

LIST307: Functional Genomics (Part I)

3/13 (Tue) : Introduction to Genomics

3/15 (Thu) : Genome Mapping

3/20 (Tue) : Human Genome Project

3/22 (Thu) : Next-Generation Sequencing

3/27 (Tue) : NGS & Sequence Alignment

3/29 (Thu) : Haplotyping & SNPs

4/3 (Tue) : Variation and Whole Genome Sequencing

4/5 (Thu) : Genome-wide association study

4/10 (Tue) : GWAS, eQTL, Phylogenetic analysis

4/12 (Thu) : Basics of Functional Genomics

4/17 (Tue) : Exome-Seq & Comparative Genomics

4/19 (Thu) : Review of functional genomics (part I)

Middle-term exam: 4/24 (Thu) or 4/26 (Tue) ?

Summary & More

: GWAS , eQTL

GWAS : Variation - Phenotype - Expression

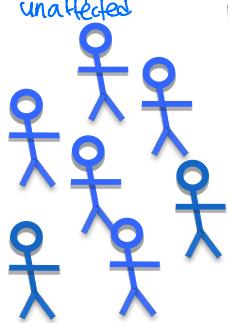
The diagram illustrates the components of GWAS. It features three main labels: "Variation" (in blue), "Phenotype" (in red), and "Expression" (in blue). "Variation" is further divided into "all kind of" (above) and "(genotype)" (below). A curved arrow points from "Variation" up towards "Phenotype". Below "Phenotype", a dashed arrow points from "Expression" back to "Variation", labeled "How affect?". Above "Phenotype", the word "which are associated" is written in blue. The label "eQTL" is placed above "Expression".

Sung Wook Chi

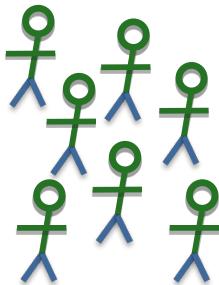
Division of Life Sciences, Korea University

GWAS

Control
Population
unaffected



Disease
Population
affected



SNP chip

disease와 관련된
유전 척도 유전체



detect Tag SNP \Rightarrow enough to figure out region
NGS sequencing

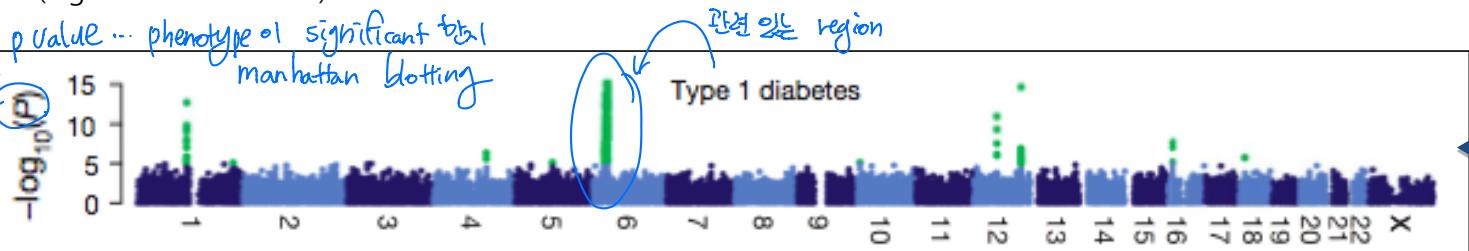


P-value

$$LOD = - \log_{10} \frac{\text{probability of birth sequence with a given linkage value (observed)}}{\text{probability of birth sequence with no linkage (expected)}}$$

(logarithm of the odds)

p value ... phenotype of significant trait
manhattan blotting



GWAS: Practice

sample	seq	type	grade
1	actgtacattc	normal	0
2	actgtacattc	normal	0
3	actgtacattc	normal	0
4	actgtacattc	normal	0
5	actgtacattc	normal	0
6	actgtacattc	tumor	2
7	actgtacattc	tumor	3
8	actctactta	tumor	2
9	actgtacattc	tumor	3
10	actctactta	tumor	4

	a	c
Normal	0.2 + D	0.3 - D
Tumor	0.2 - D	0.3 + D

$$D = 0.3 - 0.2 = 0.1 > 0$$

$$D_{\max} = \min(0.3, 0.2) = 0.2$$

$$D' = 0.1/0.2 = 0.5 \quad (?) \quad \text{significant 힘 있는 증거}$$

Chi square test (p=0.067889) : significant 힘 있는 증거

Permutation
(False-discovery rate)

observed	frequency		Total
	a	c	
Normal	3	2	5
Tumor	1	4	5
Total	4	6	10

observed	probability		
	a	c	
Normal	0.3	0.2	
Tumor	0.1	0.4	little increased → associated with significant

Expected

$P_n = \frac{\text{Normal Seq}}{5} = 0.5$	$P_a = 4/10 = 0.4$
$P_t = 5/10 = 0.5$	$P_c = 6/10 = 0.6$

Expected	a	c
Normal	0.5×0.4	0.5×0.6
Tumor	0.5×0.4	0.5×0.6

Expected	a	c
Normal	0.2	0.3
Tumor	0.2	0.3

The Genome-wide Association Study (GWAS)

① Design my set first

Affected (Disease)

affected

3 groups of id 250 each 1000 total

WGS (Variation)

→ based on population & region we see
→ pattern different

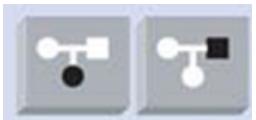
Case-control designs



Disease

More defined case

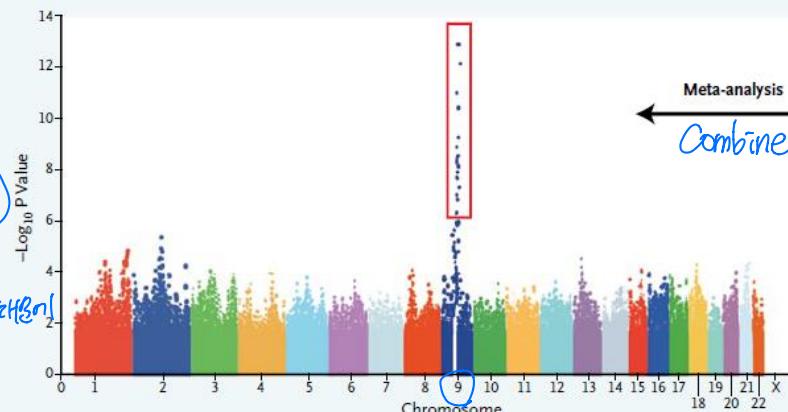
Family based designs (Trio)



Quantitative trait loci (QTLs)

Statistical significance

JY
mutation rates! p values! 3%!



but often not linked

Samples of GWAS

Case-control designs (Disease vs Normal)

easy way to grouping



- A difference in the frequency of an allele or genotype of the polymorphism under test between the two groups indicates that the genetic marker may increase risk of the disease or likelihood of the trait, or be in linkage disequilibrium with a polymorphism which does.
 - Haplotypes can also show association with a disease or trait.
- Problem: genotype and haplotype frequencies vary between ethnic or geographic populations, possibly due to different ancestry.
- > population stratification, population structure

↳ Different Sampling 대상의 발생: background frequency와 달라짐 (인종, 지역 등)

Family based designs (Trio)



↳ because all kind of SNPs are inherited from parent

안암병원 → Collect patient & their genotype

Cohort: a group of subjects with a common defining characteristic

Tracks two or more groups forward from exposure to outcome

limiting area
and relate with genotype

Population stratification

Population 1

$$\text{freq } (A) = 0.8$$

$$\text{freq } (a) = 0.2$$

Diabetics = 20% over rated

$$0.8/(0.8+0.3)$$

$$P(A) \text{ diabetic} = (0.73 \times 20\%) + 0.27 \times 10\% = \underline{\underline{17.3\%}}$$

$$P(a) \text{ diabetic} = (0.22 \times 20\%) + 0.78 \times 10\% = \underline{\underline{12.2\%}}$$

Population 2 combined case

$$\text{freq } (A) = 0.3$$

$$\text{freq } (a) = 0.7$$

Diabetics = 10% under rated

Stratification may be environmental, cultural, or genetic

-Problem: genotype and haplotype frequencies vary between ethnic or geographic populations, possibly due to different ancestry.

-> population stratification

EXERCISE 3.3 Perform a case-control association test

Two parallel association studies between a candidate gene and skin cancer are performed in New York City and Miami, Florida. The numbers of cases and controls who are homozygous for a 3-bp insertion or deletion polymorphism or who are heterozygous are as follows:

Q. Is this deletion related to skin cancer?

New York

	Insertion	Heterozygous	Deletion
w/ cancer	Case	221	198
w/o cancer	Control	279	212

Miami

	Insertion	Heterozygous	Deletion
w/ cancer	Case	171	236
w/o cancer	Control	189	244

Calculate the probability of association between the deletion polymorphism and the disease for each population separately and for the combined population, using the formula:

Chi square test

$$\chi^2 = \frac{N [N(r_1 + 2r_2) - R(n_1 + 2n_2)]^2}{RS[N(n_1 + 4n_2) - (n_1 + 2n_2)^2]} \rightarrow \text{식을 알겠는 } X$$

Discuss the possible reasons for any discrepancy in the conclusions from the two studies.

N=1,000

New York**observed**P(in)
=0.5

	474	Insertion	Heterozygous	Deletion
Case	221	198	55	490
Control	279	212	35	

Miami

	Insertion	Heterozygous	Deletion
Case	171	236	83
Control	189	244	77

P(in)
=0.36

0.45%

$$\chi^2 = N [N(r_1 + 2r_2) - R(n_1 + 2n_2)]^2 / RS[N(n_1 + 4n_2) - (n_1 + 2n_2)^2]$$

Quantity	New York	Miami	Combined
r ₁	198	236	434
r ₂	55	83	138
n ₁	410	480	890
n ₂	90	160	250
N	1,000	1,000	2,000
R	474	490	964
S	526	510	1,036
	7.6	0.8	6.9
P<0.01		↓ Deletion of Cancer 연관성 있다고 하기 어려움	P<0.01

Plugging these values into the formula results in χ^2 values of 7.6, 0.8, and 6.9 for New York, Miami, and the combined population, respectively. Prima facie this suggests little effect of the polymorphism in Miami, but suggests a significant effect at $p < 0.01$ in New York and across the two populations.

$$P(\text{in}) = ((221+279)/1000) = 0.5$$

$$P(\text{het}) = ((198+212)/1000) = 0.41$$

$$P(\text{del}) = ((55+35)/1000) = 0.09$$

$$P(\text{case}) = ((221+198+55)/1000) = 0.474$$

$$P(\text{cont}) = ((279+212+35)/1000) = 0.526$$

Expected

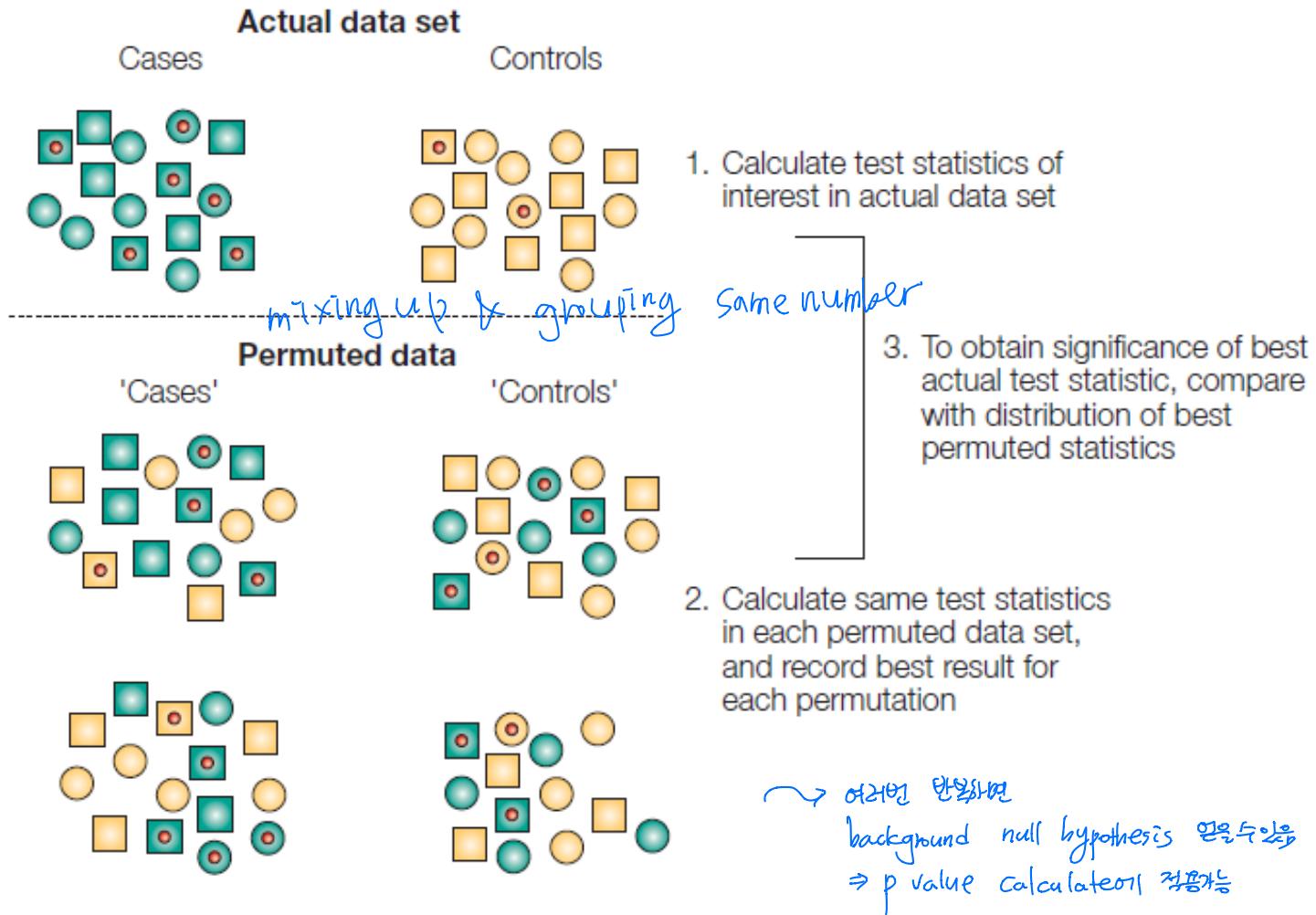
P(case)xP(in)	P(case)xP(het)	P(case)xP(del)
P(cont)xP(in)	P(cont)xP(het)	P(cont)xP(del)
237	194	42
263	216	47

$$\chi^2 = \sum \frac{(O - E)^2}{E}$$

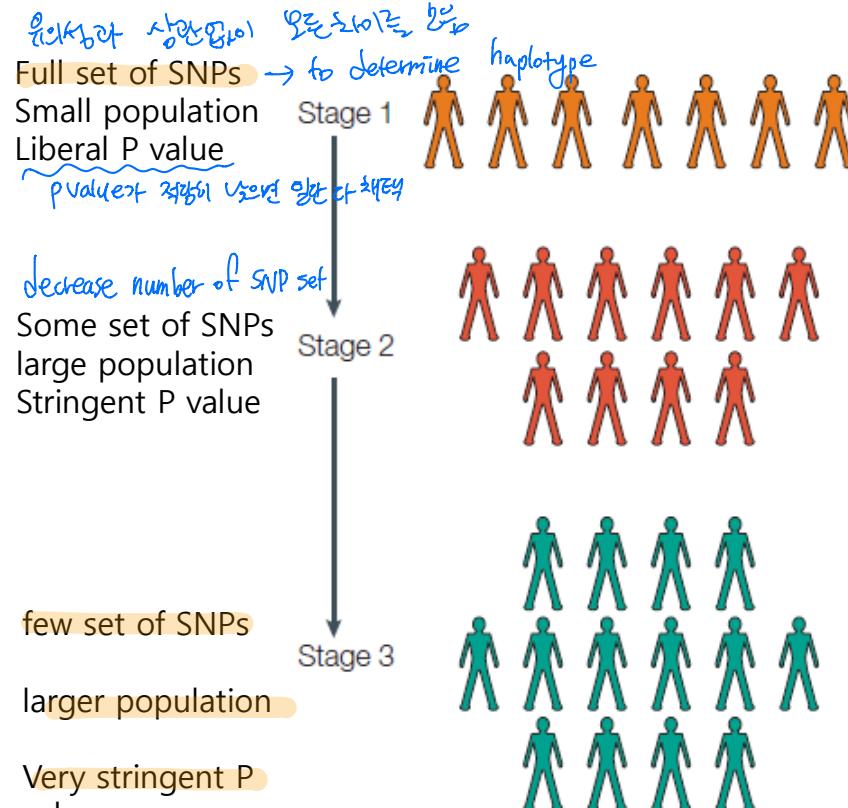
$$\chi^2 = \text{the test statistic} \quad \sum = \text{the sum of}$$

O = Observed frequencies E = Expected frequencies

Permutation testing



Using a multistage approach to minimize sample sizes.



Number of SNPs

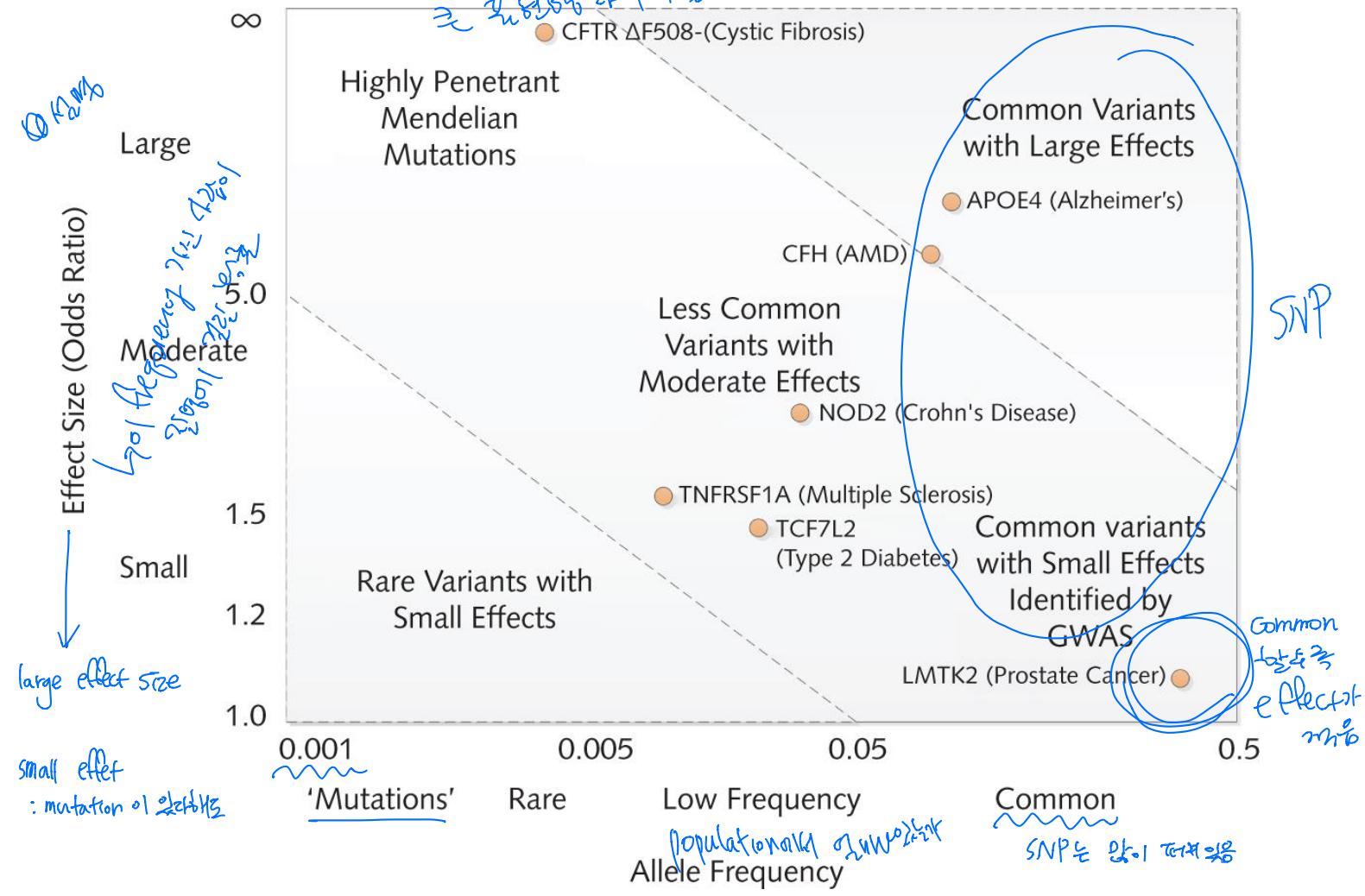
Genotype full set of SNPs in relatively small population at liberal p value

Screen second, larger population at more stringent p value

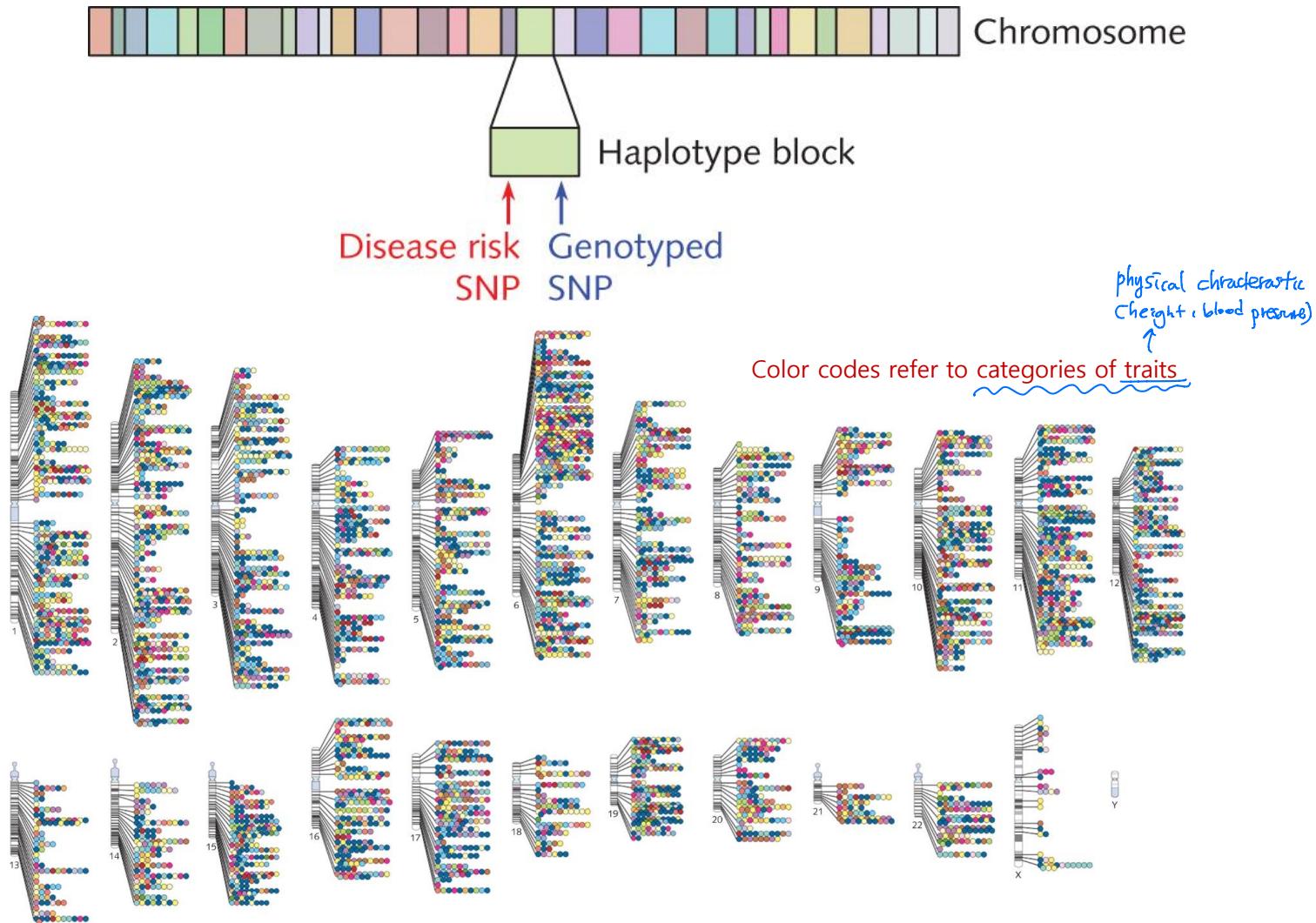
Optional third stage for increased stringency

→ more focusing on the SNP associated with disease we are interested

Allele frequency (variation rate) vs. effect size (penetrance)



GWAS : loci associated with traits



Ansan Cohort (GWAS in Korea)

Korea Association REsource (KARE)

NGS가 많은 데이터를 SNP chip을 이용해 SNP tracking



	Baseline	1 st follow-up	2 nd follow-up	3 rd follow-up	4 th follow-up	5 th follow-up
Ansung	5,018	4,580 (91.3)	3,975 (79.2)	3,434 (68.4)	3,403 (67.8)	3,186 (63.5)
Ansan	5,020	4,023 (80.1)	3,540 (70.5)	3,255 (64.8)	3,262 (65.0)	3,052 (60.8)
Total	10,038	8,601 (85.7)	7,594 (74.9)	7,594 (66.6)	6,688 (66.4)	6,688 (62.1)

- ❖ Participants (follow-up rate)
- ❖ Two cohorts: Ansung (rural area) & Ansan (urban area)

height, blood pressure = all kind of data following Study

A large-scale genome-wide association study of Asian populations uncovers genetic factors influencing eight quantitative traits

Yoon Shin Cho¹, Min Jin Go¹, Young Jin Kim¹, Jee Yeon Heo¹, Ji Hee Oh¹, Hyo-Jeong Ban¹, Dankyu Yoon², Mi Hee Lee¹, Dong-Joon Kim¹, Miey Park¹, Seung-Hun Cha¹, Jun-Woo Kim¹, Bok-Ghee Han¹, Haesook Min¹, Younghin Ahn¹, Man Suk Park¹, Hye Ree Han¹, Hye-Yoon Jang³, Eun Young Cho³, Jong-Eun Lee³, Nam H Cho⁴, Chol Shin⁵, Taesung Park^{2,6}, Ji Wan Park⁷, Jong-Keuk Lee⁸, Lon Cardon⁹, Geraldine Clarke¹⁰, Mark I McCarthy^{10,11}, Jong-Young Lee¹, Jong-Koo Lee¹², Bermseok Oh^{1,13} & Hyung-Lae Kim¹

To identify genetic factors influencing quantitative traits of biomedical importance, we conducted a genome-wide association study in 8,842 samples from population-based cohorts recruited in Korea. For height and body mass index, most variants detected overlapped those reported in European samples. For the other traits examined, replication of promising GWAS signals in 7,861 independent Korean samples identified six previously unknown loci. For pulse rate, signals reaching genome-wide significance mapped to chromosomes 1q32 (rs12731740, $P = 2.9 \times 10^{-9}$) and 6q22 (rs12110693, $P = 1.6 \times 10^{-9}$), with the latter ~ 400 kb from the coding sequence of *GJA1*. For systolic blood pressure, the most compelling association involved chromosome 12q21 and variants near the *ATP2B1* gene (rs17249754, $P = 1.3 \times 10^{-7}$). For waist-hip ratio, variants on chromosome 12q24 (rs2074356, $P = 7.8 \times 10^{-12}$) showed convincing associations, although no regional transcript has strong biological candidacy. Finally, we identified two loci influencing bone mineral density at multiple sites. On chromosome 7q31, rs7776725 (within the *FAM3C* gene) was associated with bone density at the radius ($P = 1.0 \times 10^{-11}$), tibia ($P = 1.6 \times 10^{-6}$) and heel ($P = 1.9 \times 10^{-10}$). On chromosome 7p14, rs1721400 (mapping close to *SFRP4*, a frizzled protein gene) showed consistent associations at the same three sites ($P = 2.2 \times 10^{-3}$, $P = 1.4 \times 10^{-7}$ and $P = 6.0 \times 10^{-4}$, respectively). This large-scale GWA analysis of well-characterized Korean population-based samples highlights previously unknown biological pathways.

Positional cloning of a candidate complex disease

(A) low resolution

Linkage mapping	7 cM at Chr 2q37
Construct physical contig	~1.7 Mb, 22 genes
Low resolution case-control association scan	~100 kb, 3 genes
Sequence 10 individuals over 66 kb	179 variants, including 161 SNPs
Strong association with SNP 43	Calpain 10, Intron 3
Haplotype tests → complexity	3 SNP heterozygotes
Population attributable risk	4% Northern Europeans 14% Mexican Americans

noninsulin-dependent diabetes mellitus

large chunk of genetic region

narrow down region

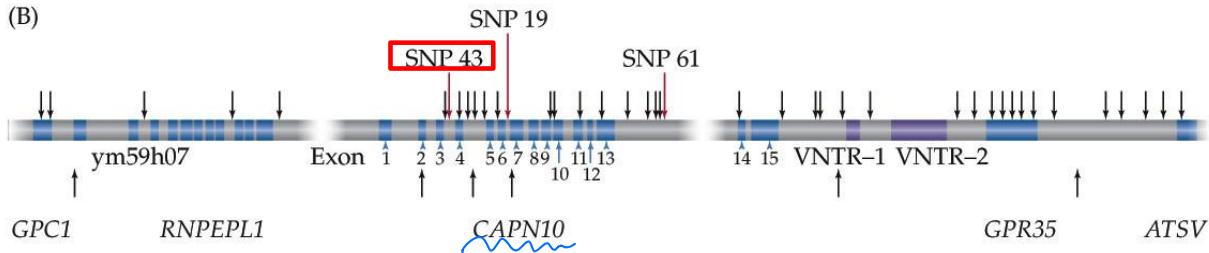
positional cloning
: classical way

SNP affects binding to nuclear factor in human pancreatic extracts
and transcriptional activation of a reporter gene in a cell line

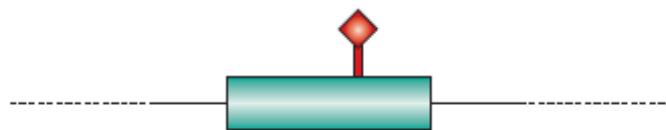
Hypothesis: Calpain 10 protease involvement in NIDDM

recent way to
identify SNPs - disease

(B)

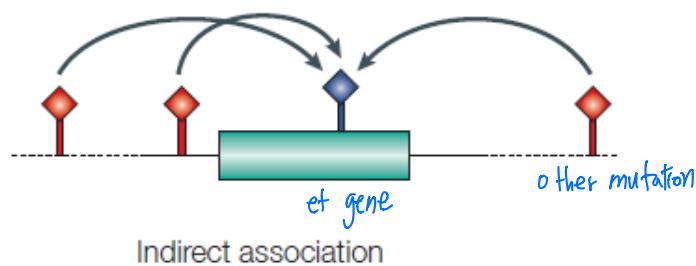


Testing SNPs for association by direct and indirect method



Exon region of mutation of SNP
Alteration of amino acids
Change in protein stability & function

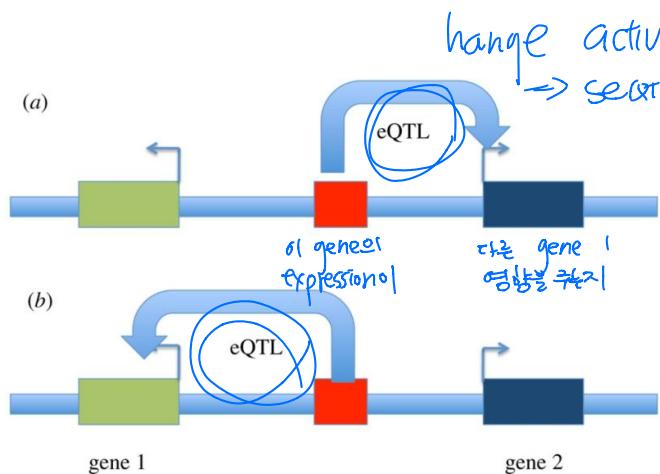
Direct association



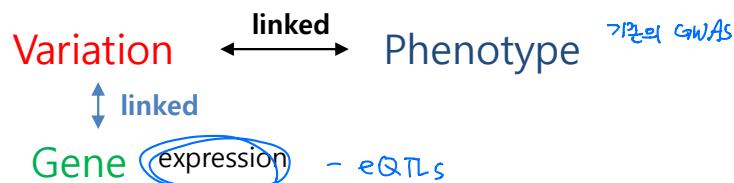
Alteration of neighboring gene expression

Alteration of splicing of neighboring genes.

⇒ difficult to explain how
of promoter

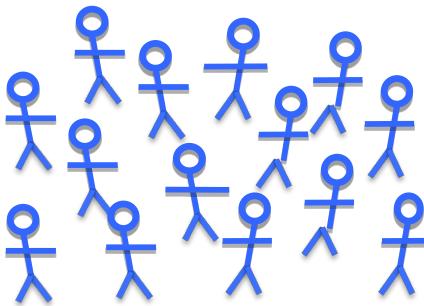


Expression quantitative trait loci (eQTLs) are genomic loci that contribute to variation in expression levels of mRNAs

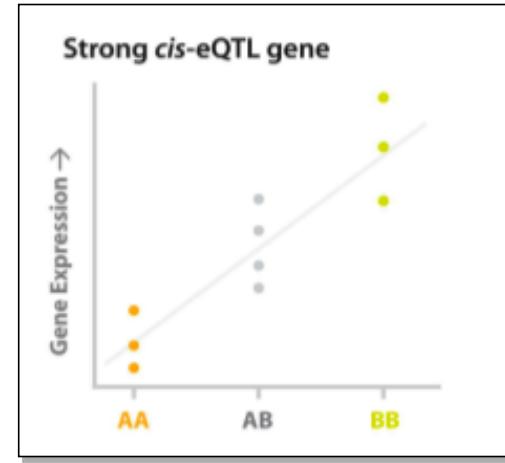


A Quantitative Gene-Expression Association (eQTL)

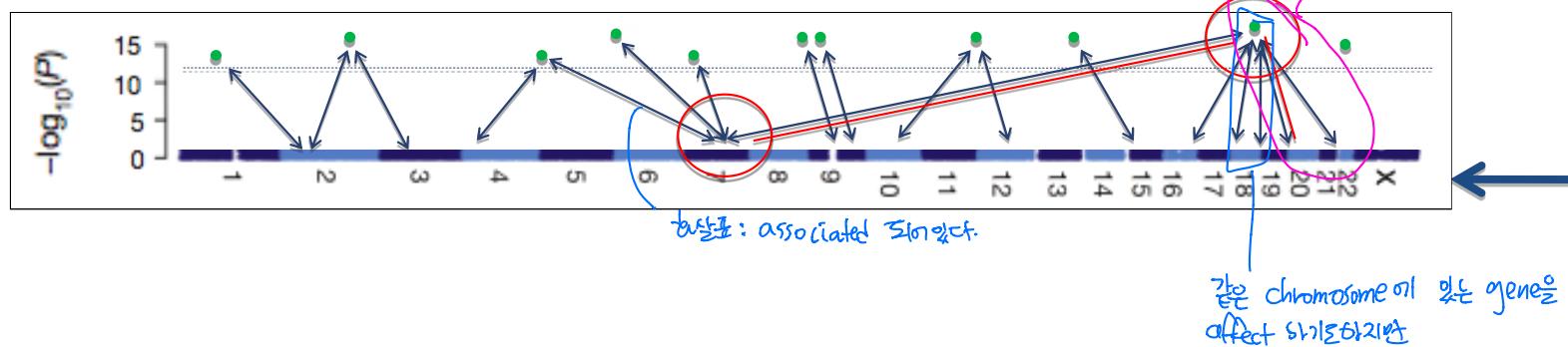
Sample Population



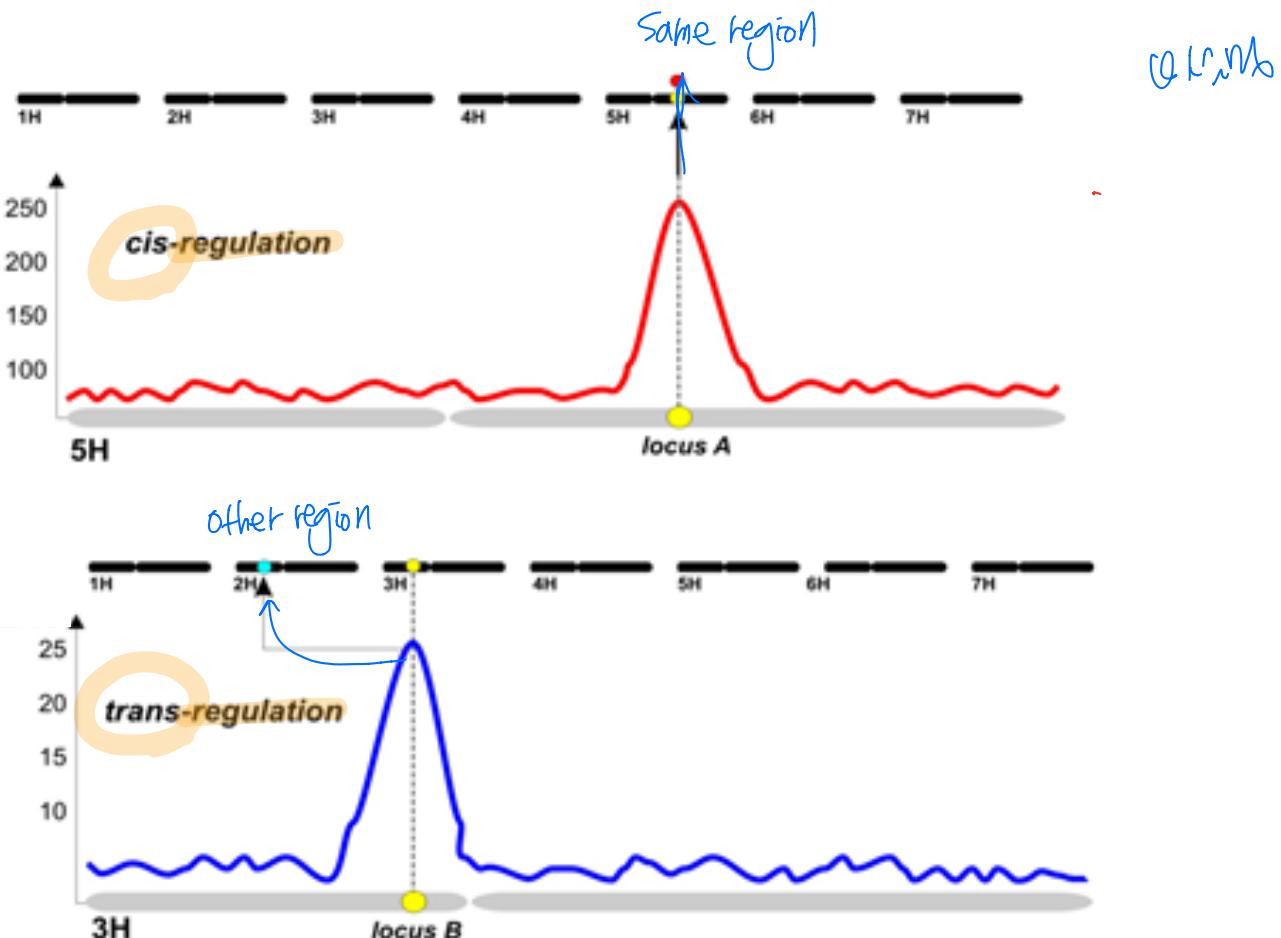
cDNA Levels



Expression Quantitative Trait Loci (eQTLs)

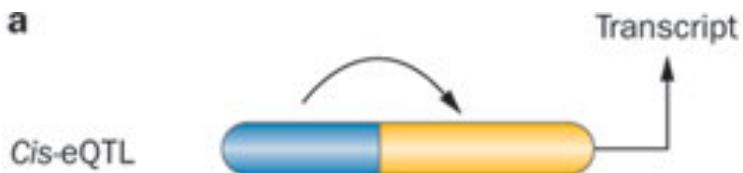


A Quantitative Gene-Expression Association (eQTL)

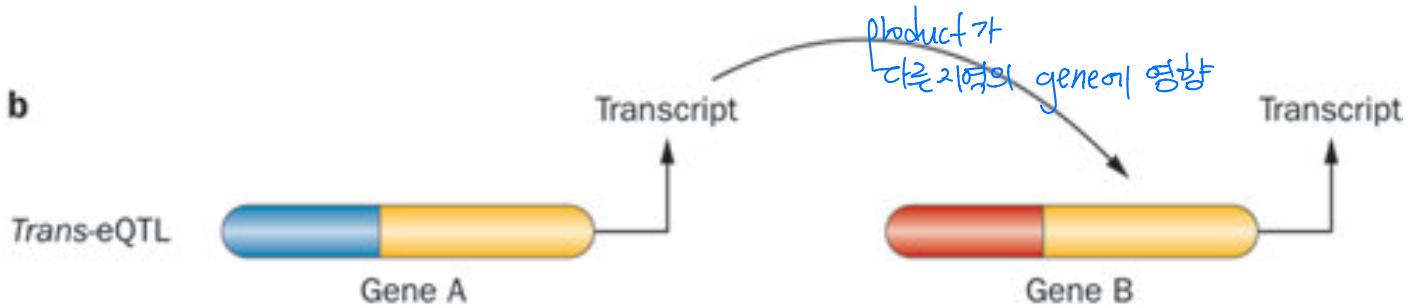


A Quantitative Gene-Expression Association (eQTL)

a

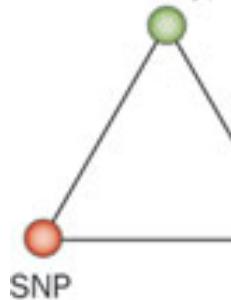


b



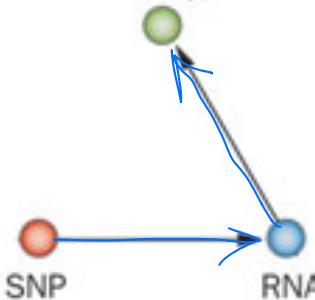
c

Phenotype



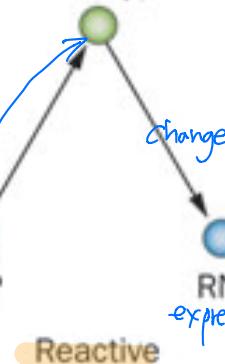
Correlation

Phenotype



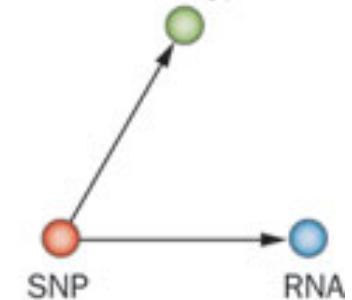
Causal

Phenotype



Reactive

Phenotype



Independent

GWA

Phylogenetic analysis

Sung Wook Chi

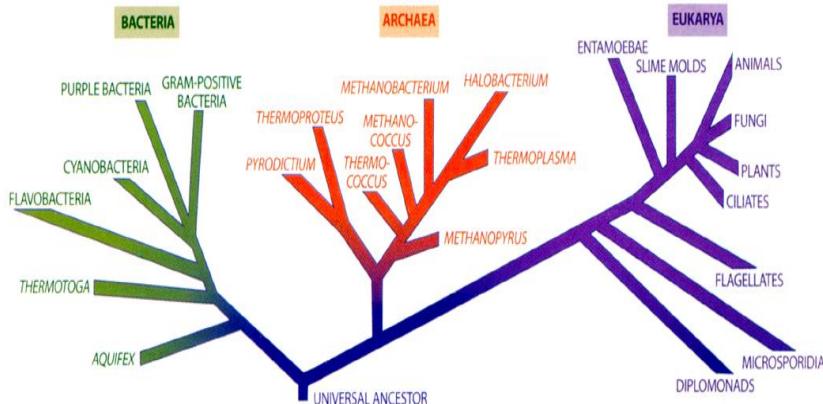
Division of Life Sciences, Korea University

Phylogenetics

phylogenetics is the study of evolutionary relationships among groups of organisms (Phylogeny), which are discovered through molecular sequencing data

use to categorize

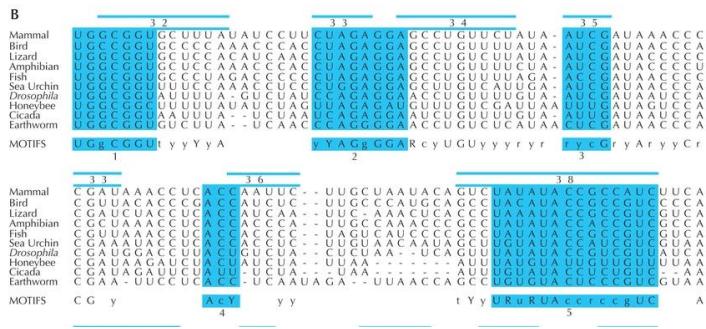
A version of the “tree of life” : Phylogenetic tree



DNA: morphological

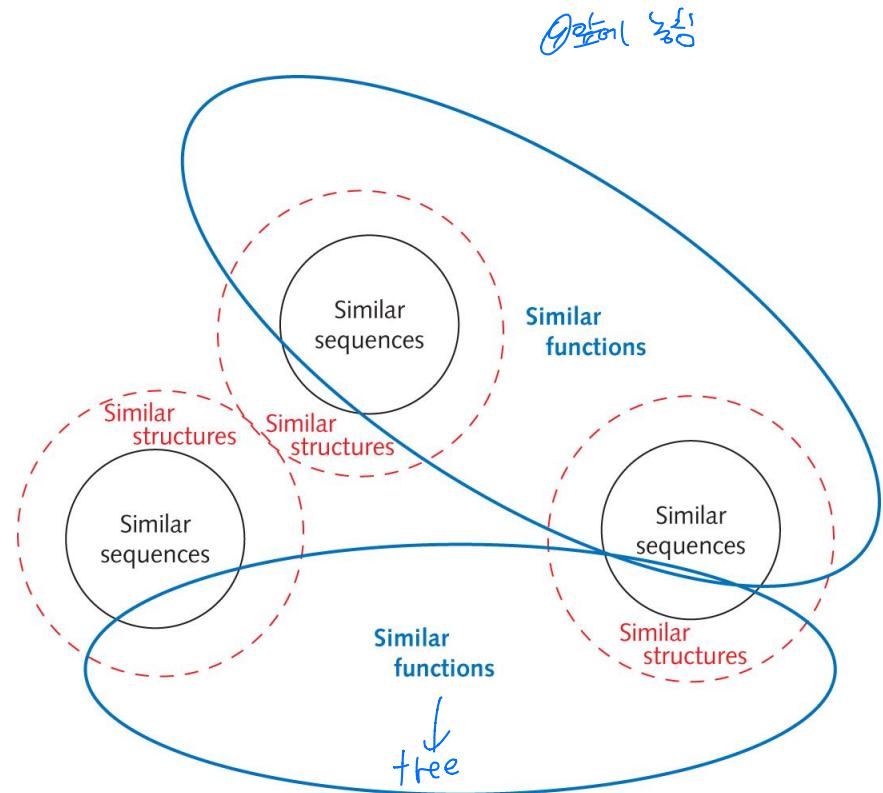
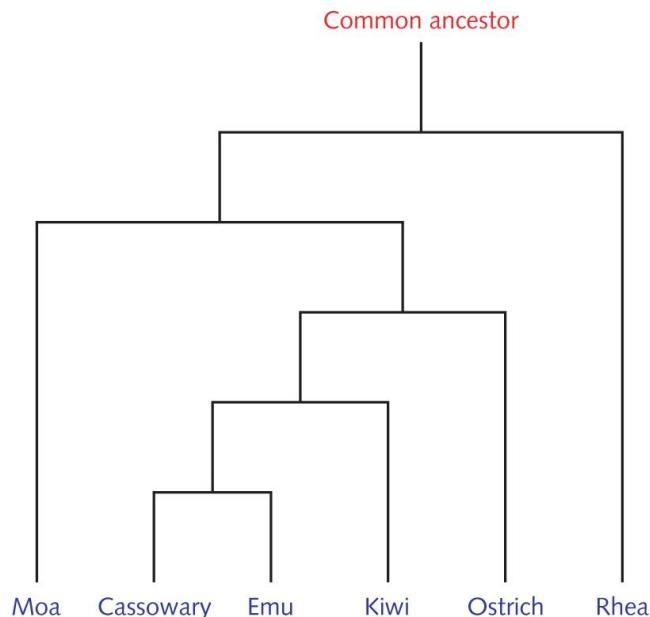
M. Madigan and B. Marrs, 1997

Obtained from aligned sequences of ribosomal RNA



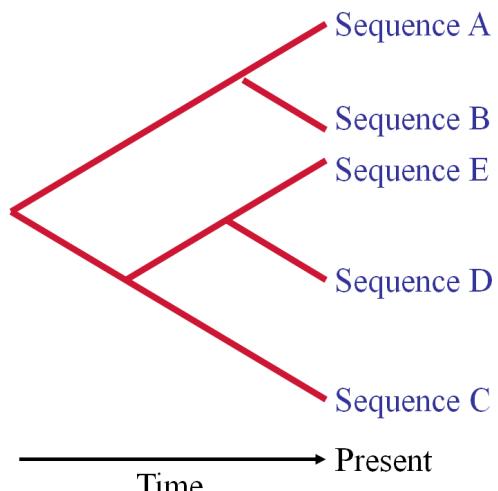
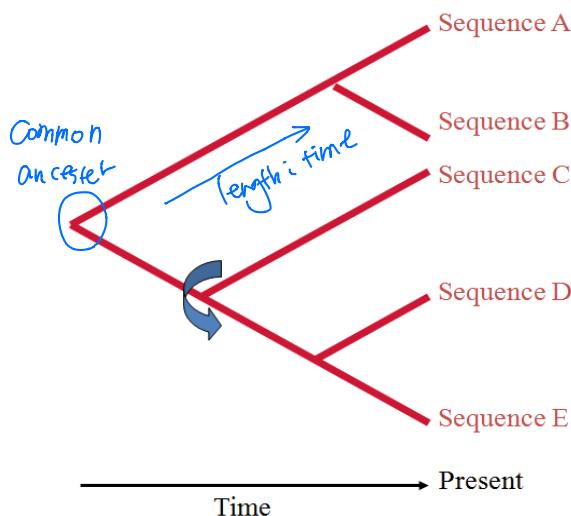
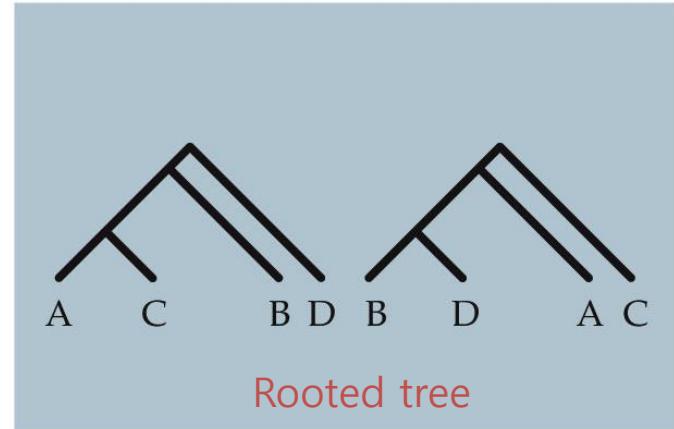
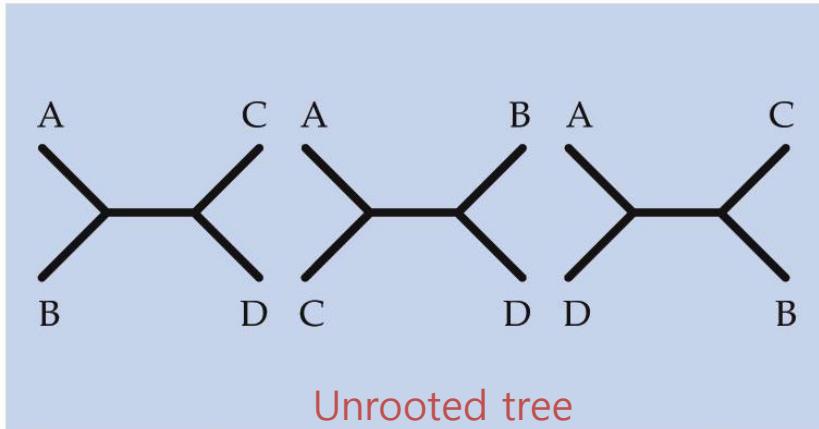
DNA information alignments

Phylogenetic analysis



Separate path of evolution but converges to the same result

Phylogenetics (evolutional tree)



Sequence Alignment, distances & tree

s1: ACCGTGAAGCCAATAC
s3: AGCGTGCAGCCAATAC

s2: A-CGTGCAACCATTAC

s4: AGGGTGCCGC-AATAC
s5: AGGGTGCCAC-AATAC

- Sequences to be aligned are phylogenetically related

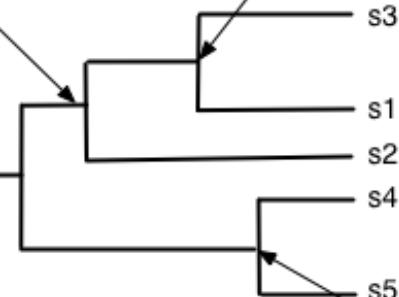
blue sequence diagram
↳ ↳
↳ ↳
↳ ↳
↳ ↳

Alignment 3

s1: ACCGTGAAGCCAATAC
s3: AGCGTGCAGCCAATAC
s2: A-CGTGCAACCATTAC

Alignment 2

s1: ACCGTGAAGCCAATAC
s3: AGCGTGCAGCCAATAC



Alignment 1

s4: AGGGTGCCGC-AATAC
s5: AGGGTGCCAC-AATAC

Perform pair-wise alignments

Measure distances

Generate tree

common ancestor

TACGGCTTTACCGA

descendants

1 TAACGGCTTTACCGA

2 TACGGCTGTACGA

aligned sequences

TAACGGCTTTACCGA

TA-CGGCTGTAC-GA

Variant
sequences

Sequence
Alignment

Phylogenetic
tree

insertion in descendant 1

mutation in descendant 2

deletion in descendant 2

↳ ↳
↳ ↳
↳ ↳
↳ ↳

Phylogenetic tree reconstruction

- Clustering Method: distance-based
 - **UPGMA**
 - Neighbor-joining
- Cladistic method: Character state-based methods
 - **Maximum parsimony**
- Maximum Likelihood

최적의
최적화

Clustering Method : UPGMA

UPGMA (unweighted pair group method with arithmetic mean)

Q(GVMB)

Distance matrix

hamming distance \Rightarrow 2 bits error

	ATCC	ATGC	TTCG	TCGG
ATCC	0	1	2	4
ATGC		0	3	3
TTCG			0	2
TCGG				0

initial grouping based rule

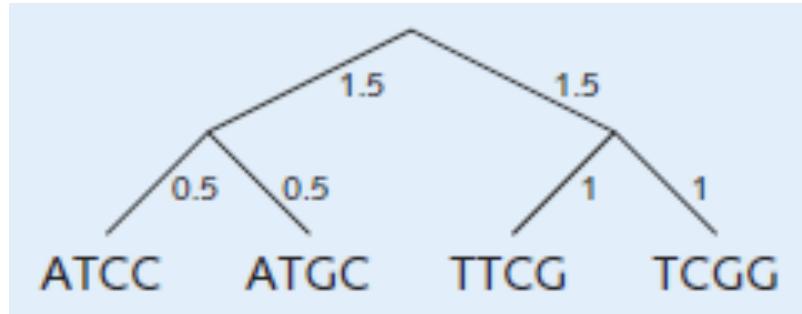
ATCC ATGC

The reduced distance matrix is:

	{ATCC, TTCG}	TCGG
{ATCC, ATGC}	0	$\frac{1}{2}(2 + 3) = 2.5$
TTCG	0	2
TCGG		0

symmetric rule

arbitrary rule



Cladistic Method : Maximum Parsimony

n=4

horizon in 4-locus tree

ATCA

A→G

A→T

ATCG

TTCA

ATCG C→G

ATGG

T→C

TCCA

TTCA

n=8

ATCG

G→A

A→T

ATCA

TTCG

A→G

A→T, T→C

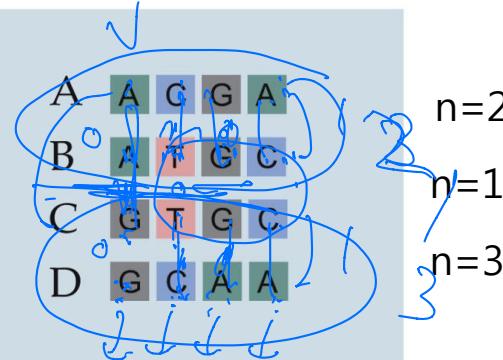
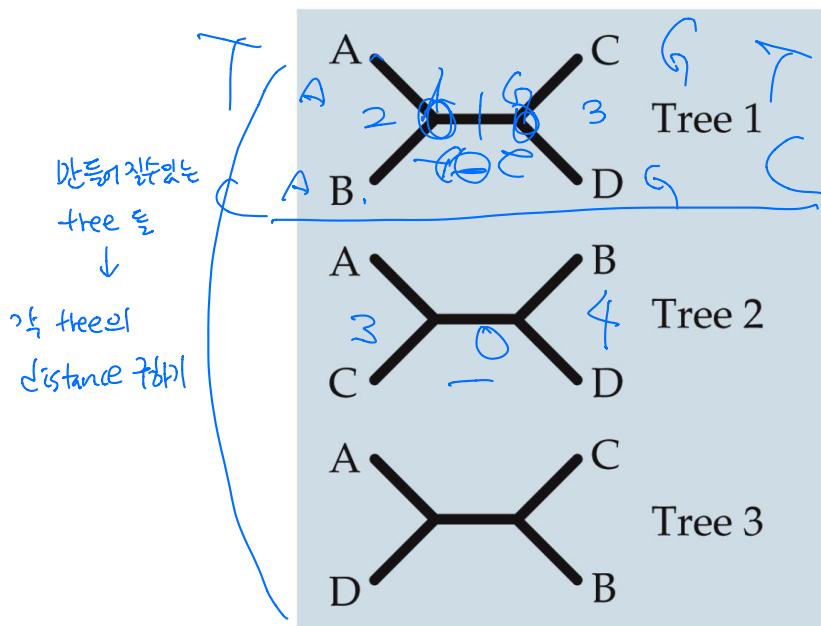
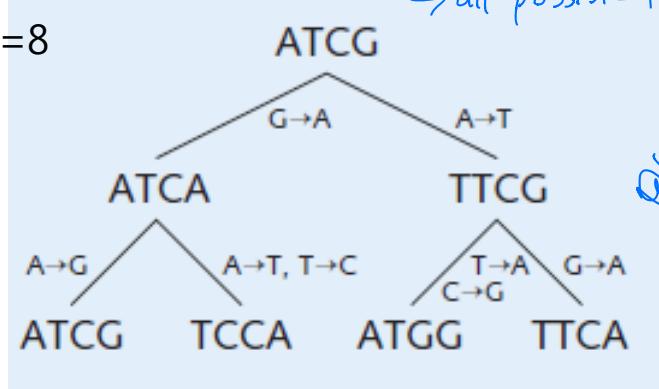
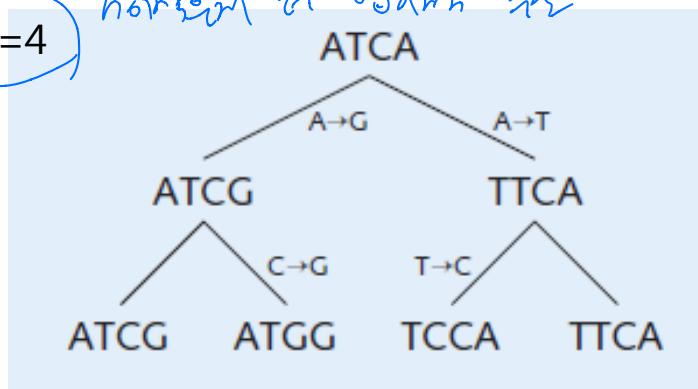
ATGG

T→A
C→G

TTCA

all possible trees

Q(?)



Tree	1	2	3	4	Tot
1	1	2	1	2	6
2	2	2	1	2	7
3	2	1	1	1	5

제일 적절한 tree

Maximum Parsimony (exercise)

EXERCISE 2.4 A simple phylogenetic analysis

A few years ago, a pressing question in the area of human evolution was the exact phylogenetic relationship between humans, chimpanzees, and gorillas. Fossil, morphological, and early molecular data provided conflicting and inconclusive results as to which pair of organisms were most closely related. Suppose the following four short DNA sequences are available, including a sequence from the outgroup organism, orangutan:

Human	CGAAATGCAT
Chimpanzee	CGGACTGTAT
Gorilla	CAAGCTGTAC
Orangutan	TAAGACGTAG

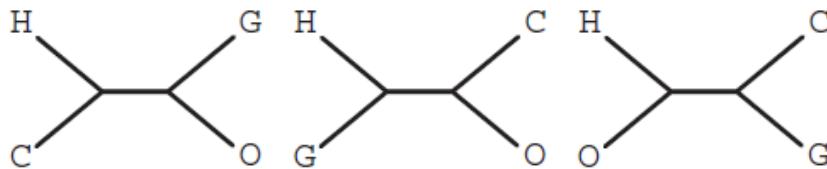
↳ most divergent
⇒ root

Carry out a phylogenetic analysis to infer the rooted evolutionary tree of human, chimpanzee, and gorilla.

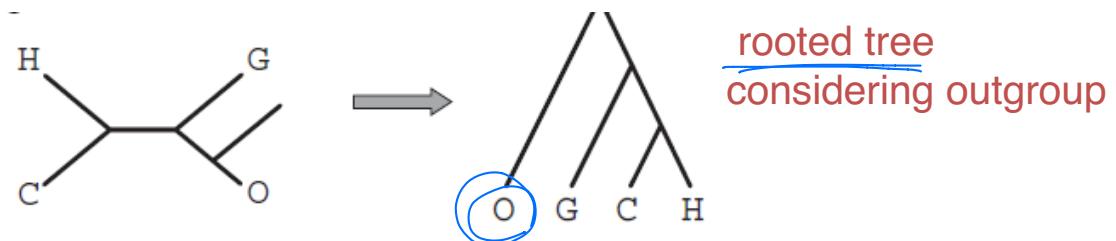
Maximum Parsimony (exercise)

<i>Human</i>	CGAAATGCAT
<i>Chimpanzee</i>	CGGACTGTAT
<i>Gorilla</i>	CAAGCTGTAC
<i>Orangutan</i>	TAAGACGTAG

all possible unrooted tree



	1234567890	Total	like to be true calculate edit distance
Tree A	1111210102	10	
Tree B	1212210102	12	
Tree C	1212110102	11	



Multiple Sequence Alignment

- Idea: Sequences to be aligned are phylogenetically related
 - these relationships are used to guide the alignment
 - Popular implementations: **CLUSTALW**,
1. **Perform pair-wise alignments** between all pairs of sequences : $(n \times (n-1)/2$ possibilities)
 2. **Generate distance matrix.**
 - Distance between a pair = number of mismatched positions in alignment divided by total number of matched positions
 3. Generate a Neighbor-Joining '**guide tree**' from distance table
 4. Use guide tree to progressively align sequences in pairs from tips to root of tree.

Multiple Sequence Alignment

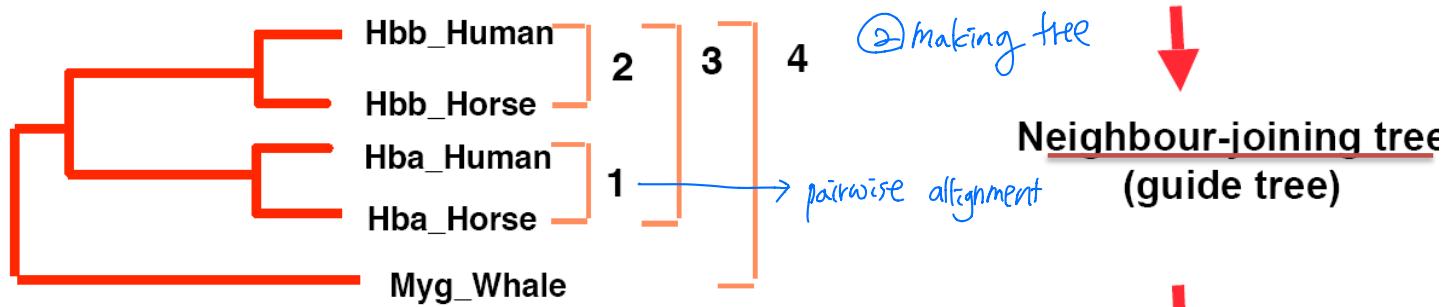
Q19/20

Hbb_Human	1	-				
Hbb_Horse	2	.17	-			
Hba_Human	3	.59	.60	-		
Hba_Horse	4	.59	.59	.13	-	
Myg_Whale	5	.77	.77	.75	.75	-

① Calculate distance

CLUSTAL W

Quick pair wise alignment
calculate distance matrix

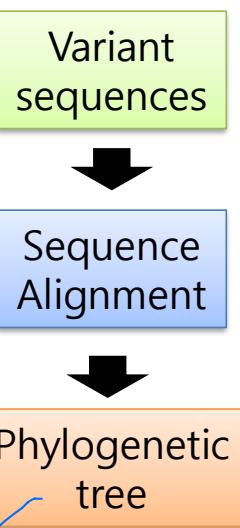


alpha-helices

1	PEEKSAVTALWGKVN	--VDEVGG	2	3	4
2	GEEKAAVLALWDKVN	--EEEVGG			
3	PADKTNVKAAWGKVGGAHAGEYGA				
4	AADKTNVKAAWSKVGGGHAGEYGA				
5	EHEWQLVLHVWAKVEADVAGHGO				

Progressive alignment
following guide tree

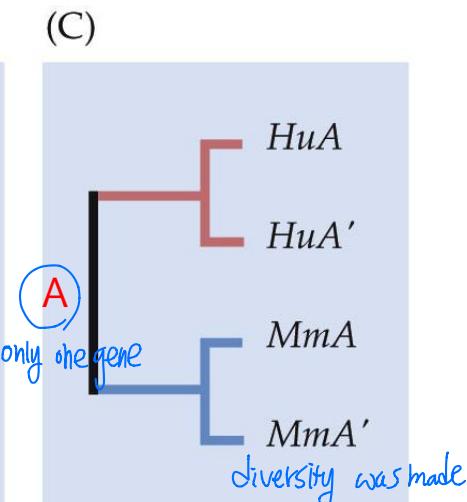
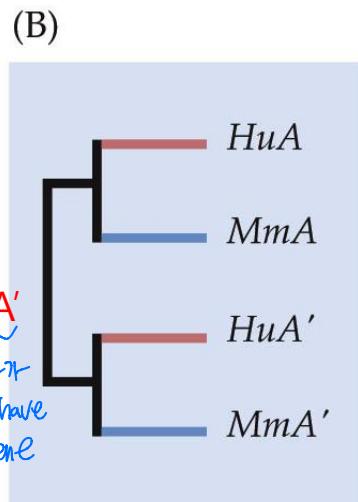
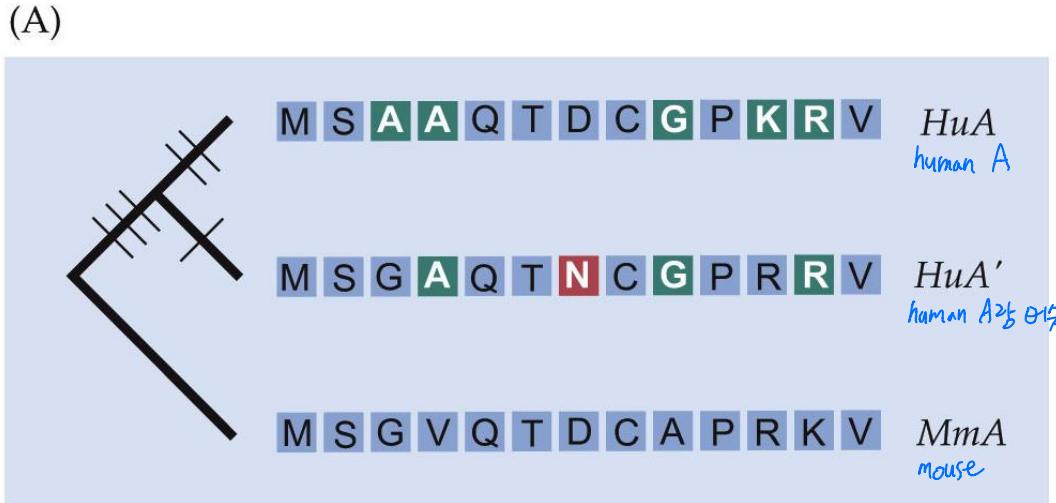
Interpretation of phylogenetic tree: Orthologs and paralogs



- Recalculate / Infer Distance or similarity b/w sequences
- Topology (order)
- Length (evolutional time)



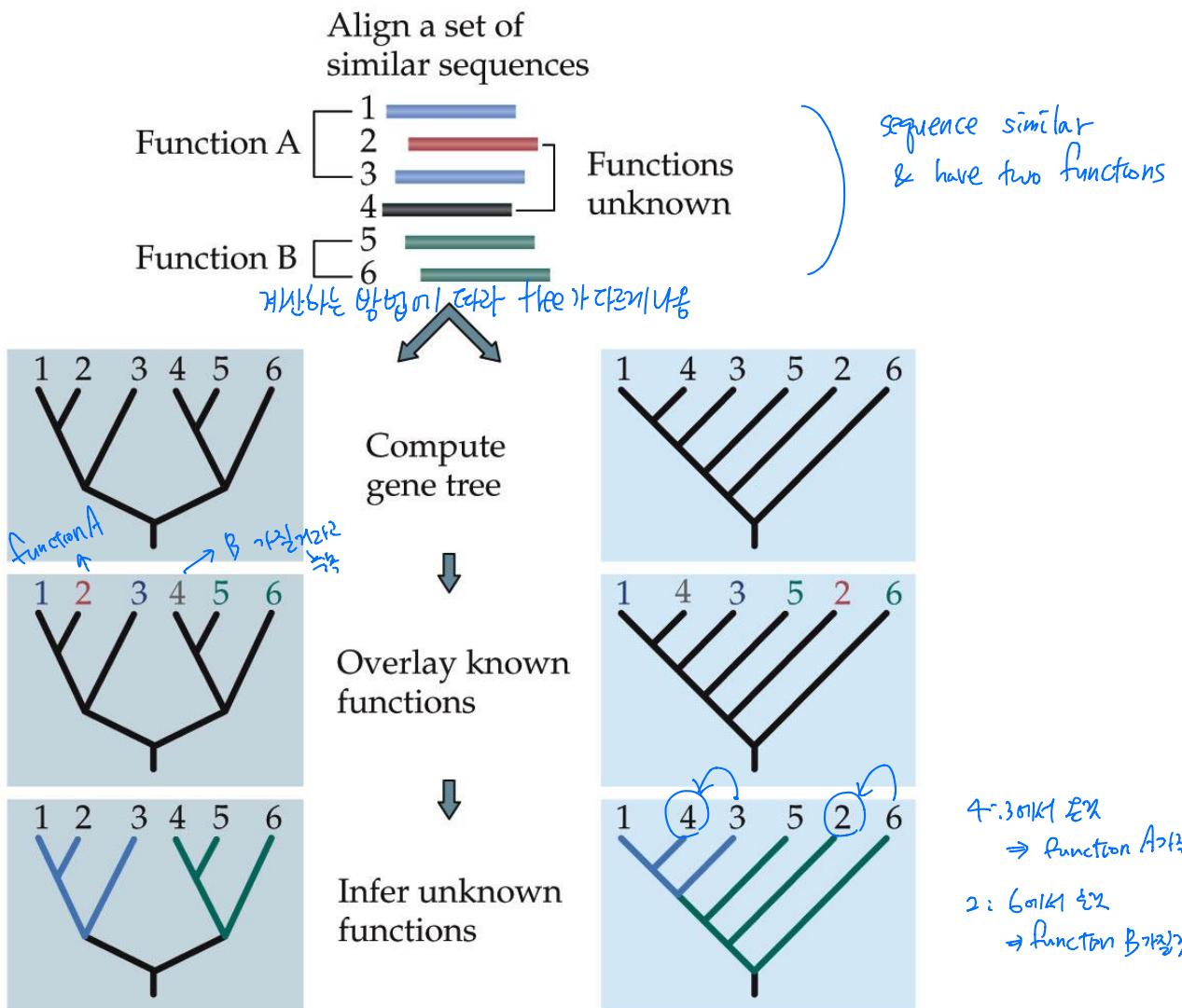
Evolutionary relationship / interpretation



Evolutionary event 이 향후에 일어남

Inferring gene function from phylogenetic analysis

①설명



Phylogenetic shadowing

Within closely related species

