

Summary & Review

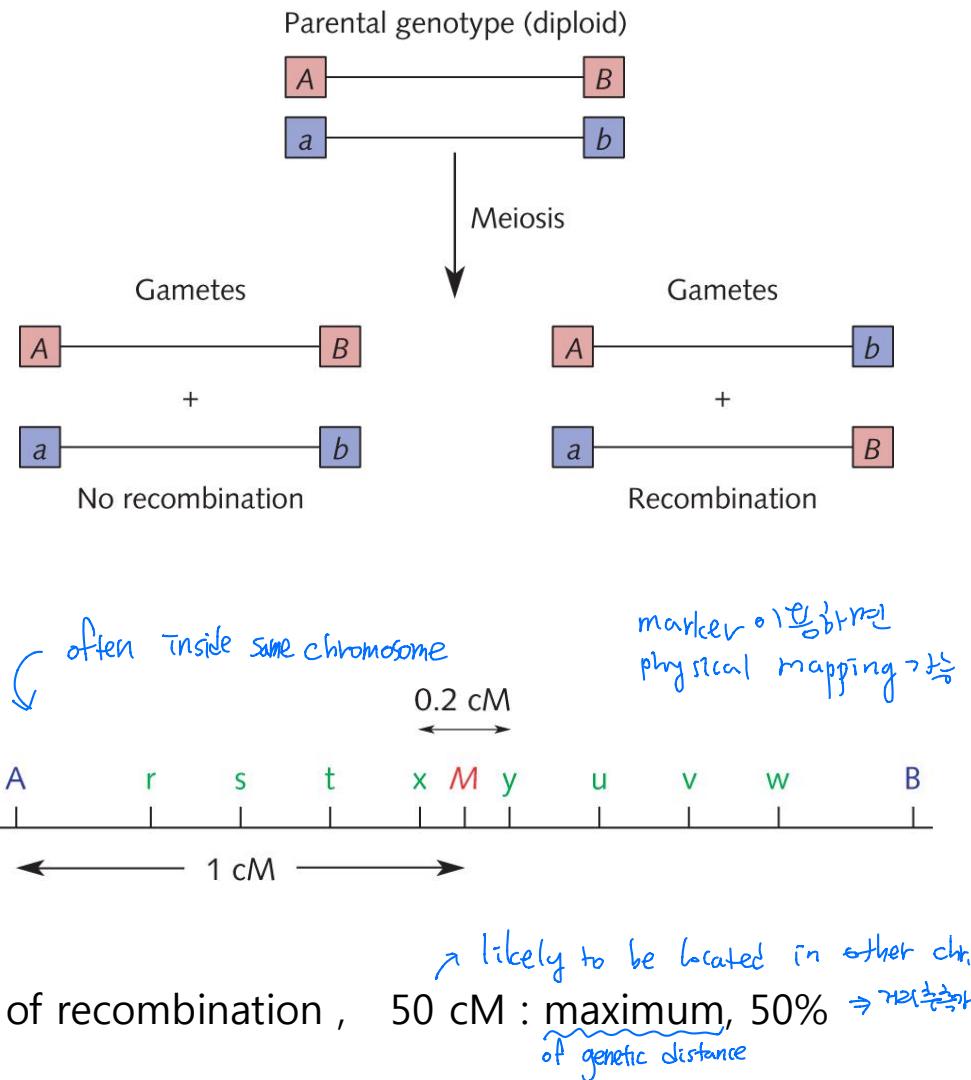
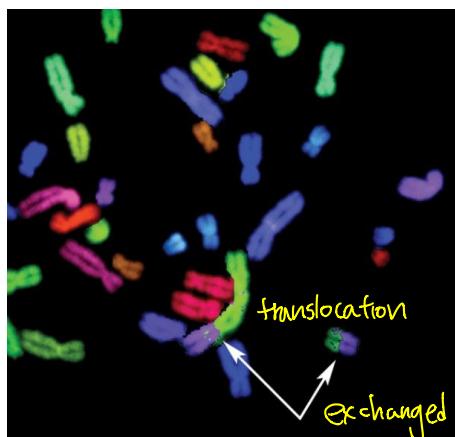
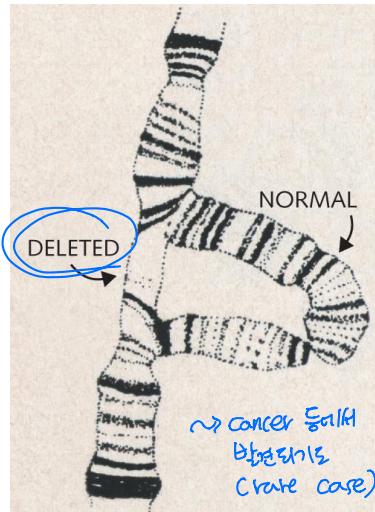
Sung Wook Chi

Division of Life Sciences, Korea University

What we learned in the previous lecture

- **Genomics**
 - The Core Aims of Genome Science
- **Mapping Genomes** ↗ sequencing 遺伝子塩基配列
 - Genetic Maps, Physical Maps, Cytological Maps
size & fragment
 - Comparative Genomics
- **The Human Genome Project**
 - Objectives, Internet Resources

Genetic Map: recombination

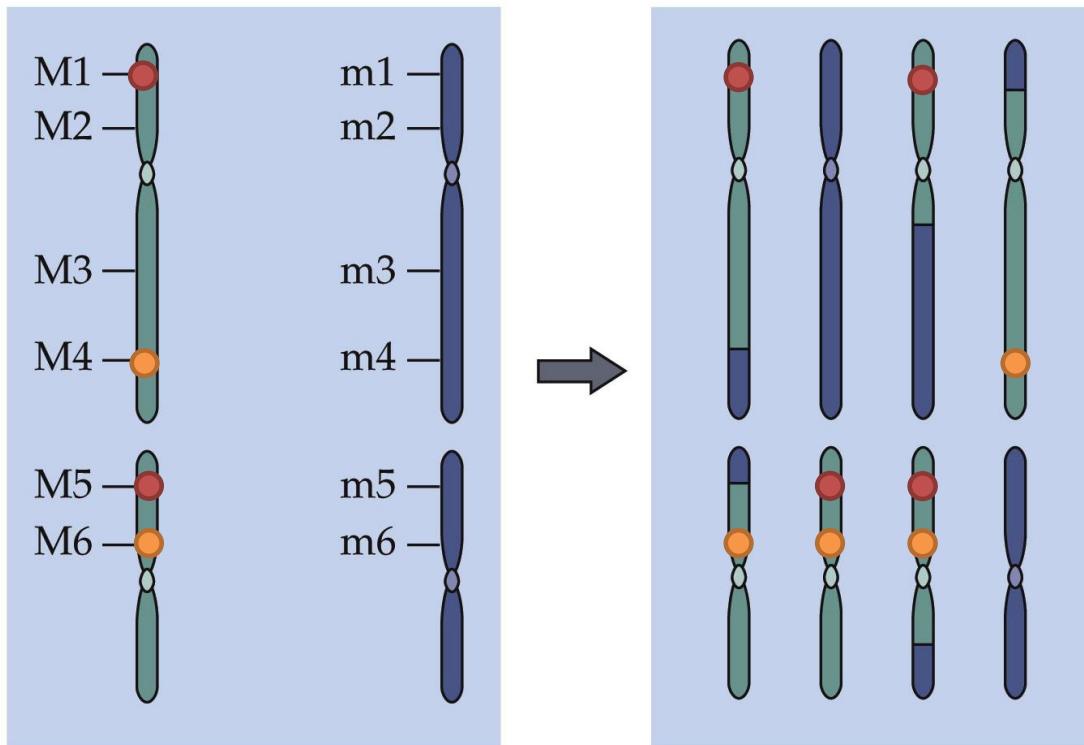


Mapping Genomes : Genetic Maps

- **Genetic Maps**

: relative order of genetics markers in linkage groups in which the distance between markers is expressed as units of recombination

↳ should be located in same chromosome



Independent segregation

Located in different chromosome

: 50% chance of co-segregation

기계마다 유전자는 함께 묶어
chunk together & phenotype linked
⇒ linked.

Genetic Maps

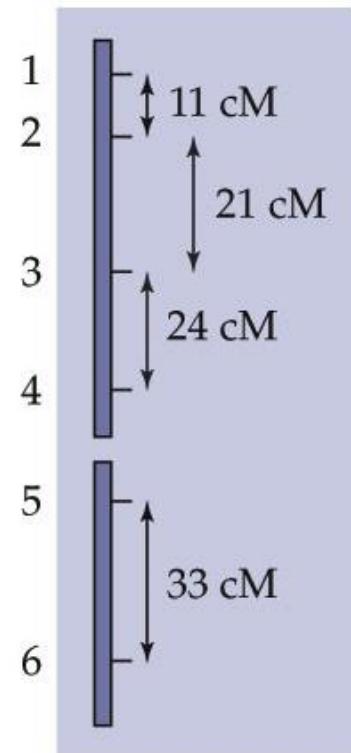
how we can construct genetic map

Linkage data -> Mapping function

(B)

	M1	M2	M3	M4	M5	M6
m1	-	.11	.31	.51	.49	.53
m2		-	.22	.46	.52	.48
m3			-	.25	.51	.50
m4				-	.49	.52
m5					-	.33
m6						-

(C)

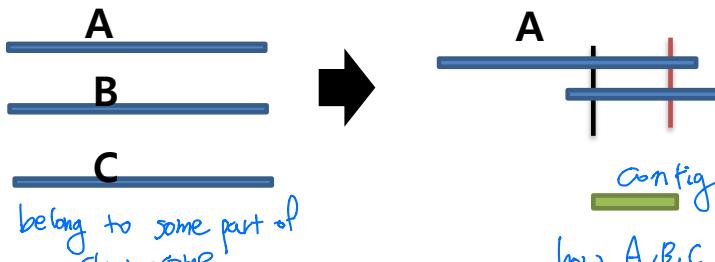


50cM – different chromosome

Mapping Genomes : Physical Maps

• Physical Maps

: an assembly of contiguous stretches of chromosomal DNA - contigs in which the distance between landmark sequences of DNA is expressed in kilobases.



• Assemble contigs

1. Alignment of randomly isolated clones based on shared restriction fragment length profiles

2. Hybridization-based approaches

: chromosome walking

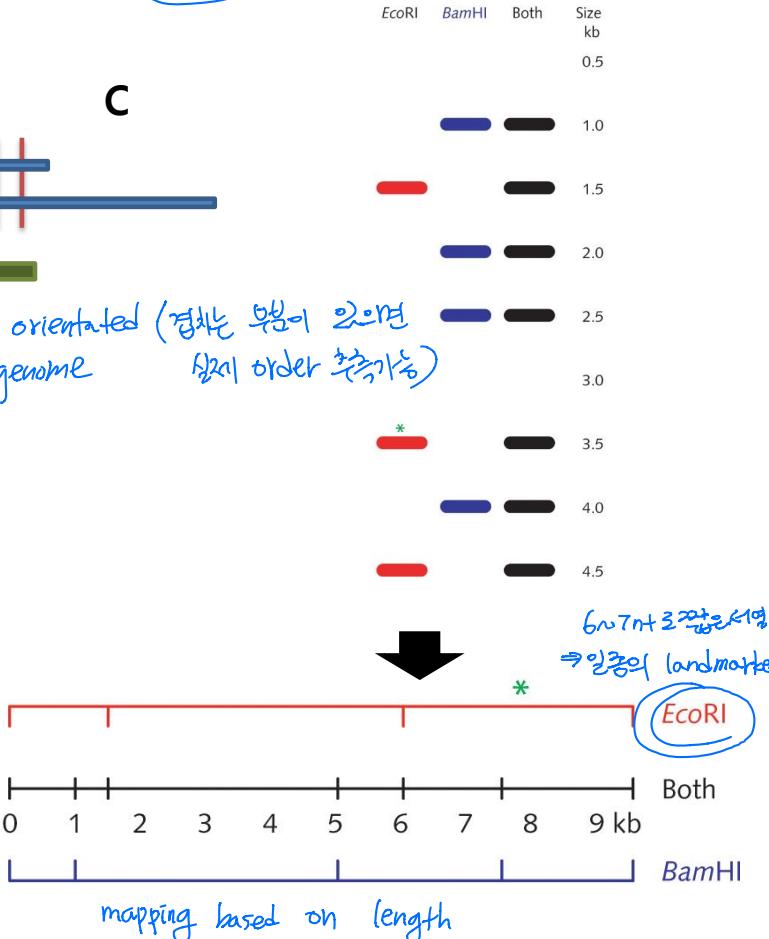
: **sequence-tagged sites (STSs)**:

single occurrence in the genome

특정 chromosome에서만 발견되는 sequence가 있으면
act as landmark → 위치 찾기 수단

Marker 간의 실제 거리를 infer하는 방법
fragment analysis \Rightarrow arrange and find relative order

contigs

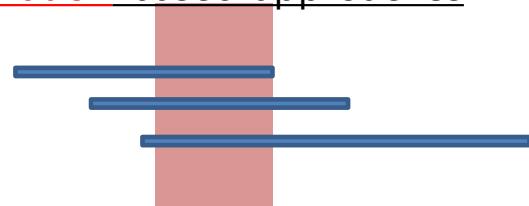


Mapping Genomes : Physical Maps

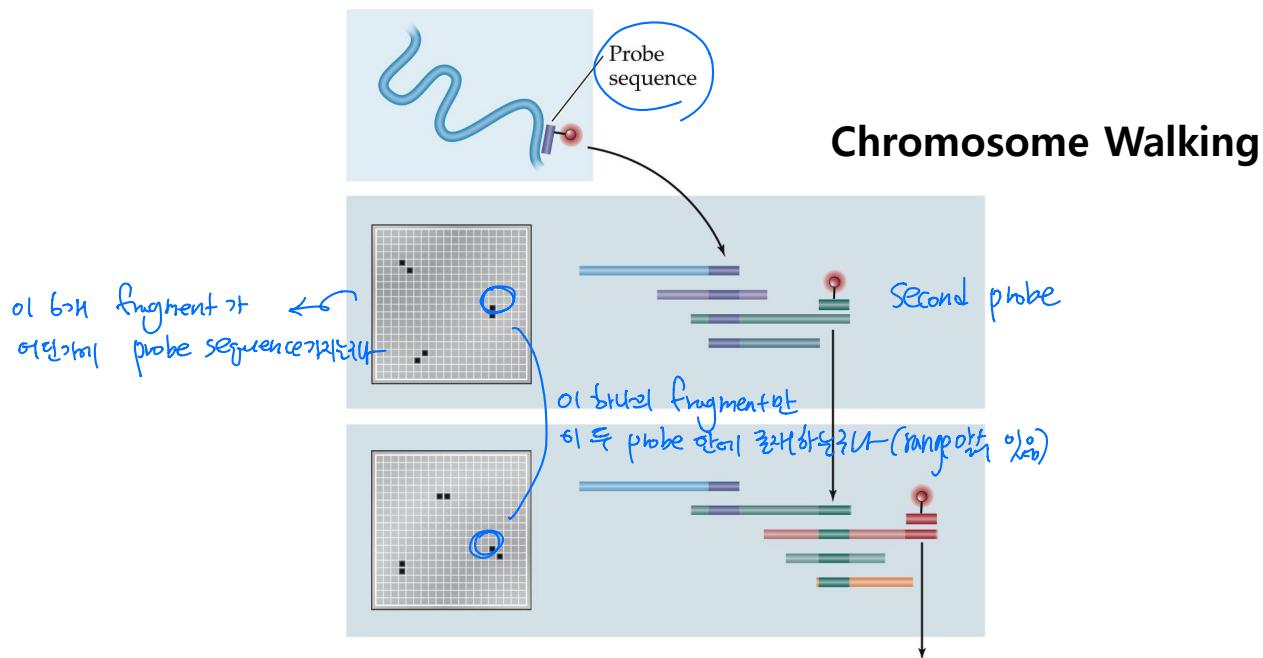
1. Assembling of contigs: Hybridization-based approaches

: chromosome walking

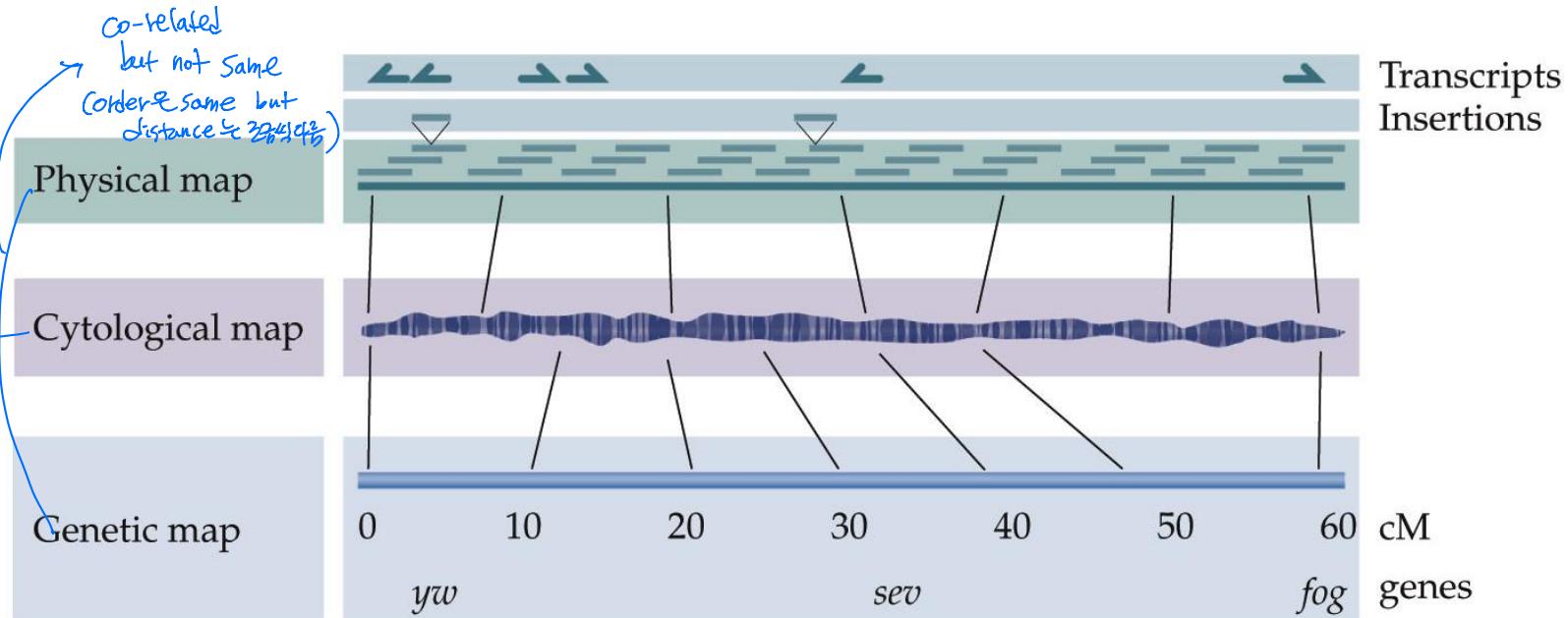
know sequence → design probe



: sequence-tagged sites (STSs): single occurrence in the genome

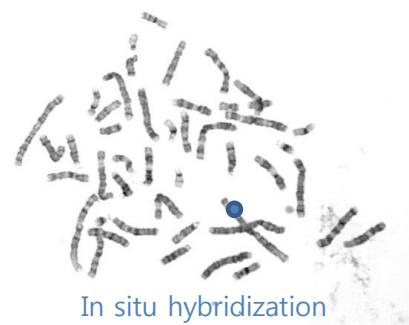


Mapping Genomes :Cytological Map

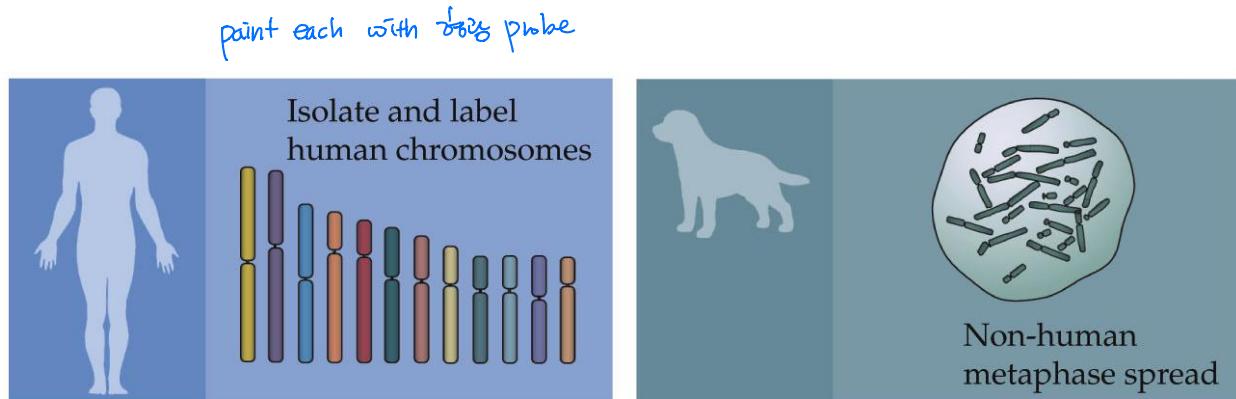


Physical map, Cytological map and Genetic map are correlated,

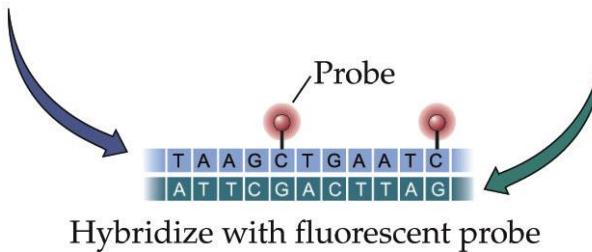
But not the same..



Comparative Genomics: Chromosome painting



Different color combination for each chromosome probe



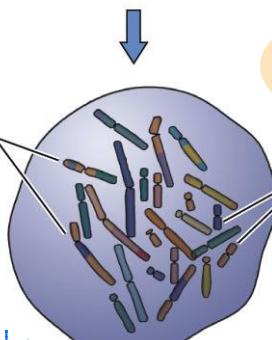
Hybridize with fluorescent probe

사람의 genome에 표지한
개의 genome에 FISH
→ 겹겹치는 부분 알 수 있음

Multicolored chromosomes indicate breakage/fusion events

Synteny

Compute large chunk of chromosome between diff. species



FISH: fluorescence in situ hybridization

Single-color chromosomes indicate complete correspondence between species

The Human Genome Project

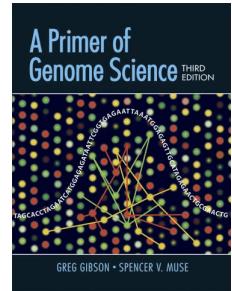
Sung Wook Chi

Division of Life Sciences, Korea University

What we will learn today

- **Sanger Sequencing / PHRED score**
- **Human genome project**
 - What we know about the human genome

genome 디지털



Chapter2. Genome
Sequencing and annotation



2001

Draft



2004

Completion



2007

Regulation



2008

Variation

1000
Genome
sequence

DNA sequencing (Sanger Sequencing)

DNA



James Watson

(1928-)

A, T, G, C

No. 4356 April 25, 1953 NATURE

MOLECULAR STRUCTURE OF NUCLEIC ACIDS

A Structure for Deoxyribose Nucleic Acid

J. D. WATSON
F. H. C. CRICK

Medical Research Council Unit for the Study of the Molecular Structure of Biological Systems, Cavendish Laboratory, Cambridge. April 2.

A detailed scientific illustration of the DNA double helix, showing the characteristic twisted ladder structure with phosphate groups on the outside and nitrogenous base pairs on the inside.

Sequencing



Fred Sanger
(1918-2013)

Protein Sequencing (1955, insulin)
Nobel prize 1958

DNA Sequencing (Sanger sequencing) (1977)
Nobel prize 1980

Sanger Sequencing: Principle

DNA polymerase reaction

template

template

ATGTGGCATGCTAGCTAGCCCTACGTATTGCAGGAT temp left

 Add primer

ATGTGGCATGCTAGCTAGCCCTACGTATTGCAGGAT

TACACCGTTACGATCG

– Primer

 Add nucleotides
and polymerase

ATGTGGCATGCTAGCTAGCCCTACGTATTGCAGGAT

TACACCGTACGATCGATCGGGATGC . . .

↓ Separate by electrophoresis

TACACCGTACGATCGATCGGGATGC

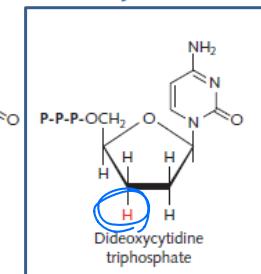
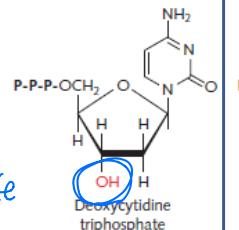
TACACCGTACGATCGATCGGGATG

TACACCGTACGATCGATCGGGAT

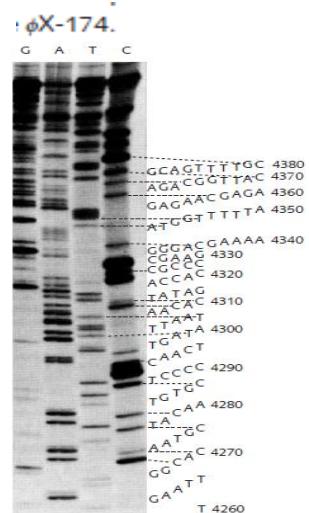
TACACCGTACGATCGATCGGGAA

▶ 헝겊색 다르게

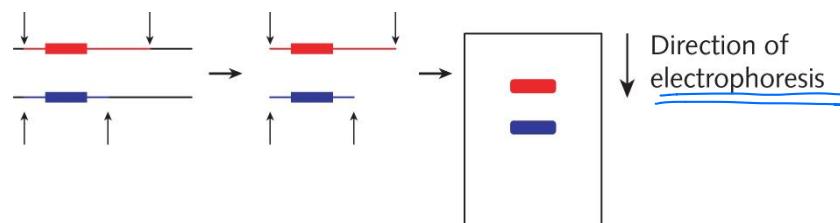
Nucleotide : Dideoxy NTP



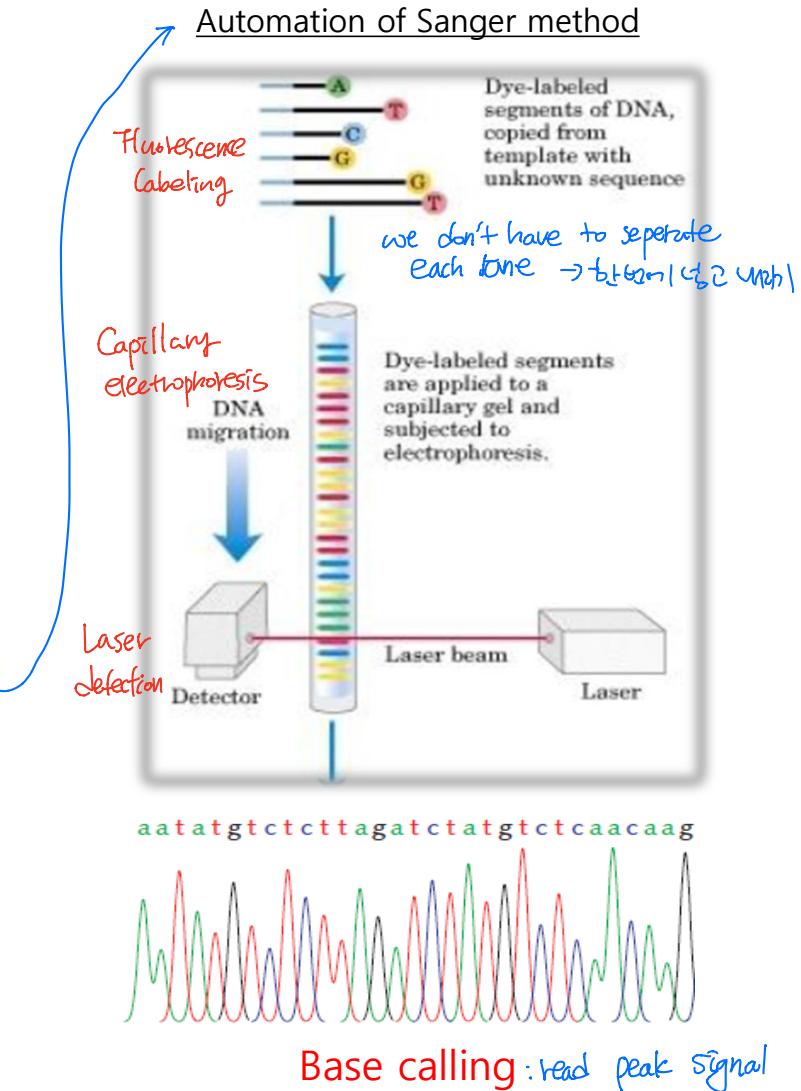
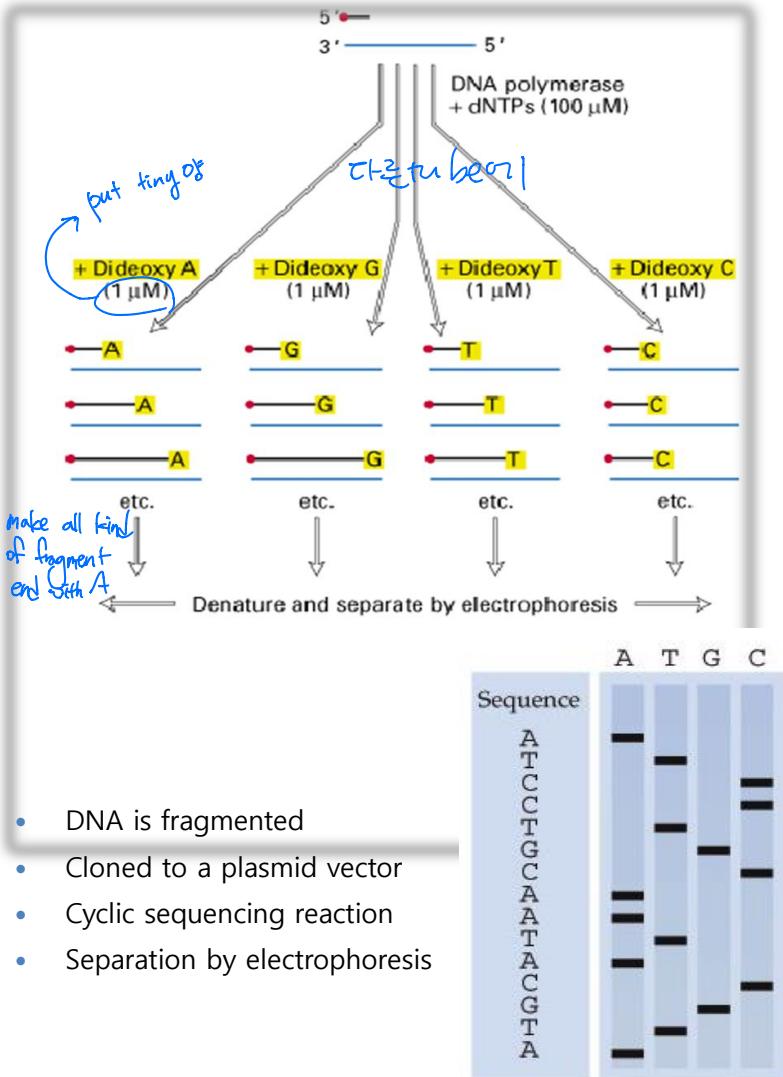
Chain termination



From: Sanger, F., Nicklen, S., & Coulson, A.R. (1977). DNA sequencing with chain-terminating inhibitors. Proc. Natl. Acad. Sci. USA **74**, 5463–5467.

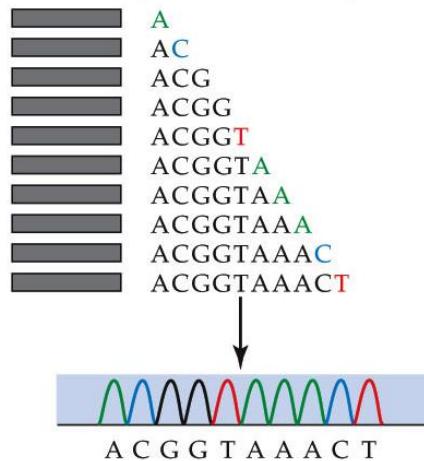


Sanger Sequencing : Chain termination method



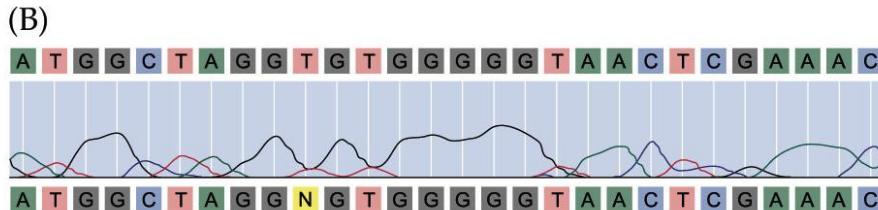
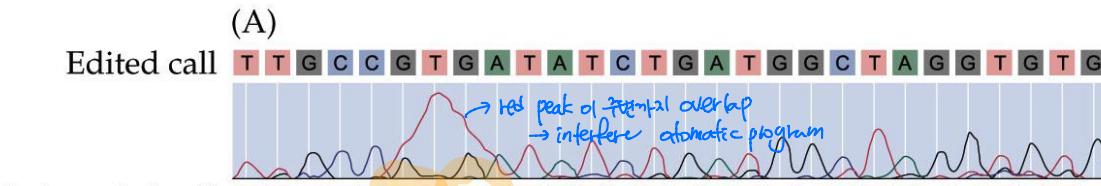
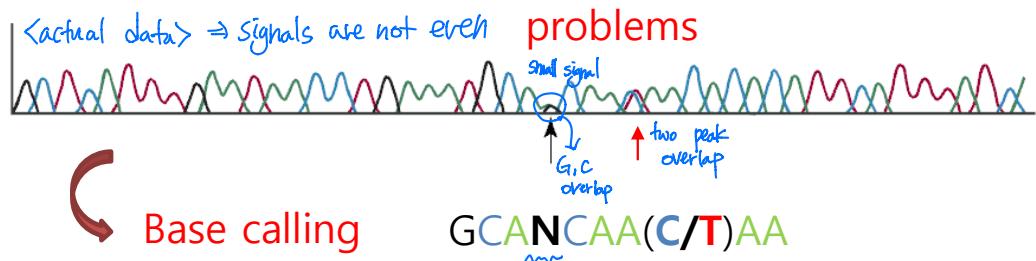
- DNA is fragmented
- Cloned to a plasmid vector
- Cyclic sequencing reaction
- Separation by electrophoresis

Sanger Sequencing : Base calling



EXERCISE 2.1 *Reading a sequence trace*

Read the following DNA sequence obtained by direct sequencing of a single individual organism, assuming that on the trace green is A, red is T, blue is C, and black is G. Remark on any ambiguities in the sequence.



- Shape of signal is import
- How to evaluate quality of base calling?

Need quality score !!

PHRED score: quality score of base calling

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



Phillip Green

Phred (Software for base-calling)
Base-calling
Algorithm (1998)

$$Q = -\log_{10} P$$

probability of error

Good

$Q > 20$,
Error rate < 0.01
Good

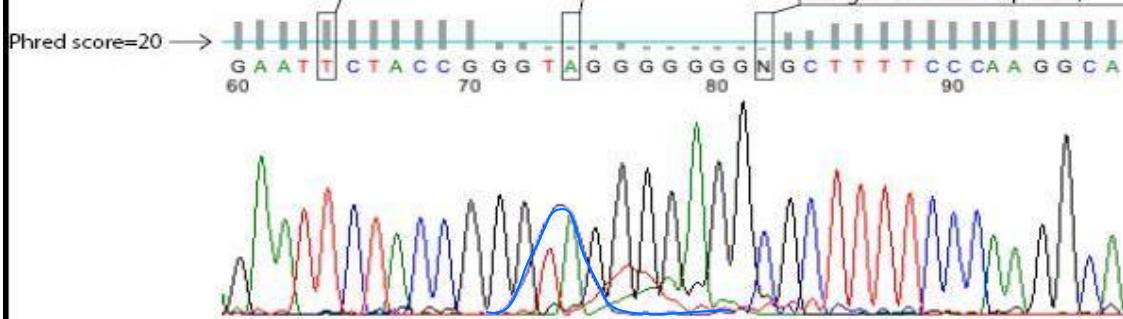
$Q = 20$ ($p=0.01$)

$Q < 20$,
Error rate > 0.01
bad

An example of a base that has been given a very high Phred score of 50, indicating that there is 99.999% probability that this base has been correctly assigned.

An example of a base that has been given a Phred score of 10, indicating that there is only a 90% probability that this base has been correctly assigned.

An example of a base for which no Phred score could be calculated, since the sequencer could not determine which base was present (therefore, an 'N' was designated in the sequence).



↳ sequencing 개발 사용

Locate predicted peaks, using Fourier methods to fit best distribution

Locate observed peaks, for which the area under the concavity exceeds 10% of the previous 10 peaks or 5% of the previous one

Match observed and predicted peaks using a three-stage shifting algorithm

Find missing peaks

Assess error probabilities of each peak according to four-parameter model

Fluorescence peaks

Phred
Base-calling
algorithm

Sequence reads

Human genome sequencing with Sanger method

Automation of sanger sequencing method



Vector – insert (~1.5kb) → need to be fragmented under this size

가지 않을 수 있음

- Sanger Sequencing
- Base-calling



Sequence reads (~1.5kb)
Quality scores

3 billion base pairs

Human genome

Fragmentation
(random)

→ 다시 조각

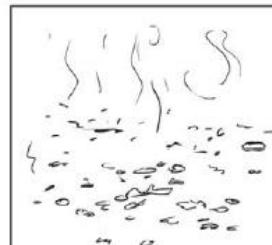
Sequencing

Shotgun Sequencing Strategy



Assembly

집어.

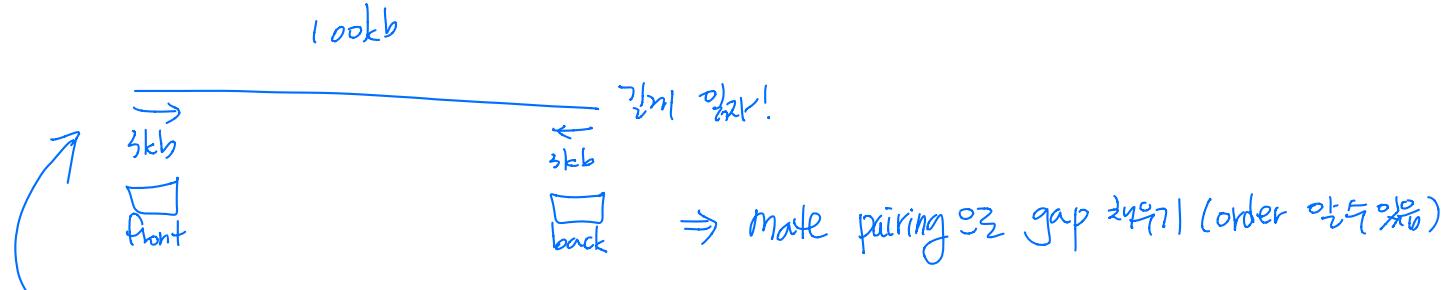


need to recombine

What we will learn today

- Sanger Sequencing / PHRED score
- Human genome project
 - Hierarchical shotgun sequencing
 - Whole genome shotgun sequencing
 - What we know about the human genome

(b) library



Human Genome Projects

Celera



Craig Venter

International Human Genome Sequencing Consortium



Francis Collins NIH MRC

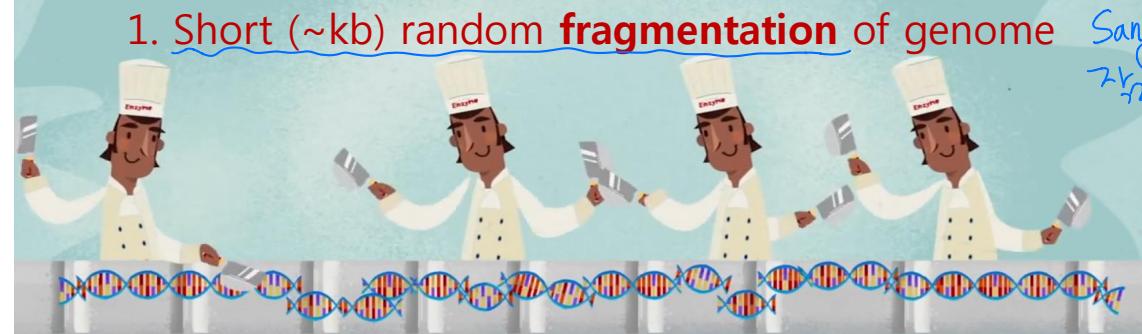
VS.
Competition

Whole-genome shotgun sequencing
break into small piece

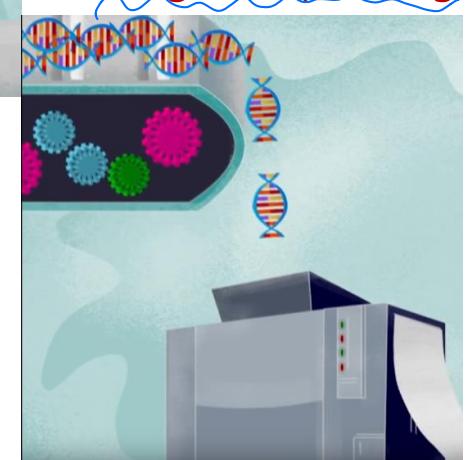
Hierarchical shotgun sequencing

1. Short (~kb) random fragmentation of genome

Sanger sequencing on gel electrophoresis
sequencing by synthesis

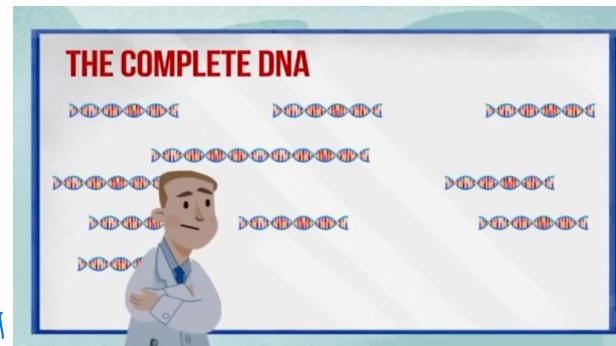


2. Sanger Sequencing



3. Assembly
(Contigs)

Small fragment
→ back into original

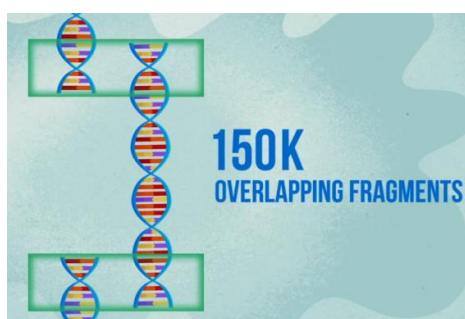


Hierarchical Shotgun Sequencing

①

Used Alu insertion fragmentation

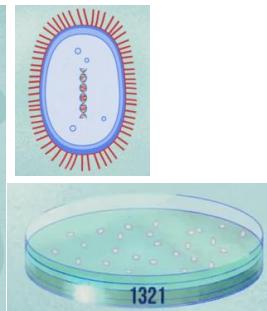
Fragmentation of genome



150kb 사이즈의 클로즈를 만들어
BAC / YAC에 넣어 관리

→ Store / amplify 가능

BAC/YAC



Scaffold (Physical maps)



Chromosome walking
named all the clones
with chromos.

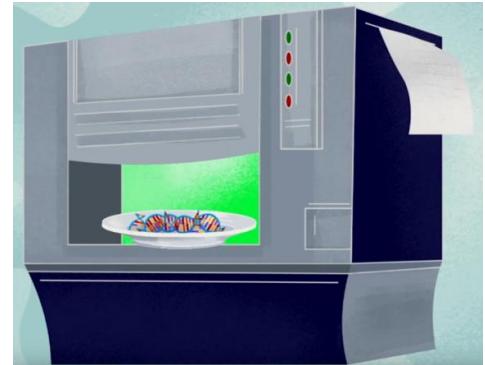
Fragmentation of BAC/YAC clone



②

각 Fragment을
contig의 연결을
작은 사이즈로 만들기 (Kb 단위)
알아내면 → 물리지도로 훨씬 이동이

Sanger Sequencing

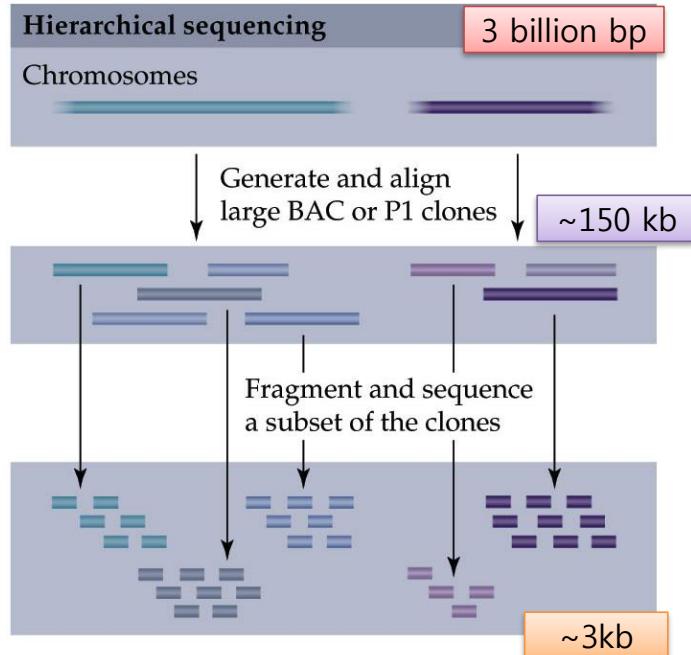
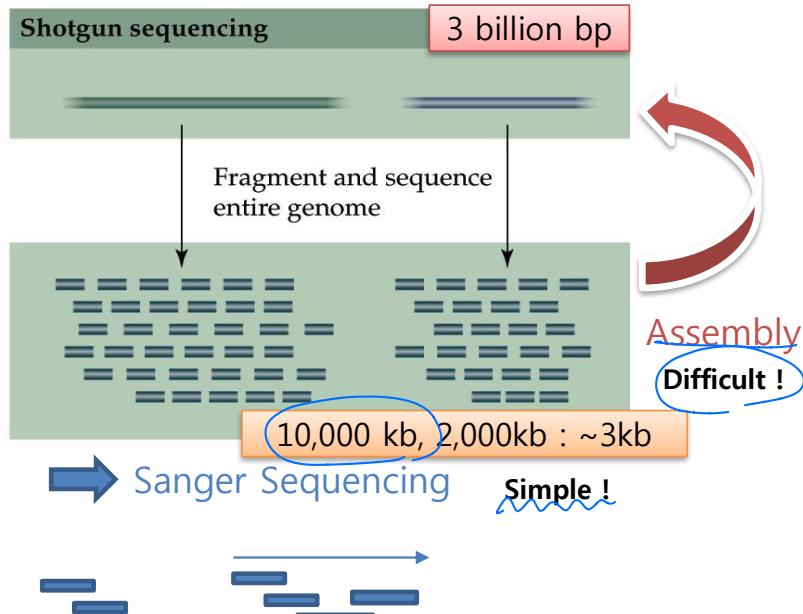


3. Assembly (Contigs)

TTTGCTCGGA
TCCTGGCTCC
AGGCCGCGTC

Hierarchical vs. whole genome shotgun sequencing

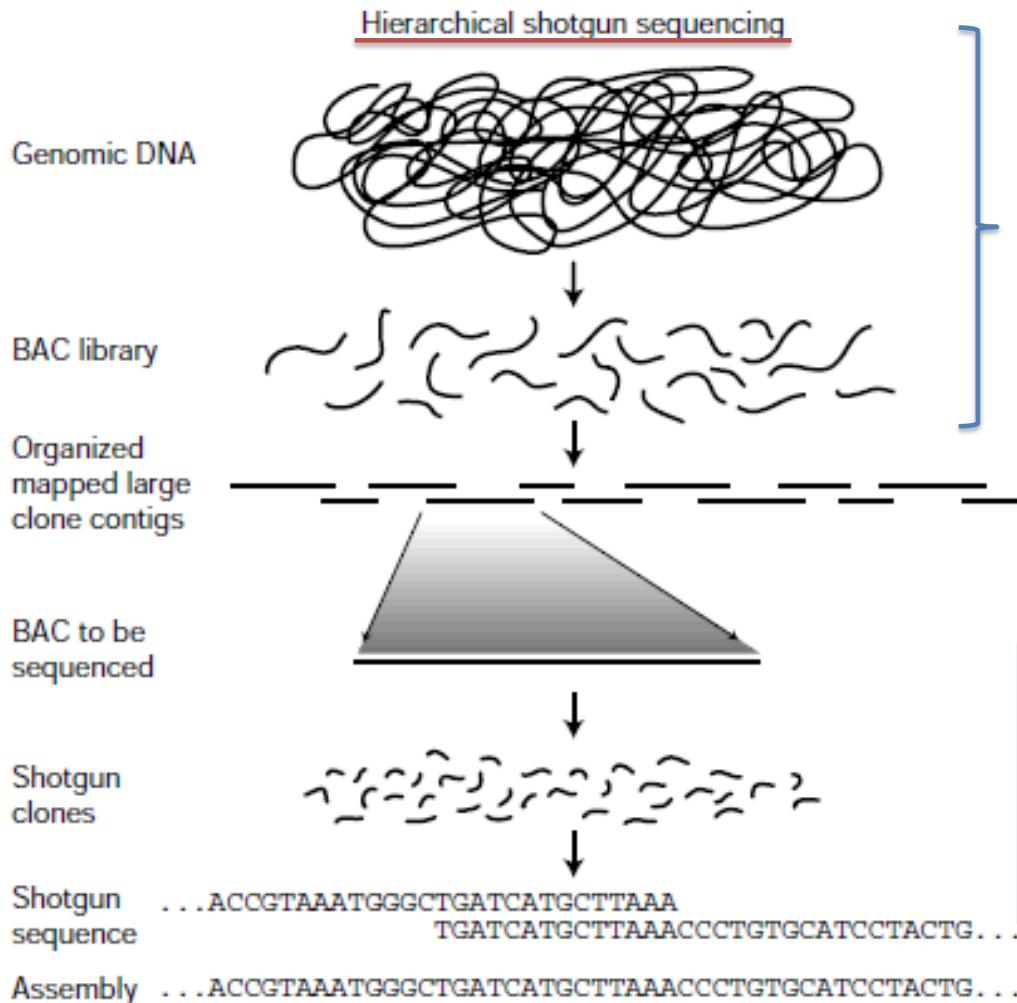
Initially attempting to fragment



- Scaffold
- Overlapper> Contigs > Supercontig

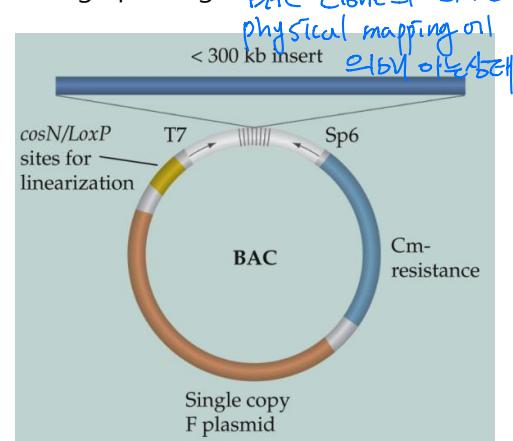
Overlap \Rightarrow gaps may be other repeat sequence

Hierarchical shotgun sequencing strategy



'BAC-to-BAC' method

- Scaffold: decide order/orientation
- Chromosome walking
- Fingerprinting



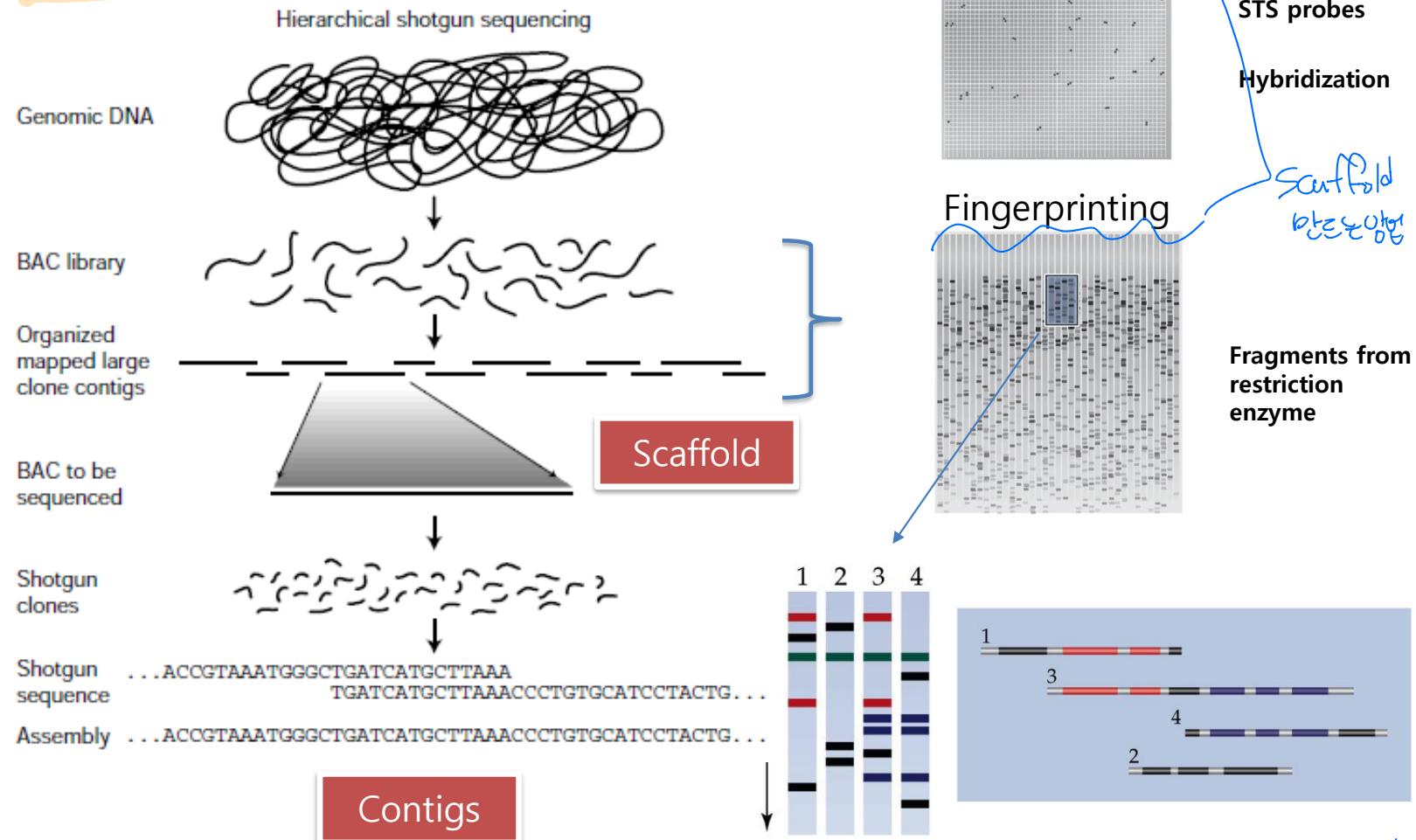
BOX
3.8

BACs: bacterial artificial chromosomes

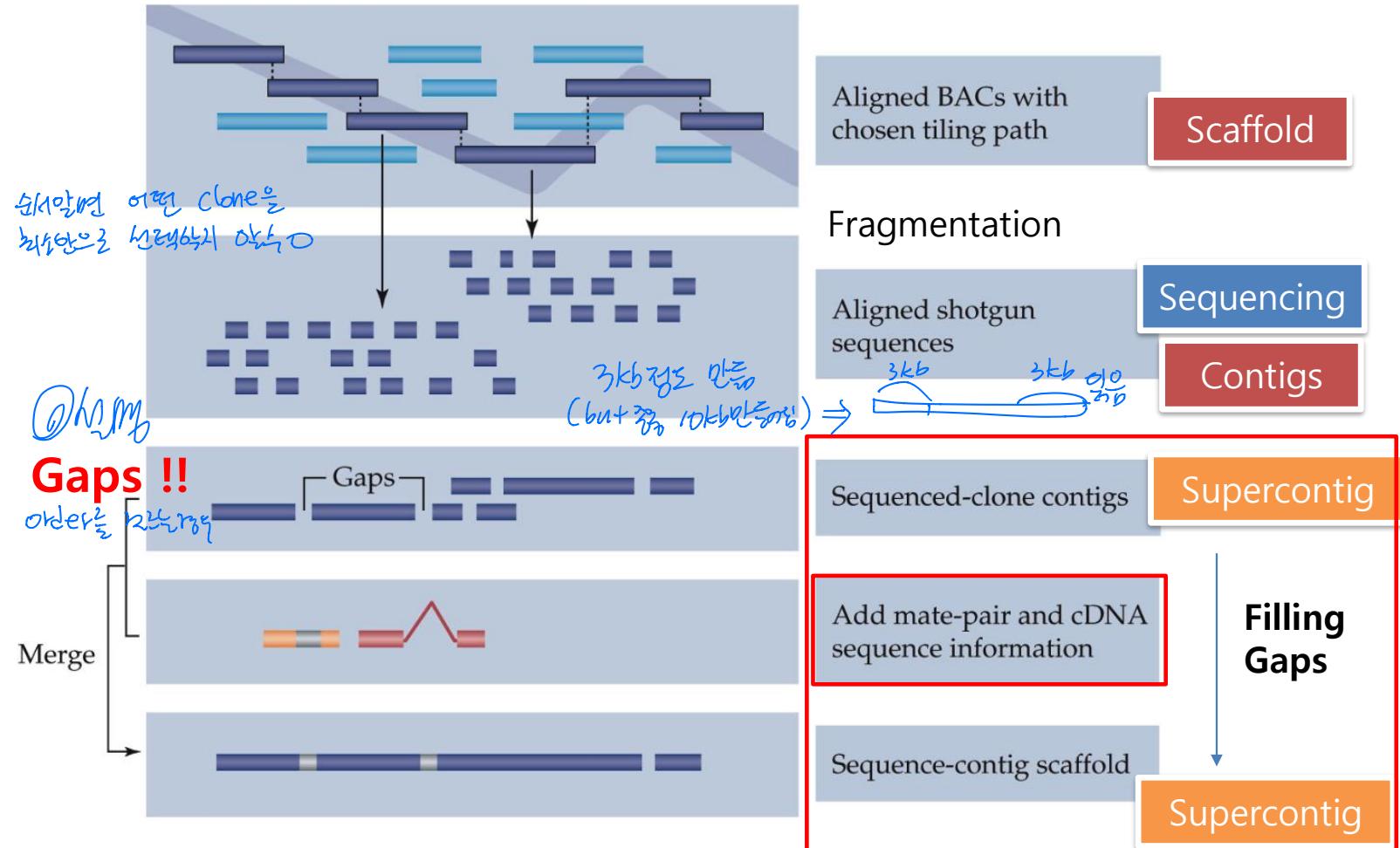
A plasmid is a small piece of double-stranded DNA in a bacterial cell, in addition to the main genome. A bacterial artificial chromosome, or BAC, is a plasmid containing foreign DNA – for instance, a fragment of the human genome. A typical BAC, in an *Escherichia coli* cell, can carry about 250 000 bp.

Hierarchical shotgun sequencing: Scaffold & Contigs

Aligning BAC clones by hybridization and fingerprinting



Hierarchical assembly of a sequence-contig scaffold (supercontig)

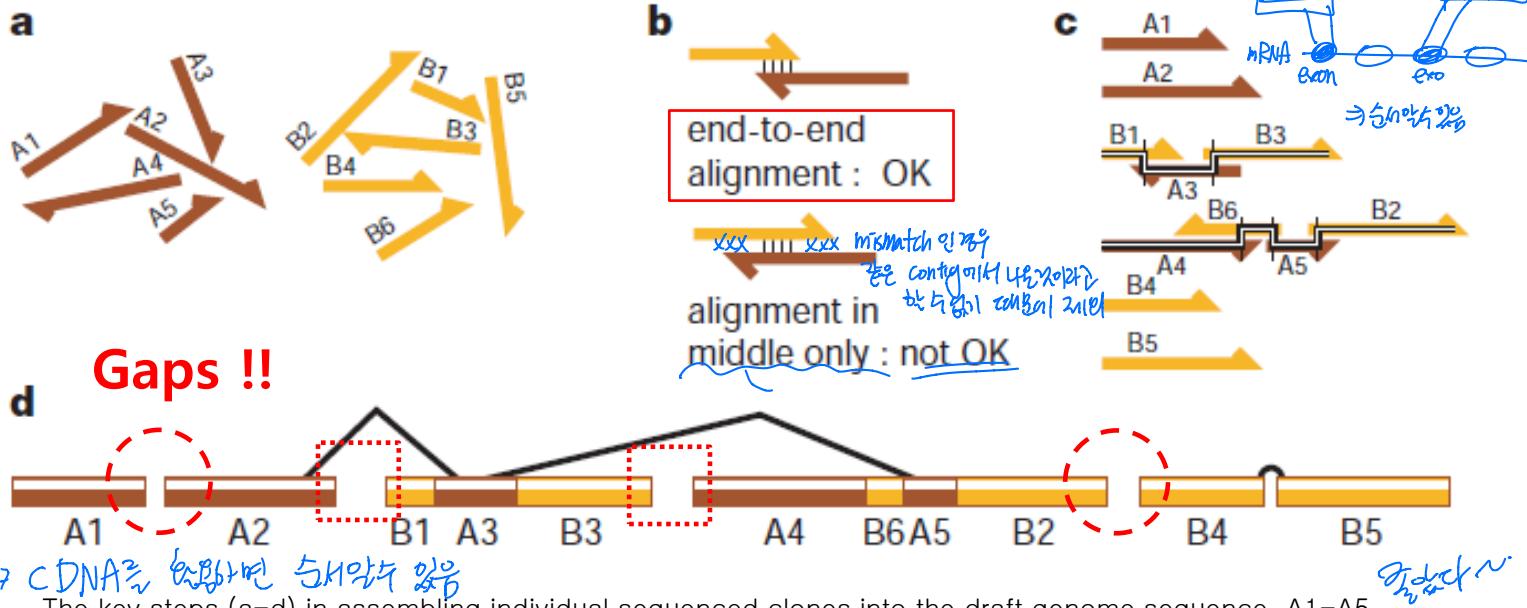


Challenge : How to merge contigs (filling gaps) to generate supercontig?

Assembling individual sequenced clones into the draft genome sequence



Assembly (Contigs > Supercontig)



The key steps (a-d) in assembling individual sequenced clones into the draft genome sequence. A1-A5 represent initial sequence contigs derived from shotgun sequencing of clone A, and 1-B6 are from clone B.

Human genome project: hierarchical shotgun sequencing

The New York Times

(Feb, 13, 2001)

READING THE BOOK OF LIFE

READING THE BOOK OF LIFE; Grad Student Becomes Gene Effort's Unlikely Hero

By NICHOLAS WADE

Published: February 13, 2001

A surprising hero helped the consortium of academic scientists decoding the human genome to avoid a drubbing by its rival, the Celera Genomics company. Scientists throughout the world are now beating an electronic path to his Web site, where they can analyze and download the human genome sequence. He is a graduate student at the University of California at Santa Cruz, and his name is not Clark but James Kent.

Methods

Assembly of the Working Draft of the Human Genome with GigAssembler

W. James Kent^{1,3} and David Haussler²

Genome Res. 2001 September; 11(9): 1541–1548

CDNA
gap filling program



can used to fill the gap

UCSC Genome Bioinformatics

UCSC genome browser

<http://genome.ucsc.edu/>

Genome Res. 2002 12 (6): 996–1006



articles

Initial sequencing and analysis of the human genome

International Human Genome Sequencing Consortium*

(First draft, 2001)

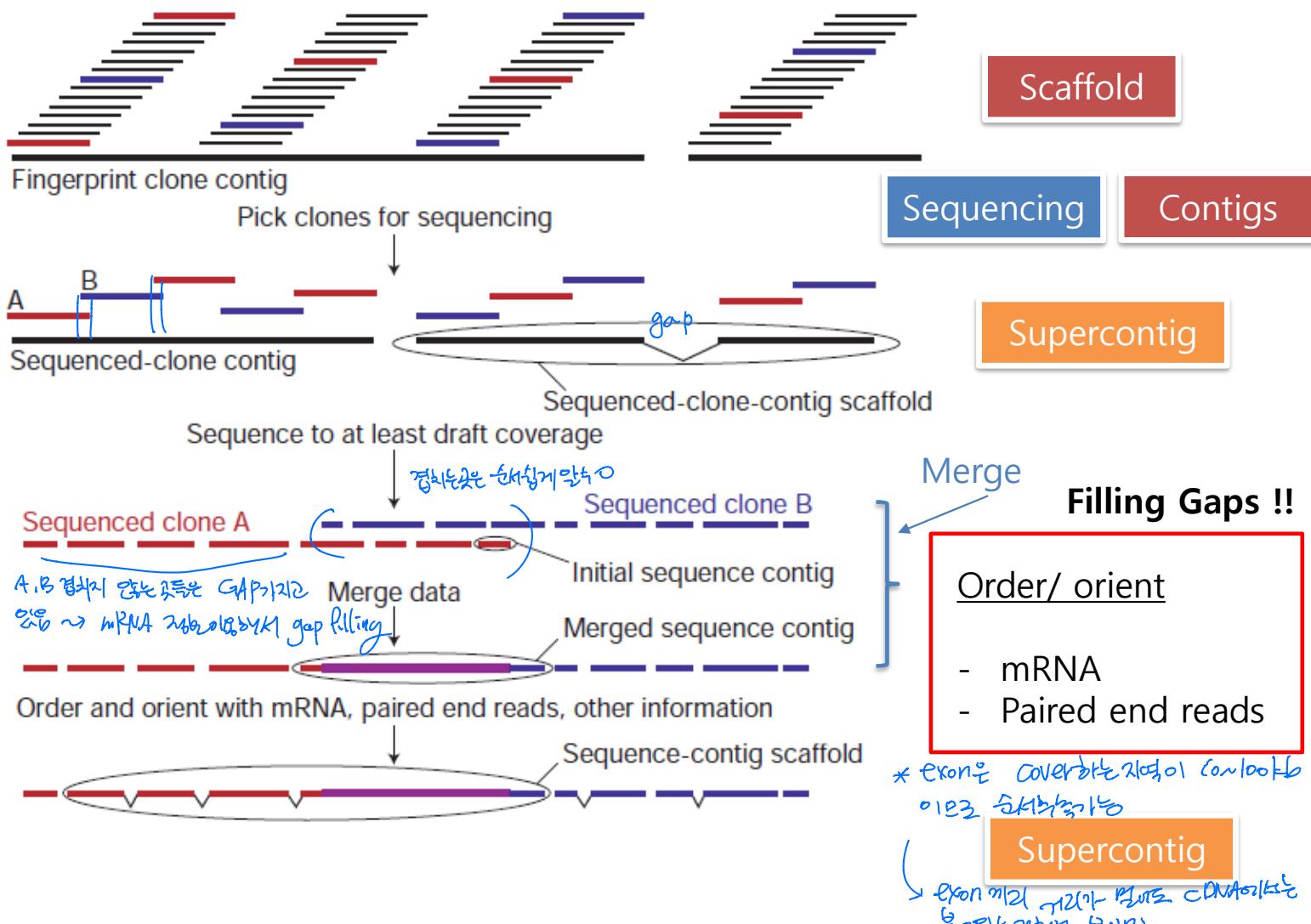


Jim Kent



David Haussler

Genome sequence assembly (Levels of clone)



Whole genome shotgun sequencing (Celera's approach)

우리방법보다 더 많은 gap



Problems: Repetitive sequences

반복서열이 너무 많아서 문제

ACGTTGTCGACTAGGATCGCTCGTGAGGG ATGCAGCAGCAGCAGTTGTGAGAATCCAC
ATCGCTCGTGAGGGGCT ATAATGCAGCAGCAGCA

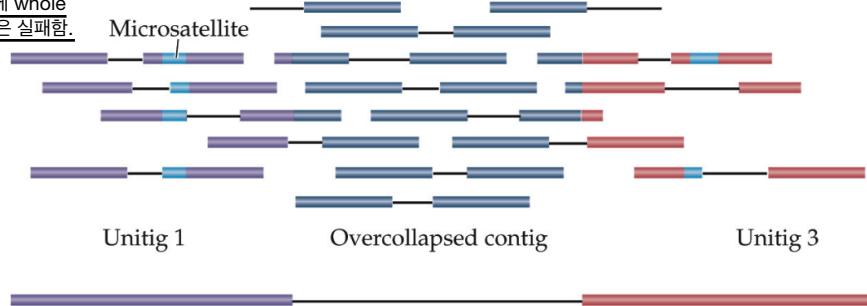


조각들 중에서 어느 방향으로든 unique한 서열 찾으려는

Unitigs and repeat resolution

Unique contiguous sequence alignments

repeats



두 region을 cover하는 다른 범위의 clone을
만들어 → 두 unitig 사이를 끊음

Solution: Mate pairs from larger fragments

3 billion bp

자른 이유 : 생어seq.로는 3kb정도밖에 못읽음 (긴
서열이라면 앞뒤로3kb정도... 그래서 길게 잘라놓으면
cDNA방법과 같이 같은 contig내에 있는 서열임
을 알 수 있음) 반복서열의 존재 때문에 whole
genome를 모두 3kb로 나누는 방법은 실패함.

1.5kb

2,000kb

Pair-end read
(Mate pairs)

1.5kb

10,000 kb

1.5kb

Completion of human genome project

BOX
3.9

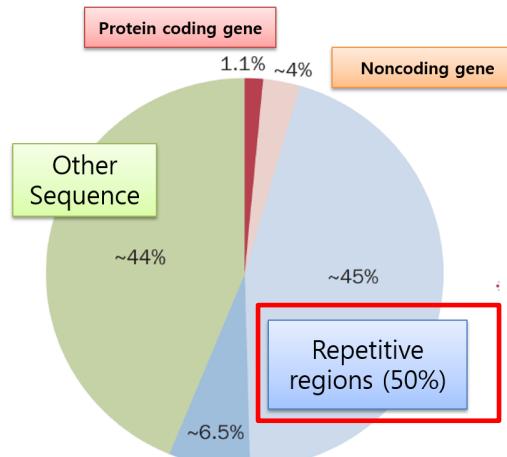
Common and different steps in 'BAC-to-BAC' and whole-genome shotgun methods

'BAC-to-BAC' method

Whole-genome shotgun method

1. Make random cuts to produce fragments of:
-150 kb ~2000 kb and 10 000 kb
 2. Make plasmid library in BACs.
 3. Fingerprint, overlap, 3. Skip this step.
and order BAC clones.
 4. Partially sequence 1500 bp subfragments
of individual clones.
 5. Assemble overlaps by computer.

to find
scaffold



2001

Draft

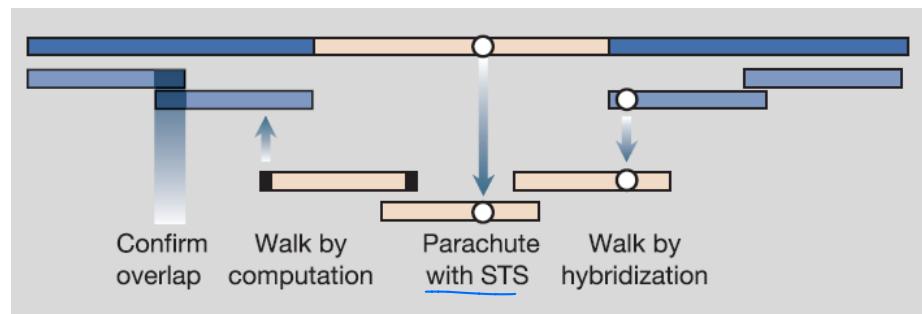


2004

Completion

2004년 1월 gap 학기 수강권 이용권으로 대체 가능

Remaining gaps > finishing the physical maps



iterative ‘walking’ from the ends of contigs

What we will learn today

- **Sanger Sequencing** / PHRED score
- **Human genome project**
 - Hierarchical shotgun sequencing
 - Whole genome shotgun sequencing
 - What we know about the human genome



2004

centromere, telomere, repeat, segment, articles

Finishing the euchromatic sequence of the human genome

International Human Genome Sequencing Consortium*

* A list of authors and their affiliations appears in the Supplementary Information

The sequence of the human genome encodes the genetic instructions for human physiology, as well as rich information about human evolution. In 2001, the International Human Genome Sequencing Consortium reported a draft sequence of the euchromatic portion of the human genome. Since then, the international collaboration has worked to convert this draft into a genome sequence with high accuracy and nearly complete coverage. Here, we report the result of this finishing process. The current genome sequence (Build 35) contains 2.85 billion nucleotides interrupted by only 341 gaps. It covers ~99% of the euchromatic genome and is accurate to an error rate of ~1 event per 100,000 bases. Many of the remaining euchromatic gaps are associated with segmental duplications and will require focused work with new methods. The near-complete sequence, the first for a vertebrate, greatly improves the precision of biological analyses of the human genome including studies of gene number, birth and death. Notably, the human genome seems to encode only 20,000–25,000 protein-coding genes. The genome sequence reported here should serve as a firm foundation for biomedical research in the decades ahead.

↳ 같은 결과를 바탕으로 유전체 개수 추정

Current state of human genome sequences

assembly
most recent assembly
↑
Dec. 2013 (GRCh38/hg38)
Feb. 2009 (GRCh37/hg19)
Mar. 2006 (NCBI36/hg18)
May 2004 (NCBI35/hg17)
July 2003 (NCBI34/hg16)

GRCh38 Genome Reference Consortium
Human Reference 38 ([GRCh 38/ hg38](#))

chr10.fa file

chr11.fa
chr11_g1000202_random.fa
chr12.fa
chr13.fa
chr14.fa
chr15.fa
chr16.fa
chr17_ctg5_hap1.fa
chr17.fa
chr17_g1000203_random.fa
chr17_g1000204_random.fa
chr17_g1000205_random.fa
chr17_g1000206_random.fa
chr18.fa
chr18_g1000207_random.fa
chr19.fa
chr19_g1000208_random.fa
chr19_g1000209_random.fa

chr1.fa
chr1_g1000191_random.fa
chr1_g1000192_random.fa
chr20.fa
chr21.fa
chr21_g1000210_random.fa
chr22.fa
chr2.fa
chr3.fa
chr4_ctg9_hap1.fa
chr4.fa
chr4_g1000193_random.fa
chr4_g1000194_random.fa
chr5.fa
chr6_apd_hap1.fa
chr6_cox_hap2.fa
chr6_dbb_hap3.fa
chr6.fa
chr6_mann_hap4.fa

chr6_mcf_hap5.fa
chr6_qbl_hap6.fa
chr6_sssto_hap7.fa
chr7.fa
chr7_g1000195_random.fa
chr8.fa
chr8_g1000196_random.fa
chr8_g1000197_random.fa
chr9.fa
chr9_g1000198_random.fa
chr9_g1000199_random.fa
chr9_g1000200_random.fa
chr9_g1000201_random.fa
chrM.fa
chrUn_g1000211.fa
chrUn_g1000212.fa
chrUn_g1000213.fa
chrUn_g1000214.fa
chrUn_g1000215.fa

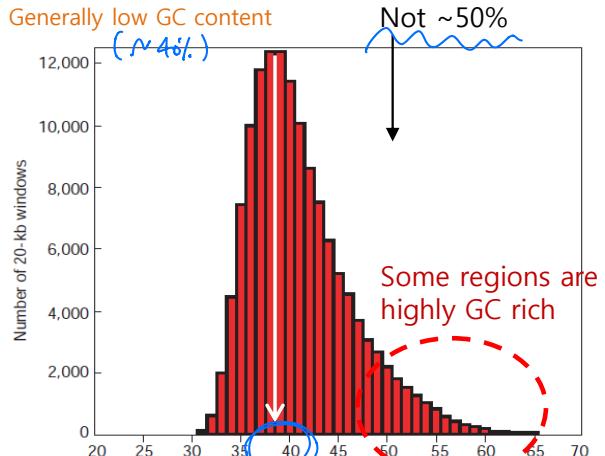
3.2 Gb text file (Fasta format)

chrUn_g1000216.fa
chrUn_g1000217.fa
chrUn_g1000218.fa
chrUn_g1000219.fa
chrUn_g1000220.fa
chrUn_g1000221.fa
chrUn_g1000222.fa
chrUn_g1000223.fa
chrUn_g1000224.fa
chrUn_g1000225.fa
chrUn_g1000226.fa
chrUn_g1000227.fa
chrUn_g1000228.fa
chrUn_g1000229.fa
chrUn_g1000230.fa
chrUn_g1000231.fa
chrUn_g1000232.fa
chrUn_g1000233.fa
chrUn_g1000234.fa

chrUn_g1000235.fa
chrUn_g1000236.fa
chrUn_g1000237.fa
chrUn_g1000238.fa
chrUn_g1000239.fa
chrUn_g1000240.fa
chrUn_g1000241.fa
chrUn_g1000242.fa
chrUn_g1000243.fa
chrUn_g1000244.fa
chrUn_g1000245.fa
chrUn_g1000246.fa
chrUn_g1000247.fa
chrUn_g1000248.fa
chrUn_g1000249.fa
chrX.fa
chrY.fa

- Unplaced contig names consist of the chromosome number, followed by the NCBI accession number, followed by "random"
- Unlocalized contig names (contigs whose associated chromosome is not known) consist of "chrUn" followed by the NCBI accession number

GC contents in human genome



CpG island

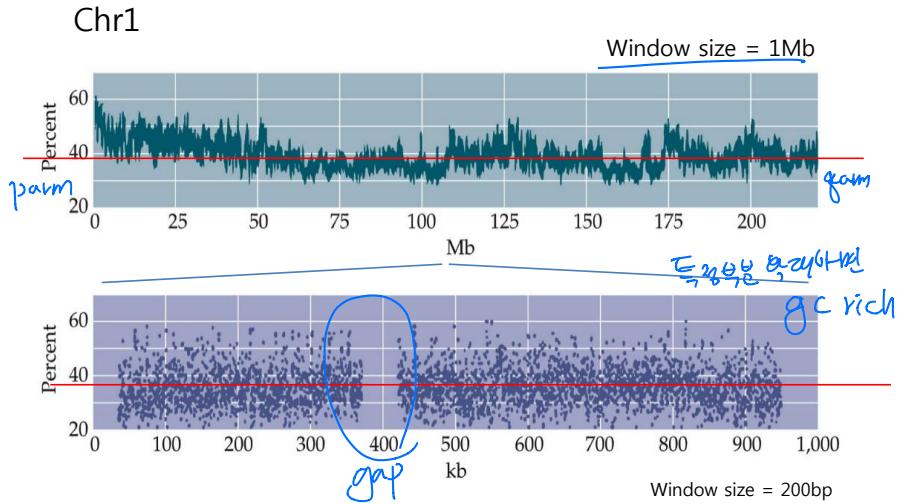
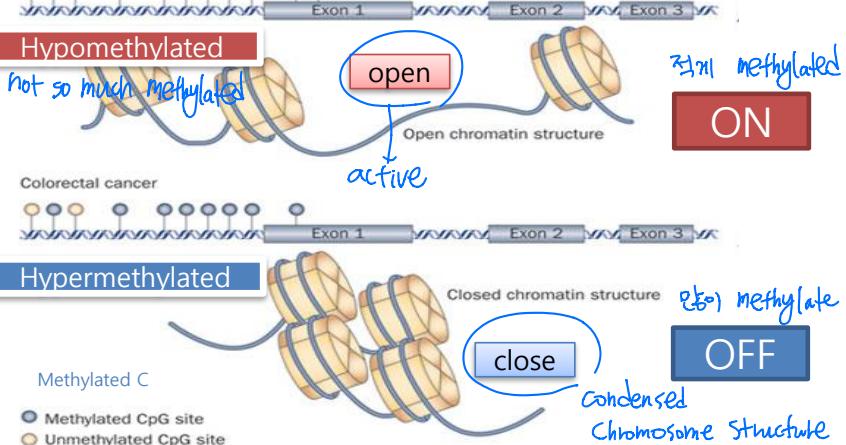
Differential regulation

GC content

region(s) A/B

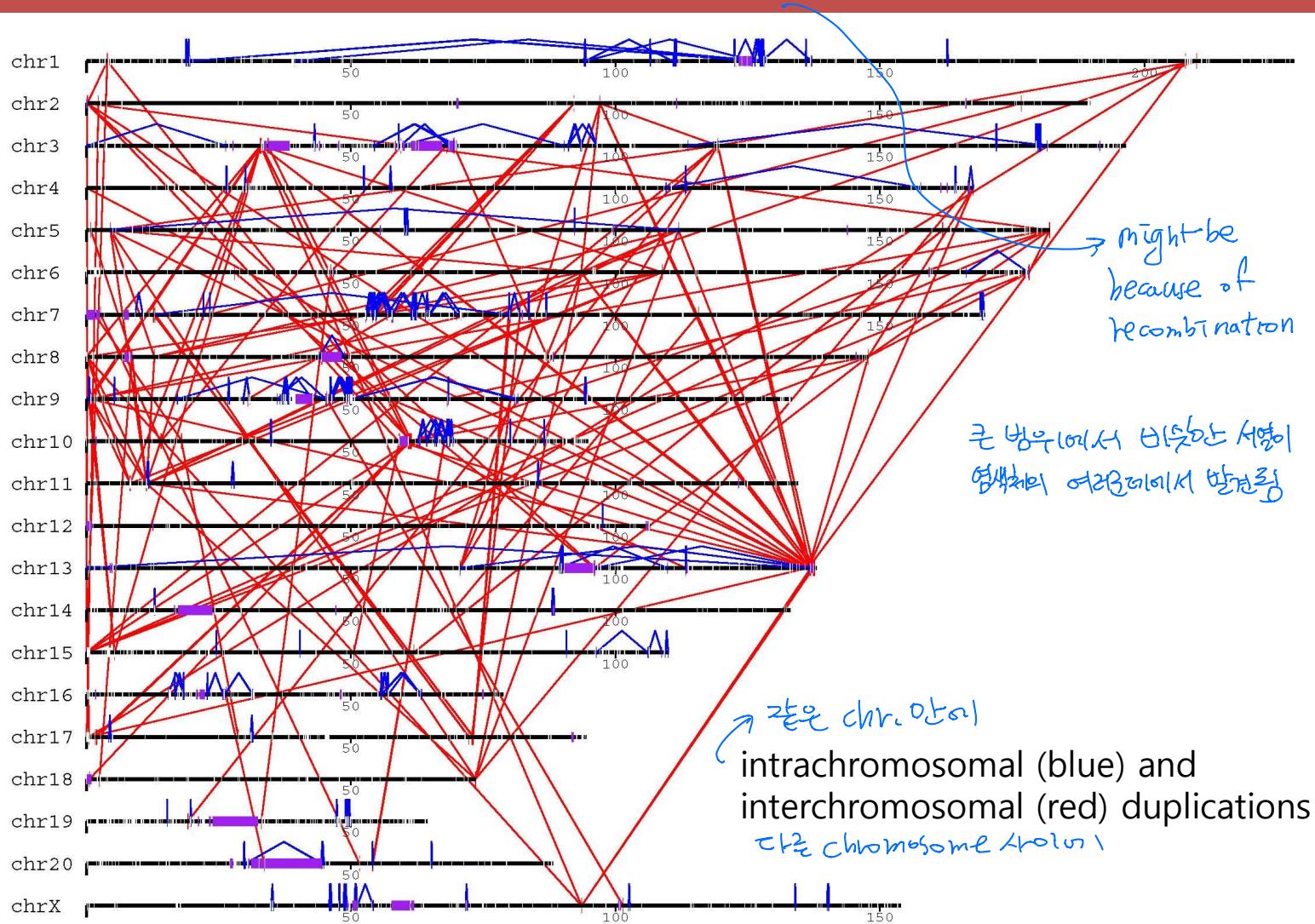
Normal Colon epithelium

shut-down / active b/w regulation



↑ G+C rich
CpG island ⇒ methylation에 의한 regulation 가능

Segmental duplication



Segmental duplication

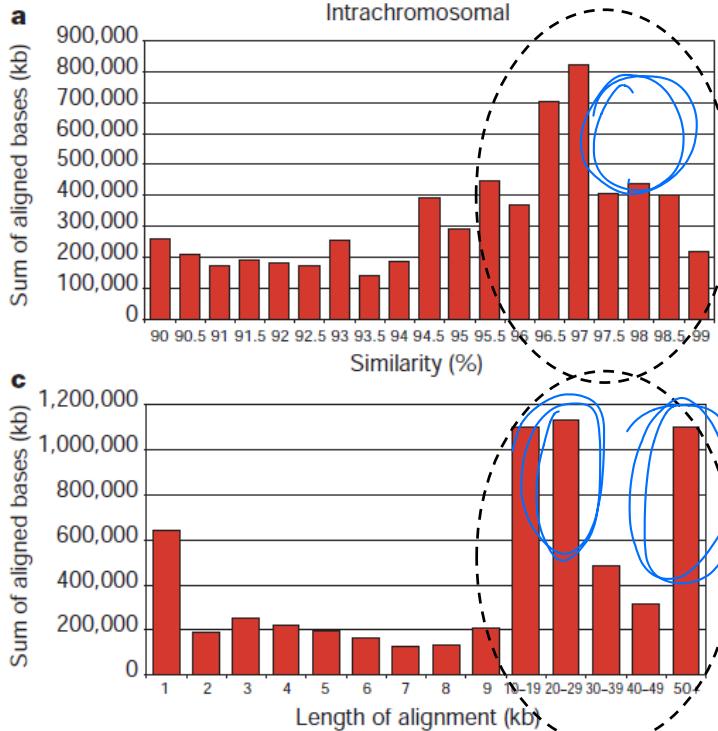
(B) TDW

Off-
within recombination

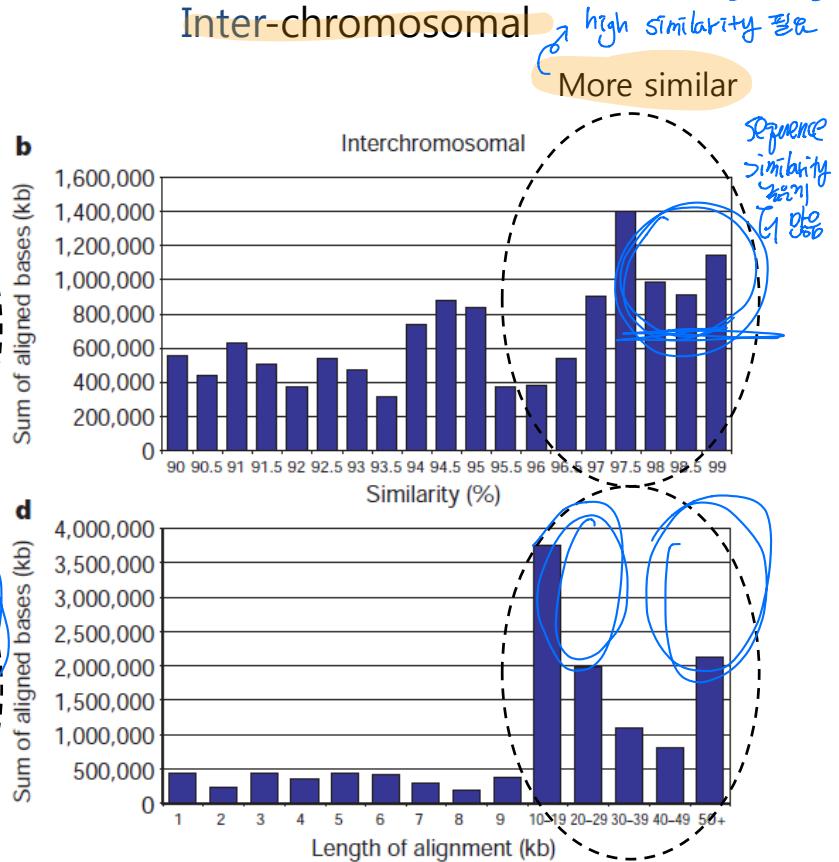
the chromosome atom
translocation

high similarity

Intra-chromosomal



Inter-chromosomal



More similar

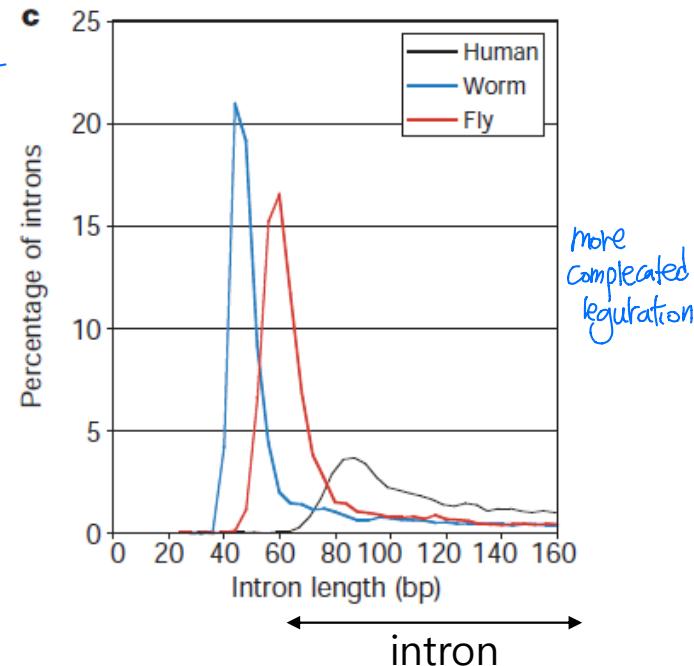
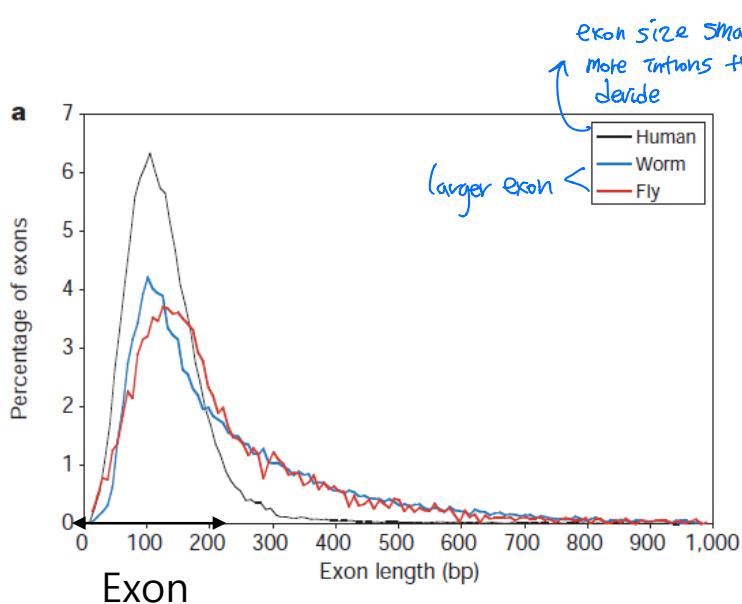
Sequence
similarity
증가
증가

shorter

Characteristics of human genes

Q1

	Median	Mean
Internal exon	122 bp	145 bp
Exon number	7	8.8
Introns	1,023 bp	3,365 bp
3' UTR	400 bp	770 bp
5' UTR	240 bp	300 bp
<u>Coding sequence</u> (CDS)	1,100 bp	1,340 bp
Genomic extent	367 aa	447 aa
	14 kb	27 kb



Summary: Human Genome Project



<https://www.youtube.com/watch?v=-gVh3z6MwdU>