

# Summary & More

: Whole genome sequencing (NGS)

Sung Wook Chi

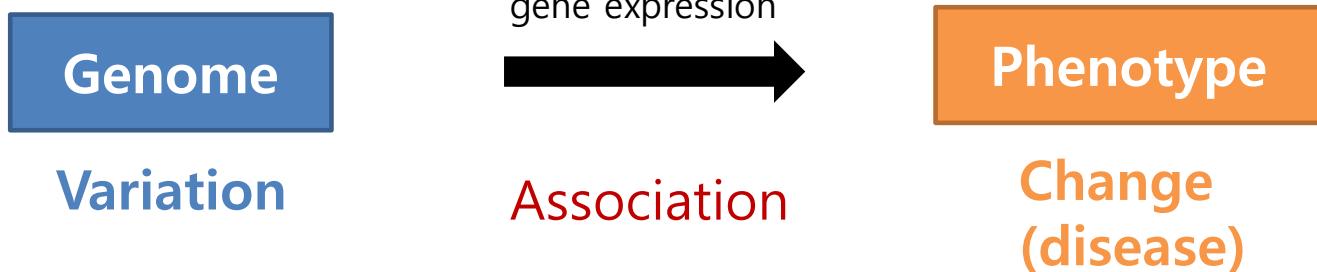
Division of Life Sciences, Korea University

# What we learned

1. Introduction ( from Genomics to Functional Genomics)
2. Genome Projects
3. Human Genome Projects
4. Next-generation Sequencing
5. NGS & Sequence Alignment

## Variation of genome sequence

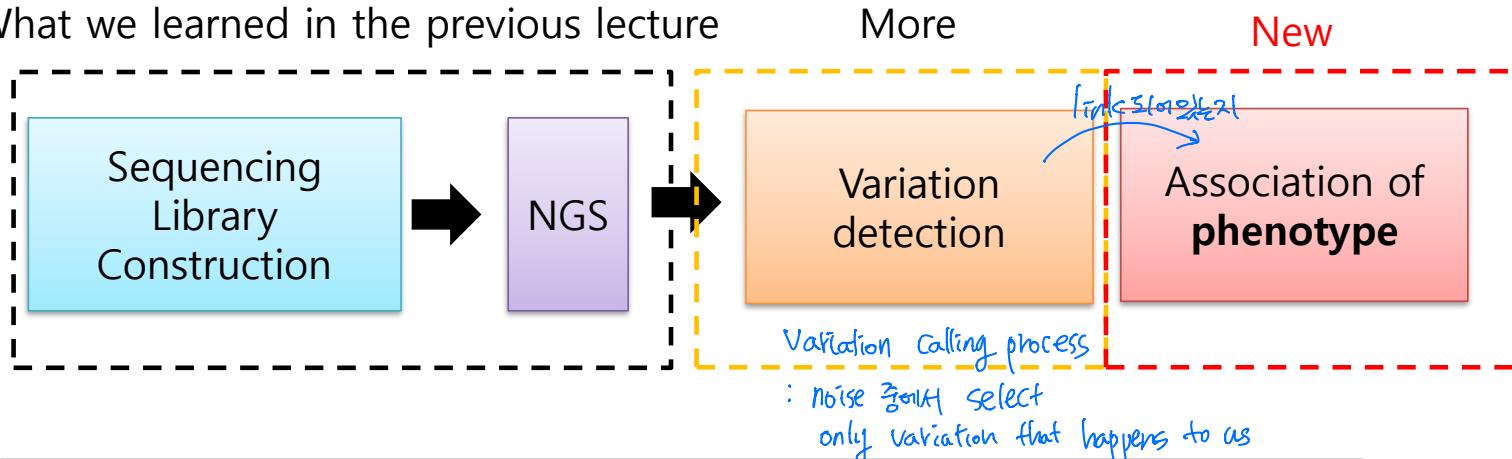
1. Haplotype & SNPs
2. Whole genome sequencing *Variation 연구를 위한 WGS의 적용*



# Analysis of genomic variation

## Variation of genome sequence

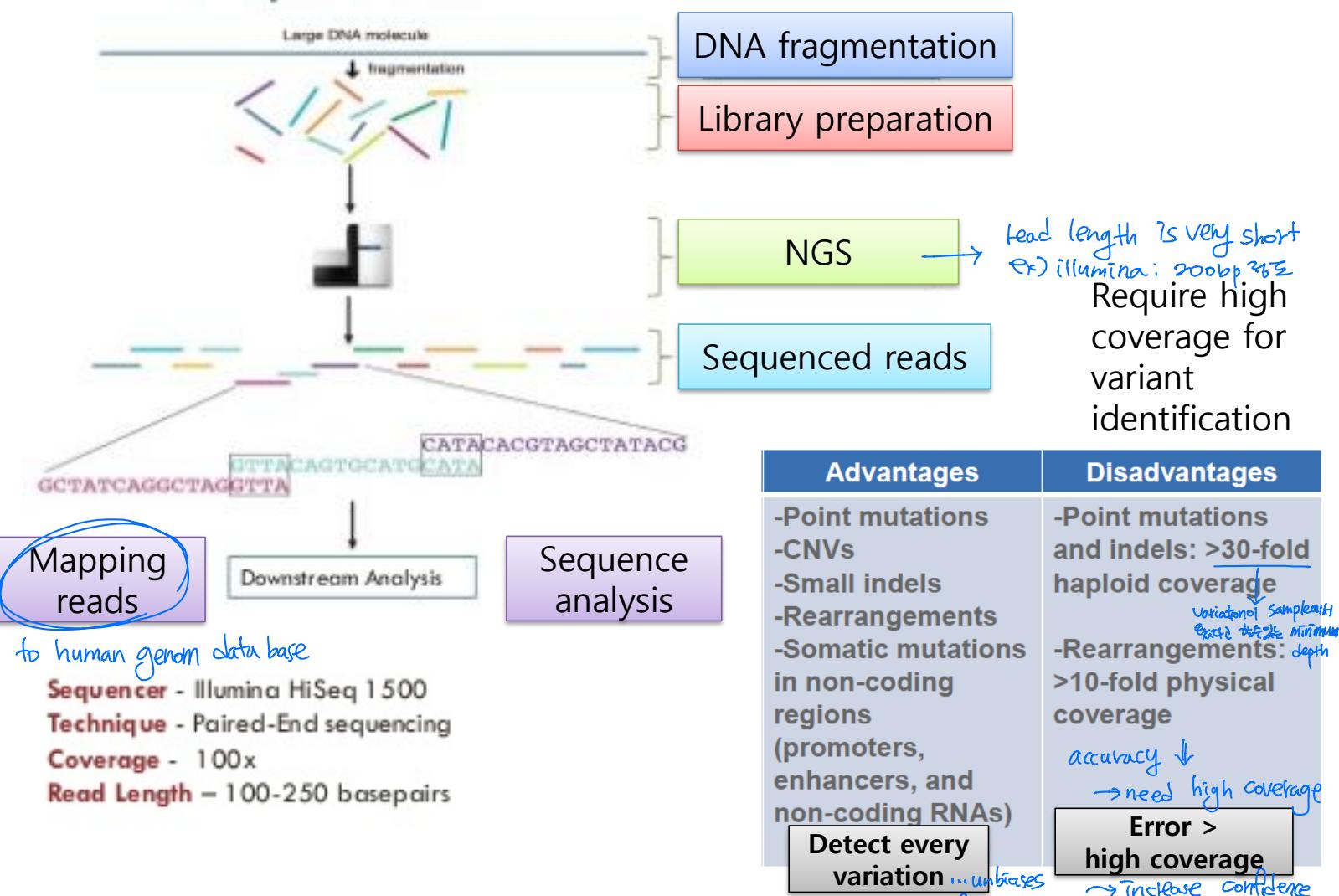
What we learned in the previous lecture



1. Whole genome sequencing (WGS by NGS)
2. Variation Analysis for NGS data
3. Genome-wide association study (GWAS)

# Whole Genome Sequencing (NGS) for variant identification

## Principle of WGS

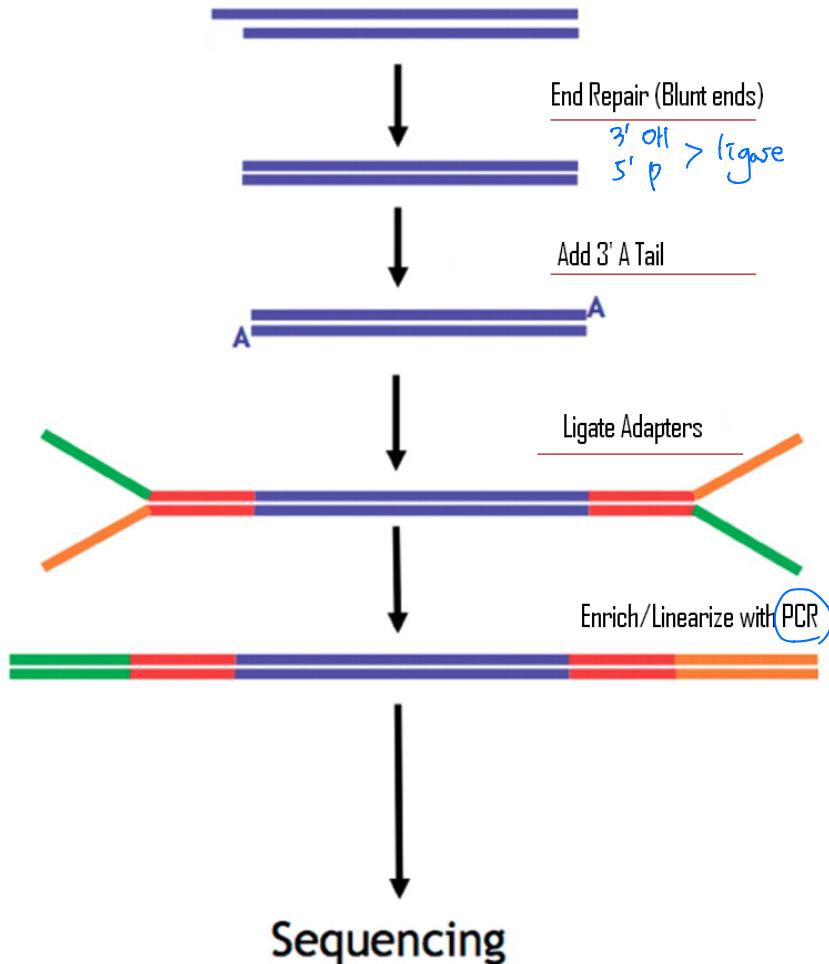


→ contain all information

# NGS Library preparation

Shear Genomic DNA or begin with cDNA

OKay



1. DNA fragmentation  
/ Target Selection

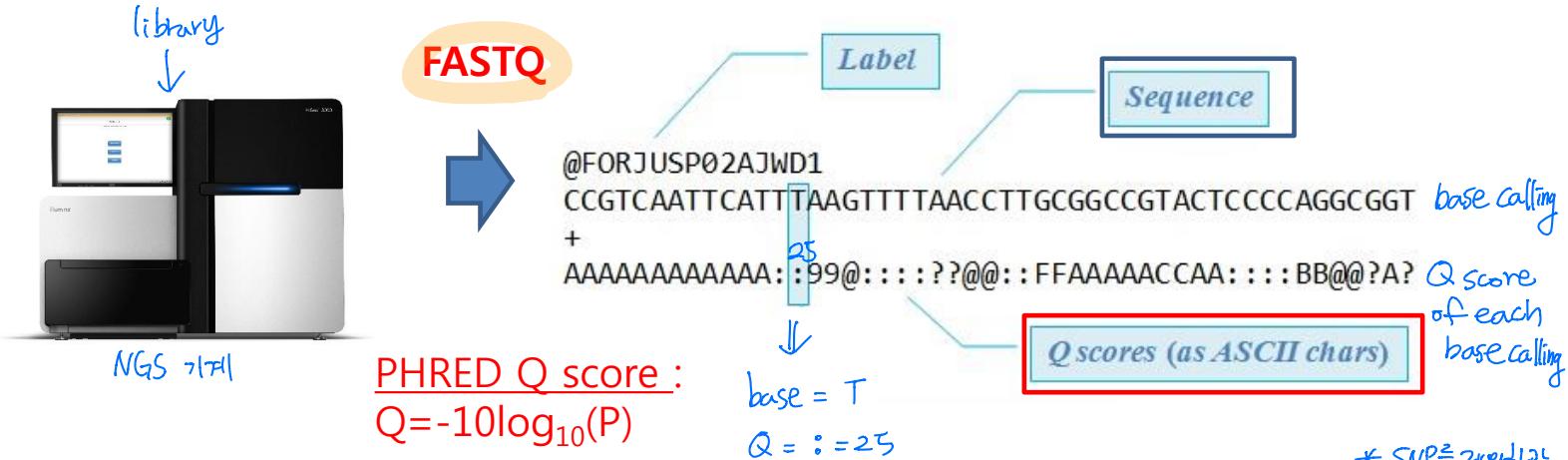
PCR amplification for specific region of interest

2. Adapter ligation

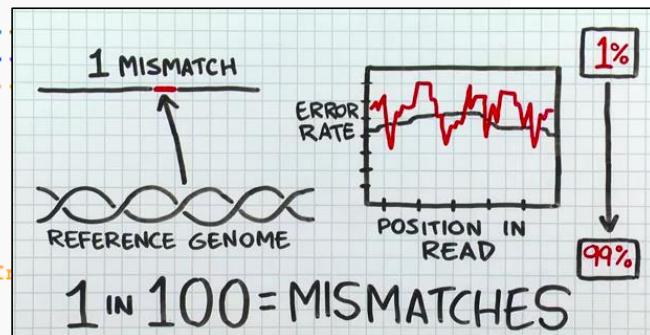
3. Size selection

4. Library quantification (QC)

## Sequencing data from NGS : FASTQ format & Q score

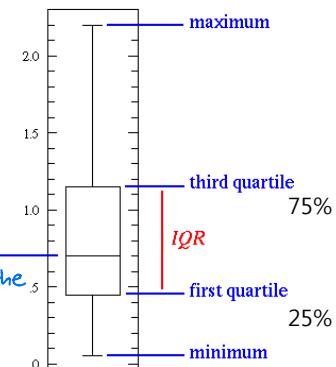
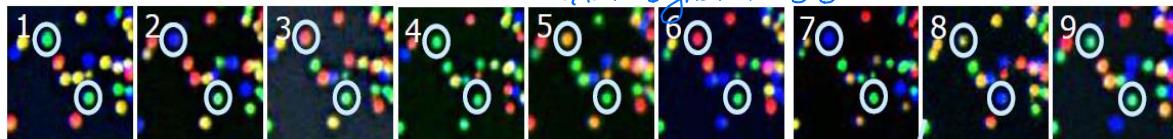


S - Sanger Phred+33, raw reads typically (0, 40)  
X - Solexa Solexa+64, raw reads typically (-5, 40)  
I - Illumina 1.3+ Phred+64, raw reads typically (0, 40)  
J - Illumina 1.5+ Phred+64, raw reads typically (3, 40)  
with 0=unused, 1=unused, 2=Read Segment Quality Control In  
(Note: See discussion above).  
L - Illumina 1.8+ Phred+33, raw reads typically (0, 41)



# Sequence Quality Control

recording same position → quenching의 제대로 되지 않으면 뛰는 signal이  
두의 Signal에 영향

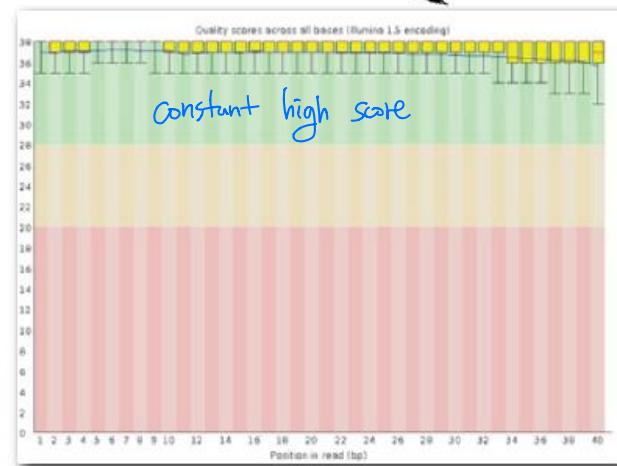


per base sequence quality

good

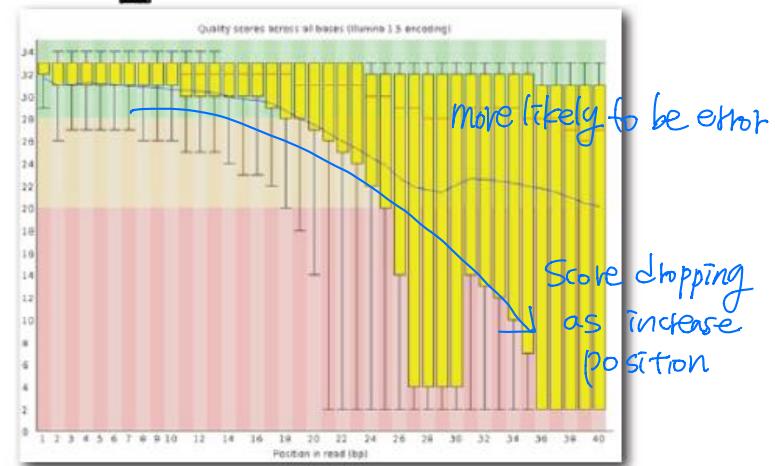
bad

Q score



Position in reads

Q score



Position in reads

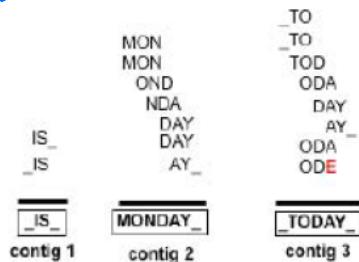
Q score tends to decrease depending on the increment of position

# NGS Sequence mapping

assemble seq newly

de novo assembly

reference genome seq!  
do de novo assembly



NDA DAY  
MON MON  
OND IS AY  
TO TO  
DAY IS ODA  
TO ODA  
DAY AY

most case

Reference Mapping

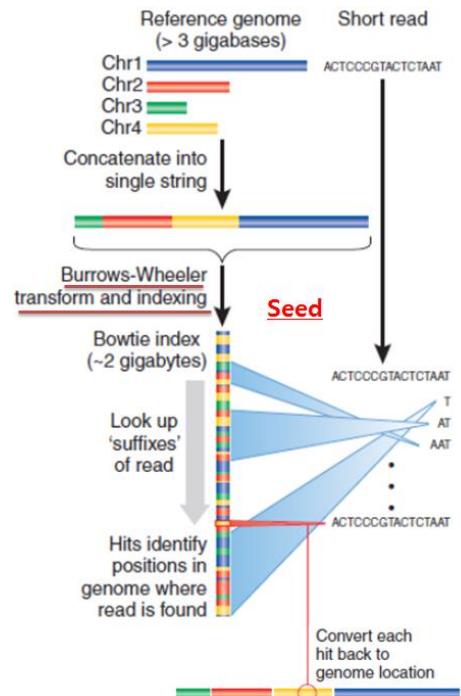
(다른 알고리즘은 레퍼런스 seq를 찾는다)

reference mapping

ODE TO IS  
TO AY OND NDA  
TO ODA MON  
TOD DAY  
DAY ODA IS MON AY  
TODAY\_IS\_MONDAY reference  
TODAY\_IS\_MONDAY consensus

(modified from Panu Somervuo)

Mapping billions of short reads on genome

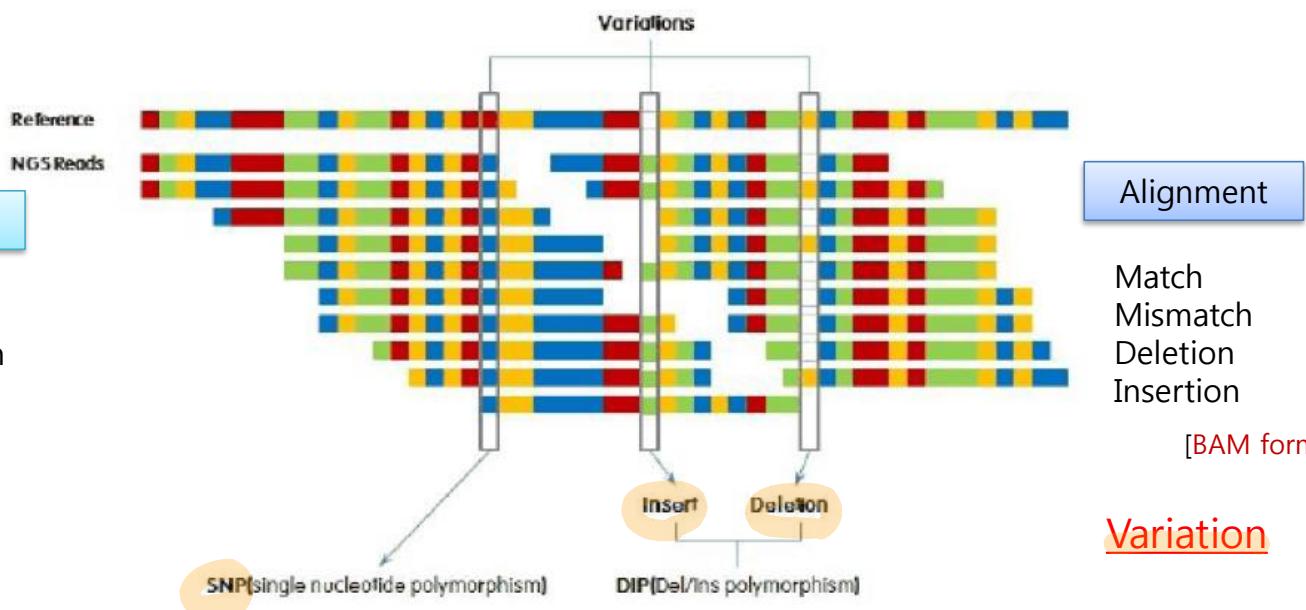
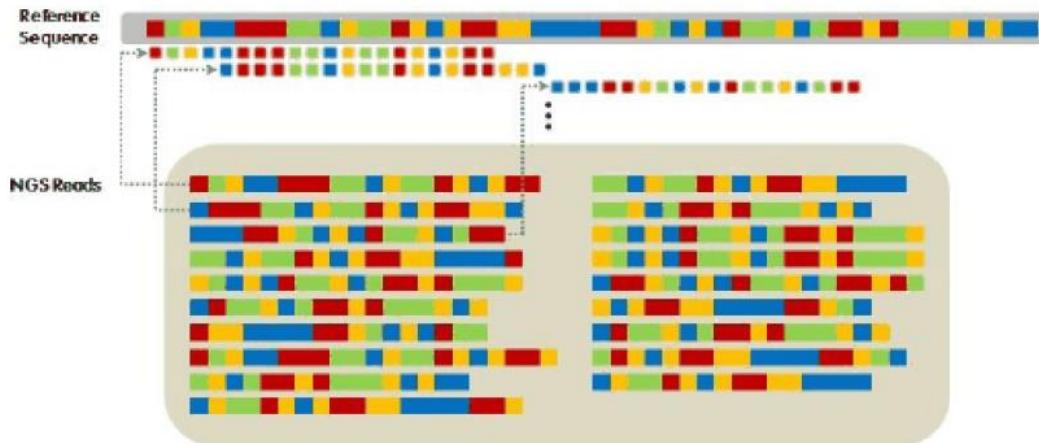


## NGS WGS reads mapping (BWT based method)

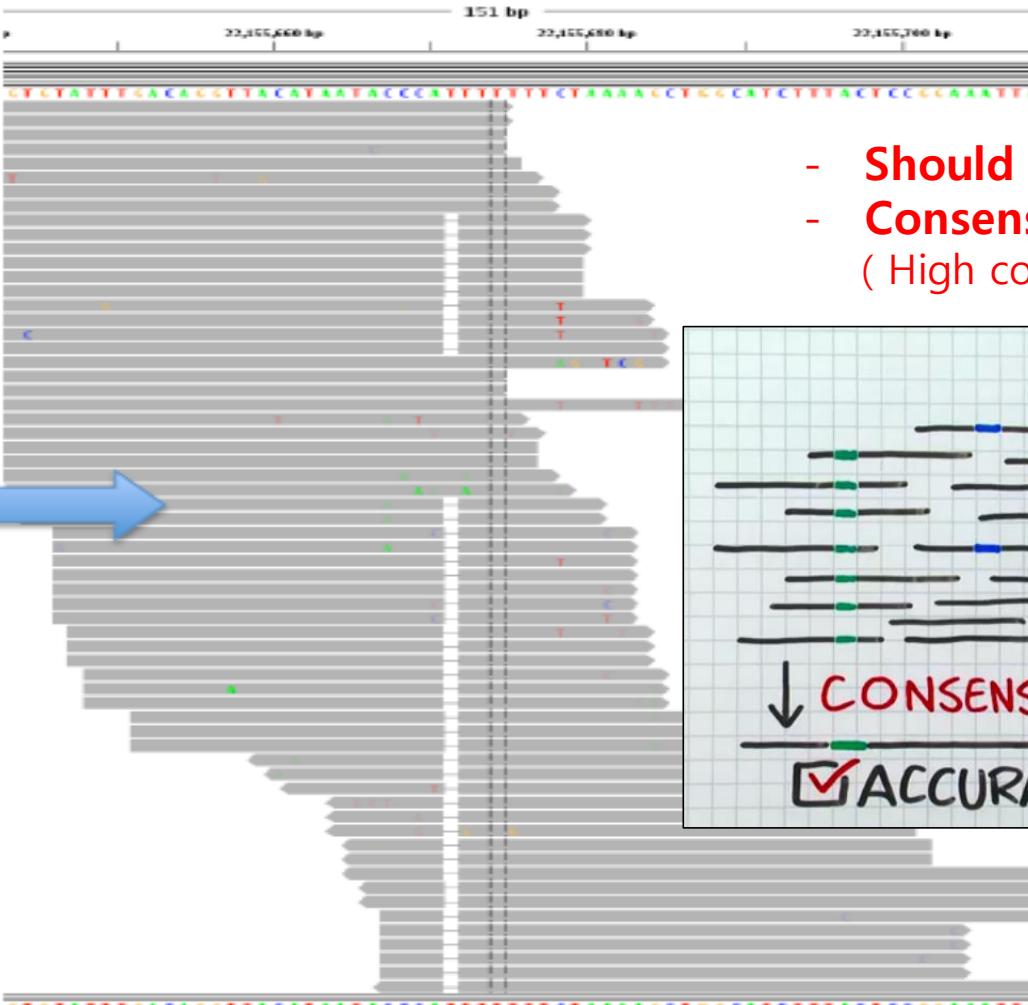
- Transform genome sequences into data structure (condensed, indexed, fast search)
- Find billions of short reads sequence there (near perfect match)
- Location of matched reads (Genomic coordination, Alignment)

BWT (Burrow-Wheeler transform) based mapping algorithm

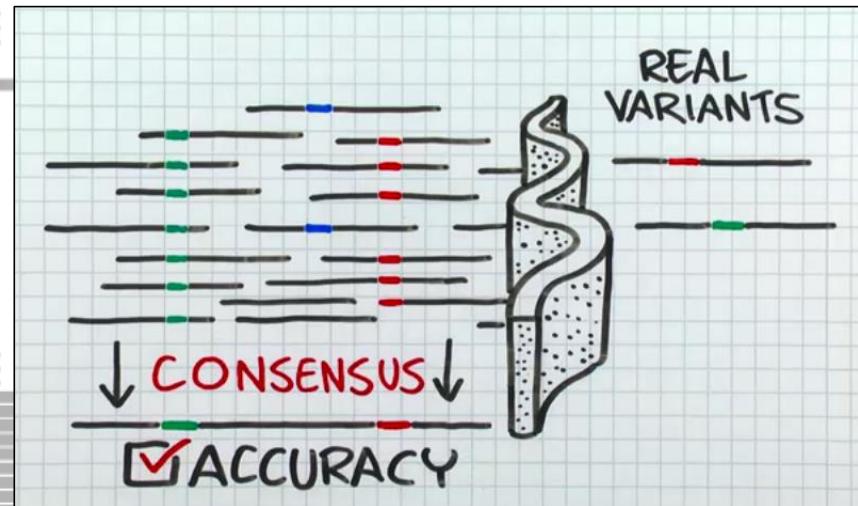
# Sequence Alignment (Mapping)



# Variation in NGS reads : Real or not ?



- Should consider Q-score
- Consensus  $\Rightarrow$  reproducibly see this event  
( High coverage > Statistics )

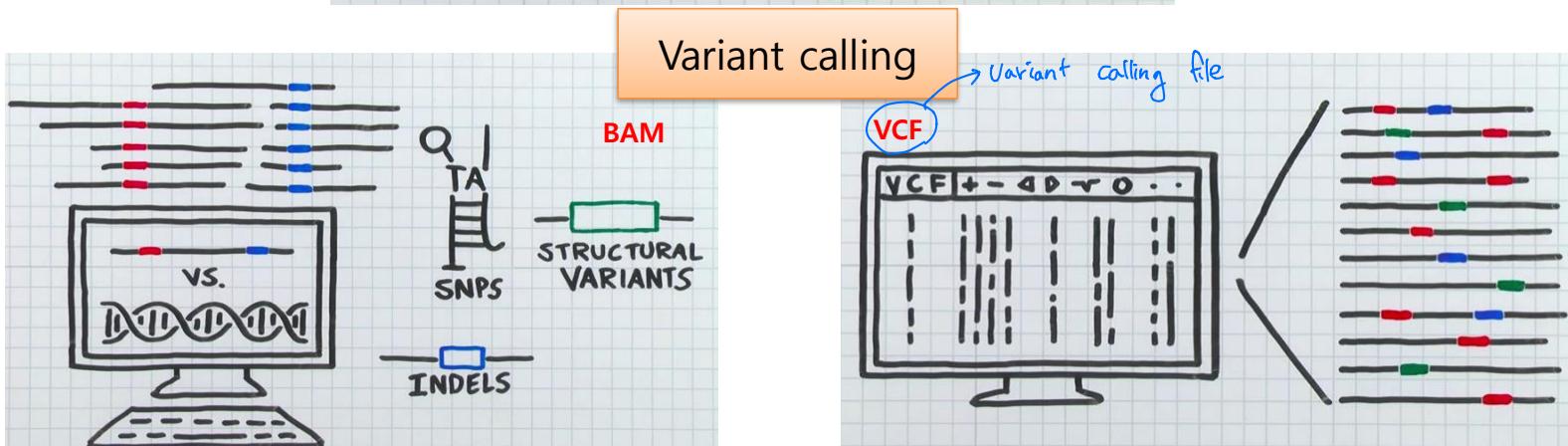
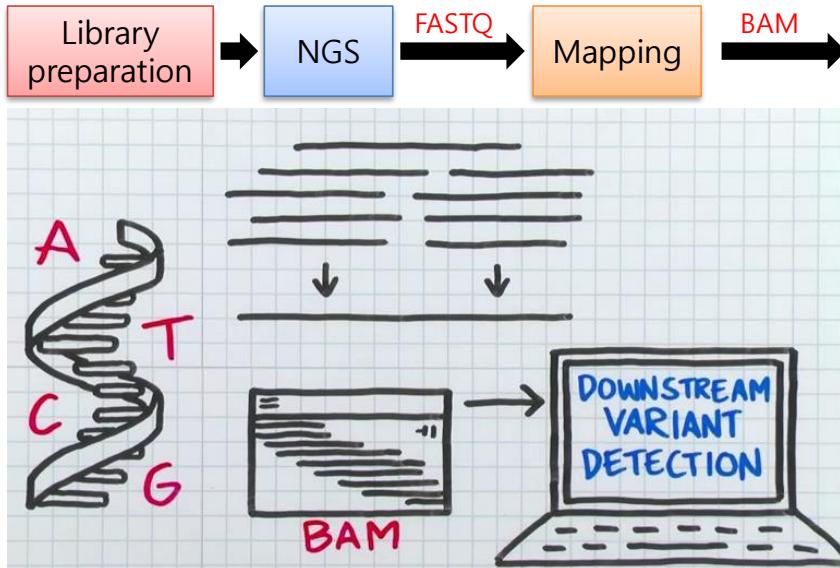


~ Statistically test

**"Variant calling"**

Seq error가 아니라 Sampleonly인 경우  
[VCF format]

# NGS WGS analysis for variation detection



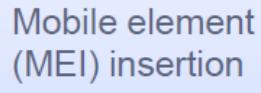
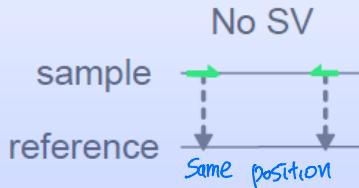
# SV Discovery with diverse resources

Structure variation

## Discovery of Structural Variation

- Challenging in mapping
- Difficult for interpretation

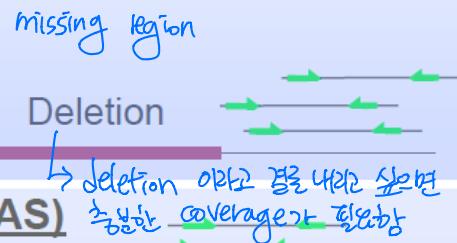
### Read Pairs (RP)



### Read Depth (RD)



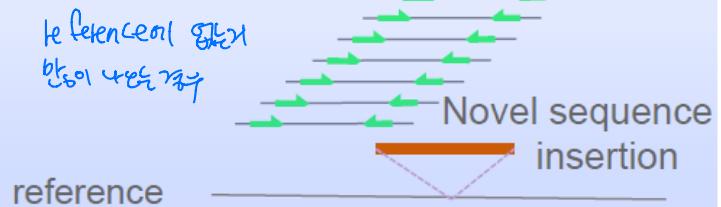
### Copy number variation



### Split Reads (SR)



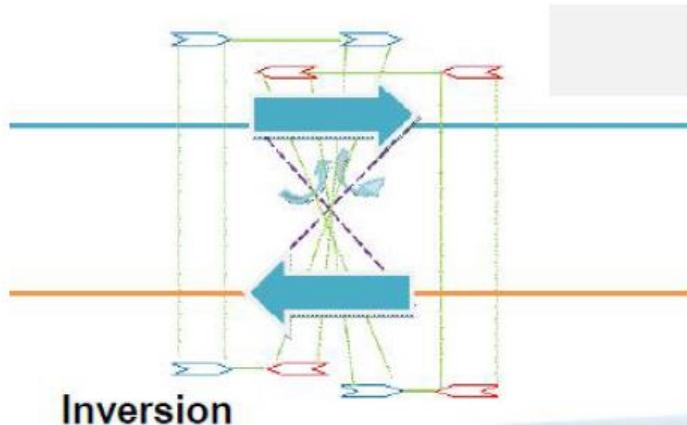
### Assembly (AS)



# SV Discovery with diverse resources

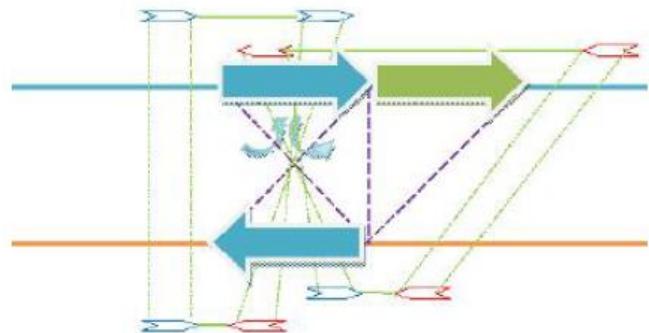
④ kym

Structure Variation의 예



Inversion

direction of



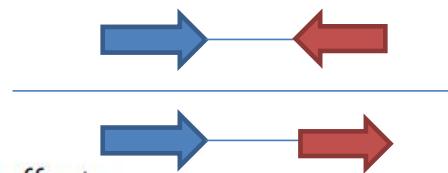
Inversion + Deletion

direction & size

Faulty mapping of pair-end reads could be interpreted as "inversion" event !!

## ❑ Inversion

- Does not involve a loss of genetic information
- Simply rearranges the linear gene sequence
- Generally considered to have no deleterious or harmful effects



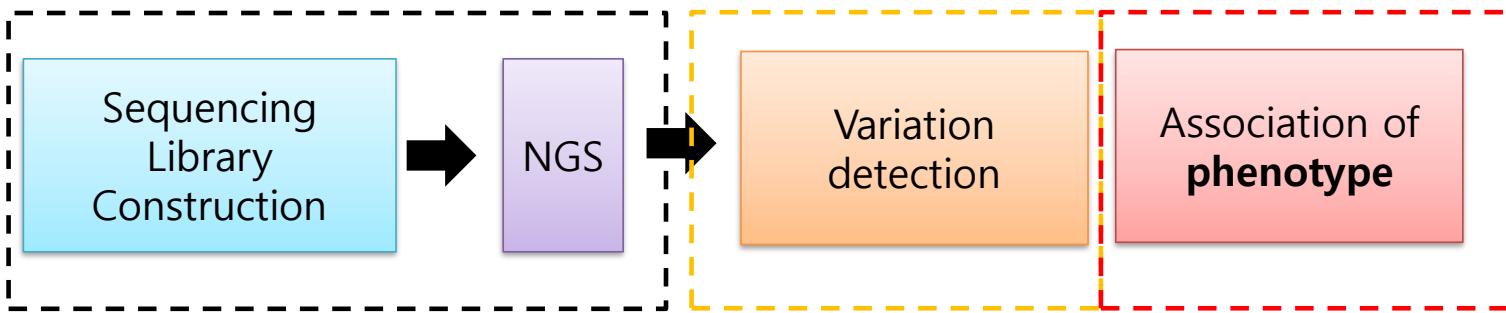
# Analysis of genomic variation

## Variation of genome sequence

What we learned in the previous lecture

More

New



1. Whole genome sequencing (WGS by NGS)
2. Variation Analysis for NGS data
3. Genome-wide association study (GWAS)

# Genome-wide association study (GWAS)

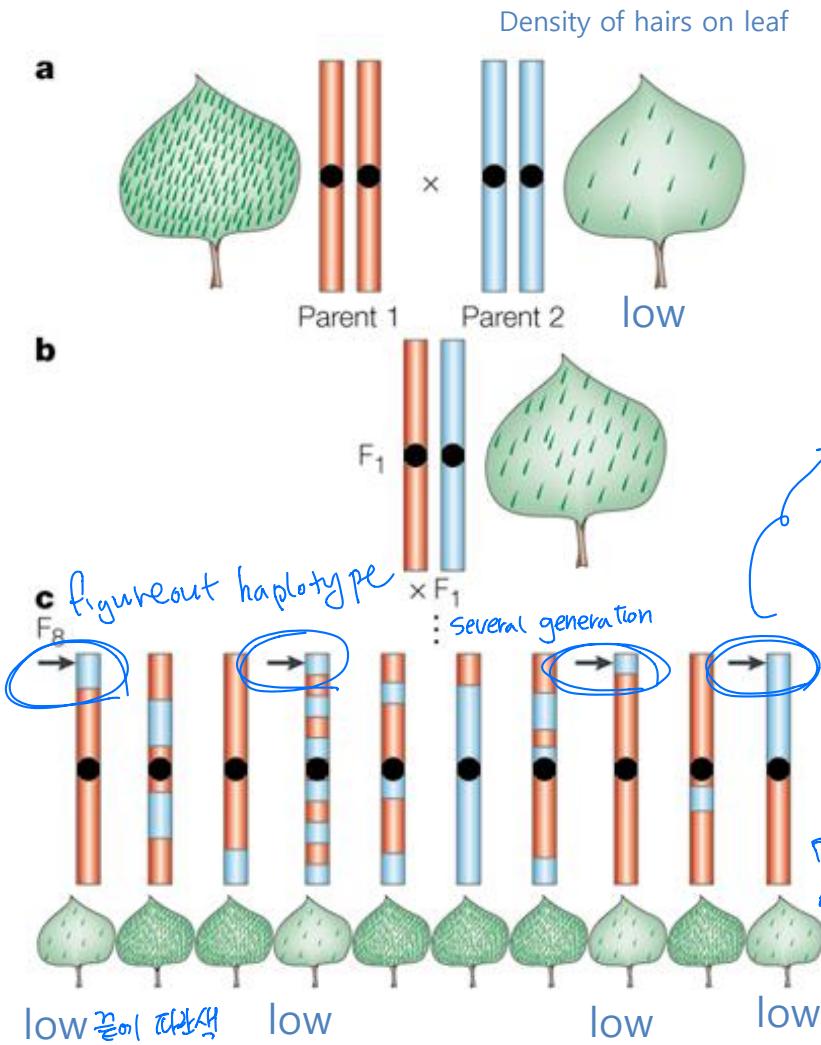
Sung Wook Chi

Division of Life Sciences, Korea University

$P \rightarrow$  phenotype or disease loci  
vary in degree

# Variation and phenotype: Quantitative trait loci (QTL)

Variation (region) vs. Phenotype (trait)



## Quantitative traits

: phenotypes (characteristics) that vary in degree and can be attributed to polygenic effects (i.e., product of two or more genes, and their environment)

## Quantitative trait loci (QTLs)

: stretches of DNA linked to, or containing, the genes or variations that underlie a quantitative trait.

## How to identify QTLs ?

QTL linkage mapping

phenotype  
SNPs  
compare

- SNPs , Haplotype
- association study with phenotypes

기본원리 100

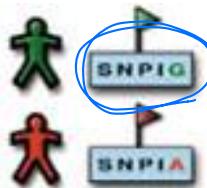
# Using SNPs to Track Predisposition to Disease and other Genetic Traits



DNA from different individuals sequenced

ATCGCTGCATGCCA  
ATCGCTGCATGCCA  
ATCGCTGCATGCCA  
ATCGCTGCATGCCA  
ATCGCTGCATGCCA  
ATCGCTGCATGCCA

Variation at a single nucleotide



of SNP > disease 관계를 찾는다  
patient > normal 일치

- Quantitative traits
- Diseases

## Hypothesis-free approach

### Sample with disease



A higher than expected incidence in a disease group suggests SNPIG is associated with a disease (or SNPIA is protective)

patient에서 더 많음

→ SNP > Disease phenotype or 원인 찾는 것

### Normal population



In a population, a certain percentage will have one version, the rest the other



### Assumption

- Bi-allelic SNPs
- Common ancestors
- Linkage disequilibrium and haplotypes
- Common disease-common variant

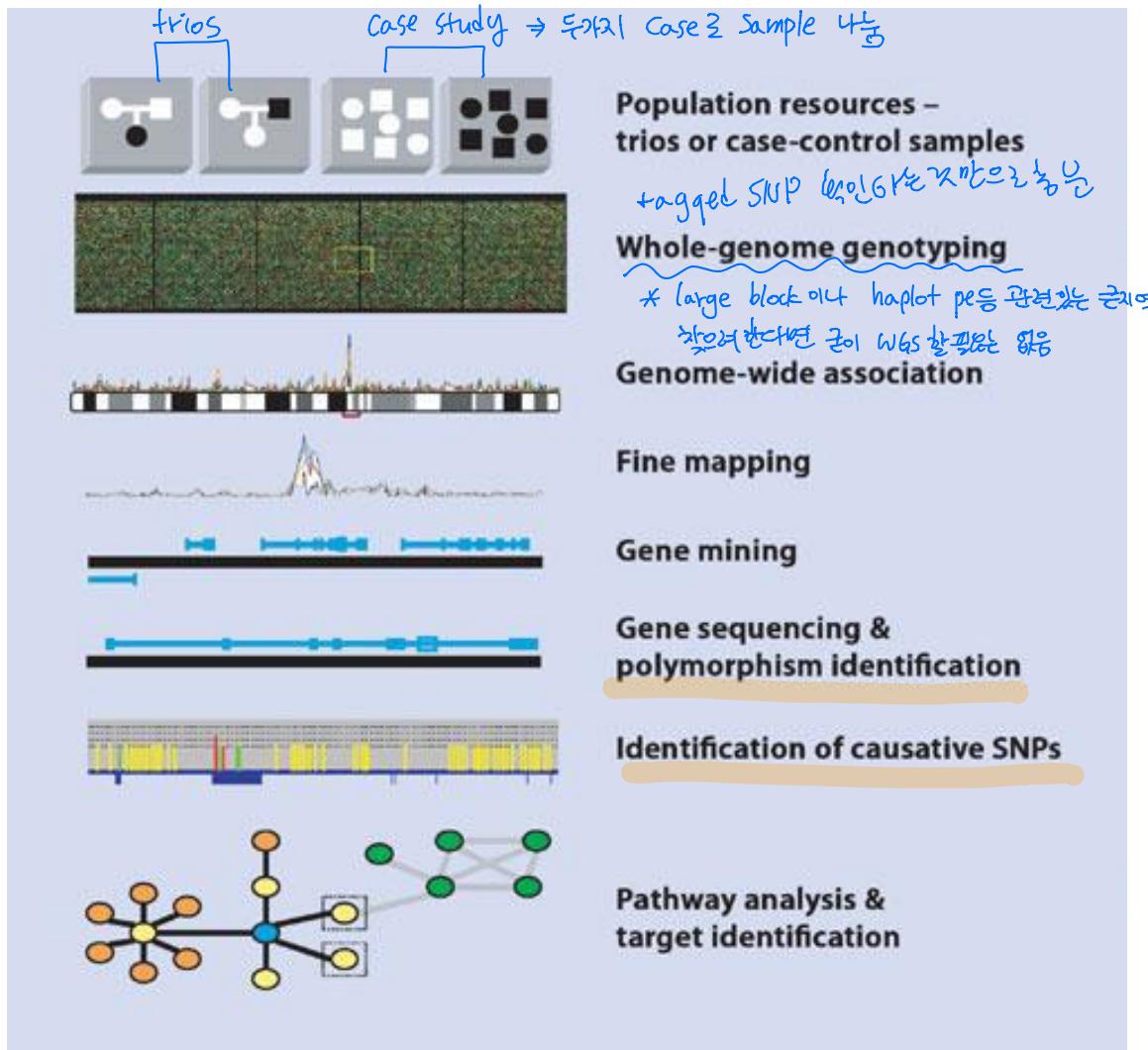
The hypothesis that genetic influences on susceptibility to common diseases are attributable to a limited number of variants present in more than 1% to 5% of the population.

이정도의 낮은 association 가능

→ 이정도가 significant 짜리로 되어야

polygenic event  
→ nothing is tightly linked  
→ 낮은 Association 일정도 중요!

# genome-wide association study (GWAS)

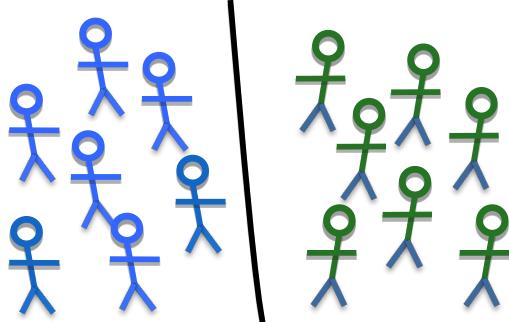


## GWAS:

An examination of genetic variation across a given genome whether it associated with phenotypes (quantitative traits, diseases)

# GWAS

Control Population



Haplotyping  
(tagged SNPs)



Whole genome sequencing (WGS)

NGS sequencing

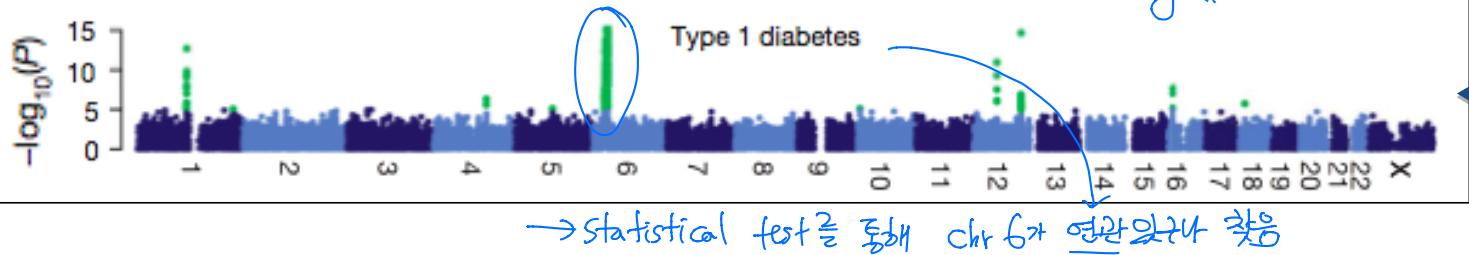
or



D' 등이 매우 낮아 이용하기 어렵기 때문에  
p-value 를 사용

Statistics (p-value) > low p-value > significant  
= Significant

Linkage disequilibrium analysis ( SNP & Phenotype) > Statistics (p-value) > low p-value > significant



# SNP variation & phenotypic change

- **Classification of SNPs**

: transition, transversion, single-base indels

**1. Noncoding SNPs**, 5'UTR, 3'UTR, intron, intergenic regions;

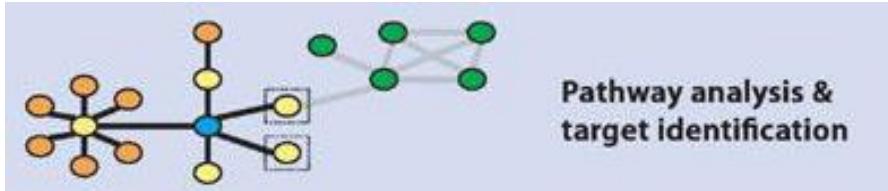
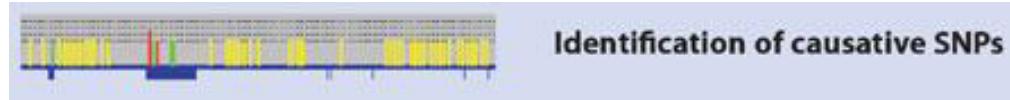
**2. Coding SNPs**,

- nonsynonymous or missense or replacement polymorphism

- Synonymous or sense polymorphism → a.a가 동일한 <sup>기능을 가진다</sup> same function

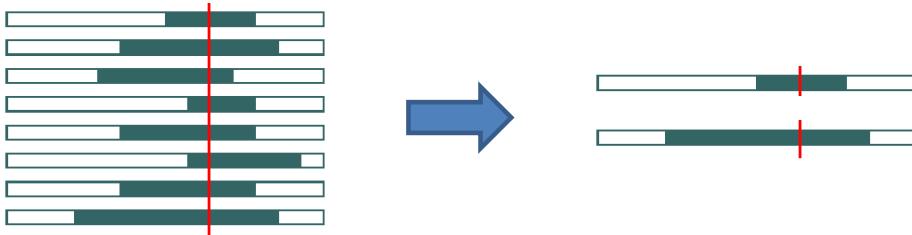
- Regulatory polymorphism : Synonymous and noncoding polymorphism

**3. haplotype**, a distinct combination of single nucleotide types on a single chromosome at a locus. <sup>used to narrow down region</sup>



Need to explain the mechanism how SNP affects the phenotype !!  
(gene expression, regulation)

# Haplotyping : linkage disequilibrium



n=10

catcag aag  
tatcag aac  
tatcag aac  
catcag aag  
catcag aag  
tatcag aac  
catcag aag  
catcag aag  
tatcag aac  
catcag aag

1. Segregation sites ?



6 x c      4 x t

SNP1

cg  
tc  
cg  
cg  
tc  
cg  
cg  
cg  
tc  
cg

2. Minor allele frequencies?

$$\begin{aligned} \text{SNP1} &= 0.4 \\ \text{SNP2} &= 0.5 \end{aligned}$$

3. Expected average heterozygosity (H)?

$$\frac{(0.4 \times 1 + 0.5 \times 1 + 0 \times 7)}{9} = 0.1$$

SNP2

5 x g      5 x c

9 nucleotides

that your sequence results in this block can discriminate haplotypes from your ancestor

# Haplotyping : linkage disequilibrium

A	B
1	1
1	1
1	1
1	1
1	1
1	1
1	1
1	1
1	1
1	1
1	1
1	1
1	1
1	1
1	1
1	1

1	2
2	1
2	1
2	2
2	2
2	2
2	2
2	2
2	2
2	2
2	2
2	2
2	2
2	2
2	2

2	2
2	2
2	2
2	2
2	2

$$D = \frac{24}{16} = 0.15$$

$$D > 0$$

	$B_1$	$B_2$	Total
$A_1$	$p_{11} = p_1q_1 + D$	$p_{12} = p_1q_2 - D$	$p_1$
$A_2$	$p_{21} = p_2q_1 - D$	$p_{22} = p_2q_2 + D$	$p_2$
Total	$q_1$	$q_2$	1

Expected

$$p_1 = 10/16 = 0.625$$

$$q_1 = 11/16 = 0.629$$

$$p_2 = 0.375$$

$$q_2 = 0.371$$

$$p_1q_1 = 0.39$$

$$\frac{6}{16}$$

$$\frac{5}{16}$$

$$\frac{52}{16}$$

observed

	$B_1$	$B_2$	Total
$A_1$	$9/16 = p_1q_1 + D$	$1/16 = p_1q_2 - D$	$p_1$
$A_2$	$2/16 = p_2q_1 - D$	$4/16 = p_2q_2 + D$	$p_2$
Total	$q_1$	$q_2$	1

$$\frac{16}{16^2} = \frac{50}{16^2} - D$$

$$\frac{1}{16} = \frac{50}{16} \cdot \frac{1}{16}$$

$$D = \frac{24}{16}$$

$$D = 0.133$$

$$D_{\max} = \min(0.625 \times 0.371, 0.375 \times 0.629) = 0.23$$

link 되어 있는지 알아내기 위해 random에서 나온 값이 아님을 확인

Whether it is significant ?

# Statistical significance

Null hypothesis ( $H_0$ )



not associated

Independent  $\Rightarrow$  (내가 다른 게 아니면  
특정 결과는 가능)

Random process  
(Permutation)

If we know the null distribution, we just simply used it  
(t-test, chi-square test...)

무언가 있는 D' 가 뉴턴 같음  
있는 것이라면 뉴턴

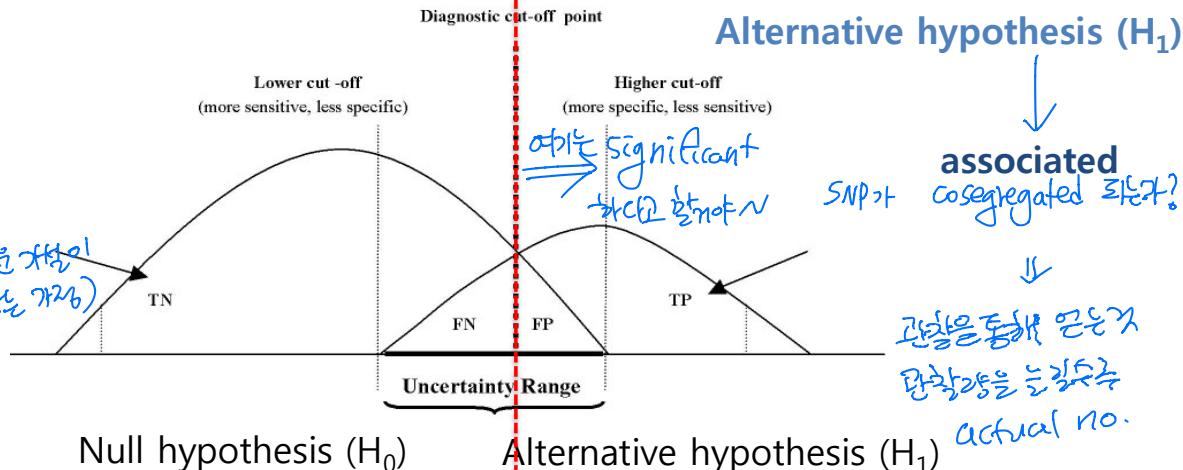
D' 가 낮다고 해도

random으로 있어야 한다면

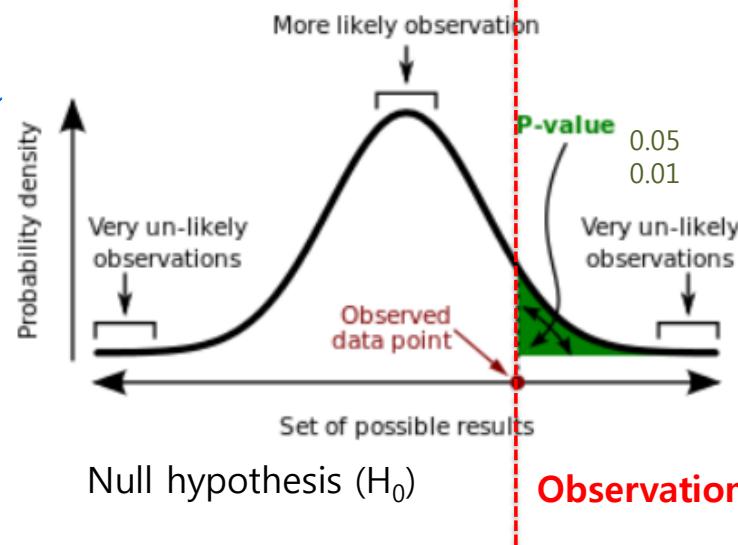
Significant

$\Rightarrow$  disease 두 개

SNP가 있는



Null hypothesis ( $H_0$ )      Alternative hypothesis ( $H_1$ )



# Chi square test

카이 스퀘어

observed

Expected

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

$\chi^2$  = the test statistic      $\sum$  = the sum of

O = Observed frequencies    E = Expected frequencies

catcag aag  
tatcag aac  
tatcag aac  
catcag aag  
catcag aag  
tatcag aac  
catcag aag  
catcag aag  
tatcag aac  
catcag aag

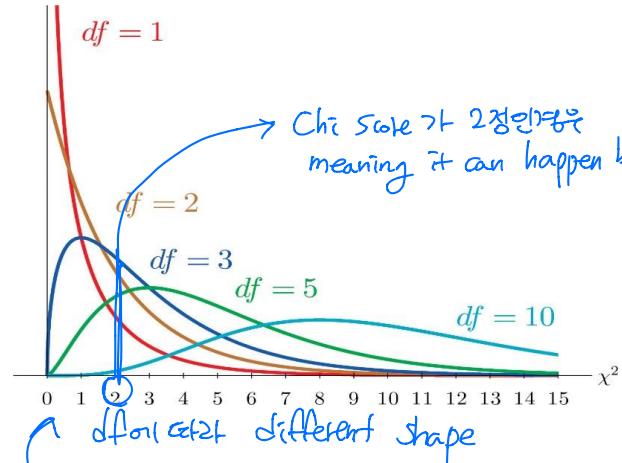
observed

	c	t	total
g	P <sub>cg</sub>	P <sub>tg</sub>	P <sub>g</sub>
c	P <sub>cc</sub>	P <sub>tc</sub>	P <sub>c</sub>
	P <sub>c</sub>	P <sub>t</sub>	1

Expected

	c	t	total
g	P <sub>c</sub> X P <sub>g</sub>	P <sub>t</sub> X P <sub>g</sub>	P <sub>g</sub>
c	P <sub>c</sub> X P <sub>c</sub>	P <sub>t</sub> X P <sub>c</sub>	P <sub>c</sub>
	P <sub>c</sub>	P <sub>t</sub>	1

Chi square test ( p=0.067889)

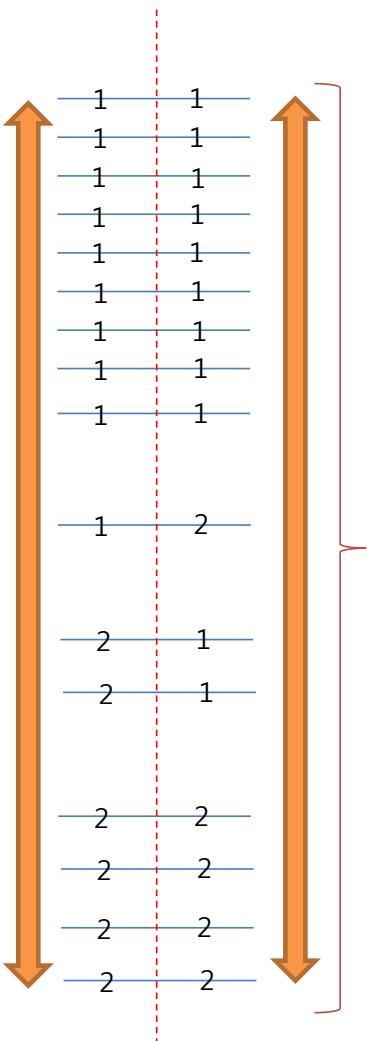


Degrees of freedom (df) = n-1 = 4-1 = 3

large P' value → 무의미  
Chi square test →  
Significant 되지 않았다면  
→ random 브랜드 일정한 결과

# Fisher's exact test

review  
gratul.



observed

	c	t	total
g	P <sub>cg</sub>	P <sub>tg</sub>	P <sub>g</sub>
c	P <sub>cc</sub>	P <sub>tc</sub>	P <sub>c</sub>
	P <sub>c</sub>	P <sub>t</sub>	1

	1	2
1	9	1
2	2	4

Hypergeometric distribution

Shuffling  
(Enumerate)

Table B

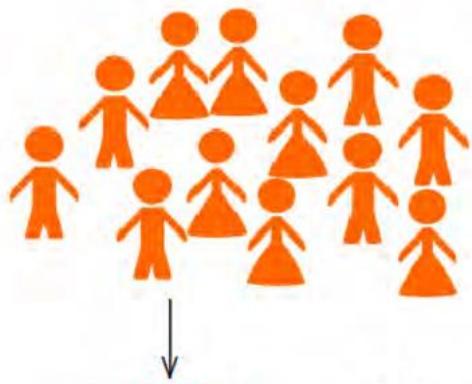
$n_{11}$	$n_{12}$	$n_{21}$	$n_{22}$	$D$	Probability	Cumulative probability
10	0	1	5	0.195	0.001	0.001
9	1	2	4	0.133	0.034	0.035
8	2	3	3	0.070	0.206	0.241
7	3	4	2	0.008	0.412	0.653
6	4	5	1	-0.055	0.288	0.941
5	5	6	0	-0.117	0.058	1.000

D<sub>9,1,2,4</sub>

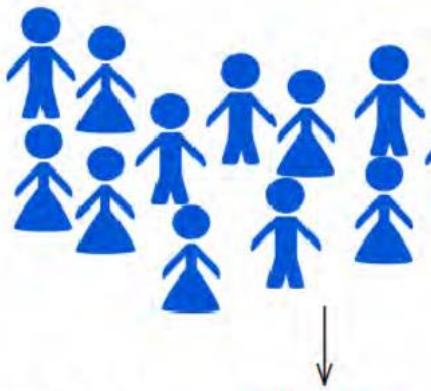
Sampled  
obtained  
2x2

# GWAS: Genome-wide association study

## Affected Individuals

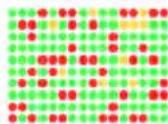


## Unaffected Individuals

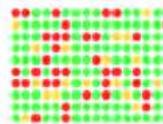


Q/Hyp

sample	seq	type	grade
1	actgtacttc	normal	0
2	actgtacttc	normal	0
3	actgtacttc	normal	0
4	actgtccttc	normal	0
5	actgtccttc	normal	0
6	actgtccttc	tumor	2
7	actgtccttc	tumor	3
8	actctaactta	tumor	2
9	actgtccttc	tumor	3
10	actctaactta	tumor	4



**SNPs analyzed  
and compared  
statistically**



## Distribution of SNPs : Population genetics

- SNPs – phenotypic variation (Quantitative genetics)
- QTL (quantitative trait loci) = polygenic phenotypic variation

\* normal 2%  
patient 3%  $\Rightarrow$  significant  
 $\rightarrow p$  value

Linked?

SNPs



Phenotype

# GWAS (practice)

Observed frequency

patient	A	a	Total
Affected	9	1	10
Unaffected	2	4	6
Total	11	5	16

$$p_1 = 10/16 = 0.625 \text{ Affected}$$

$$p_2 = 0.374 \text{ unaffected}$$

$$q_1 = 11/16 = 0.6875 \text{ A}$$

$$q_2 = 0.3125 \text{ a}$$

$$D = 9/16 - \frac{p_1}{p_1} \times 0.6875 = 0.1328$$

$$D_{\max} = \min(0.625 \times 0.3125, 0.374 \times 0.6875) = 0.195313$$

	$A_1$	$A_2$	Total
$B_1$	$x_{11} = p_1 q_1 + D$	$x_{21} = p_2 q_1 - D$	$q_1$
$B_2$	$x_{12} = p_1 q_2 - D$	$x_{22} = p_2 q_2 + D$	$q_2$
Total	$p_1$	$p_2$	1

	A	a
Affected	$e_{11}$	$e_{21}$
Unaffected	$e_{12}$	$e_{22}$

expected frequency

$$e_{11} = 16 \times p_1 \times q_1 = 6.875$$

$$D > 0 \Rightarrow P_1 q_2, P_2 q_1$$

	A	a	Total
Affected	6.875	3.125	10
Unaffected	4.125	1.875	6
Total	11	5	16

$$D' = 0.1328 / 0.195313 = 0.679934$$

6.6이므로 P Value 가 안좋으면  
그는 linkage가 아님이 알 수가 없음 (Significant하지 않음)

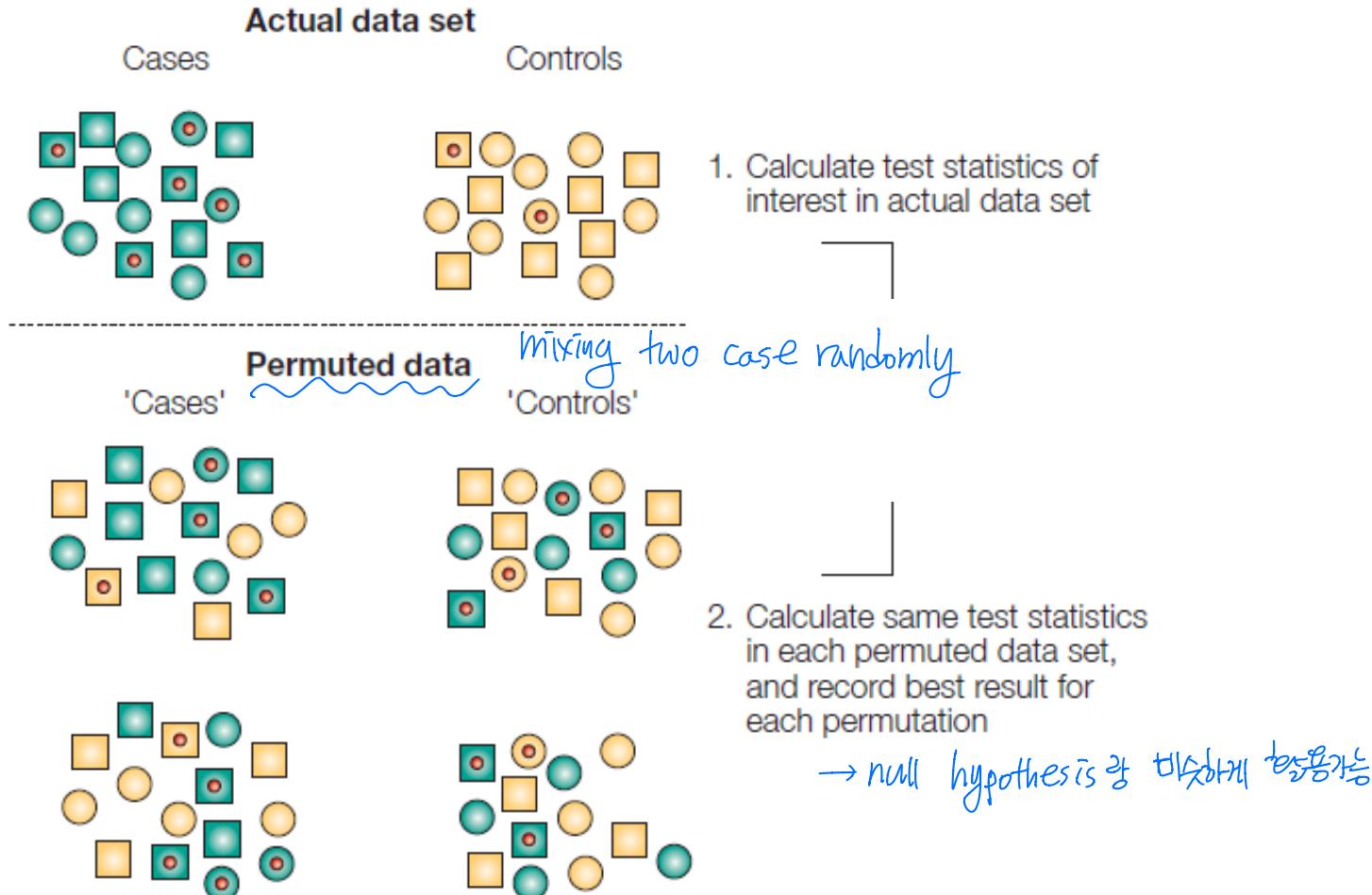
Chi square test (  $p=0.017911$  )  
Fisher's exact test

Significant Linkage  
disequilibrium

(A is significantly  
associated with disease)

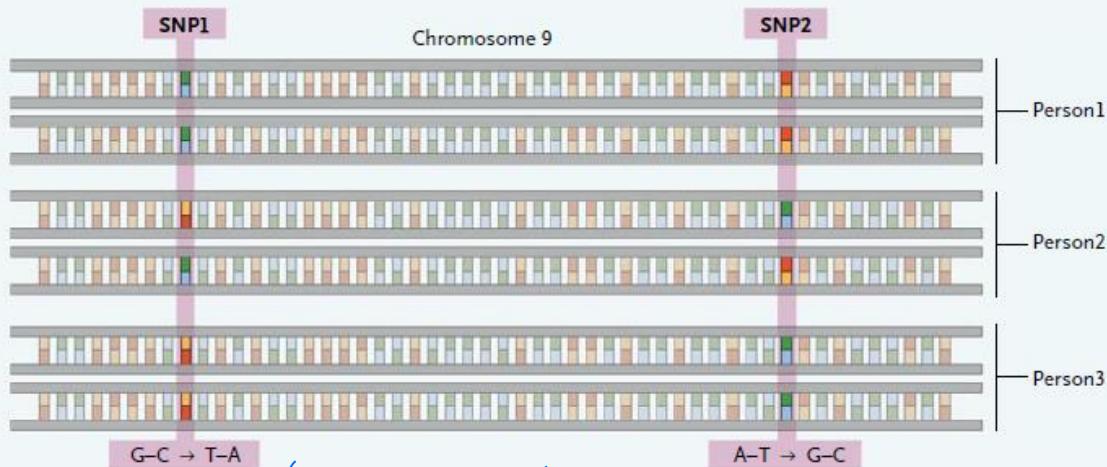
Permutation  
(False-discovery rate)

# Permutation testing

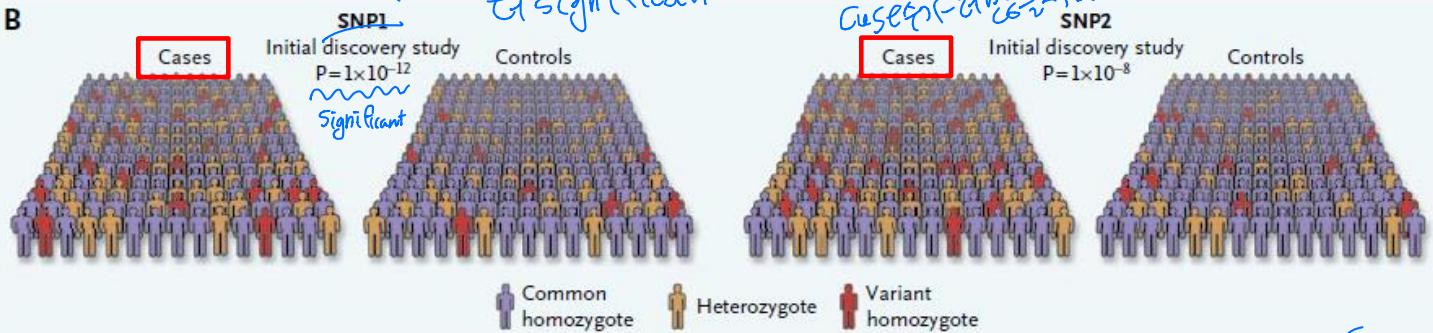


# The Genome-wide Association Study (GWAS)

A



B



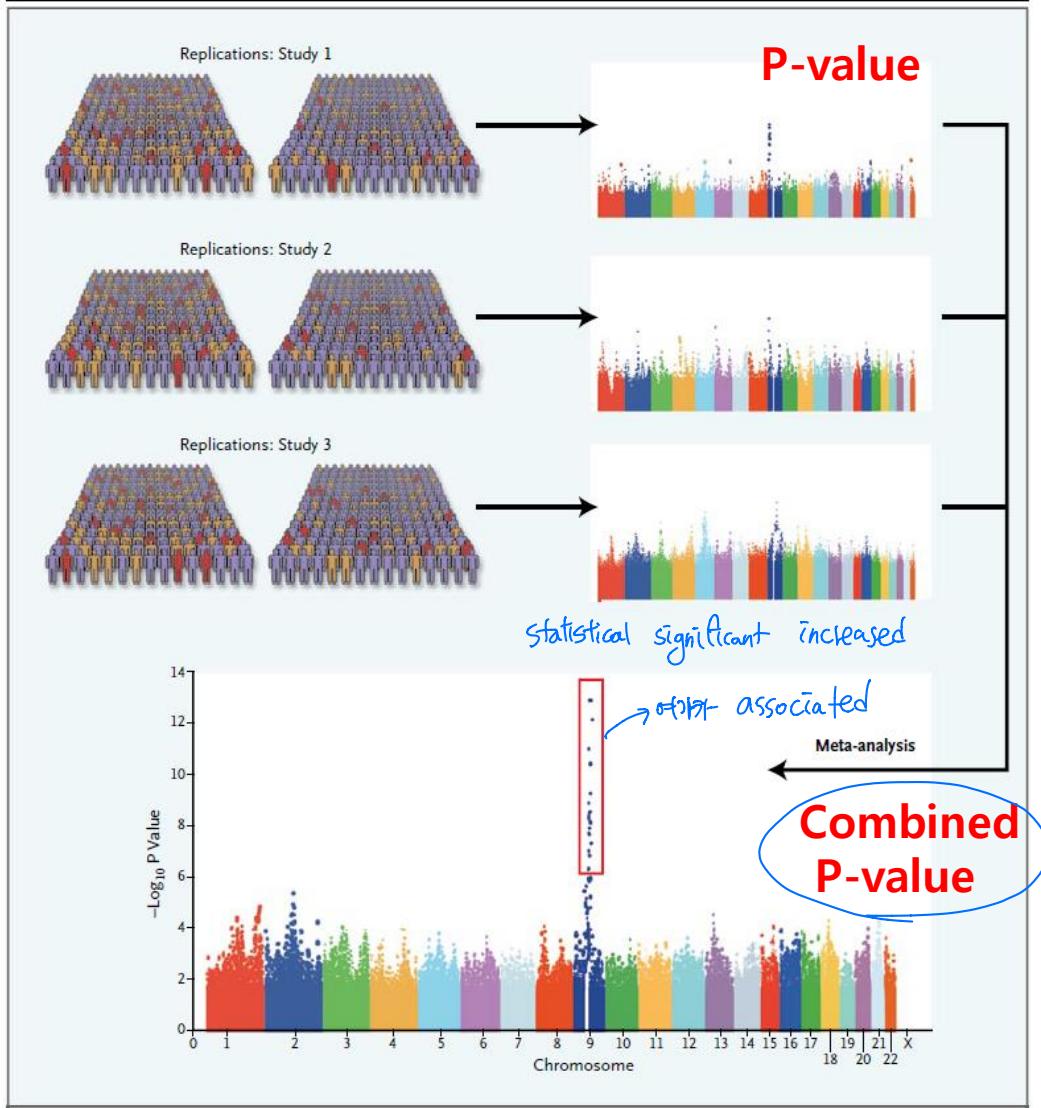
- An examination of genetic variation across a given genome
- Designed to identify genetic associations with observable traits –Such as blood pressure or weight, –or why some people get a disease or condition

frequency /afe 4cf  
associationl 퍼 품 은

# The Genome-wide Association Study (GWAS)

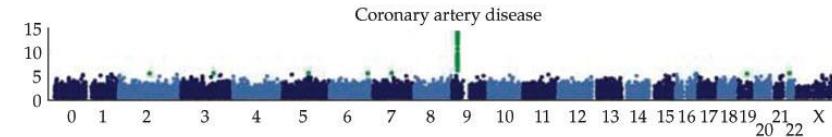
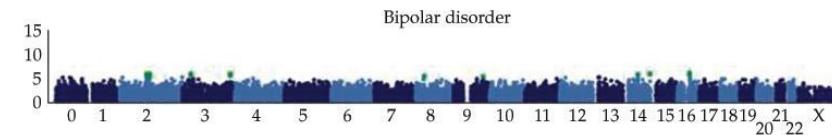
Meta-analysis is required

Statistical significance



# Genome-wide Association Studies (example)

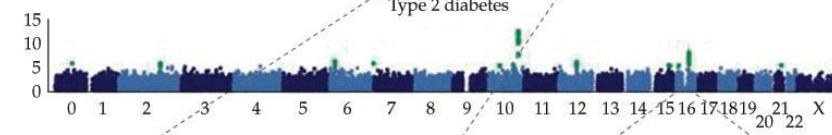
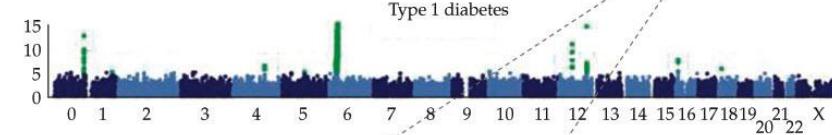
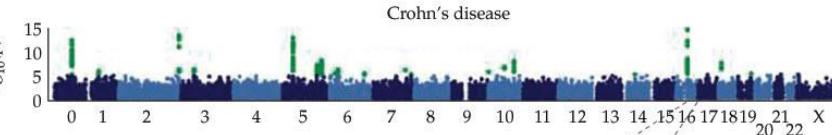
(A)



$$Y^2_i = -(10 \log_{10}(P))$$

Statistical significance

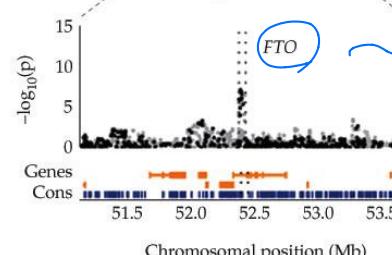
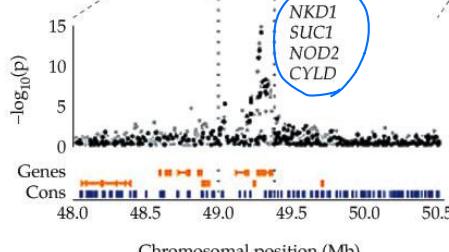
$-\log_{10}(p)$



(B)

→ p value 가 낮을  
수록 보정  
= associated

each dot : SNP  
& have own p-value



각 병에 연관된 Variation