

# Summary & Review

## Genome Sequence Analysis

### 1) Sequence alignment

- Purpose of sequence alignment
- Principles of pair-wise sequence alignment
- Sequence alignment for database search: BLAST

### 2) Sequence variation analysis *haplotyping*

- Pre-required basic knowledge of human genome

Sung Wook Chi

Division of Life Sciences, Korea University

# Biological question from sequence alignment

ACGCTGA

Common  
Ancestor

ACTGT

Evolutional Changes

ACGCTGA

ACGCTGA  
A--CTGT

Sequence  
alignment 1

ACTGT

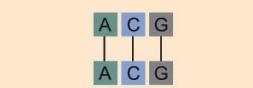
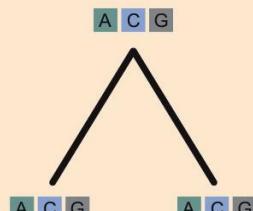
ACGCTGA

ACGCTGA  
ACTGT--

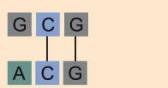
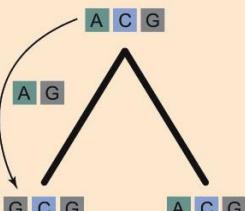
Sequence  
alignment 2

ACTGT

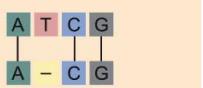
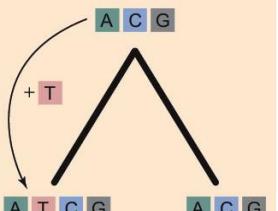
(A) Identity



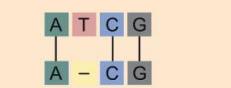
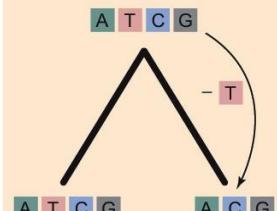
(B) Substitution



(C) Insertion



(D) Deletion



What is the best  
optimal alignment?

# Sequence Alignment

## Pairwise sequence alignment

|            | match   | Alignment 1 | Alignment 2 |
|------------|---|-------------|-------------|
| Sequence 1 | ACGCTGA   | ACGCTGA     |             |
| Sequence 2 | A--CTGT   |             | ACTGT--     |
|            |  | mismatch    |             |

Gap opening, Gap extension

## Optimal sequence alignment

1. Evaluation of sequence similarity : Distance, Score
2. Performing optimal sequence alignment search: Dynamic Programming

### Global Alignment

```
TCAG -- T - G T C G A A G T - T A  
| | | | | | | | | | | | | | | | | | | |  
T - A G G C T A G - C - A - G T G T A
```

### Global alignment

Conserved region of sequence  
> Functional domain / element

ex) kinase

### Local Alignment

```
TCAG T G T C G A A G T T A  
| | | | | | | | | | | | | | | | | | | |  
T A G G C T A G C A G T G T A
```

### Local alignment

High sequence similarity > Homolog > same function  
Similar seq and

# How to measure sequence similarity : Score

**Score** = (match or mismatch penalty) + gap penalty

| C | S  | T  | P  | A  | G  | N  | D  | E  | Q  | H  | R  | K  | M  | I  | L  | V  | F  | Y | W |   |
|---|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|---|---|---|
| C | 9  |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |   |   |   |
| S | -1 | 4  |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |   |   |   |
| T | -1 | 5  |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |   |   |   |
| P | -1 | -1 | 7  |    |    |    |    |    |    |    |    |    |    |    |    |    |    |   |   |   |
| A | 1  | 0  | -1 | 4  |    |    |    |    |    |    |    |    |    |    |    |    |    |   |   |   |
| G | 0  | -2 | -2 | 0  | 6  |    |    |    |    |    |    |    |    |    |    |    |    |   |   |   |
| N | -3 | 1  | 0  | -2 | -2 | 0  | 6  |    |    |    |    |    |    |    |    |    |    |   |   |   |
| D | -3 | 0  | -1 | -1 | -2 | -1 | 1  | 6  |    |    |    |    |    |    |    |    |    |   |   |   |
| E | -4 | 0  | -1 | -1 | -1 | -2 | 0  | 2  | 5  |    |    |    |    |    |    |    |    |   |   |   |
| Q | -3 | 0  | -1 | -1 | -1 | -2 | 0  | 0  | 2  | 5  |    |    |    |    |    |    |    |   |   |   |
| H | -3 | -1 | -2 | -2 | -2 | 1  | -1 | 0  | 0  | 8  |    |    |    |    |    |    |    |   |   |   |
| R | -3 | -1 | -1 | -2 | -1 | -2 | 0  | -2 | 0  | 1  | 0  | 5  |    |    |    |    |    |   |   |   |
| K | -3 | 0  | -1 | -1 | -1 | -2 | 0  | -1 | 1  | -1 | 2  | 5  |    |    |    |    |    |   |   |   |
| M | -1 | -1 | -1 | -2 | -1 | -3 | -2 | -3 | -2 | 0  | -2 | -1 | 1  | 5  |    |    |    |   |   |   |
| I | -1 | -1 | -3 | -1 | -4 | -3 | -3 | -3 | -3 | -3 | -3 | -3 | 1  | 4  |    |    |    |   |   |   |
| L | -1 | -2 | -1 | -3 | -1 | -4 | -3 | -4 | -3 | -2 | -3 | -2 | 2  | 2  | 4  |    |    |   |   |   |
| V | -1 | -2 | 0  | -2 | 0  | -3 | -3 | -2 | -2 | -3 | -3 | -2 | 1  | 3  | 1  | 4  |    |   |   |   |
| F | -2 | -2 | -2 | -4 | -2 | -3 | -3 | -3 | -3 | -3 | -3 | -3 | 0  | 0  | -1 | 6  |    |   |   |   |
| Y | -2 | -2 | -2 | -3 | -2 | -3 | -2 | -1 | -2 | -2 | -1 | -1 | -1 | -1 | -1 | 3  | 7  |   |   |   |
| W | -2 | -3 | -2 | -4 | -3 | -2 | -4 | -4 | -3 | -2 | -2 | -3 | -3 | -1 | -3 | -2 | -3 | 1 | 2 |   |
|   | C  | S  | T  | P  | A  | G  | N  | D  | E  | Q  | H  | R  | K  | M  | I  | L  | V  | F | Y | W |

Match = 3 , Mismatch = -1

|   | A  | G  | T  | C  |
|---|----|----|----|----|
| A | 20 | 10 | 5  | 5  |
| G | 10 | 20 | 5  | 5  |
| T | 5  | 5  | 20 | 10 |
| C | 5  | 5  | 10 | 20 |

transition : 표준  $\rightarrow$  표준  
transversion : 표준  $\rightarrow$  비표준

DOROTHY-----HODGKIN

↳ amino acid

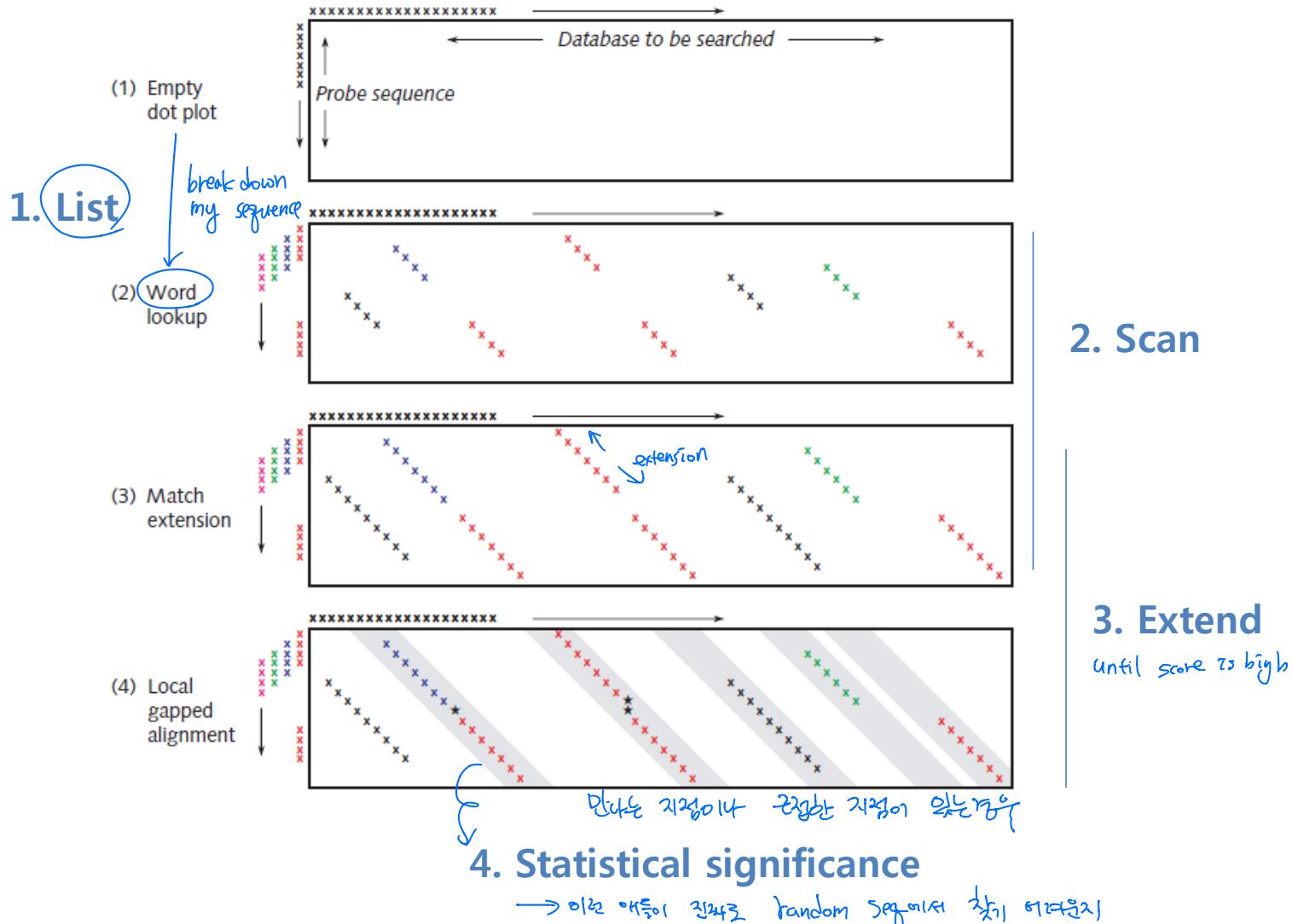
- Dot Plot : seq 형태를 비교하는데 유용



Repetitive sequence

| A | B | R | A | C | A | D | A | B | R | A | C | A | D | A | B | R | A | C | A     | D | A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|-------|---|---|
| B | A | R | A | C | A | D | A | B | R | A | C | A | D | A | B | R | A | C | A     | D | A |
| R | B | A | R | A | C | A | D | B | R | A | C | A | D | B | R | A | C | A | D     | B | R |
| A | R | B | A | R | A | C | A | B | R | A | C | A | D | B | R | A | C | A | D     | B | R |
| C | A | R | B | A | R | A | C | B | R | A | C | A | D | B | R | A | C | A | D     | B | R |
| A | C | A | R | B | A | R | A | C | B | R | A | C | A | D | B | R | A | C | A     | D | B |
| D | A | C | A | R | B | A | R | C | A | R | A | C | A | D | B | R | A | C | A     | D | B |
| A | B | R | A | C | A | D | A | B | R | A | C | A | D | B | R | A | C | A | D     | B | R |
| B | R | A | C | A | D | A | B | R | A | C | A | D | B | R | A | C | A | D | B     | R | A |
| R | A | C | A | D | A | B | R | A | C | A | D | B | R | A | C | A | D | B | R     | A | C |
| A | C | A | D | B | R | A | C | A | D | B | R | A | C | B | R | A | C | A | D     | B | R |
| C | A | D | B | R | A | C | A | D | B | R | A | C | B | R | A | C | A | D | B     | R | A |
| A | D | B | R | A | C | B | R | A | C | B | R | A | C | B | R | A | C | B | R     | A | C |
| D | A | C | B | R | A | D | B | R | A | C | B | R | A | C | B | R | A | C | B     | R | A |
| A | C | B | R | A | D | B | R | C | A | D | B | R | A | C | B | R | A | D | B     | R | A |
| B | R | A | C | B | R | A | D | C | R | A | C | B | R | A | C | B | R | A | D     | B | R |
| R | A | C | B | R | A | D | C | R | A | C | B | R | A | C | B | R | A | D | B     | R | A |
| A | C | B | R | A | D | C | R | A | C | B | R | A | C | B | R | A | D | B | R     | A | C |
| C | A | D | B | R | A | C | B | R | C | A | D | B | R | A | C | B | R | A | D     | B | R |
| A | D | B | R | A | C | B | R | C | R | A | D | B | R | A | C | B | R | A | D     | B | R |
| B | R | A | C | B | R | A | D | C | R | A | C | B | R | A | C | B | R | A | D     | B | R |
| R | A | C | B | R | A | D | C | R | A | C | B | R | A | C | B | R | A | D | B     | R | A |
| A | C | B | R | A | D | C | R | A | C | B | R | A | C | B | R | A | D | B | R     | A | C |
| D | A | C | B | R | A | D | C | R | A | C | B | R | A | C | B | R | A | D | B     | R | A |
| A | C | B | R | A | D | C | R | C | R | A | D | B | R | A | C | B | R | A | D     | B | R |
| C | A | D | B | R | A | C | B | R | C | R | A | D | B | R | A | C | B | R | A     | D | B |
| A | D | B | R | A | C | B | R | C | R | A | D | B | R | A | C | B | R | A | D     | B | R |
| B | R | A | C | B | R | A | D | C | R | A | D | B | R | A | C | B | R | A | D     | B | R |
| R | A | C | B | R | A | D | C | R | A | D | B | R | A | C | B | R | A | D | B     | R | A |
| A | C | B | R | A | D | C | R | C | R | A | D | B | R | A | C | B | R | A | D     | B | R |
| D | A | C | B | R | A | D | C | R | C | R | A | D | B | R | A | C | B | R | A     | D | B |
| A | C | B | R | A | D | C | R | C | R | A | D | B | R | A | C | B | R | A | D     | B | R |
| C | A | D | B | R | A | C | B | R | C | R | A | D | B | R | A | C | B | R | A     | D | B |
| A | D | B | R | A | C | B | R | C | R | A | D | B | R | A | C | B | R | A | D     | B | R |
| B | R | A | C | B | R | A | D | C | R | A | D | B | R | A | C | B | R | A | D     | B | R |
| R | A | C | B | R | A | D | C | R | A | D | B | R | A | C | B | R | A | D | B     | R | A |
| A | C | B | R | A | D | C | R | C | R | A | D | B | R | A | C | B | R | A | D     | B | R |
| D | A | C | B | R | A | D | C | R | C | R | A | D | B | R | A | C | B | R | A     | D | B |
| A | C | B | R | A | D | C | R | C | R | A | D | B | R | A | C | B | R | A | D     | B | R |
| C | A | D | B | R | A | C | B | R | C | R | A | D | B | R | A | C | B | R | A     | D | B |
| A | D | B | R | A | C | B | R | C | R | A | D | B | R | A | C | B | R | A | D     | B | R |
| B | R | A | C | B | R | A | D | C | R | A | D | B | R | A | C | B | R | A | D     | B | R |
| R | A | C | B | R | A | D | C | R | A | D | B | R | A | C | B | R | A | D | B     | R | A |
| A | C | B | R | A | D | C | R | C | R | A | D | B | R | A | C | B | R | A | D     | B | R |
| D | A | C | B | R | A | D | C | R | C | R | A | D | B | R | A | C | B | R | A     | D | B |
| A | C | B | R | A | D | C | R | C | R | A | D | B | R | A | C | B | R | A | D     | B | R |
| C | A | D | B | R | A | C | B | R | C | R | A | D | B | R | A | C | B | R | A     | D | B |
| A | D | B | R | A | C | B | R | C | R | A | D | B | R | A | C | B | R | A | D     | B | R |
| B | R | A | C | B | R | A | D | C | R | A | D | B | R | A | C | B | R | A | D     | B | R |
| R | A | C | B | R | A | D | C | R | A | D | B | R | A | C | B | R | A | D | B     | R | A |
| A | C | B | R | A | D | C | R | C | R | A | D | B | R | A | C | B | R | A | D     | B | R |
| D | A | C | B | R | A | D | C | R | C | R | A | D | B | R | A | C | B | R | A</th |   |   |

# Sequence Database Search: BLAST



## EXERCISE 2.3 Perform a BLAST search

IL13

Using the gene that you identified in Exercise 1.1 perform a BLASTn search for homologs in other species. What conclusions can you reach regarding whether this gene is part of a gene family, or regarding its phylogenetic distribution?

**ANSWER:** Continuing with our example of the IL13 gene from Chapter 1, copy the transcript sequence of this gene either from the Ensembl or GenBank page. (For Ensembl, click on the “Transcript info” link and the sequence appears at the bottom of the page. For GenBank, first select “Nucleotide”; then type “IL13” into the search box and follow the link to the mRNA annotation for NM\_002188.) Return to the NCBI home page, click on “BLAST,” bring up the “Nucleotide-nucleotide (blastn)” page, and paste your sequence into the search box. Choose some search options (for example, “all RefSeq mRNA sequences”) and submit the request. Depending on the number of queries in the queue, this could take a minute or two. In the case of IL13, alignments are shown for human, chimp, macaque, canine, and equine matches to the 1280 nt search. If you select optimization for somewhat similar sequences rather than highly similar ones, very high match scores are also seen for several other mammalian IL13 loci. The closest match for the Drosophila genome, however, has an E-value of just 0.17, and even this score is due to a single perfect match—a 23-bp element—that is probably not indicative of a homologous gene. Thus IL13 appears to be unique to vertebrates (and possibly even to mammals, since no fish or bird sequences have significant E-values), and to be present in a single copy per genome.

NCBI 웹사이트  
mRNA 주제  
FASTA 형식으로  
Web BLAST  
검색 키워드.

homolog 같은 종의 생물

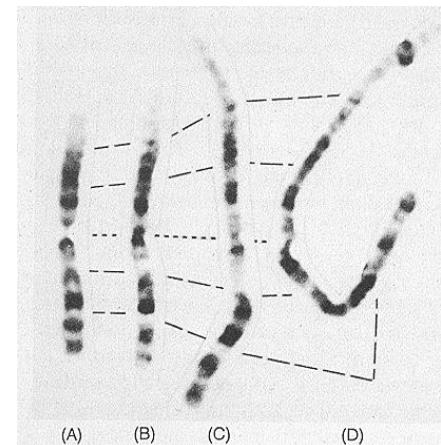
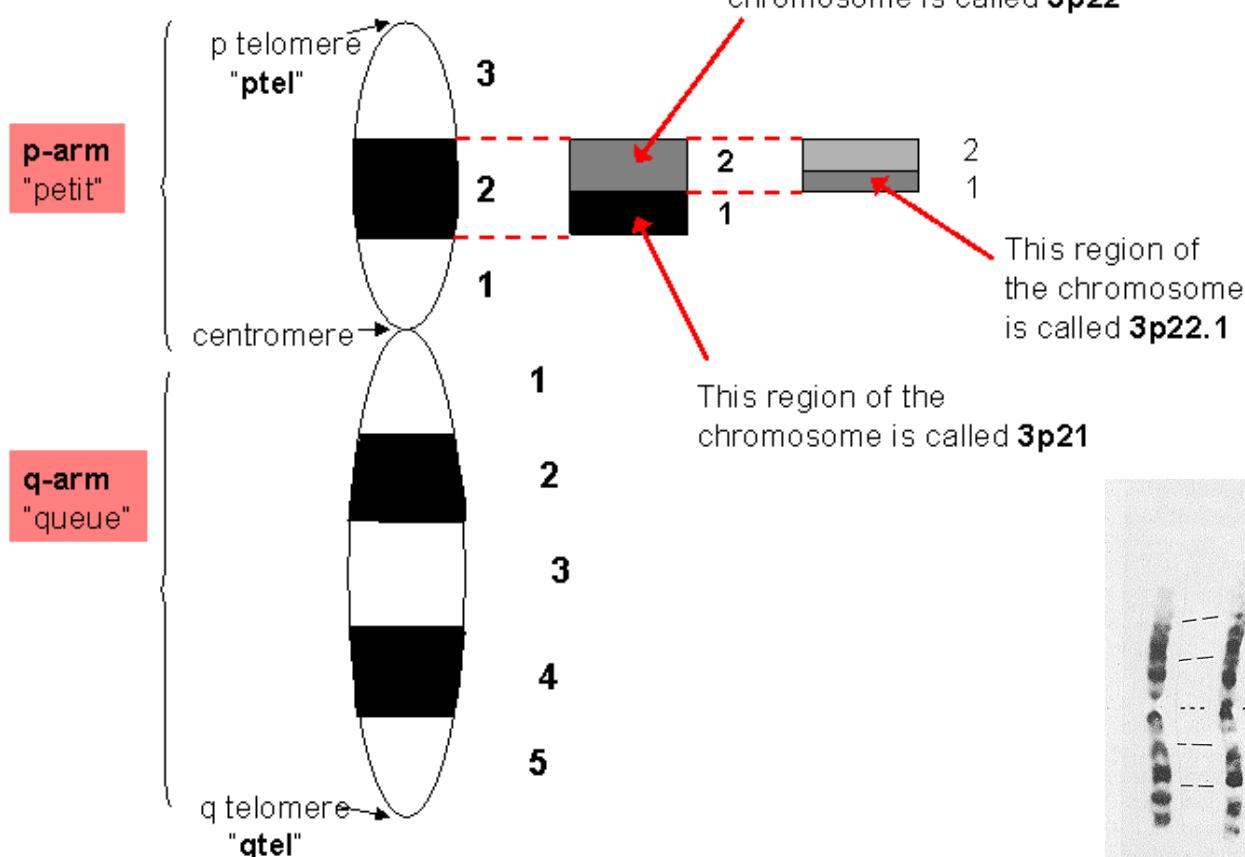
# What we will learn today

## Genome Annotation

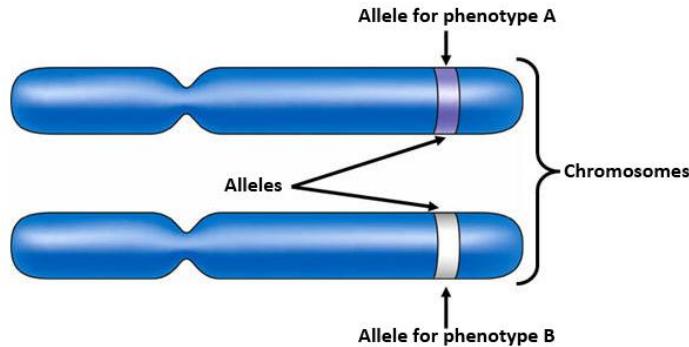
- Sequence alignment
- Linkage disequilibrium (LD), Haplotype, SNP
  - => 표현형을 놓고서 다른 세기들 간에 linkage 를 찾는다 (SNP 활용)
- > Study of genomic variation  
(will be covered in next lecture)

# Nomenclature of human chromosome

## Chromosome 3:



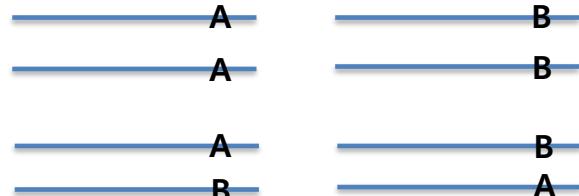
# Human Genetics : recombination, linkage



- Allele, Genotype, Phenotype, Trait

- Homozygote, Heterozygote

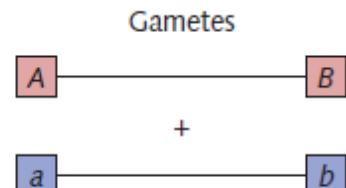
$AA / BB$        $AB / BA$



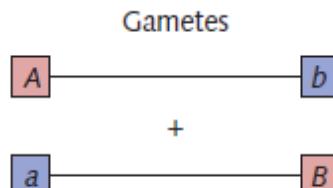
Parental genotype (diploid)



Meiosis



No recombination



Recombination

**Segregation** The separation of corresponding alleles during the reproductive process.

**Independent assortment** The uncorrelated choices of genes for different characters that each parent transmits to children.

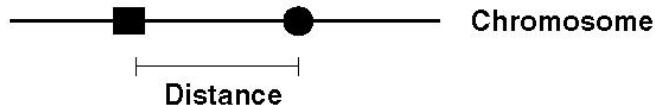
**Linkage** Absence or reduction of independent assortment of parental genes, which are usually transmitted together because they lie on the same chromosome.

In Same Segregation block (Chromosome)

서로 다른 짐수의 allele 들의 축정된 분포가 같았으  
→ allele 조건으로 전달됨. 기반으로, 텐트, 레이, 기린을 가리킨다. FBT

# Linkage disequilibrium (LD)

RECOMBINATION RATE:



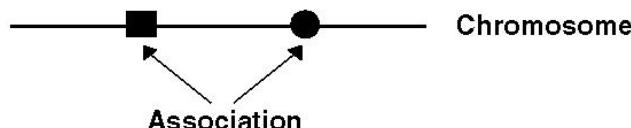
genetic map 2126 association

Mendelian inheritance

:alleles segregate independently

→ to understand gene block

LINKAGE DISEQUILIBRIUM:

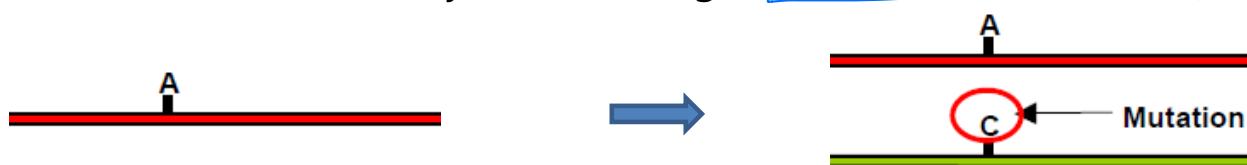


Hardy-Weinberg disequilibrium

Two markers are in **linkage disequilibrium** when alleles at two or more loci **do not segregate independently**

Alleles that exist today arose through ancient mutation events

(at least inherited from parent)

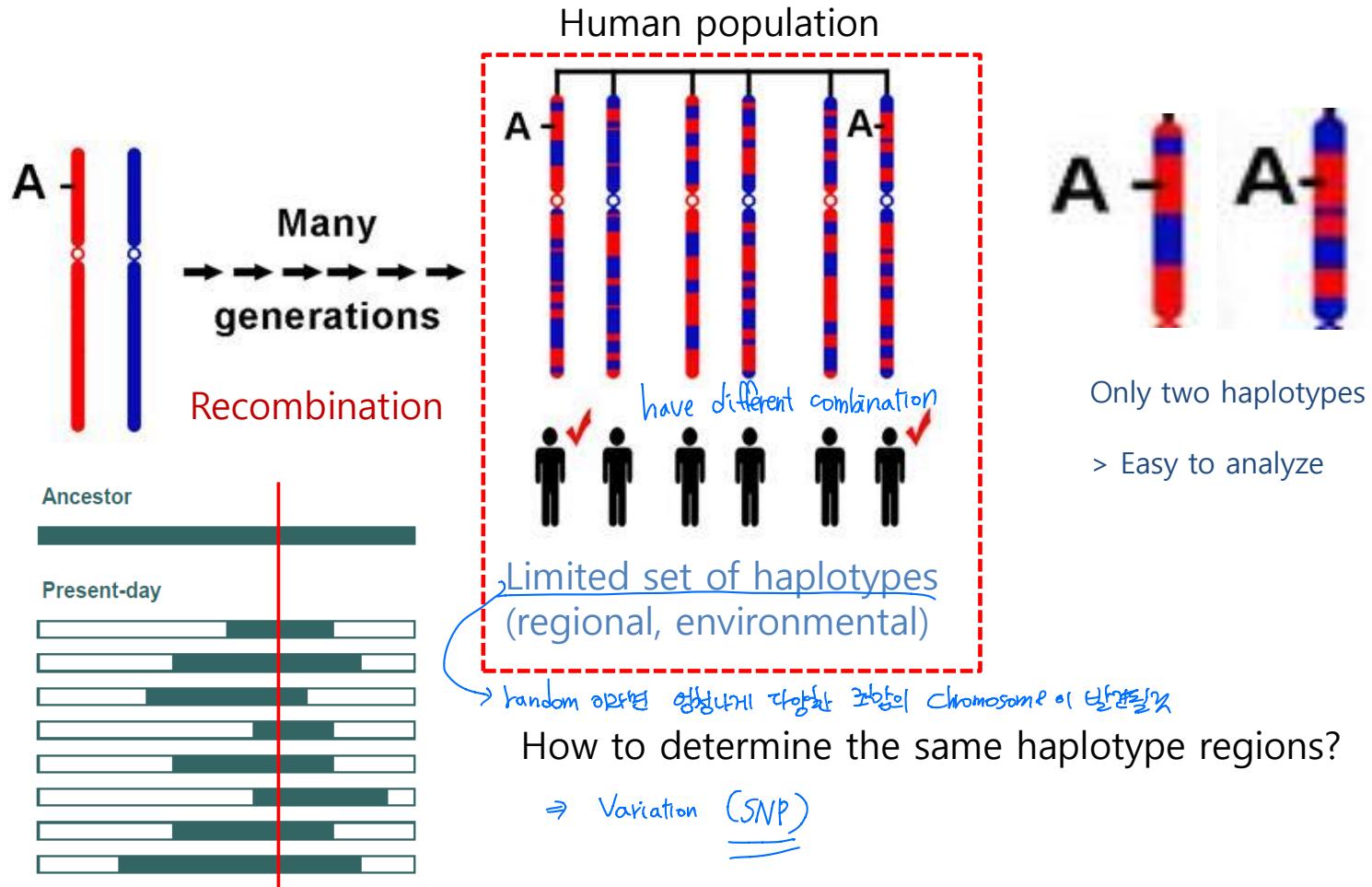


Recombination generates new arrangements for ancestral alleles



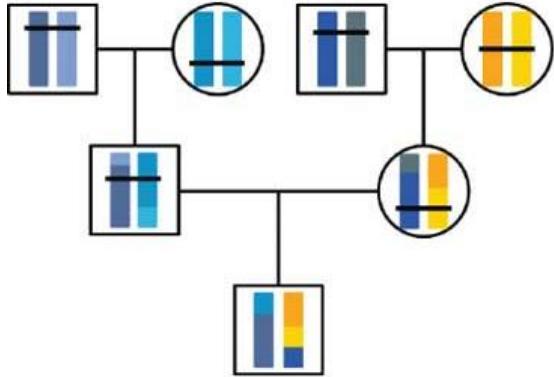
# haplotype

Haplotype: group of genes (DNA regions) in chromosome that are inherited (segregated) together from a single parent during recombination

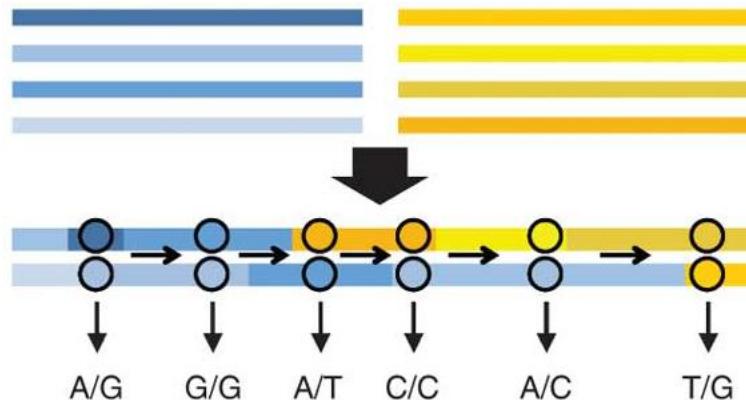


# Tagged SNP

## Ancestors of variable ancestry



## Sampled admixed individual



0.1%의 SNP의 밀도 -> 약 1kb 당 1개의 SNP

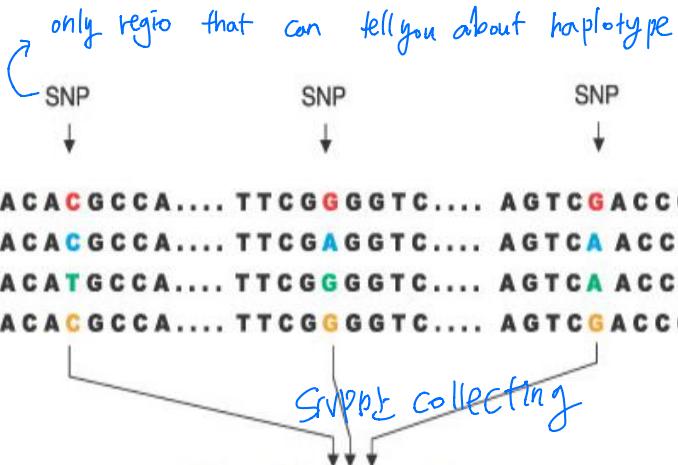
특정 지역은 100kb당 5개 이내의 SNP -> recombination이 거의 일어나지 않는 지역 => haplotype

### a SNPs

Chromosome 1: AACAC**G**CCA.... TTG**G**GGGTc.... AGTC**G**ACCG....  
 Chromosome 2: AACAC**G**CCA.... TTG**G**AGGTc.... AGTC**A**ACCG....  
 Chromosome 3: AACAT**G**CCA.... TTG**G**GGGTc.... AGTC**A**ACCG....  
 Chromosome 4: AACAC**G**CCA.... TTG**G**GGGTc.... AGTC**G**ACCG....

### b Haplotypes

Haplotype 1: **C**TCAAAAGTACGGTT**C**AGGCA  
 Haplotype 2: **T**TGATT**T**GCGCAACAGTAATA  
 Haplotype 3: **C**CCGAT**T**GTGATA**T**CTGGTG  
 Haplotype 4: **T**CGATT**C**CGCGGGTT**C**AGACA



SNP만 있는 것

focus on a few  
SNP

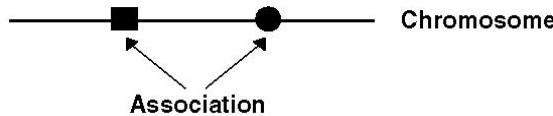
### Tagged SNP

(많은 SNP 중에 haplotype을  
figure out 하기 힘든 매듭)

Only 4 haplotypes  
: discriminate by 3 tagged SNPs

# Linkage disequilibrium (LD)

LINKAGE DISEQUILIBRIUM:



Sequencing results (population)

Observed frequency  
of co-occurrence

Two SNPs

SNP[1] haplotype은 어떤지 알수  
→ SNP[1]과 같은 recombination이  
있을 때

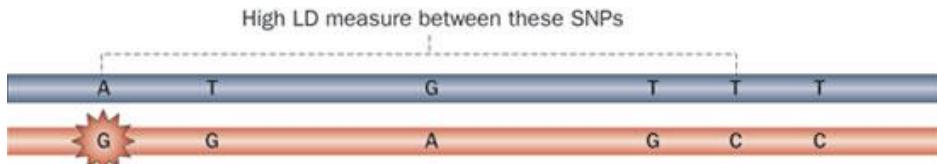
Expected frequency  
of co-occurrence

linkage disequilibrium  
가지 않음

How much it becomes  
disequilibrium ?

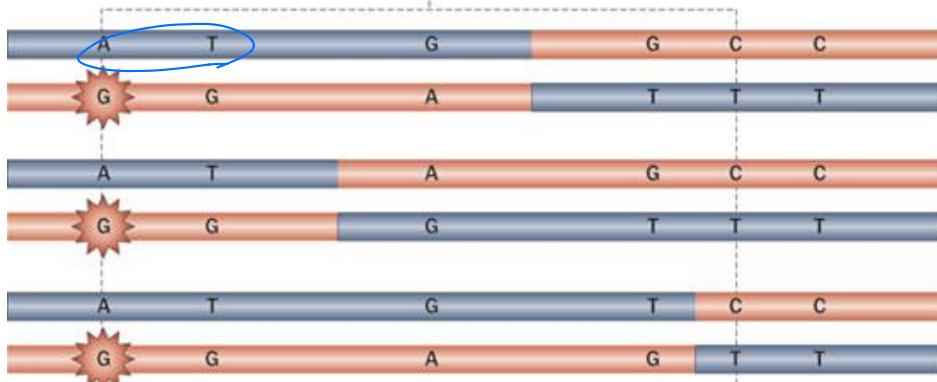
Completely linked

Linked



Recombination events in the population,  
due to chromosomal crossing-over

Lower LD measure between the SNPs  
Reduction in size of haplotype blocks



Completely independent

observed

= expected

# Calculation of linkage disequilibrium

To understand the calculation of linkage disequilibrium consider following example

Suppose there are two genes on Chromosome 5 of apple, each with two alleles

linkage disequilibrium?

randomly  $\leftarrow$  ACTGGTAT.....GATCAACCAG  
G-A linkage?  $\rightarrow$  ACTCGTAT.....GATCAACCAG  
ACTCGTAT.....GATCATCCAG

if independent  
 $\Rightarrow \frac{1}{3}(G) \times \frac{2}{3}(A)$

$$= \frac{2}{9}$$

G-A allele frequency  
is not independent

SNP1                    SNP2

Showing only alleles for both SNPs

| Alleles  | SNP1 | SNP2 |
|----------|------|------|
| Allele 1 | G    | A    |
| Allele 2 | C    | T    |

SNP1                    SNP2  
① 1/2 p1 p2  
② 1/2 q1 q2

## 1. Observed haplotype frequencies

| Haplotype | Frequency |
|-----------|-----------|
| $A_1B_1$  | $x_{11}$  |
| $A_1B_2$  | $x_{12}$  |
| $A_2B_1$  | $x_{21}$  |
| $A_2B_2$  | $x_{22}$  |

| Allele | Frequency               |
|--------|-------------------------|
| $A_1$  | $p_1 = x_{11} + x_{12}$ |
| $A_2$  | $p_2 = x_{21} + x_{22}$ |
| $B_1$  | $q_1 = x_{11} + x_{21}$ |
| $B_2$  | $q_2 = x_{12} + x_{22}$ |

| SNP1   |                     | SNP2   |                     |
|--------|---------------------|--------|---------------------|
| Allele | Frequency           | Allele | Frequency           |
| G      | $p_1 = \frac{1}{3}$ | A      | $q_1 = \frac{2}{3}$ |
| C      | $p_2 = \frac{2}{3}$ | T      | $q_2 = \frac{1}{3}$ |

# Calculation of linkage disequilibrium

| SNP1   |           | SNP2   |           |
|--------|-----------|--------|-----------|
| Allele | Frequency | Allele | Frequency |
| G      | $p_1$     | A      | $q_1$     |
| C      | $p_2$     | T      | $q_2$     |

| Haplotype | Frequency |
|-----------|-----------|
| $A_1B_1$  | $x_{11}$  |
| $A_1B_2$  | $x_{12}$  |
| $A_2B_1$  | $x_{21}$  |
| $A_2B_2$  | $x_{22}$  |

When haplotype frequencies are equal to the product of their corresponding allele frequencies, it means the loci are in linkage equilibrium

observed

expected

| Haplotype frequency     |   | Product of allelic frequency |
|-------------------------|---|------------------------------|
| $p_{11} GA \frac{1}{3}$ | = | $p_1 q_1 \frac{2}{9}$        |
| $p_{12} GT 0$           | = | $p_1 q_2 \frac{1}{9}$        |
| $p_{21} CA \frac{1}{3}$ | = | $p_2 q_1 \frac{4}{9}$        |
| $p_{22} CT \frac{1}{3}$ | = | $p_2 q_2 \frac{2}{9}$        |

$$\begin{aligned} D &= \frac{1}{9} \\ D &= \frac{1}{9} \\ D &= \frac{1}{9} \\ D &= \frac{1}{9} \end{aligned}$$

observed frequency가 0 이될 수 X  
minimum

|       | $A_1$                  | $A_2$                  | 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     |

$$0.2+D$$

$$0.2-D$$

$$0.3-D$$

$$0.3+D$$

$$\Rightarrow D \max 0.2$$

| Allele | Frequency               |
|--------|-------------------------|
| $A_1$  | $p_1 = x_{11} + x_{12}$ |
| $A_2$  | $p_2 = x_{21} + x_{22}$ |
| $B_1$  | $q_1 = x_{11} + x_{21}$ |
| $B_2$  | $q_2 = x_{12} + x_{22}$ |

SNP1

SNP2

$$\begin{aligned} G & A = \frac{2}{4} \Rightarrow 0.5 \\ G & A \Rightarrow \frac{2}{4} \times \frac{3}{4} = 0.375 \\ \text{최소가 } 0.125, & \text{이미 어지럽이나니 } \\ 0.2 \text{랑 } 0.125 \text{랑 } & 100 \\ 100 & \downarrow \text{한계점} \\ 0 & 0 \text{이 같았을} \end{aligned}$$

# Linkage Disequilibrium(LD)

$D_{\max}$ : maximum possible value of D

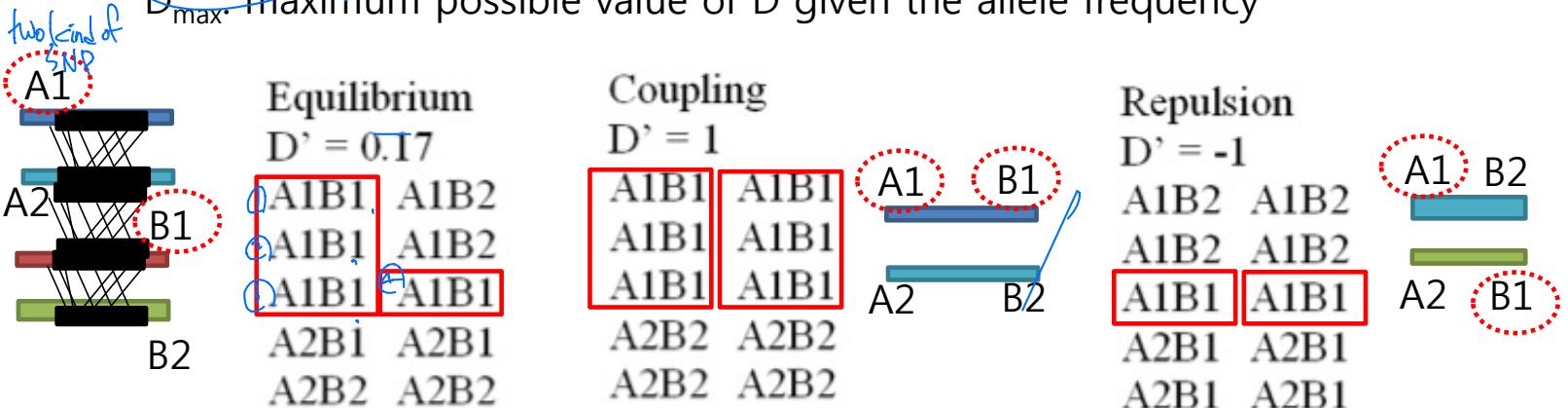
$$D' = \frac{D}{D_{\max}}$$

$0.125$   
 $0.2$

$$D_{\max} = \begin{cases} \min(p_1q_1, p_2q_2) & \text{when } D < 0 \\ \min(p_1q_2, p_2q_1) & \text{when } D > 0 \end{cases}$$

$\Sigma p_i q_i$   
 $D \propto \Sigma p_i q_i$

$D_{\max}$ : maximum possible value of D given the allele frequency



$$p_{11} = p_1q_1 + D$$

$$D = 0.24$$

$$D = -0.16$$

$$D = 6/10 - 6/10 \times 6/10$$

$$D = 2/10 - 6/10 \times 6/10$$

$$\begin{aligned} D &= p_{11} - p_1q_1 \\ \text{observed} &= \underline{\underline{4/10}} - \underline{6/10} \times \underline{6/10} \\ &= \underline{0.4} - \underline{0.36} = \underline{0.04} \end{aligned}$$

$$D_{\max} = \min(p_1q_2, p_2q_1)$$

$$= \min(6/10 \times 4/10, 4/10 \times 6/10) = \underline{0.24}$$

$$D' = 0.04 / 0.24 = \sim 0.17$$

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

# Exercise



## EXERCISE 3.1 *Quantifying heterozygosity and LD*

Using the following ten sequences:

|    |            |            |            |            |     |         |
|----|------------|------------|------------|------------|-----|---------|
| 1  | gctgcatcag | aagaggccat | caagcgcatc | actgtacttc | tgc | catggcc |
| 2  | gctgtatcag | aacaggccat | caagcgcatc | actgtacttc | tgc | catggcc |
| 3  | gctgtatcag | aacaggccat | caagcacatc | actgtacttc | tgc | catggac |
| 4  | gctgcatcag | aagaggccat | caagcacatc | actgtccttc | tgc | catggcc |
| 5  | gctgcatcag | aagaggccat | caagcacatc | actgtccttc | tgc | catggcc |
| 6  | gctgtatcag | aacaggccat | caagcgcatc | actgtccttc | tgc | catggcc |
| 7  | gcggcatcag | aagaggcgat | caagcacatc | actgtccttc | tgc | catggac |
| 8  | gctgcatcag | aagaggccat | caagcacatc | actctacttc | tgc | catggcc |
| 9  | gctgtatcag | aacaggccat | caagcgcatc | actgtccttc | tgc | catggcc |
| 10 | gctgcatcag | aagaggccat | caagcgcatc | actctccttc | tgc | catggcc |

(a) Count the number of segregating sites. SNPs

(b) Calculate the expected average heterozygosity per nucleotide.

(c) Determine the level of linkage disequilibrium between the common polymorphisms.

SNPs

## 9 SNPs

1 gatgcatacg aagaggccat caagggcatac actgtacttc tgccatggcc (n=50)  
 2 gctgtatcag aacaggccat caagggcatac actgtacttc tgccatggcc  
 3 gctgtatcag aacaggccat caagggcatac actgtacttc tgccatggac  
 4 gatgcatacg aagaggccat caagggcatac actgtacttc tgccatggcc  
 5 gatgcatacg aagaggccat caagggcatac actgtacttc tgccatggcc  
 6 gatgtatcag aacaggccat caagggcatac actgtacttc tgccatggcc  
 7 gatgcatacg aagaggccat caagggcatac actgtacttc tgccatggac  
 8 gatgcatacg aagaggccat caagggcatac actgtacttc tgccatggcc  
 9 gatgtatcag aacaggccat caagggcatac actgtacttc tgccatggcc  
 10 gatgcatacg aagaggccat caagggcatac actgtacttc tgccatggcc

0.1 # 3  
 0.4 # 3

**ANSWER:** (a) There are 9 segregating sites, at positions 3, 5, 13, 18, 26, 28, 34, 36, and 49.

0.5 # 1  
 0.2 # 2

(b) The minor allele frequencies at these 9 sites are, respectively, 0.1, 0.4, 0.4, 0.1, 0.5, 0.1, 0.2, 0.4, and 0.2. Since the expected heterozygosity of each nucleotide is given by  $2pq$  where  $p$  is one allele frequency and  $q$  ( $= 1 - p$ ) is the other allele frequency, the expected average heterozygosity  $H$  (including the 41 nonpolymorphic sites) is

$$H = [(41 \times 0) + (3 \times 0.18) + (3 \times 0.48) + (2 \times 0.32) + 0.50]/50 = 0.0624$$

This means that, on average, just over three site differences are expected between any pair of randomly chosen alleles.

0.48  
 0.32  
 0.5

$$0.0624 \times 50 = \sim 3.1$$

→ 0.18 (m21) 7/10  
 2/50

(c) There are four common polymorphisms, at positions 5, 13, 26, and 36.

cgga  
tcga  
tcha  
cgac  
cgac  
tcgc  
cgaa  
cgaa  
tcgc  
cgac

If we extract these sites, it is easier to see how they are related. Following the procedure in Box 3.1, draw a table of haplotype frequencies for each pairwise combination. For example, for sites 5 and 13:

Site 5

t ( $p_1 = 0.4$ )  
c ( $p_2 = 0.6$ )

Site 13

|                   |                   |
|-------------------|-------------------|
| c ( $q_1 = 0.4$ ) | g ( $q_2 = 0.6$ ) |
| $p_{11} = 0.4$    | $p_{12} = 0.0$    |
| $p_{21} = 0.0$    | $p_{22} = 0.6$    |

observed

Since  $D = p_{11} - p_1 q_1$ , for this pair  $D_{5,13} = 0.4 - (0.4 \times 0.4) = 0.24$ .

expected

not zero  $\rightarrow$  likely to be linked

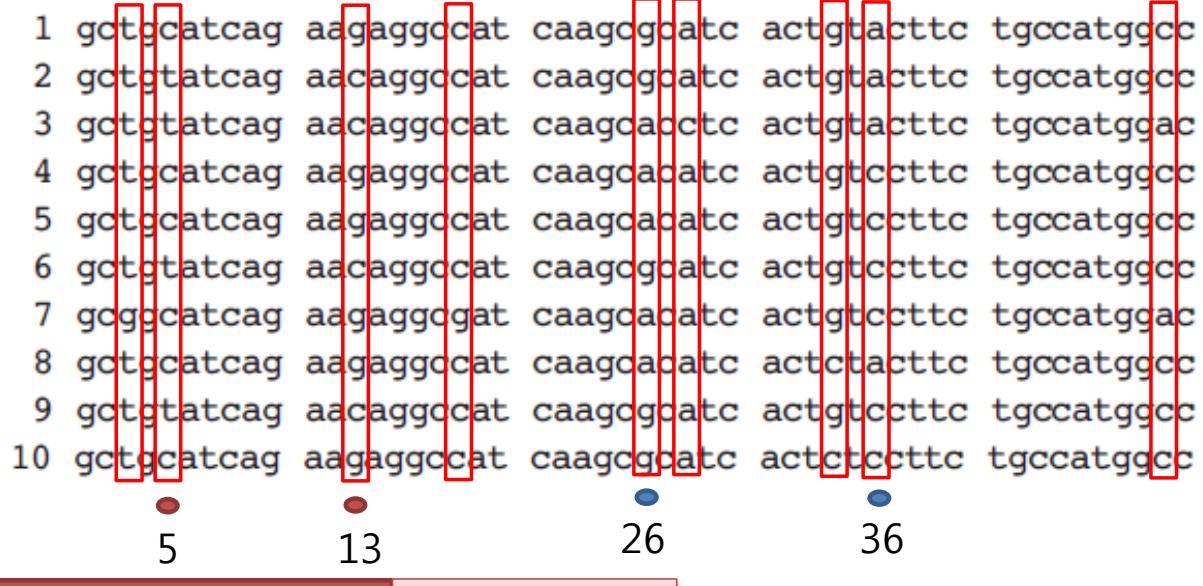
This is also the maximal value,  $D_{max}$ , since all of the less common alleles at both sites always segregate together. Consequently,  $D'$  is equal to 1.

You should be able to calculate the following table for the other linkage disequilibrium estimates.  $D_{max}$  is just the maximum value  $p_{11}$  could take, minus  $p_1 q_1$ .

$\min(p_1 q_2, p_2 q_1)$  when  $D > 0$

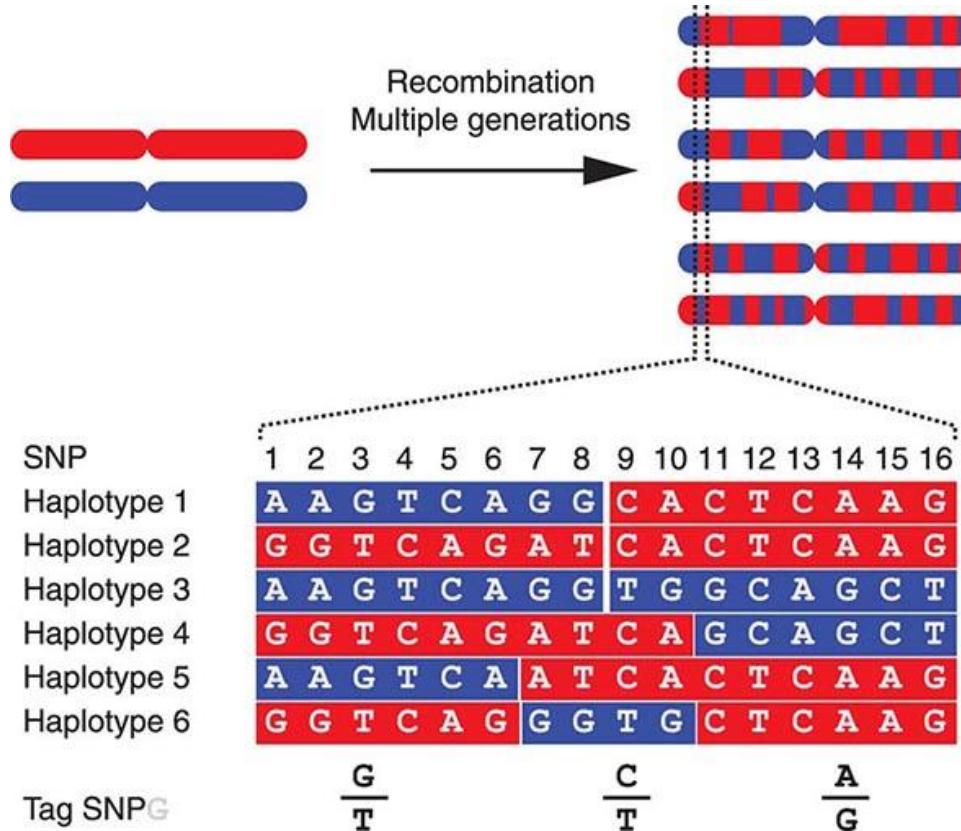
| Allele pair | $p_1$ | $q_1$ | $p_{11}$ | $D$  | $D_{max}$ | $D'$ |                                  |
|-------------|-------|-------|----------|------|-----------|------|----------------------------------|
| 5t, 13c     | 0.4   | 0.4   | 0.4      | 0.24 | 0.24      | 1.00 | $0.4 \times 0.6, 0.6 \times 0.4$ |
| 5t, 26g     | 0.4   | 0.5   | 0.3      | 0.10 | 0.20      | 0.50 |                                  |
| 5t, 36a     | 0.4   | 0.4   | 0.2      | 0.04 | 0.24      | 0.17 |                                  |
| 13c, 26g    | 0.4   | 0.5   | 0.3      | 0.10 | 0.20      | 0.50 |                                  |
| 13c, 36a    | 0.4   | 0.4   | 0.2      | 0.04 | 0.24      | 0.17 |                                  |
| 26c, 36a    | 0.5   | 0.4   | 0.2      | 0.00 | 0.20      | 0.00 |                                  |

This means that there is complete LD between the first two sites, but linkage equilibrium between the last two sites, with partial LD for all other comparisons.

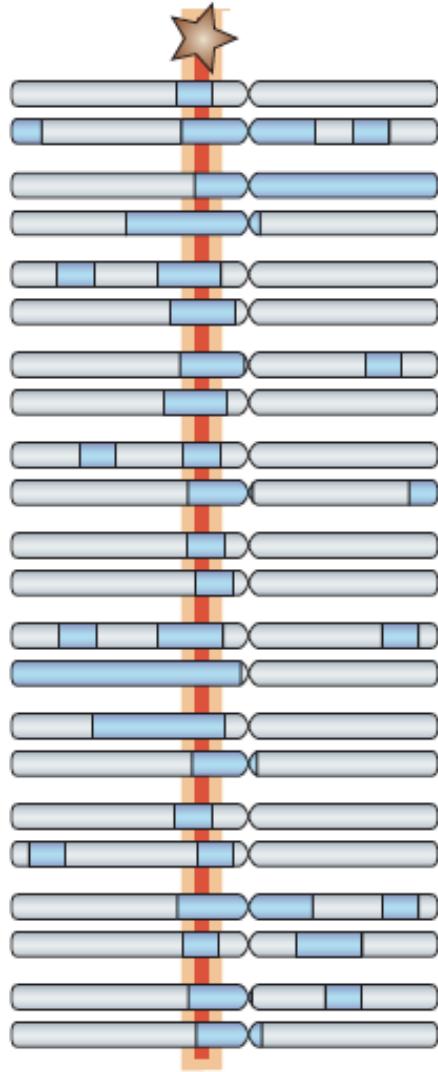


| Allele pair | $p_1$ | $q_1$ | $p_{11}$ | D    | $D_{max}$ | $D'$ |
|-------------|-------|-------|----------|------|-----------|------|
| 5t, 13c     | 0.4   | 0.4   | 0.4      | 0.24 | 0.24      | 1.00 |
| 5t, 26g     | 0.4   | 0.5   | 0.3      | 0.10 | 0.20      | 0.50 |
| 5t, 36a     | 0.4   | 0.4   | 0.2      | 0.04 | 0.24      | 0.17 |
| 13c, 26g    | 0.4   | 0.5   | 0.3      | 0.10 | 0.20      | 0.50 |
| 13c, 36a    | 0.4   | 0.4   | 0.2      | 0.04 | 0.24      | 0.17 |
| 26c, 36a    | 0.5   | 0.4   | 0.2      | 0.00 | 0.20      | 0.00 |

# Haplotype: identification of associated regions

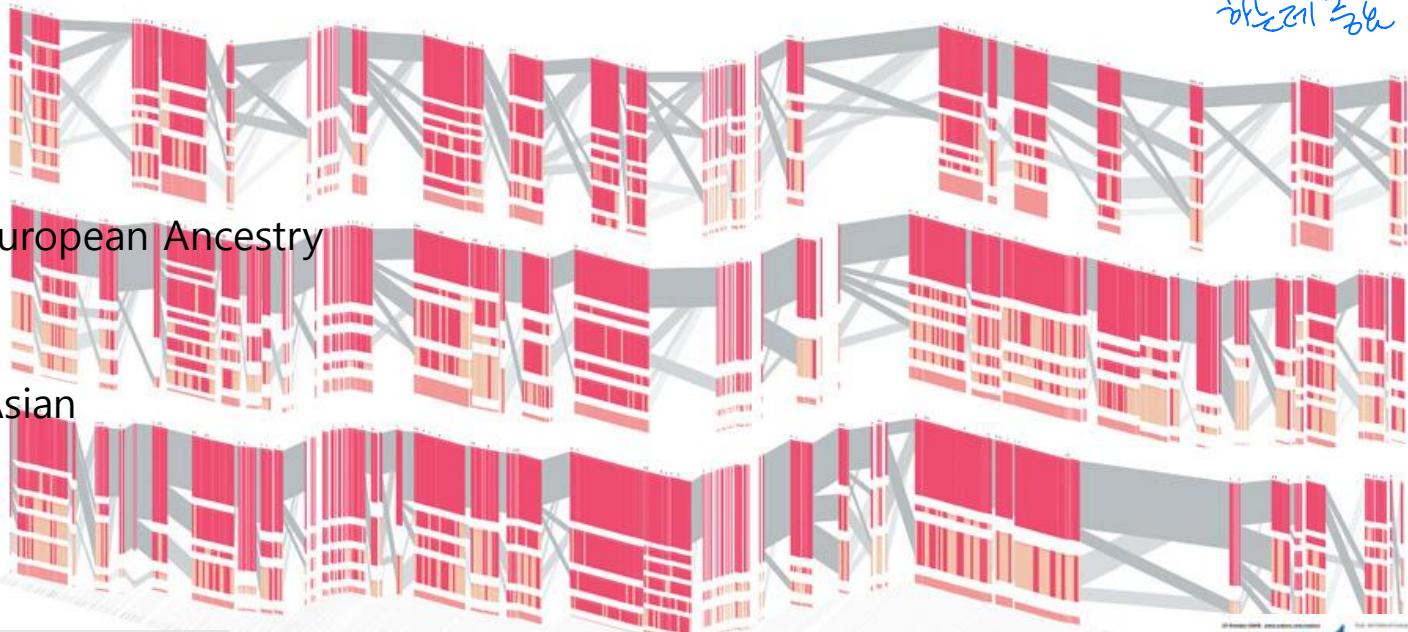


~ 정기적으로 찾기



# Haplotype Map (HapMap)

African



European Ancestry

Asian



(2008)

질병유전 / 개별개체의 차이를 파악  
하는 데 활용



(2005)

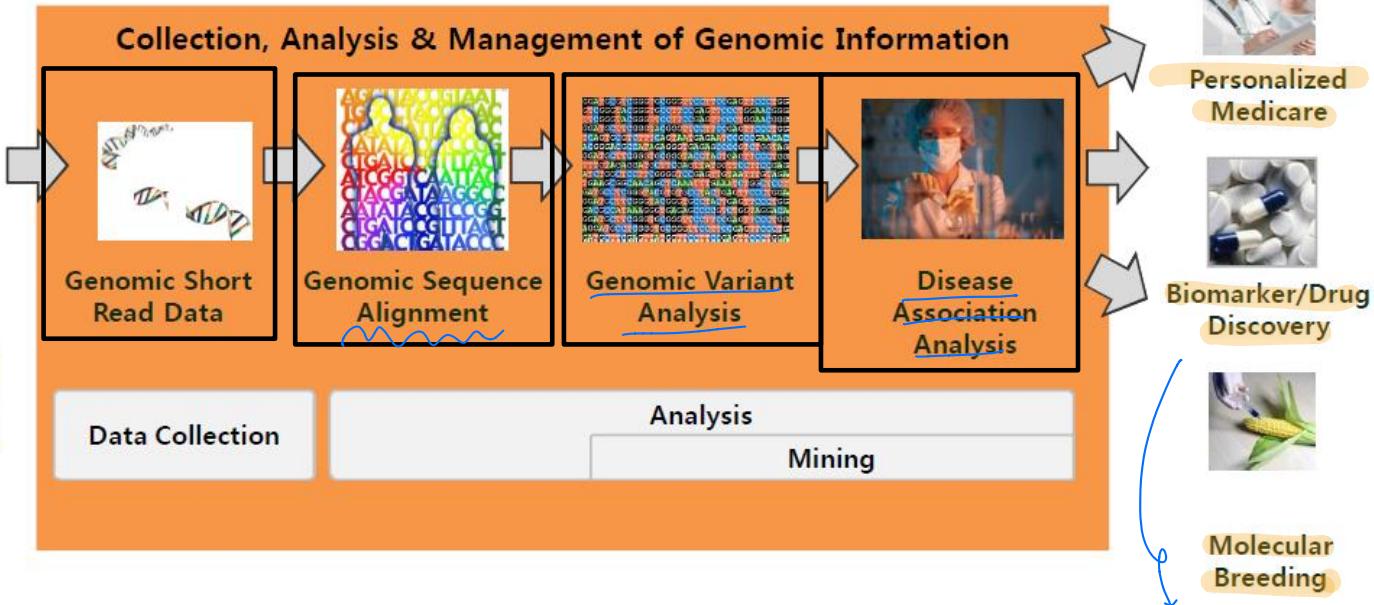
# Genome Variation Analysis

## Introduction to genomic variation

Sung Wook Chi

Division of Life Sciences, Korea University

# Genomic Variant Analysis



Variation

Genome

Gene expression

Change

Phenotype

SNP关联分析  
基因表达  
药物筛选  
作物育种

# Genetic Variations



## Sizable

Chromosome numbers *Some cancer case*

Segmental duplications,

**Copy Number Variation (CNV)**

*Some region duplicated / deleted*

Translocations

Inversion

Sequence Repeats

Transposable Elements

*repeat*

Short deletions and insertions

Tandem Repeats

Nucleotide Insertions and Deletions (Indels)

*one or two insertion / deletion*

**Single Nucleotide Polymorphisms (SNPs)**

*more frequent in population*

Mutations *frequency low*

## Structural

## Sequence

## Minor

# Sequence Variations (SNP, mutation)

## POLYMORPHISM

Widespread

Accepted mutation

Quantitative trait

(자연선택)

93% C T A A G T A

7% C T A C G T A

## MUTATION

99.7% C T A A G T A

0.3% C T A C G T A

Rare

Diseases  
(more lethal)

**Polymorphism:** Single DNA base change found in >1% of population

**Mutation:** Single DNA base change found in <1% of population

*Genetic mutations are one type of genetic polymorphism*

### Sequences

. . . T C A A G T C A A G C G A T C A T G . . .  
. . . T C A A G T C A A G C G A T C A [G] G . . .  
. . . T C A [G] G T C A A G [T] G A T C A T G . . .  
. . . T C A [G] G T C A A G [T] G A T C A T G . . .  
. . . T C A A G T C A A G C G A T C A [G] G . . .  
. . . T C A A G T C A A G C G A [A] C A [G] G . . .

### Haplotypes

### Tagging SNPs

DNA sequence

| Block 1 | Block 2   | Block 3 |
|---------|-----------|---------|
| A G G   | T C A C T | T A G   |
| C A T   | A C A C T | G C C   |
| A G G   | A C G T T | T A G   |
| A G G   | A C G T T | T A G   |
| A G G   | T C A C T | G A G   |
| A G G   | T T A C A | G C C   |

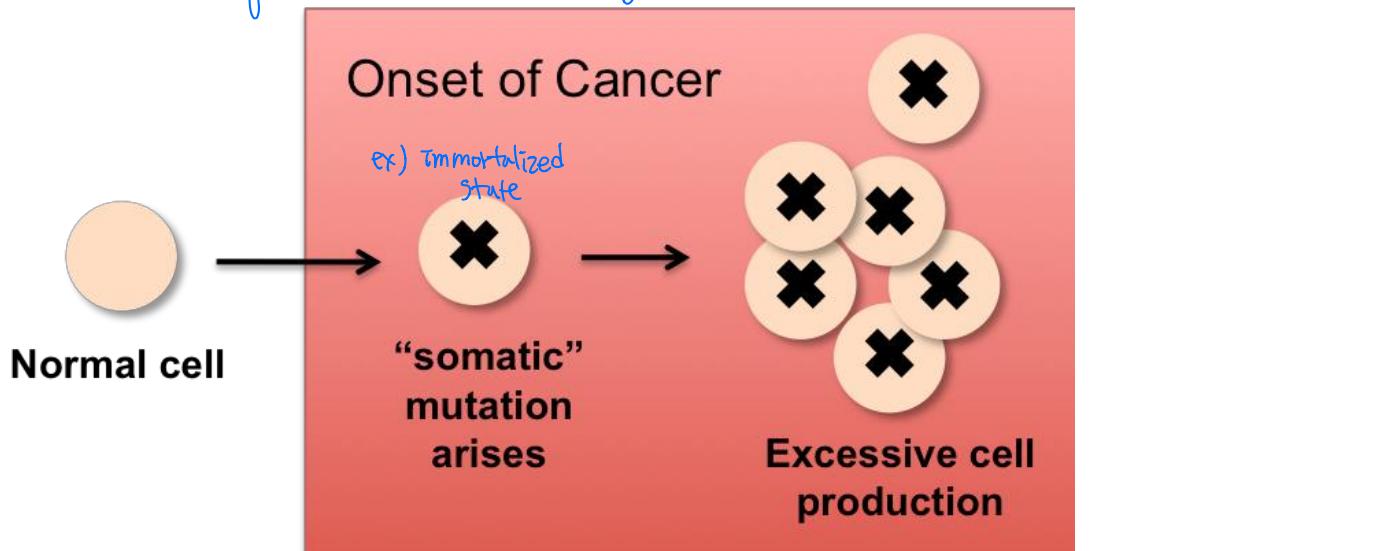
↓      ↓      ↓      ↓

G/A    C/T    C/T    A/C

↳ large haplotype region

# Genetic Variations

- Germline variation (Blood, normal tissue)  
: GWAS, Targeted genetic study  
*↳ all inherited & every tissue of my body have same variation*
- Somatic variation (Tumor tissue) *→ variation of tumor tissue only*  
: Cancer genomics  
*↳ acquired mutation during cell division*

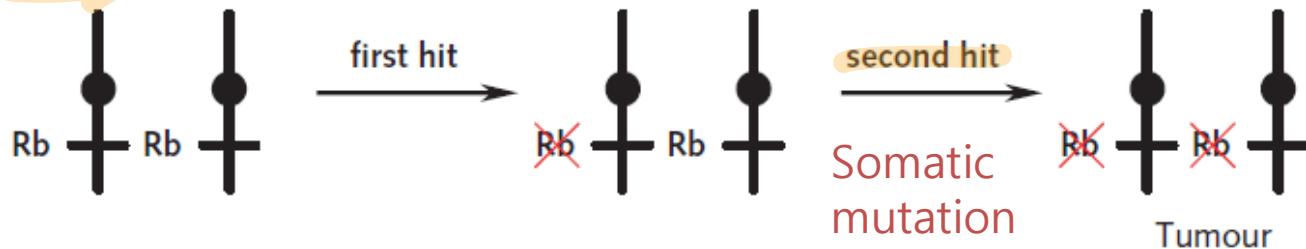


# Mutation and Cancer

## Somatic mutation

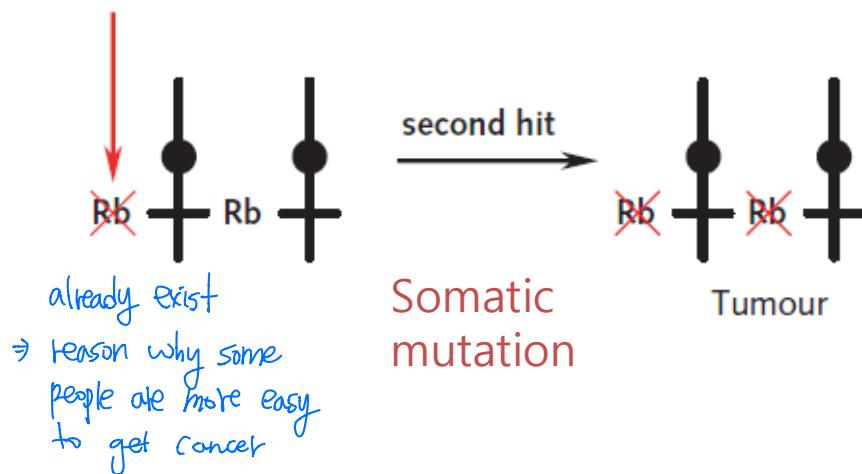
산발적인, 이전에 발생하는

Sporadic retinoblastoma: two hits required



Two hits theory

## Germ line mutation



already exist  
⇒ reason why some  
people are more easy  
to get cancer

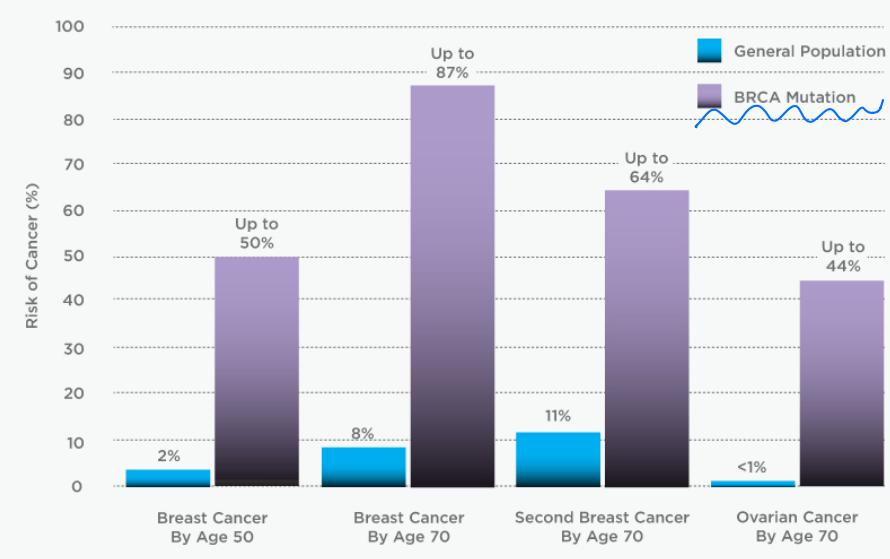
# Mutation and Cancer

- (1) Mutations detectable in the genome indicate propensity for development of cancers. Mutations in BRCA1 and BRCA2, as indicators for likelihood of breast and ovarian cancer development, are probably the best known.
- (2) Sequence analysis can predict disease progression and outcome.
- (3) Sequence analysis can help choose optimal treatment.
- ↗ BRCA էսու բարեկայի օլդին էտին էլք

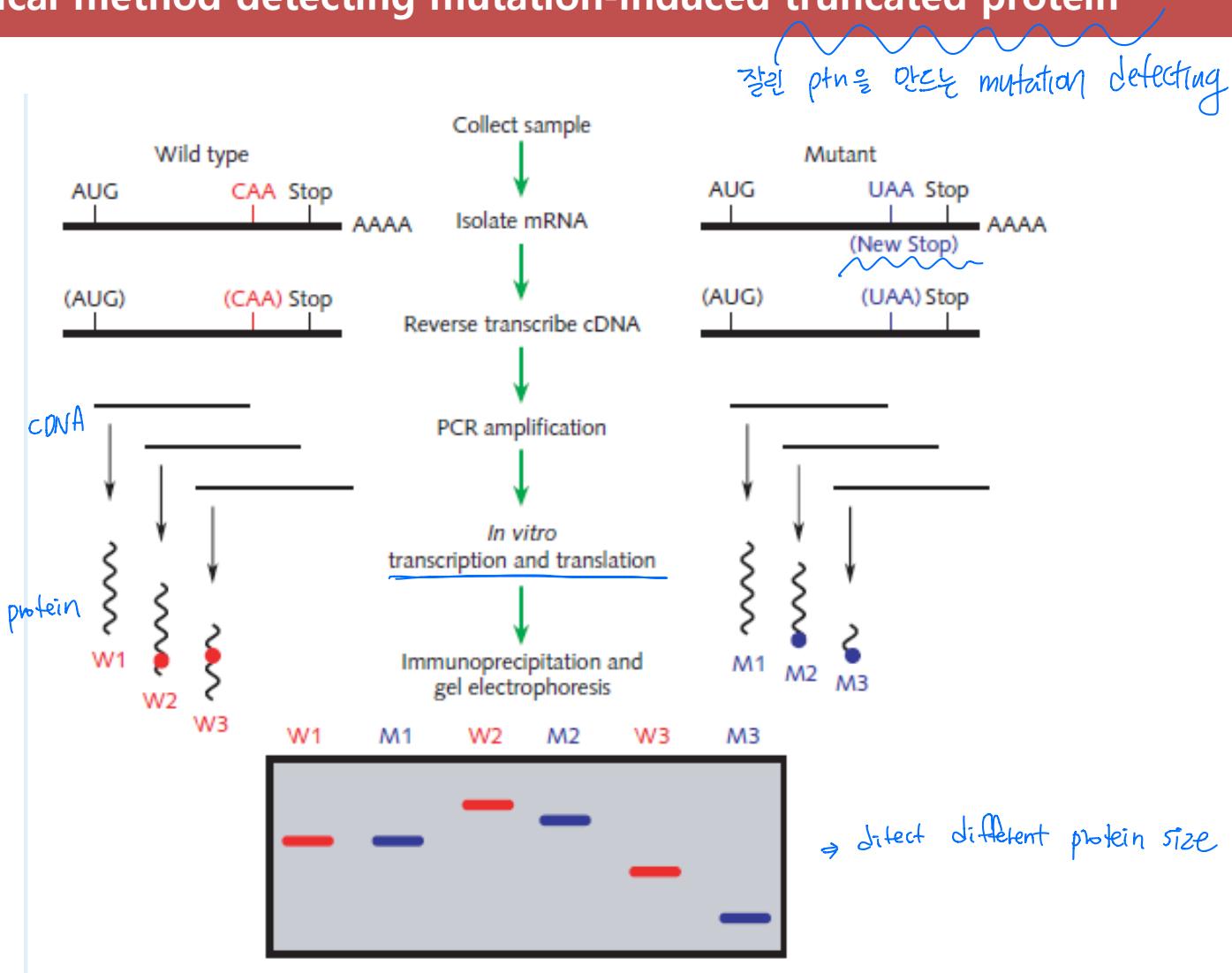


Table 2.1 Common *BRCA1* and *BRCA2* mutations

| Population     | Common<br><i>BRCA1</i><br>mutations                 | Common<br><i>BRCA2</i><br>mutations |
|----------------|---|-------------------------------------|
| Ashkenazi Jews | 185delAG, 5382insC                                  | 6174delT                            |
| Iceland        |   | 999del5                             |
| Denmark        | 2594delC, 5208T→C                                   |                                     |
| Lithuania      | 4153delA, 5382insC,<br>61G→C                        |                                     |
| China          | 589delCT, IVS7-27del10,<br>1081delG, 2371-2372delTG | 3337C→T                             |

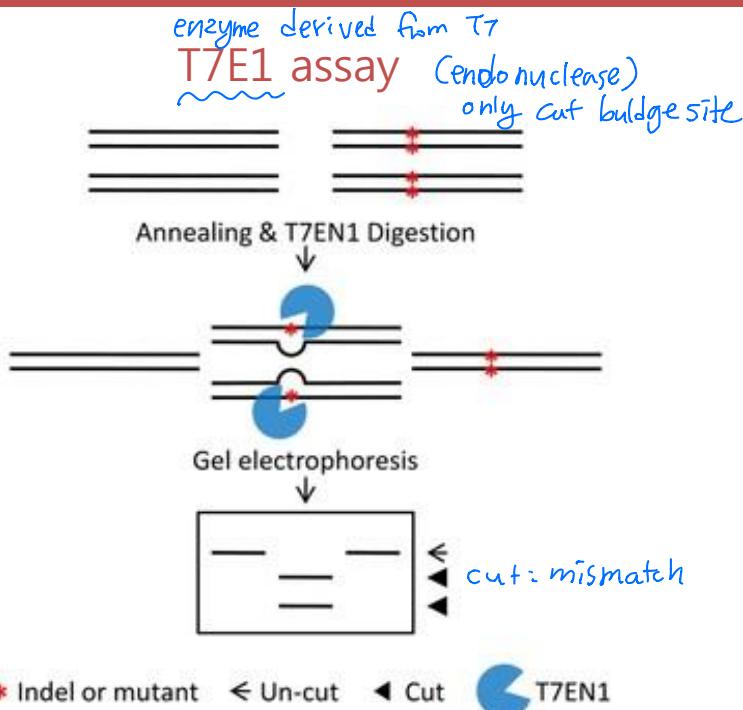
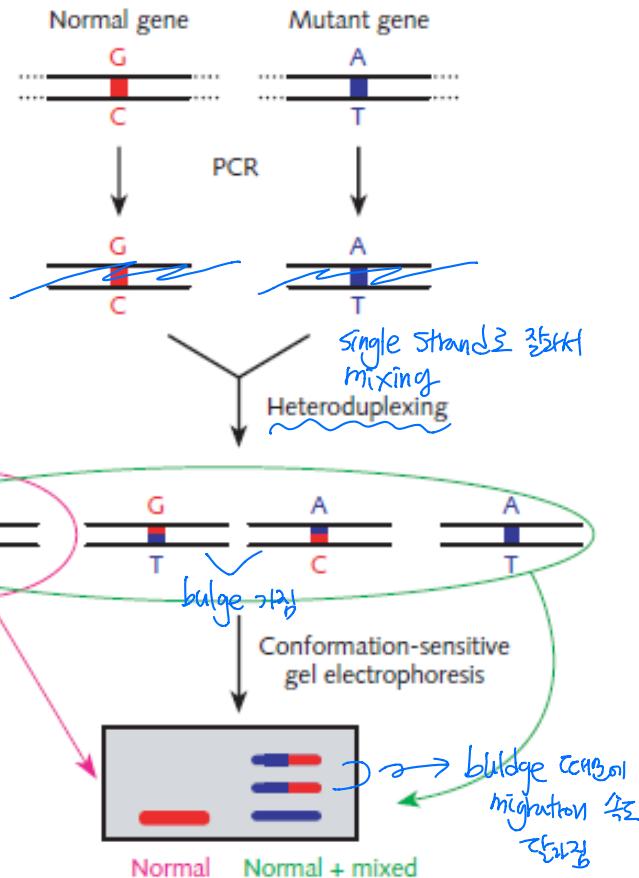


# Classical method detecting mutation-induced truncated protein



# How to experimentally detect mutation?

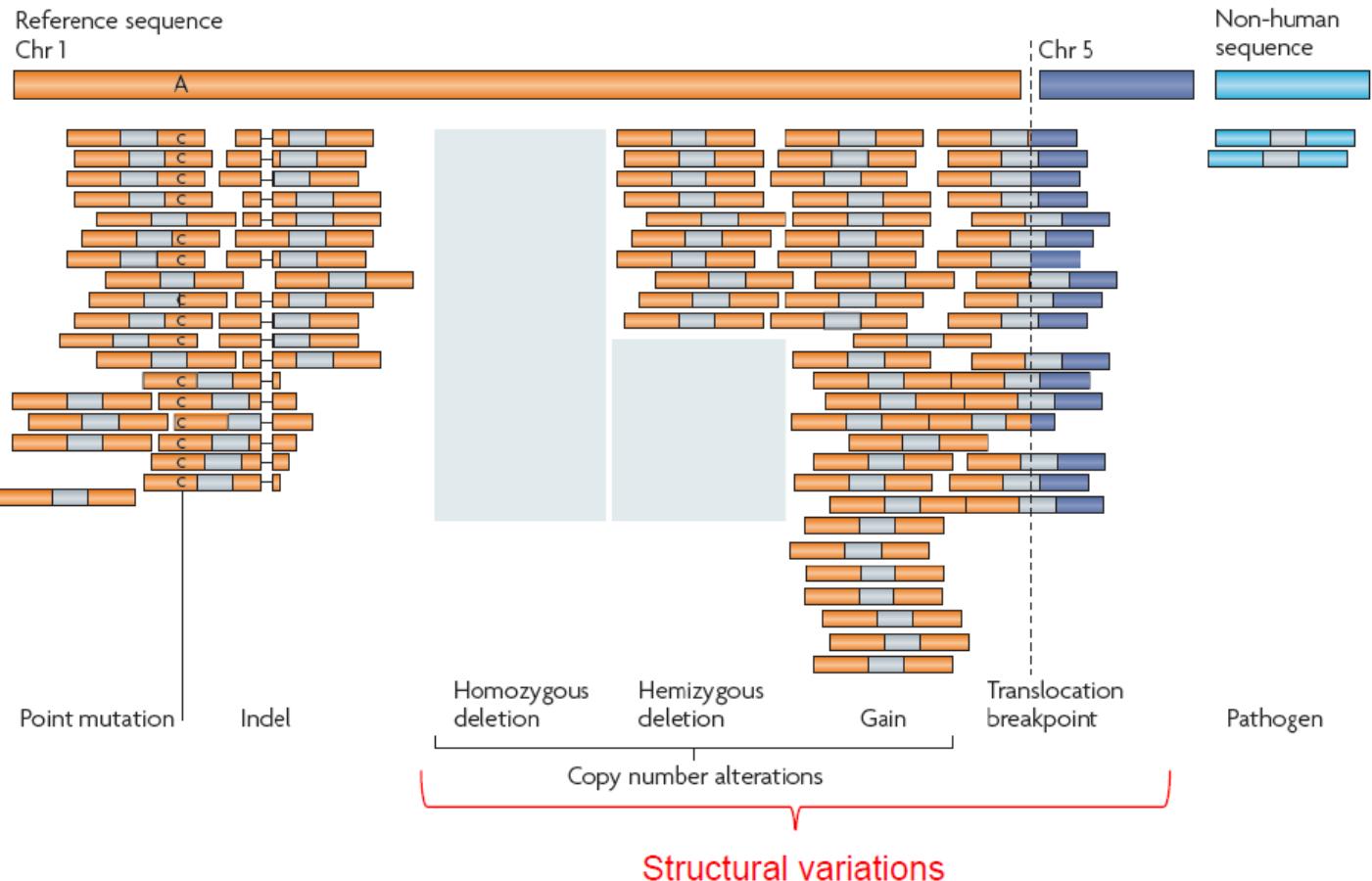
## Conformation-sensitive gel electrophoresis



**T7E1 enzyme:**  
cut mismatched DNA duplex site

Now days we are using sequencing (NGS)

# Structural Variations (SVs)



# Oncogenes vs. Tumor Suppressor genes



- Oncogenes
  - Growth signals
  - Cell multiplication
  - Activated in cancer
- Tumor Suppressor genes
  - Growth suppressive signals
  - Cell stop dividing
  - Inactivated in cancer

# Corrupted Genomes in Cancer Cells



Amplification



Point mutation



Translocation

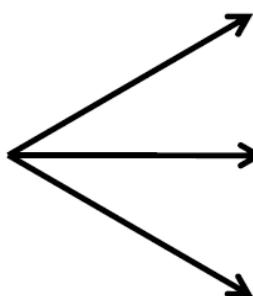
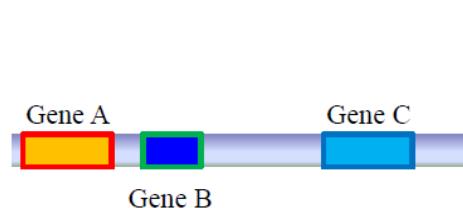


Increased signal

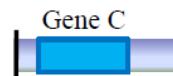
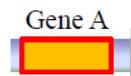


Abnormal signal

## Different Types of CNVs

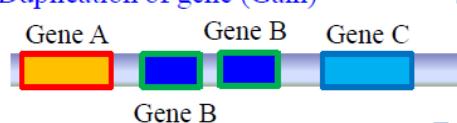


Deletion of gene (Loss)



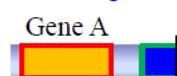
Deletion (loss)

Duplication of gene (Gain)



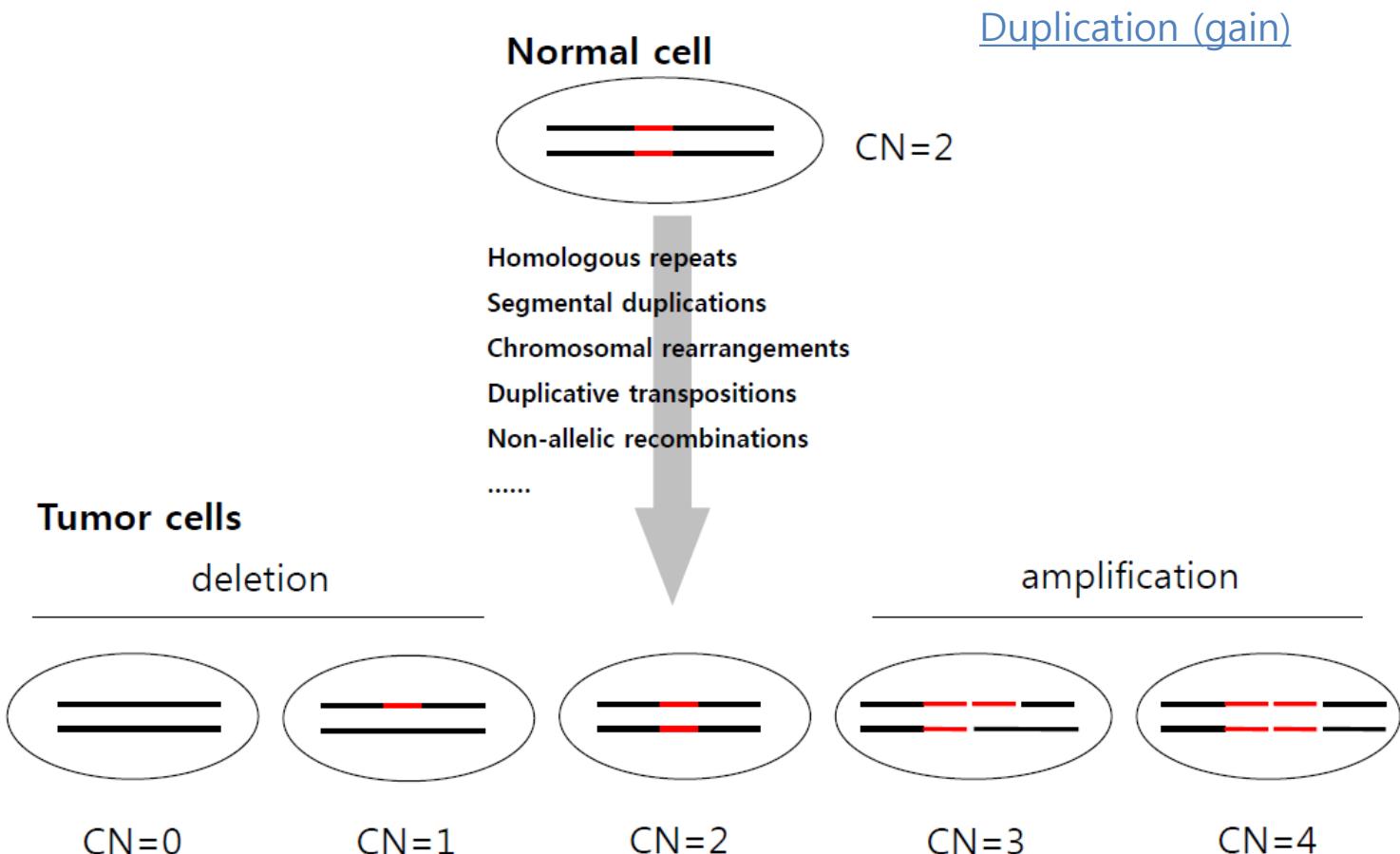
Duplication (gain)

Deletion generated new fused gene



Fusion (deletion, new function)

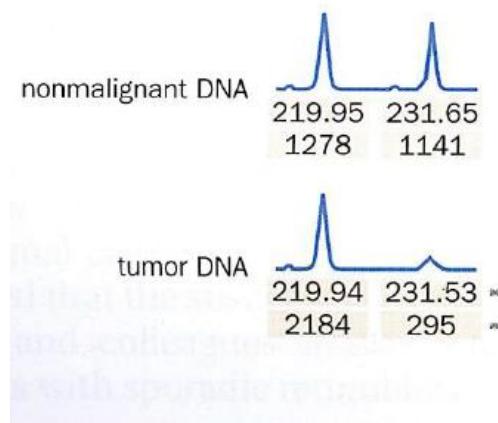
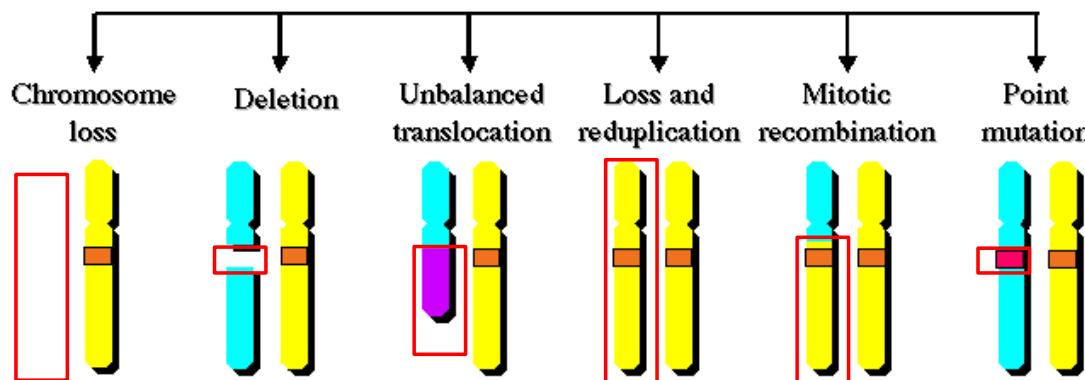
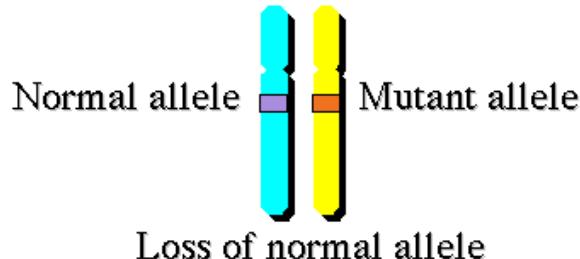
# Genetic Alterations in Tumor (Copy number changes)



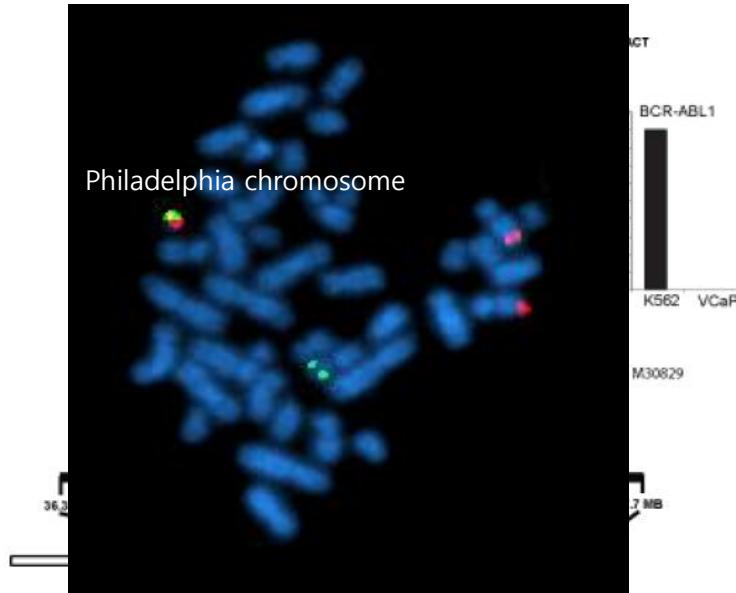
# Loss of heterozygosity (LOH)

## Loss of heterozygosity as a marker to locate tumor suppressor genes

- Somatic genetic changes in retinoblastoma caused loss of heterozygosity (LOH) at markers close to the RB1 locus
- By screening paired blood and tumor samples with markers spaced across the genome, we may discover the locations of tumor suppressor genes

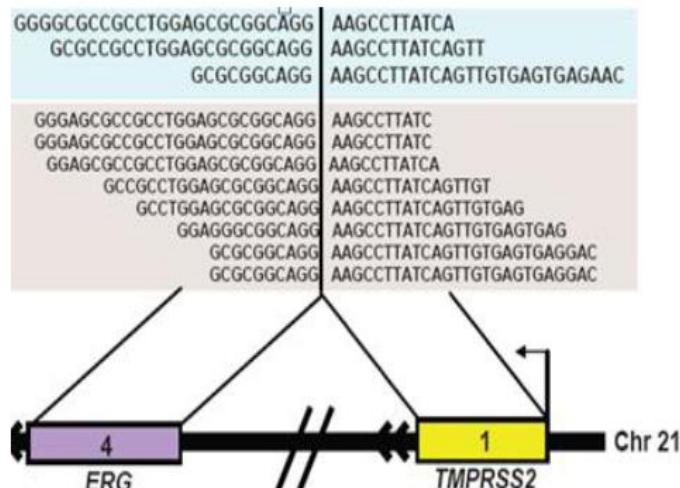


# Fusion Gene in Cancer

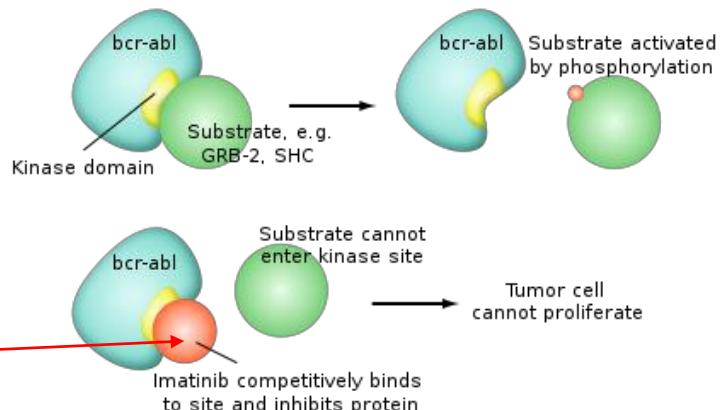


BCR-ABL1 caused by translocation  
Chronic myeloid leukemia

## Translocation



TMPRSS2-ERG caused by deletion  
Prostate cancer



# Genomics and Medicine : 100,000 Genome Project in England



Cancer and rare diseases

About Us ▾ 100,000 Genomes Project ▾ Taking Part ▾ For Healthcare Professionals ▾ Research ▾ Industry Partnerships ▾ News & Events ▾



Genomics England is delivering  
the **100,000 Genomes Project.**

We are creating a new genomic medicine service with the NHS – to support  
**better diagnosis and better treatments** for patients. We are also enabling medical research.

[More information about the 100,000 Genomes Project](#)

Start in late 2012, aim to finish by 2017

<https://www.youtube.com/watch?v=hxou7ayQSZQ>