



APPLIED GENOMICS

DNA SEQUENCING

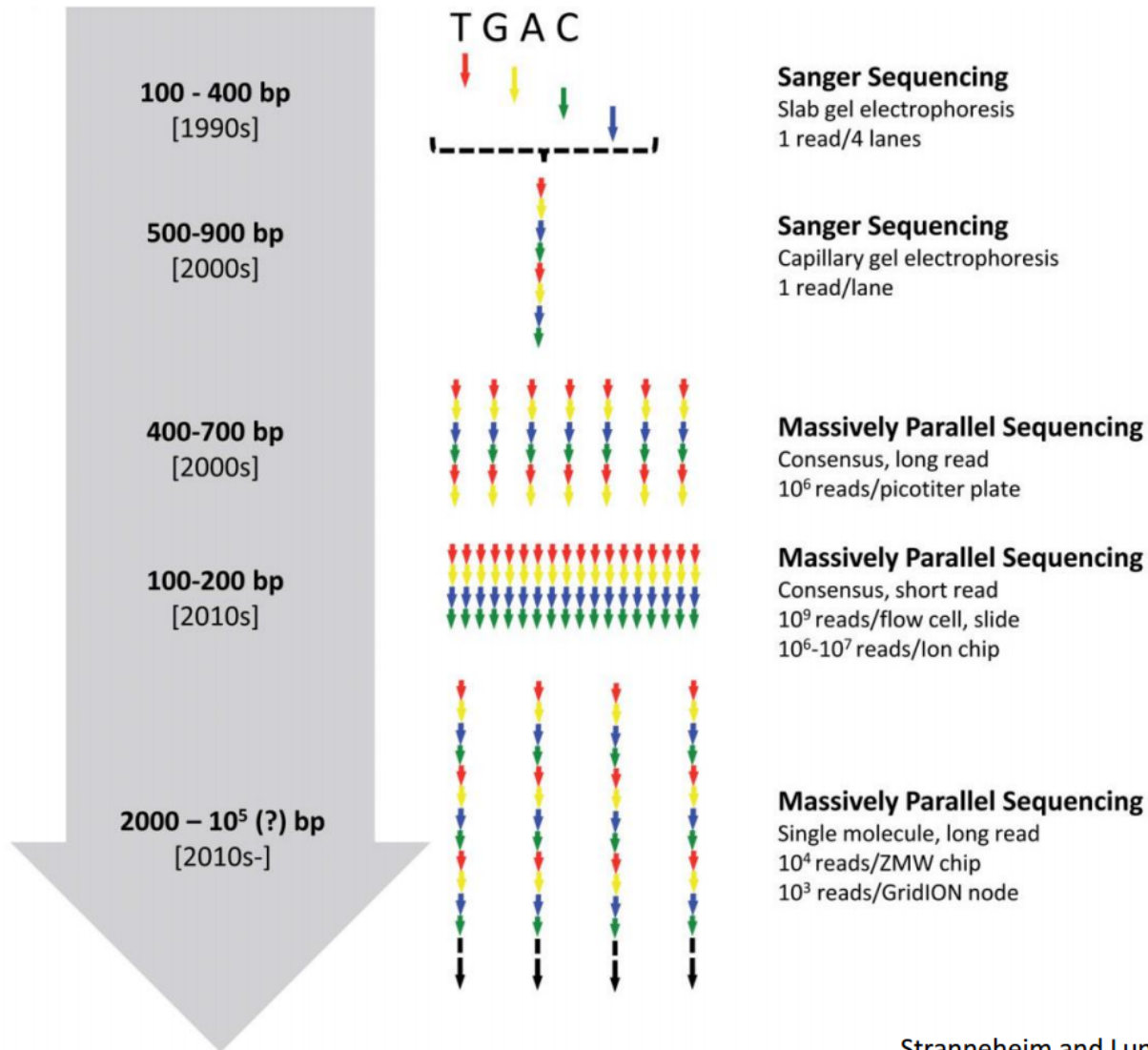
Prof. Alberto Pallavicini
pallavic@units.it

GENOMICS ANALYSIS

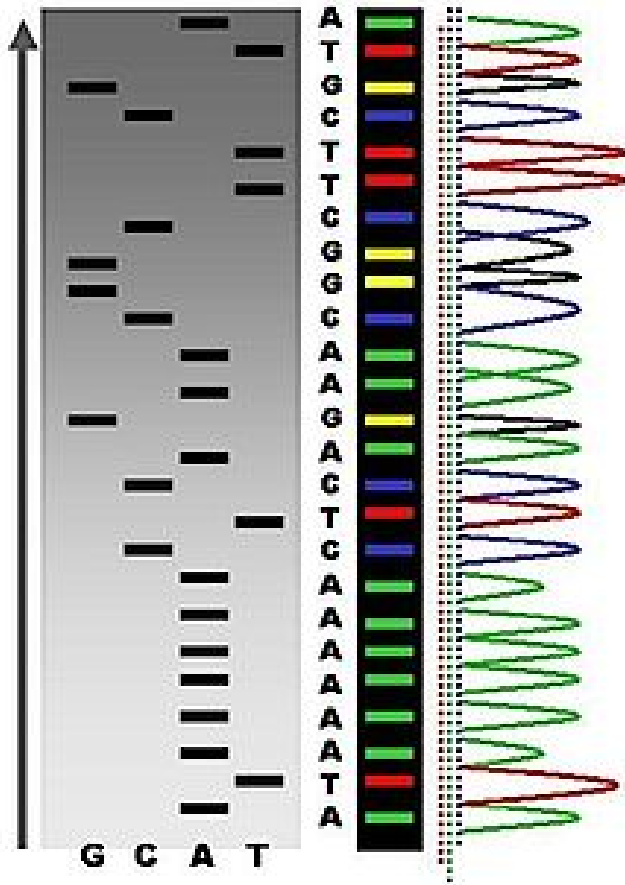
- Genomics is highly technology driven. The enormous impact of genomics research on medical and agro-technological sciences has inspired commercial life-science companies to develop innovative genomic tools at a tremendously high speed.



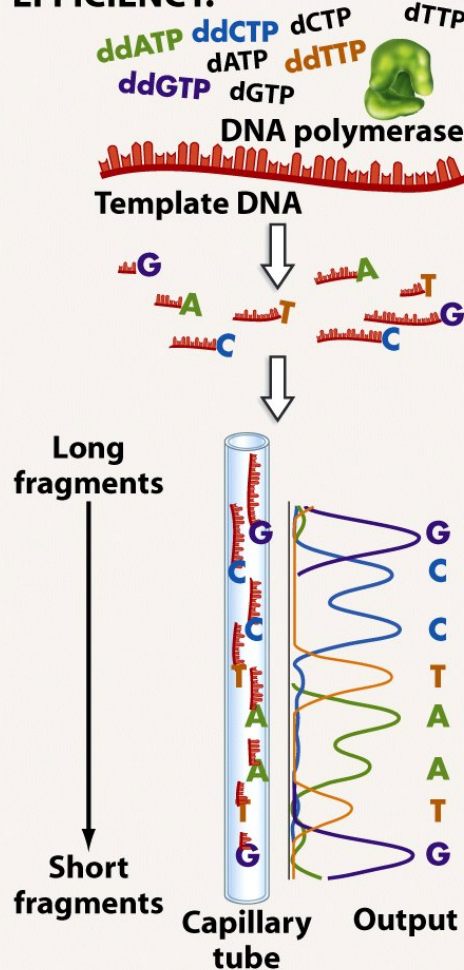
GENOMICS ANALYSIS



SANGER SEQUENCING



FLUORESCENT MARKERS IMPROVE SEQUENCING EFFICIENCY.



1. Do one sequencing reaction instead of four. Reaction mix contains ddATP, ddTTP, ddGTP, ddCTP with distinct fluorescent markers. (With radioactive labels, four reactions are needed—one labeled ddNTP at a time.)

2. Fragments that result have distinctive labels.

3. Separate fragments via electrophoresis in mass-produced, gel-filled capillary tubes. Automated sequencing machine reads output.

SANGER SEQUENCING

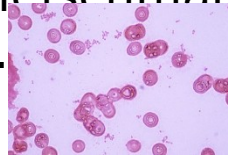
- Advantages
 - Long reads (~900bps)
 - Suitable for small projects

- Disadvantages
 - Low throughput
 - Expensive



SANGER SEQUENCING

1994: *H. Influenzae*
1.8 Mbp
(Fleischmann et al.)



2007: Global Ocean Sampling Expedition
~3,000 organisms,
7Gb (Riesenfeld et al.)



1980

1990

2000



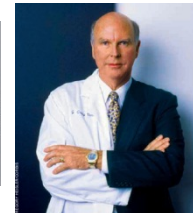
1982: *lambda virus*

DNA stretches up to
30-40Kbp
(Sanger et al.)



2001: *H. Sapiens, D. Melanogaster*

3 Gbp
(Venter et al.)



NEXT GENERATION SEQUENCING: WHY NOW?

- **Motivation:** HGP and its derivatives, personalized medicine
- **Short reads applications:** (re-)sequencing, other methods (e.g. gene expression)
- Advancements in technology



HIGH PARALLELISM IS ACHIEVED IN POLONY SEQUENCING

Sanger

Polony

Cyclic array sequencing ($>10^6$ reads/array)

Cycle 1



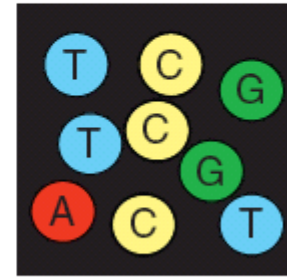
What is base 1?

Cycle 2



What is base 2?

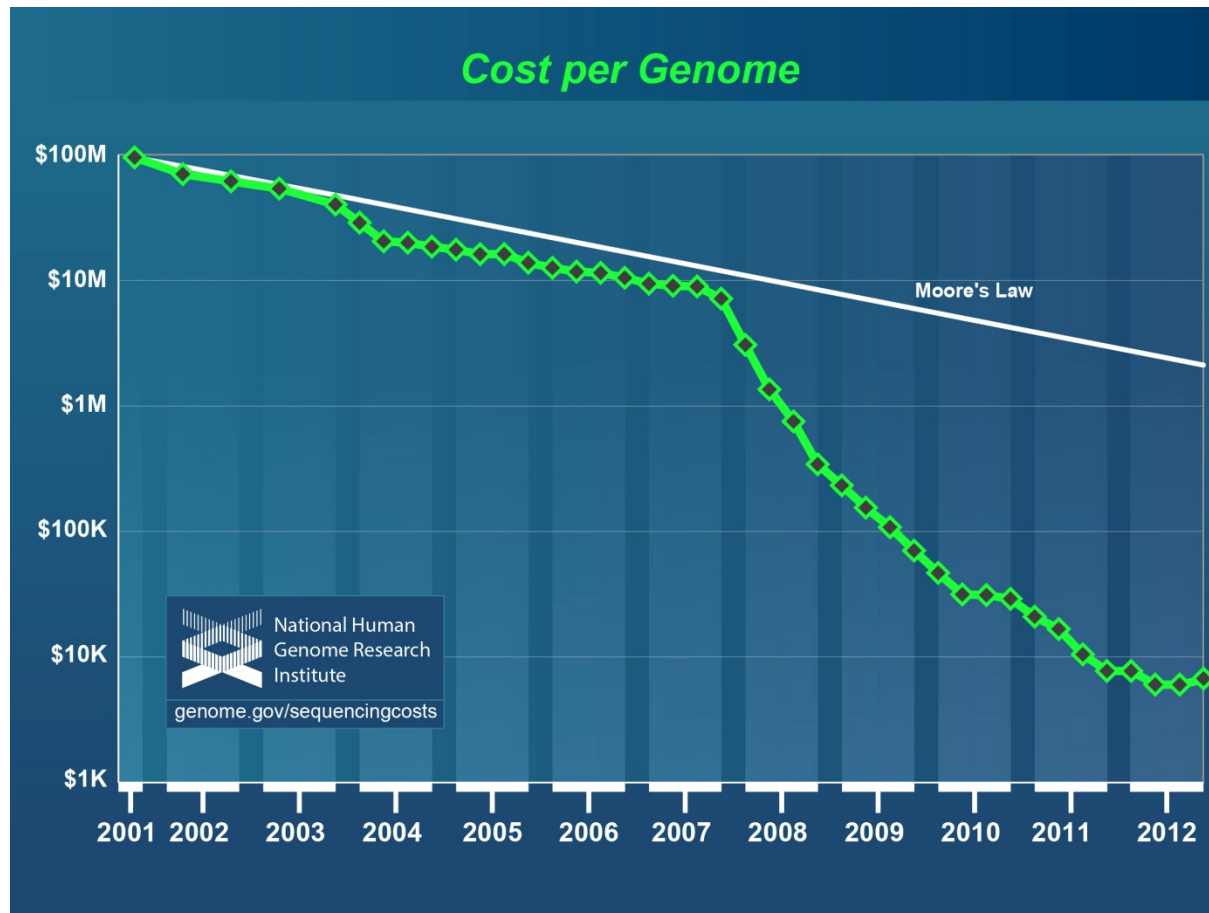
Cycle 3



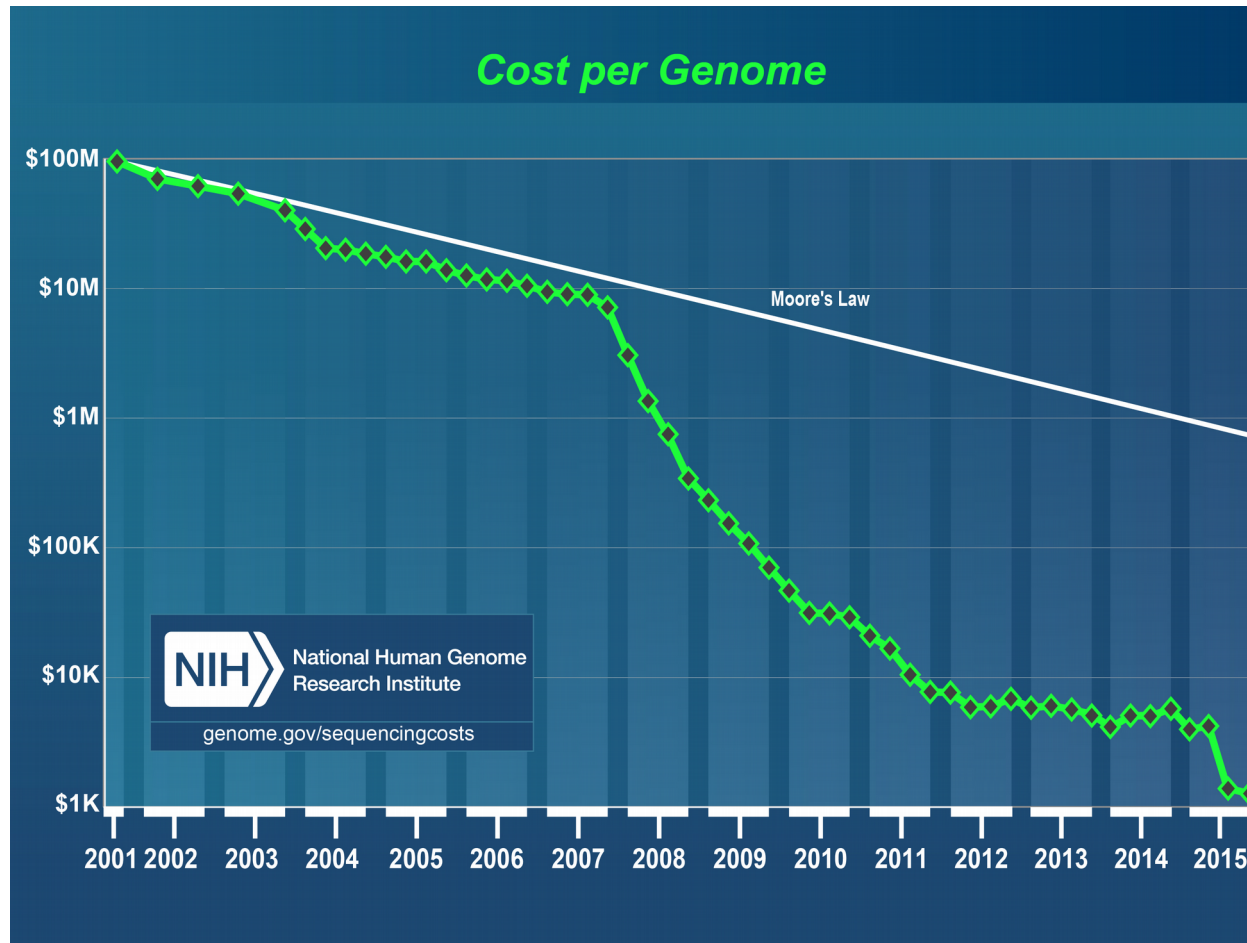
What is base 3?



Sequencing costs have fallen



Sequencing costs have fallen



HIGH PARALLELISM IS ACHIEVED IN POLONY SEQUENCING

Perchè Next Generation Sequencing






Si possono generare centinaia di milioni di corte sequenze (35bp-250bp) in una sola corsa in un tempo breve con un basso prezzo per base sequenziata.

2000

- Illumina HiSeq 2500, MiSeq, Next seq 500
- Life Technologies Ion Proton/Ion PGM
- Applied Biosystems SOLiD e Roche/454 FLX, Titanium



ILLUMINA MACHINES

	 MiniSeq System	 MiSeq Series	 NextSeq Series	 HiSeq Series	 HiSeq X Series*
Key Methods	Amplicon, targeted RNA, small RNA, and targeted gene panel sequencing.	Small genome, amplicon, and targeted gene panel sequencing.	Everyday exome, transcriptome, and targeted resequencing.	Production-scale genome, exome, transcriptome sequencing, and more.	Population- and production-scale whole-genome sequencing.
Maximum Output	7.5 Gb	15 Gb	120 Gb	1500 Gb	1800 Gb
Maximum Reads per Run	25 million	25 million [†]	400 million	5 billion	6 billion
Maximum Read Length	2 × 150 bp	2 × 300 bp	2 × 150 bp	2 × 150 bp	2 × 150 bp
Run Time	4–24 hours	4–55 hours	12–30 hours	<1–3.5 days (HiSeq 3000/HiSeq 4000) 7 hours–6 days (HiSeq 2500)	<3 days
Benchtop Sequencer	Yes	Yes	Yes	No	No
System Versions	<ul style="list-style-type: none"> MiniSeq System for low-throughput targeted DNA and RNA sequencing 	<ul style="list-style-type: none"> MiSeq System for targeted and small genome sequencing MiSeq FGx System for forensic genomics MiSeqDx System for molecular diagnostics 	<ul style="list-style-type: none"> NextSeq 500 System for everyday genomics NextSeq 550 System for both sequencing and cytogenomic arrays 	<ul style="list-style-type: none"> HiSeq 3000/HiSeq 4000 Systems for production-scale genomics HiSeq 2500 Systems for large-scale genomics 	<ul style="list-style-type: none"> HiSeq X Five System for production-scale whole-genome sequencing HiSeq X Ten System for population-scale whole-genome sequencing



ION TORRENT MACHINES



Ion S5 Systems



Ion PGM System

<https://www.thermofisher.com/it/en/home/brands/ion-torrent.html>





Next generation sequencing

	Run Time	Read Length	Quality	Total nucleotides sequenced	Cost /MB
454 Pyrosequencing	24h	700 bp	Q20-Q30	1 GB	\$10
Illumina Miseq	27h	2x300bp	> Q30	15 GB	\$0.15
Illumina Hiseq 2500	1 - 10days	2x250bp	>Q30	3000 GB	\$0.05
Ion torrent	2h	400bp	>Q20	50MB-1GB	\$1
Pacific Biosciences	30m - 4h	10kb - >40kb	>Q50 consensus >Q10 single	500 - 1000MB /SMRT cell	\$0.13 - \$0.60

<http://www.hindawi.com/journals/bmri/2012/251364/>
<http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3431227>

circa 2015

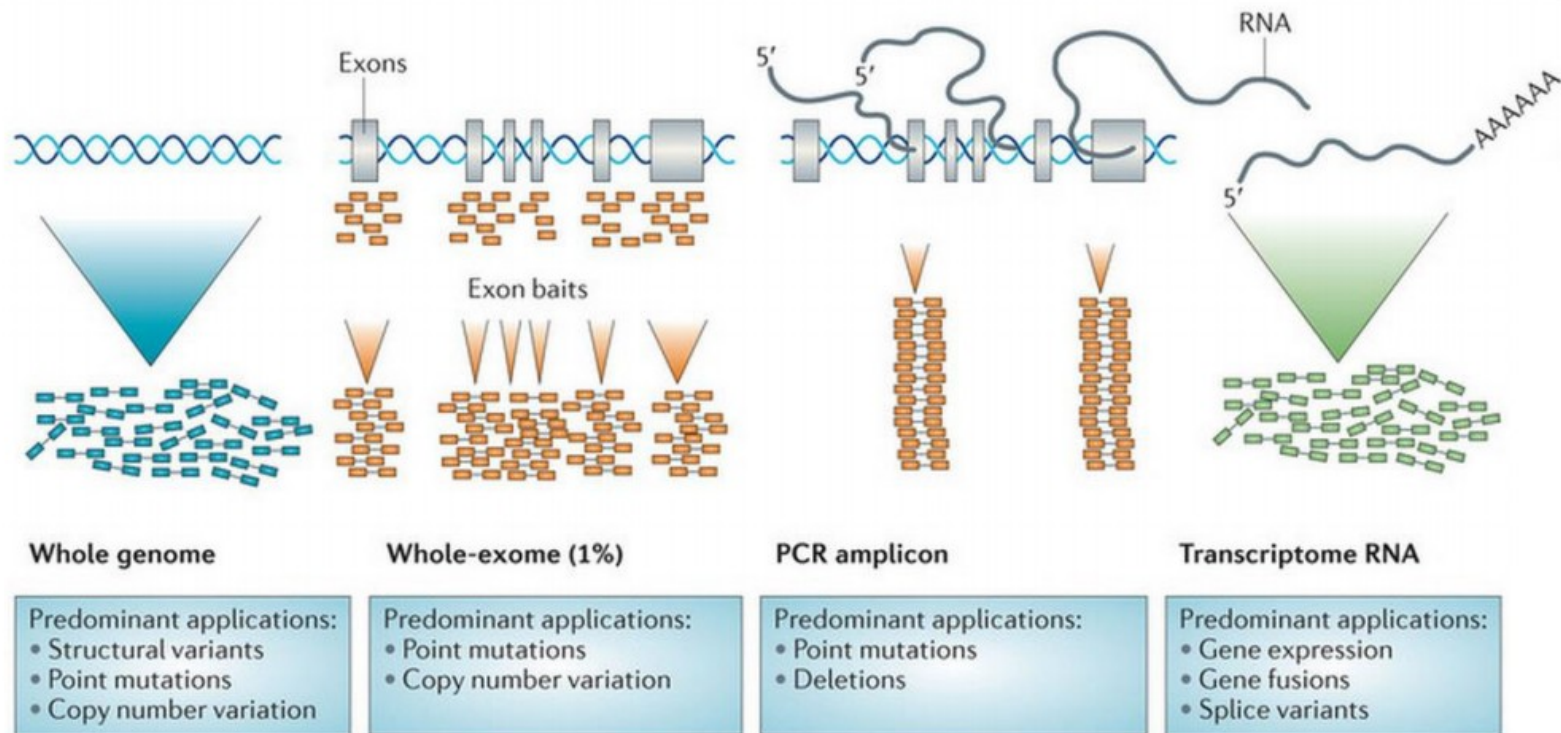


Terminology

- **Coverage (depth):** the number of times a nucleotide is read during the sequencing process
- **Quality Score:** Each called base comes with a quality score which measures the probability of base call error.
- **Paired-End Sequencing:** Both end of the DNA fragment is sequenced, allowing highly precise alignment.
- **Multiplex Sequencing:** "barcode" sequences are added to each sample so they can be distinguished in order to sequence large number of samples on one lane.
- **Mapping:** Align reads to reference to identify their origin.
- **Assembly:** Merging of fragments of DNA in order to reconstruct the original sequence.
- **Duplicate reads:** Reads that are identical.
- **Multi-reads:** Reads that can be mapped to multiple locations equally well.



Applications: genomes, exomes, transcriptomes



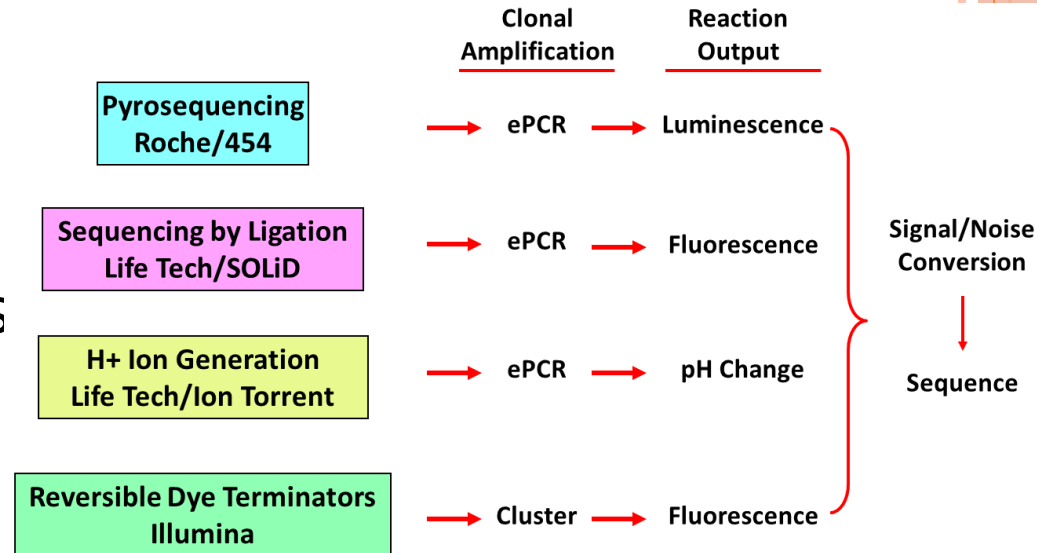
Implementing personalized cancer genomics in clinical trials

Richard Simon & Sameek Roychowdhury

Nature Reviews Drug Discovery 12, 358–369 (2013) | doi:10.1038/nrd3979

NGS PLATFORMS OVERVIEW

- Differ in design and chemistries
- Fundamentally related-sequencing of thousands to millions of clonally amplified molecules in a massively parallel manner
- Orders of magnitude more information-will continue to evolve



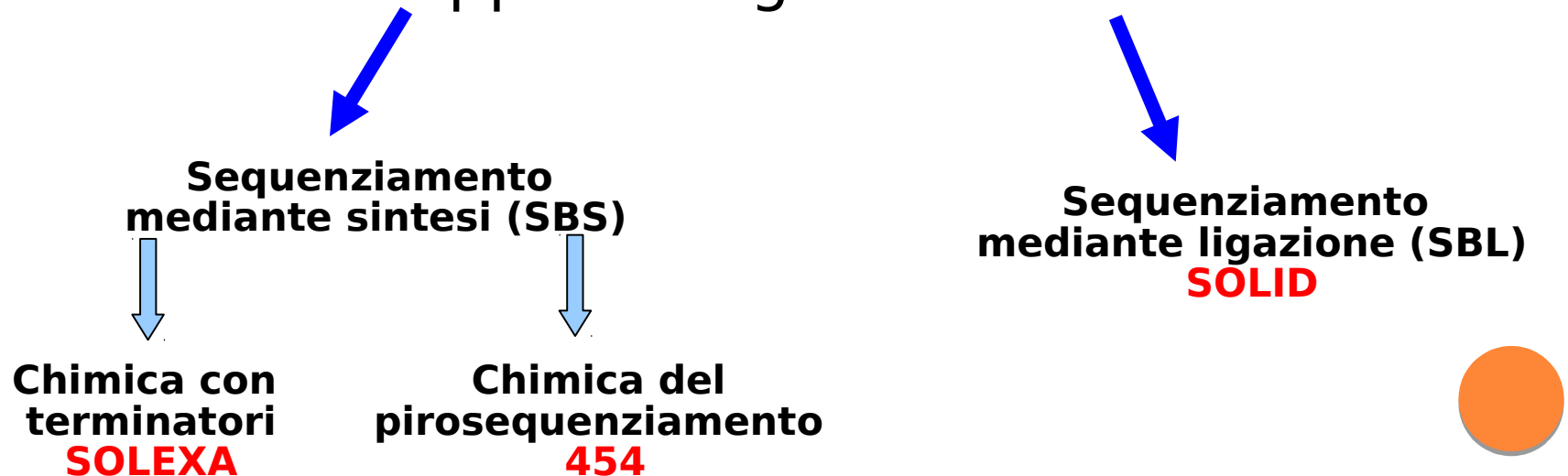
Pacific Biosciences
Helicos Biosciences
NABsys
VisiGen Biotechnologies
Complete Genomics
Oxford Nanophore
Technologies



SEQUENZIAMENTO DI NUOVA GENERAZIONE

Si basano sul principio del sequenziamento di *'cluster'* clonali

Il processo, che incomincia con una singola molecola target, prevede la creazione di targets clonali durante un processo intermedio di amplificazione. Copie multiple identiche sono infatti necessarie per avere un alto rapporto segnale-rumore



SEQUENZIAMENTO SANGER AD ALTA PROCESSIVITÀ

PREPARAZIONE DELLA LIBRERIA

Frammentazione casuale del DNA genomico
clonazione e trasformazione in batteri

7-10 giorni

assumendo di possedere
una piattaforma robotica
per alta processività

Raccolta delle colonie

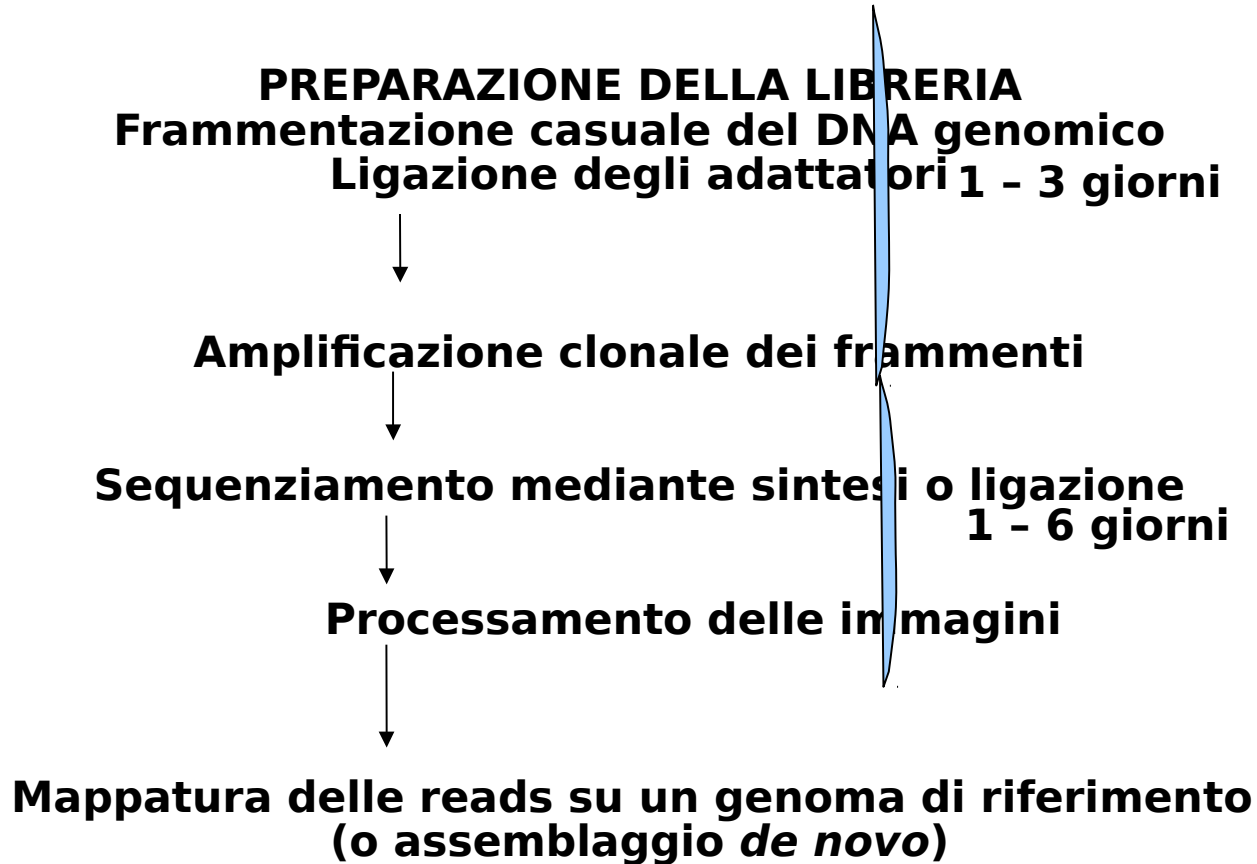
Purificazione del DNA dalle colonie
Sequenziamento Sanger
Elettroforesi capillare

Settimane-anni (!),
dipendentemente
dalla dimensione del
genoma (e copertura
richiesta), dal numero
di sequenziatori
capillari


Whole genome *de novo* assembly or mapping
to a reference (re-sequencing)



SEQUENZIAMENTO DI NUOVA GENERAZIONE



Vantaggi delle piattaforme di nuova generazione

- Non sub-clonazione, non utilizzo di cellule batteriche *E. coli*
 - abolizione di *bias* di clonazione
 - rapidità nel preparare le librerie (non c'è colony picking!)
 - Ciascuna sequenza proviene da una molecola di DNA unica.
 - quantificazione attraverso 'conta' digitale
 - aumento del range dinamico
 - rilevazione di varianti rare
 - Fornisce una eccezionale risoluzione per molti tipi di esperimenti (es. analisi di espressione, sequenziamento di DNA immunoprecipitato, di RNA piccoli, analisi di medie/grandi inserzioni-delezioni nei genomi....)
 - Rivoluzionaria diminuzione del costo e del tempo per generare dati di sequenza (lavorano in multi-parallelo)
 - Richiesta meno robotica nelle fasi precedenti al caricamento sul sequenziatore
- 

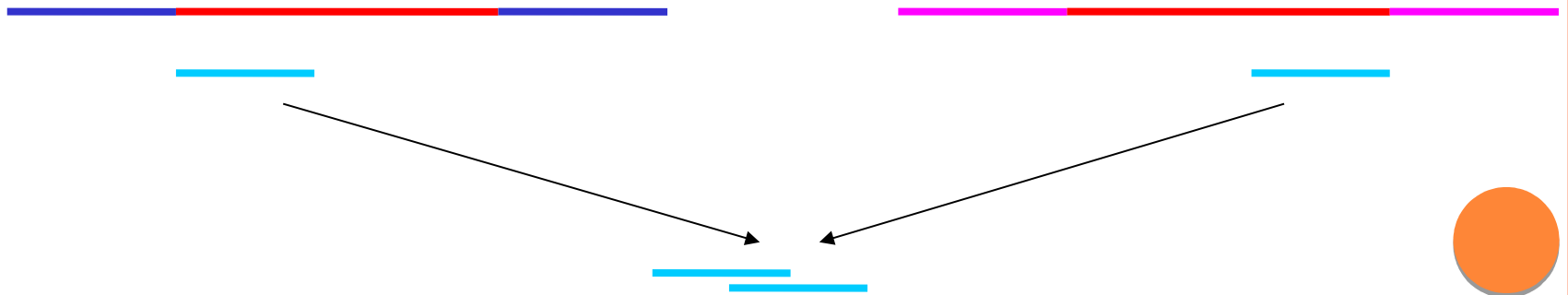
Svantaggi delle piattaforme next-gen

- Sono prodotte sequenze più corte
 - relativamente alle sequenze da sequenziatori capillari (metodo Sanger)
 - è necessario ri-parametrizzare l'accuratezza della procedura di chiamata delle basi
 - enorme difficoltà nell'analisi dei dati; richiesto un grande sforzo di programmazione per costruire nuovi algoritmi.
- La mole enorme di dati 'traumatizza' le infrastrutture informatiche.
 - da 10 Gb a diversi Tb di dati grezzi prodotti per corsa (dipende dalla piattaforma)
 - il processamento delle *read* tramite *pipeline* informatiche richiede molta capacità di calcolo (CPU)
 - è necessario prendere accurate decisioni su cosa salvare e cosa cancellare



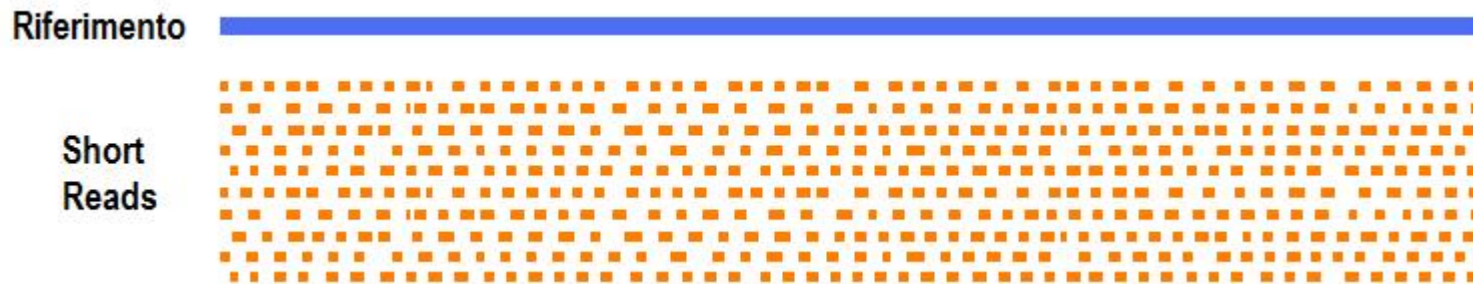
SEQUENZE CORTE

- Sequenze corte, ma tecnologia in continua evoluzione:
 - 454: 100 basi → 200 → **400-500** → ?
 - Solid: 25 basi → 35 → **50** → 100 → ?
 - Illumina: 32 → 36 → 75-100 → **125 → 150** → ?
- Difficoltà di assemblare sequenze corte *de novo*, soprattutto per il problema delle sequenze ripetute complicato ancora di più rispetto a Sanger (lunghezza media 700-900bp)

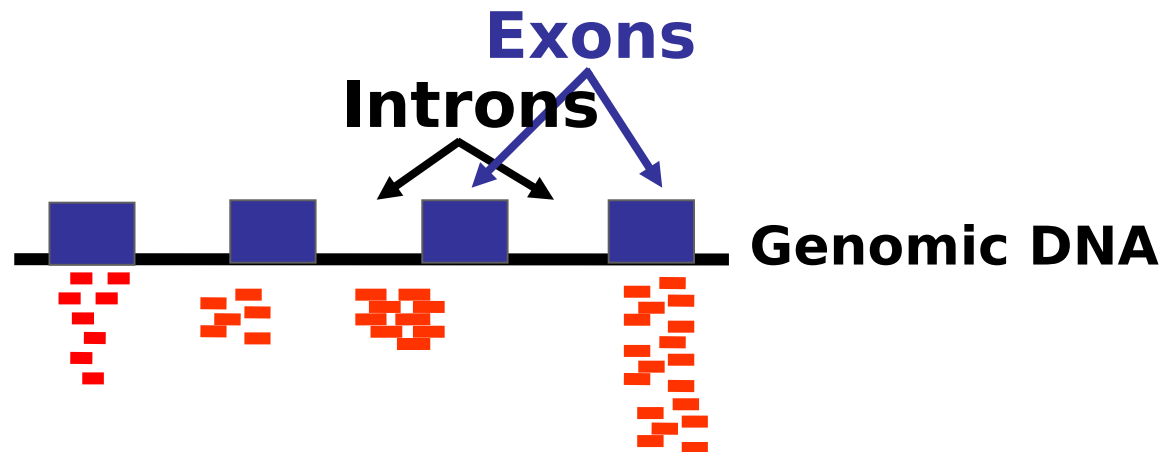


RISEQUENZIAMENTO

- In presenza di un genoma di riferimento di buona qualità posso effettuare un ri-sequenziamento e allineare tutte le reads ottenute:



• Non solo del genoma, ma anche del trascrittoma



VERSO IL GENOMA DA MILLE DOLLARI....

- **Costo 1 anno Sanger, reads 700bp:**
 - 1 anno, 1 sequenziatore a pieno regime=260 Mbp
 - 260 Mbp=circa 370.000 sequenze (lunghezza media 700bp)=370.000 EUR
 - EUR/base=0,0014
 - Sequenziamento di un genoma batterico (es E. coli, 4.5Mbp) con copertura 10x=64.000 EUR
 - 1 genoma umano (dimensione 3.6 Gbp), copertura 1x=60 anni (!) =5M EUR

Costo 1 corsa 454, reads 300-400bp:

10 ore, 1 sequenziatore=fino a 0.6 Gbp

0.6 Gbp=10.000 EUR

EUR/base= 0,000016

Sequenziamento di un genoma batterico (es E. coli, 4.5Mbp) con copertura 10x= 9.600 EUR

1 genoma umano (dimensione 3.6 Gbp), copertura 10x=almeno 60 corse (più di 1 mese)=576K EUR



VERSO IL GENOMA DA MILLE DOLLARI....

Costo 1 corsa, Illumina 2x75bp:

10 giorni, 1 sequenziatore=fino a 18 Gbp

18 Gbp=10,000 EUR

EUR/base= 0,00000055

Sequenziamento di un genoma batterico con copertura 10x= 25 EUR

1 genoma umano (dimensione 3.6 Gbp), copertura 10x=2 corse=20K EUR

- **Costo 1 corsa Solid, reads 2x50bp:**

- **12 giorni, 1 sequenziatore= 20 Gbp**

- **20 Gbp = 8.000 EUR**

- **EUR/base=0,00000044**

- **Sequenziamento di un genoma batterico con copertura 10x= 18 EUR**

- **1 genoma umano (dimensione 3.6 Gbp), copertura 10x= circa 2 corse (1 mese) = 16K EUR**



**SEQUENZIAMENTO CON
LA TECNOLOGIA 454**

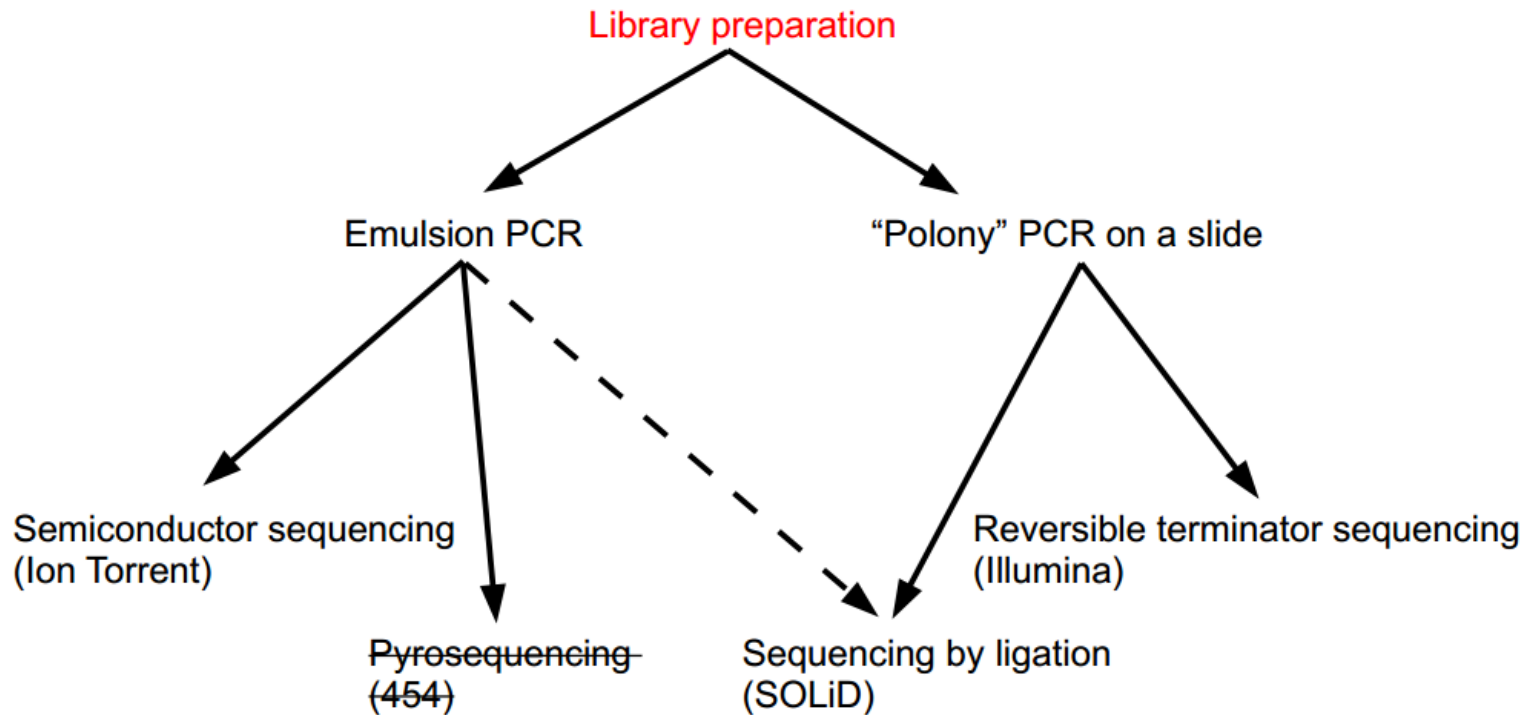
454 LIFE
SCIENCES

First to the Finish

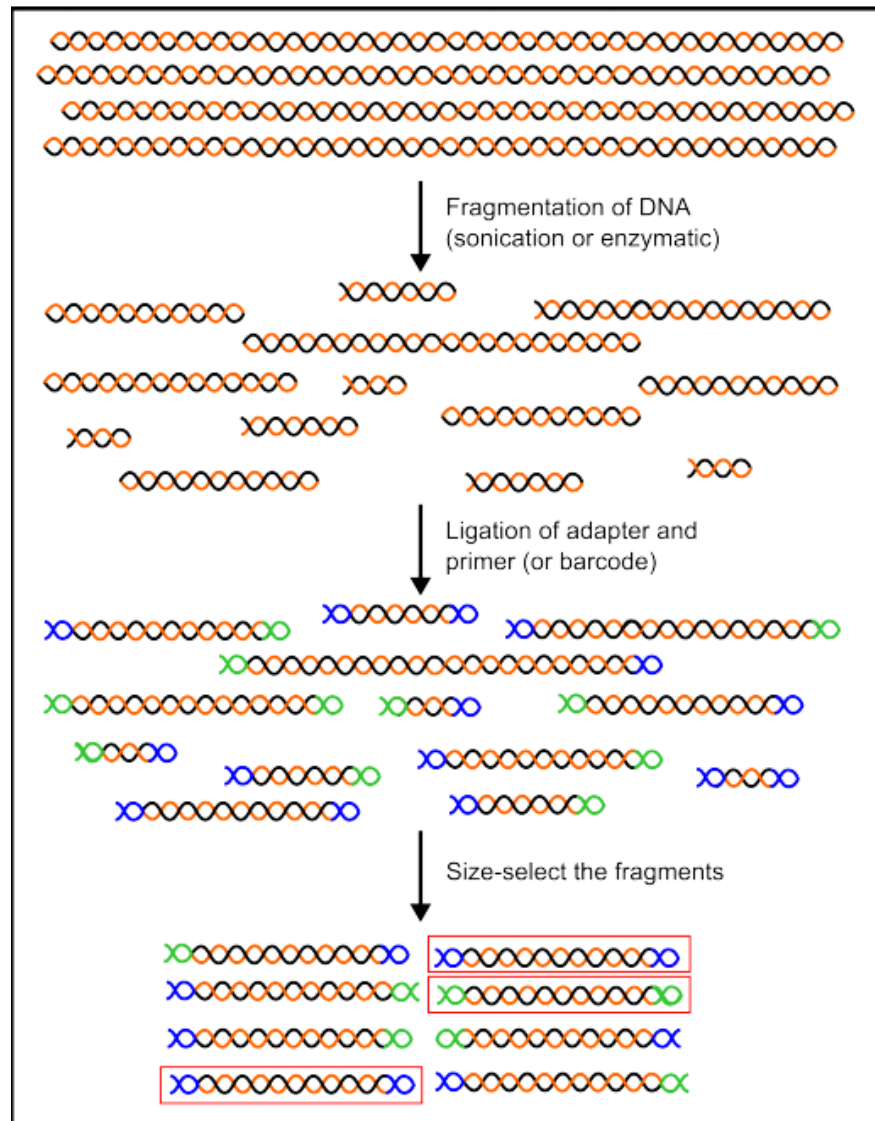


WORKFLOW

Next Generation Sequencing: Amplified Single Molecule Sequencing

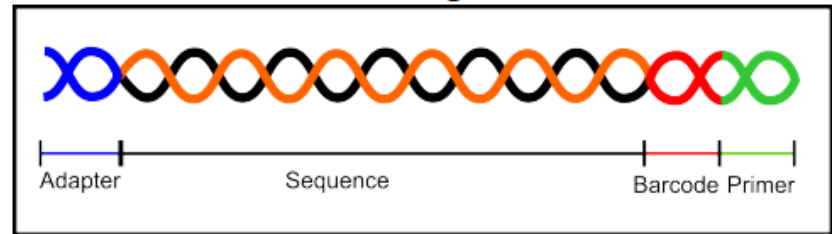


Next Generation Sequencing: Amplified Single Molecule Sequencing

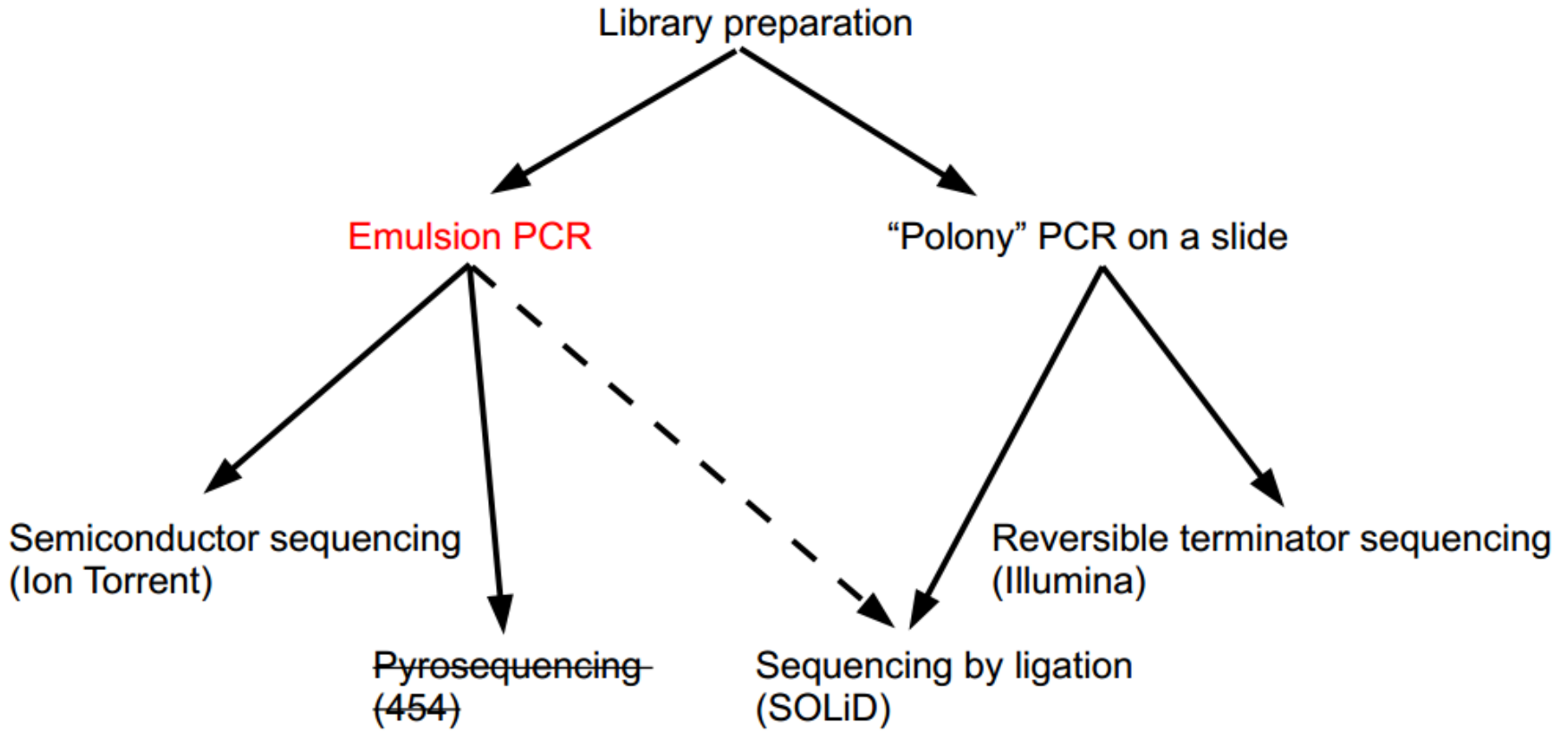


Library preparation

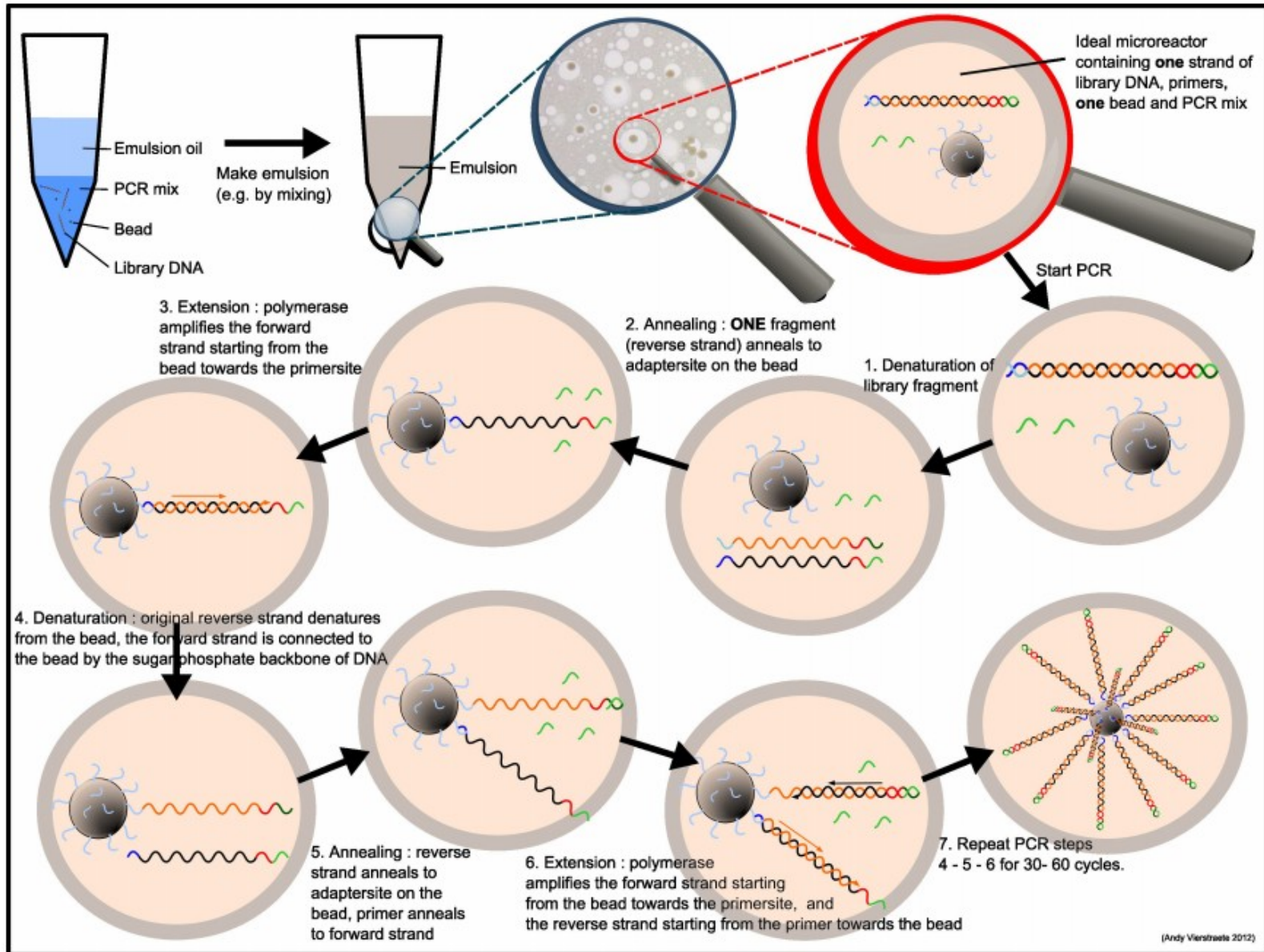
Good fragments:



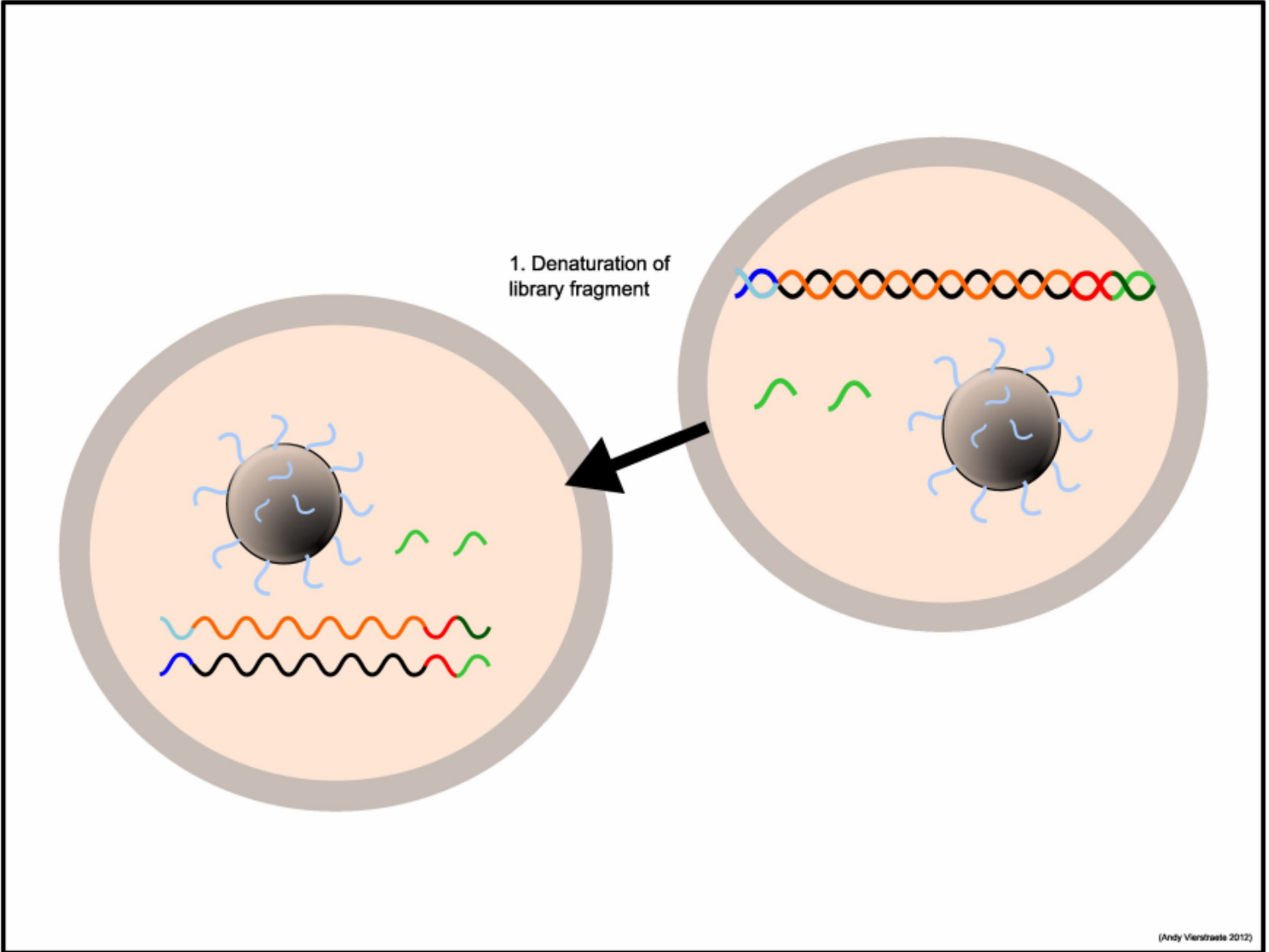
Next Generation Sequencing: Amplified Single Molecule Sequencing



Next Generation Sequencing: Amplified Single Molecule Sequencing Emulsion PCR



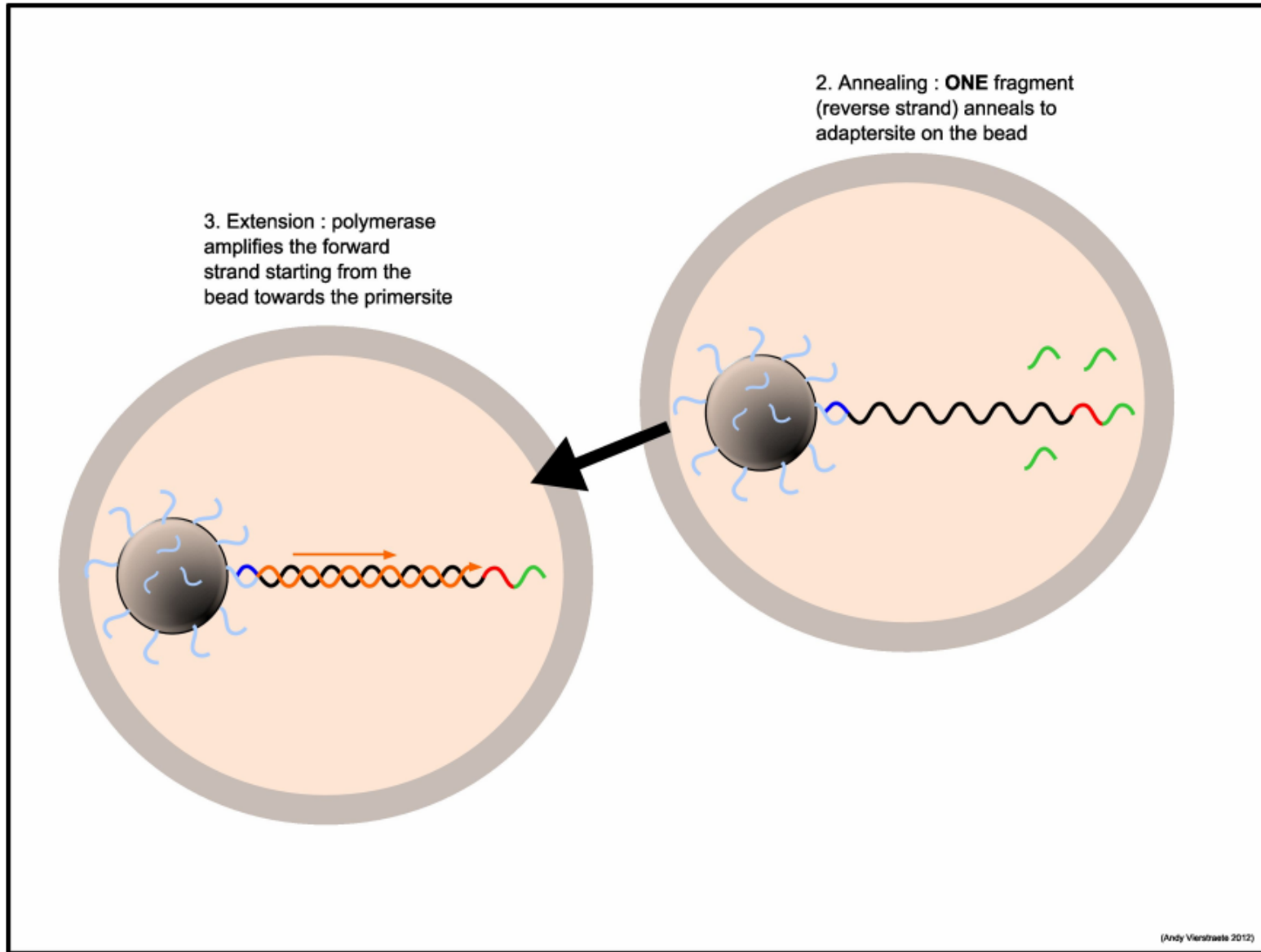
Next Generation Sequencing: Amplified Single Molecule Sequencing Emulsion PCR



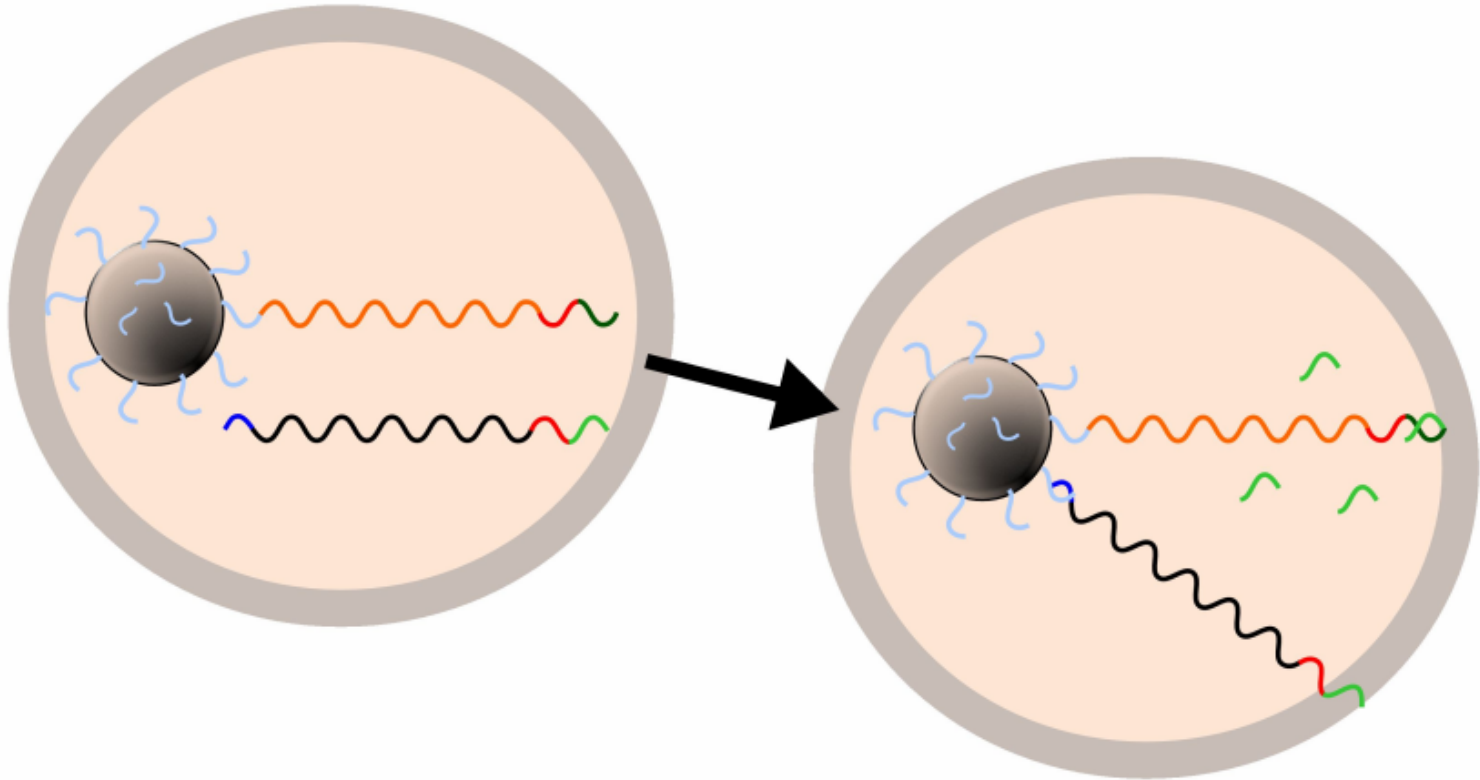
(Andy Venkates 2012)



Next Generation Sequencing: Amplified Single Molecule Sequencing Emulsion PCR

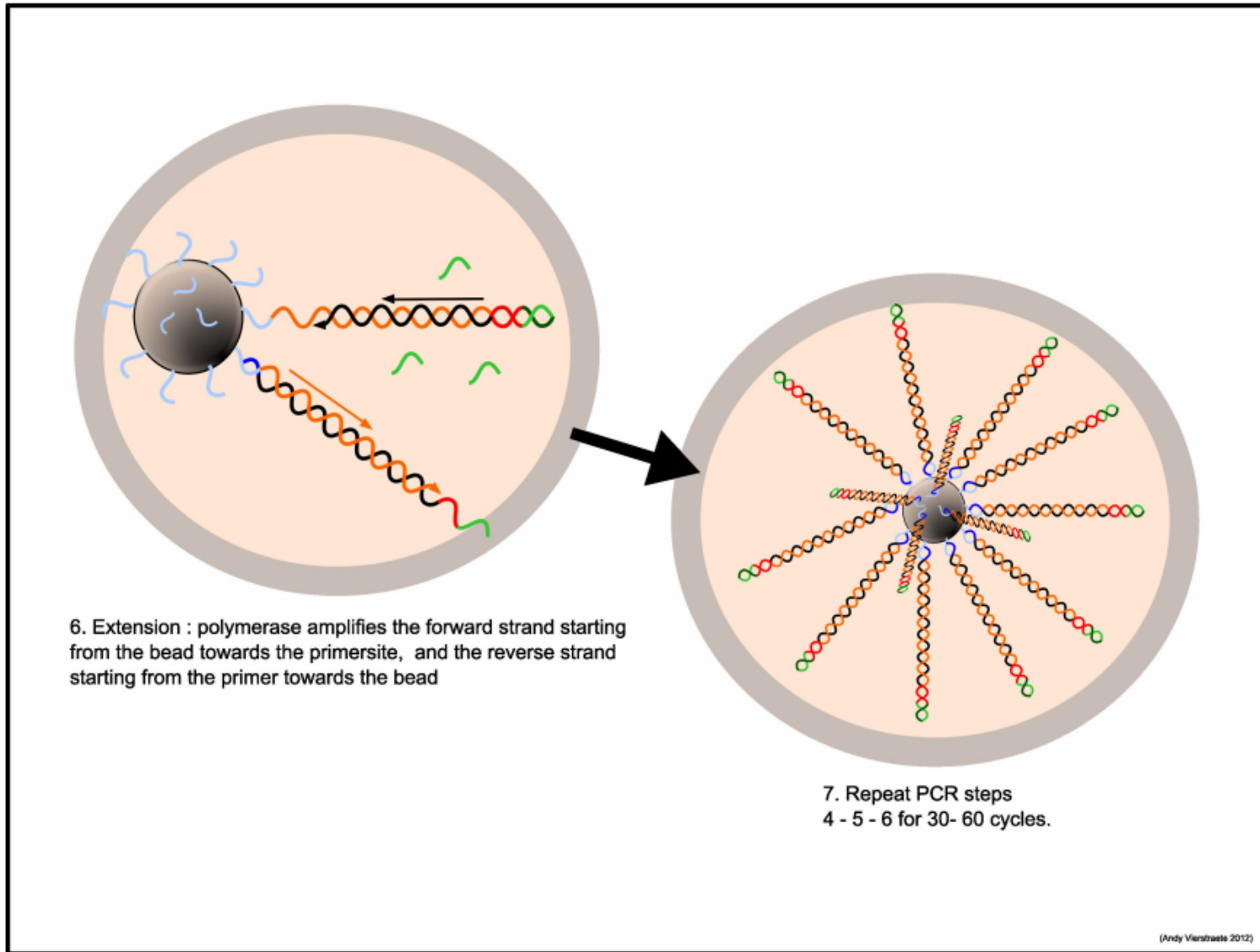


4. Denaturation : original reverse strand denatures from the bead, the forward strand is connected to the bead by the sugar phosphate backbone of DNA



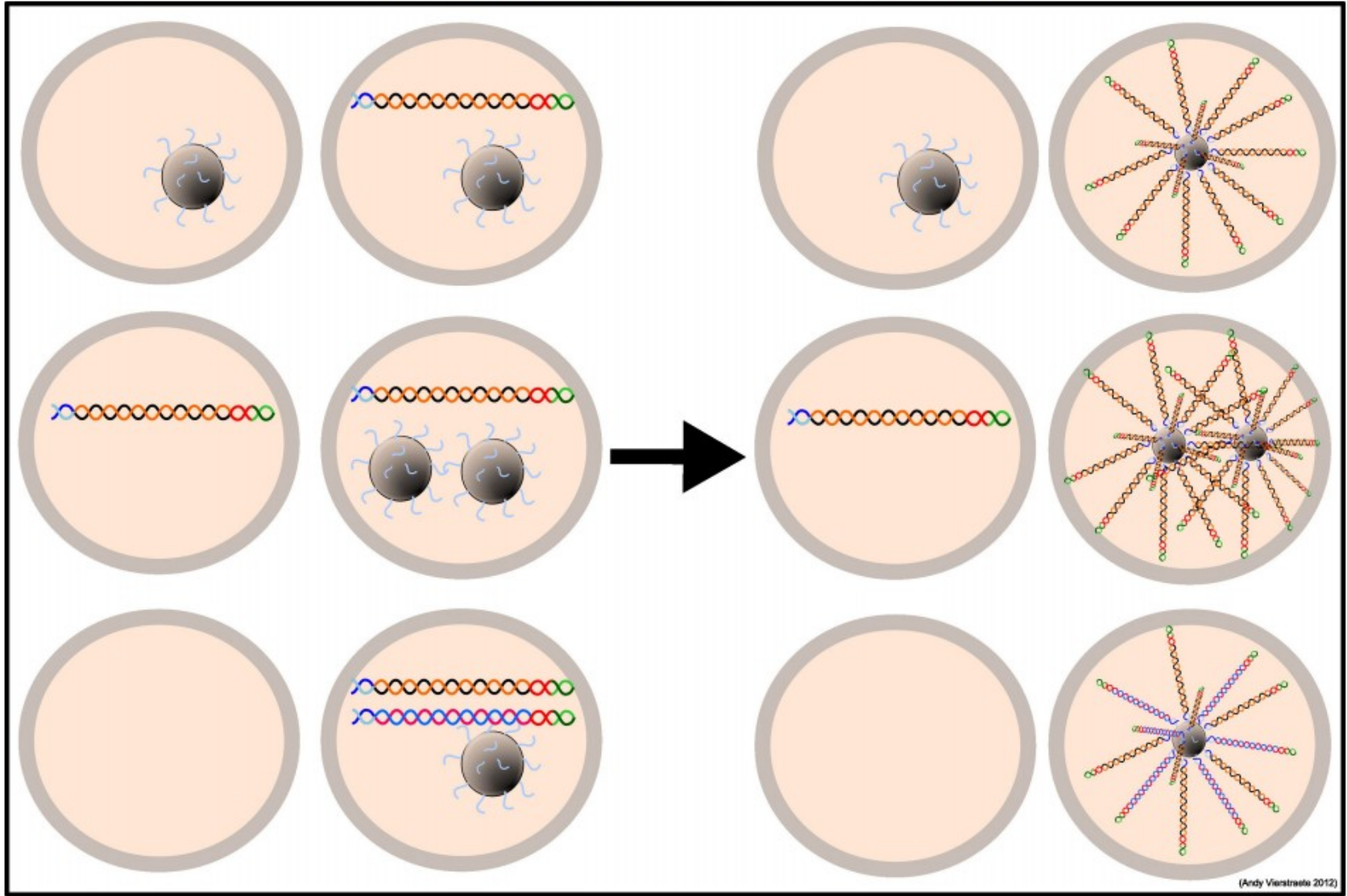
5. Annealing : reverse strand anneals to adaptersite on the bead, primer anneals to forward strand

Next Generation Sequencing: Amplified Single Molecule Sequencing Emulsion PCR

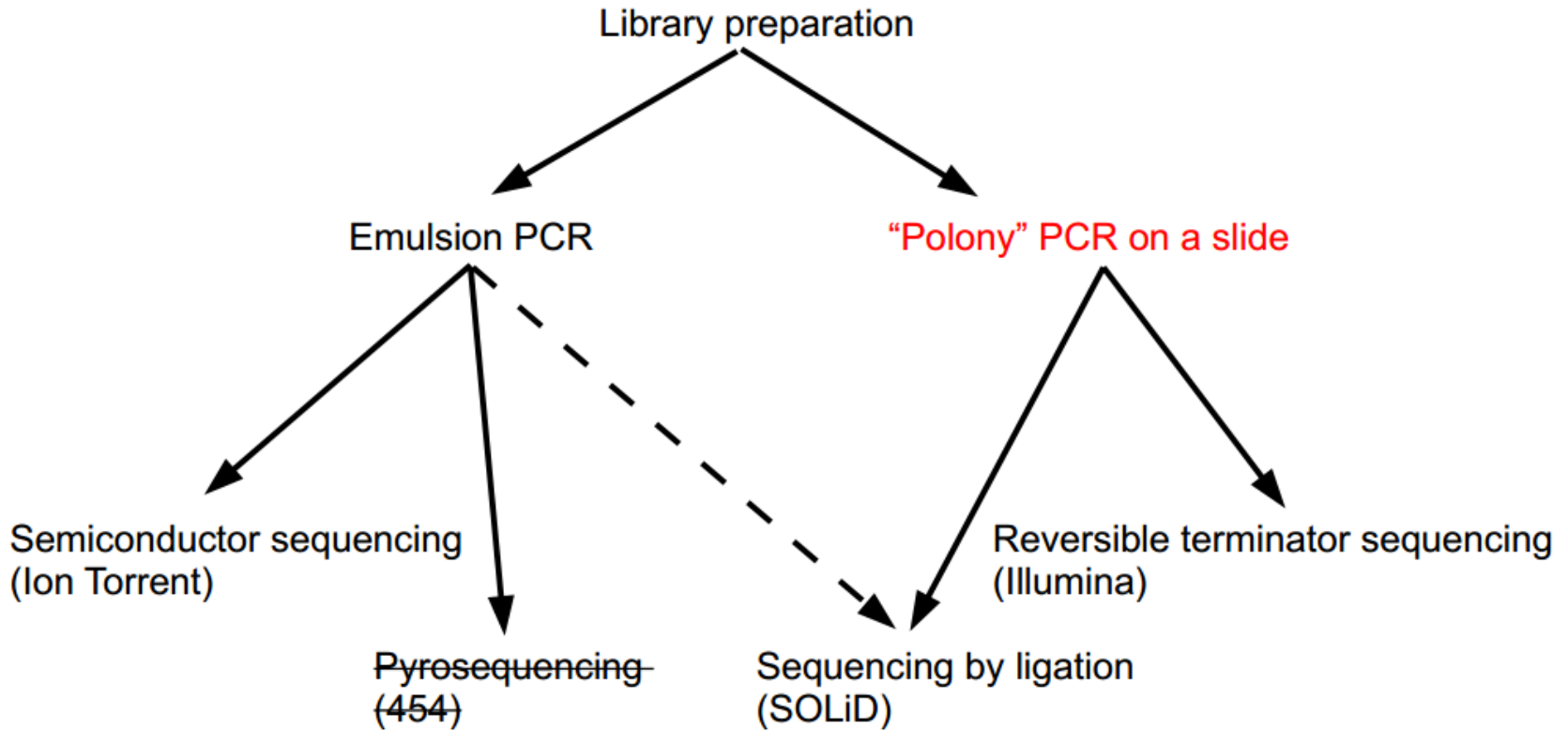


Next Generation Sequencing: Amplified Single Molecule Sequencing Emulsion PCR

different micro reactors: only 15 % are good ones

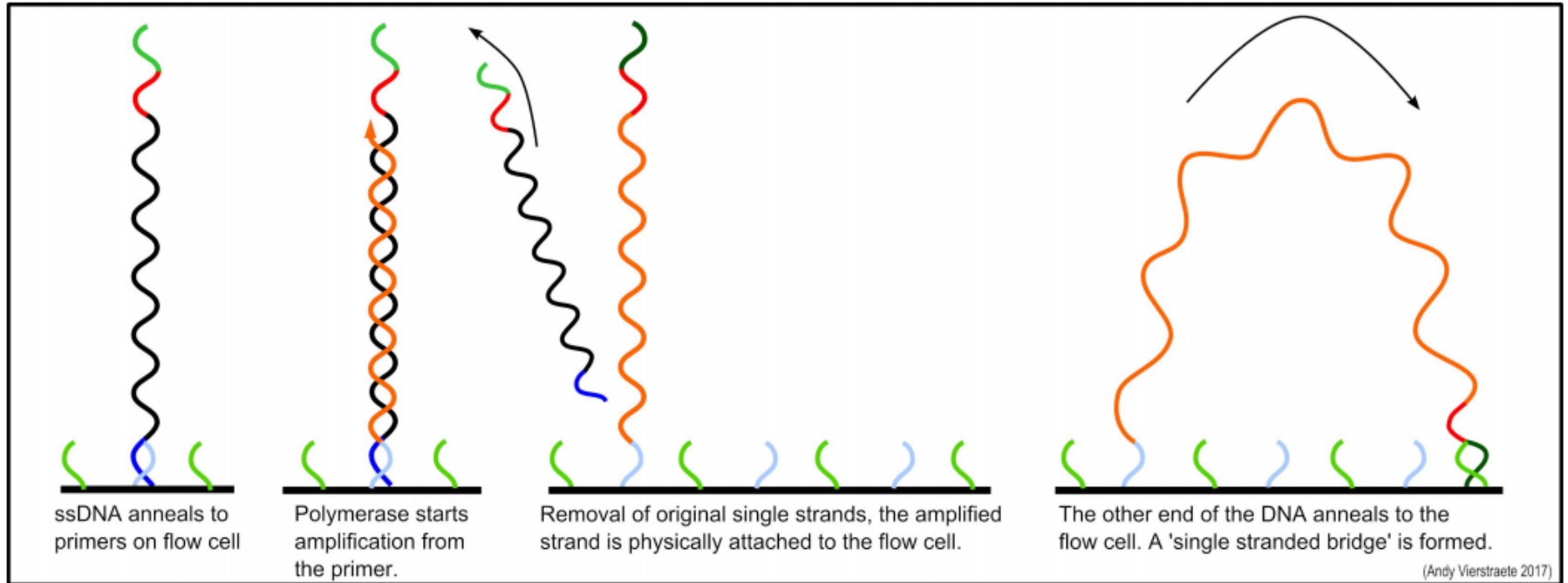


Next Generation Sequencing: Amplified Single Molecule Sequencing



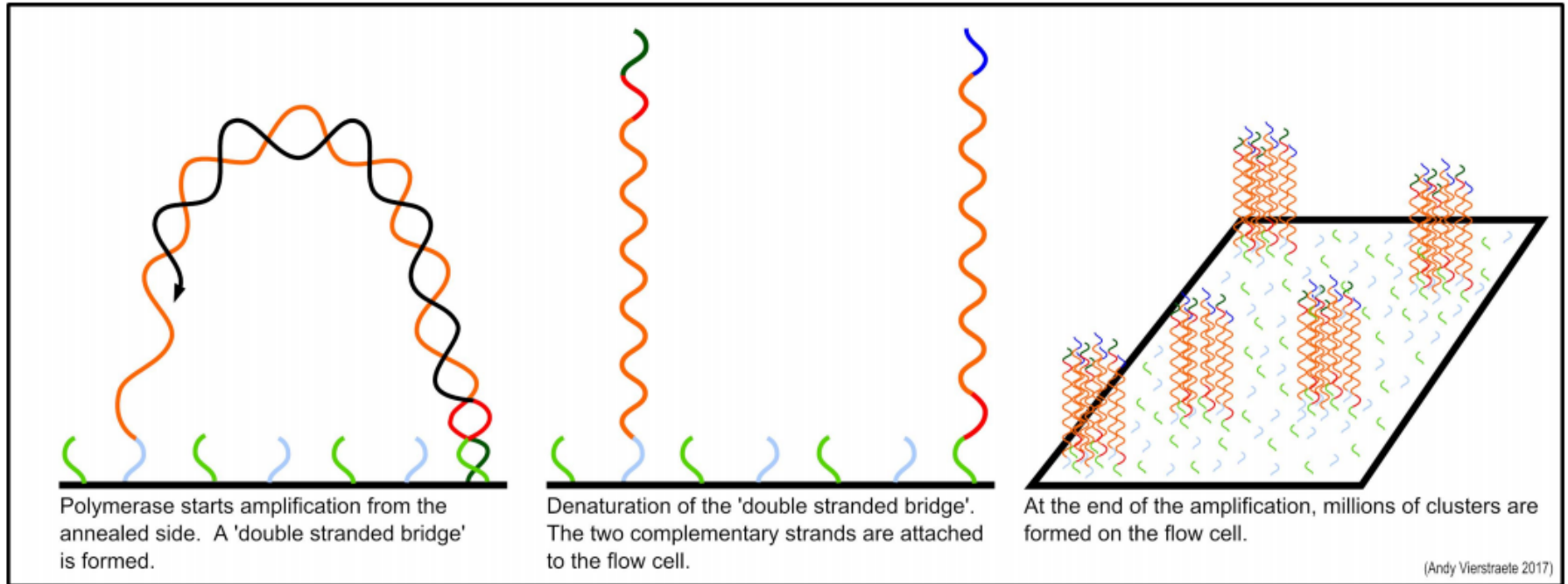
Next Generation Sequencing: Amplified Single Molecule Sequencing “Polony” PCR

Bridge amplification: Illumina

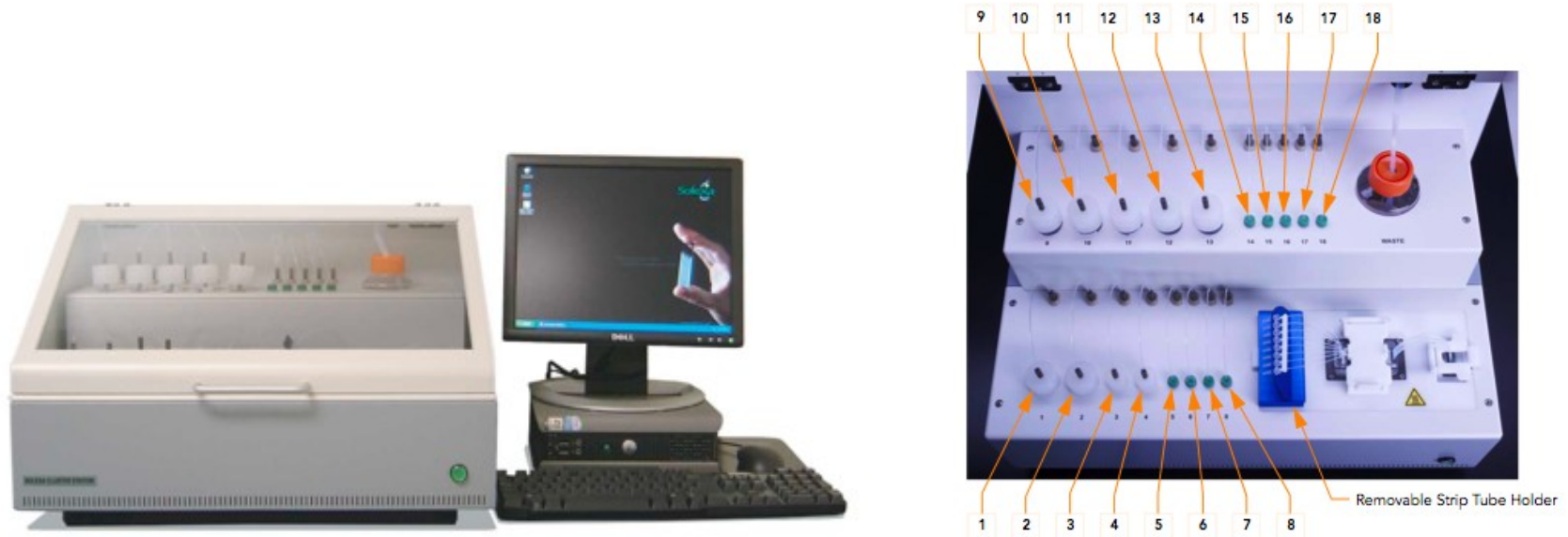


Next Generation Sequencing: Amplified Single Molecule Sequencing “Polony” PCR

Bridge amplification: Illumina



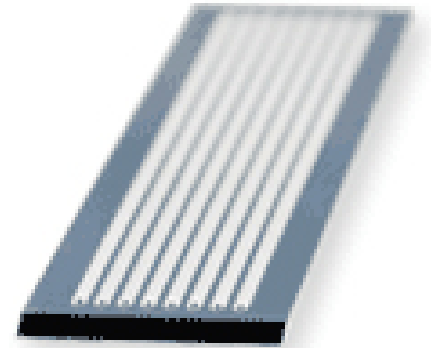
Cluster Station



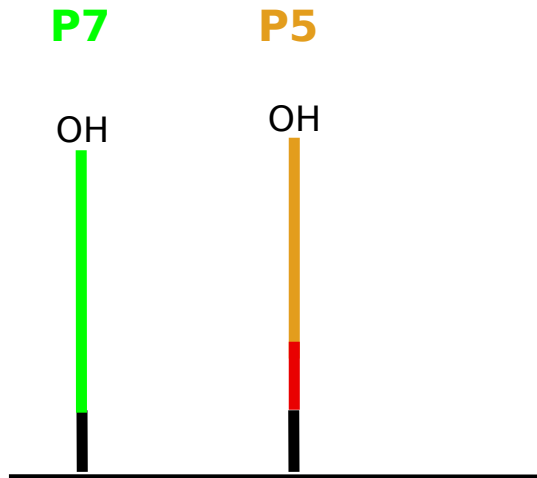
Strumento che permette di preparare la *flow-cell* (=supporto di vetro su cui i frammenti della libreria verranno sequenziati in parallelo)

Cluster Generation

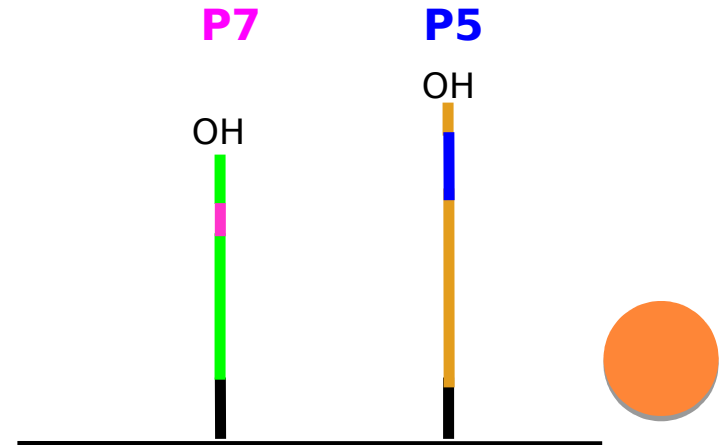
Grafted Flow Cells

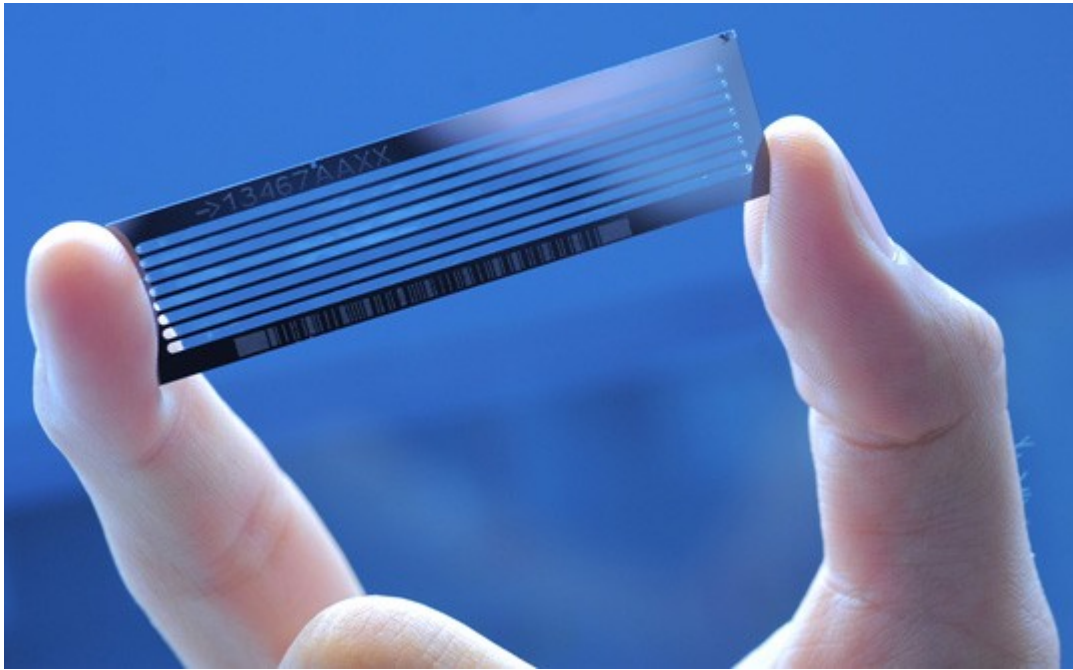


Single pass

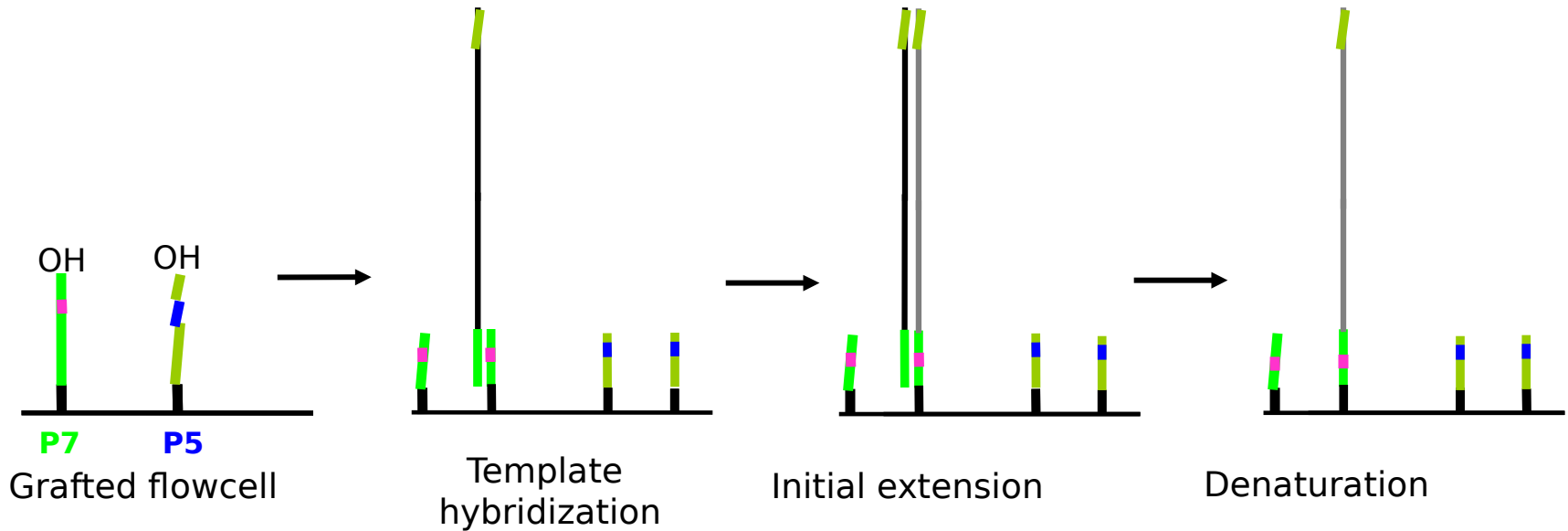


Paired end

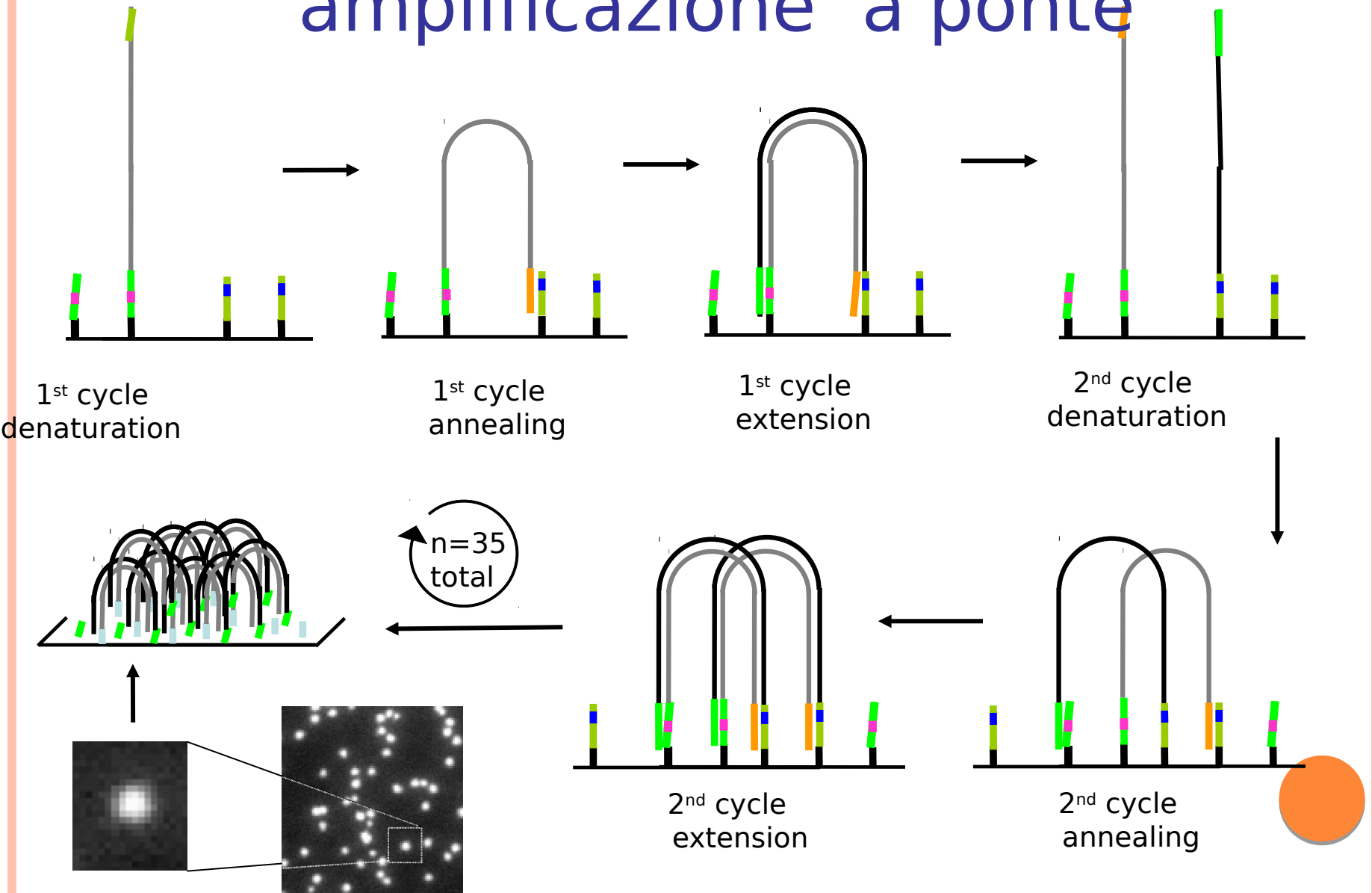




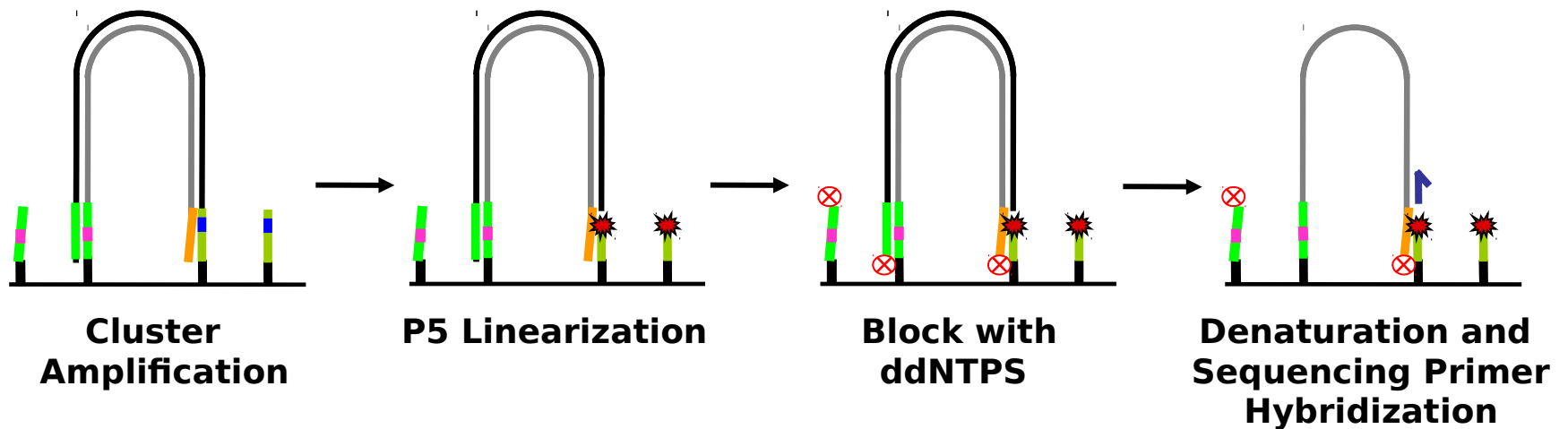
Cluster Generation



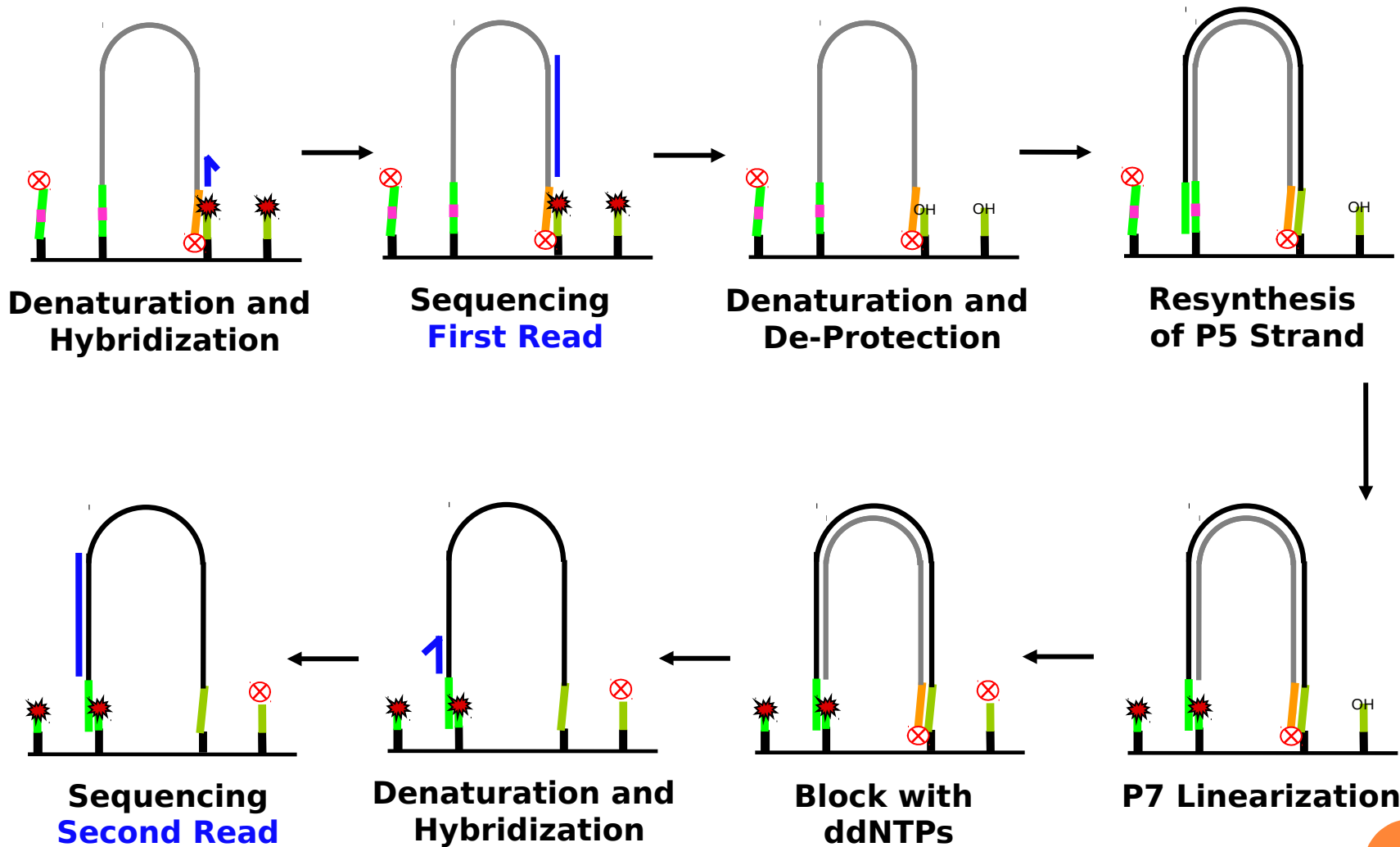
La generazione dei cluster: amplificazione a ponte



La generazione dei cluster

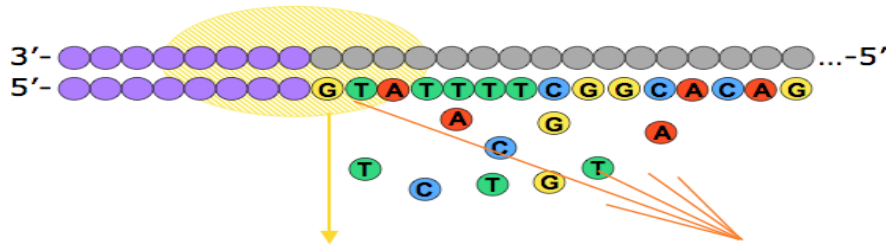


Sequenziamento pair-end

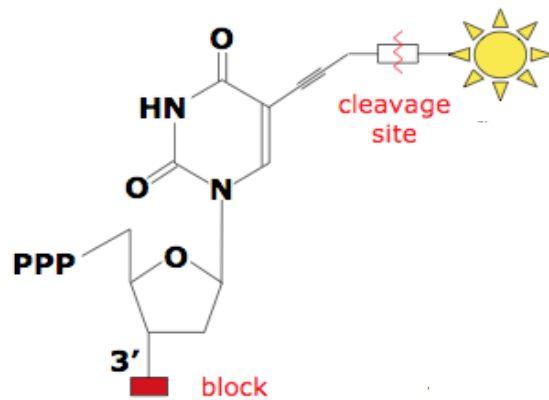
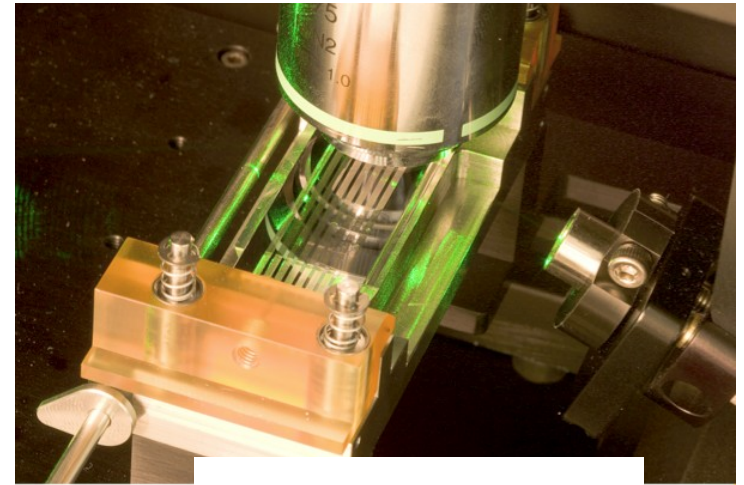


Sequencing Chemistry

Sequencing by Synthesis



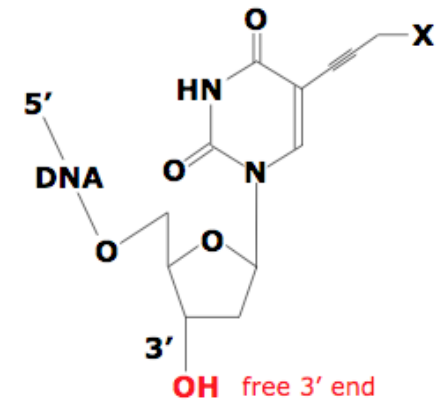
- Cycle 1: Add sequencing reagents
First base incorporated
Remove unincorporated bases
Detect signal
- Cycle 2-n: Add sequencing reagents and repeat



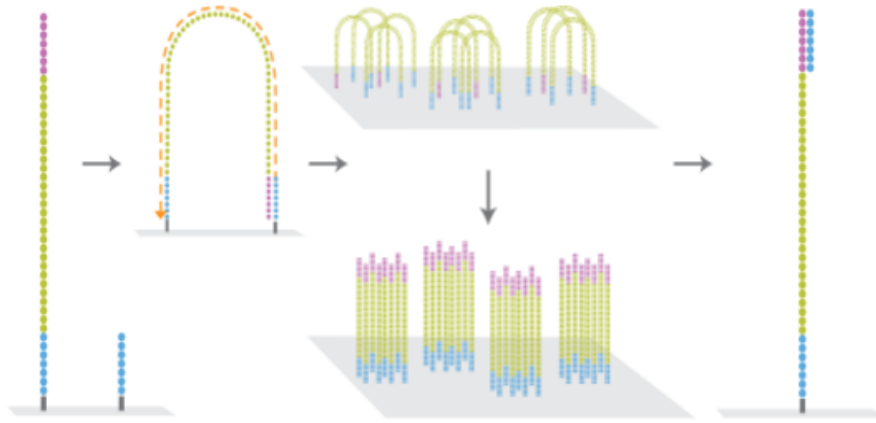
Cleave fluorophore



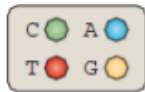
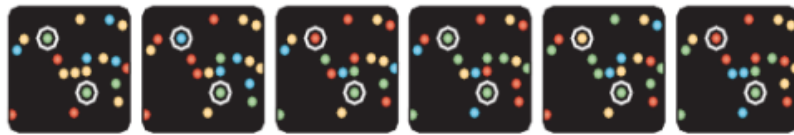
De-block 3' terminus



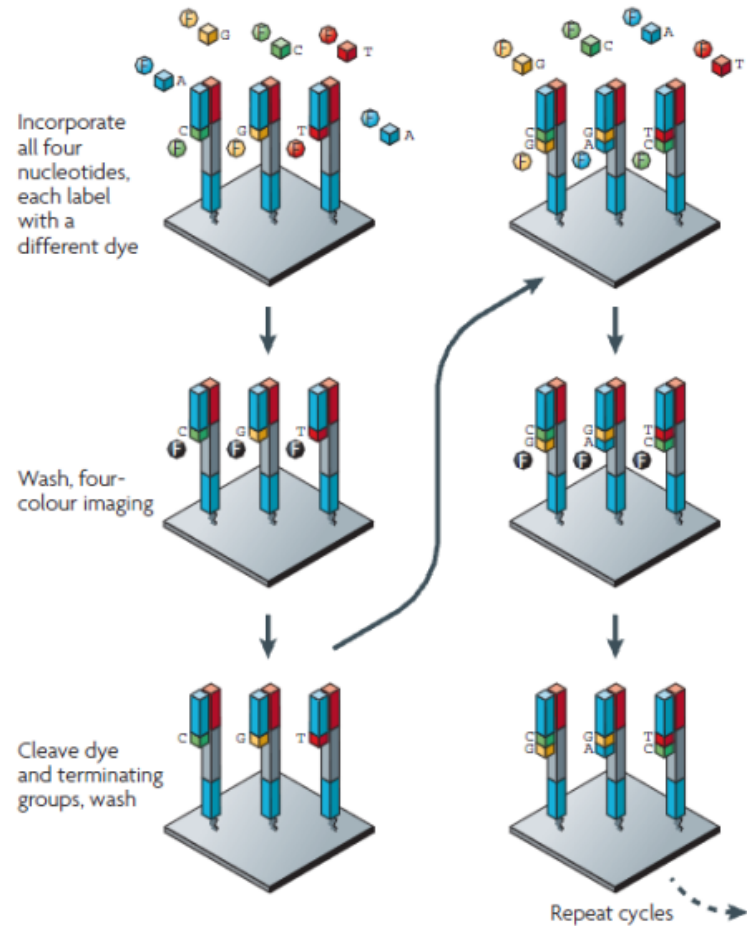
Illumina: Reversible termination sequencing



4 nucleotides with different dye flow simultaneously



Top: CATCGT
Bottom: CCCCCC



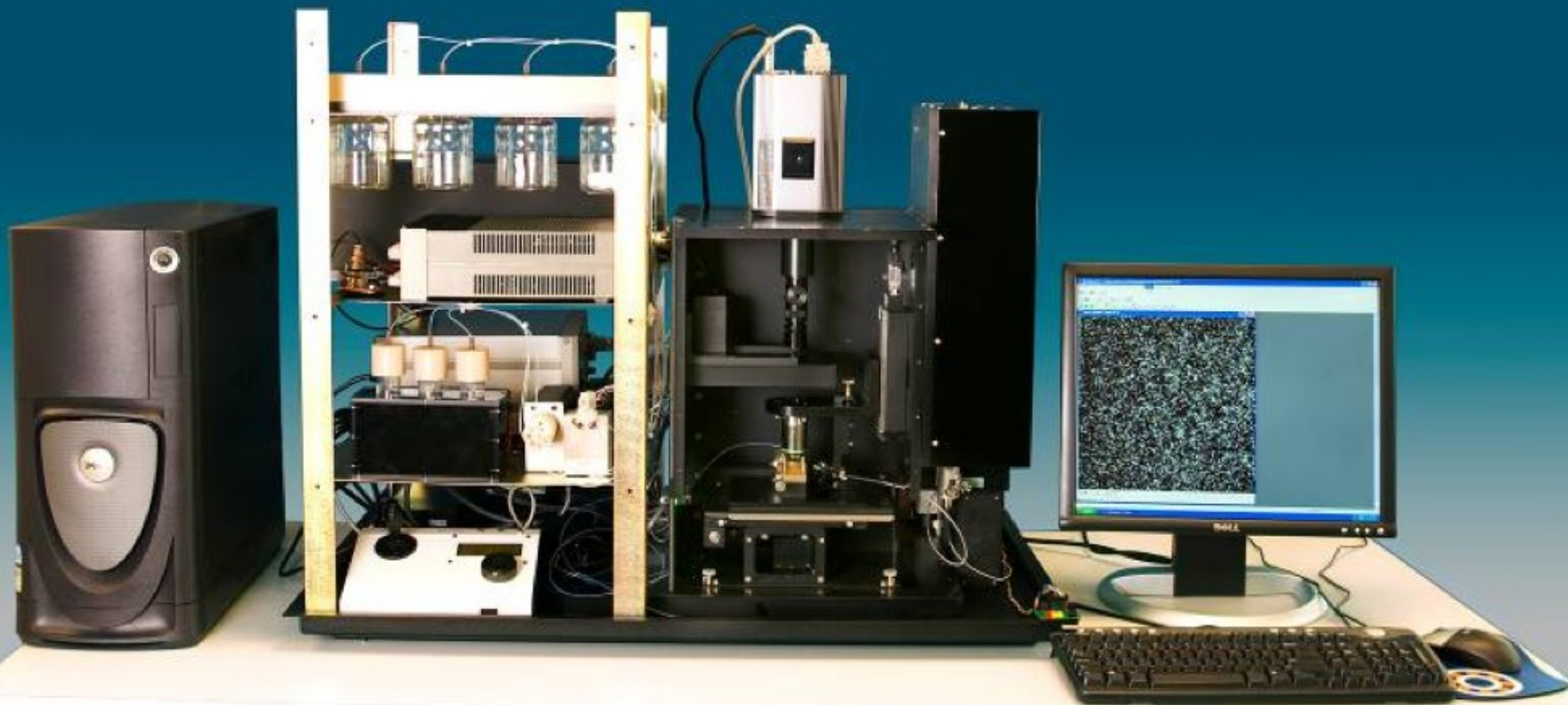
Illumina 2- and 4-channel SBS (sequencing by synthesis) sequencing technology



Fluidics & electronics

Flow cell & detection

Laser optics



ILLUMINA IMPROVEMENT



- Longer read length (250 bp)
- 3-fold more reads (15 M)
- Higher throughput (5-7 Gb)
- Faster run time



Two run configurations

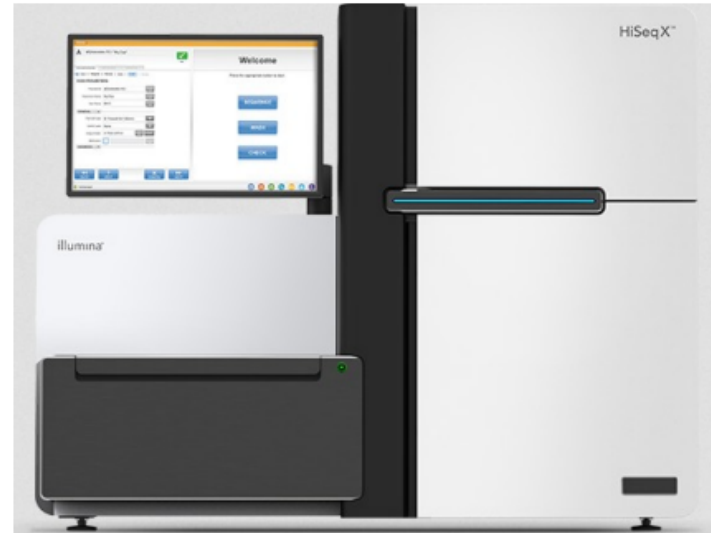
- Fast run config can be done in 27 hours and produce 120 Gb
- Standard run config remains the same (600 Gb in 17 days)



NovaSeq 5000 / 6000



HiSeq X



Next Generation Sequencing: Amplified Single Molecule Sequencing

Illumina

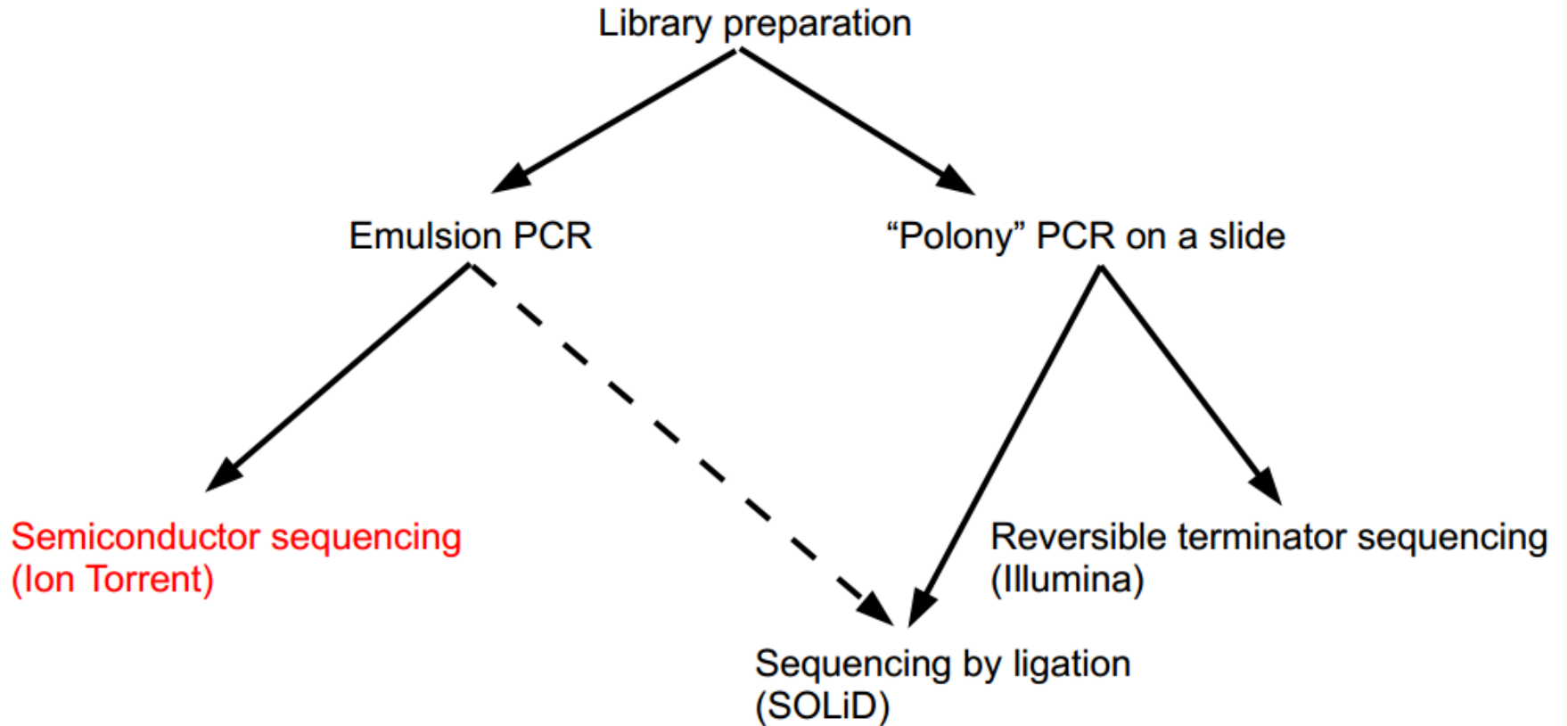
	MiniSeq	MiSeq	NextSeq 500
Read Length	2 x 150 bp	2 x 300 bp	2 x 150 bp
Throughput	7.5 Gb	15 Gb	120 Gb
Reads per run	50 million	50 million	800 million
Accuracy	99,9 % (>80%)	99,9 % (>70%)	99,9 % (>80% of the bases)
Run Time	24 hours	55 hours	29 hours

	HiSeq 2500 / 3000 / 4000	NovaSeq 5000 / 6000	HiSeq X
Read Length	2 x 125 / 2 x 150 / 2 x 150 bp	2 x 150 bp	2 x 150 bp
Throughput	1000 / 750 / 1500 Gb	500 - 1000 / 500 - 3000 Gb	1800 Gb
Reads per run	4 / 2,5 / 5 billion	1.6 - 3.3 / 1.6 - 10 billion	6 billion
Accuracy	99,9 % (>80% of the bases)	99,9 % (>80% of the bases)	99,9 % (>75%)
Run Time	6 / 3,5 / 3,5 days	19 - 40 hours	< 3 days

Workflow: Library preparation → Bridge amplification → Reversible termination sequencing



Next Generation Sequencing: Amplified Single Molecule Sequencing



Jonathan M. Rothberg

- 1999: founded
454 Life Sciences
- 2007: founded
Ion Torrent



Ion Torrent

PGM
(Personal Genome Machine)



S5 / S5 XL



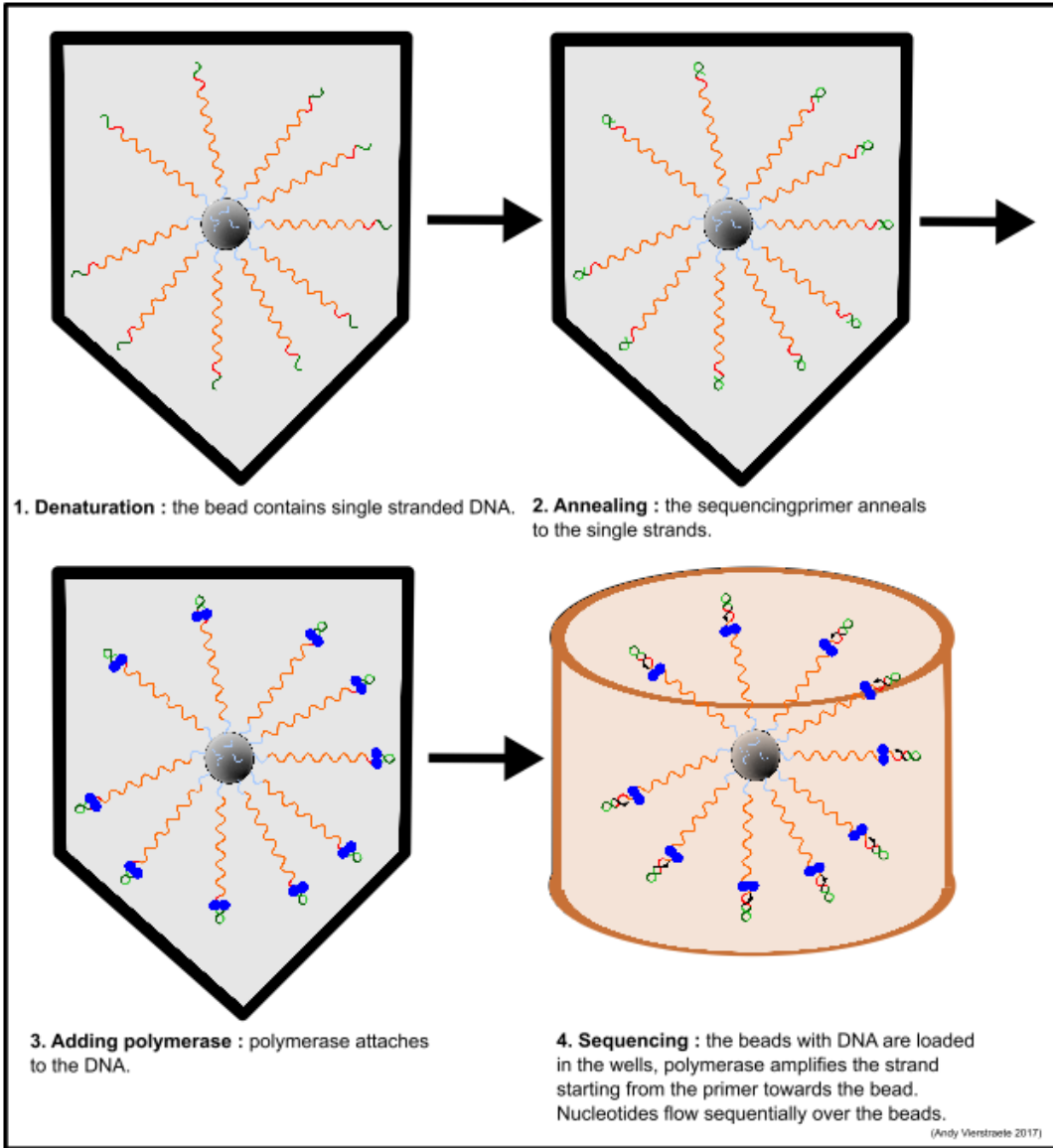
Proton



	PGM	S5 / S5 XL	Proton
Chip	314 – 316 - 318	520 – 530 – 540	PI – PII
Read length	400 bp	400/600 – 400/600 – 200	200 bp - ?
Throughput	0,1 – 0,6 – 2 Gb	2 – 8 - 15 Gb	10 -100 Gb
Reads per run	0,5 – 3 – 6 million	5 – 20 – 80 million	80 - 250 million
Accuracy	99 % (raw read)	99 % (raw read)	99 % (raw read)
Run Time	4 – 5 – 7 hours	8 – 17 - 17 hours (4 times faster for XL)	4 hours



Ion Torrent



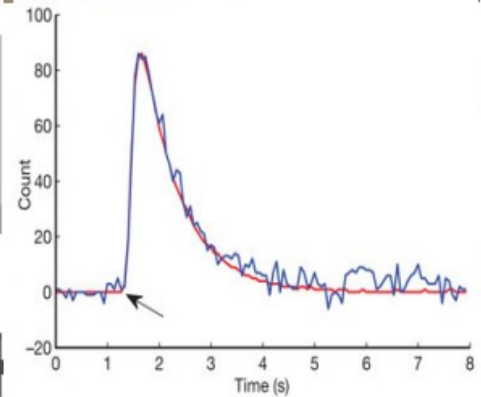
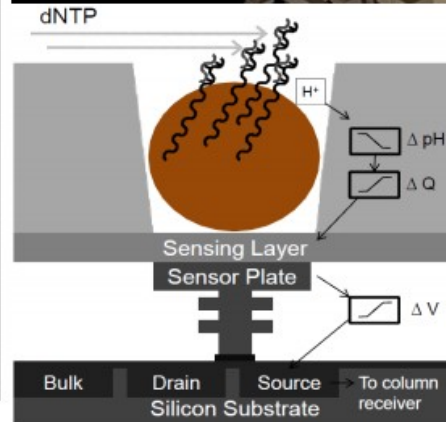
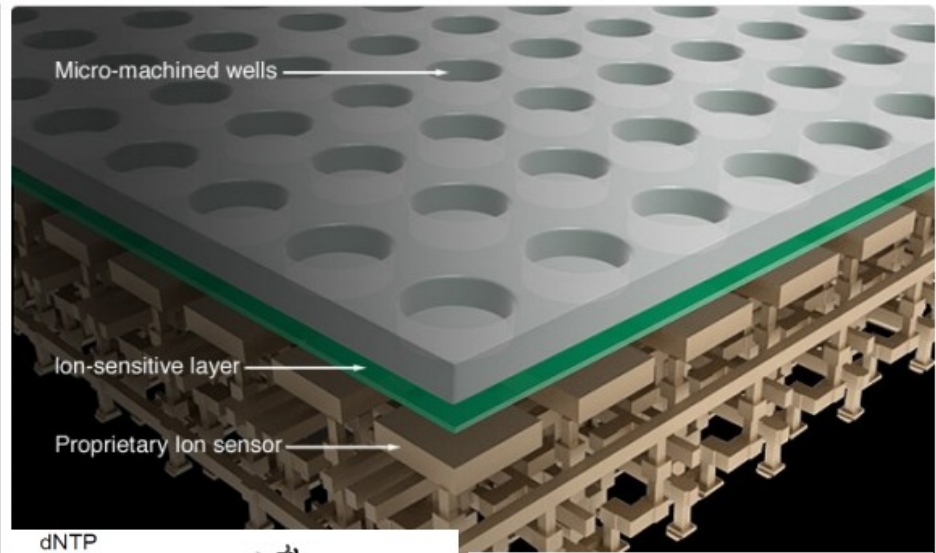
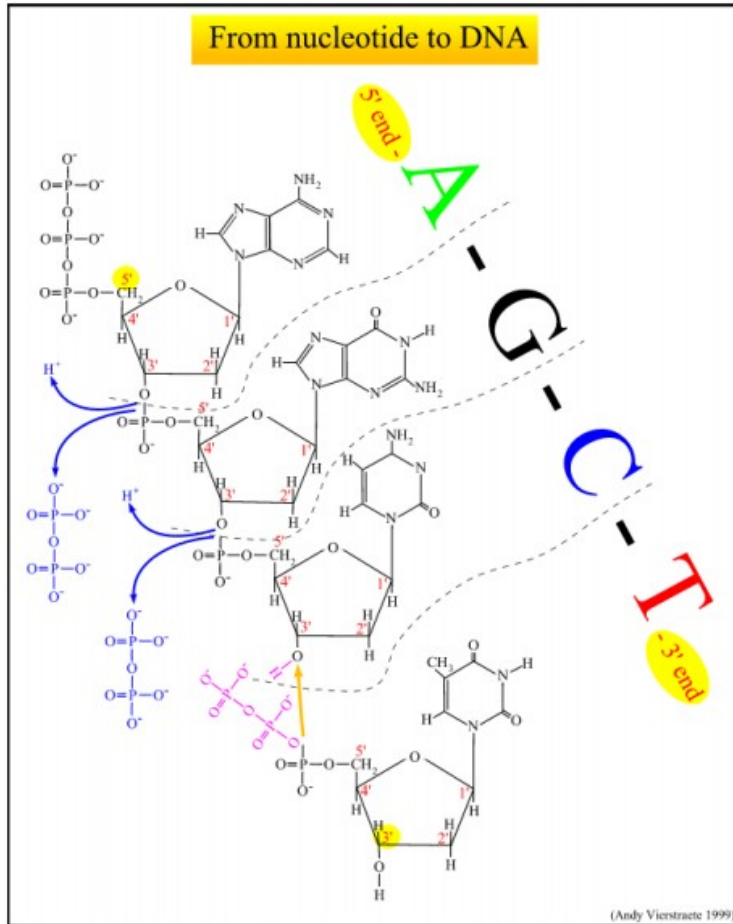
Sequencing



Next Generation Sequencing: Amplified Single Molecule Sequencing

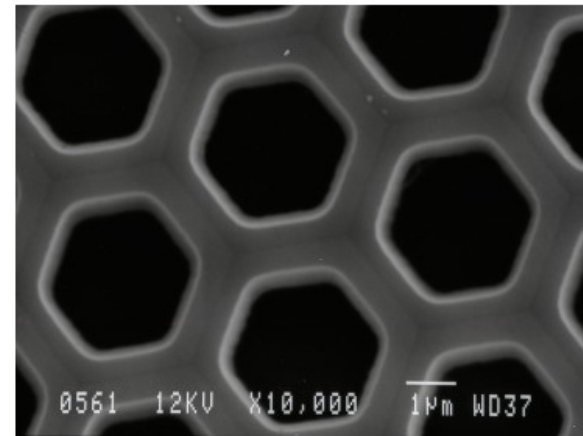
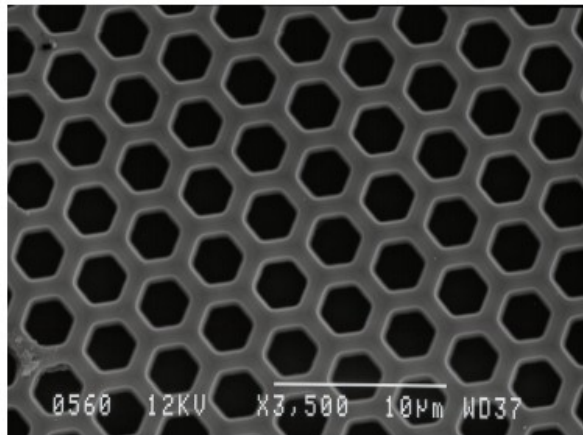
Ion Torrent

Workflow: Library preparation → Emulsion PCR → Semiconductor Sequencing

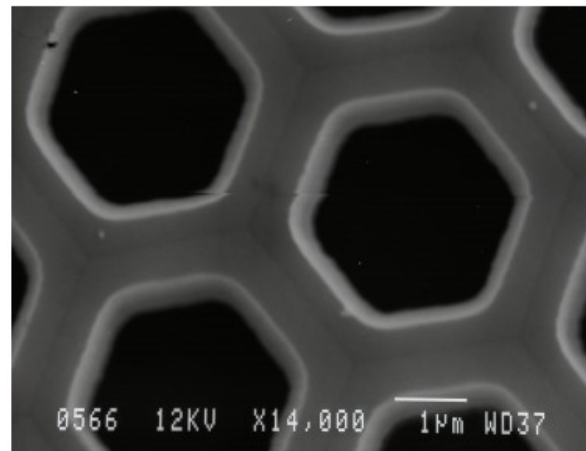


Next Generation Sequencing: Amplified Single Molecule Sequencing Ion Torrent

Workflow: Library preparation → Emulsion PCR → Semiconductor Sequencing

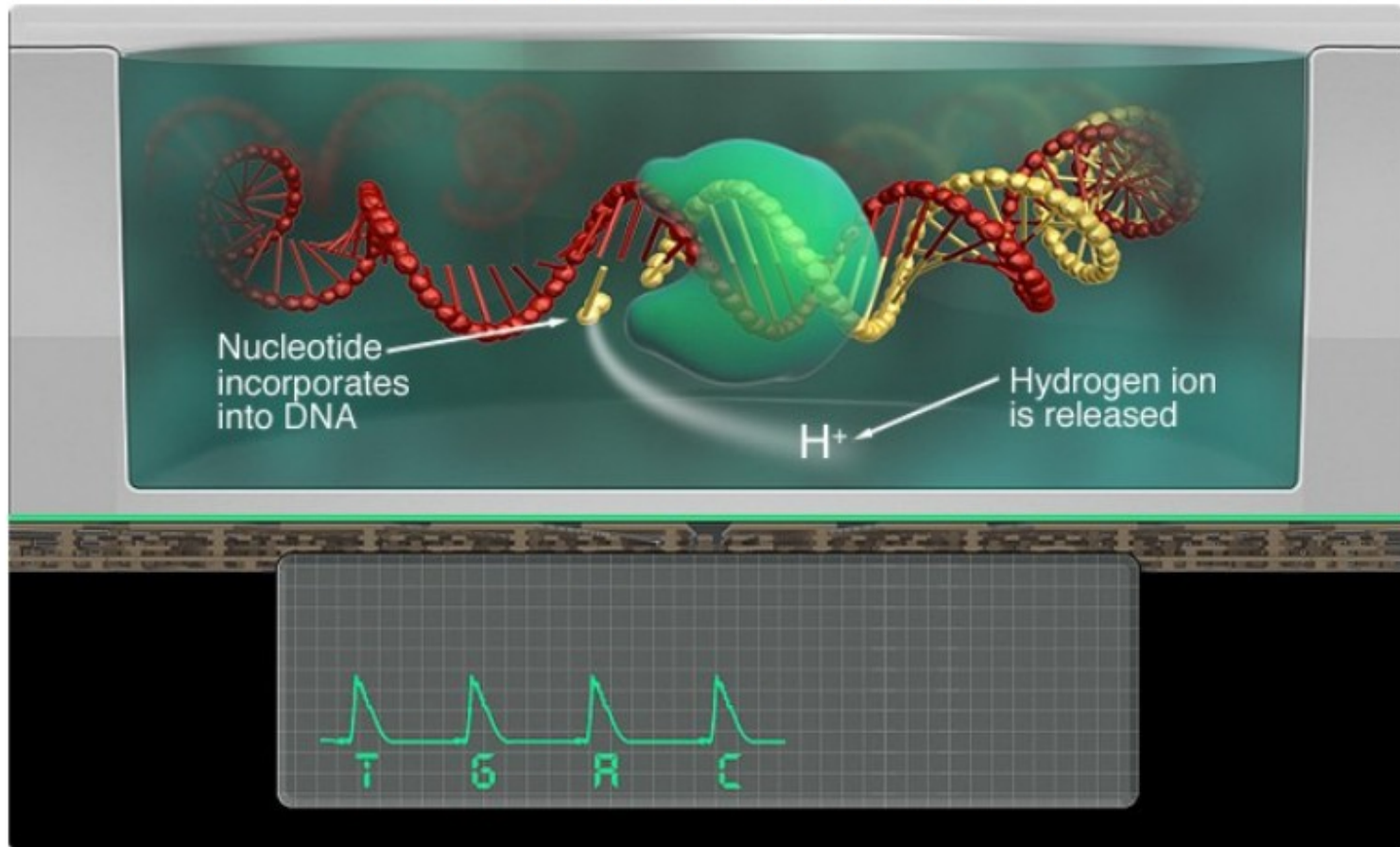


Picture of the wells in a
318 chip



Next Generation Sequencing: Amplified Single Molecule Sequencing Ion Torrent

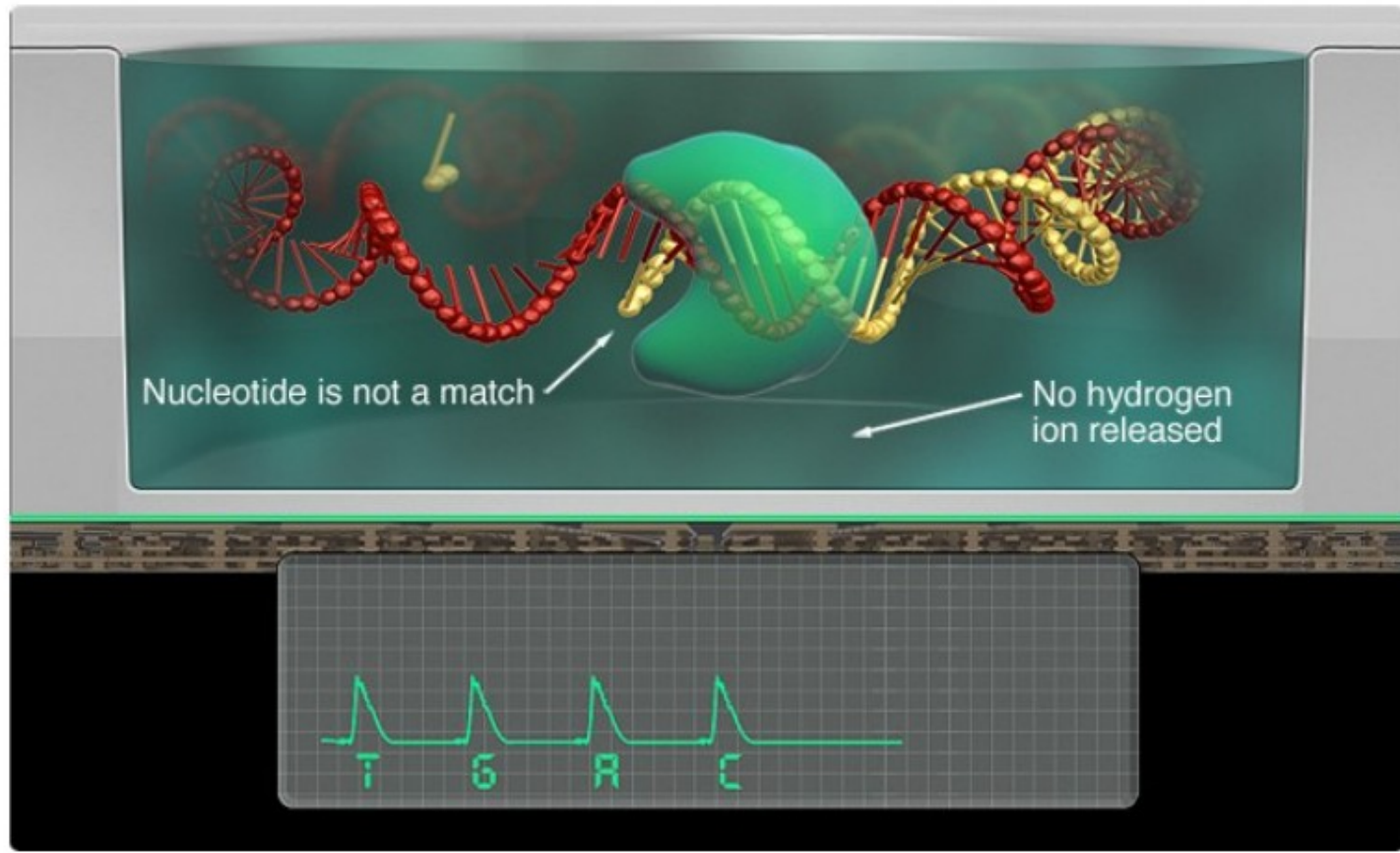
4 nucleotides flow sequentially



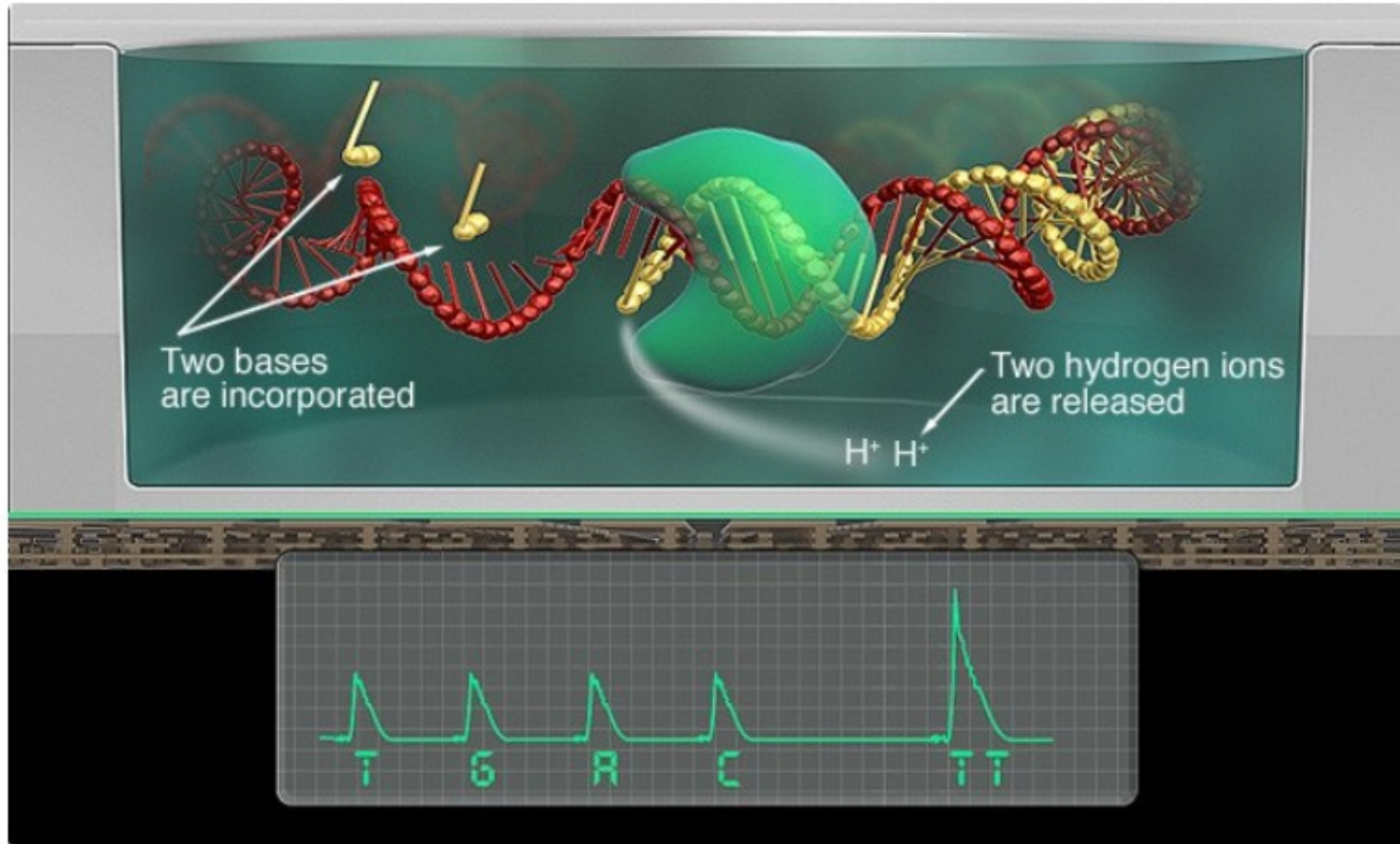
No camera, just a pH sensor



Next Generation Sequencing: Amplified Single Molecule Sequencing Ion Torrent

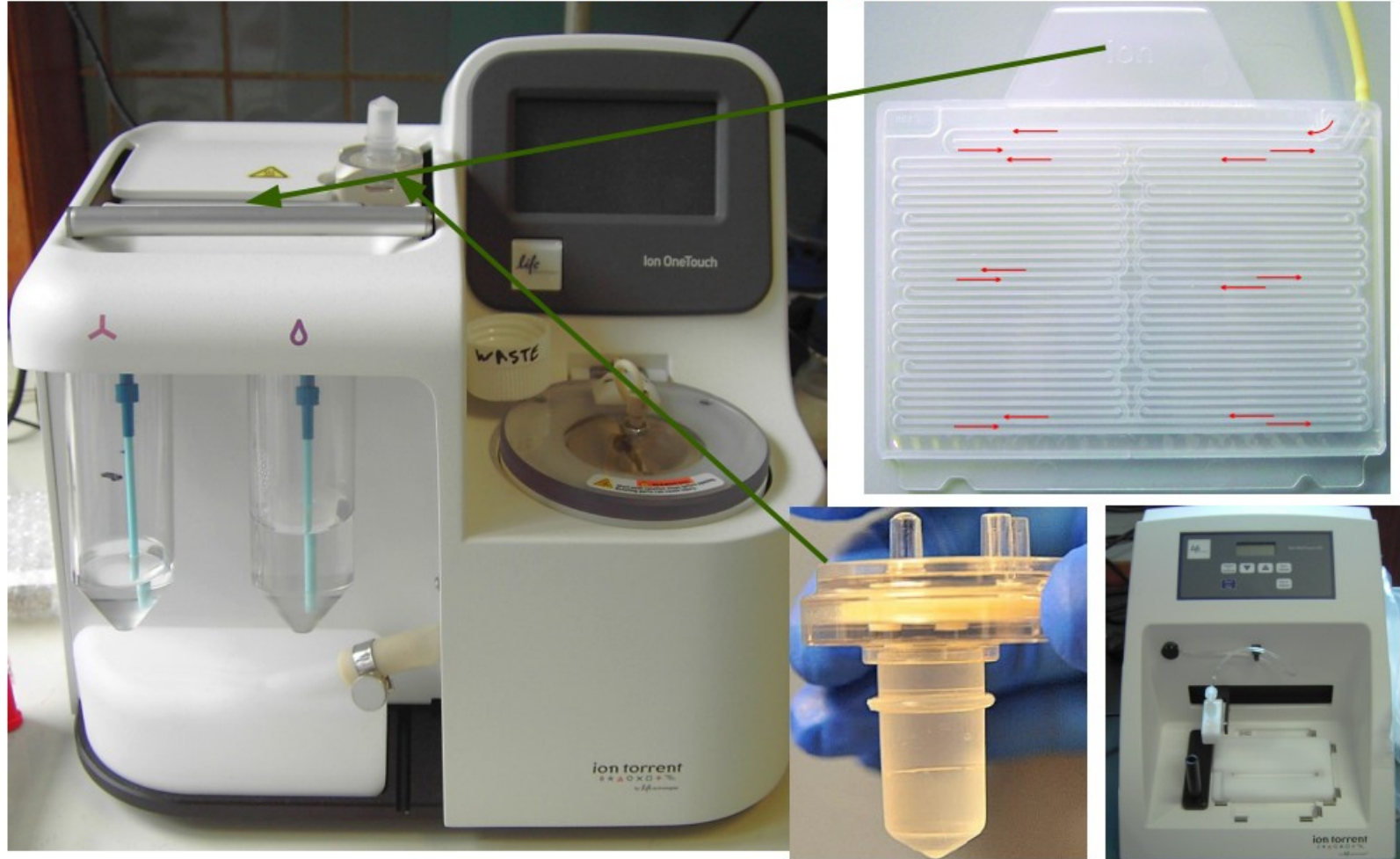


Next Generation Sequencing: Amplified Single Molecule Sequencing Ion Torrent



Next Generation Sequencing: Amplified Single Molecule Sequencing Ion Torrent

Emulsion PCR: OneTouch Instrument

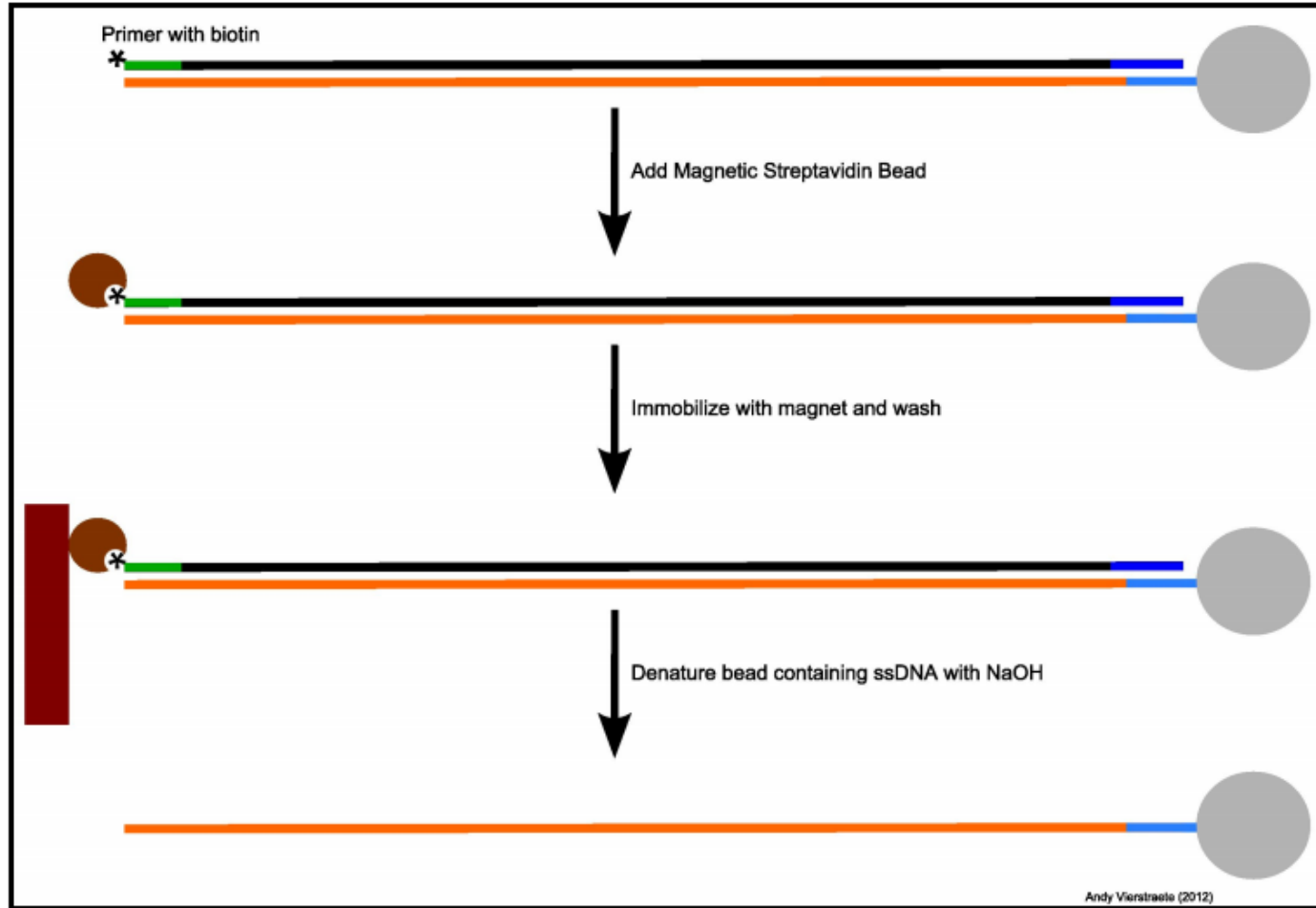


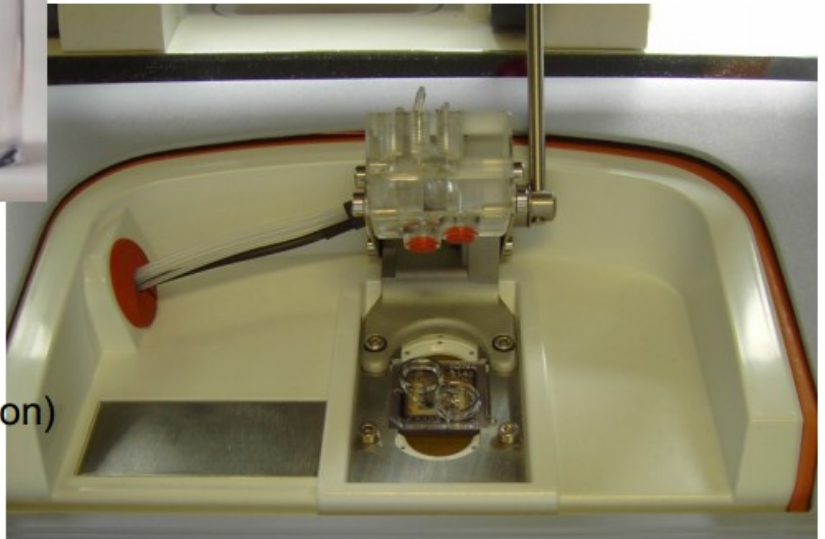
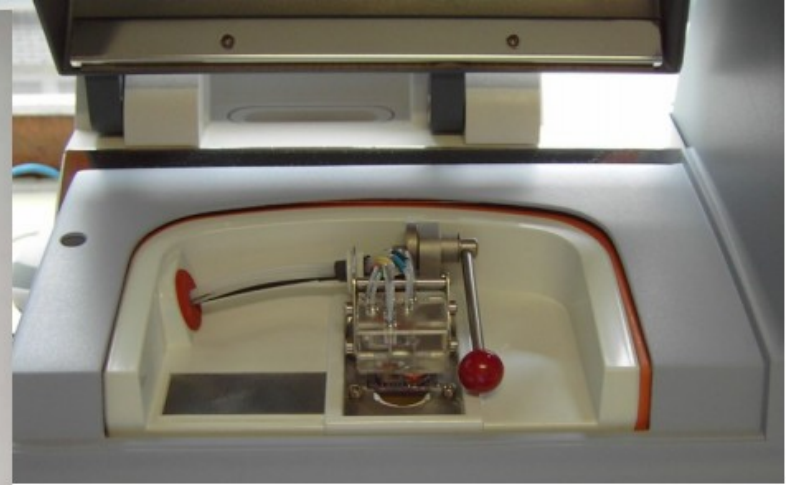
15 min hands-on; 4-8 hours amplification; 35 min enrichment



Next Generation Sequencing: Amplified Single Molecule Sequencing Ion Torrent

Enrichment: select only the beads that contain DNA
-> maximizing sequencing yield





W1 bottle: 350 μ l 100 mM NaOH
W2 bottle: 2 liter W2 solution (contains pcr solution)
W3 bottle: 50 ml W3 (= buffer with known pH)
4 tubes with 20 μ l of the different dNTP's



Pacific Biosciences: il futuro, il sequenziamento massivo di singole molecole

Published online before print January 23, 2008, 10.1073/pnas.0710982105

PNAS | January 29, 2008 | vol. 105 | no. 4 | 1176-1181

[◀ Previous Article](#) | [Table of Contents](#) | [Next Article ▶](#)

BIOLOGICAL SCIENCES / BIOPHYSICS

Selective aluminum passivation for targeted immobilization of single DNA polymerase molecules in zero-mode waveguide nanostructures

Jonas Korlach, Patrick J. Marks, Ronald L. Cicero, Jeremy J. Gray, Devon L. Murphy, Daniel B. Roitman, Thang T. Pham, Geoff A. Otto, Mathieu Foquet, and Stephen W. Turner*

Pacific Biosciences, 1505 Adams Drive, Menlo Park, CA 94025

Communicated by Watt W. Webb, Cornell University, Ithaca, NY, November 20, 2007 (received for review August 6, 2007)

Optical nanostructures have enabled the creation of subdiffraction detection volumes for single-molecule fluorescence microscopy. Their applicability is extended by the ability to place molecules in the confined observation volume without interfering with their biological function. Here, we demonstrate that processive DNA synthesis thousands of bases in length was carried out by individual DNA polymerase molecules immobilized in the observation volumes of zero-mode waveguides (ZMWs) in high-density arrays.



Third Generation Sequencing: Single Molecule Sequencing

Pacific Biosciences

Pacbio RS

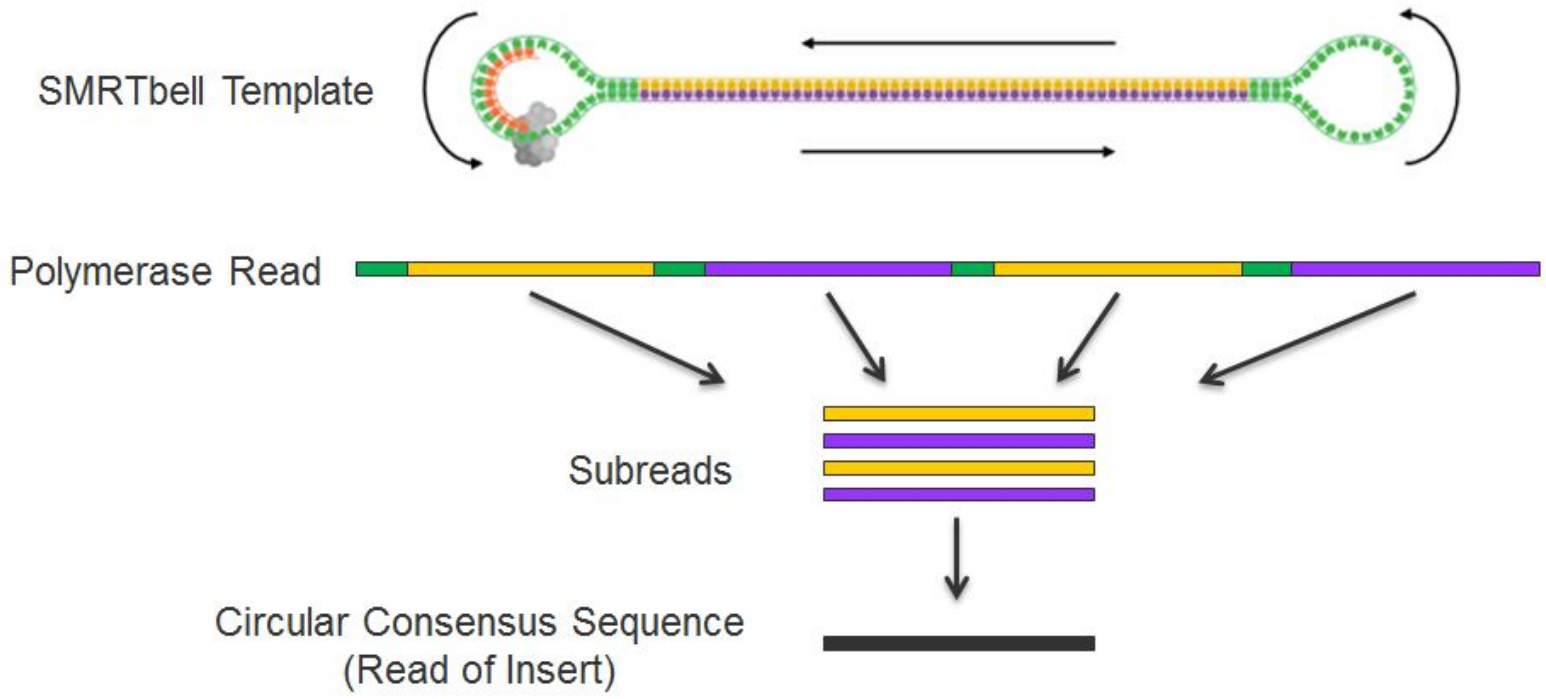


Sequel System



	Pacbio RS	Sequel System
Read Length	50 % > 20 kb (max > 60 kb)	50 % > 20 kb (max > 60 kb)
Throughput	1 Gb/SMRT cell (max 16/run)	5-8 Gb/SMRT cell (max 16/run)
Reads per run	55,000	365,000
Accuracy	86 %	86 %
Run Time	30 minutes – 6 hours/ SMRT cell	30 minutes – 10 hours/ SMRT cell

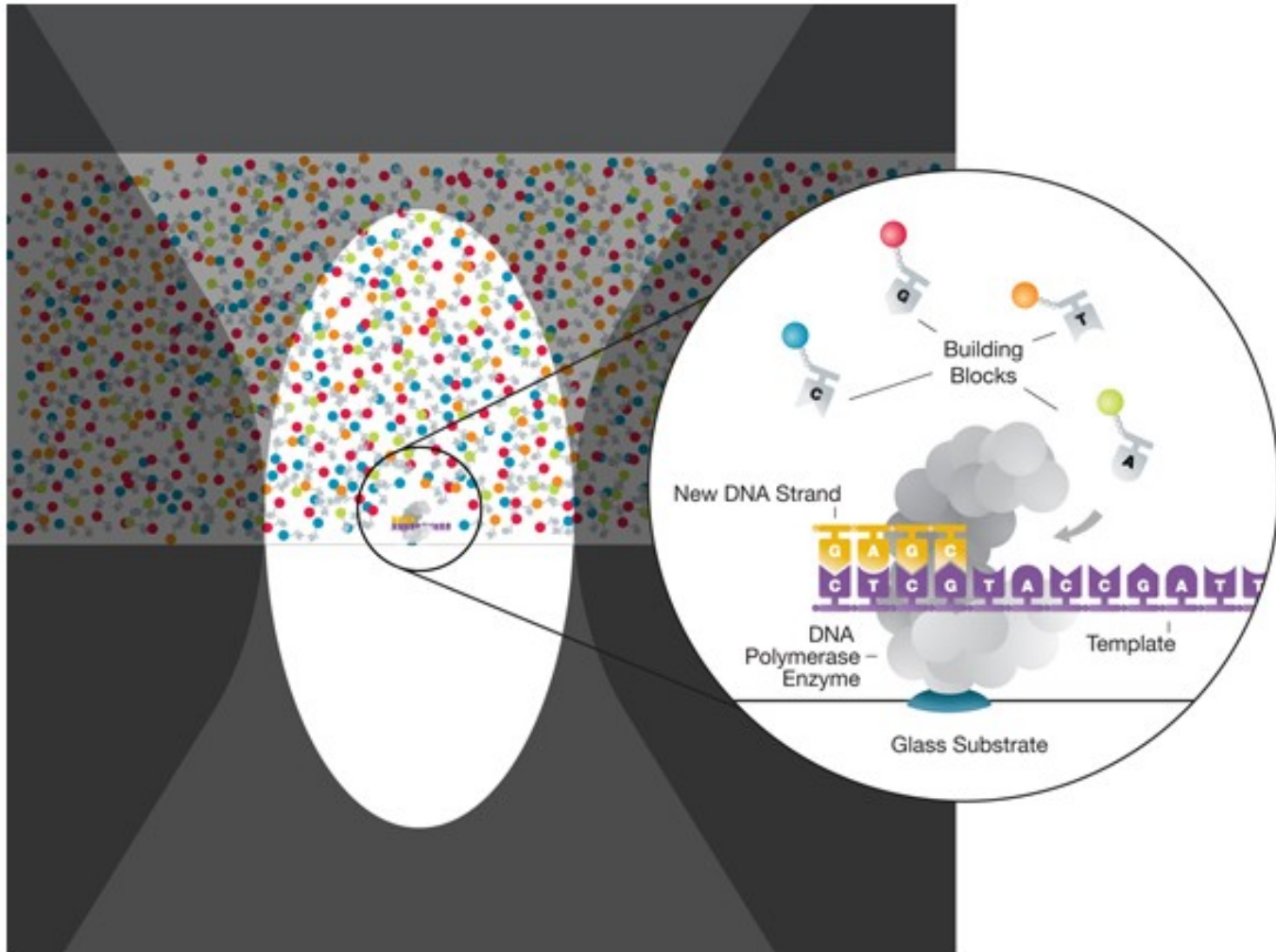




SMRT SEQUENCING

- **SMRT=Single Molecule Real-Time**
- **una ZMW è un foro, di diametro di qualche decina di nanometri, fabbricato in un film di metallo di 100nm depositato su un substrato di diossido di silicene**
- **Ciascun ZMW diventa una camera di visualizzazione nanofotonica che fornisce un volume di rilevazione di solo 20 zeptolitri (10^{-21} litri).**
- **In un tale volume, l'attività di una singola molecola può essere rilevata in un background di migliaia di nucleotidi marcati in fluorescenza**

Pacific Biosciences ZMW



ZMW='Zero Mode Waveguide'



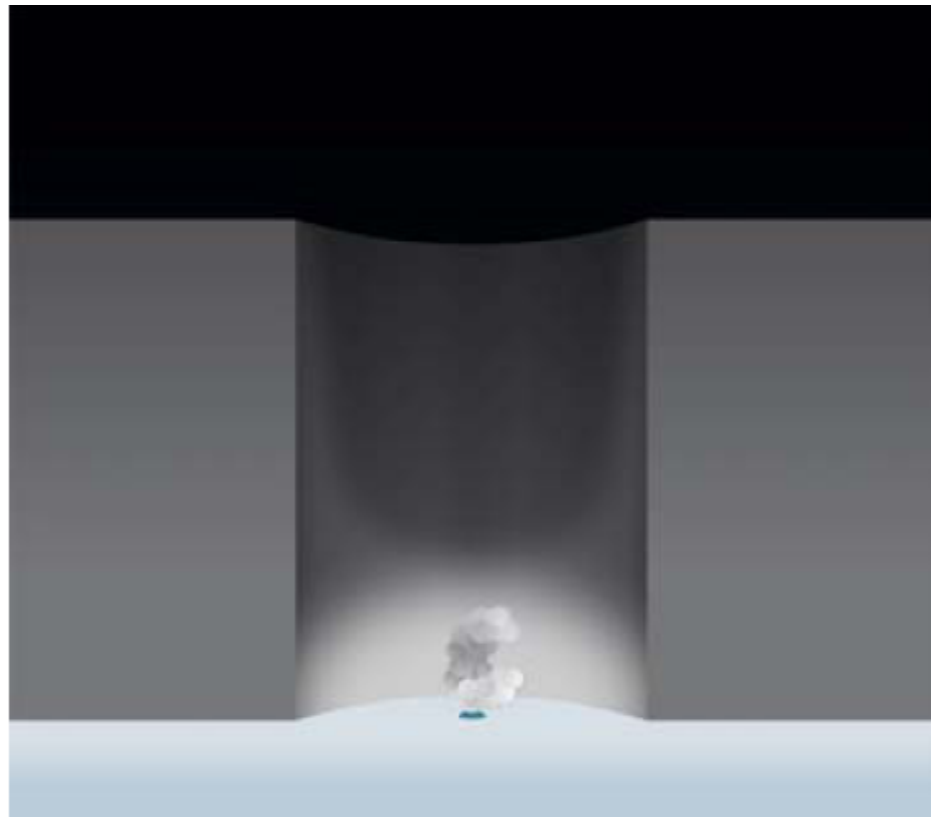
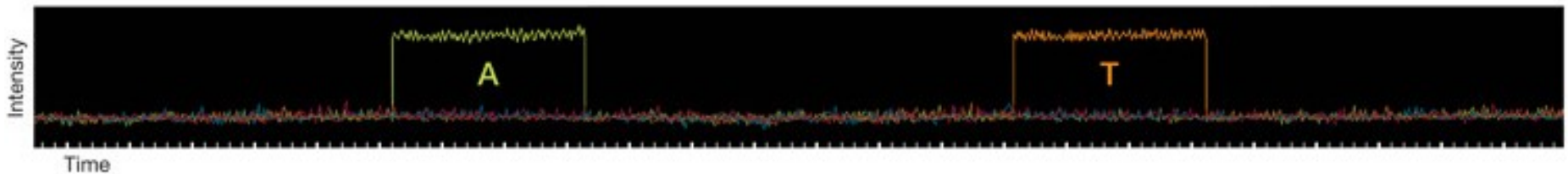


Figure 5. *ZMW with DNA polymerase*

A single DNA polymerase molecule is attached to the bottom of the ZMW using a proprietary biased immobilization process.



Pacific Biosciences



- ✓ Sequenziamento di singole molecole
- ✓ Incorporazione di molecole fluorescenti (sequenziamento mediante sintesi)
- ✓ Monitoraggio in tempo reale dell'attività della polimerasi – (eccitazione/rilevazione)



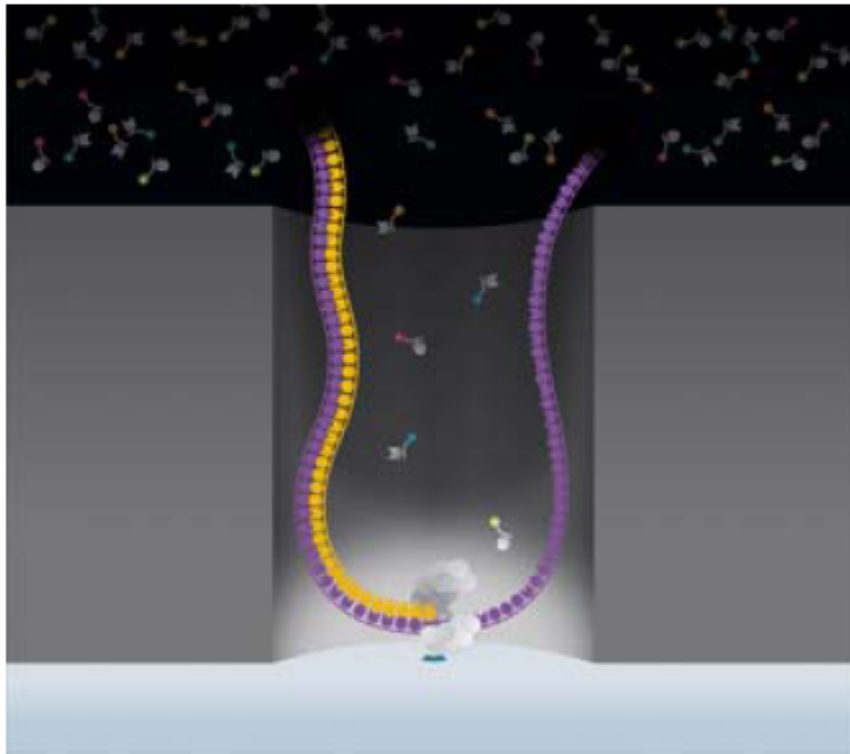
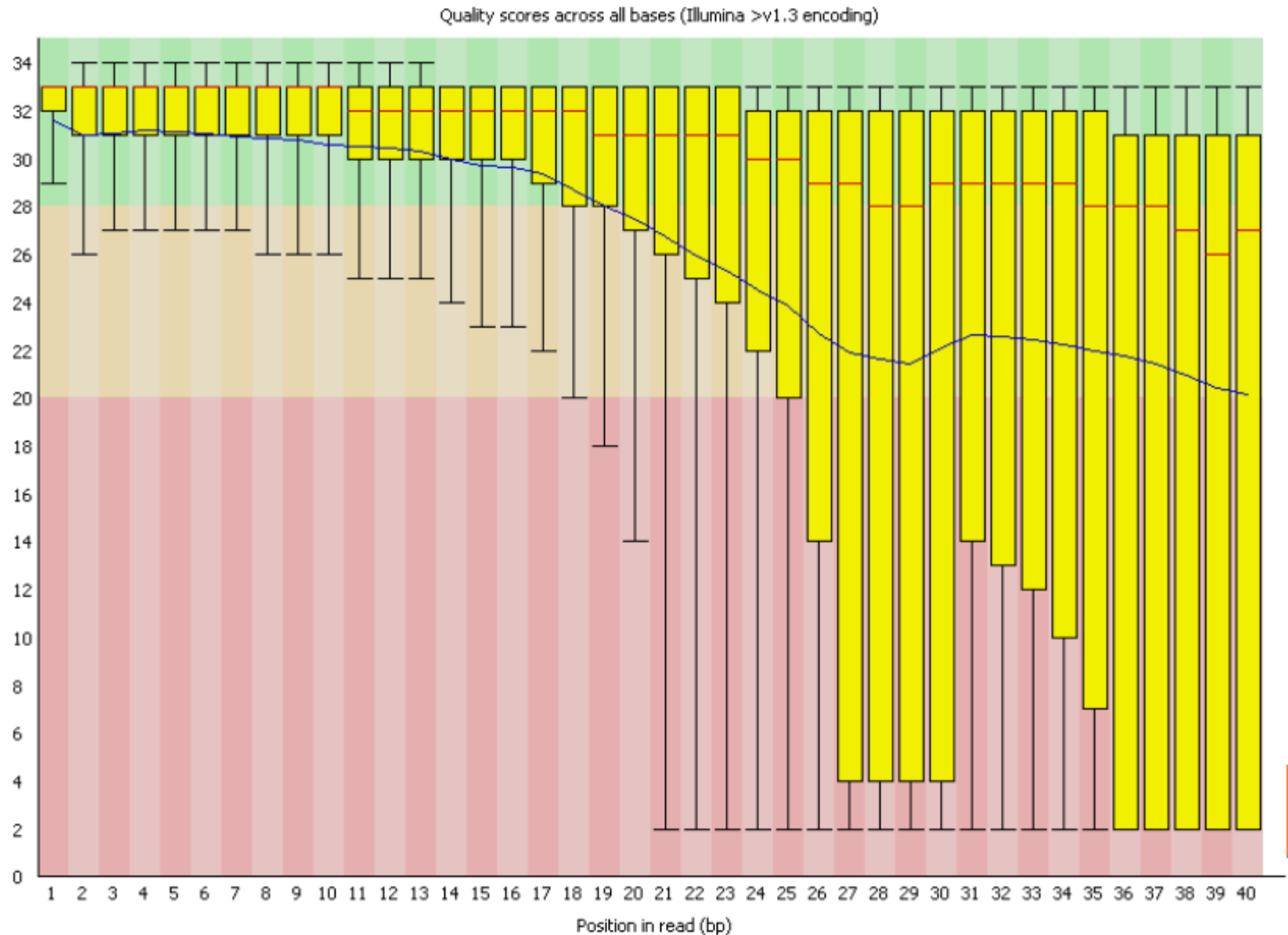


Figure 11. Synthesis of long DNA.

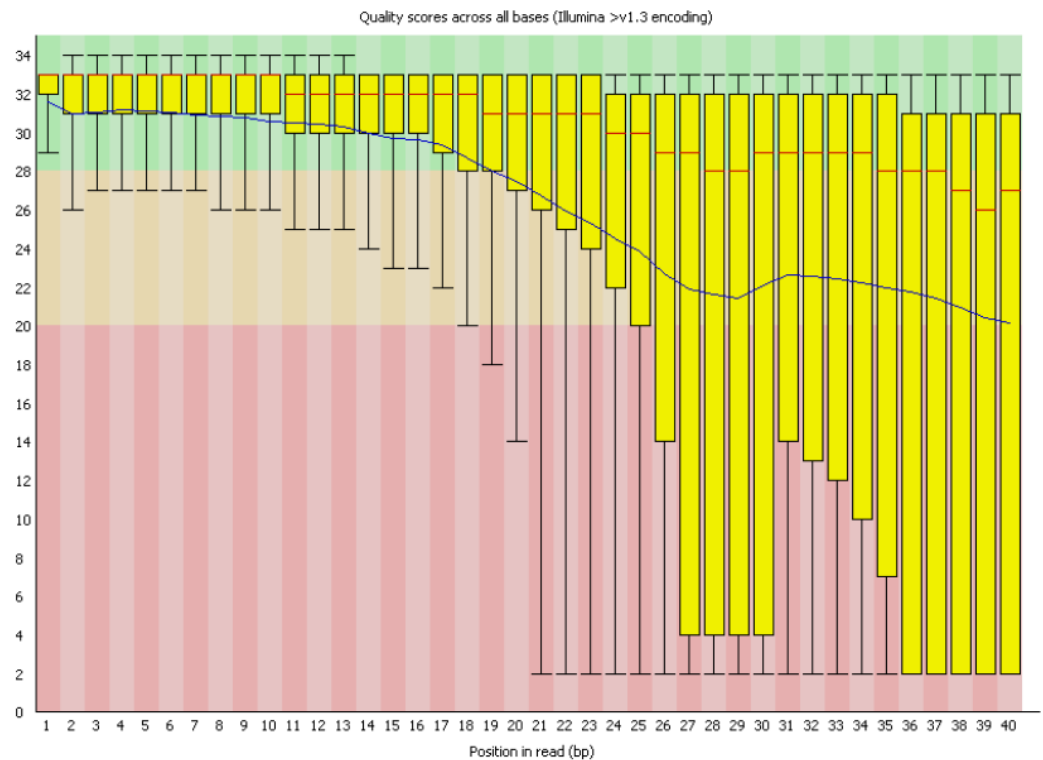
DNA polymerase processively incorporates nucleotides producing long, natural DNA.



FastQC



FastQC



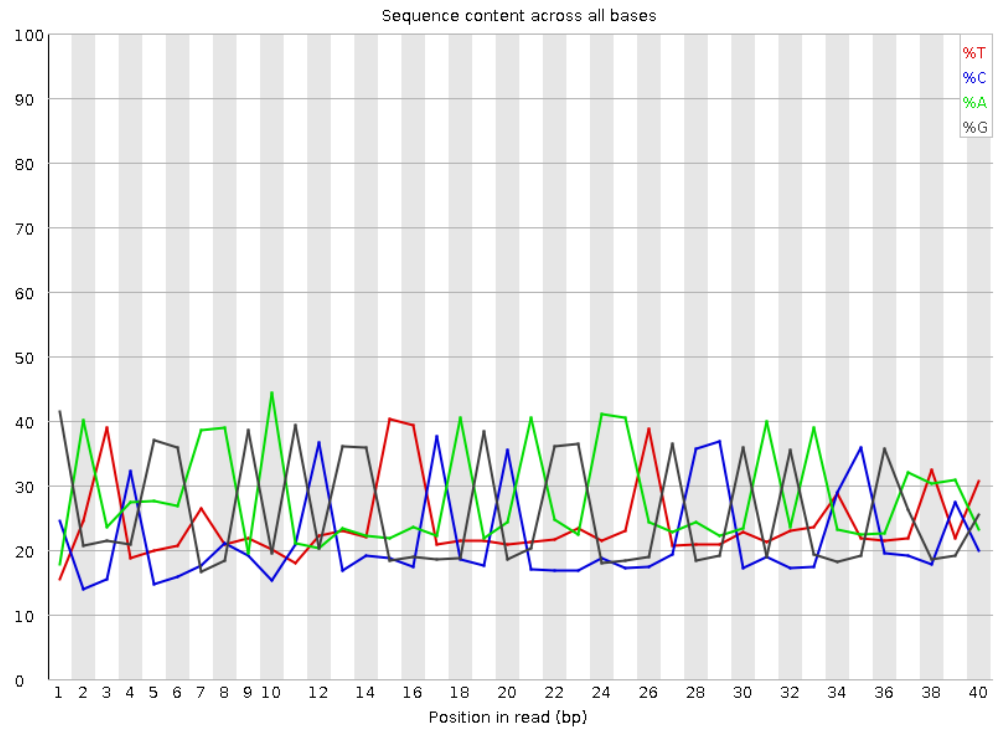
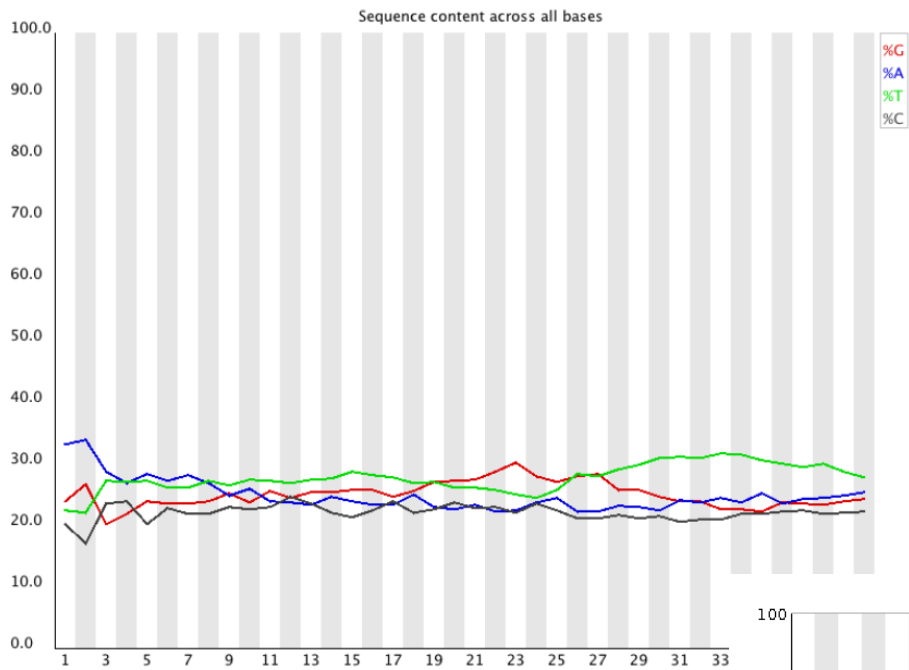
The central red line is the median value

The yellow box represents the inter-quartile range (25-75%)

The upper and lower whiskers represent the 10% and 90% points

The blue line represents the mean quality





Third Generation Sequencing: Single Molecule Sequencing

Oxford Nanopore

SmidgION



MinION



GridION X5



PromethION

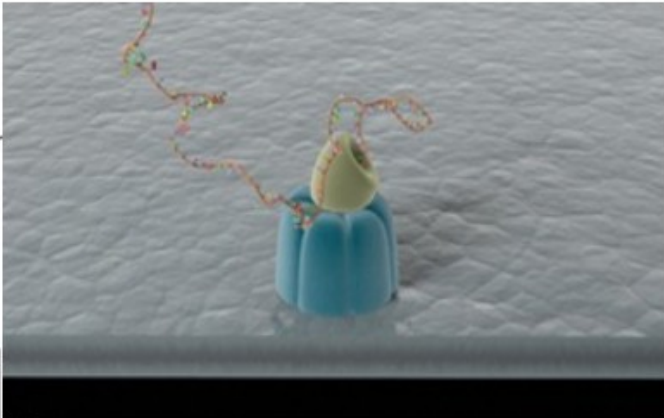
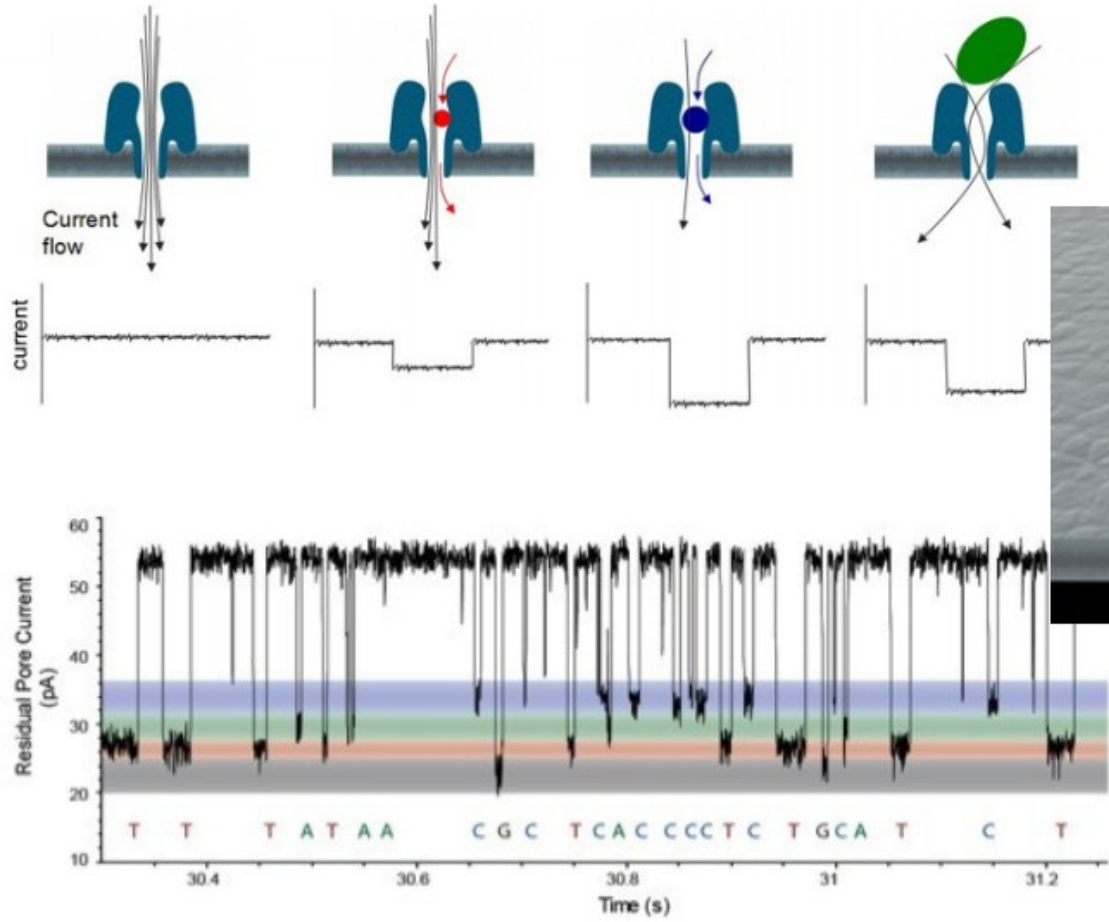


	SmidgION	MinION	GridION X5	PromethION
Read Length	?	> 200 kb		
Throughput	1 Gb (1 flow cell with 256 pores) ?	10-20 Gb (1 flow cell with 512 pores)	100 Gb (5 flow cells with 512 pores/cell) 2560 pores	50 – 250 Gb per flow cell/48hours ? (48 flow cells, 3000 pores/cell) 144,000 pores
Reads per run	?	10,000 – >300,000		
Accuracy		90 % (1D) – 96 % (1D ²)		
Run Time	1 – 4 hours	1 - 48 (70) hours	1 – 48 hours	1 - 48 hours

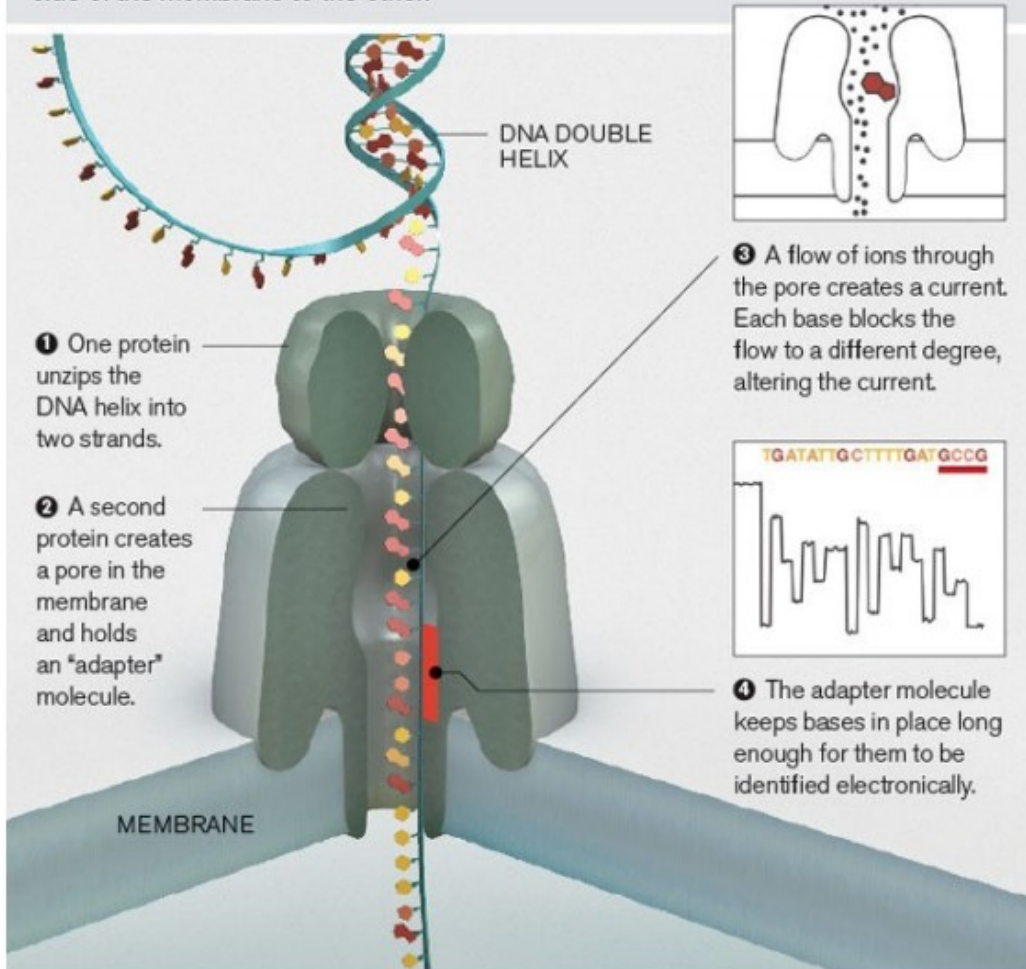
consensus accuracy improved to 99.5% at 30× coverage



Oxford Nanopore



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.



R9.4: 450 nucleotides/second



VoITRAX

Rapid, programmable, portable, disposable sample processor

