

Summary & Review

: Human Genome Project

Sung Wook Chi
Division of Life Sciences, Korea University

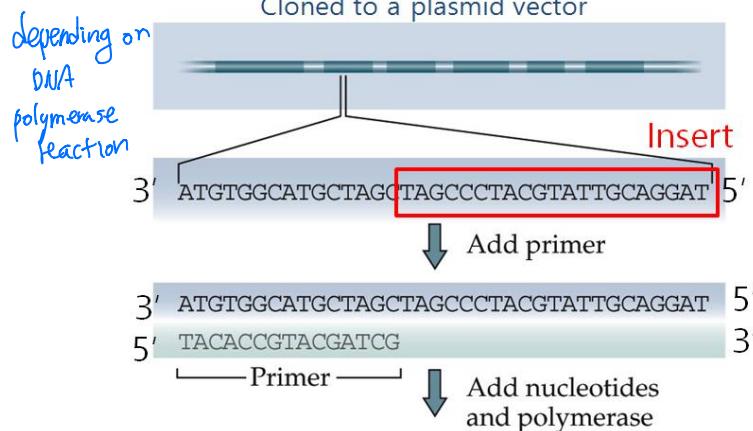
What we learned in the previous class

Sanger Sequencing / PHRED score
quality score

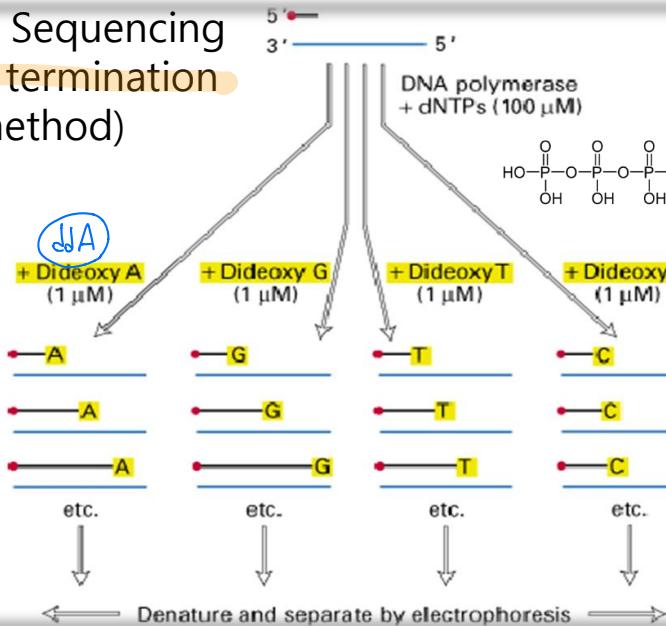
Human Genome Project

- Hierarchical vs. whole-genome shotgun sequencing
- Human genome

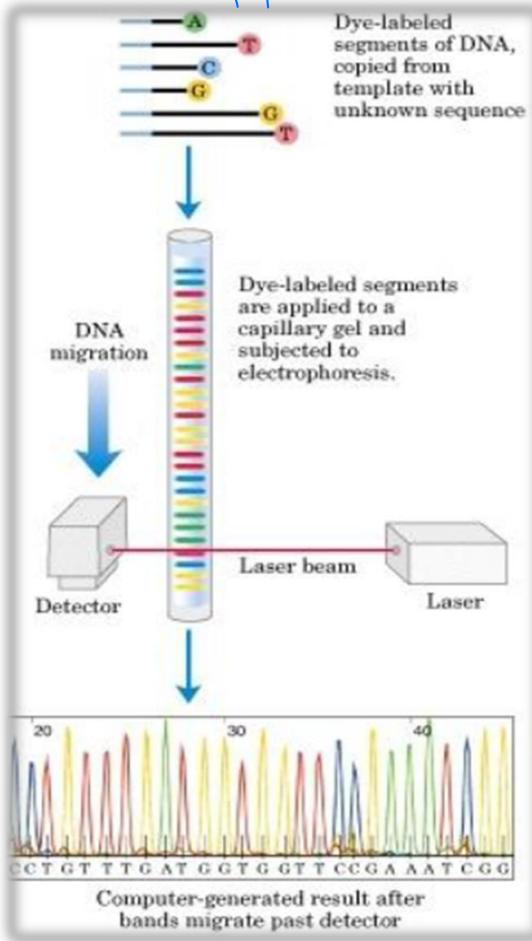
Sanger Sequencing : Chain termination method



Sanger Sequencing
(Chain termination method)



Dideoxy nucleotide
→ Stop point



Phred Scores: quality of base calling

Quality score
(PHRED)
: base calling

$$q = -10 \log_{10}(p)$$

decrease error rate
→ Increase of score

Error rate
in base calling process
(estimated based on
shape/height of peaks)

BOX
3.7

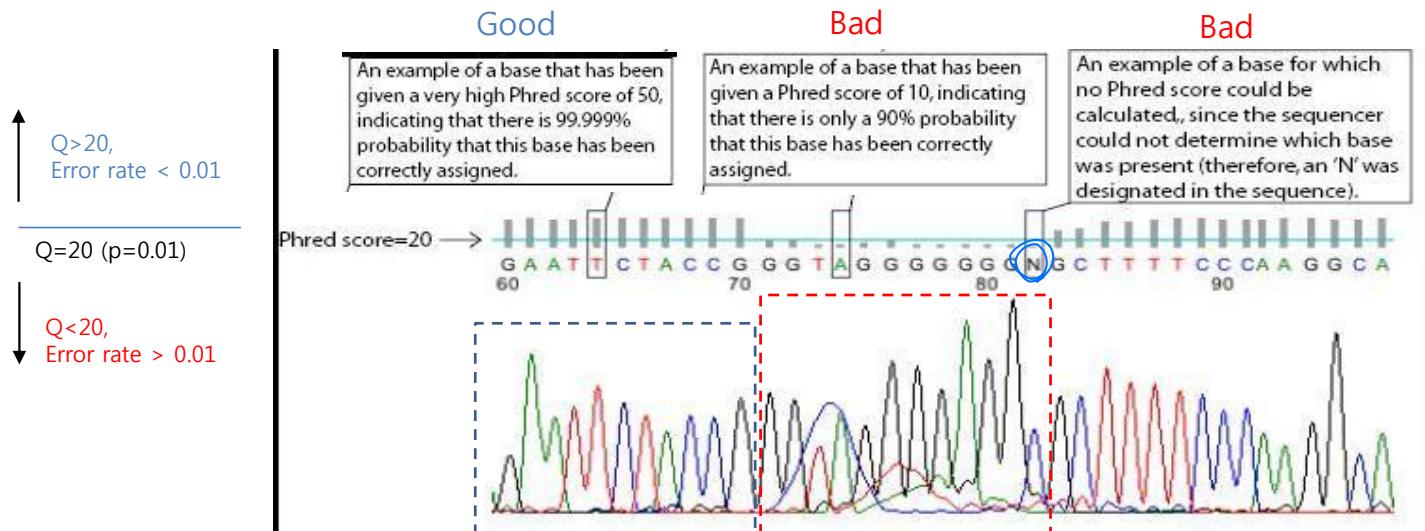
Phred scores: a measure of quality
of sequence determination

The phred score of a sequence determination is a measure of sequence quality. It specifies the probability that the base reported is correct.

If p = the probability that a base is in error, then
the corresponding phred score $q = -10 \log_{10}(p)$.

Here is a short table:

Quality score q	Probability of error	Error rate
10	0.1	1 base in 10 wrong
20	0.01	1 base in 100 wrong
30	0.001	1 base in 1000 wrong
40	0.0001	1 base in 10 000 wrong



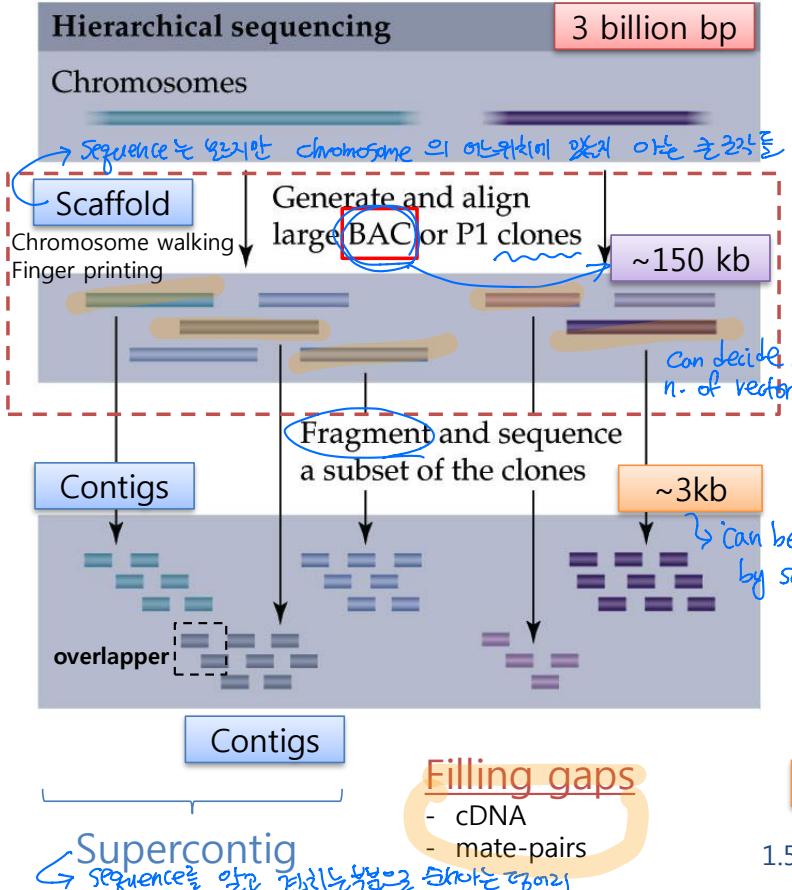
Hierarchical vs. shotgun sequencing

Human genome → **Fragmentation** → **Sequencing** → **Assembly (De novo)**

Hierarchical shotgun sequencing

'BAC-to-BAC' method

(International consortium)



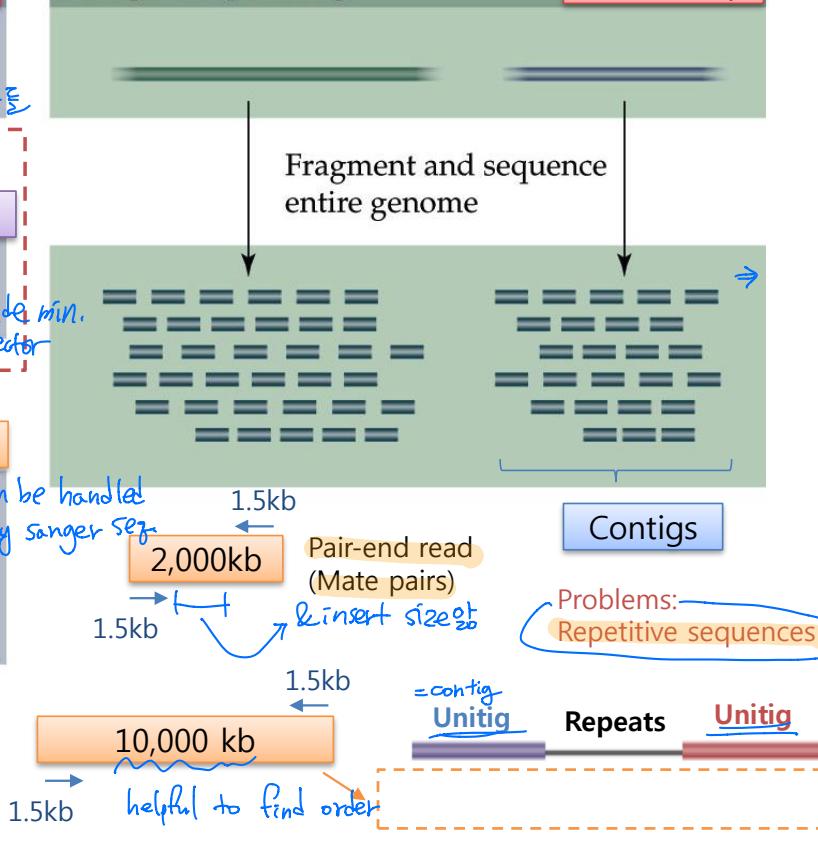
Whole-genome shotgun sequencing

nively fragment everything

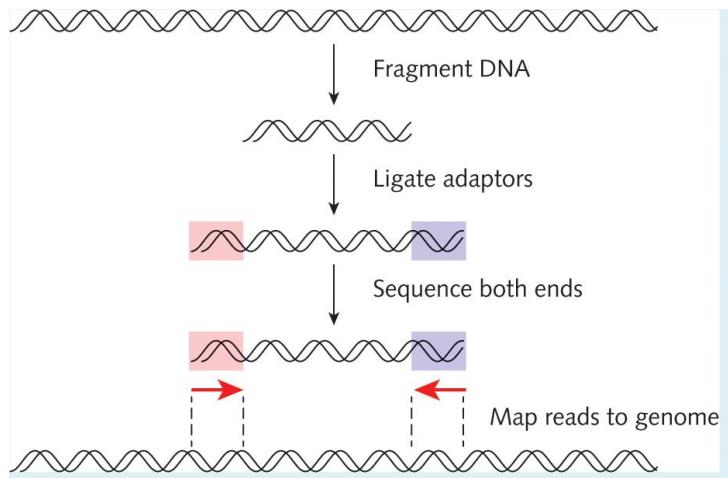
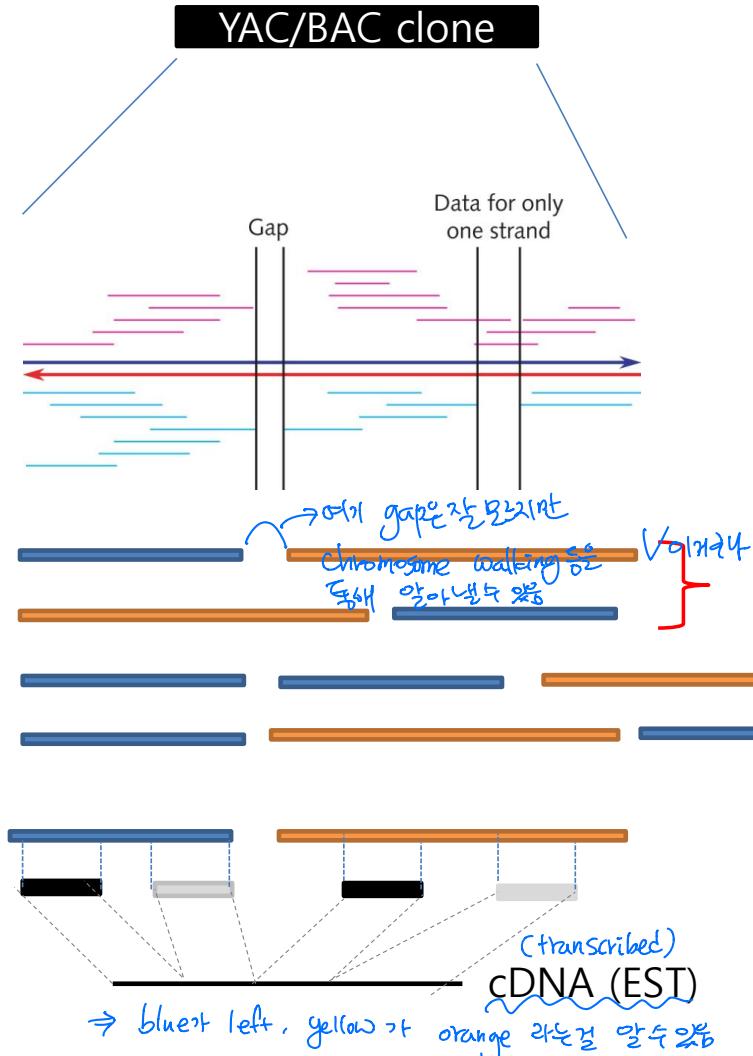
(Cerela)

Shotgun sequencing

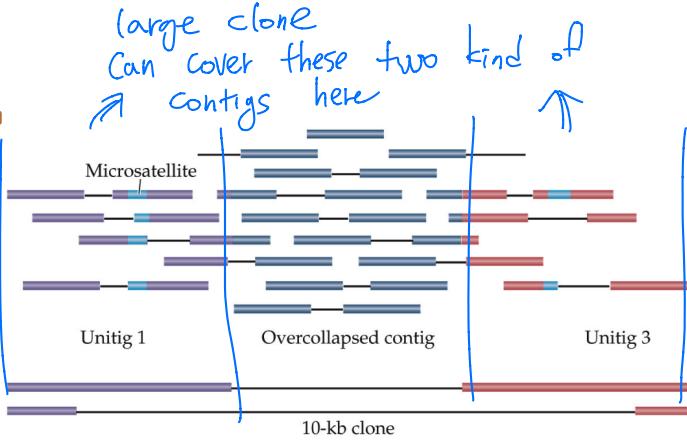
3 billion bp



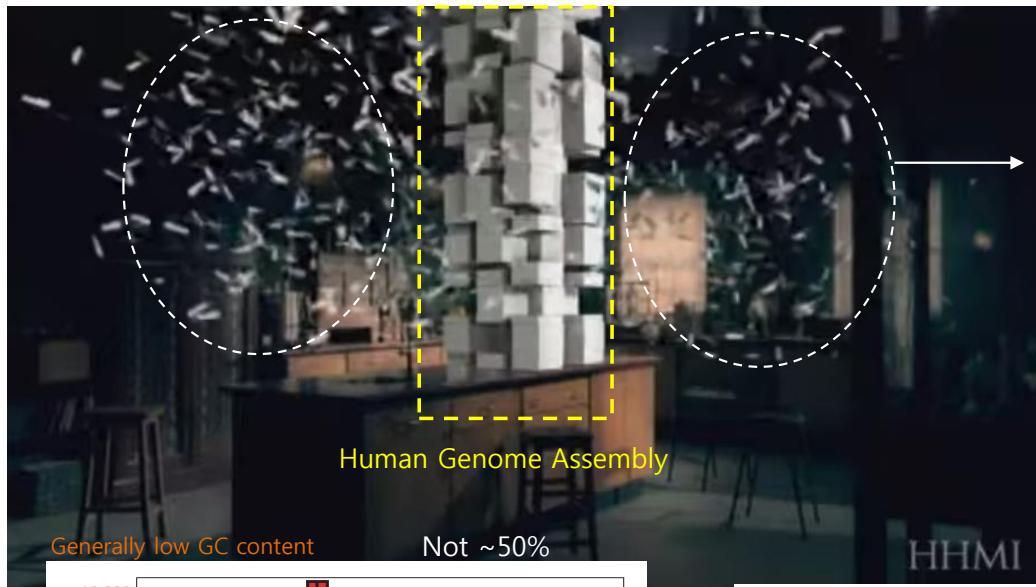
Assembly of sequence reads: issue of gaps caused by repeat sequences



Pair-end reads (mate pair)

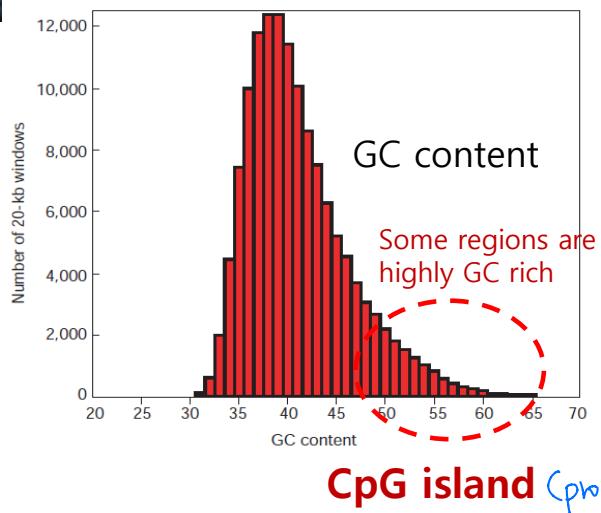


Human genome



Generally low GC content

Not ~50%

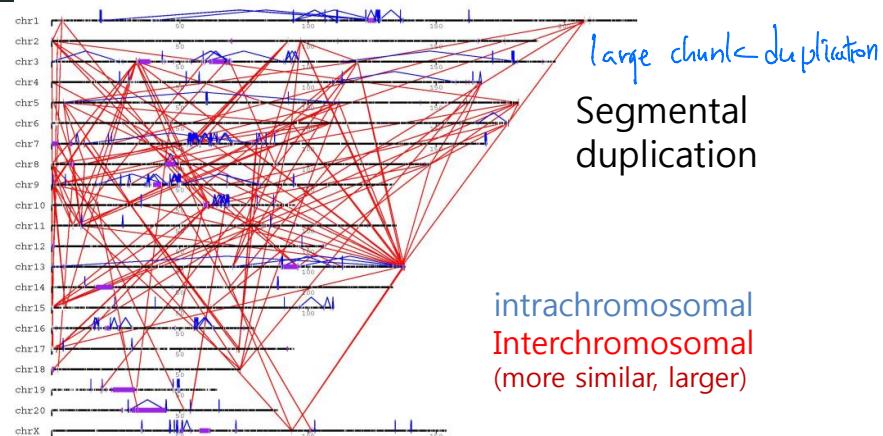


어느 크로모솜에도 위치한
Unlocalized contigs

"chrUn"
Consisted of repeated seq.

크로모솜은 암시지만 위치 모르겠음
Know the chromosome number,
but unplaced

"Chr1 xxxx random"



Next-Generation Sequencing (NGS)

: High-throughput Sequencing

parallelized

병렬화된 секью就越 = 高通量 секью就越

Sung Wook Chi

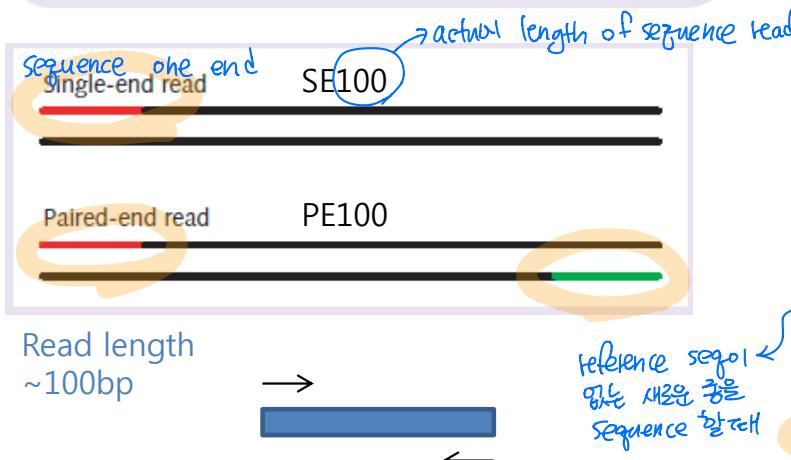
Division of Life Sciences, Korea University

Terminology for DNA sequencing

Fragment: a small piece of genomic DNA – typically several hundred bp in length – subject to an individual partial sequence determination, or *read*.

Single-end read: technique in which sequence is reported from only one end of a fragment (see Figure 1.7).

Paired-end read: technique in which sequence is reported from both ends of a fragment (with a number of undetermined bases between the reads that is known only approximately).



Paired-end sequence Raw sequence obtained from both ends of a cloned insert in any vector, such as a plasmid or bacterial artificial chromosome.



Coverage (or depth) The average number of times a nucleotide is represented by a high-quality base in a collection of random raw sequence. Operationally, a 'high-quality base' is defined as one with an accuracy of at least 99% (corresponding to a PHRED score of at least 20).

Read length: the number of bases reported from a single experiment on a single fragment.

Assembly: the inference of the complete sequence of a region from the data on individual fragments from the region, by piecing together overlaps.

Contig: a partial assembly of data from overlapping fragments into a contiguous region of sequence.

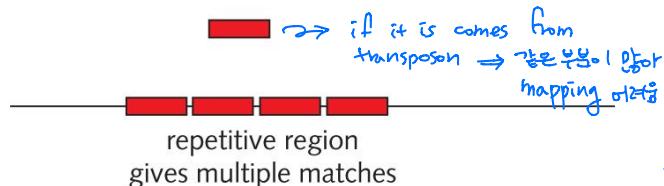
De novo sequencing: determination of a full-genome sequence without using a known reference sequence from an individual of the species to avoid the assembly step.

Resequencing: determination of the sequence of an individual of a species for which a reference genome sequence is known. The assembly process is replaced by mapping the fragments onto the reference genome.

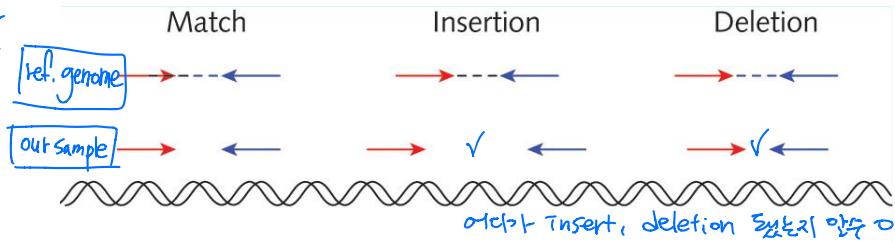
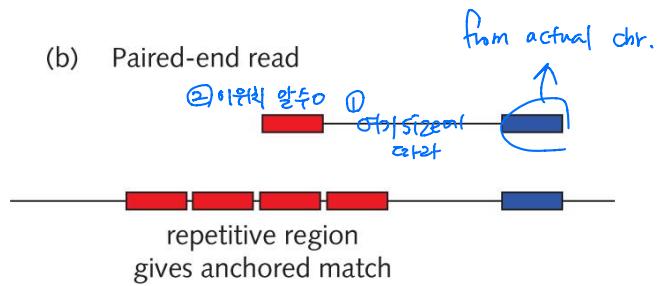
Single-end vs. pair-end read

<PEI 이용해 알수 있는 장점들>

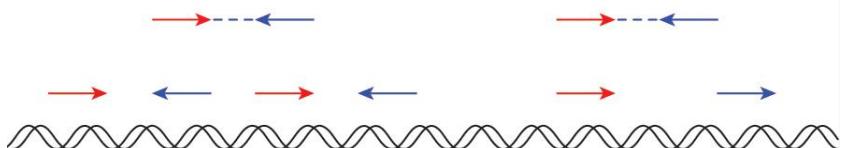
(a) Single-end read



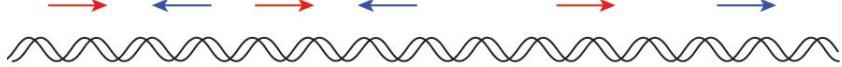
(b) Paired-end read



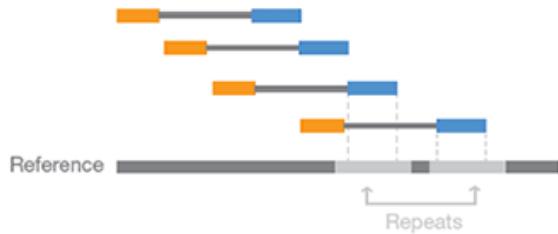
Tandem duplication



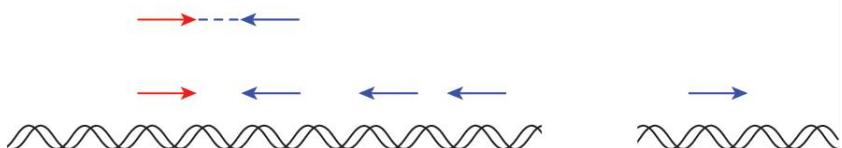
Inversion



Alignment to the Reference Sequence



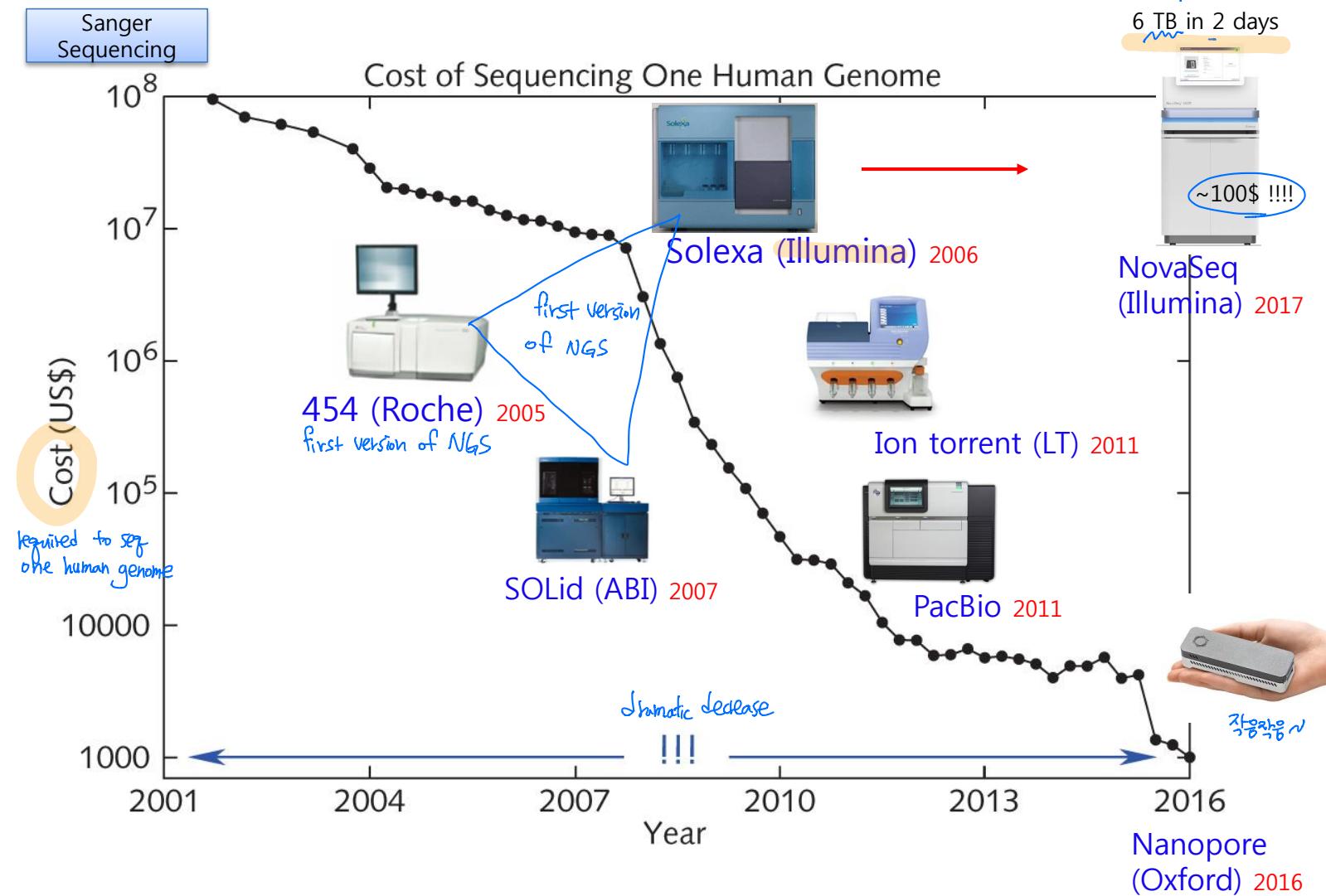
Repeat element insertion



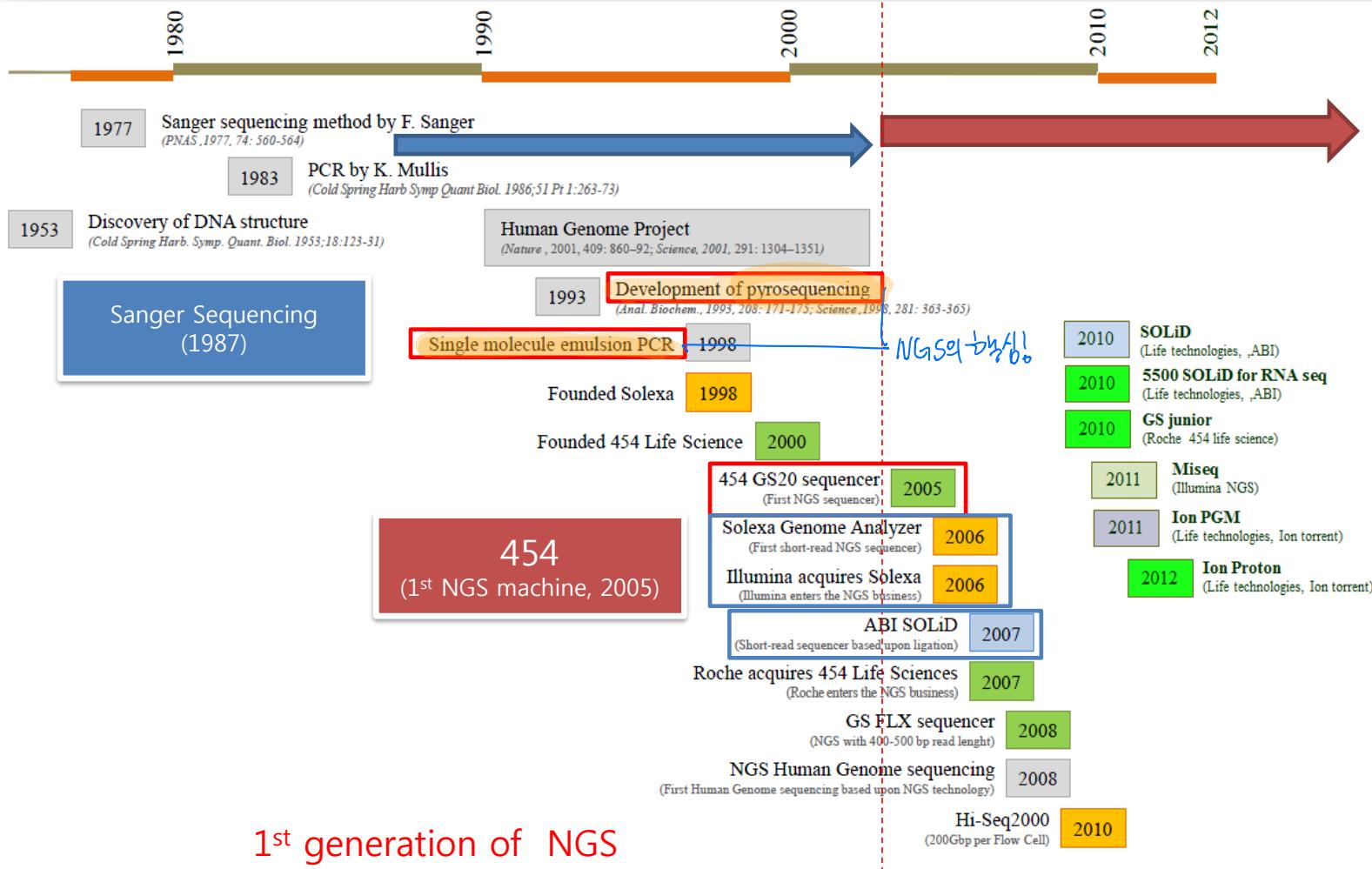
Next-generation Sequencing (NGS)

genome: ~3Gb
→ 10x depth need → ~30Gb data

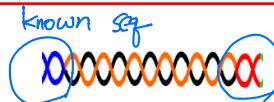
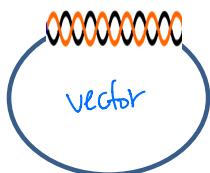
6 TB in 2 days



Sequencing Technology : History



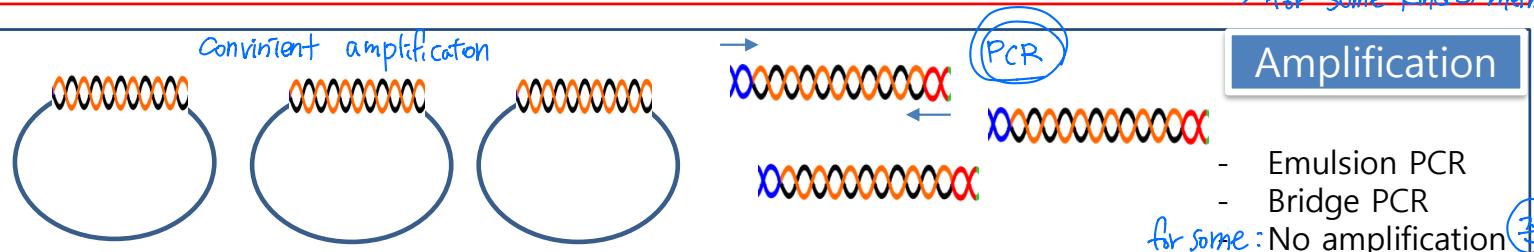
DNA Sequencing



for all fragment
→ γ^32 endonuclease (S1) + sequencing vector or gel electrophoresis (Sanger seq)

Library Preparation

- No preparation
- ↳ for some kind of method



DNA polymerase reaction like Sanger sequencing

Chain termination method

- Pyrosequencing
- SBS (Sequencing by synthesis)
- : reversible chain termination

Sequencing Chemistry

- Ligase reaction
- translocation

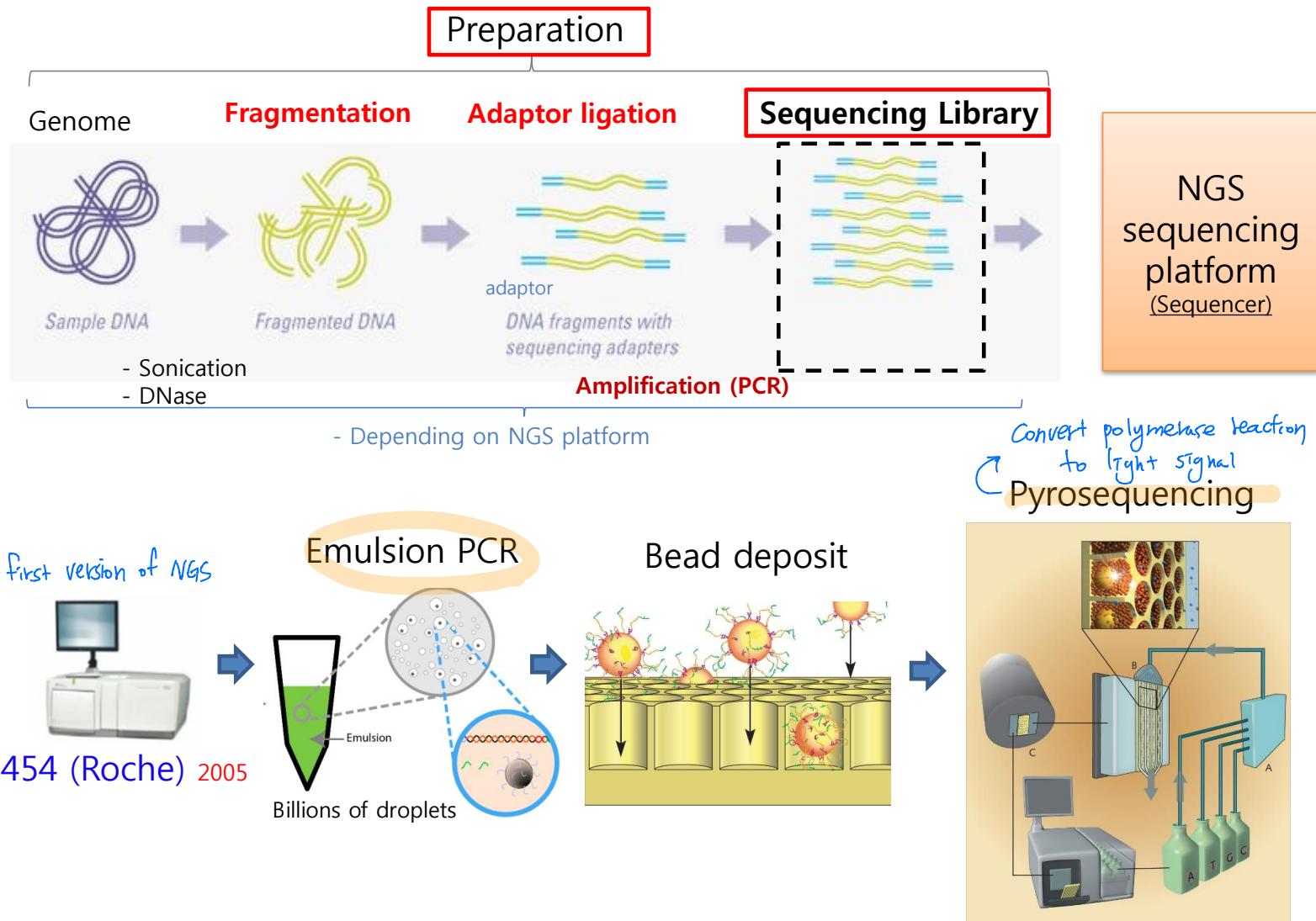
Fluorescence color

- Light
- Fluorescence color
- Electric signal (without light)
→ directly receive signal

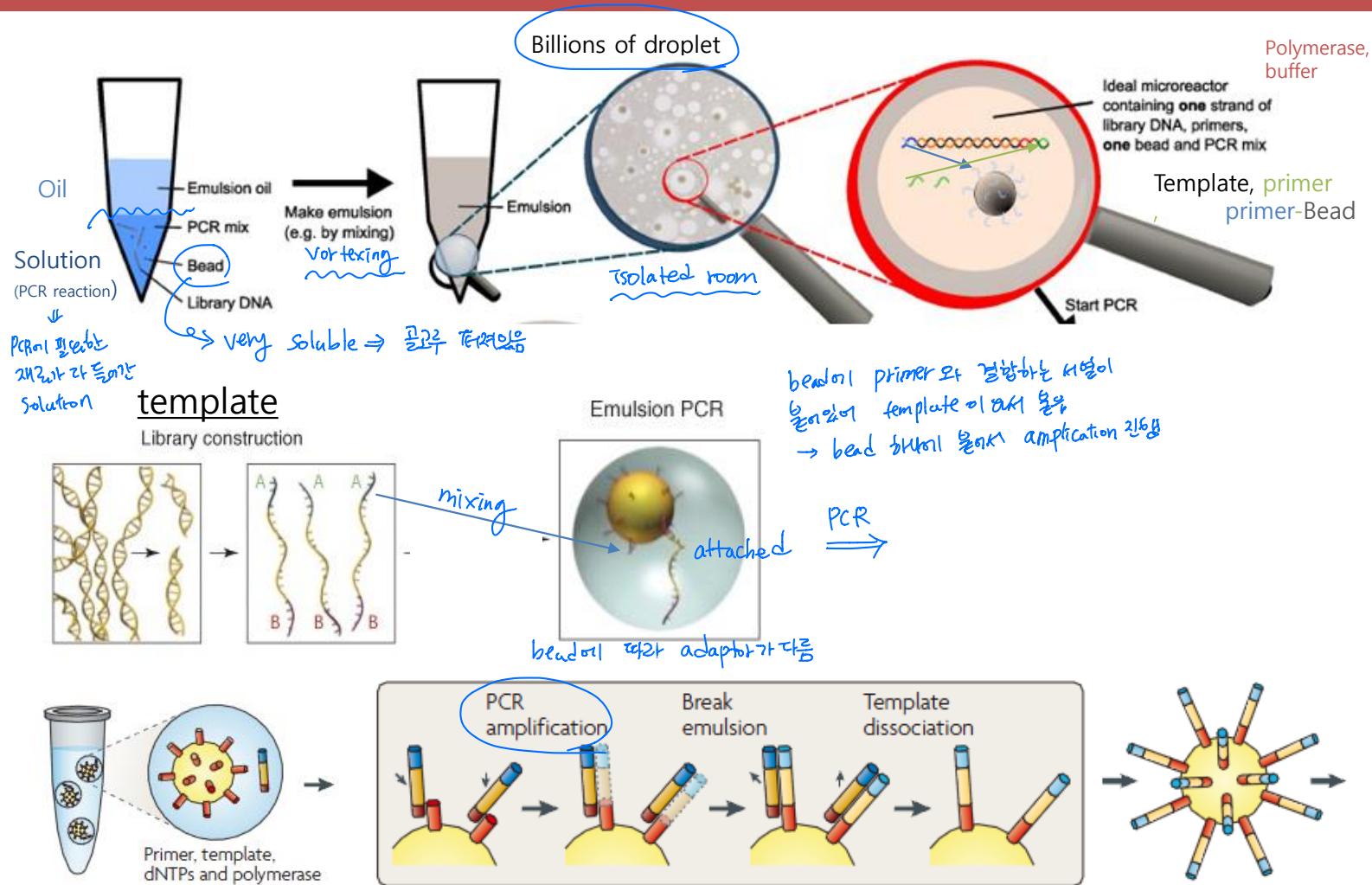
Signal detection

↳ determine A,T,G,C

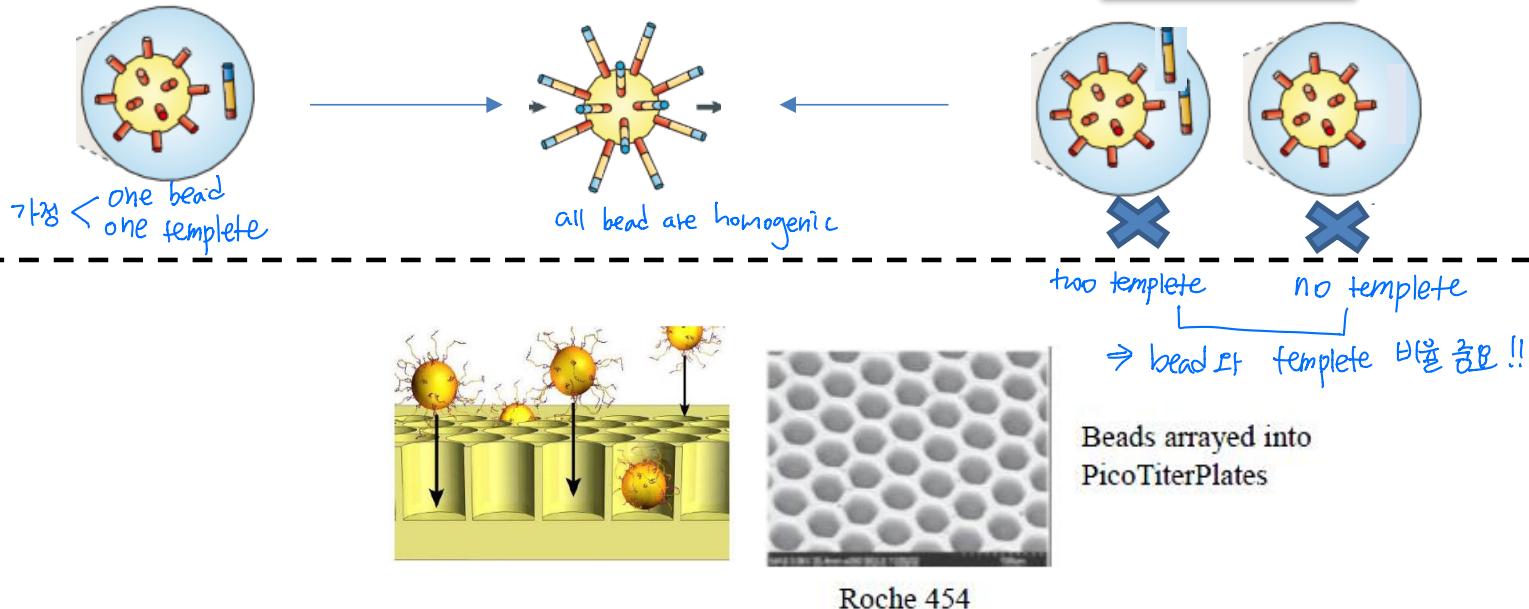
DNA Sequencing by NGS (flow & principle)



Emulsion PCR (454/Rocher, ABI SOLiD)



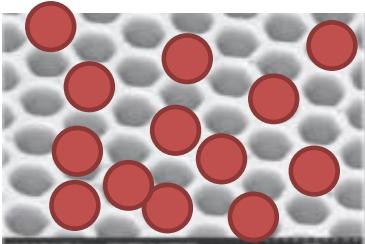
Emulsion PCR & bead loading : critical issues



Not 100% loading



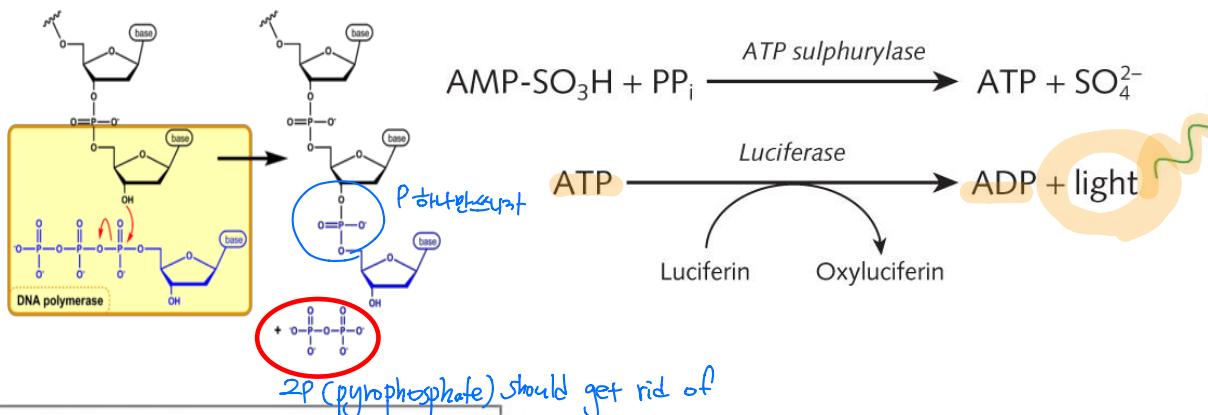
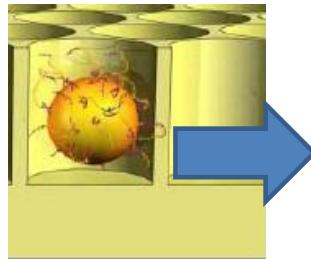
problems



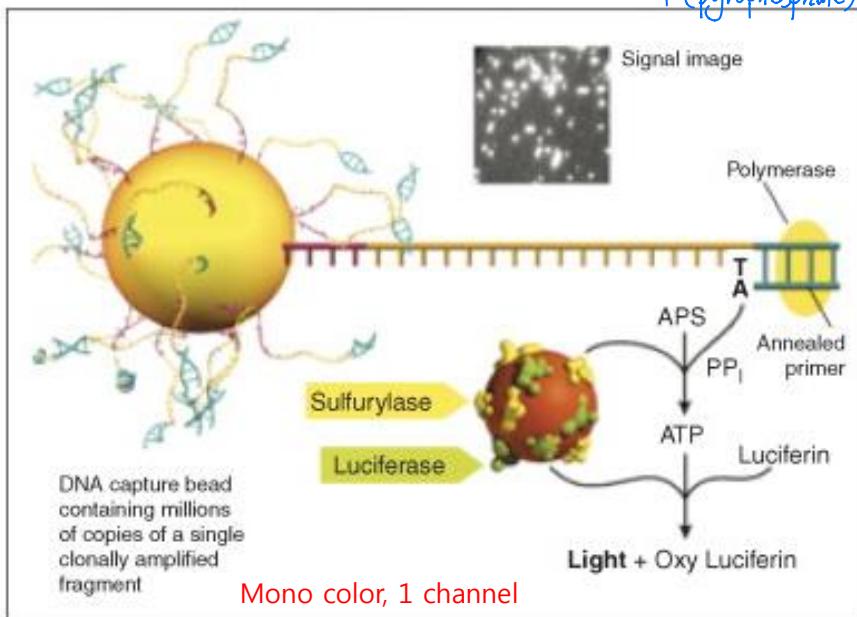
또는 구멍에 bead 허나눠
넣기 어려움
(속면도가 커서 well에 넣는게
어려워짐)

Pyrosequencing (454/Roche)

Polymerase reaction

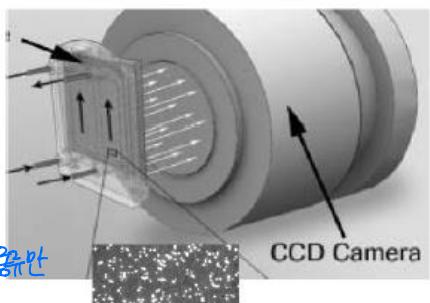


2P (pyrophosphate) should get rid of



PP_i → light

① 한 번에 dNTP 투여



these positions have T

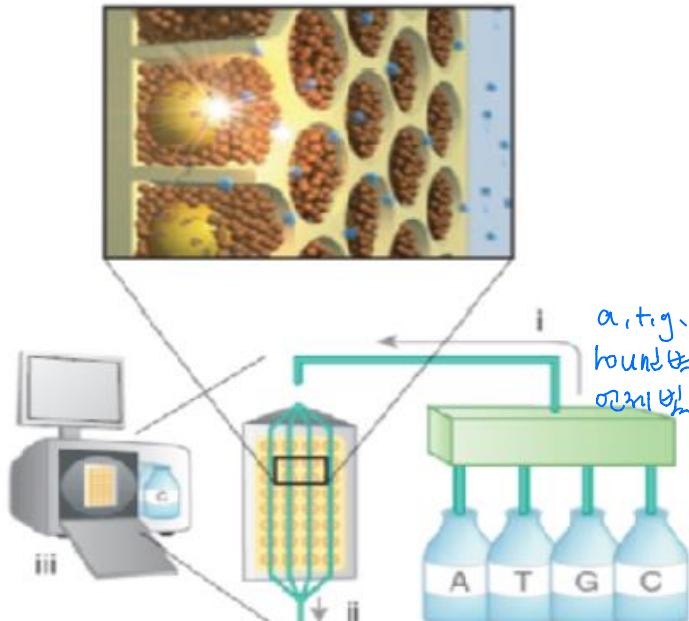
Pyrosequencing (454, Roche) : Problem for homopolymer

③ *

② *

dATP > dTTP > dGTP > dCTP

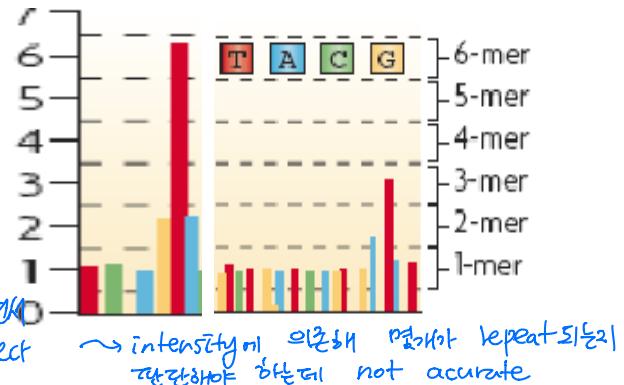
① *



<https://www.youtube.com/watch?v=rsJoG-AuINE>

TCA GG TTTTTT AA

High error rate for reading homopolymer

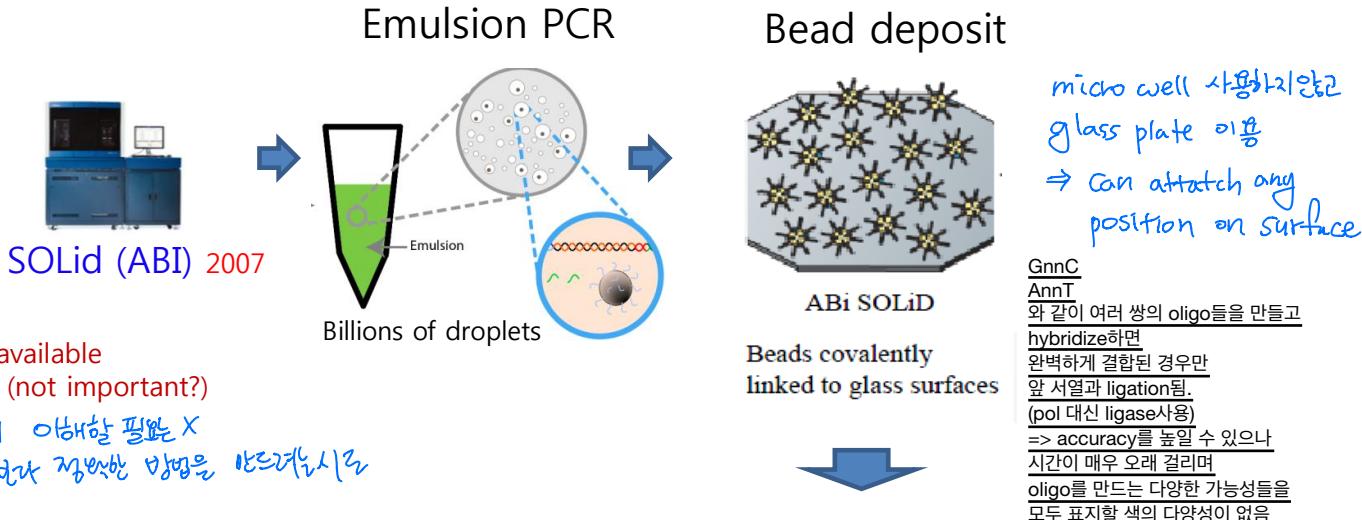


- Long read (up to 700bp),
Low throughput (1 million)
↳ empty well 도 있음 100% 활용률 X
- High accuracy (99.997%)
- Low accuracy for
homopolymer sequence.

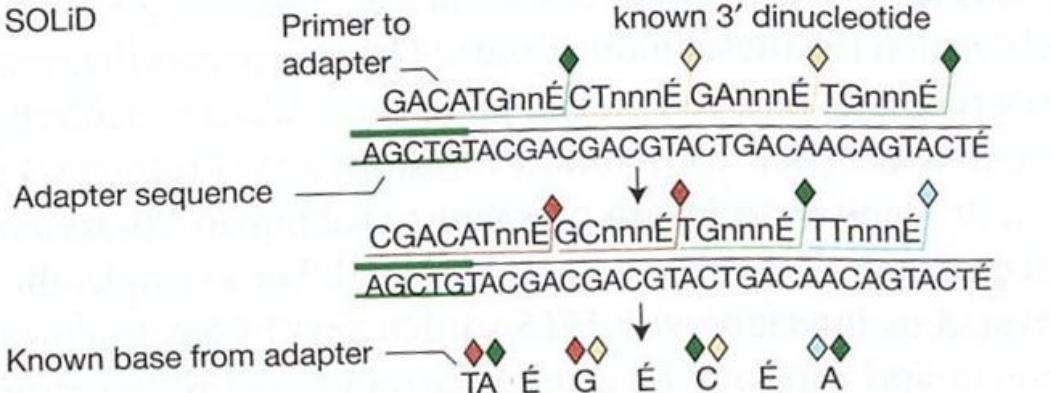
AAAAAA

Sequencing	Amplif.	Chemistry	Read lenght (bp)	Run time (d)	Gbp/day	DNA required (μg)
Roche 454 GS FLX Titanium	emPCR	Pyrosequencing	250-400	0.35 *	1.3	3-5

SOLiD (ABI): ligation based sequencing method



Ligation based sequencing



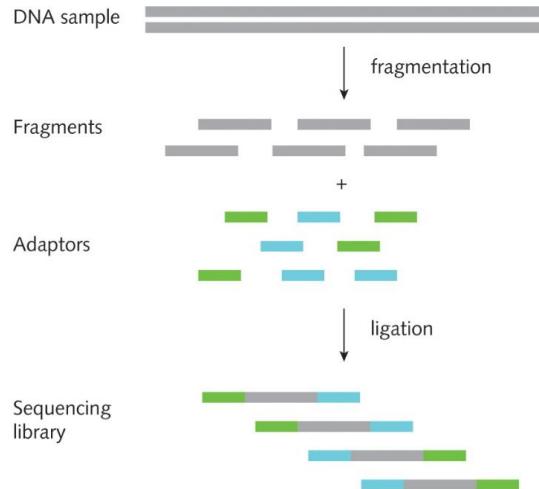
Sequencing by Ligation

- High accuracy
- Large computing time

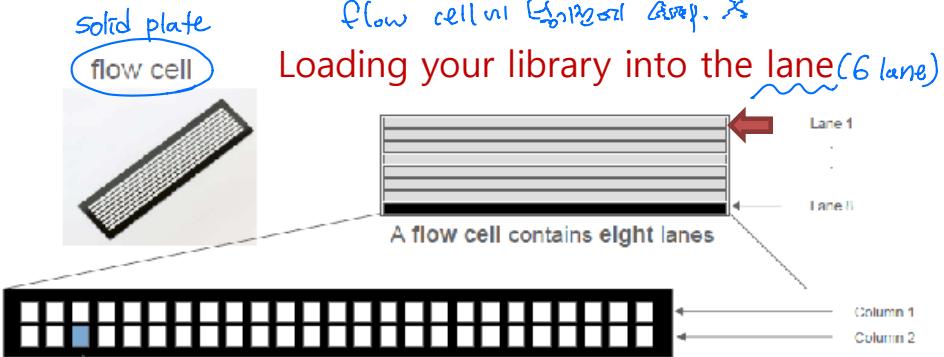
Coding scheme
second base

First base	ACGT
A	◊
C	◊
G	◊
T	◊

Illumina NGS platform (Solexa)



↳ → 광량이 이용되는 platform



flow cell이 네이버로 가요. X

Loading your library into the lane (6 lane)

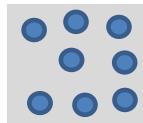
Same as polymerase track



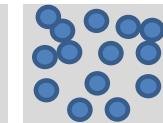
SBS
(sequencing by synthesis)



Sparse
(low number of reads)
도 아닙....

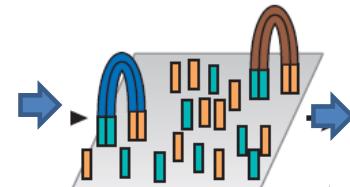


Optimal

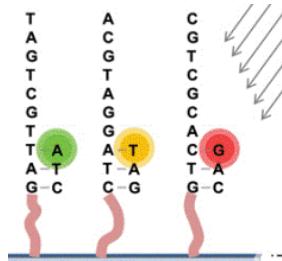


crowded
(Interfering of signal;
low quality)

Bridge amplification



Solid phase amplification

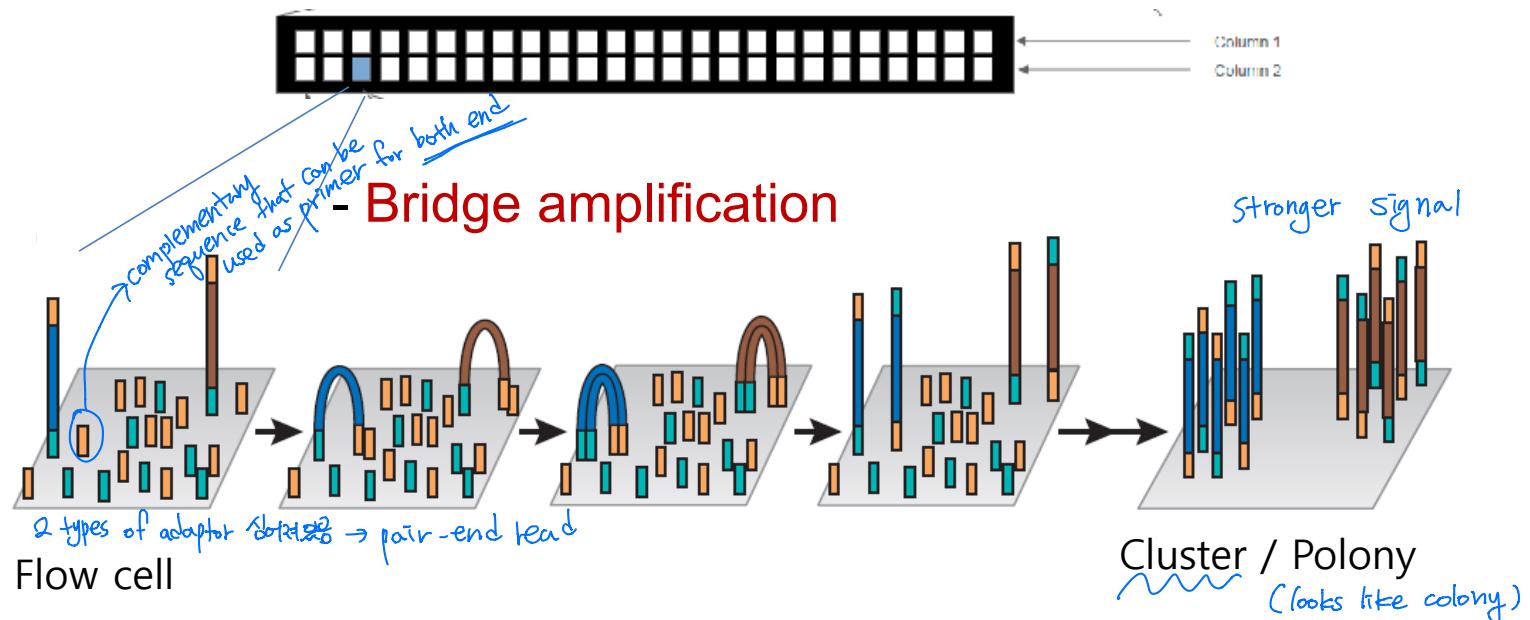


Reversible termination
chemistry

Loading (spread) your library
on glass plate (flow cell)

Solexa (Illumina)
2006

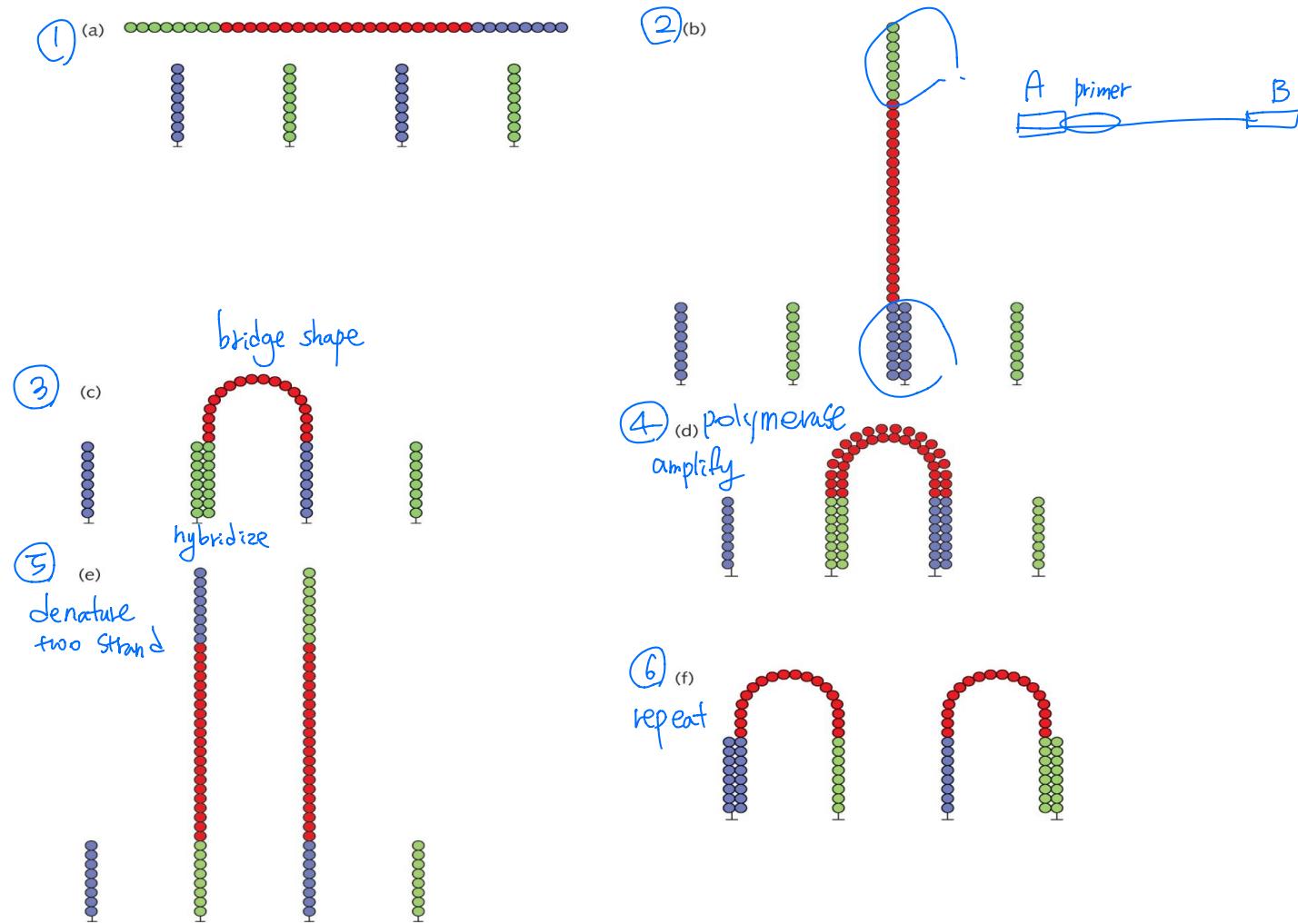
Bridge amplification : Library amplification in Illumina platform



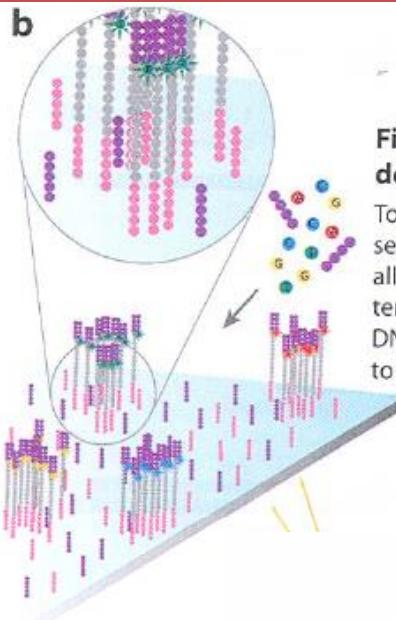
- DNA fragments are flanked with adaptors.
- A flat surface coated with two types of primers, corresponding to the adaptors.
- Amplification proceeds in cycles, with one end of each bridge tethered to the surface.

Cluster (polony) generation on solid-phase

Bridge amplification : Library amplification in Illumina platform



Sequencing chemistry in Illumina platform

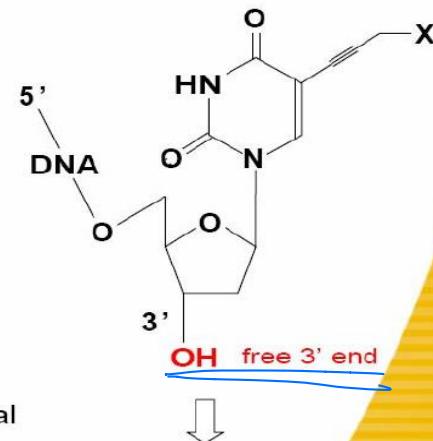
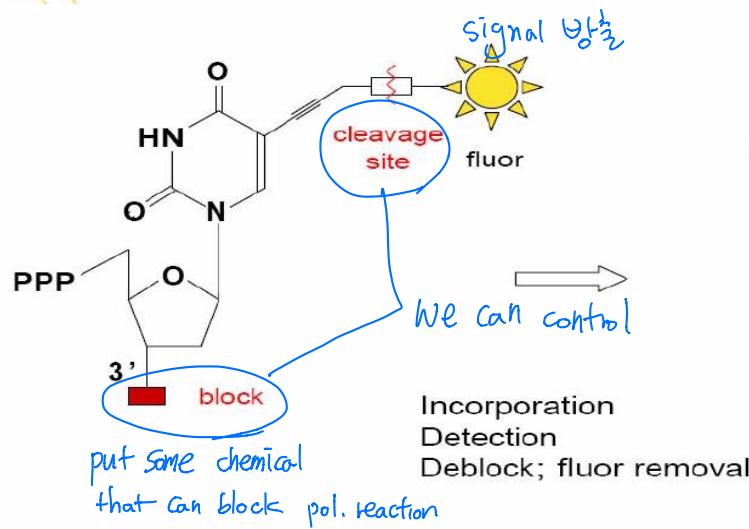


Reversible Terminator Chemistry

Solexa

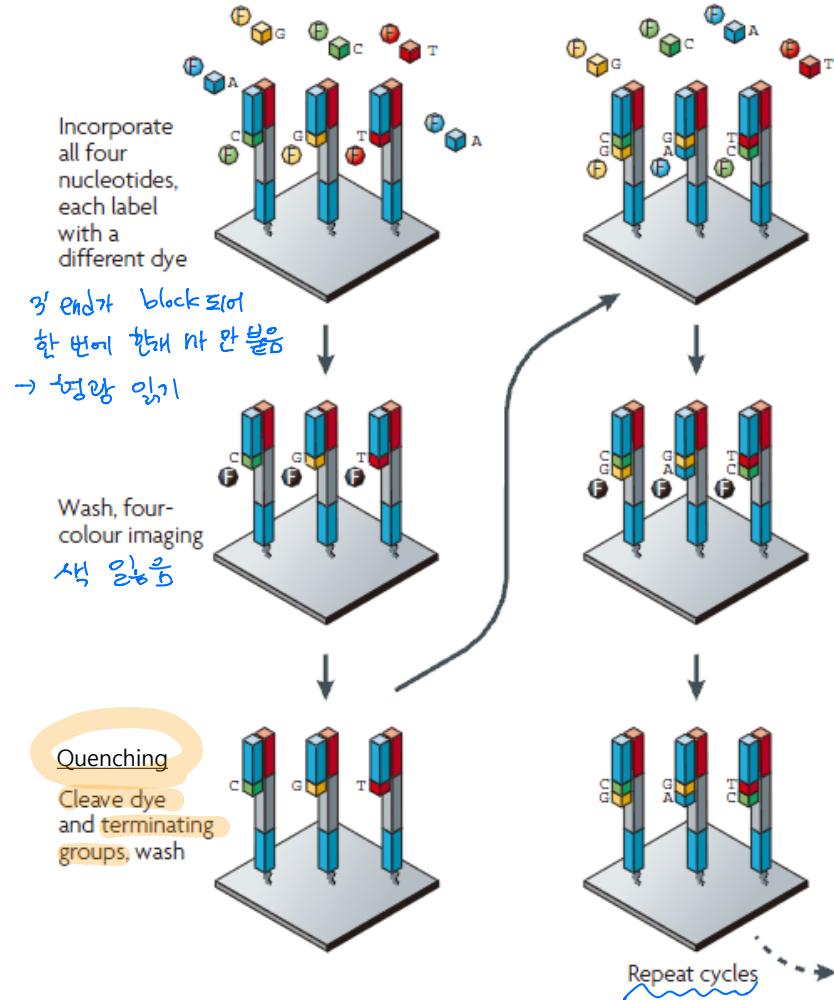
- All 4 labelled nucleotides in 1 reaction
- Higher accuracy
- No problems with homopolymer repeats

SBS (Sequencing by synthesis)



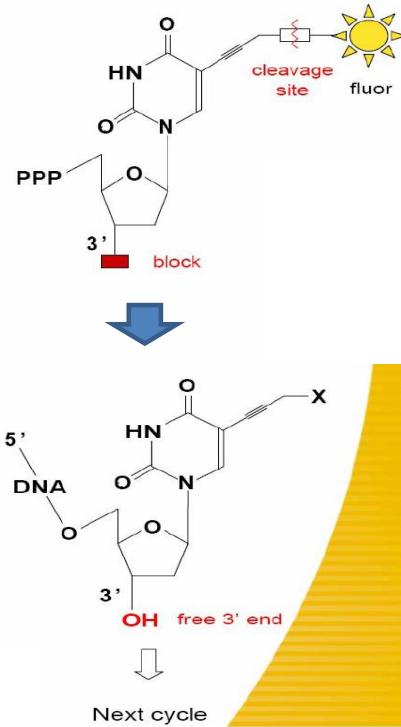
Sequencing chemistry in Illumina platform

a Illumina/Solexa — Reversible terminators



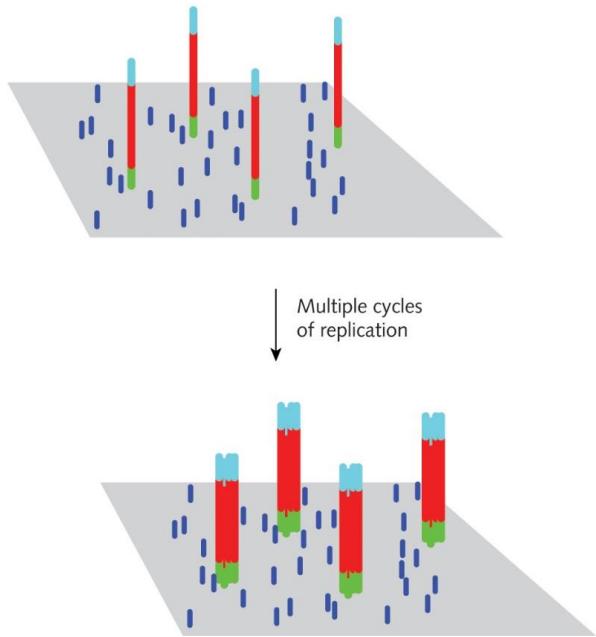
one cycle → detect one nt

- SBS (Sequencing by synthesis)
Reversible termination method

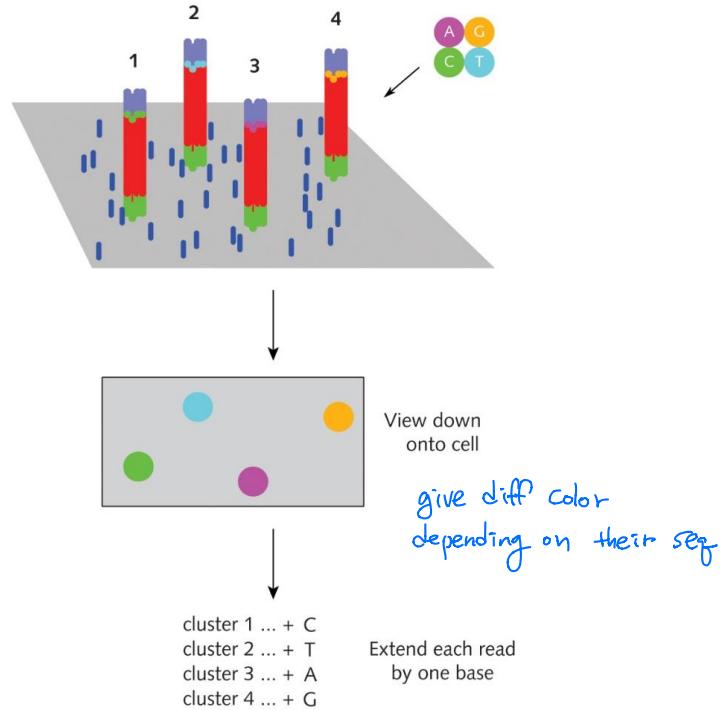


Solexa (Illumina) NGS method

Bridge amplification



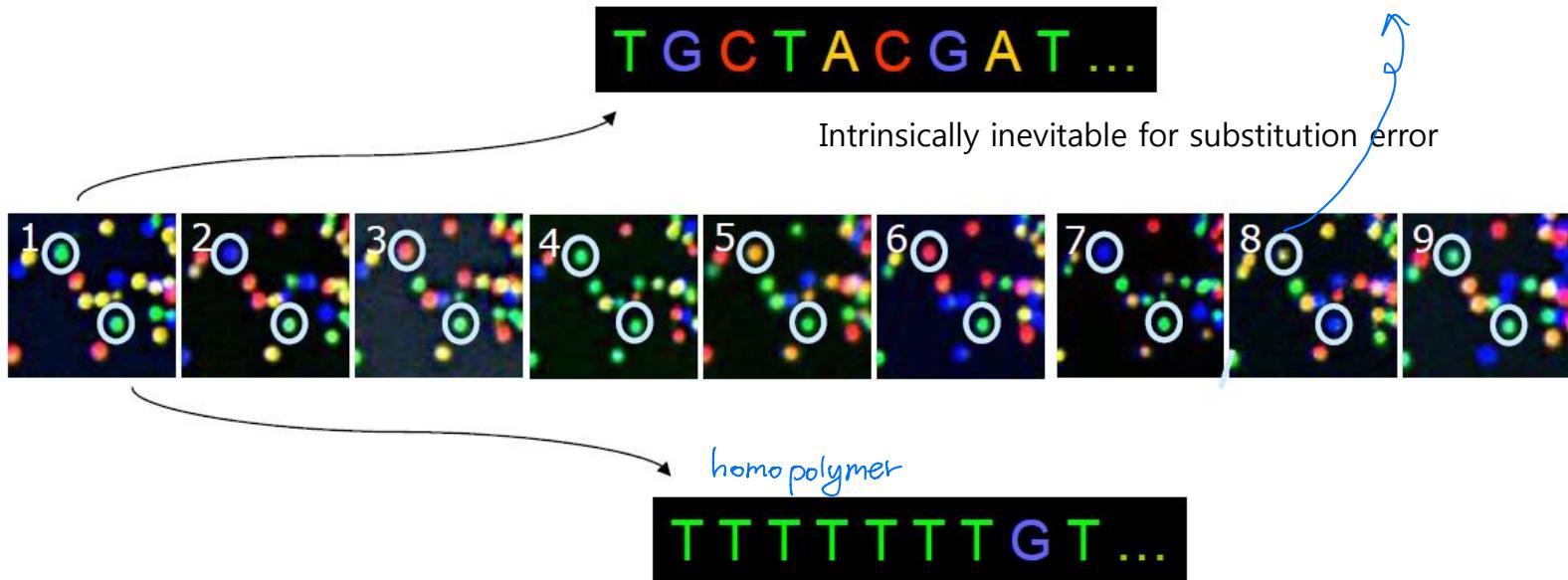
SBS (sequencing by synthesis)



Reversible termination chemistry

Base calling from image (Illumina)

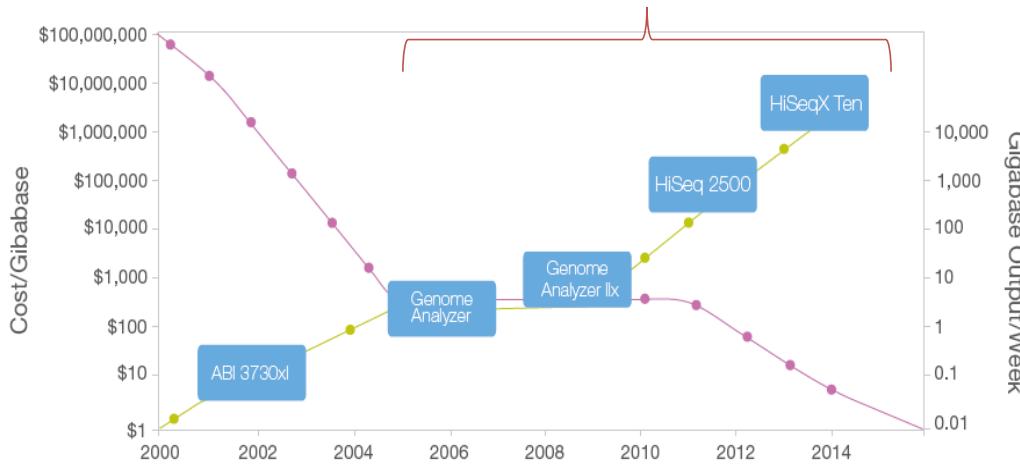
looks like colors are mixed
Quenching of dyes
signals interfere



⇒ not able to read long length because of reading error

Illumina NGS platform (Solexa)

Winner of NGS market



Illumina (NovaSeq)
6 TB in 2 days

Small size



MiniSeq System

Power and simplicity for targeted sequencing.

MiSeq Series

Small genome and targeted sequencing.

NextSeq Series

Everyday genome, exome transcriptome sequencing, and more.

HiSeq Series

Production-scale genome, exome, transcriptome sequencing and more.

HiSeq X Series

Population- and production-scale human whole-genome sequencing.

<https://www.youtube.com/watch?v=qtVM18Twi0>

HiSeq 2500 System (High-Output Mode)

HiSeq 2500 System (Rapid Run Mode)

29 hours–6 days	7–60 hours
1000 Gb	300 Gb
4 billion	600 million
2 × 125 bp	2 × 250 bp

limit of length

HiSeq platform: 300 million reads, 150 PE

1st generation NGS platforms (2005-)

Solexa / Illumina



Illumina GAIL

Market winner

454 /Roche



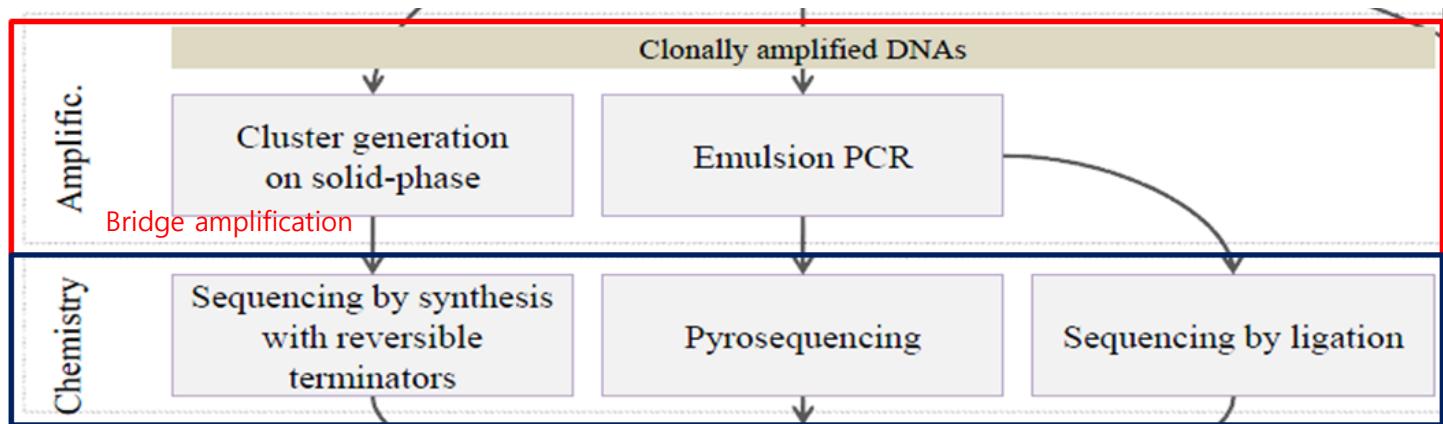
Roche 454

First NGS sequencer



ABi SOLiD

Terminated



4 channel
fluorescence

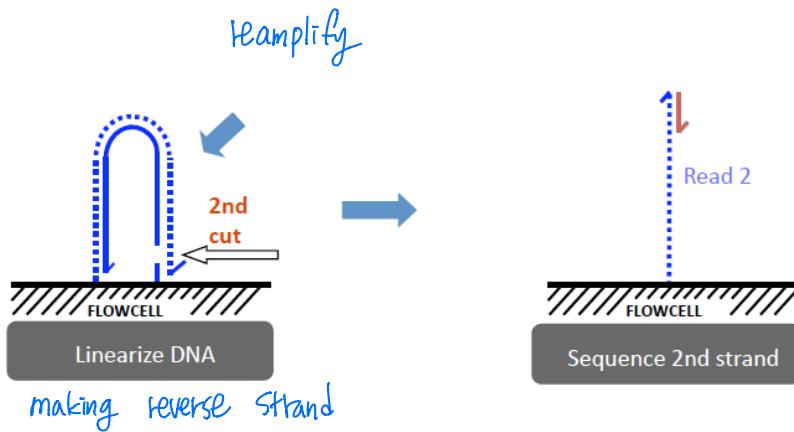
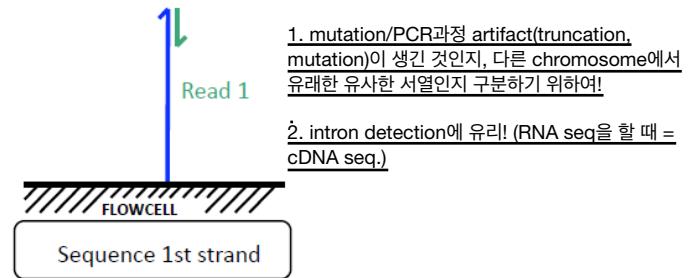
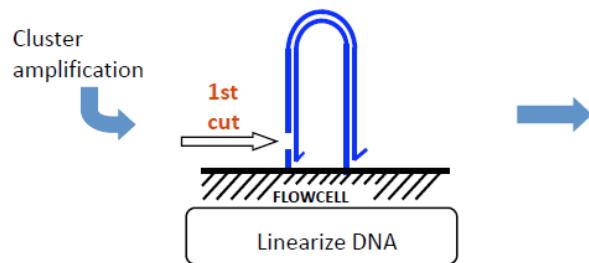
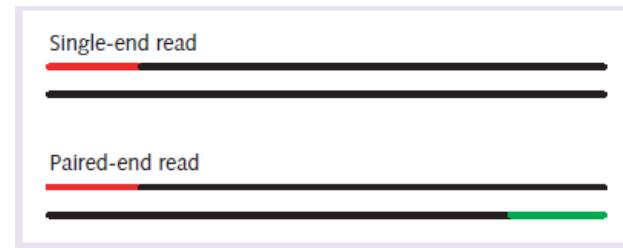
1 channel light

4 channel
fluorescence

Paired-End Sequencing (Illumina)

Illumina/Solexa sequencing

Single-end or Pair-end read



Multiplexing with Barcode: Sample Indexing

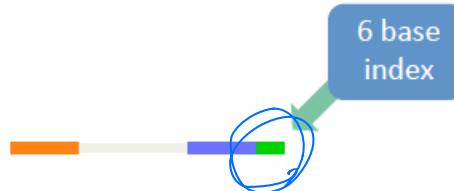
Illumina/Solexa sequencing

put multiple sample in one lane
→ 같은 샘플 다른 바코드로 구분

Sequencing multiple samples in a single channel, reducing cost/sample

Simple construct design

Based on paired end process

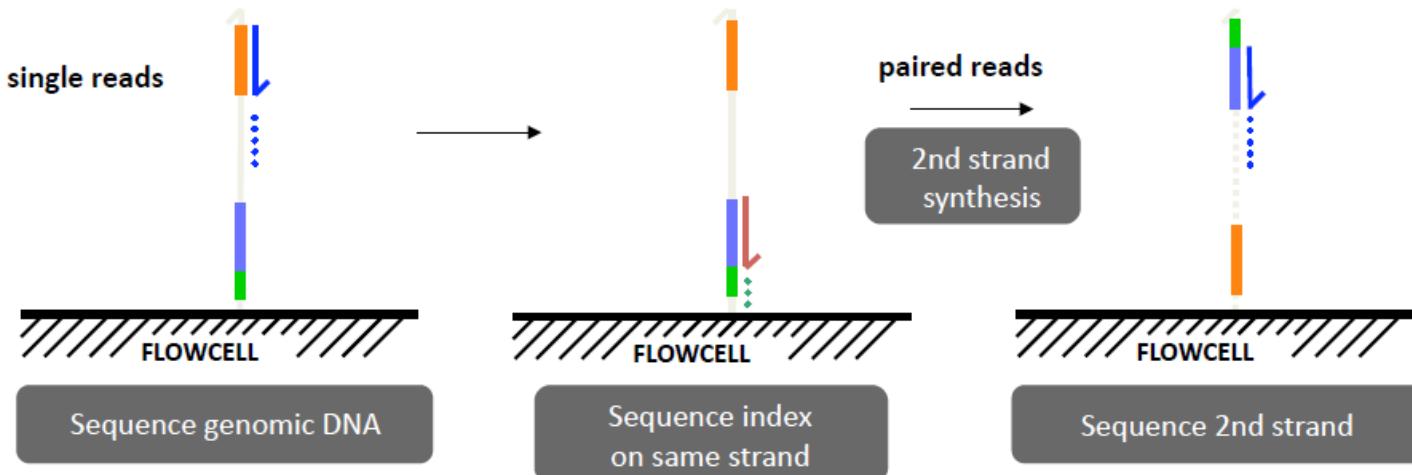


1 lane : 240 million reads

12 samples x 20 million reads

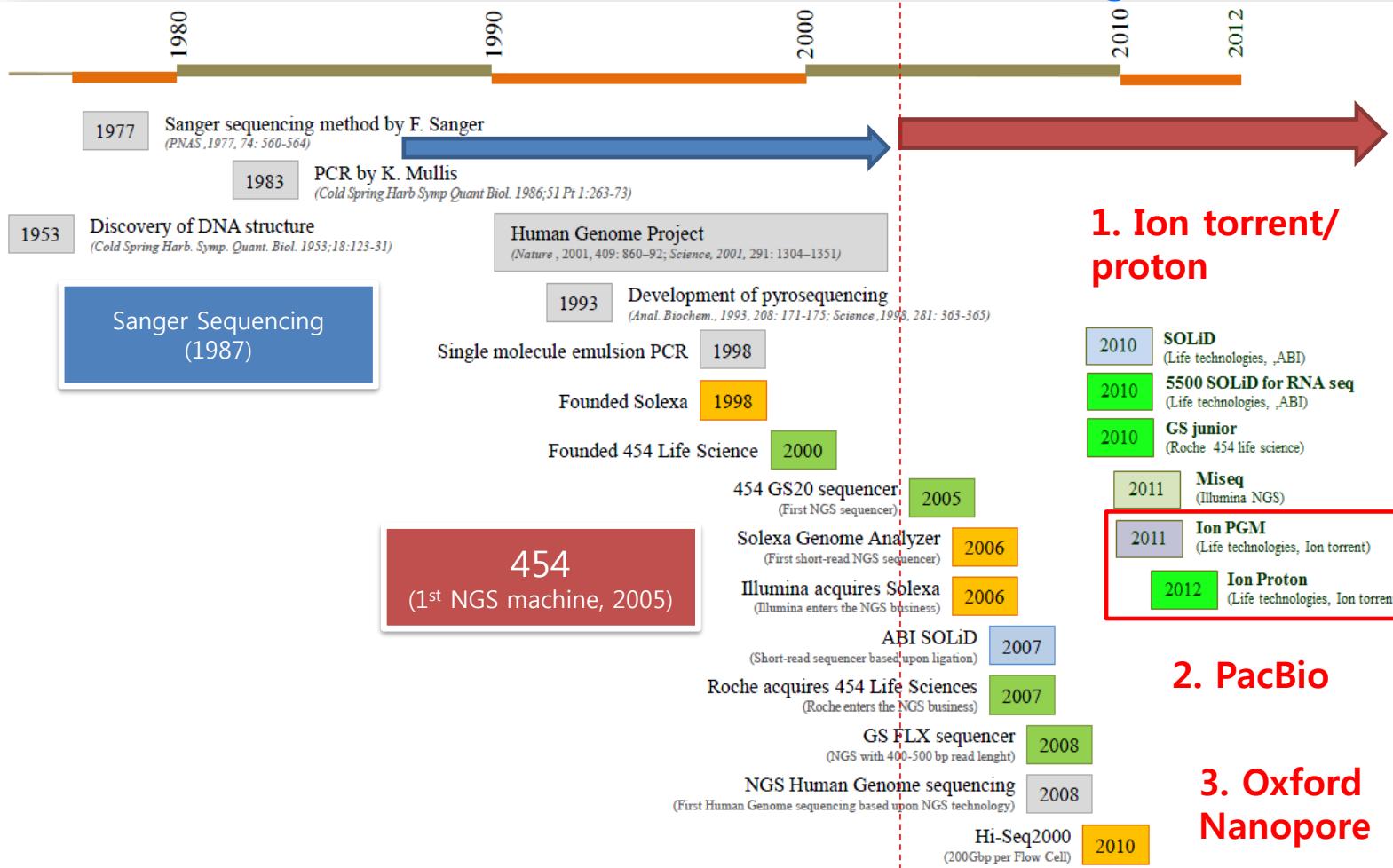
Sequencing protocol

Automated processing of indexing read

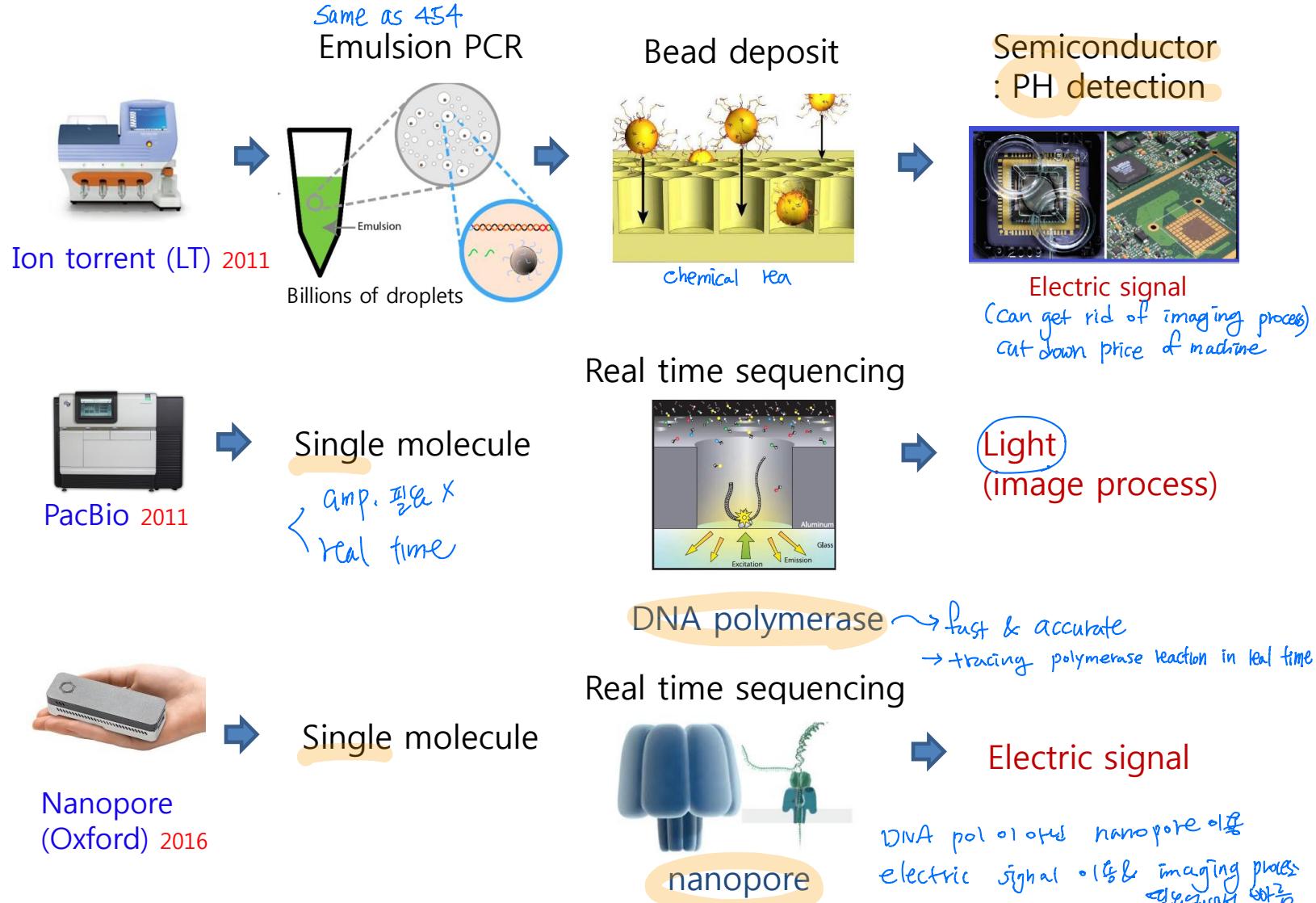


2nd Generation of NGS sequencer

⇒ second generation

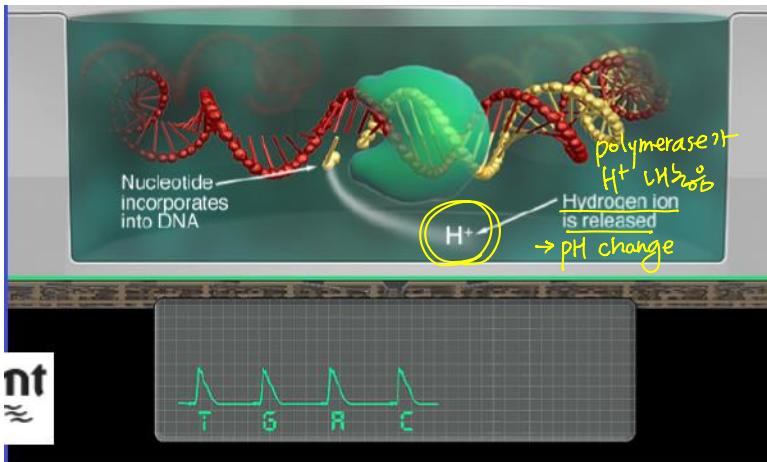
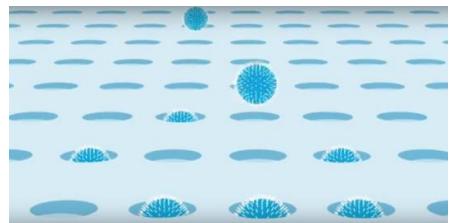


2nd Generation of NGS sequencer



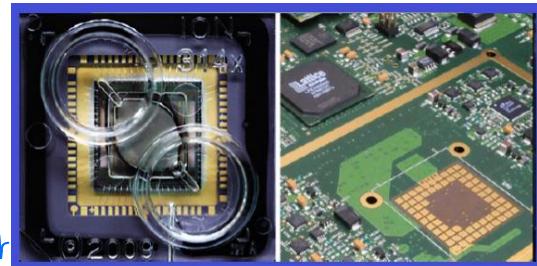
Ion torrent/proton: semiconductor sequencing

Emulsion PCR

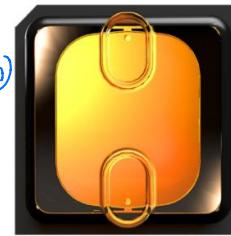


Electric signal
(No camera !!, low price)

1.5 million pH meters semiconductor chip



Ion
torrent
use chips
with pH detector



80 million pH meters semiconductor chip

각 패셀을 pH 미터로 별도로

DNA polymerase reaction

Release of proton (H^+) > change in pH > detecting as electric signal

well을 100만 이용 X

The same problems in

- Emulsion PCR, Pyrosequencing

homopolymer의
연속된碱 number X

Pacific Biosciences

No amplification (Single molecule)

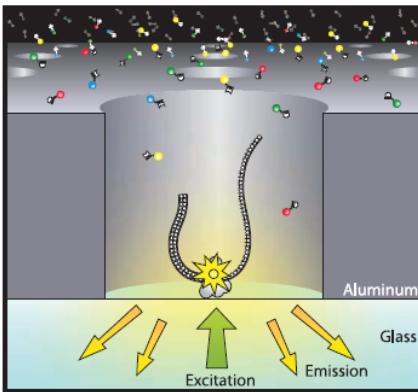
Long reads (3kb ~ 15kb),

Fast (real time, 30 min)

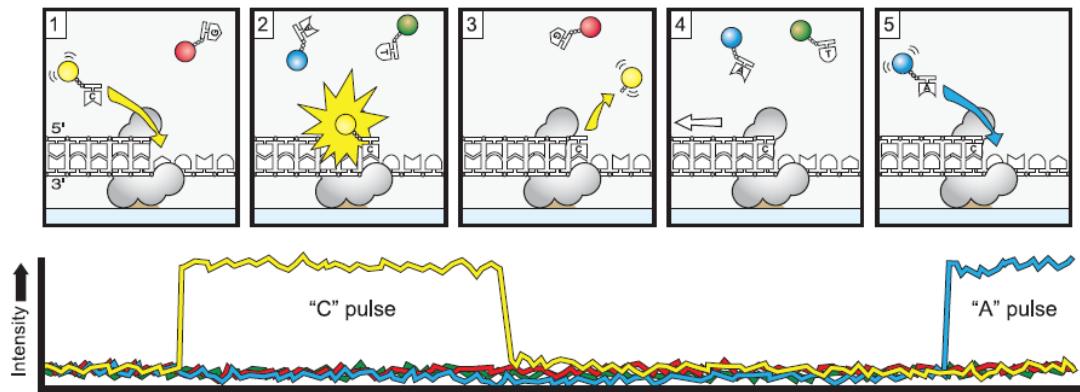
↳ but too fast to detect with our technology

tiny amount of signal & too fast \Rightarrow hard to detect
but low accuracy (95%)
low throughput (70,000 reads)

A



B

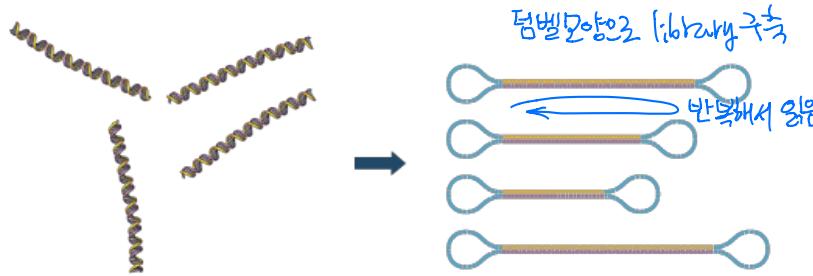


Single Molecule, Real Time Sequencing

<https://www.youtube.com/watch?v=NHCJ8PtYCFc>

Useful for long read sequencing
(resequencing genomes)

SMRT (Single Molecule RealTime) Sequencing : PacBio



Library Preparation

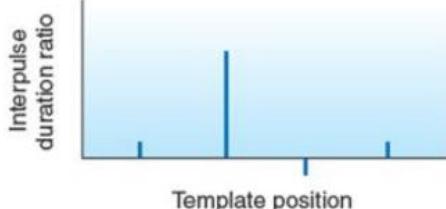
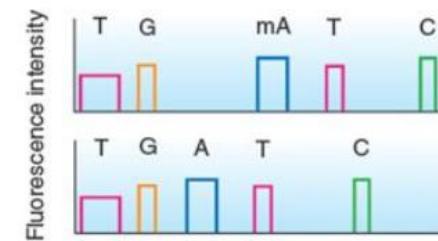
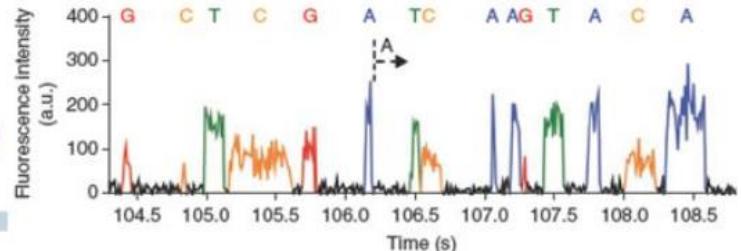
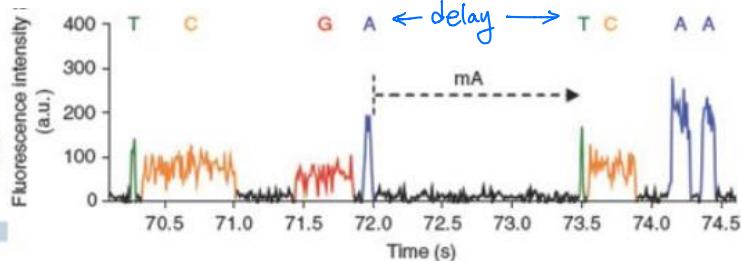
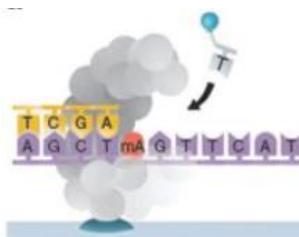
- Universal SMRTbell template accepts insert sizes from 250 bp to 40 kb

Redundant Sequencing Reading

Consensus accuracies > 99.999%
repeat several time → reduce error



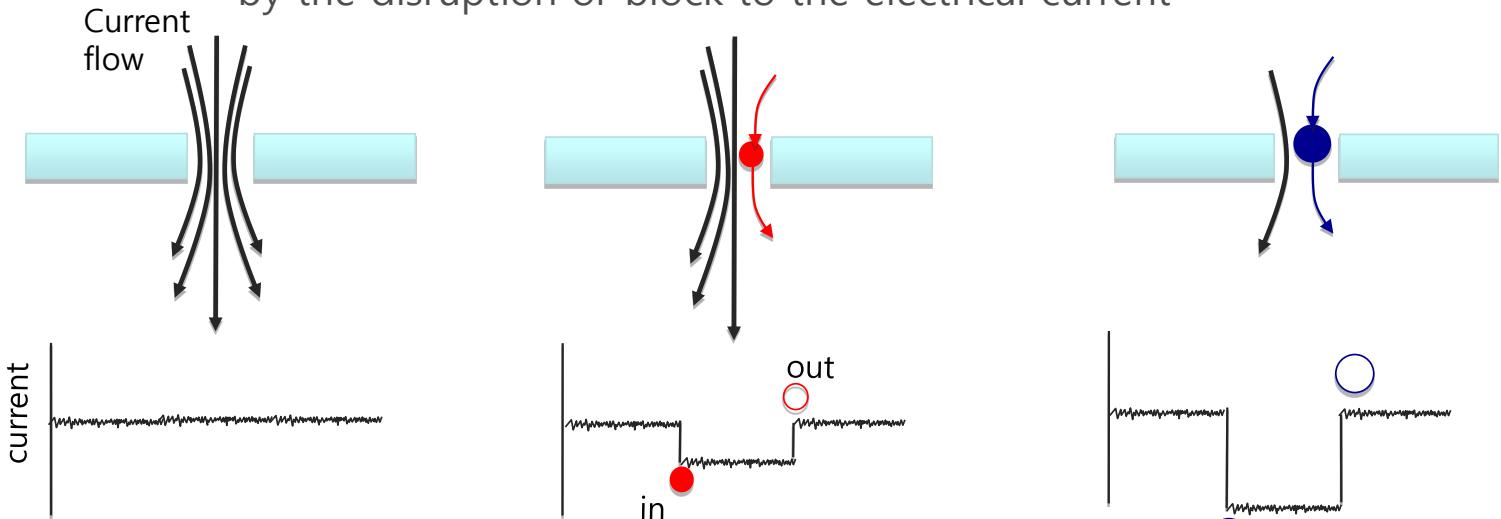
Realtime direct detection of modified DNA



Oxford Nanopore

Electronic &
Single molecule
Real time

- Nanopore = 'very small hole'
- Electrical current flows through the hole
- Introduce analyte of interest into the hole → identify "analyte" by the disruption or block to the electrical current



nature
nanotechnology

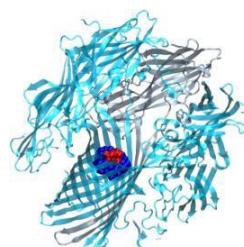
PUBLISHED ONLINE: XX XX 2009 | DOI: 10.1038/NNANO.2009.12

ARTICLES

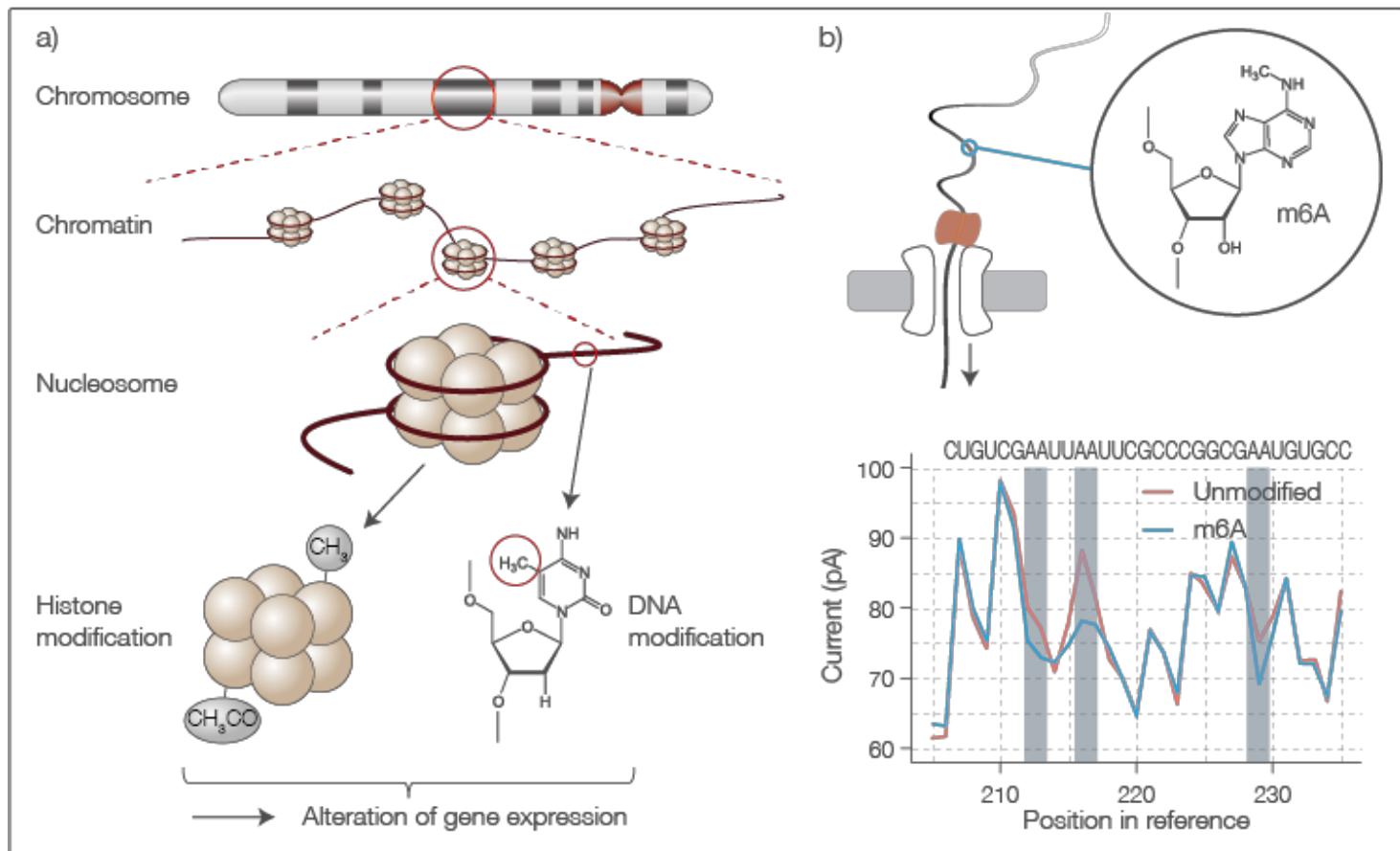
Continuous base identification for single-molecule nanopore DNA sequencing

James Clarke¹, Hai-Chen Wu², Lakmal Jayasinghe^{1,2}, Alpesh Patel¹, Stuart Reid¹ and Hagan Bayley^{2*}

<https://www.youtube.com/watch?v=BNz880V52rQ>



Direct detection of modified DNA/RNA by Oxford Nanopore

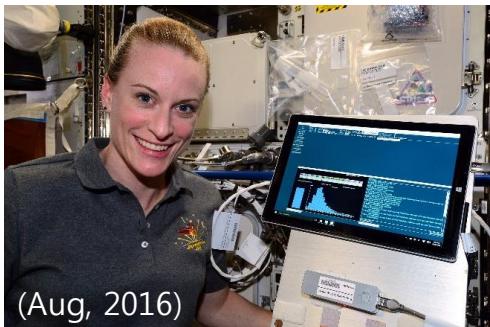


Realtime Portable Sequencer: Oxford Nanopore

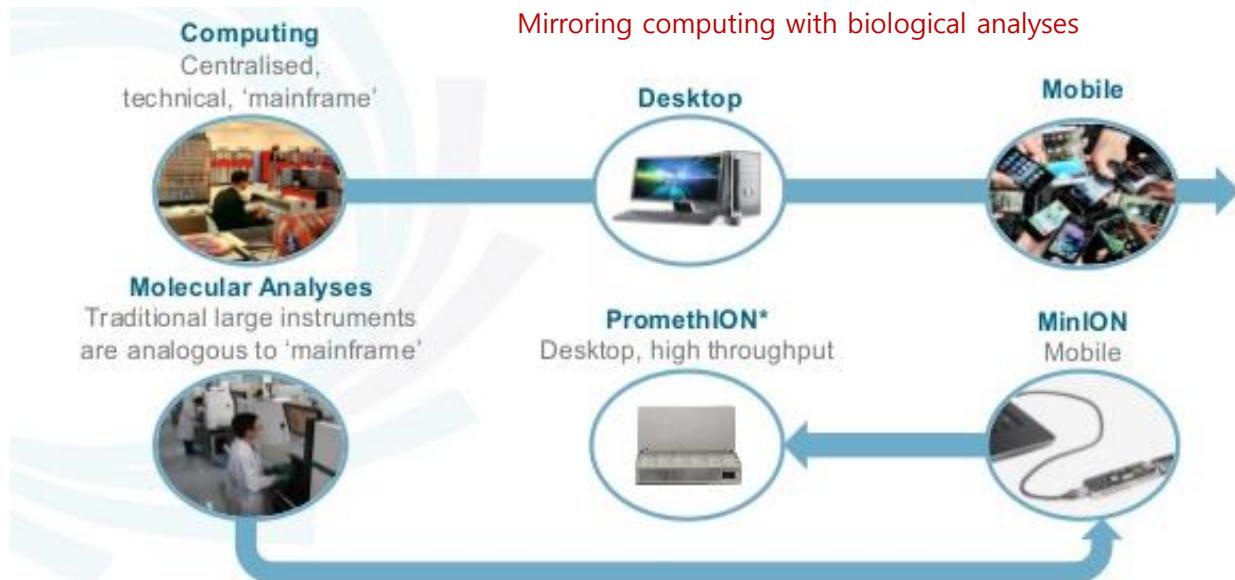
Oxford Nanopore (**MinION**)



9th NASA/SpaceX commercial cargo

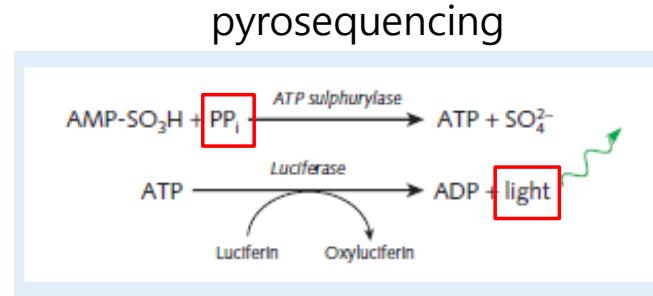
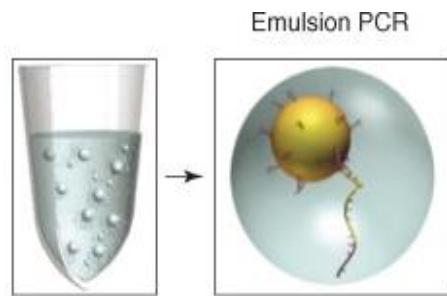


Creation of internet of living things (DNA)

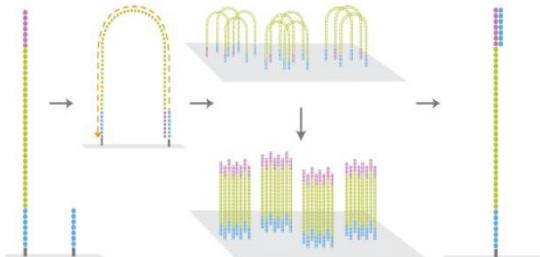


NGS platform comparison

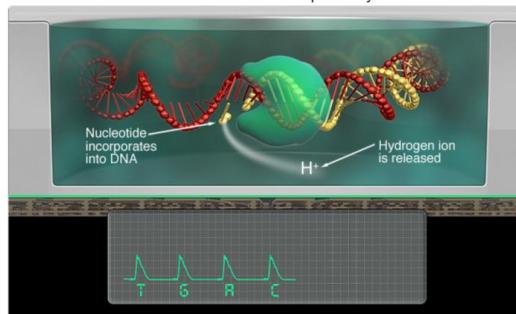
Platform	Amplification	Chemistry	Detection	description
454	emPCR	Pyrosequencing	Image (mono)	Low throughput, long read (~400), ~0.35 day
Abi SOLiD	emPCR	Seq by ligation	Image (color)	High throughput, short read (~50), ~2 weeks
Illumina /Solexa	Solid-phase	Reversible terminator	Image (color)	High throughput, short read (~50), ~1 week
Ion torrent /proton	emPCR	Pyrosequencing	H ⁺ (pH)	low-High throughput, short read (~50), 5 hours
Pacific Bioscience	Single molecule	realtime, polymerase	Image (color)	longer read (3-15kb), Accuracy (~95%), 30min
Oxford Nanopore	Single molecule	realtime, polymerase	H ⁺ (pH)	?



Solid-phase (bridge amplification)



Proton (pH meter)



NGS platform comparison

Method	Generation	Read length (bp)	Single pass error rate (%)	No. of reads per run	Time per run	Cost per million bases (US\$)
Sanger ABI 3730x1	First	600–1000	0.001	96	0.5–3 h	500
Ion Torrent	Second	200	1	8.2×10^7	2–4 h	0.1
454 (Roche) GS FLX+	Second	700	1	1×10^6	23 h	8.57
Illumina HiSeq 2500 (High Output)	Second	2 × 125	0.1	8×10^9 (paired)	7–60 h	0.03
Illumina HiSeq 2500 (Rapid Run)	Second	2 × 250	0.1	1.2×10^9 (paired)	1–6 days	0.04
SOLiD 5500x1	Second	2 × 60	5	8×10^8	6 days	0.11
PacBio RS II: P6-C4	Third	1.0–1.5 × 10 ⁴ on average	13	$3.5\text{--}7.5 \times 10^4$	0.5–4 h	0.40–0.80
Oxford Nanopore MinION	Third	2–5 × 10 ³ on average	38	$1.1\text{--}4.7 \times 10^4$	50 h	6.44–17.90

From: Rhoads, A. and Au, K.F. (2015). PacBio Sequencing and Its Applications. *Genomics, Proteomics & Bioinformatics*, **13**, 278–289. (This article gives citations of sources of data.)

Application of Next-Generation Sequencing

DNA

Genome
Resequencing

Methylation
Analysis
(Bisulfite sequencing)

Functional
Elements
(ChIP-Seq, DNAse-Seq)



RNA

mRNA Tag
Profiling
(HITS-CLIP)

Small RNA
Identification

Transcriptome
Sequencing
(RNA-Seq)