

Summary & Review

: introduction to genomics

Sung Wook Chi
Division of Life Sciences, Korea University

What we learned in the previous lecture

1. Introduction

- What is **genomics** and functional genomics?

2. Architecture of Genome (Human, Eukaryote)

- Chromosome, Gene Structure (genome browser)

3. Flow of genomic information

- Gene expression, ESTs, epigenomic information

4. Today's Genomics

- Human genome project, NGS, personal genomics

5. Human genome

- Repeats (retrotransposons)

6. Open access of genome sequences

- Comparative genomics, Variations (CNV)

Functional Genomics

What I introduced in the previous lecture

Genome

Gene expression



Phenotype

1. Genome
2. Human Genome Project
3. Genome Annotation
4. Genomic Variation
5. Comparative Genomics

Genome-wide
methods

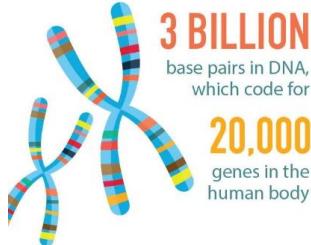
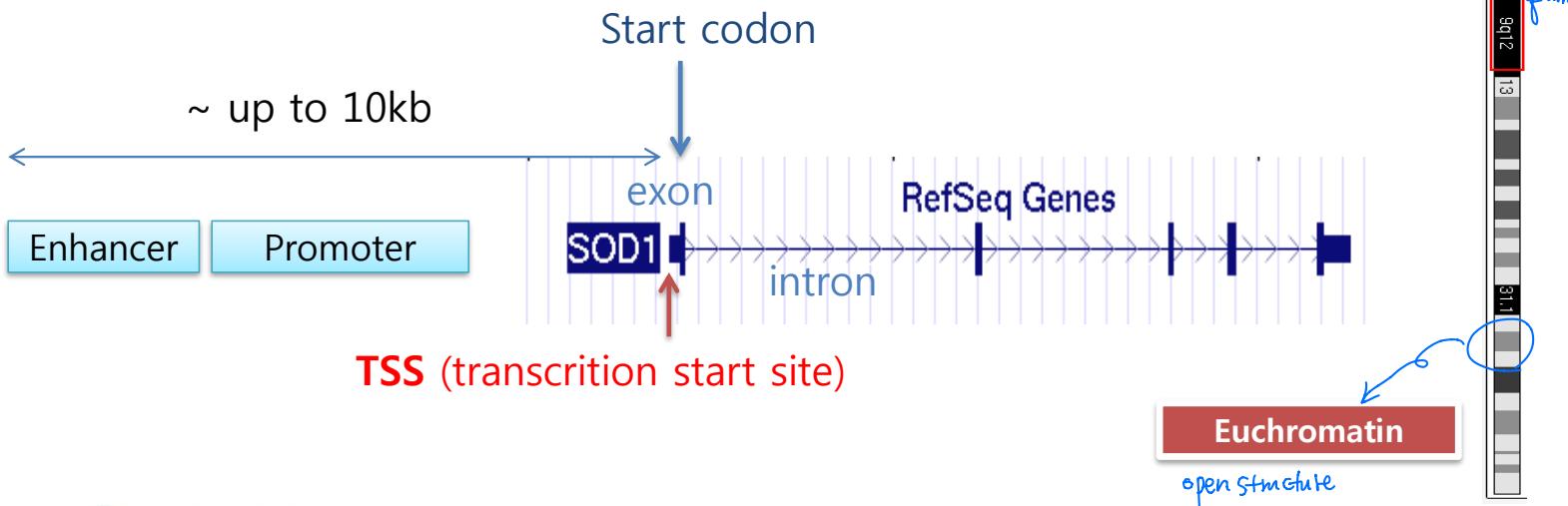
1. Next-Generation Sequencing (NGS)
(WGS, Exome-Seq)
2. Microarray
3. RNA-Seq
4. ChIP-Seq
5. Ribosomal Profiling
6. Bioinformatics

1. Biological function
2. Diseases
3. Quantitative trait /population

Genome and gene structure

Genome: whole DNA sequences (Gene + Chromosome)

- DNA + Histone (2A, 2B, 3, 4) x 2 = Nucleosome
 - Nucleosome > Chromatin > Chromosome



Human Genome Project (1990-2003)

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

Flow of genomics information

Phenotype = **genotype** + (environment + life history + **epigenetics**)

Central Dogma

Genome

Transcriptome

Proteome

The diagram illustrates the process of gene expression. It starts with DNA, which undergoes Transcription to produce Pre-mRNA. This Pre-mRNA then undergoes Processing to become mRNA. Finally, mRNA undergoes Translation to produce protein.

```
graph LR; DNA[DNA] -- "Transcription" --> PremRNA[Pre-mRNA]; PremRNA -- "Processing" --> mRNA[mRNA]; mRNA -- "Translation" --> Protein[protein]
```

regulation point

- 1) initiation
- 2) elongation

3) Splicing / Poly (A)

4) Stability

5) Repression

Complexity !

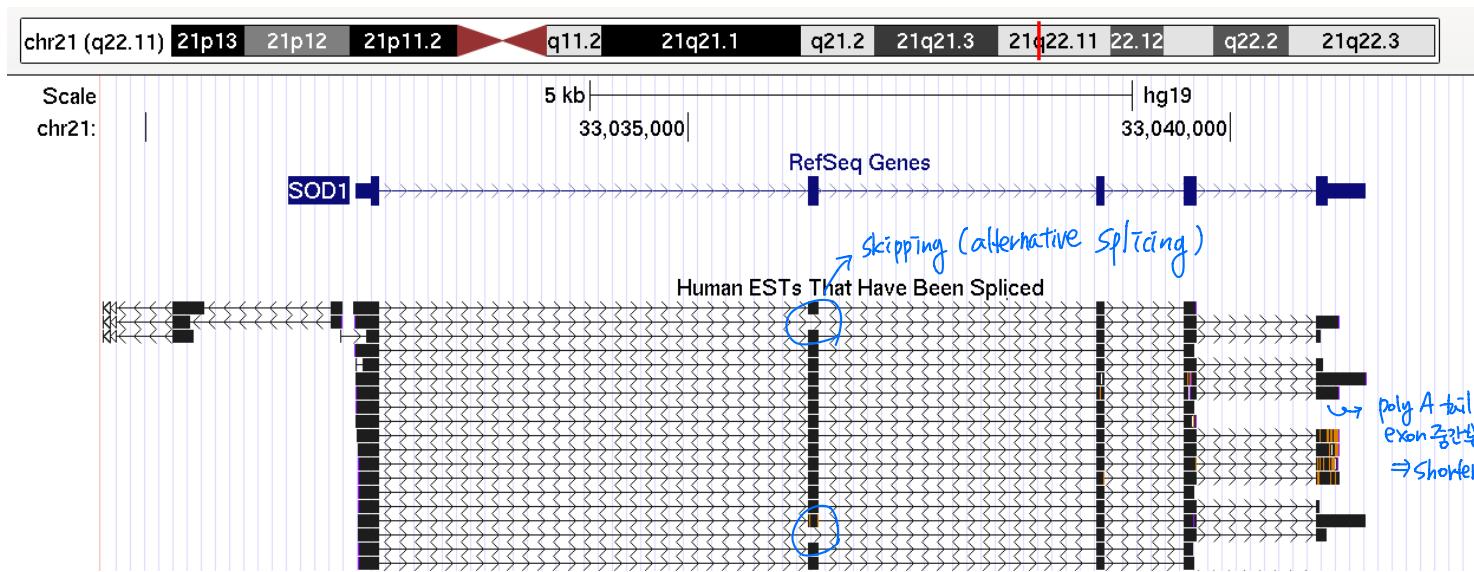
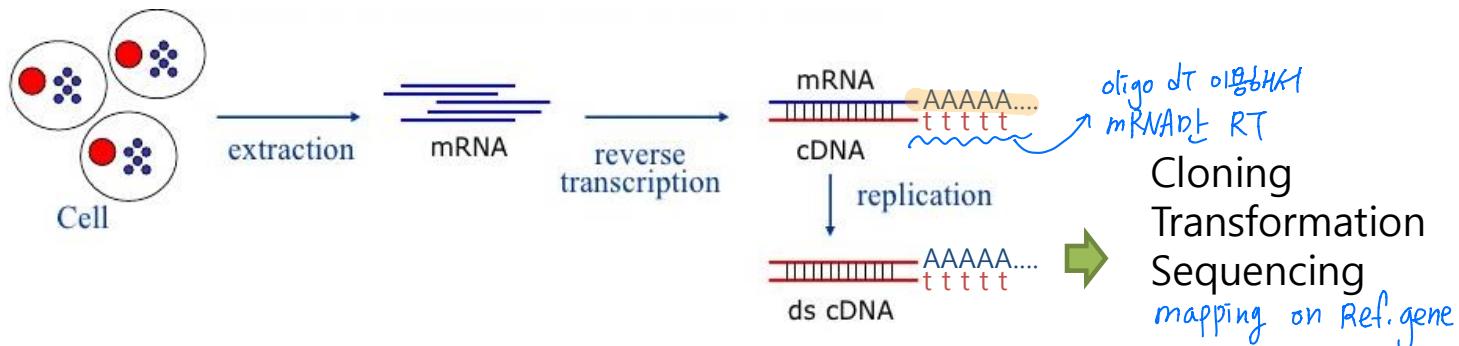
RNA polymerase II

Chromatin structure

mRNA
Ribosome

How to figure out expression? Expressed Sequence Tag (EST)

Gene expression
Genome → mRNAs
How?
EST clones



Human genome: ~50% is repeat sequence

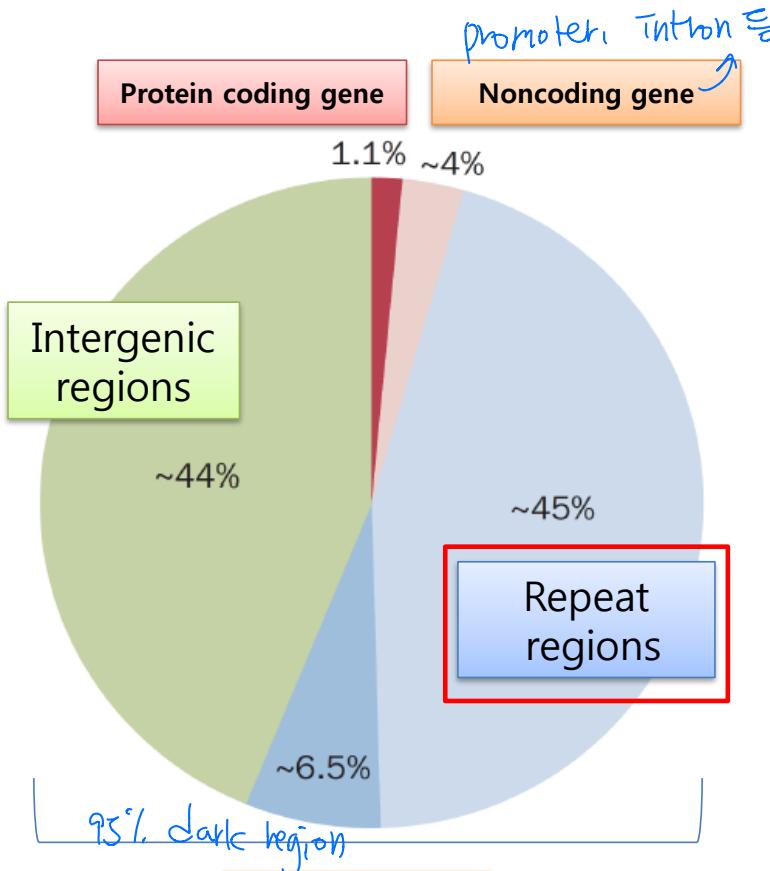


Table 1.2 Transposable elements in the human genome

Element	Estimated number	% of total genome
SINE + LINE	2.4×10^6	33.9
LTR Retrotransposon	0.3×10^6	8.3
Transposons	0.3×10^6	2.8
Total	3.0×10^6	~45

Data from: Bannert, N. & Kurth, R. (2004). Retroelements and the human genome: New perspectives on an old relation. Proc. Natl. Acad. Sci. USA 101, 14572–14579.

~40% of repeat regions is "retrotransposon"

Majority of retrotransposon is LINE or SINE \Rightarrow Some of them are still active \rightarrow can jump around

\rightarrow Cause genomic variation (disease)

~75% of intergenic regions are transcribed !!!



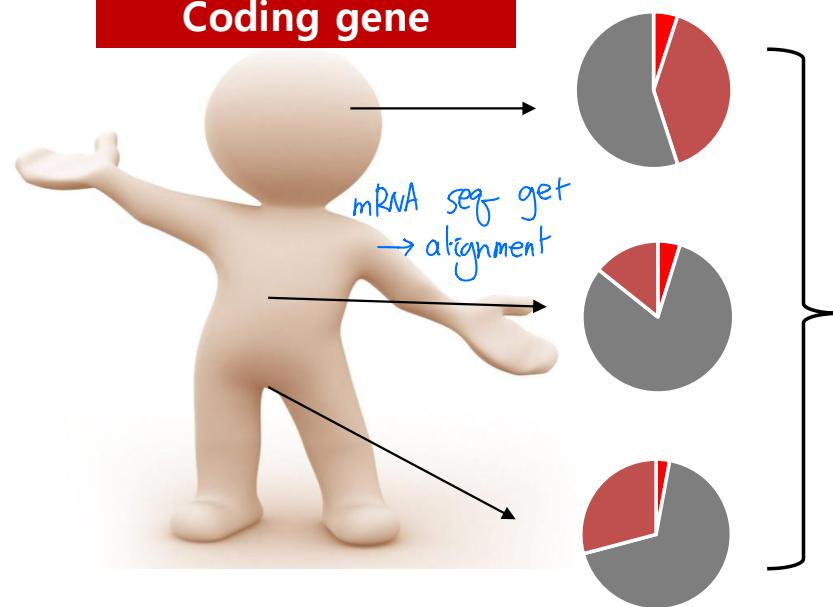
ENCODE (Encyclopedia of DNA Elements) Projects

aim to sequence transcribed gene

functional elements in the human genome
(regulatory elements, non-coding RNAs.....)

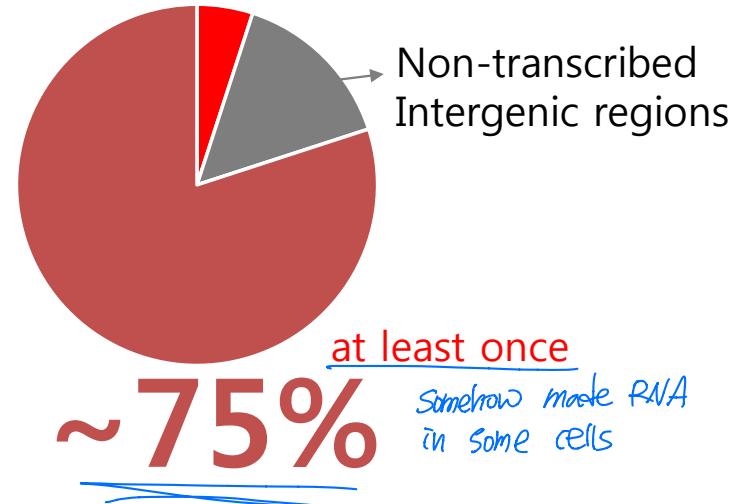
(2012)

Coding gene



Non-coding genes

Human genome



Importance of non-coding genes (RNAs) in human genome

Regulatory function ?

NGS & personal genomics

1990 - 2003

Human genome project

- Sequencing (3 billion bp)
- 13 years, 3 billion \$

2015~ (under develop)



Oxford Nanopore (MinION)



Portable, real-time

~2008

Next-generation sequencing

- 2 month, <1 million \$

Now

NGS

- 1 week , ~1,000 \$

2017, Jan



- 2 days
- ~100 \$
- !!!

Illumina (NovaSeq)

6 TB in 2 days

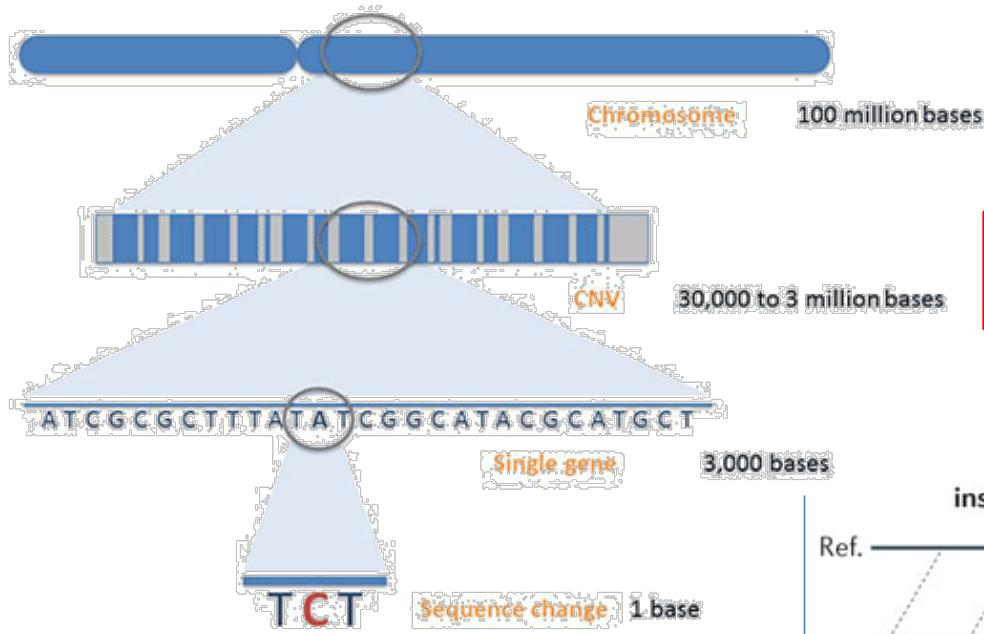
- Personal Genomics



- Single Cell Genomics

질병의 경우 밝히는 동시에
증상을 억제

Genomic Variation



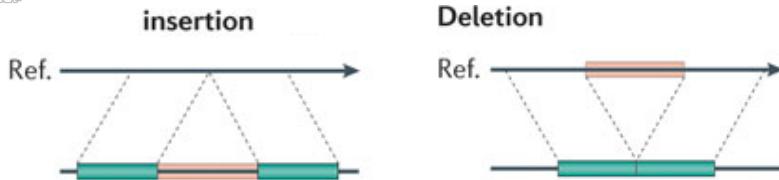
Structural variation

DNA recombination

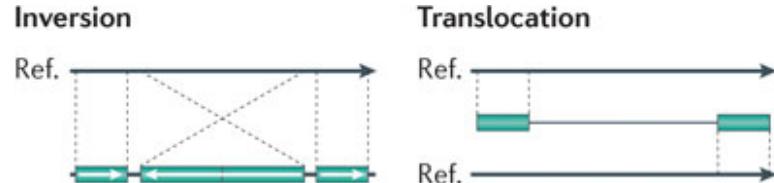
→ meiosis 중인 일정

1. Gene duplication (CNV: Copy number variation)

2. Large insertion / deletion



3. Inversion, Translocation



Sequence variation → phenotype 이 영향

1. Single nucleotide change

: SNP (Single Nucleotide Polymorphism)
↗ mutation 유전적 변화

↗ 인증적 변화

2. Small nucleotide deletion or insertion

: Indel

Genome Projects

: Organization and Objectives

Sequencing은 두 유전자의 거리를 알아내는 법

sequencing은 업로드 시장인 어떤かい genome을 알아내는 법

Sung Wook Chi

Division of Life Sciences, Korea University

What we will learn today

• Genomics

- ## - The Core Aims of Genome Science

• Mapping Genomes

조각내서다 연결

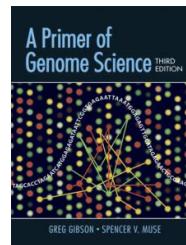
- Genetic Maps, Physical Maps, Cytological Maps
genetic recombination marker → order 현미경통해 Staining을 보고 패턴을 통해 유전자 위치
 - Comparative Genomics

- The Human Genome Project

- #### - Objectives, Internet Resources

Details about human genome project will be on next lecture

figure out location & distance
w/o sequencing



Chapter 1. Genome Projects: Organization and Objectives

Genomics

Genome

Sequencing

phenotype

Mapping

- Identification of genes
- Location of genes
- Relation of genes
- Mapping
- Assembly 조립
- Annotation 주석

장내에 geno이 같은
chromosome가 있는가?
if so, how far?

Variation

- Variation (sequence, structure)
- Comparative genomics probe 등 seq- seq 기법 이용
DNA-DNA hybridization 등 기술 이용
- Variant identification (SNP)
(Translocation)
- Phylogenetic analysis

Functional genomics
Gene expression analysis

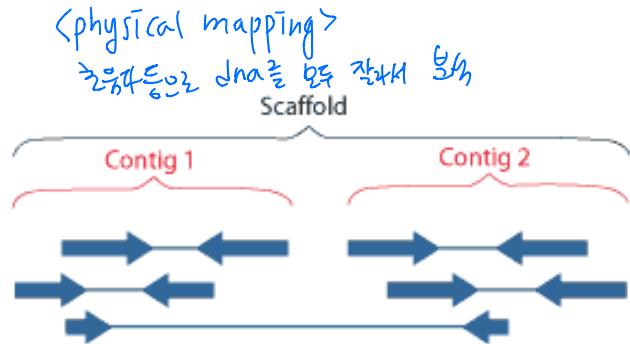
Analysis (Bioinformatics)

The Core Aims of Genome Science

Mapping

1. Assembly of genetic and physical maps

- linkage, recombination frequency
- contig, scaffold



2. Mapping expressed gene sequences in the order of genome

- cDNA clones, ESTs

3. Annotate complete set of genes

- Sequence search, annotation
- ORF (open reading frame)

Variation

mapping 키워드 Variation 정의

4. Sequence variation

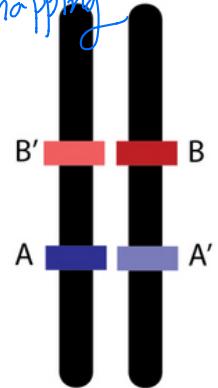
- SNPs, haplotype, linkage disequilibrium

Linked (A and B') (A' and B)
⇒ 연관되는 genetic mapping

For two genes,
each with two
alleles:

A	■
A'	■
B'	■
B	■

The possible
haplotypes
are:
AB
AB'
A'B
A'B'



Phased haplotypes

5. Comparative genomics

- To provide the resources for comparison with other genomes (synteny)

The Core Aims of Genome Science

Functional genomics

6. Functional genomics Screening

- accumulate **functional data**, including biochemical and phenotypic properties of genes
 - : genomics (*near-saturation mutagenesis, high-throughput reverse genetics*)
 - : transcriptomics
 - : proteomics, structural genomics)

7. Compiling gene expression

- SAGE (serial analysis of gene expression)
- transcription profiling (microarray, RNA-Seq)

Analysis

8. Bioinformatics (Computational genomics)

- To establish an **integrated** Web-based **database** and research **interface**

9. Systems biology (Integrated genomics)

- Understanding whole interactions of biological components (genes, proteins.....)
- Understanding emerging properties of network interactions : Network biology

What we will learn today

- **Genomics**

- The Core Aims of Genome Science

- **Mapping Genomes**

- Genetic Maps, Physical Maps, Cytological Maps
- Comparative Genomics

- **The Human Genome Project**

- Objectives, Internet Resources

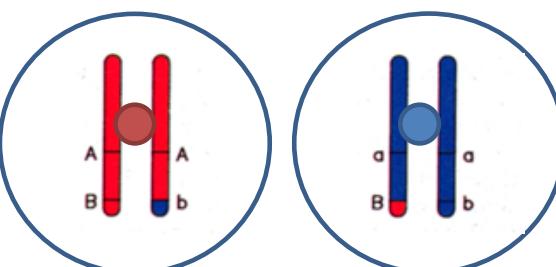
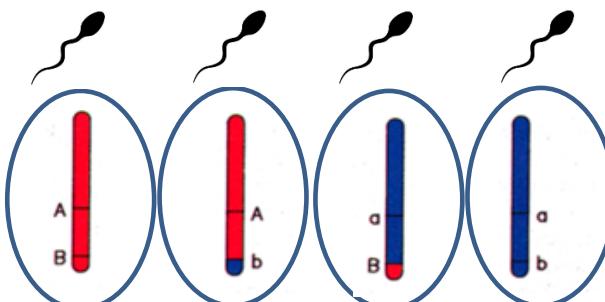
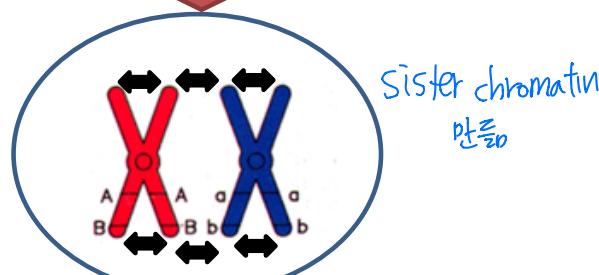
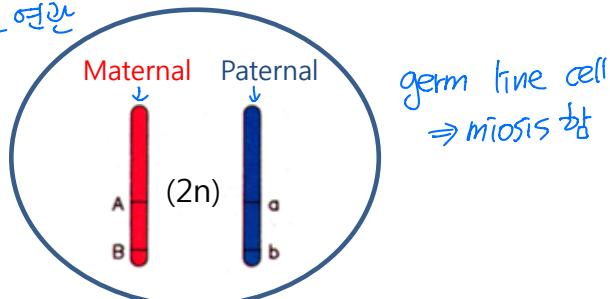
Recombination: Genetic maps

1. location에 따른 유전자 없을 때

① 가교를 거쳐 연관

② 유전자는 거쳐 연관

③ 드립

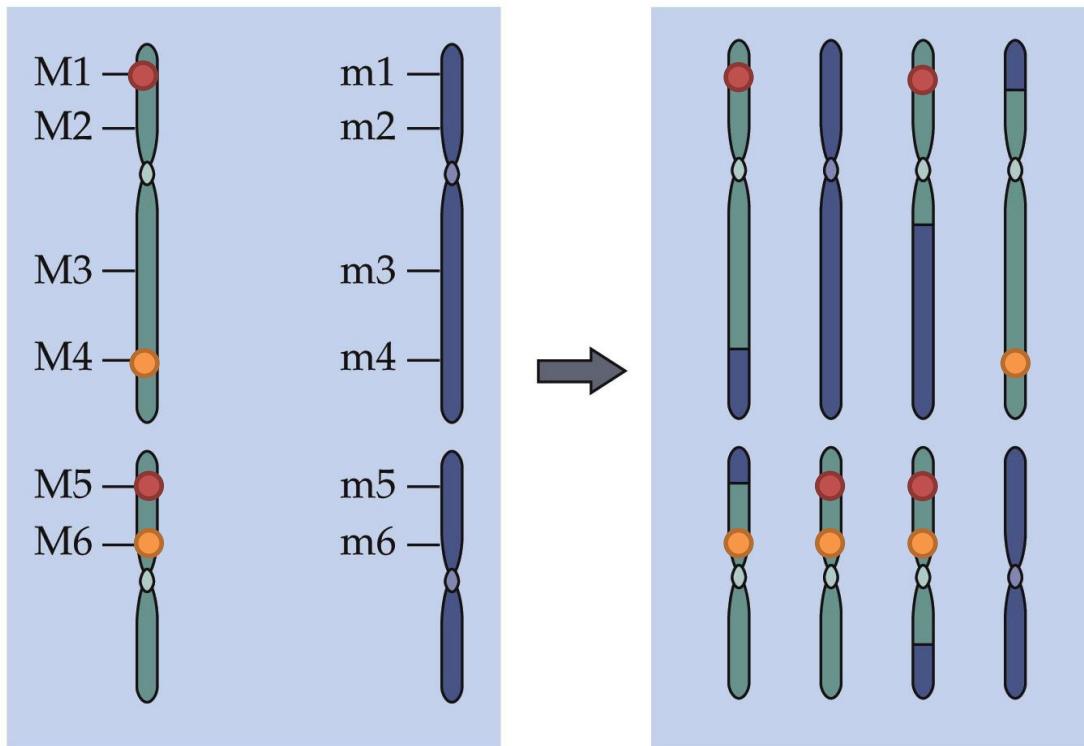


Recombination

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
- centiMorgan (cM) : $1\text{cM} = \text{recombination frequency of } 0.01$

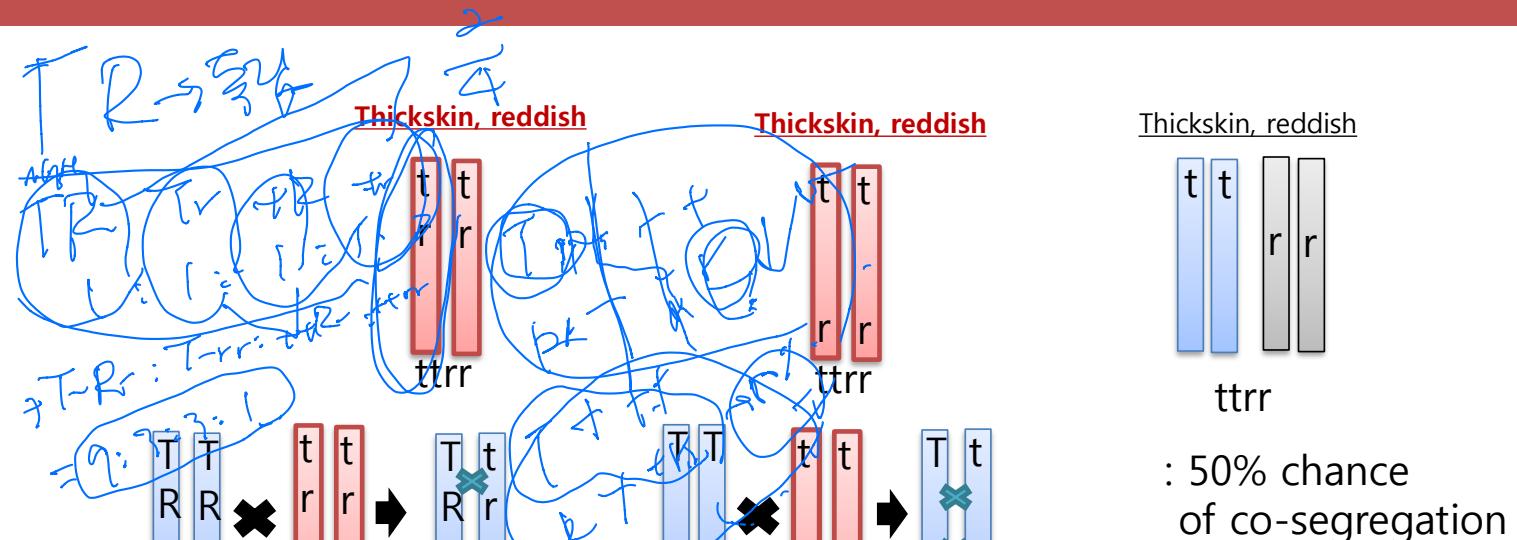


Located in
different
chromosome

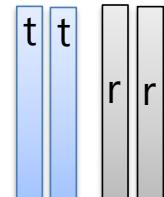
: 50% chance of
co-segregation

52% 경우

Genetic distance



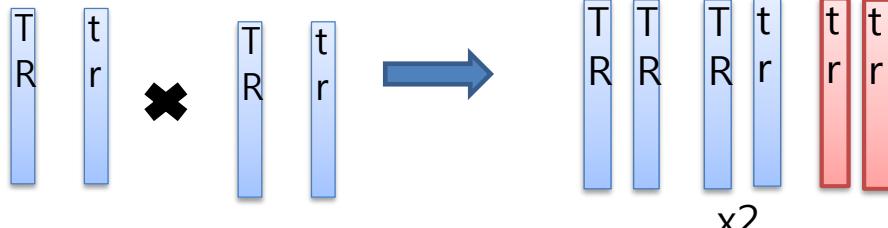
Thickskin, reddish



ttrr

: 50% chance
of co-segregation

TTRR ttrr TtRr TTRR ttrr TtRr



Recombination rate !!!!

Genetic Maps /

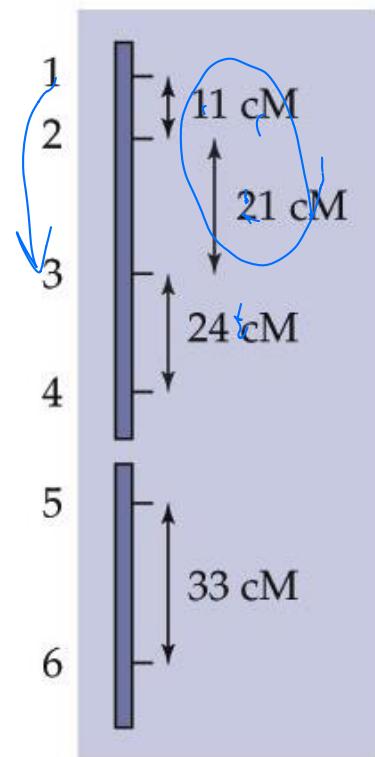
Linkage data -> Mapping function

(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						-

각각의 상관관계 풍선하여 거리를 측정

(C)



50cM – different chromosome

random packing DNA 두 유전자 사이에는 50cM를 두고

2/21 4/21

EXERCISE 1.1 Constructing a genetic map

Suppose that a breeder of orange trees begins to assemble a genetic map based on four recessive loci—thickskin, reddish, sour, and petite—named after the fruit phenotypes of homozygotes. After identifying two true-breeding trees that are either completely wild-type or are mutant for all four loci, the breeder crosses them and plants an orchard of the resultant F_2 trees. Based on the following frequencies of mutant classes, determine which loci are likely to be on the same chromosome and which are the most closely linked.

		Total (n=968)	
Normal	402	Thick	18
Petite	127	Red, sour, and petite	12
Sour	115	Red	11
Thick and red	108	Thick and petite	10
Thick, red, and petite	42	Thick, red, sour, and petite	8
Thick, red, and sour	41	Thick and sour	7
Sour and petite	38	Red and petite	3
Red and sour	24	Thick, sour, and petite	2

Mendelian ratio $\frac{52}{968} \times 100 = 5.3\%$

$$\text{recessive } (1/4) \longrightarrow 968 \times 1/4 = 242$$

WT x Mut

ww mm

wm wm

F_2

ww wm wm mm

Thick

$$108+42+41+18+10+8+7 + 2 = \underline{\underline{236}}$$

$$\text{Reddish} = \underline{\underline{249}}$$

$$\text{Sour} = \underline{\underline{247}}$$

$$\text{Petite} = \underline{\underline{242}}$$

$$\frac{4 \times 3}{2} \Rightarrow 6 \times 8 =$$

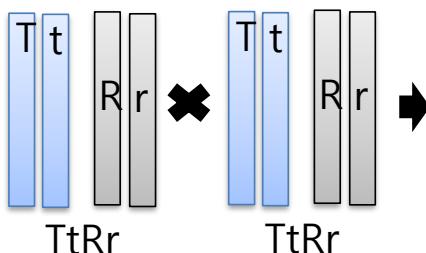
				Total (n=968)
Normal	402	Thick	18	
Petite	127	Red, sour, and petite	12	
Sour	115	Red	11	
Thick and red	108	Thick and petite	10	
Thick, red, and petite	42	Thick, red, sour, and petite	8	
Thick, red, and sour	41	Thick and sour	7	
Sour and petite	38	Red and petite	3	
Red and sour	24	Thick, sour, and petite	2	

B/W two phenotypes (genes)

Unlinked

(Segregate independently)

Thickskin, reddish



$$968 \times 1/4 \times 1/4 = \sim 60$$

B/W two phenotypes (genes)

Unlinked

(Segregate independently)

독립적(독립)인 경우로 해서 면밀히 계산해보면 ⇒ 애를 기운으로 연관성이 있다.

$$968 \times 1/4 \times 1/4 = \sim 60$$

Thick and red	199	Thick and petite	62	Thick and sour	58
Red and petite	65	Red and sour	85	Sour and petite	60

highly linked

likely linked

likely independent
or unlinked



Mapping Genomes : Physical Maps

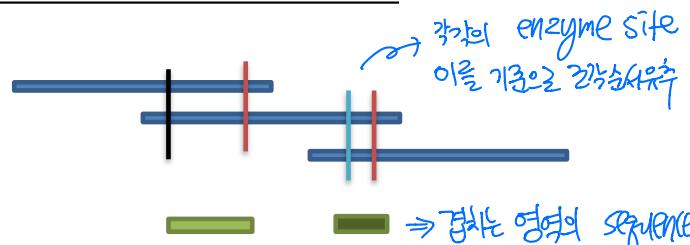
- **Physical Maps**

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

인접한, 접두사
특정 위치 / 이미 알고 있는 sequence

- Assemble contigs

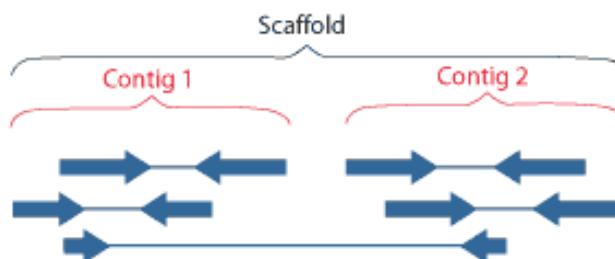
1. Alignment of randomly isolated clones based on shared restriction fragment length profiles



2. Hybridization-based approaches

: chromosome walking

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

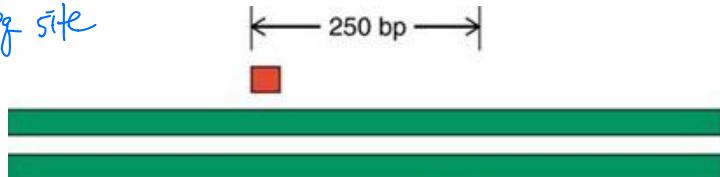


Sequence-Tagged Sites (STSs) 서열표지부위

- Short (200 to 500 base pair) DNA sequence

Can be used as marker of location

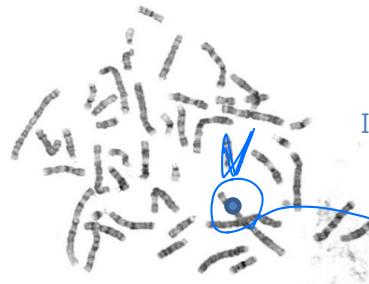
- Marker 역할을 해야하므로 딱 한번만 있는 STS site
- Single occurrence in the genome whose location and base sequence are known.



- Detectable by PCR
- Design short primers

: Hybridize few hundred bp apart

: Amplify a predictable length of DNA



서열표지 혼화학
In situ hybridization
→ 염색체 위치결정

find one signal

① 독특한 블록의 유전체 DNA 시퀀스
약 200 - 500bp 길이로 2종
사기로 PCR

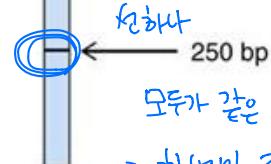
PCR

특정 제한효소를 사용하여
marker로 사용



Electrophoresis
gel에서 내림

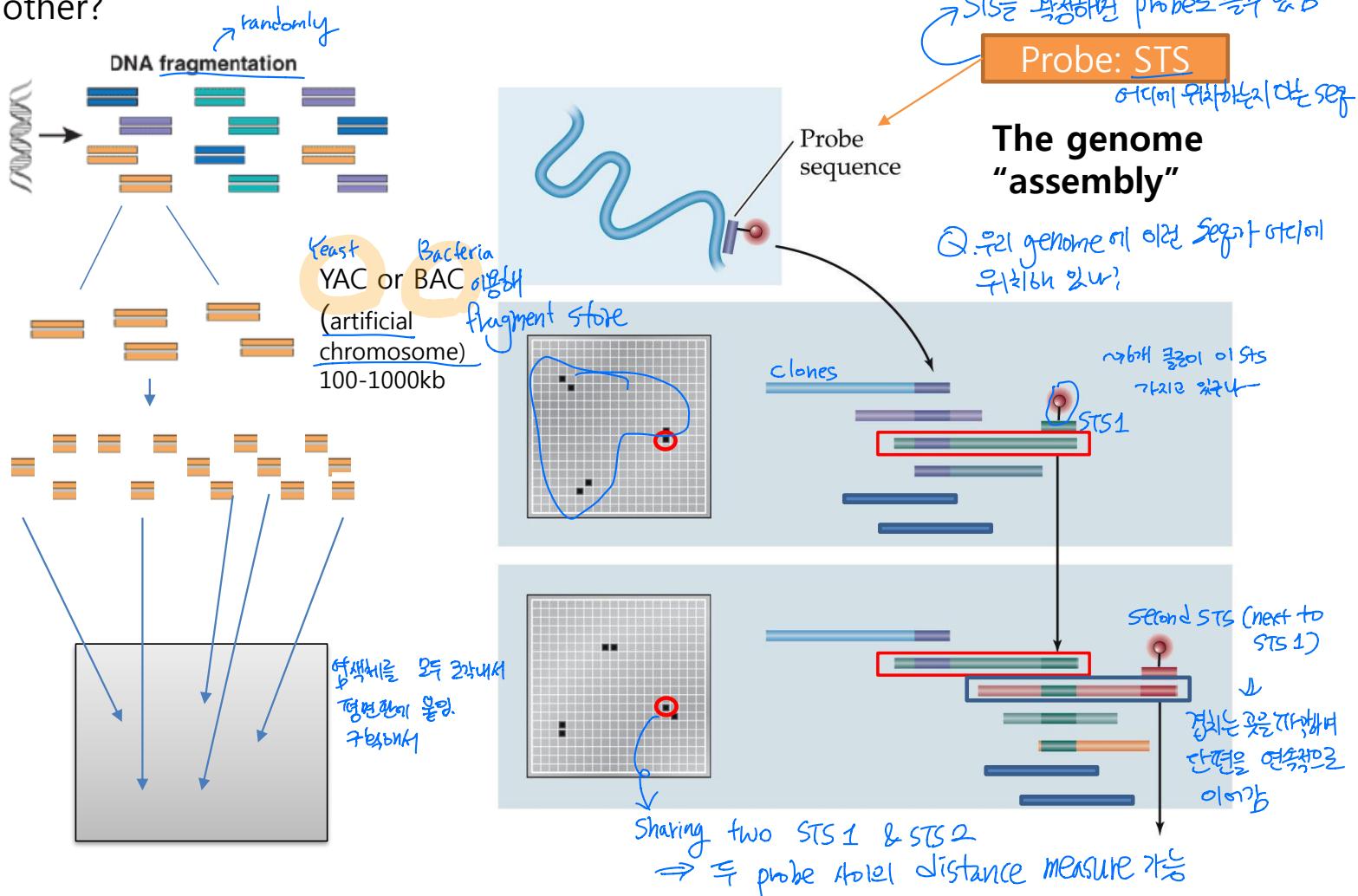
② 노란색



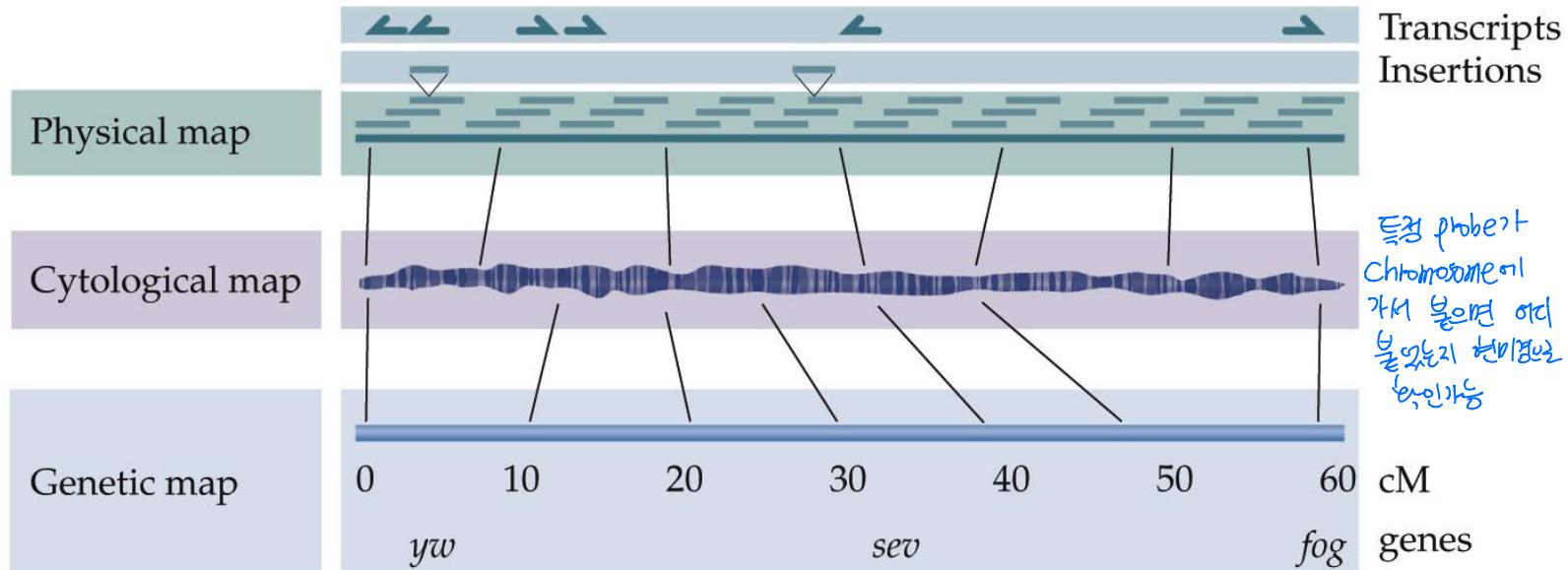
모두가 같은 bp
→ 한번만 존재하는 STS
⇒ probe로 사용

Chromosome Walking

How are individual clones in a genomic library positioned relative to each other?
→ STS를 확장하면 probe



Mapping Genomes :Cytological Map



Physical map, Cytological map and Genetic map are correlated,

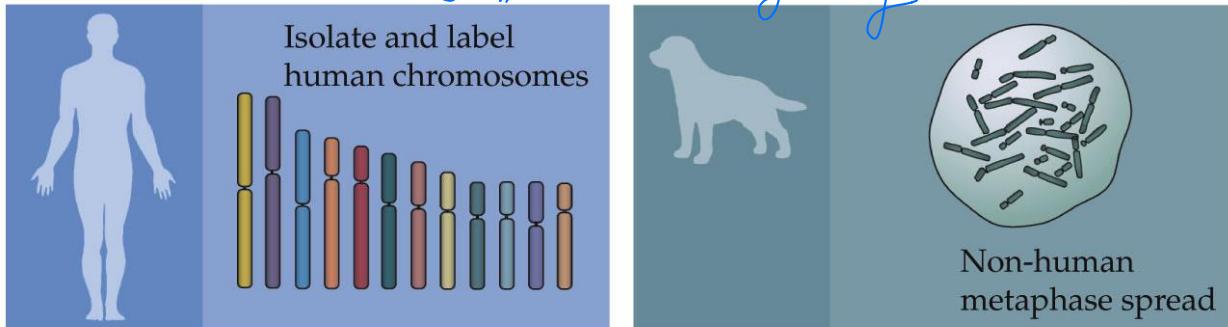
But not the same..

Comparative Genomics

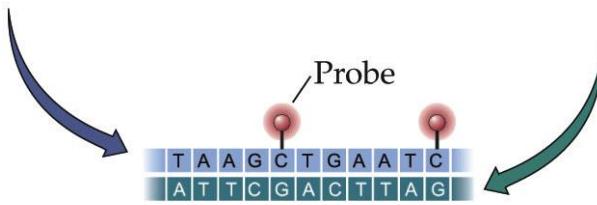
- **Synteny** 두 종 이상의 유전자가 다른 종에 의해 공유되는 동일한 염색체에 있는 경우
: conservation of gene order between chromosome segments of two or more organisms.or Species
*Some kind of genes have same sequence
→ same function*
- **Homologs**
→ comparative genomics only 관심
: orthologs, paralogs
- **Chromosome painting**
 - Rearrangement of gene order: Chromosome fusions and splits, reciprocal translocations and inversions.
 - each chromosome of one species is separately labeled with a set of fluorescent dye

Comparative Genomics: Chromosome painting

각 chr. 대표하는 STS probe를 design 할 수 있음 → 각 chromosome 대표하는 probe를
가의 유전체에 융여됨 ⇒ 연관성 / synteny



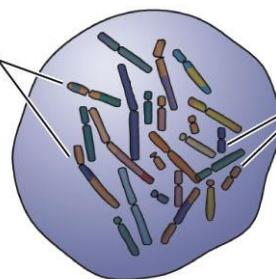
Different color combination for each chromosome probe



Hybridize with fluorescent probe

FISH: fluorescence in situ hybridization

Multicolored chromosomes indicate breakage/fusion events



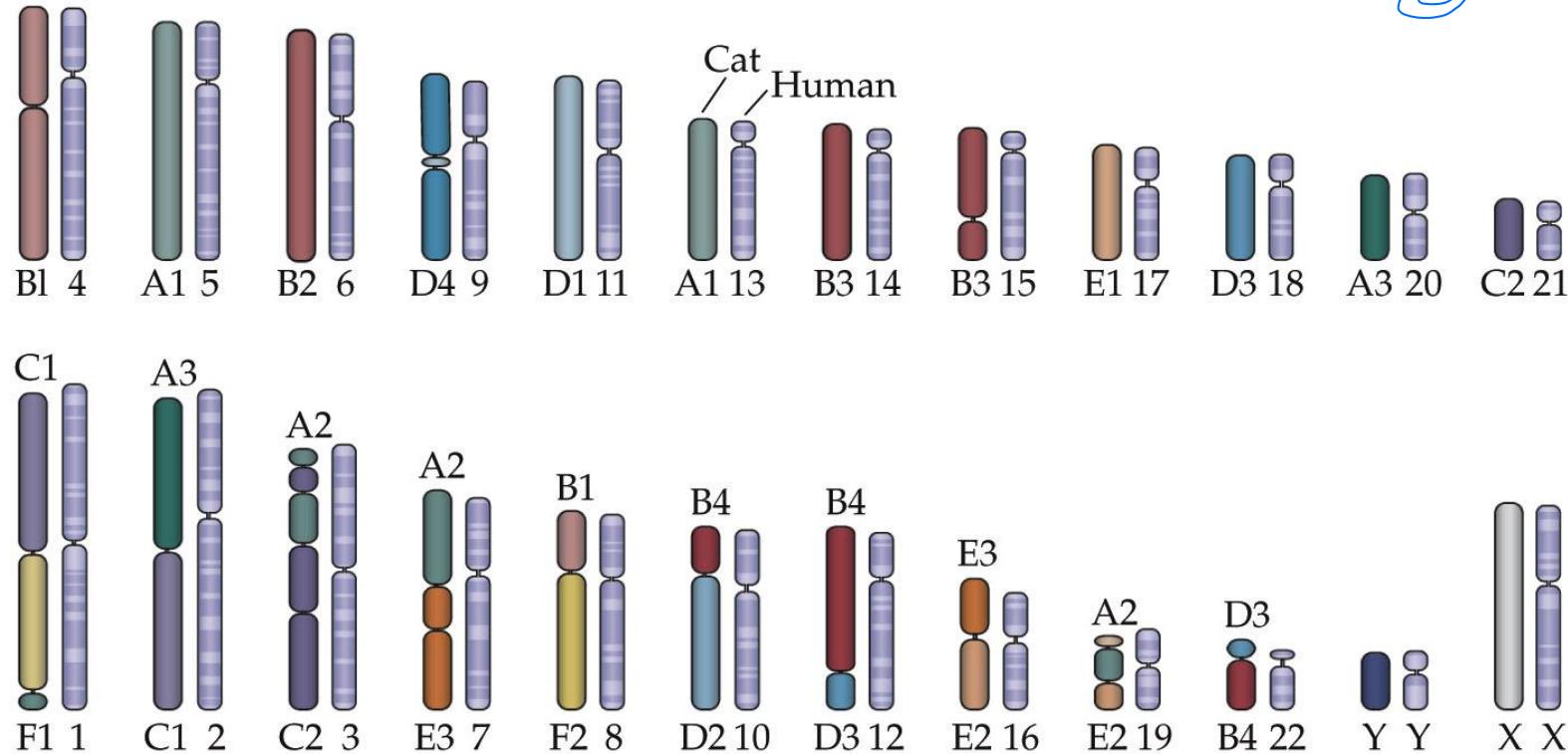
Single-color chromosomes indicate complete correspondence between species

Synteny b/w cat & human genomes

Old school: Chromosome painting (Labelled probe)

→ 가정학 관찰

6



New school: Genome sequence alignment across multi-species

What we will learn today

- **Genomics**

- The Core Aims of Genome Science

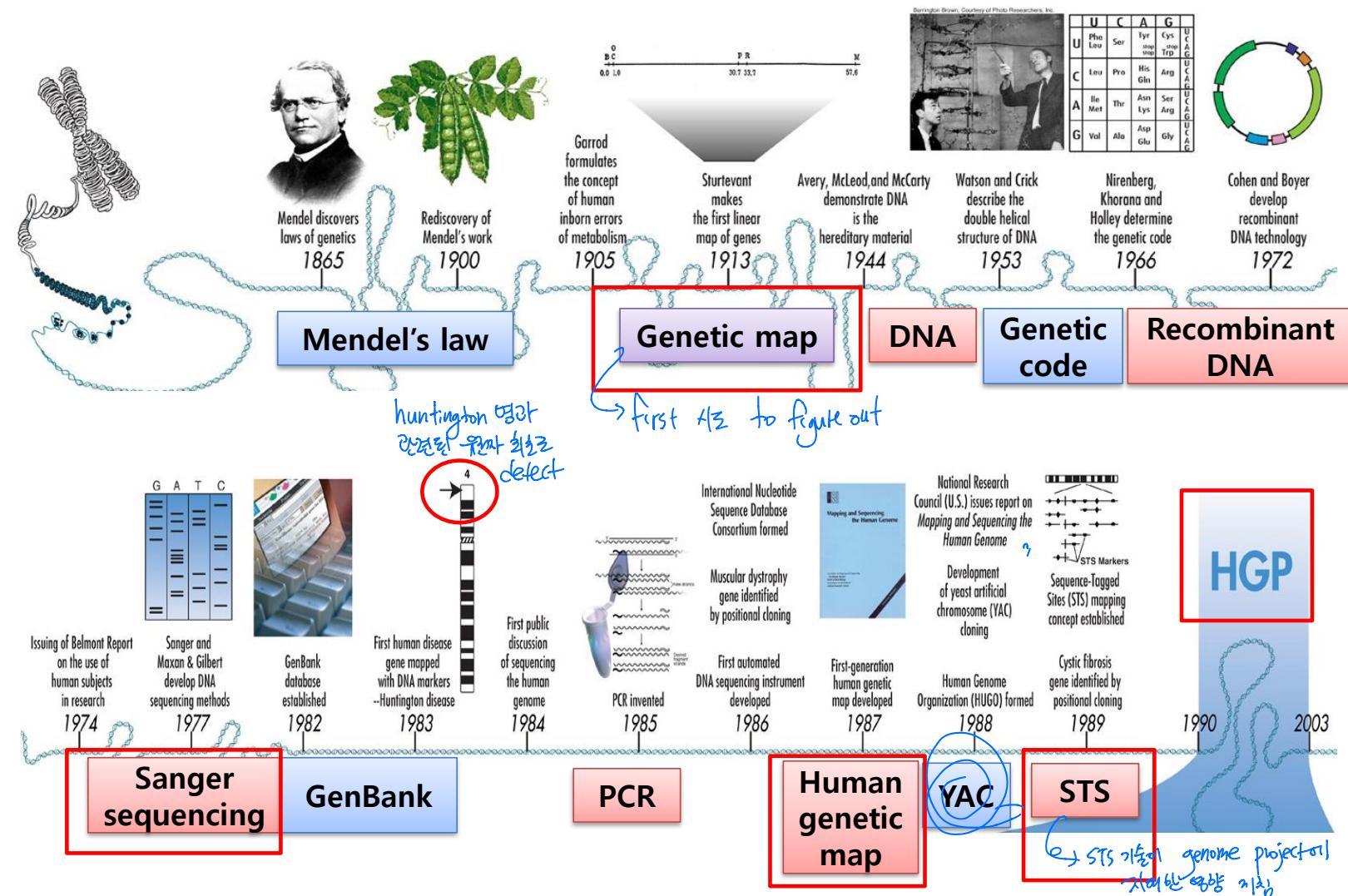
- **Mapping Genomes**

- Genetic Maps, Physical Maps, Cytological Maps
- Comparative Genomics

- **The Human Genome Project**

- Objectives, Internet Resources

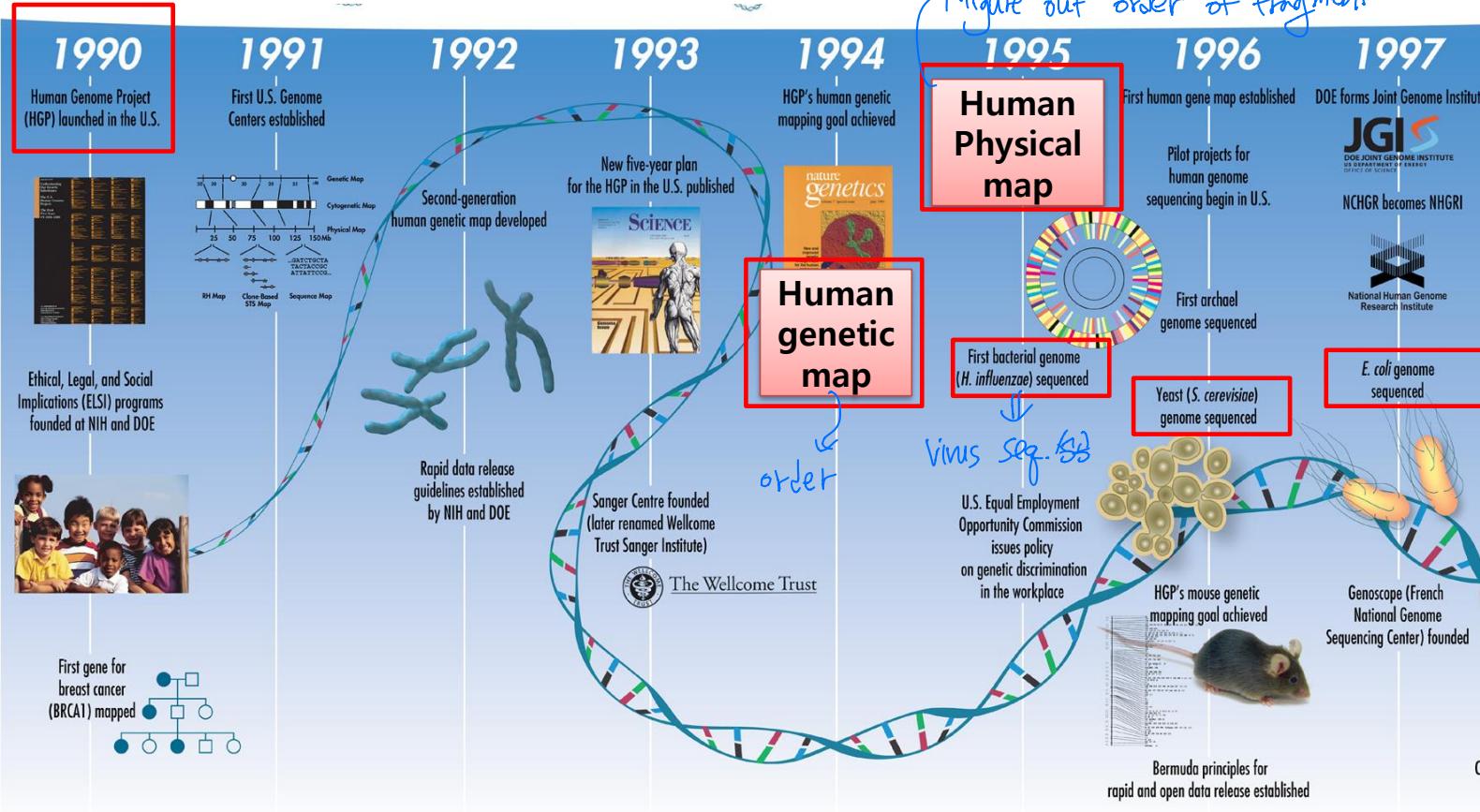
Human Genome Projects : History



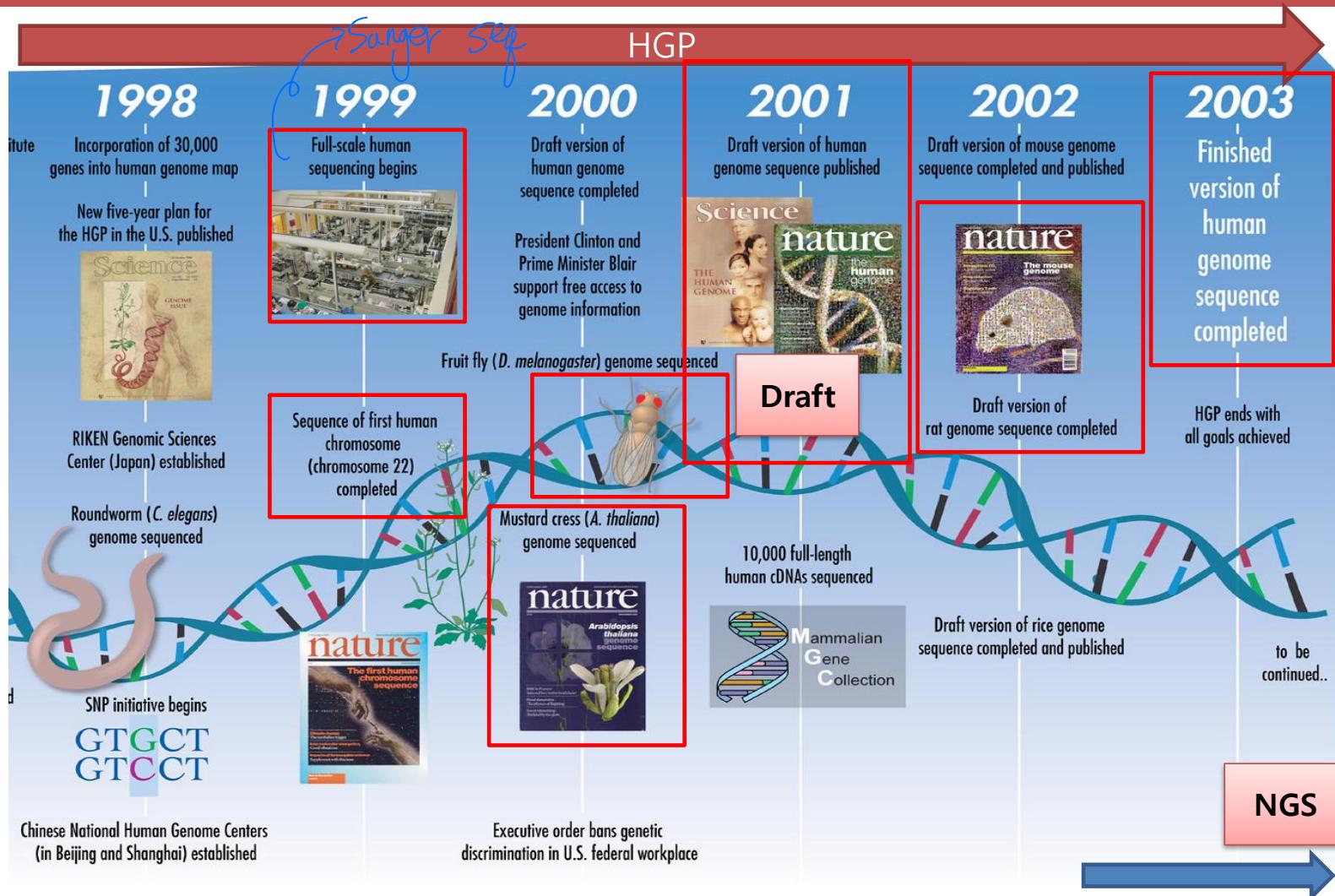
Human Genome Projects : History

HGP

figure out order of fragment



Human Genome Projects : History



The Human Genome Projects : Objectives

TABLE 1.1 Initial Goals of the Human Genome Project (Part 1)

From the First 5-Year Plan: 1993–1998

- ① **1. THE GENETIC MAP**
Complete 2 to 5 cM map by 1995
Develop new technology for rapid and efficient genotyping
- ② **2. THE PHYSICAL MAP**
Complete STS map to 100 kb resolution
- ③ **3. DNA SEQUENCING** *Sanger seq*
Develop approaches to sequence highly interesting regions on Mb scale
Develop technology for automated high throughput sequencing
Attain sequencing capacity of 50 Mb per year; sequence 80 Mb by 1998
- ④ **4. GENE IDENTIFICATION**
Develop efficient methods for gene identification and placement on maps
- 5. TECHNOLOGY DEVELOPMENT**
Substantially expand support for innovative genome technology research
- 6. MODEL ORGANISMS** *to other organisms*
Finish STS map of mouse genome to 300 kb resolution
Obtain complete sequence of biologically interesting regions of mouse genome
Finish sequences of *E. coli* and *S. cerevisiae* genomes
Substantial progress on complete sequencing of *C. elegans* and *D. melanogaster*

The Human Genome Projects :objectives

TABLE 1.1 *Initial Goals of the Human Genome Project* (Part 2)

From the First 5-Year Plan: 1993–1998

7. INFORMATICS

- Continue to create, develop, and operate databases and database tools
- Consolidate, distribute, and develop software for genome projects
- Continue to develop tools for comparison and interpretation of genome information

8. ETHICAL, LEGAL, AND SOCIAL IMPLICATIONS (ELSI)

- Continue to identify and define issues and develop policy options
- Develop and disseminate policy regarding genetic testing
- Foster greater understanding of human genetic variation
- Enhance public and professional education programs on sensitive issues

9. OTHER

- Training of interdisciplinary genome researchers
- Technology transfer into and out of genome centers
- Outreach

How to access genome sequencing data

1. NCBI Genome

The screenshot shows the NCBI Genomes & Maps interface for the species *Homo sapiens*. At the top, there's a search bar with dropdown menus for 'All Databases (Entrez)' and 'for' (with a text input field), and buttons for 'Go' and 'Clear'. To the left, there's a sidebar titled 'Browse your genome' with the sub-instruction 'Click on a chromosome to show' and a dropdown menu set to 'Genes'. Below this is a small graphic of a DNA helix. The main area features a large image of a DNA double helix and text that reads 'Human Genome Resources'.

2. Ensembl

The screenshot shows the Ensembl homepage. The header features the 'Ensembl' logo with 'ASIA' written next to it. A horizontal navigation bar includes links for 'BLAST/BLAT', 'BioMart', 'Tools', 'Downloads', 'Help & Documentation', 'Blog', and 'Mirrors'.

This is a zoomed-in view of the Ensembl search interface. It shows a search bar with 'Search: All species' and a dropdown menu for 'for', followed by a 'Go' button. Below the search bar is an example query: 'e.g. BRCA2 or rat 5:62797383-63627669 or coronary heart disease'.

3. UCSC Genome

The screenshot shows the UCSC Genome Bioinformatics homepage. The header has a yellow background with the text 'UCSC Genome Bioinformatics'. Below the header is a blue navigation bar with links for 'Home', 'Genomes', 'Genome Browser', 'Tools', 'Mirrors', 'Downloads', 'My Data', 'Help', and 'About Us'. On the left side, there's a vertical sidebar with buttons for 'Genome Browser', 'Blat', and 'Table'. The main content area has a section titled 'About the UCSC Genome Bioinformatics Site' with a welcome message. The footer contains a copyright notice: '© 2012 University of California, Santa Cruz'.

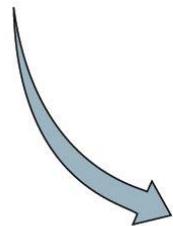
Internet Resources

NCBI
National Center for Biotechnology Information
National Library of Medicine
PubMed All Databases BLAST OMIM Books
Search All Databases for
SITE MAP Alphabetical List Resource Guide
What does NCBI do?
Established in 1988 as a national resource

⑥ Үзүүлэх
үзүүлэх

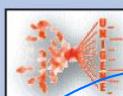
NCBI

(National Center for
Biotechnology
Information)



<http://www.ncbi.nlm.nih.gov>

- ▶ Online Mendelian Inheritance in Man
- ▶ Map viewer
- ▶ Basic Local Alignment Search Tool
- ▶ Links to literature
- ▶ Entrez data-mining tools
- ▶ Cancer Genome Anatomy Project
- ▶ Gene annotation
- ▶ GenBank



Энэ узүүлэхийн төв
disease

Segregation
Хүчинч болуулж
жадалыг

EST
dbSNP

Online Mendelian Inheritance in Man (OMIM) Web site

The screenshot shows the OMIM homepage within a web browser window. The URL in the address bar is <http://www.ncbi.nlm.nih.gov/sites/entrez?db=OMIM>. The page features the NCBI logo and the Johns Hopkins University logo. A search bar at the top has "OMIM" entered. Below the search bar are buttons for "Limits", "Preview/Index", "History", "Clipboard", and "Details". To the left, a sidebar includes links for Entrez, OMIM, Search OMIM, Search Gene Map, and Search Morbid Map. A large blue curved arrow points from the sidebar towards the main content area.

<http://www.ncbi.nlm.nih.gov>

- ▶ Gene Map (cytogenetic location of human Mendelian loci)
- ▶ Genes and Disease
- ▶ FAQs and Statistics
- ▶ Morbidity Map: Alphabetical listing of diseases and corresponding loci
- ▶ Allied Resources

*cytosol information,
disease*



Use genome browser to examine a human disease gene

특정질병관련 유전자를 웹에서 찾기

EXERCISE 1.2 Use the NCBI and Ensembl genome browsers to examine a human disease gene

Choose a human disease of interest to you, and then use the OMIM site to identify a gene that is implicated in the etiology of the disease. Then use the NCBI, UCSC, and Ensembl genome browsers to answer the following questions about the gene:

Asthma / IL13

- a. What are the various identifiers (aliases) for your gene?
- b. Where is the gene located on the chromosome (cytologically and physically)?
- c. What is the Reference Sequence (RefSeq) for the gene?
- d. How many exons are there in the major transcript, and how long is it?
- e. What is known about the function of the gene?
- f. Do the three annotations agree? Which browser do you prefer, and why?

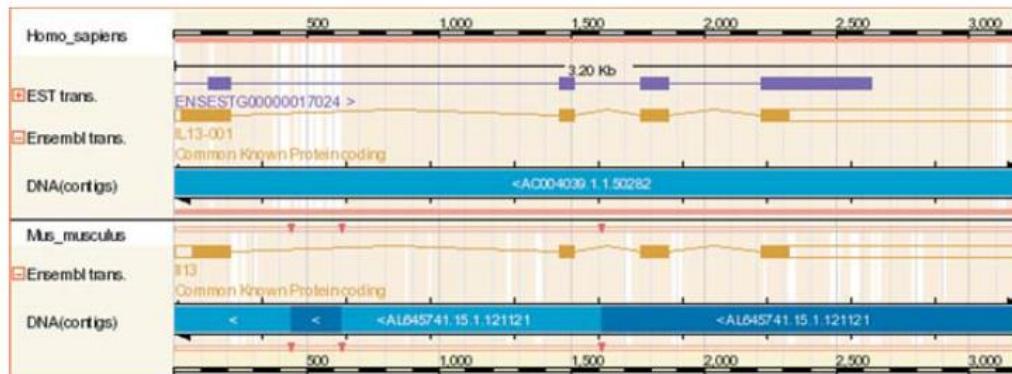
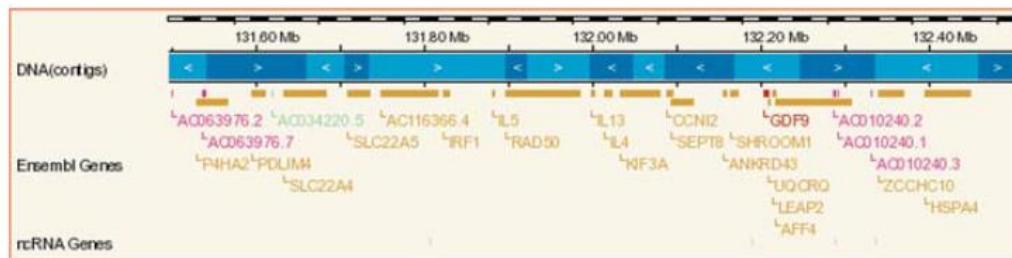
Compare the structure of a gene in a mouse and a human

EXERCISE 1.3 Compare the structure of a gene in a mouse and a human

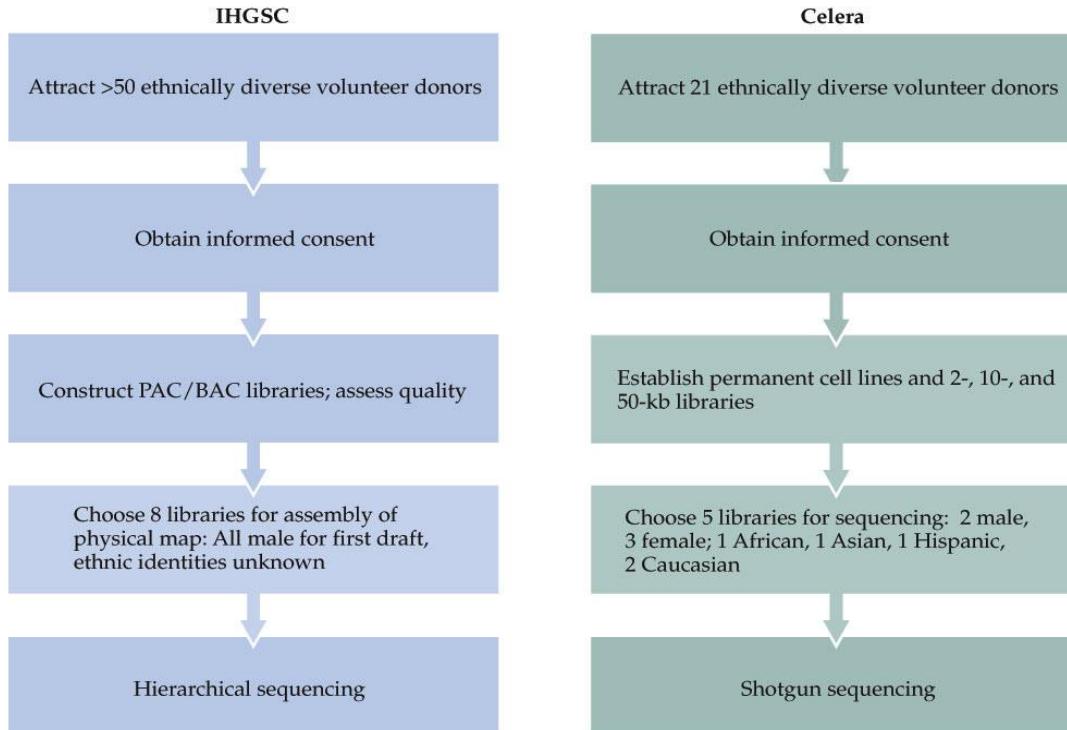
Using either the NCBI or Ensembl Browser, explore the structure of the gene you used in Exercise 1.2 in a mouse and a human (and if possible, in other vertebrates).

IL13

genome 202301K1
IL13

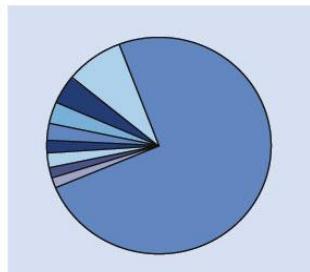


Human genome project race

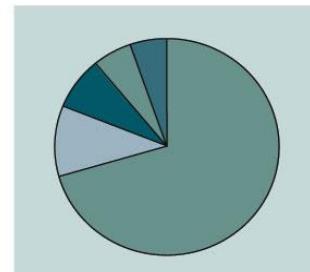


Francis Collins

International
Human Genome
Sequencing
Consortium



Vs.



Craig Venter

Celera

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