

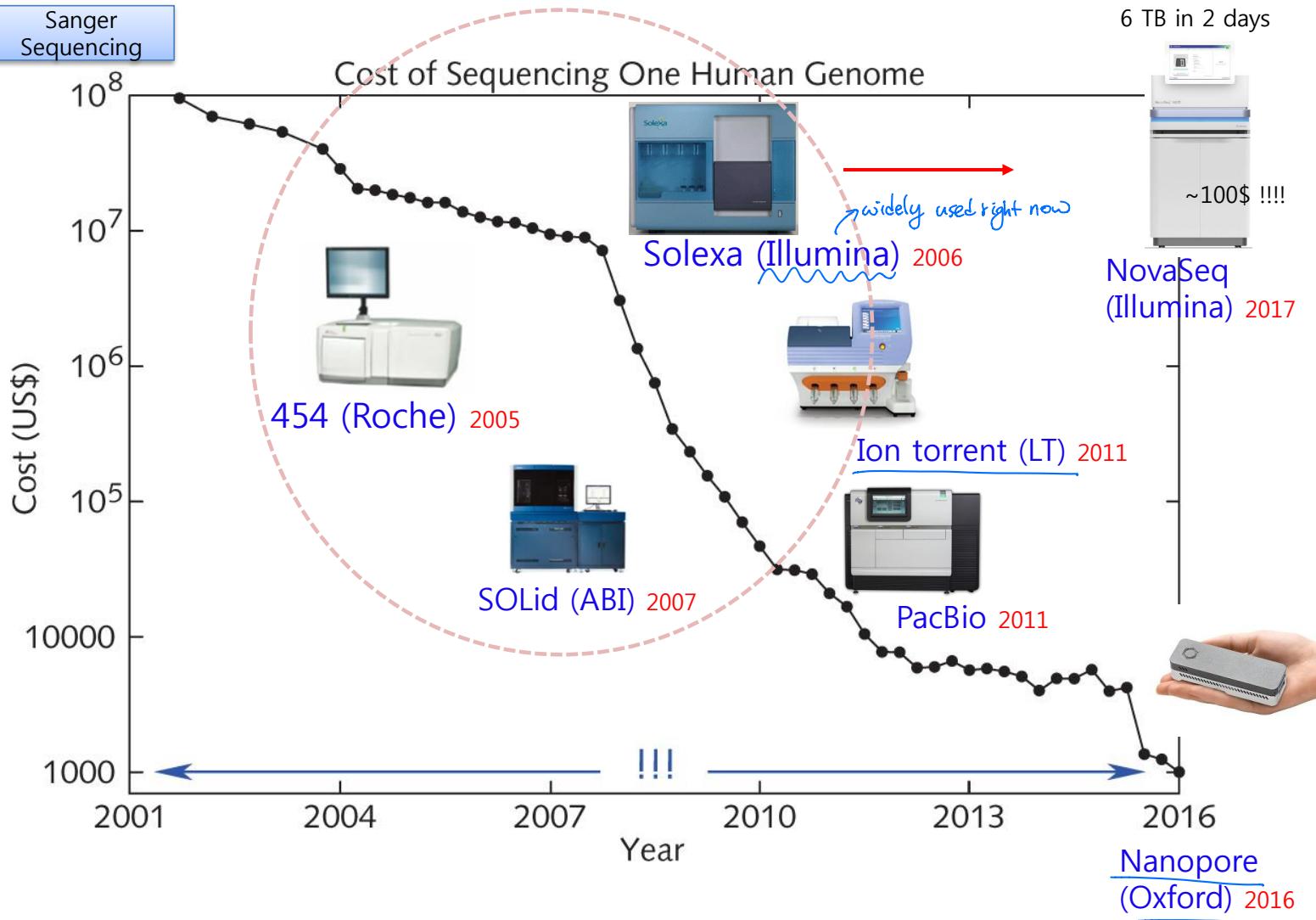
Summary & More

• Next-Generation Sequencing (NGS)

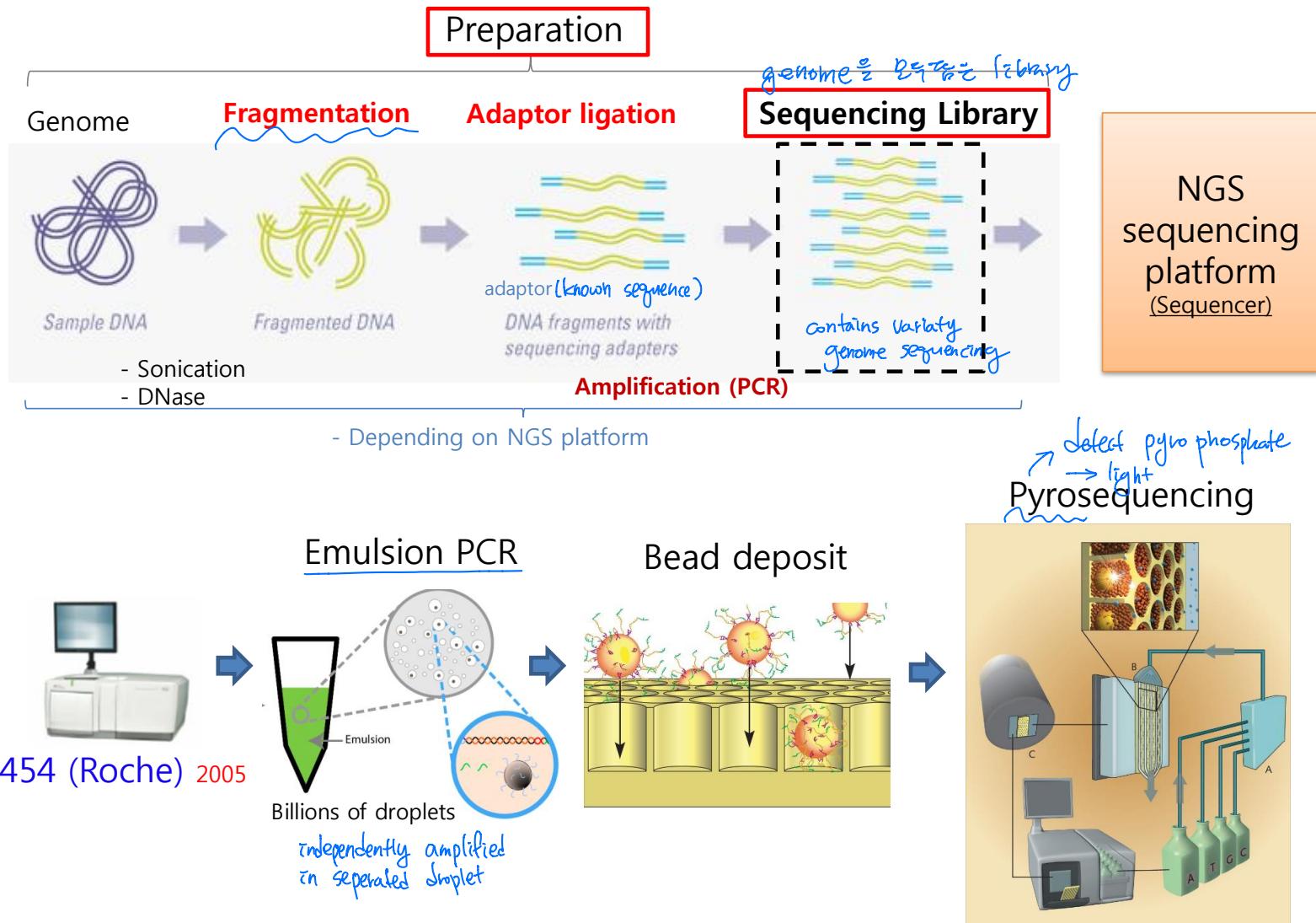
Sung Wook Chi

Division of Life Sciences, Korea University

Next-generation Sequencing (NGS)



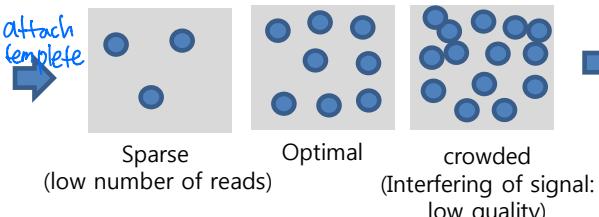
DNA Sequencing by NGS (flow & principle)



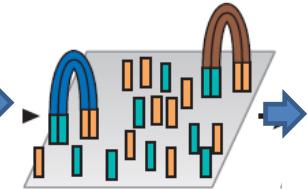
Illumina NGS platform (Solexa)



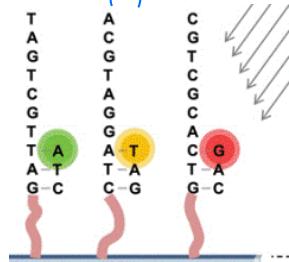
Library loading



Bridge amplification



(sequencing by synthesis)
utilize polymerase



Solexa (Illumina)
2006

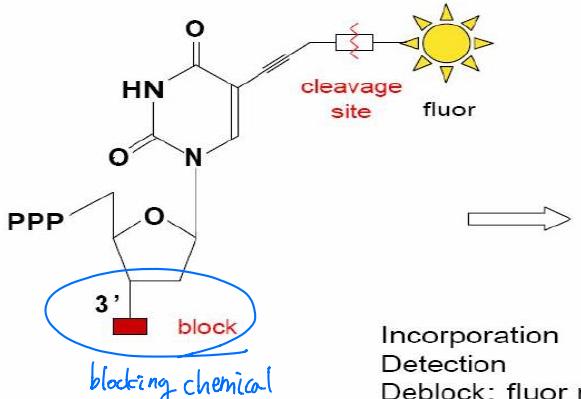
Loading (spread) your library
on glass plate (flow cell)

Solid phase amplification

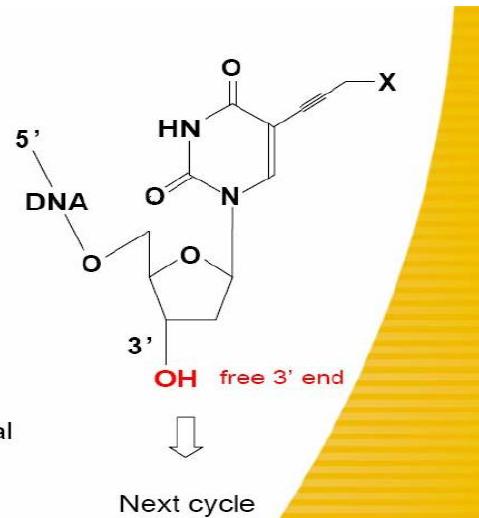
*(primer is already
attached on plate)*

Reversible termination
chemistry

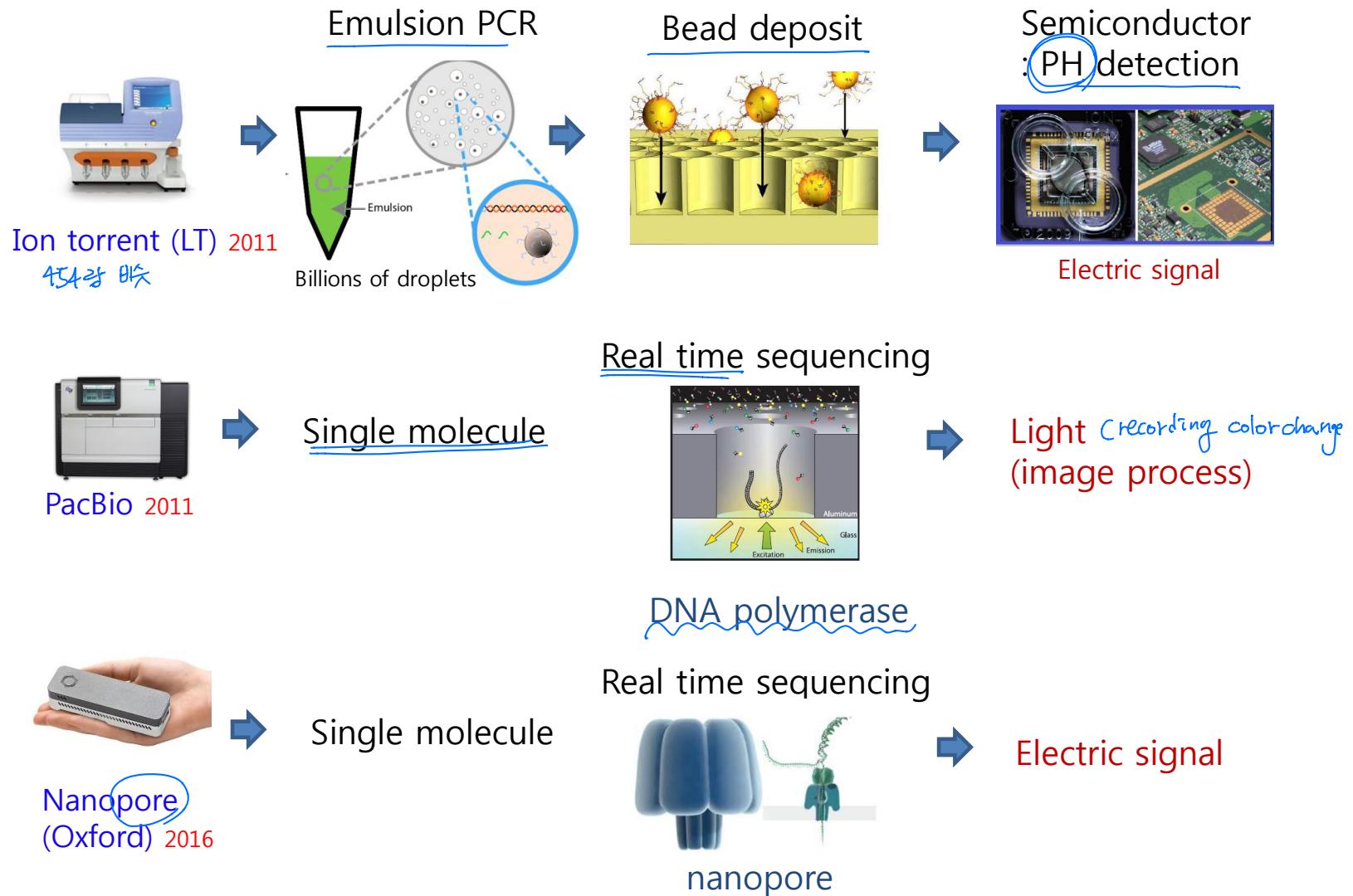
Reversible terminator chemistry



Incorporation
Detection
Deblock; fluor removal
Quenching

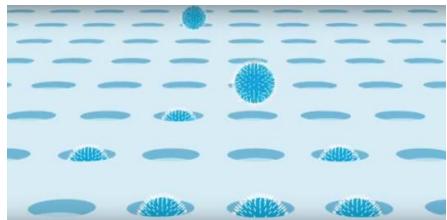


2nd Generation of NGS sequencer

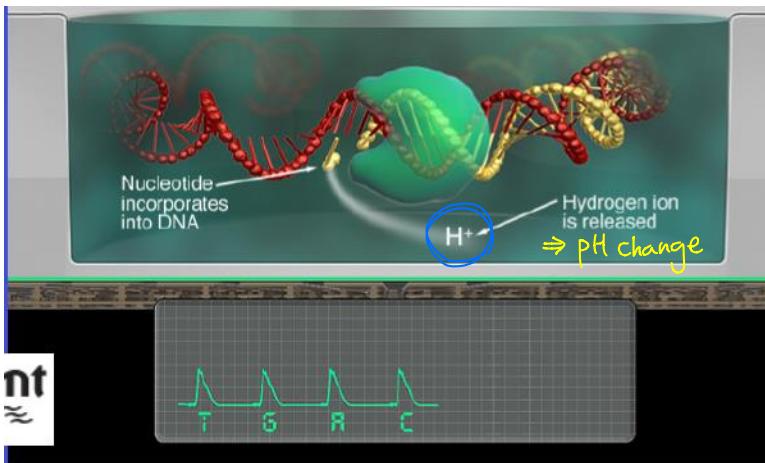


Ion torrent/proton: semiconductor sequencing

Emulsion PCR



**Electric signal
(No camera !!, low price)**

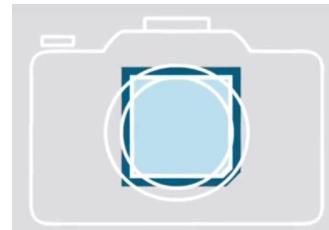


**Ion
torrent**

1.5 million pH meters semiconductor chip



**Ion
proton**



80 million pH meters semiconductor chip

→ camera → pixel로 defect → signal
하늘고마 같은 chip이지만,
pH 정보 defect 하도록 개발된 chip

monocolor

DNA polymerase reaction

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

454랑 같음 (A만 넣고 signal 나오는지 봐, ⇒ washing → T만 넣고 signal 나오는지 봐)
AA처럼 pH change signal이 2개!

장점: 많은정보를
빠르게수집가능

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

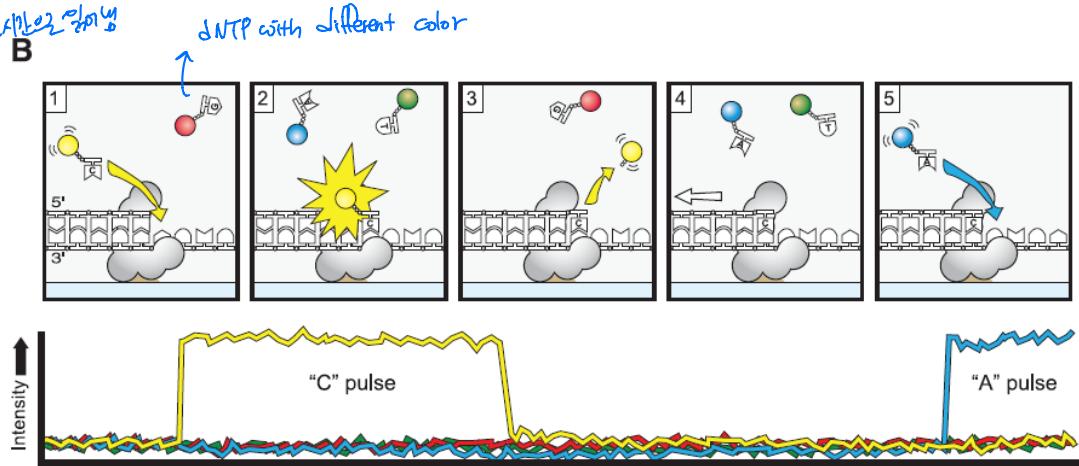
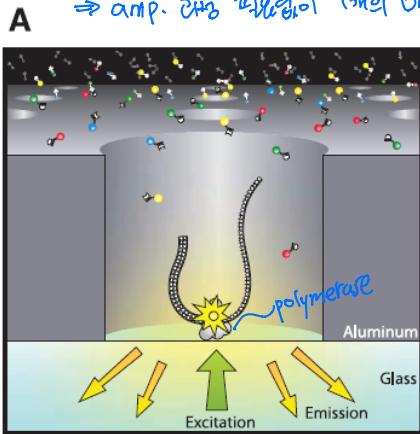
The same problems in
- Emulsion PCR, Pyrosequencing

Pacific Biosciences

No amplification (Single molecule)
Long reads (3kb ~ 15kb),
Fast (real time, 30 min)

작은 well 안에 pol. 불가 있음
⇒ amp. 그는 필요하지 않아서 DNA 한개로 충분

→ 높은 정확도로 defect 찾기 어렵지만
but low accuracy (95%)
low throughput (70,000 reads)



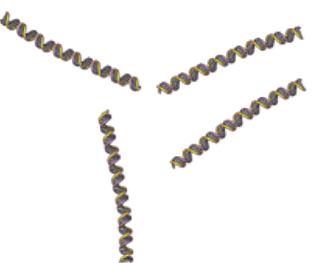
Single Molecule, Real Time
⇒ don't need amplification
Sequencing

as fast as polymerase read DNA

no limit

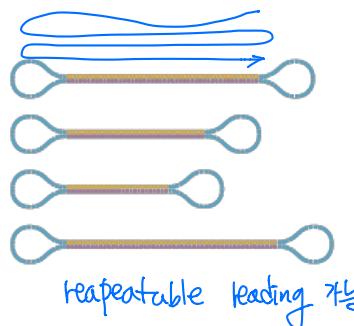
Useful for long read sequencing
(resequencing genomes)

SMRT (Single Molecule RealTime) Sequencing : PacBio



modified

Can read one molecule as many as we want



repeatable reading ↗

Library Preparation

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

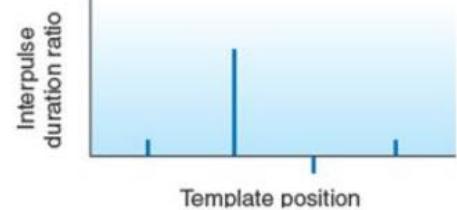
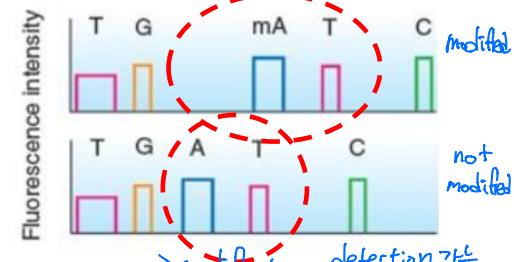
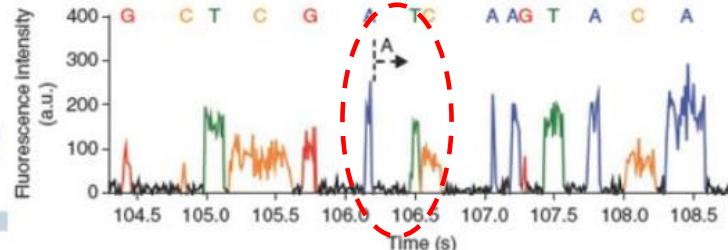
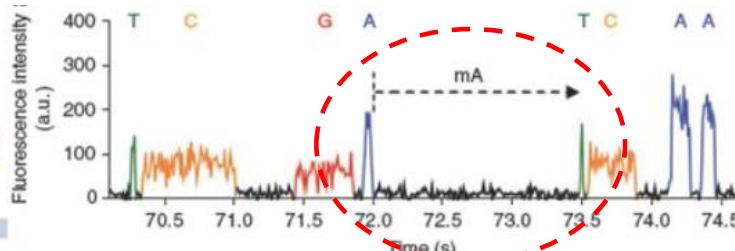
Redundant Sequencing Reading

long read ↗

Consensus accuracies > 99.999%



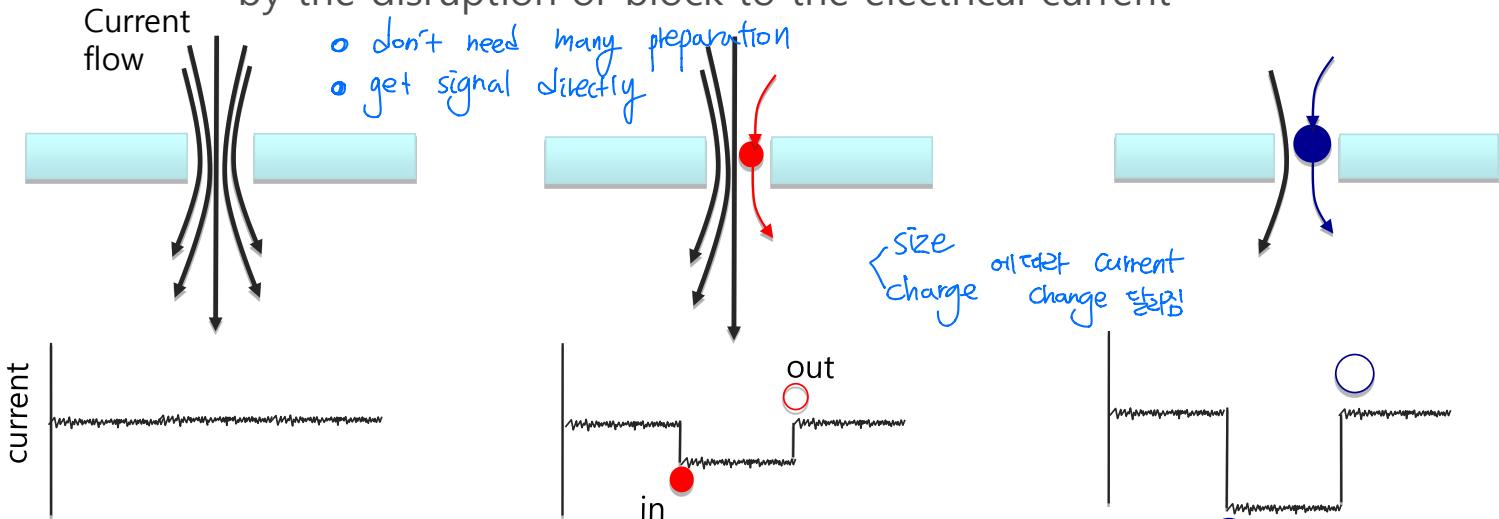
Realtime direct detection of modified DNA



Oxford Nanopore

Electronic &
Single molecule
Real time

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



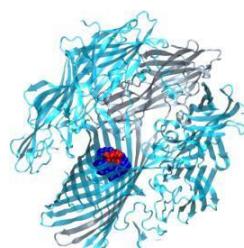
nature
nanotechnology

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

ARTICLES

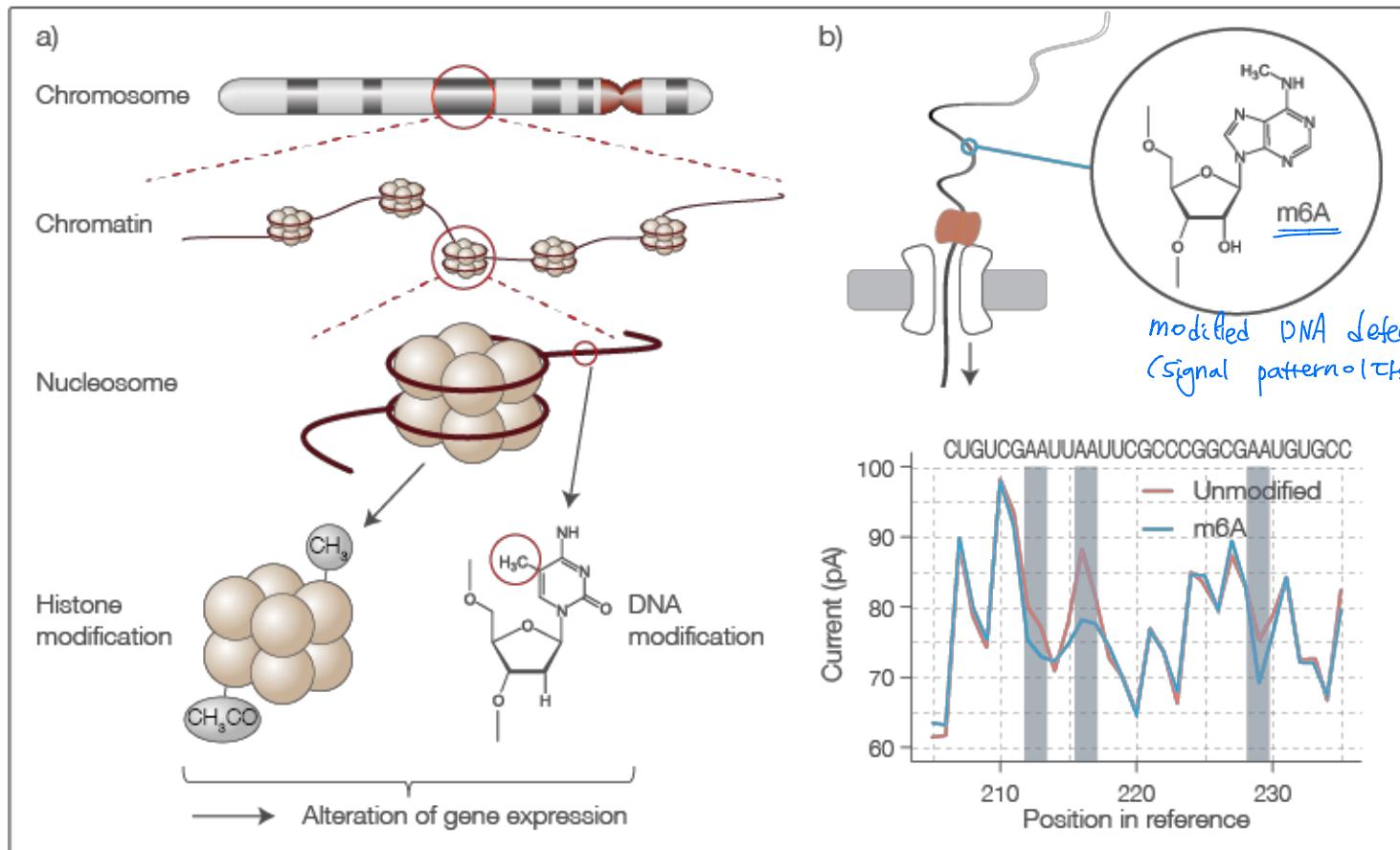
Continuous base identification for single-molecule nanopore DNA sequencing

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



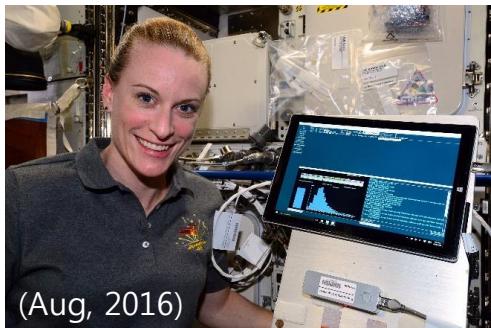
Realtime Portable Sequencer: Oxford Nanopore

Oxford Nanopore (MinION)



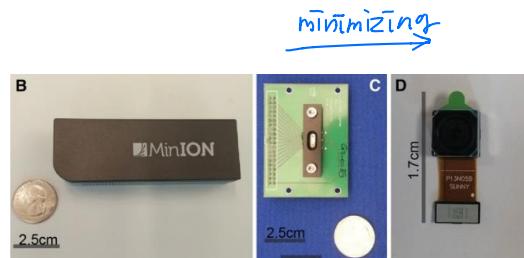
아시아의 등에서 사용

9th NASA/SpaceX commercial cargo



Space ship에 실내서 유전체를 사용

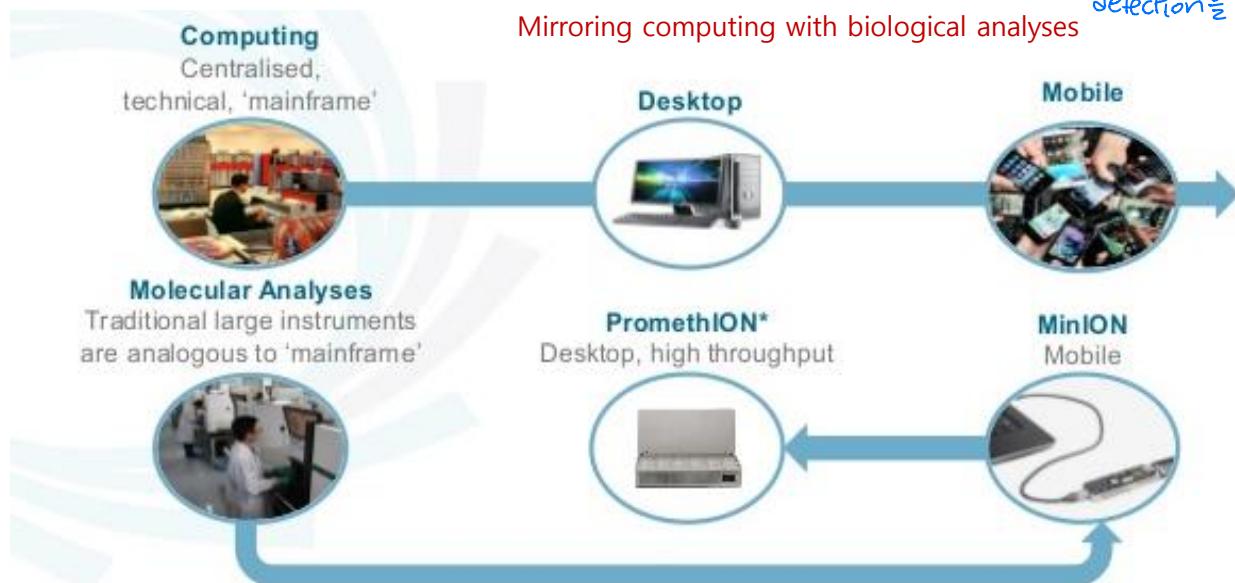
Creation of internet of living things (DNA)



IoT (Internet of Thing)

⇒ IoLT (Internet of Living Thing)

ex) micro - organism 을 통한 real time detection을 통한 생물학적 분석



NGS platform comparison

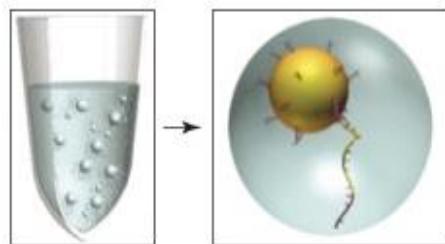


Emulsion PCR

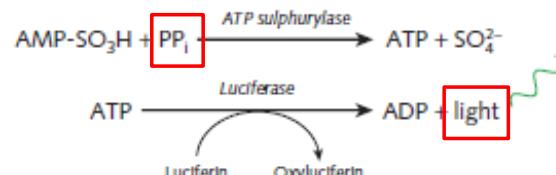
Platform	Amplification	Chemistry	Detection	description
454	emPCR	Pyrosequencing	Image (mono)	Low throughput, long read (~400), ~0.35 day
Abi SOLiD	emPCR	Seq by ligation	Image (color)	High throughput, short read (~50), ~2 weeks
Illumina /Solexa	Solid-phase	Reversible terminator	Image (color)	High throughput, 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)	?

no amp.

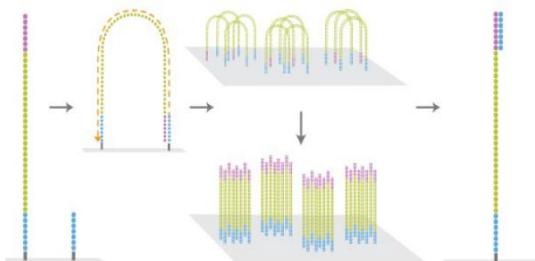
Emulsion PCR



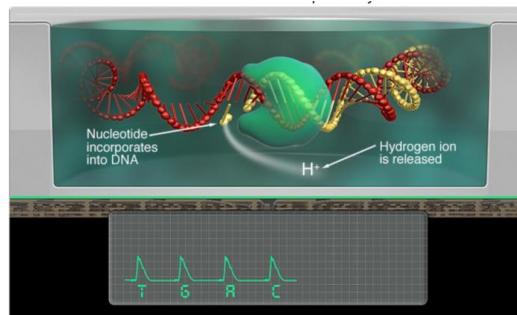
pyrosequencing



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 <i>low</i>	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 <i>pair end</i>	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 <i>high error rate</i>	$3.5\text{--}7.5 \times 10^4$	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.)

long read
 A T G C G I Signal이
 distinctive 특징 있는 짧은 읽음
 ⇒ high error rate

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)

Genome Annotation

- Sequence alignment

Sung Wook Chi

Division of Life Sciences, Korea University

Mapping reads (NGS) : Sequence Alignment

<mapping process>



↳ mapping our sequence into known sequence

fragment

fragment



Sequencing (NLG & Co)

alignment



ATGTTCCGATTAGGAAACCTATCTGTAACTGTTTCATTCACTAAAAGGGAGGAAA

Pairwise sequence alignment



Sequence Alignment : history

Amino acid sequence insulin (1951)

first sequenced amino acid



Fred Sanger
(1918-2013)

A chain	B chain
Gly	Phe 1
Ile	Val
Val	Asn
Glu	Gln
Gln	His 5
Cys	Leu
Cys	Cys
Ala	Gly
Ser	Ser
Val	His 10
Cys	Leu
Ser	Glu
Leu	Ala
Tyr	Leu 15
Gln	Tyr
Leu	Leu
Glu	Tyr
Asn	Leu
Tyr	Val
Cys	Cys 20
Gly	Gly
Glu	Arg
Arg	Gly
Gly	Phe
Phe	Phe 25
Tyr	Tyr
Thr	Thr
Pro	Pro
Lys	Lys
Ala	Ala 30

- DNA sequence (Sanger Sequencing)
(bacteriophage, 1977)



Fred Sanger



Protein sequence atlas (1960)

Mother of bioinformatics



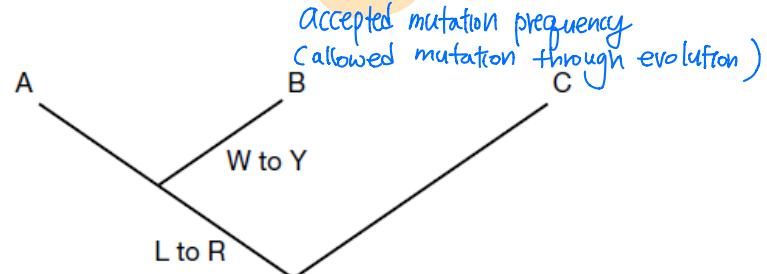
Margaret Dayhoff
(1925-1983)

Cytochrome protein (Families & Superfamilies)

↳ So many isoform & spread out through species
→ homolog 등을 보면서 유전적 차이를 찾는다
→ sequence 암호를 찾으면 같은 job
→ biology system의 암호의 mutation을 알수

A	A	W	T	V	A	S	A	V	R	T	S	I
B	A	Y	T	V	A	A	A	V	R	T	S	I
C	A	W	T	V	A	A	A	V	L	T	S	I

- Phylogenetic tree
:Evolutional and functional relationship
:Substitution matrix (PAM), one letter code



Sequence analysis > Sequence alignment

Biological question from sequence alignment

ACGCTGA

(Do same function)

Common
Ancestor

ACTGT

Changed through evolution process

Don't exactly know what happened

Evolutional Changes

alignment can tell you
about history

ACGCTGAA

2 insertion

1 substitution

ACGCTGA

A--CTGT

Sequence
alignment 1

ACGCTGA

2 substitution

2 insertion

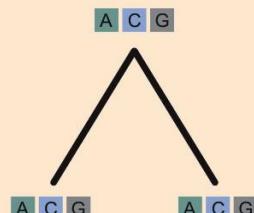
ACGCTGA

ACTGT--

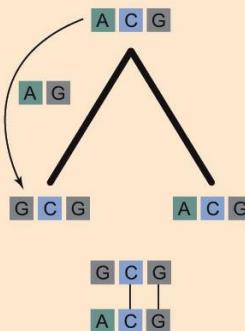
Sequence
alignment 2

ACTGT

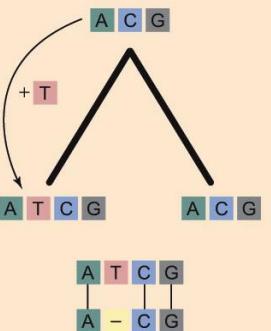
(A) Identity



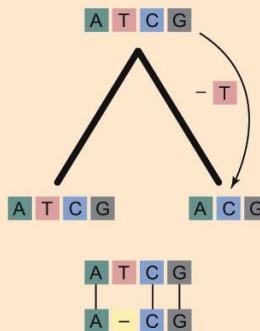
(B) Substitution



(C) Insertion



(D) Deletion



tell you most likely
change

What is the best
optimal alignment?

Which alignment is
biologically relevant?

: purpose of alignment

Two types of sequence alignment : local & global alignment

- **Sequence Alignment**

- Pairwise alignment
- Multiple Sequence alignment

Global Alignment

Same start/last position

- Suite for similar sequences
- Nearly equal length
- Overall similarity is detected

Local Alignment

regardless of start / last position
→ local part of very similar seq align

- Isolate regions in sequences
- Suitable for database searching
- Easy to detect repeats

High sequence similarity

> Homolog >
same function

Sequence 1
Sequence 2

Alignment 1

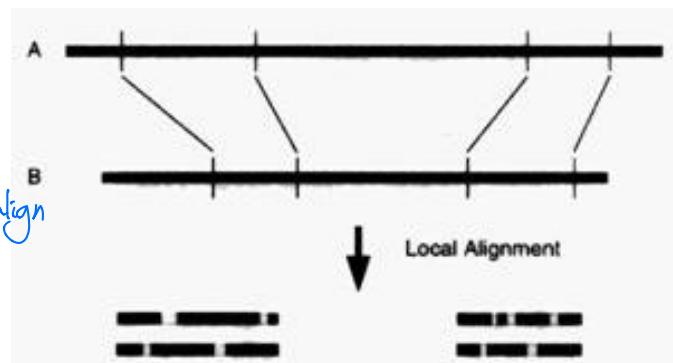
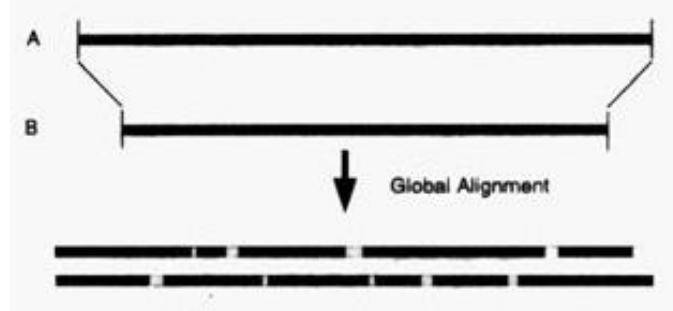
ACGCTGA
||
A--CTGT

Alignment 2

ACGCTGA
ACTGT--

Global alignment

Local alignment



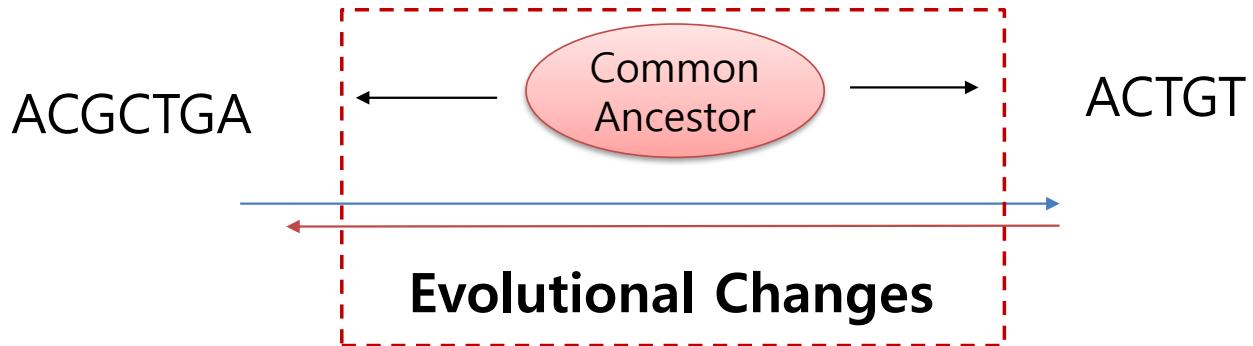
only some part

ex) kinase → 같은 업무를 같은 일자에
→ same sequence but different (local on high similarity.)

Conserved region
of sequence

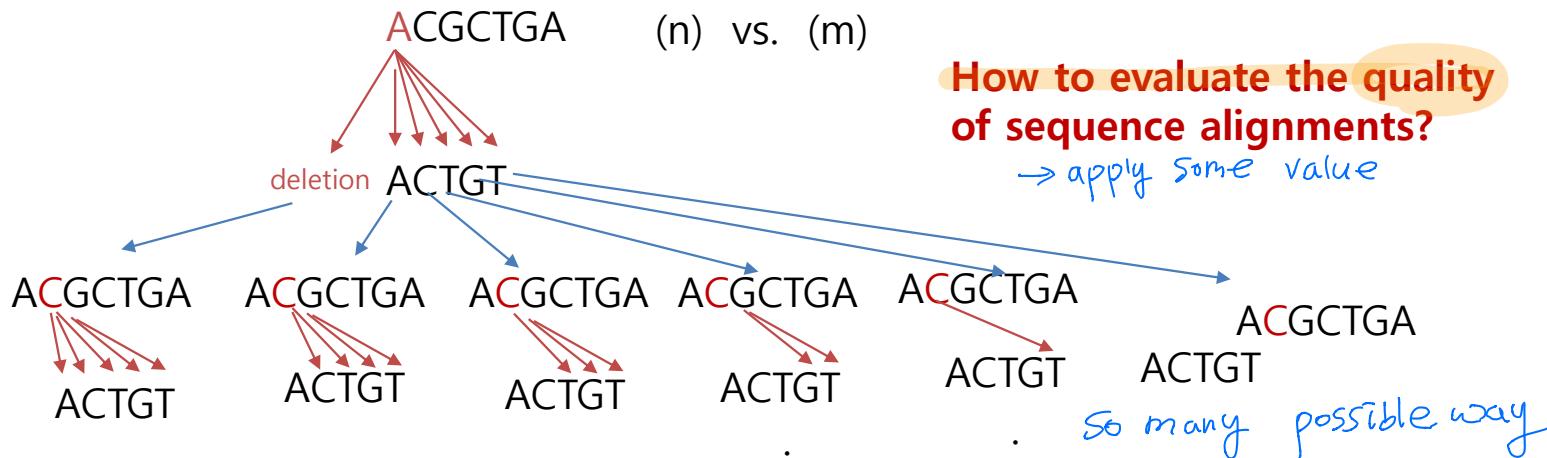
> Similar function

Sequence Alignment : Bioinformatics methods



All combination of sequence alignments: Which one is the optimal ?

all possible combination 2^{n+m}



Need to have similarity or distance value for evaluating sequence alignments.

How to measure sequence similarity : distance

how different these two seq

- (Method 1) Counting identical letters on each position

→ **Hamming distance**

: Number of position with mismatch characters

distance ↗ good
scale ↳ good

- (Method 2) Inserting gaps to maximize the number of identical letters

→ **Edit distance** (Levenshtein distance)

: How many operations are required?

A T C C G A T
| | | |
T G C - A T A T

A T C C G A T
| | | |
T G C A T A T

- x="TGCATAT" (m=7), y="ATCCGAT" (n=7)

TGCATAT insertion of "A" ①
ATGCATAT substitution of "G" with "C" ②
ATCCATAT insertion of "G" ③
ATCCGATAT deletion of "A" ④
ATCCGATI deletion of "T" ⑤

- x="TGCATAT" (m=7), y="ATCCGAT" (n=7)

TGCATAT insertion of "A" ①
ATGCATAT deletion of "T" ②
ATGCAAAT substitute of "G" with "C" ③
ATCCAAT substitute of "A" with "G" ④

✓ Hamming distance = 3 better
Edit distance = 5

✓ Hamming distance = 4 better
Edit distance = 4

How to measure sequence similarity : Score

give score minus score
Score = (match or mismatch penalty) – gap penalty

Match = 3 , Mismatch = -1

C	S	T	P	A	G	N	D	E	Q	H
C	9									
S	-1	4								
T	-1	1	5							
P	-3	-1	-1	7						
A	0	1	0	-1	4					
G	-3	0	-2	-2	0	6				
N	-3	1	0	-2	-2	0	6			
D	-3	0	-1	-1	-2	-1	1	6		
E	-4	0	-1	-1	-1	-2	0	2	5	
Q	-3	0	-1	-1	-1	-2	0	0	2	5
H	-3	-1	-2	-2	-2	-2	1	-1	0	0
R	-3	-1	-1	-2	-1	-2	0	-2	0	1
K	-3	0	-1	-1	-1	-2	0	-1	1	1
M	-1	-1	-1	-2	-1	-3	-2	-3	-2	0
I	-1	-2	-1	-3	-1	-4	-3	-3	-3	-3
L	-1	-2	-1	-3	-1	-4	-3	-2	-2	-3
V	-1	-2	0	-2	0	-3	-3	-2	-3	-2
F	-2	-2	-2	-4	-2	-3	-3	-3	-3	-3
Y	-2	-2	-2	-3	-2	-3	-2	-2	-2	-1
W	-2	-3	-2	-4	-3	-2	-4	-4	-3	-2
C	S	T	P	A	G	N	D	E	Q	H

	A	G	T	C
A	20	10	5	5
G	10	20	5	5
T	5	5	20	10
C	5	5	10	20

G→A : 푸른→푸른 Gap = -2
higher score than 푸른→다른색인

Affine gap penalty

(gap opening + (gap extension) x (n-1))

ACGCTGA
A--CTGT



Gap opening, Gap extension penalty

ORGANISM A	A	W	T	V	A	S	A	V	R	T	S	I
ORGANISM B	A	Y	T	V	A	A	A	V	R	T	S	I
ORGANISM C	A	W	T	V	A	A	A	V	R	T	S	I

↳ estimated by frequency ⇒ more frequently observed
Substitution matrix in same family ⇒ higher score

Log-odds score

rare observed "⇒ high penalty"

$$S_{i,j} = \log \frac{p_i \cdot M_{i,j}}{p_i \cdot p_j} = \log \frac{M_{i,j}}{p_j} = \log \frac{\text{observed frequency}}{\text{expected frequency}}$$

-BLOSUM62: Pairs frequencies were counted between segments less than 62% identical.

The PAM1 matrix estimates what rate of substitution would be expected if 1% of the amino acids had changed

Sequence Alignment

Pairwise sequence alignment

	match	Alignment 1	Alignment 2
Sequence 1	ACGCTGA	ACGCTGA	
Sequence 2	A--CTGT		ACTGT--
↑↑ mismatch			
Gap opening, Gap extension			

Optimal sequence alignment

1. Evaluation of sequence similarity : Distance, Score

depending on scale, program can decide which alignment is better

2. Performing optimal sequence alignment search: Dynamic Programming

Global Alignment

T C A G - - T - G T C G A A G T - T A
| | | | | | | | | | | | | | | | | | | |
T - A G G C T A G - C - A - G T G T A

alignment → global local
program →

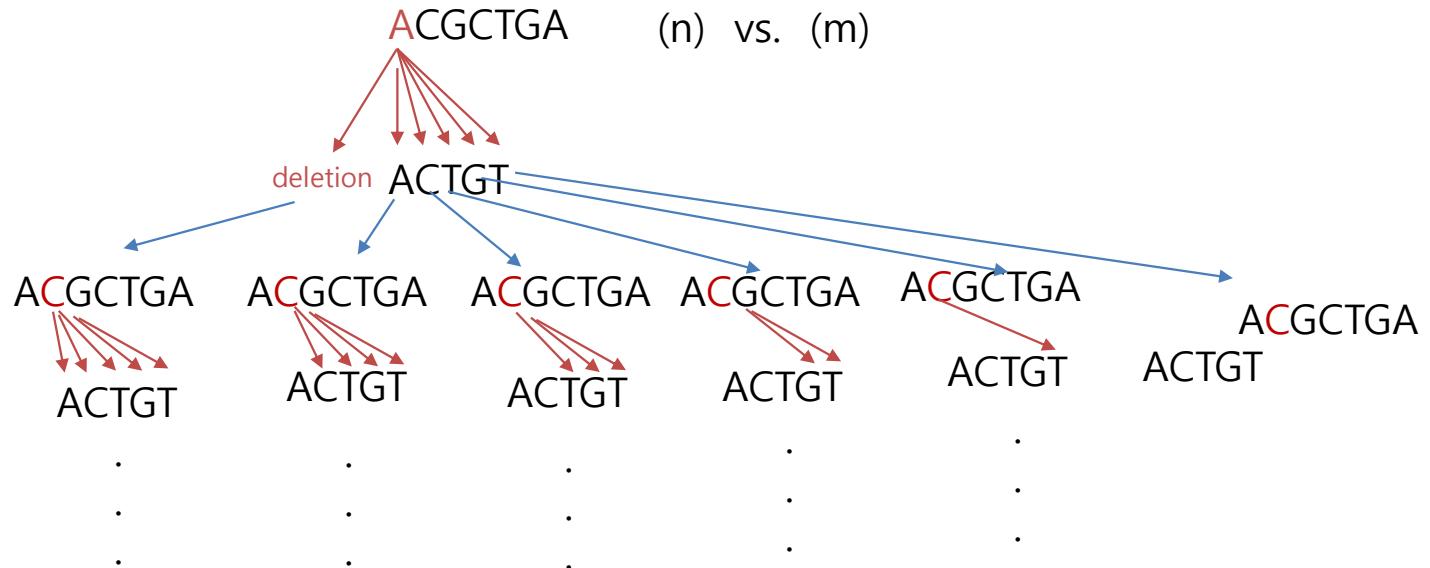
Local Alignment

T C A G T G T C G A A G T T A
| | | | | |
T A G G C T A G C A G T G T A

Sequence Alignment : Find the optimal alignment

All combination of sequence alignments:

Too many operations, and storage !!!

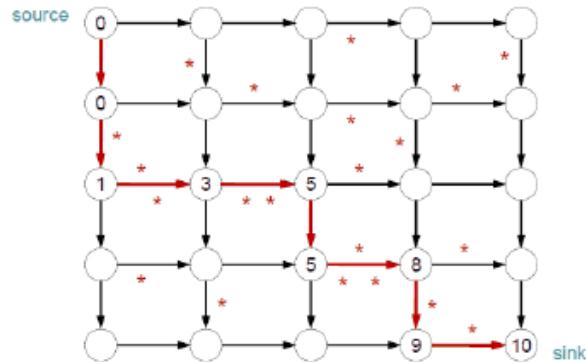
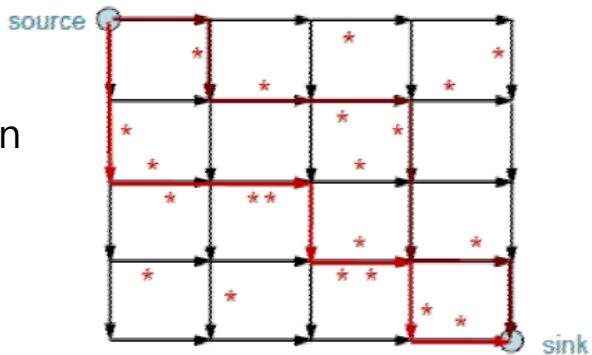


Solution :

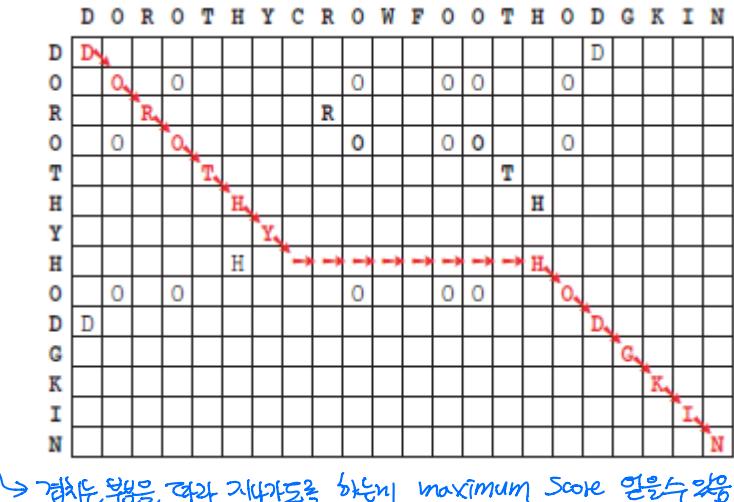
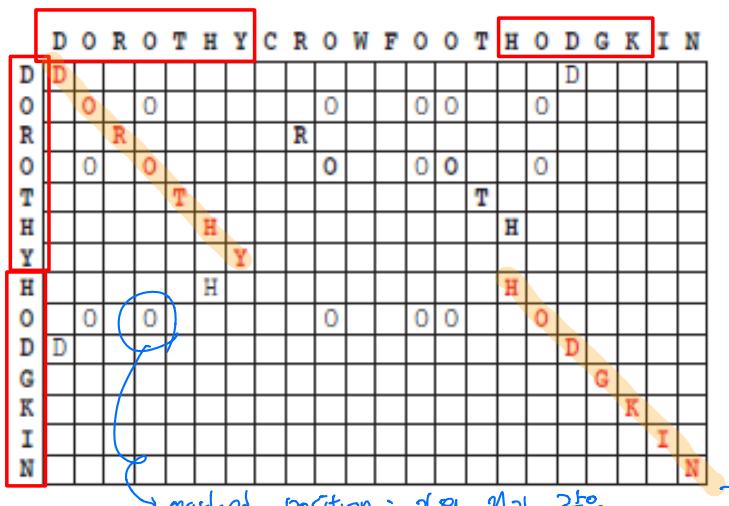
The longest Common Path Problem
Dynamic Programming

Finding optimal alignment: finding shortest distance

Manhattan
Tourist
Problem



DOROTHY-----HODGKIN
DOROTHYCROWFOOTHODGKIN



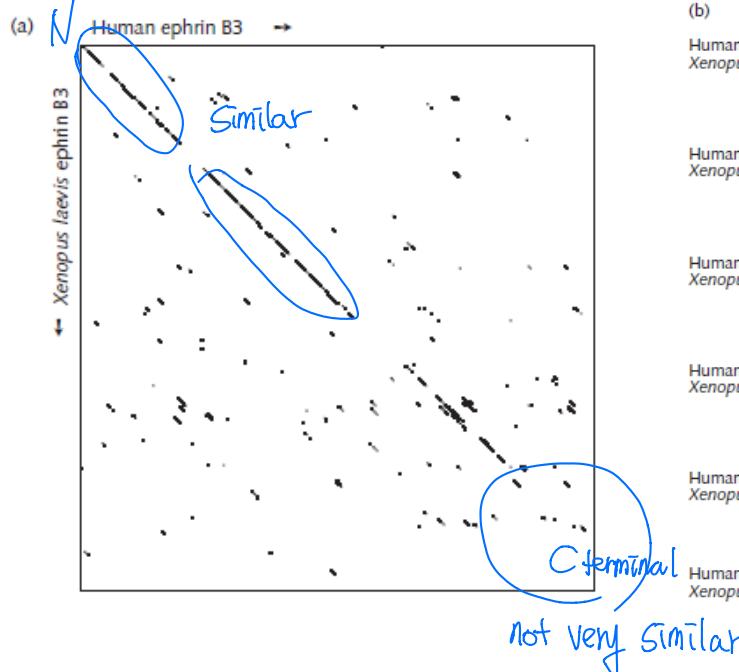
- A dotplot shows perspicuously the quality and distribution of the pattern of similarity between two sequences. Each possible alignment of the two sequences corresponds to a path through the dotplot, from upper left to lower right.

Sequence Alignment : Dot plot

zipped map

& match → mark as dot

Relationships between the sequences of ephrin B3 proteins from human and *Xenopus laevis*.



(b)

	10	20	30	40	50	60
Human	I S L E P V Y W N S A N K R F Q A E G C Y V L Y P Q I G D R L D L L C P F A R P P G P H S S P N Y E F Y K L Y L V G G A					
Xenopus	S E D P I Y W N S S N K R F D E T E G V Y L Y P Q I G D R L D L L C P F S E P Q G P F S S S P Y E Y Y K L Y L V G G T K	S L P Y W N S N K R F G Y V L Y P Q I G D R L D L L C P R P G P S S Y E Y K L Y L V G				
	70	80	90	100	110	120
Human	Q C H R - C E A P P A P N N L L T C D R F D L D L R F T I K F Q E Y S P N L W G H E F R S H H D Y Y I I A T S D G T R E					
Xenopus	E E M S S C S I L R T P N L L T C D R F S Q D L R F T I K F Q E Y S P N L W G H E F R S H H D Y Y I I A T S D G T R E	C P N L L T C D R P D L R F T I K F Q E S P N L W G H E F R S H H D Y Y I I A T S D G T R E				
	130	140	150	160	170	180
Human	G L E S L Q G G V C L T R G M K V L L R V G Q S P R G G A V P R K P V S E N P M E R D R G A A H S L E P C H E N L P G D					
Xenopus	G I E T L Q G G V C E T K G M K V T L K V G Q S P N G A T P P R R P S S A G -- K D S G I S P S V P N P D I P N V G D	G E L Q G G V C T G M K V L V G Q S P N G A T P P R R P S S A G D G S S D Y Y I I A T S D G T R E				
	190	200	210	220	230	240
Human	P T S N A T S R G A E G P L P P F S H P A V A G A A G G L A L L L L G V A G A G G G A M C W R R R R A K H S E S P R P G P					
Xenopus	T S G N A T K T G E N G P L P I S H V P L V A G A A G G A L L L L V F G V V G W V C H R R R Q A K H S D T R P P P	N A T G G P L P P P V A G A A G G G A L L L L V F G V V G W V C H R R R Q A K H S D T R P P P				
	250	260	270	280	290	300
Human	G S F G - - - R G G S L S E G S I T S P K R G G N N N G H E F S D I I M P L R P S E A G A F C H E A E P G E L G - - - I A L K G G G A A D P					
Xenopus	L S E G S I T S P K R G G N N N G H E F S D I I M P L R P S E A G A F C H E Y E K V S G D Y G H P V Y I Y Q D M A S Q S	S G R G G L G G G G G M C P R E A E P G E L G - - - I A L K G G G A A D P A				
	P F C P H Y E -					
	P A N I Y Y K V	P Y				

⇒ easy to visualize similarity

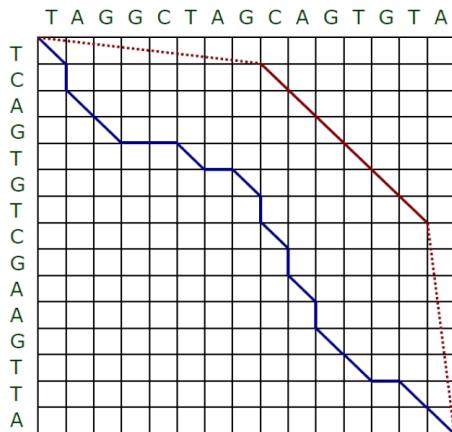
Dynamic Programming

- Dot Plot



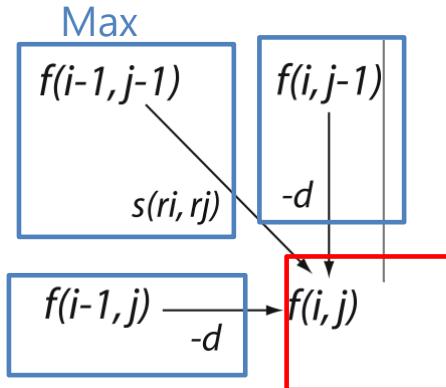
trace path to get high score **-Dynamic Programming:**

-Dynamic Programming:



/ Divide the problem into subproblems.

Previous state > Max Score > **current state (Max Score)**



➤ Global Alignment

Needleman-Wunsch

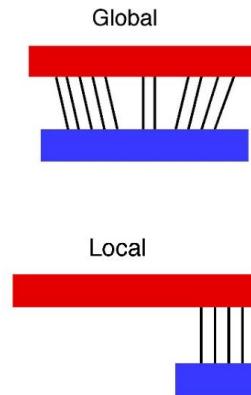
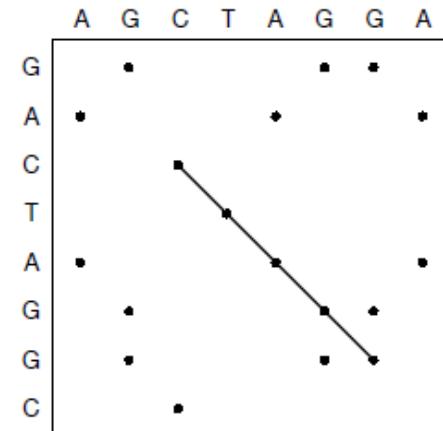
T C A G - - T - G T C G A A G T - T A
| | | | | | | | | |
T - A G G C T A G - C - A - G T G T A

➤ Local Alignment

TCAGTGTCGAAGTTA
|||
TAGGCTAGCAGTGTAA

Smith-Waterman

Sequence alignment and database search



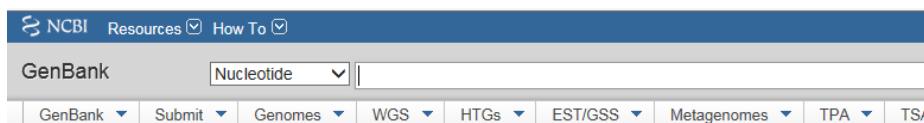
Needleman-Wunsch algorithm(1970)

Global FTFTALILLAVAV
F--TAL-LLA-AV

Smith-Waterman algorithm (1981)

Local FTFTALILL-AVAV
--FTAL-LLAAV--

GenBank (1982-) → first database collecting sequence data



Temple Smith

Michael Waterman



Walter Goad



David Lipman



BLAST (1981)
Database search

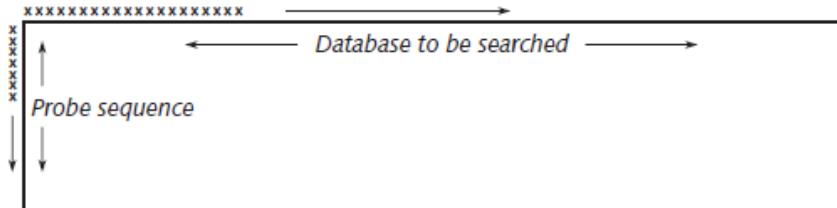
Alignment of sequence
to sequences in
GenBank database

Sequence Database Search: BLAST

like finding
text from textbook
with index

BLAST: Basic Local Alignment Search Tool,

(1) Empty dot plot

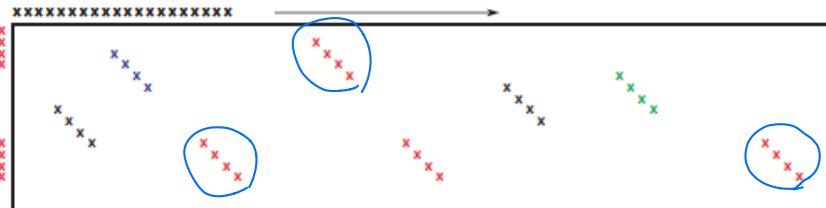


book: NCBI genebank

1. List

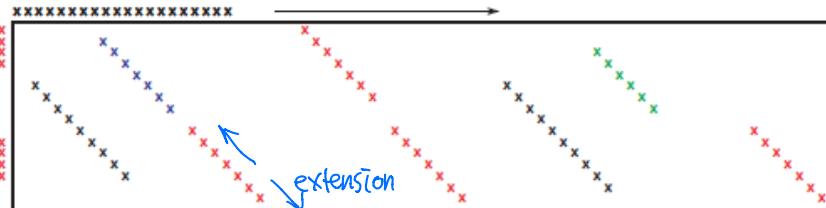
UH sequence
3 n t breakdown
⇒ match

(2) Word lookup

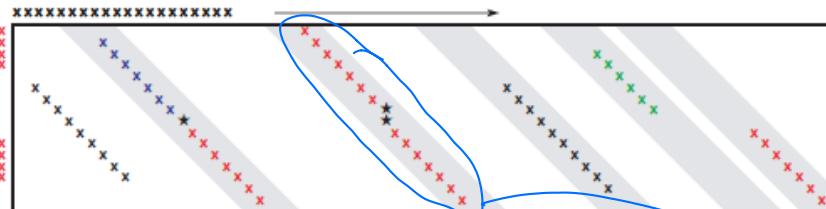


2. Scan

(3) Match extension



(4) Local gapped alignment



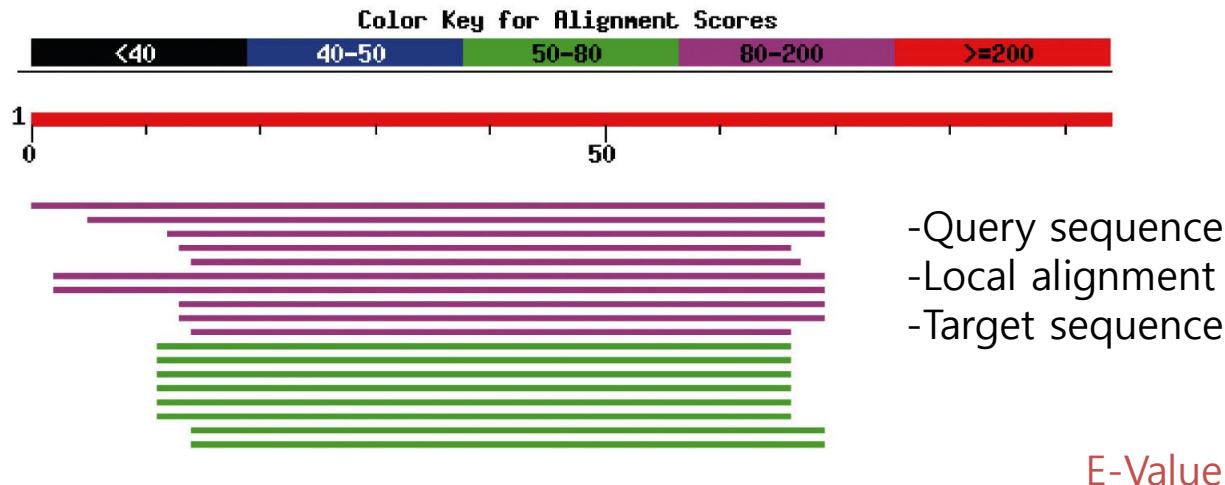
이정표에 있는 score가 진짜 significant한지
∴ random에 의한건 아님

4. Statistical significance

UH sequence와 match 할 확률
seq sequence

if mismatch가 계속
发生在 score drops
⇒ Stop extension

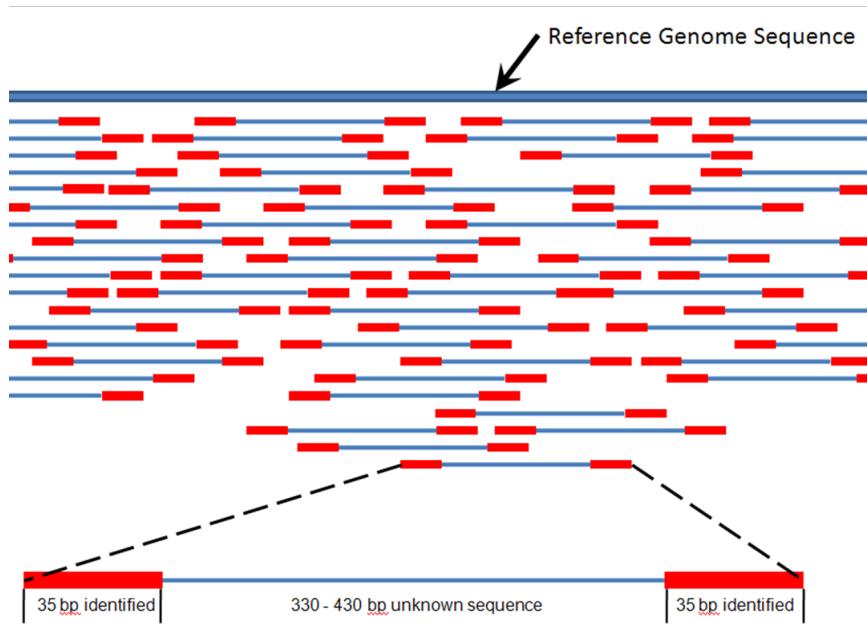
BLAST



Sequences producing significant alignments:

	Score (bits)	E Value
gi 256517 gb S45649.1 S45649 16S rRNA [Mastotermes electrod...	139	5e-33
gi 12005612 gb AF246514.1 AF246514 Drosophila ornatipennis ...	86	7e-17
gi 11119031 gb AF304735.1 AF304735 Sphyracephala bipunctipe...	84	3e-16
gi 3552018 gb AF086859.1 AF086859 Mystacinobia zealandica l...	84	3e-16
gi 256518 gb S45650.1 S45650 16S rRNA [Mastotermes darwinie...	84	3e-16
gi 15341487 gb AF403473.1 AF403473 This canus 16S ribosomal...	82	1e-15
gi 15341483 gb AF403469.1 AF403469 Icaridion debile 16S rib...	82	1e-15
gi 13435200 ref NC_002697.1 Chrysomya chloropyga mitochond...	82	1e-15
gi 13384216 gb AF352790.1 AF352790 Chrysomya chloropyga mit...	82	1e-15
gi 3552016 gb AF086857.1 AF086857 Calliphora quadrimaculata...	82	1e-15
gi 15341485 gb AF403471.1 AF403471 Malacomysia sciomyzina 16...	80	4e-15
gi 15341463 gb AF403449.1 AF403449 Helcomyza mirabilis 16S ...	80	4e-15

Mapping millions of short reads (NGS) on genome



reduce size of file

BWT (Burrows-Wheeler transform)
algorithm

