

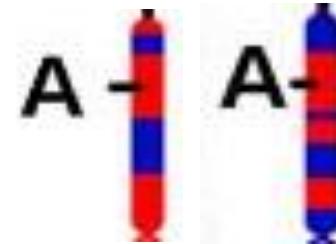
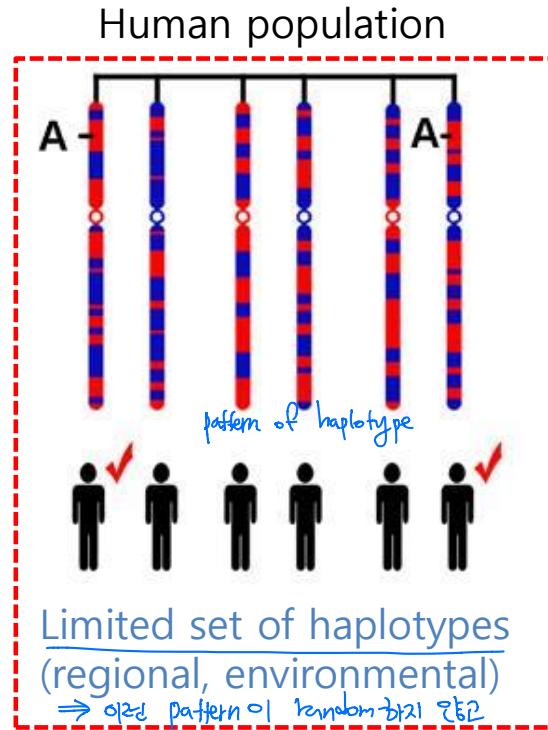
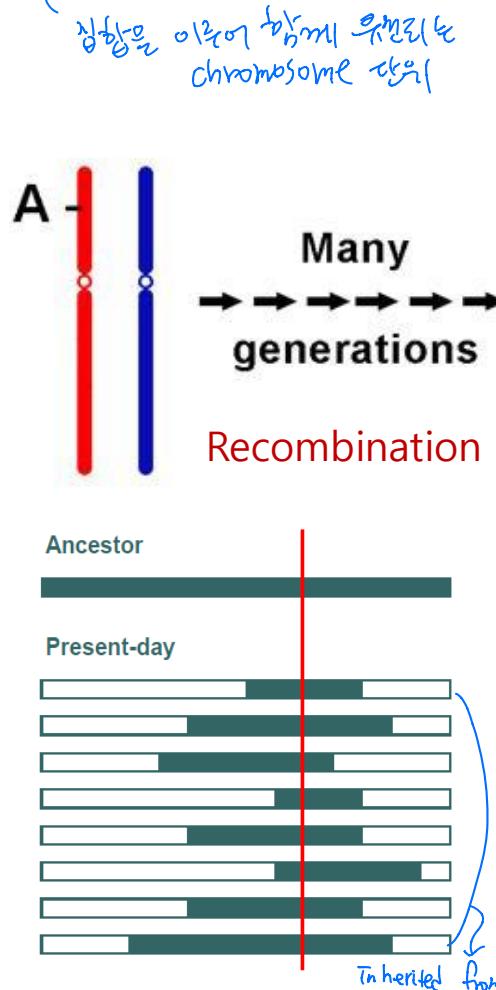
Summary & Review

Basic knowledges for human genomic variation

- How to determine the haplotypes (tagged SNP) ?
- What is the method for haplotype maps (HapMap) ?
: **Linkage disequilibrium (LD)**
- Why haplotype is important for genomic variation analysis?

haplotype

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



Only two haplotypes

> Easy to analyze

유전자는 지역을 제한하거나
haplotype 유사성

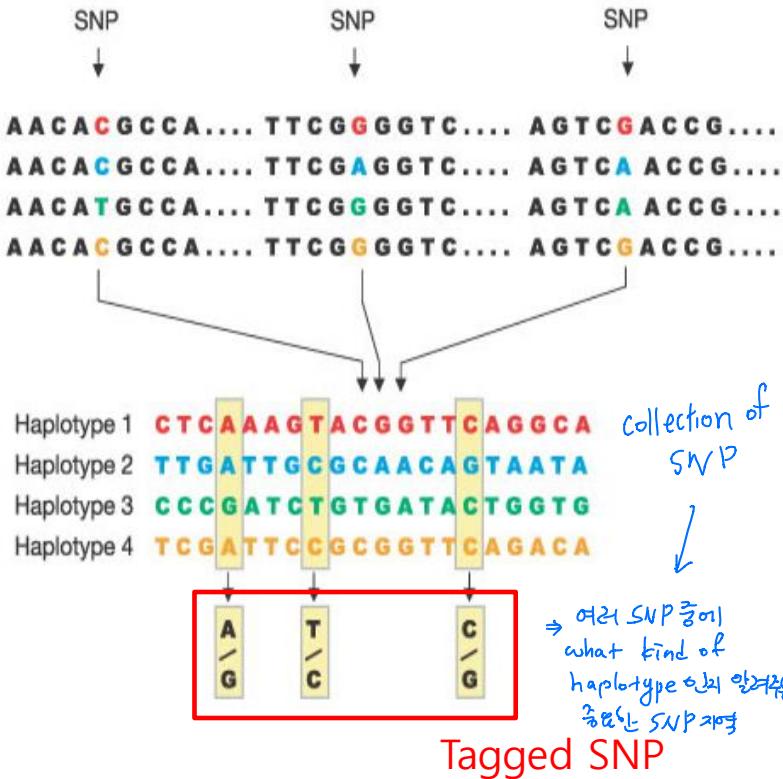
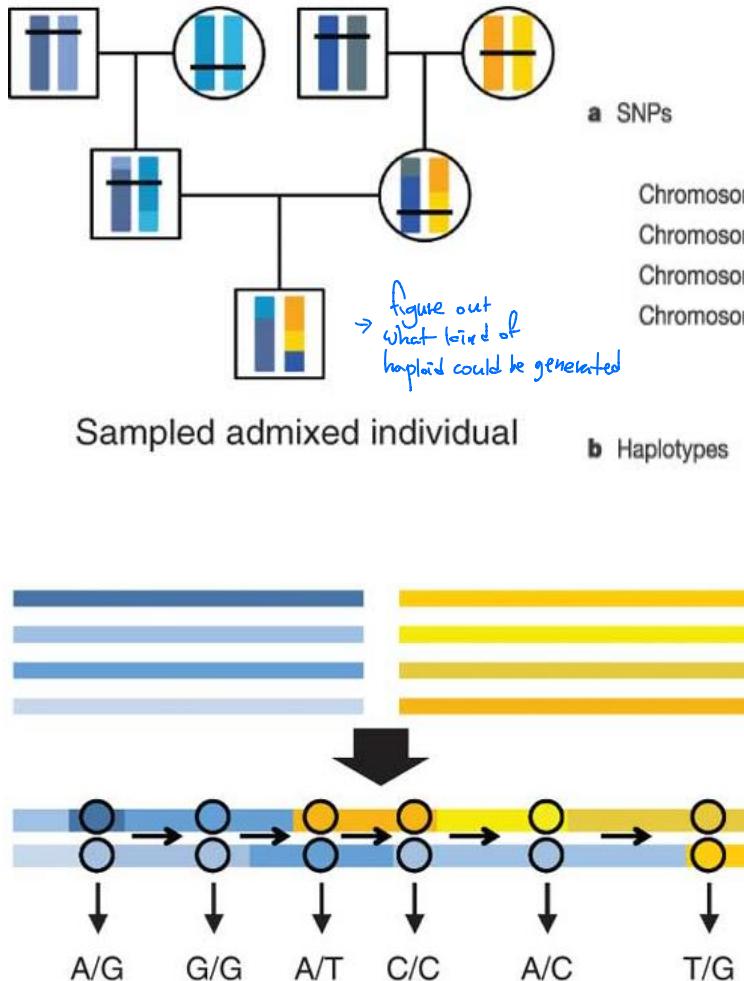
How to determine the same haplotype regions?

Variation (SNP: single nucleotide polymorphism)

SNP를 재现出한 나머지 개인은 같음

Tagged SNP

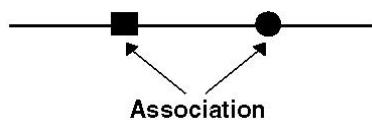
Ancestors of variable ancestry



Only 4 haplotypes
: discriminate by 3 tagged SNPs

Linkage disequilibrium (LD)

LINKAGE DISEQUILIBRIUM:



Sequencing results (population)

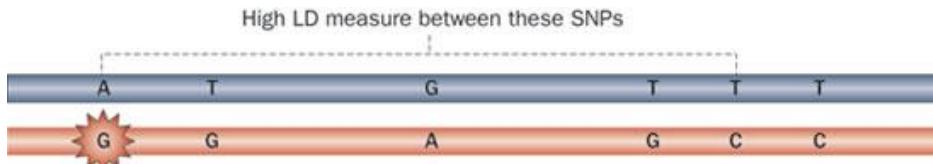
Observed frequency
of co-occurrence

Two SNPs

\neq
linked

Expected frequency
of co-occurrence

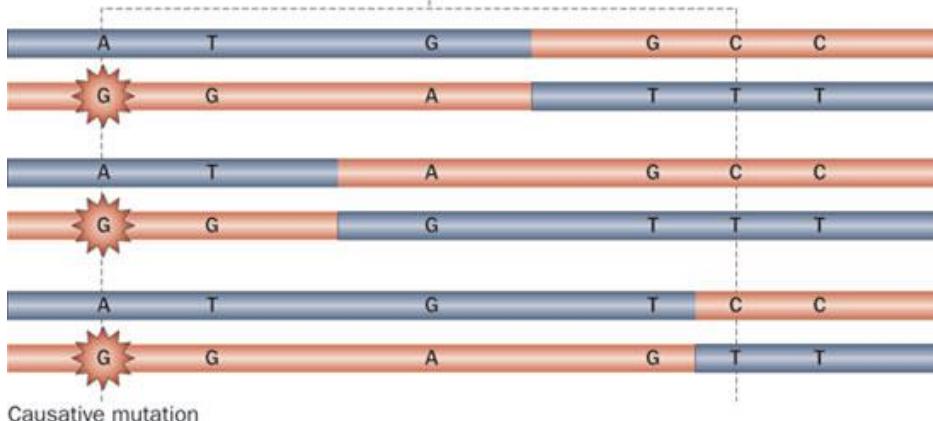
How much it becomes
disequilibrium ?



Causative mutation

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

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



Completely linked

설정되는 텐트
설정되는 텐트
설정되는 텐트
설정되는 텐트
설정되는 텐트

Linked (located in same
cosaggregated region)

A T

Completely independent

observed

=

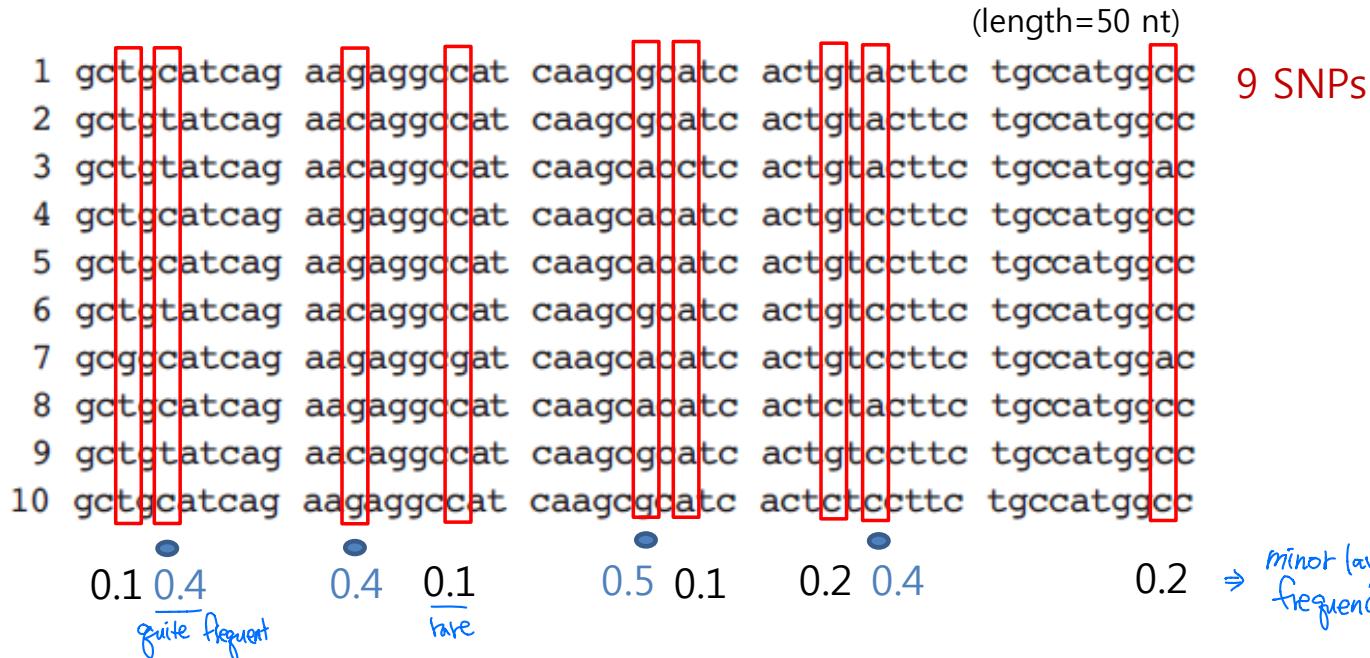
expected

Do given sequence results contain enough variations to determine haplotype ?

diversity^를 설정하는
value가 높을 \rightarrow SNP

10 sequence results
(human)
each indicate one haplotype

10 sequence of haploid



Diversity $[(41 \times 0) + (3 \times 0.1) + (3 \times 0.4) + (2 \times 0.2) + 0.5] / 50 = 0.048$

is this diverse?!

혹은

혹은

→ 2개의 염색체(diploid)가 서로 다른 정보를
가지고 있을 때 예상되는 대로 계산

actual length of sequence
 $0.048 \times 50 = 2.4$

random하게 짧았을 때
다를 경우

Expected average heterozygosity per nucleotide (H)

Select two haploid sequences from those above to make diploid,
then see how much their sequence is expected to be different to form heterozygosity

→ 0.9 X 0.1 이 될 수도 있으나 (각각의 개수 x diploid)

$$[(41 \times 0) + (3 \times 0.1 \times 0.9 \times 2) + (3 \times 0.4 \times 0.6 \times 2) + (2 \times 0.2 \times 0.8 \times 2) + 0.5 \times 0.5 \times 2] / 50 = 0.0624$$

$0.0624 \times 50 = 3.12$ → in population에서 2개의 haplotype 를 갖을 때
최초 2개 location은 heterozygote 일 것이라 예상

Calculation of linkage disequilibrium (LD)

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

observed frequency

	c	t	total
g	P_{cg}	P_{tg}	P_g
c	P_{cc}	P_{tc}	P_c
	P_c	P_t	1

≠

자유변수, 조건부로
→ 다른 변수

expected

확률적으로 계산한 것

	c	t	total
g	$P_c \times P_g$	$P_t \times P_g$	P_g
c	$P_c \times P_c$	$P_t \times P_c$	P_c
	P_c	P_t	1

parameter of diverse

D : linkage disequilibrium coefficient

	c	t	total
g	$P_{cg} = P_c \times P_g + D$		P_g
c			P_c
	P_c	P_t	1

if two loci are in linkage equilibrium
→ $D=0$

if two loci are in linkage disequilibrium
→ $D \neq 0$

probability of recombination이 작아질 수 있음 X

	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

질량 테이블과 같도록

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

$$D_{\max} = \begin{cases} \min(p_1 q_1, p_2 q_2) & \text{when } D < 0 \\ \min(p_1 q_2, p_2 q_1) & \text{when } D > 0 \end{cases}$$

Completely linked

$$(D=D_{\max})$$

$$(D'=1)$$

Linked

$$\min \{x_{21}, x_{12}\}$$

Completely independent

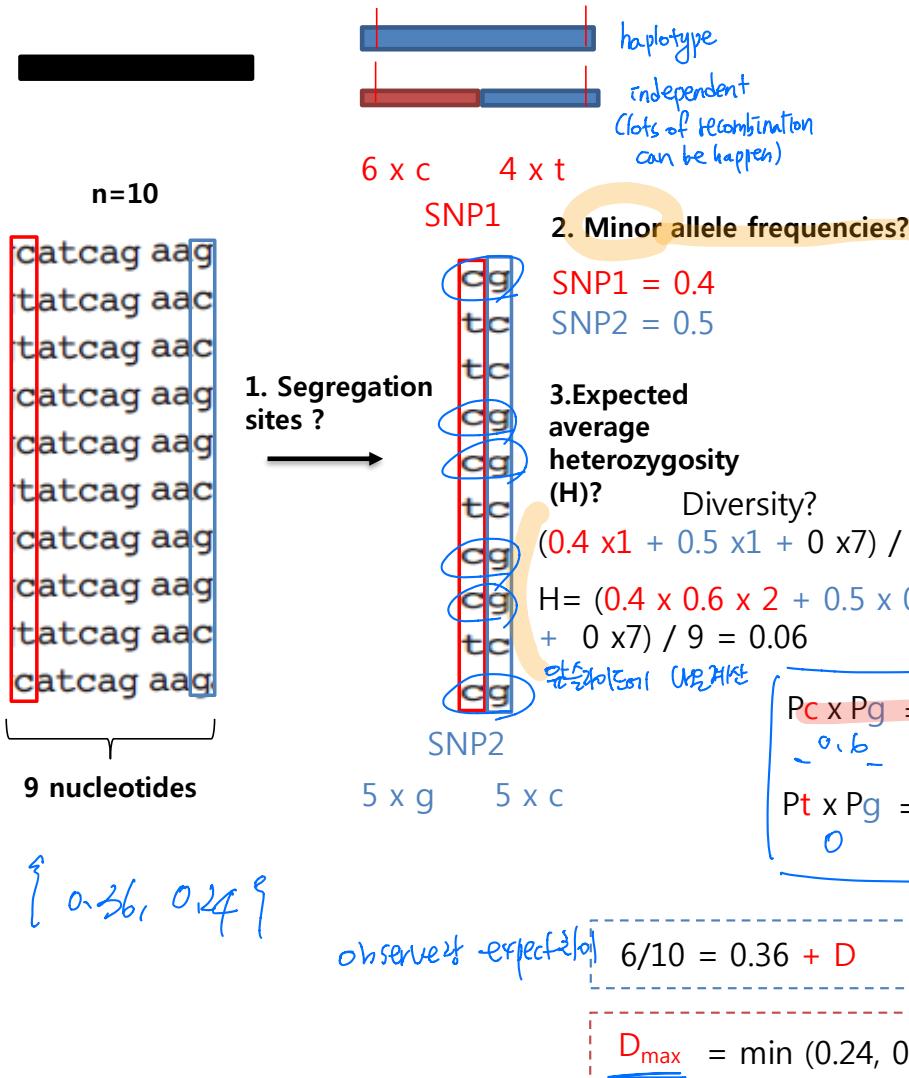
$$(D=0)$$

$$(D'=0)$$

observed

= expected

Linkage Disequilibrium (LD) analysis > haplotyping



4. Level of LD in SNP?

SNP1 (c) Vs. SNP2 (g)

- Observed freq

$$\begin{array}{l|l} P_{cg} = 6/10 & P_{cc} = 0 \\ P_{tg} = 0 & P_{tc} = 4/10 \end{array}$$

- Expected freq

$$\begin{array}{l|l} P_c = 6/10 & P_t = 4/10 \\ P_g = 6/10 & P_c = 4/10 \end{array}$$

$$\begin{array}{l|l} P_c \times P_g = 0.6 \times 0.6 = 0.36 & P_c \times P_c = 0.6 \times 0.4 = 0.24 \\ -D & -D \\ P_t \times P_g = 0.4 \times 0.6 = 0.24 & P_t \times P_c = 0.4 \times 0.4 = 0.16 \\ O & 4.0 \\ -D & +D \\ 0.24 & 0.24 \end{array}$$

$$D' = 0.24/0.24 = 1$$

Complete LD !!!

(region 1 ~ 4)

1 gatgcatacg aagaggccat caagggcatac actgtacttc tgccatggcc
 2 gctgtatcag aacaggccat caagggcatac actgtacttc tgccatggcc
 3 gctgtatcag aacaggccat caagggcatac actgtacttc tgccatggac
 4 gatgcatacg aagaggccat caagggcatac actgtacttc tgccatggcc
 5 gctgcatacg aagaggccat caagggcatac actgtacttc tgccatggcc
 6 gatgtatcag aacaggccat caagggcatac actgtacttc tgccatggcc
 7 gatgcatacg aagaggccat caagggcatac actgtacttc tgccatggac
 8 gatgcatacg aagaggccat caagggcatac actgtacttc tgccatggcc
 9 gctgtatcag aacaggccat caagggcatac actgtacttc tgccatggcc
 10 gatgcatacg aagaggccat caagggcatac actgtacttc tgccatggcc

5 13

26 independent 36

frequency가 높은 information
 높은 없음. 그 대신에 32 번째가
 높아서 → frequency는 SNP가 아니고
 to analyze haplotype

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 → best
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 → independent

(3)	26	observed		expected	
g	g	g	a	g	a
c	g	g	0.2	0.4	g 0.30
c	a	c	0.3	0.1	c 0.20
g	a				$D > 0 \Rightarrow \min \{0.3, 0.2\}$
g	a		$g = 0.6$	$g = 0.5$	
c	g		$c = 0.4$	$a = 0.5$	
g	a				
g	a	Finding D			
c	g		g		a
g	g	$g = 0.2 = 0.3 - D$		$0.4 = 0.3 + D$	
c		$0.3 = 0.2 + D$		$0.1 = 0.2 - D$	
		$D_{max} = \min \{0.1, 0.1\} = 0.1$			
		$D = 0.1$		$\Rightarrow \dots ?$	

(5)	26	g	a		
c	g				
t	g	t 0.3	0.1	0.2	0.2
t	a	c 0.2	0.4	0.3	0.3
c	a				
c	a				
t	g				
c	a				
t	g				
c	g				
t	g				
c	g				
		$D_{max} = 0.1$?			
		$t = 0.4$	$g = 0.5$		
		$c = 0.6$	$a = 0.5$		

(st) (ba)

observed

C	a		a	c
t	a	c	0.2	0.4
t	a	t	0.2	0.2
c	c			
c	c			
t	c			
c	c			
c	a			
t	c			
c	c			

$$c = 0.6 \quad a = 0.4$$

$$t = 0.4 \quad c = 0.6$$

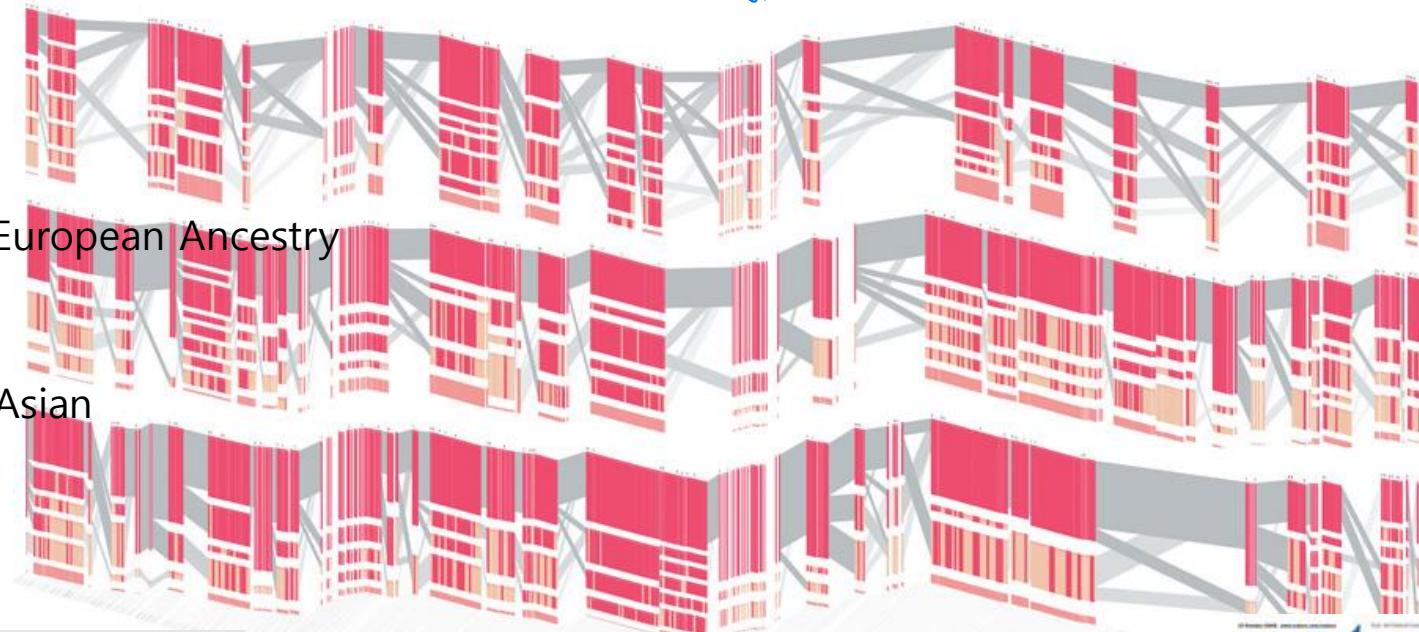
Haplotype Map (HapMap)

African

ancestor가 2번이 되었을 haplotype도 2개

European Ancestry

Asian



(2008)



(2005)

Genome Variation Analysis (review and more)

Introduction to genomic variation

focus on cancer thing

Sung Wook Chi

Division of Life Sciences, Korea University

Genetic Variations



Sizable

Chromosome numbers
Segmental duplications,
Copy Number Variation (CNV)

Translocations
Inversion
Sequence Repeats
Transposable Elements
Short deletions and insertions
Tandem Repeats

Nucleotide Insertions and Deletions (Indels)

Single Nucleotide Polymorphisms (SNPs)

1%
Mutations → occur less than 1%

Minor

Structural

Sequence

→ widespread among population
& inherited from parent
(mutation is acquired)

현장에 의해 얻어진거나 아니면
전해온 경우

Oncogenes vs. Tumor Suppressor genes



- Oncogenes
 - Growth signals
 - Cell multiplication
 - Activated in cancer

tumor : want to amplify

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

want to loss or stop

Corrupted Genomes in Cancer Cells



Amplification



Point mutation
↗ loss function



Translocation



Increased signal

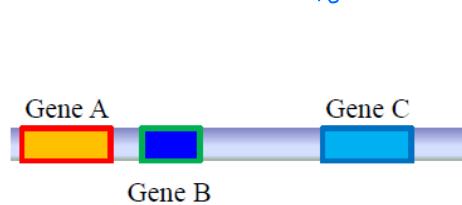


Abnormal signal

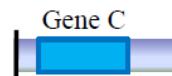
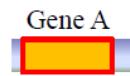
different chromosome
or 염색체로

Different Types of CNVs

Copy number variation



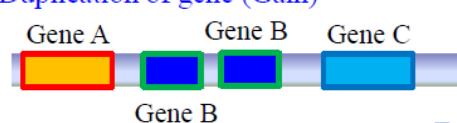
Deletion of gene (Loss)



Deletion (loss)

ex. tumor suppressor

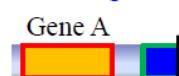
Duplication of gene (Gain)



Duplication (gain)

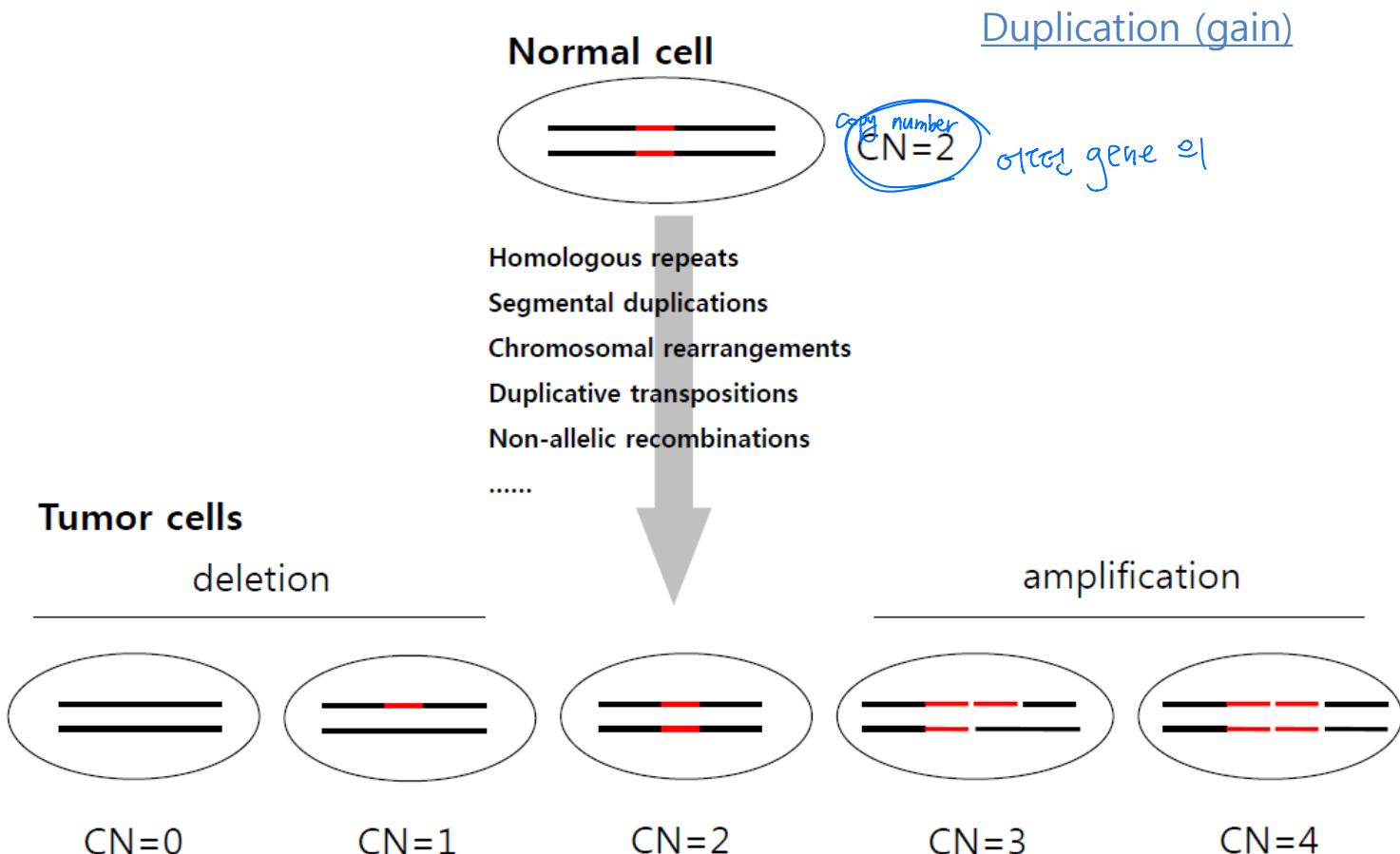
ex oncogene

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

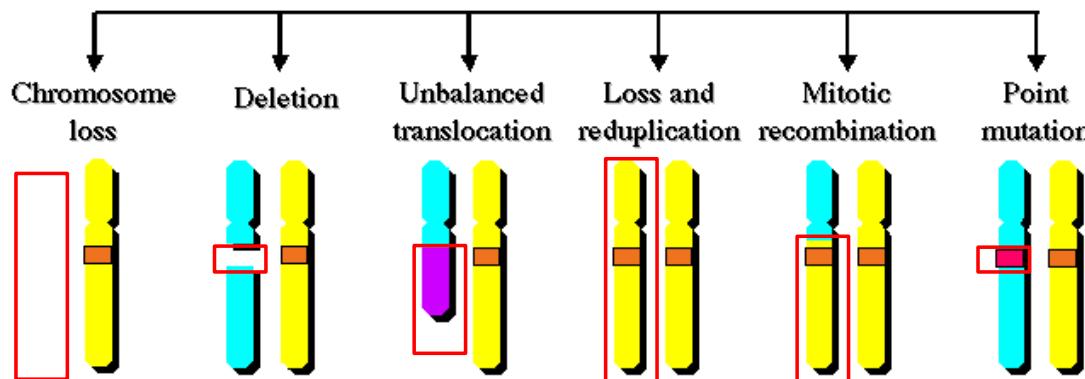
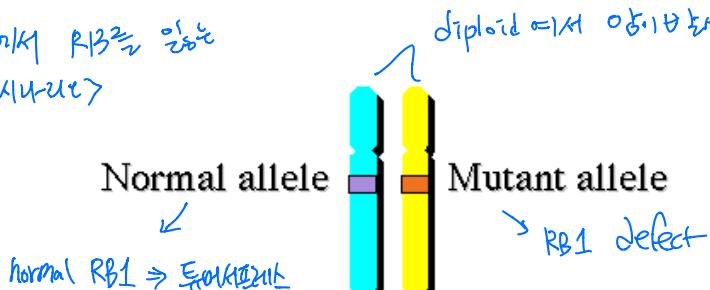
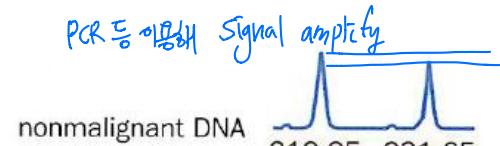
- Somatic genetic changes in retinoblastoma caused loss of heterozygosity (LOH) at markers close to the RB1 locus (well known tumor suppress gene)
- By screening paired blood and tumor samples with markers spaced across the genome, we may discover the locations of tumor suppressor genes

(tumor only RB1는 있는
아니거나...)

diploid ⇒ 2개의 다른 형질을 갖는

여러 가지가 있음

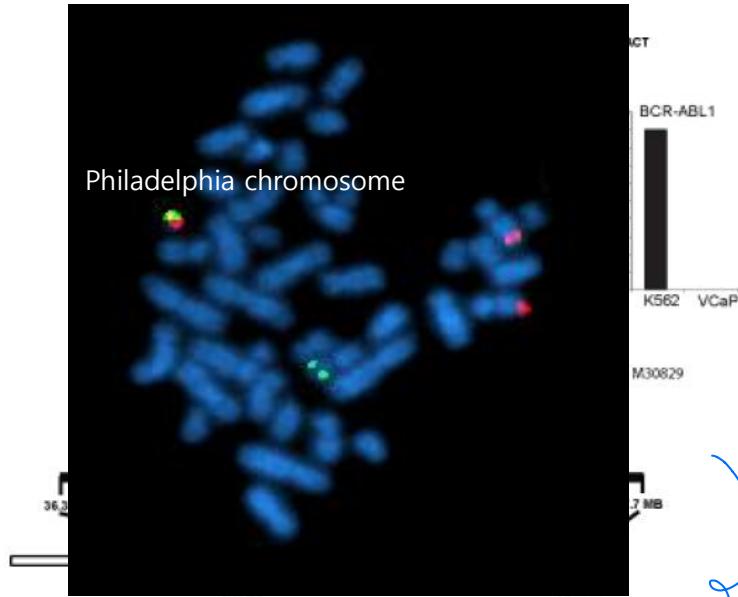
PCR 등 이용해 Signal amplify



{

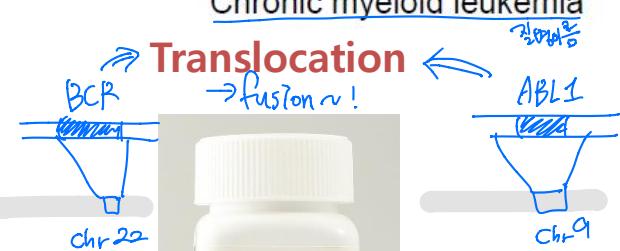
Sequencing ⇒
특정 region의 발현 정도
이상적 tumor 표본
gene 등의 유전적 특징 분석

Fusion Gene in Cancer

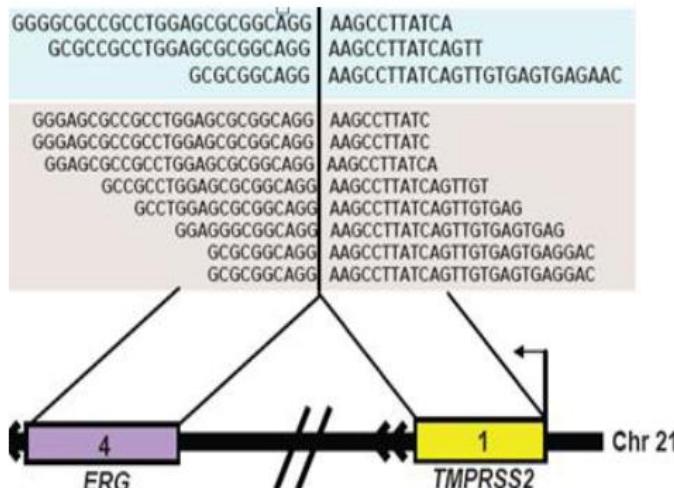


BCR-ABL1 caused by **translocation**
Chronic myeloid leukemia

Translocation
→ fusion ↗!

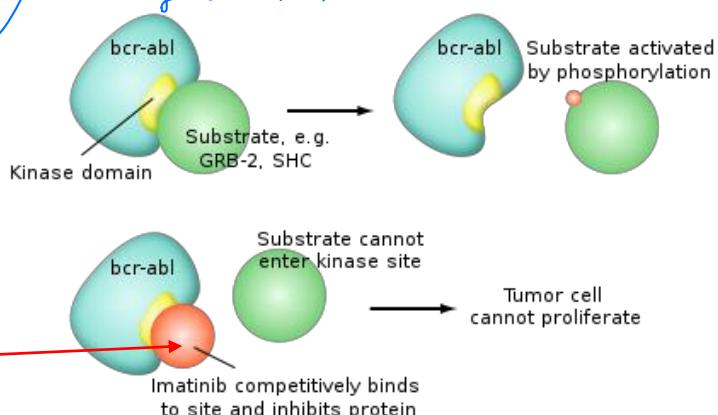


anti Cancer drug cure GRB



TMPRSS2-ERG caused by **deletion**
Prostate cancer

working as unusual

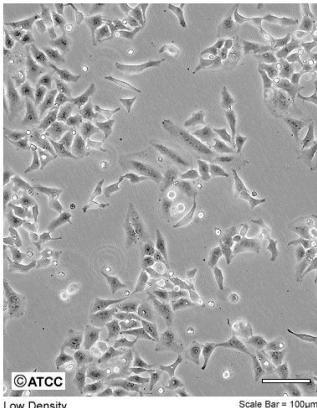


Analysis of Cancer Genome (example of HeLa cell)

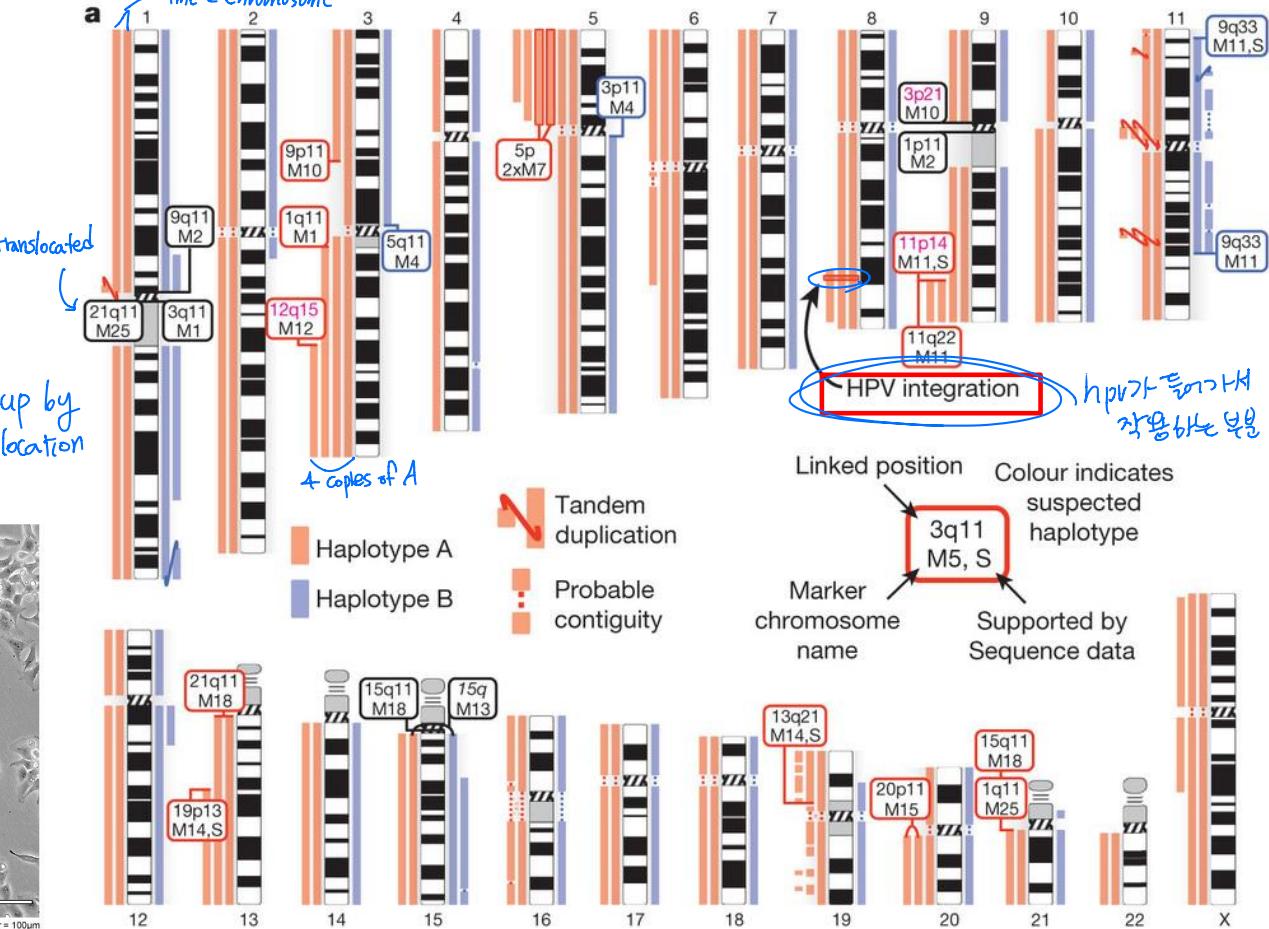


Henrietta Lacks
1945–1951

ATCC Number: CCL-2
Designation: HeLa



NGS3 cervical cancer Chromosome 8p11 normal 9p12
line = chromosome



HPV integration > Cervical cancer

Nature Vol: 500, Pages:207–211, 2013

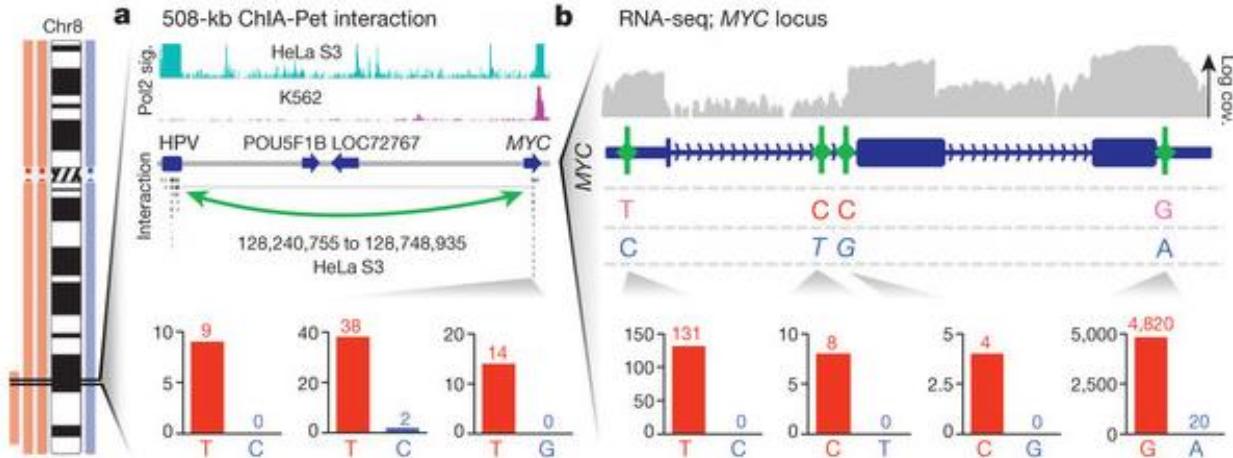
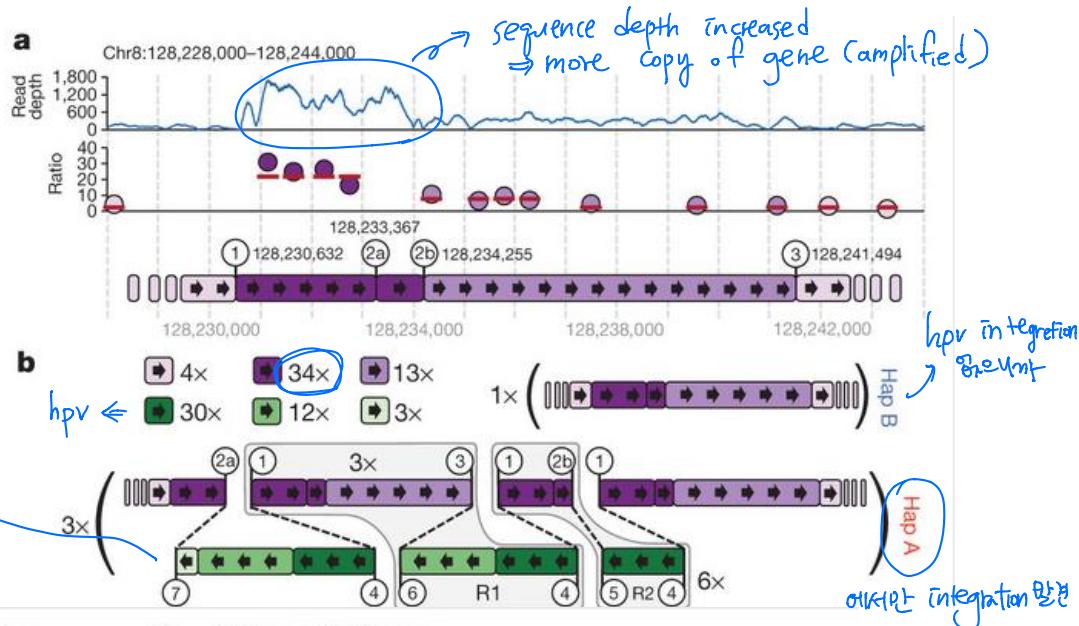
Analysis of haplotype, HPV CNVs, effect on Myc oncogene expression

HPV integration

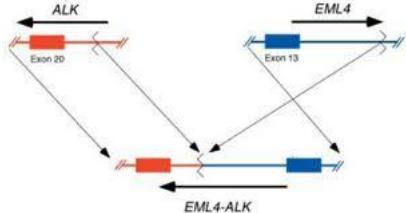
virus는 강력한 promoter 가지고 있음
여러 copy가 integrated
= 강력한 promoter의 copy
= 그로 gene의 expression↑

> Myc oncogene overexpression

> Cervical cancer (HeLa)

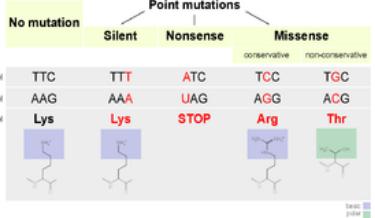


Application of NGS for translational genomics



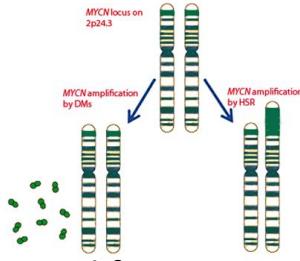
Translocation

BCR-ABL
EML4-ALK
RET/PTC



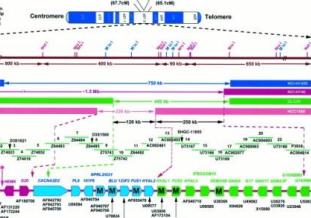
Point mutation

TP53
KRAS
BRAF



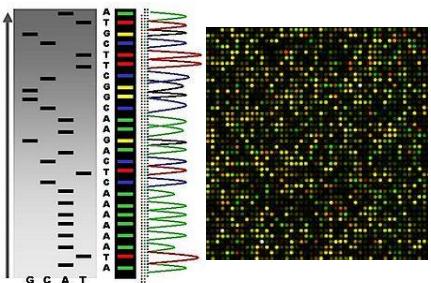
Amplification

EGFR
Myc

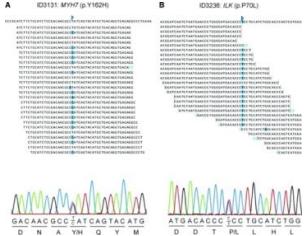


Deletion

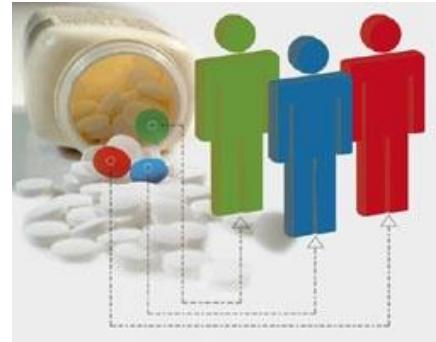
Rb
TP53
PTEN



Clinical sequencing



Personalized cancer care



Genomics and Medicine : 100,000 Genome Project in England

need lots of database



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

Cancer and rare diseases



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](#)

70000 정도는
유전증과 질환과 같은

Start in late 2012, aim to finish by 2017

still working on

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

Whole genome sequencing (NGS)

for analyzing sequence variations

Sung Wook Chi

Division of Life Sciences, Korea University

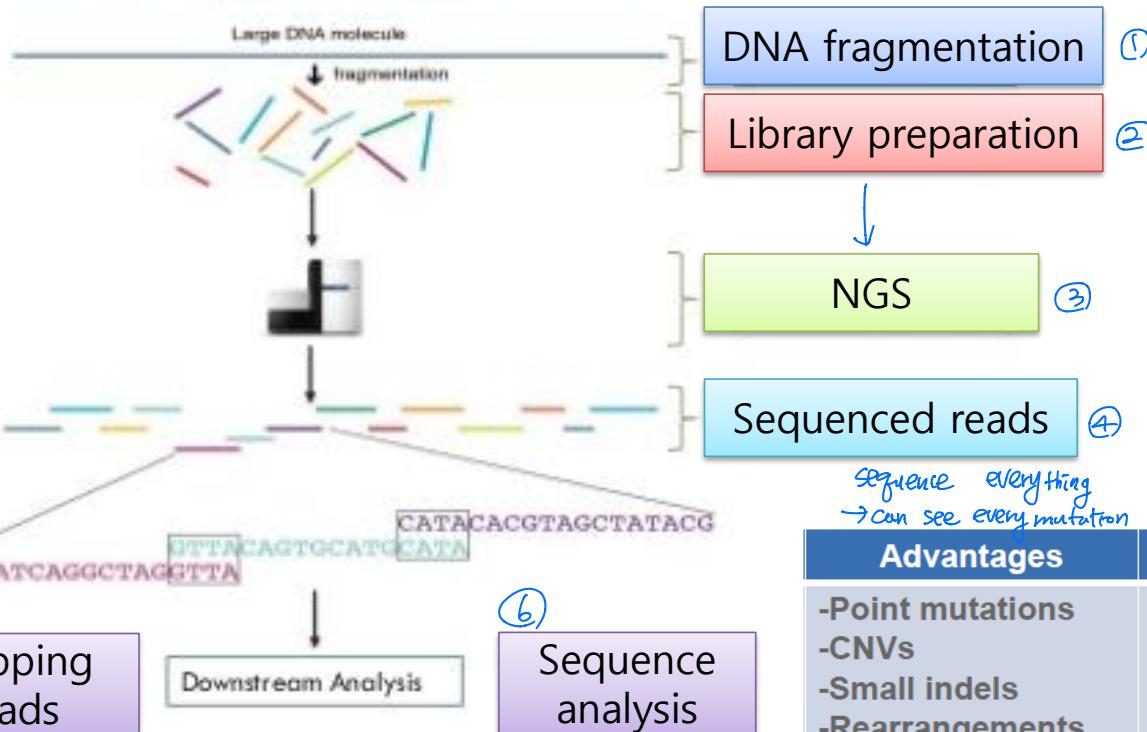
What we will learn today

Whole genome sequencing (NGS)
for analyzing sequence variations

- 1. Whole genome sequencing (WGS by NGS)**
- 2. Variation Analysis for NGS data** *how we can analyze*

Whole Genome Sequencing (NGS) for variant identification

Principle of WGS



Sequencer - Illumina HiSeq 1500
Technique - Paired-End sequencing
Coverage - 100x
Read Length - 100-250 basepairs

sequence everything
→ can see every mutation

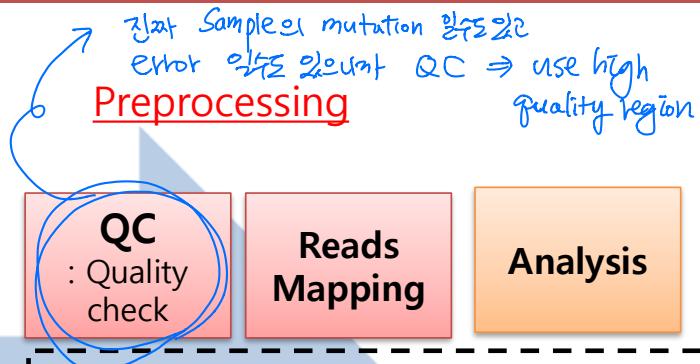
Advantages	Disadvantages
<ul style="list-style-type: none">-Point mutations-CNVs-Small indels-Rearrangements-Somatic mutations in non-coding regions (promoters, enhancers, and non-coding RNAs)	<ul style="list-style-type: none">-Point mutations and indels: >30-fold haploid coverage at least 30X depth-Rearrangements: >10-fold physical coverage

Error → high coverage
error = mutation 일치 잘 등록
→ coverage 높아야 함

Detect every variation

Genomic variation analysis by NGS

→ Increase price



Library preparation

Sample preparation

- DNA
- RNA
- Small RNA

Library preparation

- Target enrichment
- Sample barcoding

Sequence generation

Primary Data analysis

- Image processing
- Signal processing

Raw sequence data QC

Secondary data analysis

- Assembly
- Mapping

Application specific analysis

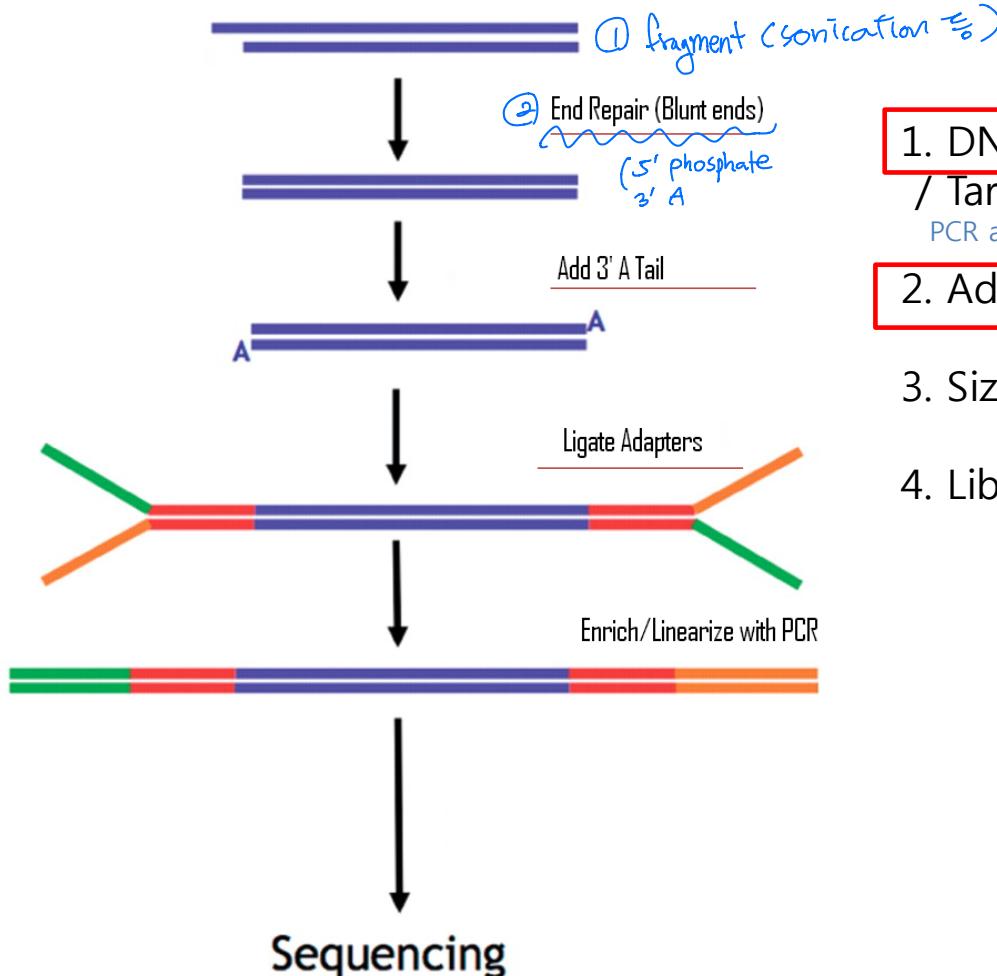
- SNP calling
- Structural variation
- Read counting



NGS data analysis for variant identification

NGS Library preparation

Shear Genomic DNA or begin with cDNA



1. DNA fragmentation / Target Selection

PCR amplification for specific region of interest

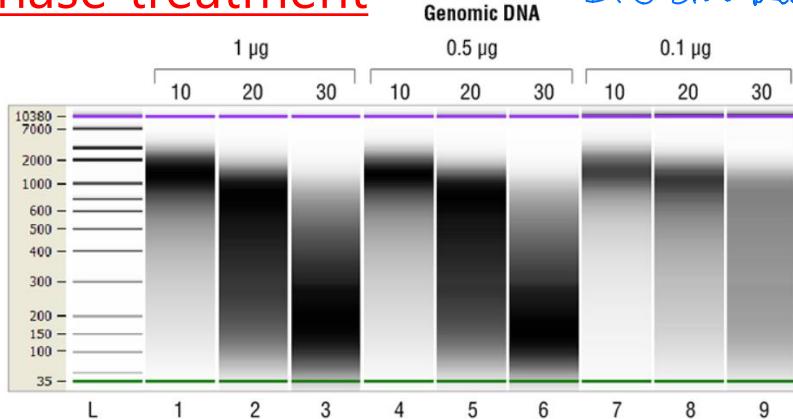
2. Adapter ligation

3. Size selection

4. Library quantification (QC)

DNA fragmentation for library construction

Dnase treatment



DNA Input
Digestion Time
Minutes



or



or



Nebulized,
sonicated or
sheared DNA



Repair

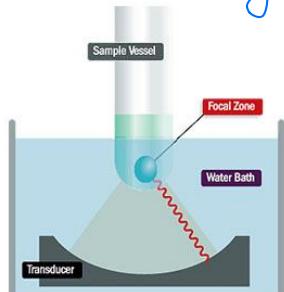
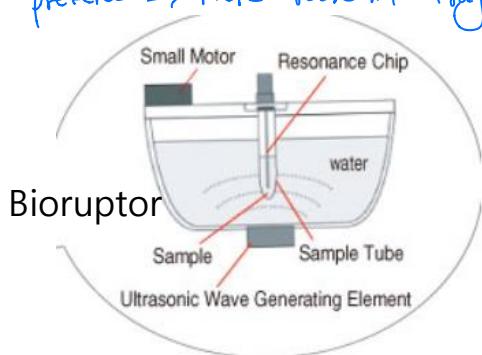
Blunt-ended,
5'-phosphorylated DNA



adaptor
ligation

Double strand :
Polymerase (Klenow fragment)

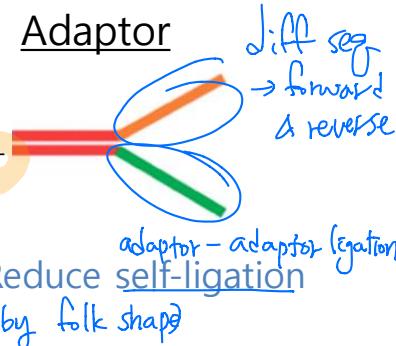
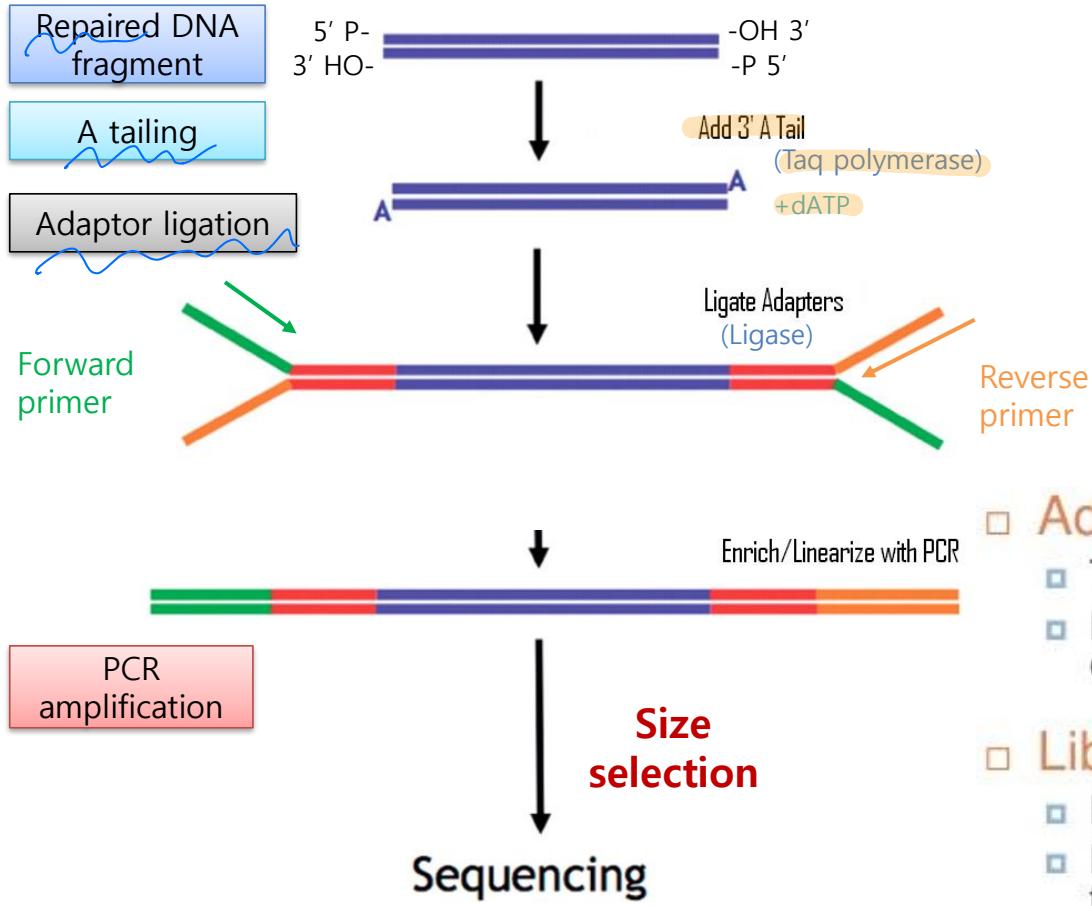
5'P : PNK (T4 Polynucleotide Kinase)
→ attach phosphate to 5'



Sonication (shearing)

preferred → more random fragment than Dnase enzyme

A-tailing, ligation, PCR amplification : Library preparation



- **Adapter ligation**
 - T-overhangs
 - Forked structure controls orientation
- **Library amplification**
 - Few cycles
 - Enrich for correctly-adapted fragments

NGS library preparation & Illumina Sequencing



Size selection

- Gel electrophoresis
- Bead based method

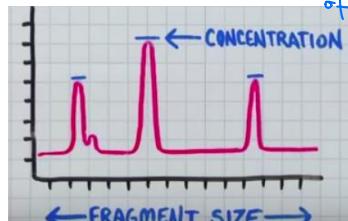
DNA quantification

: Accurate concentration
(Mole number)

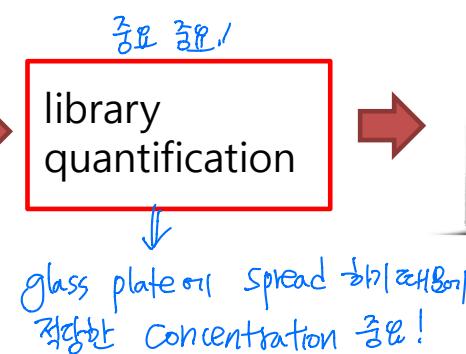
1) Spectrophotometer

- Absorbance : amount of DNA

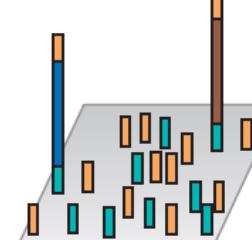
2) Bioanalyzer → measure concentration of library as accurate as possible



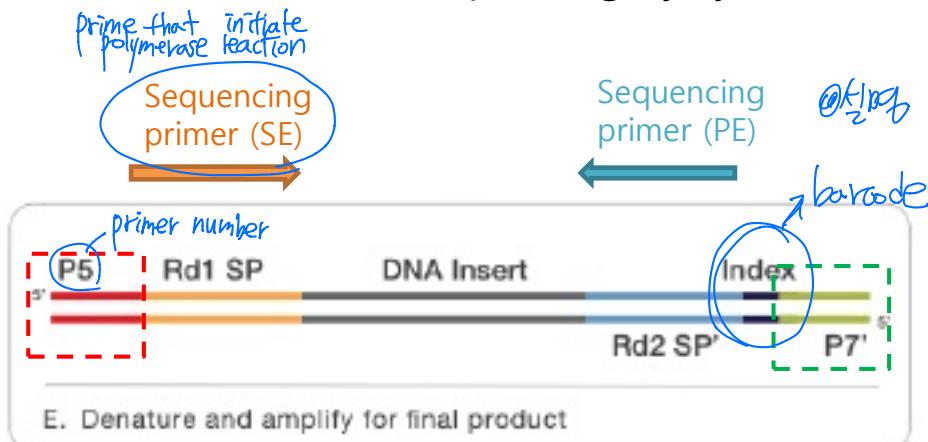
3) qPCR



Illumina NGS sequencing



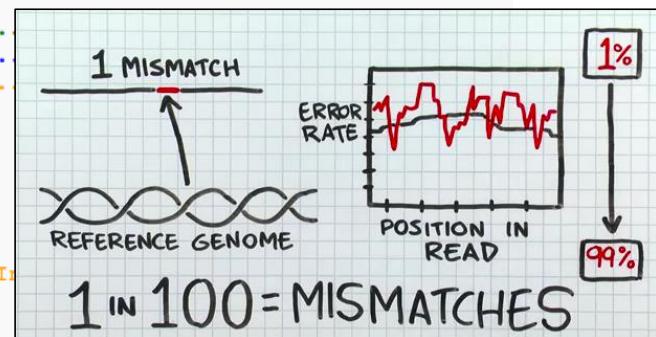
SBS (Sequencing by synthesis)



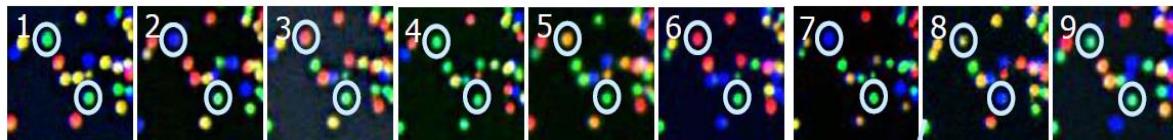
Sequencing data from NGS : FASTQ format & Q score



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 I
(Note: See discussion above).
L - Illumina 1.8+ Phred+33, raw reads typically (0, 41)



Sequence Quality Control

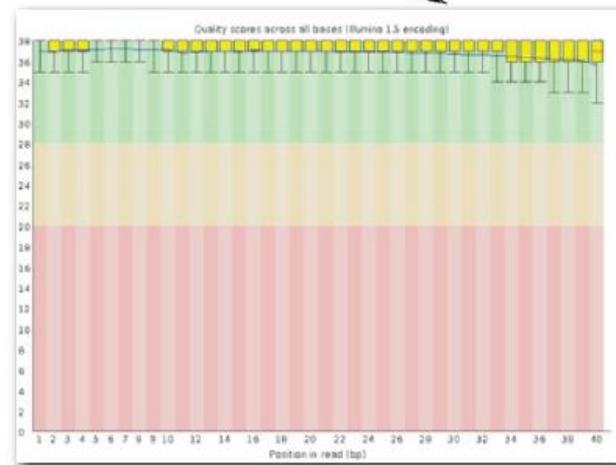


per base sequence quality

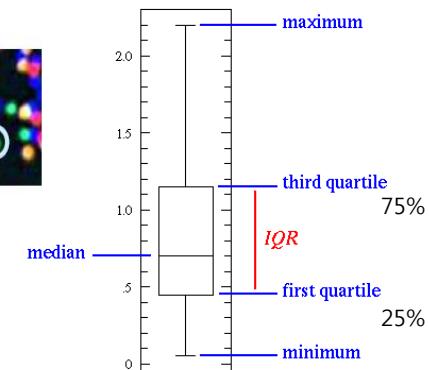
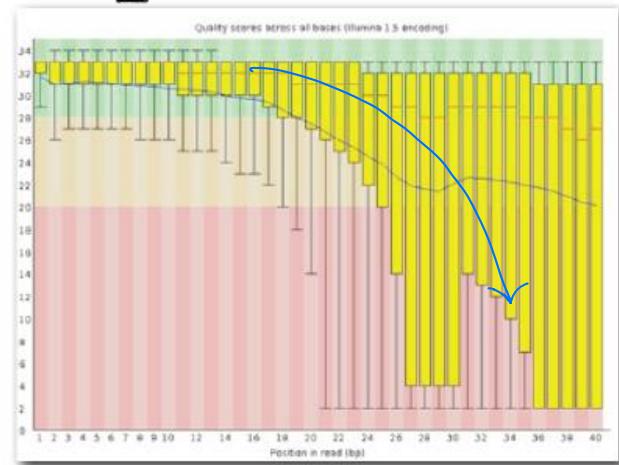
good

bad

Q score



Q score

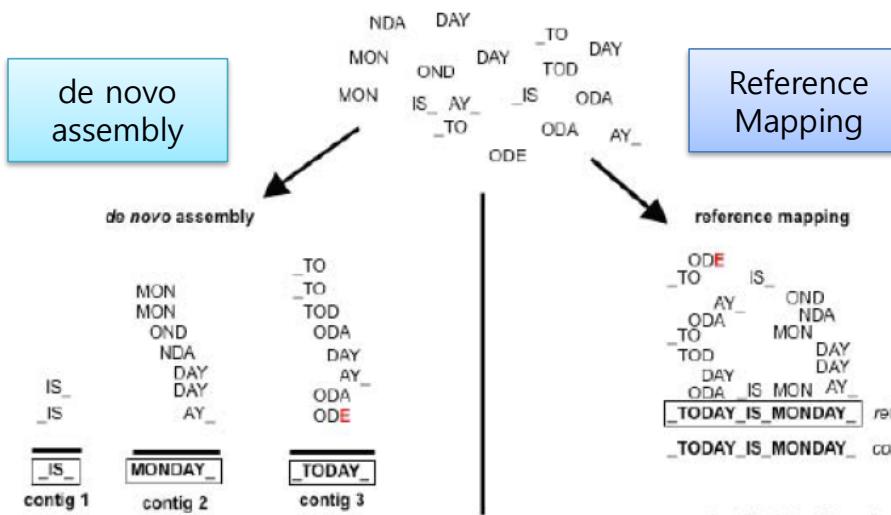


Position in reads

Position in reads

Q score tends to decrease depending on the increment of position

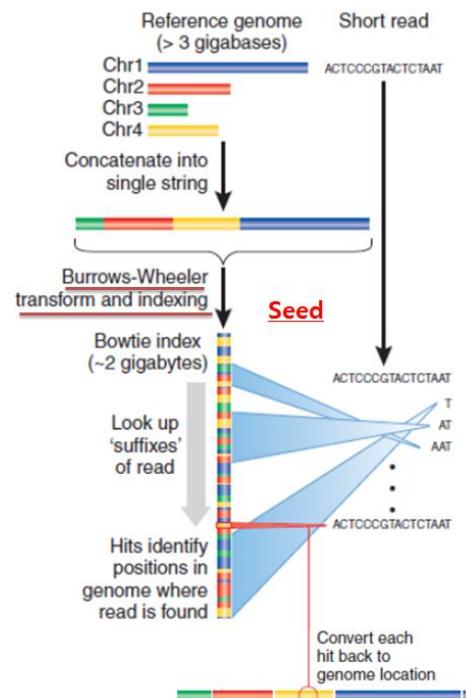
NGS Sequence mapping



NGS WGS reads mapping (BWT based method)

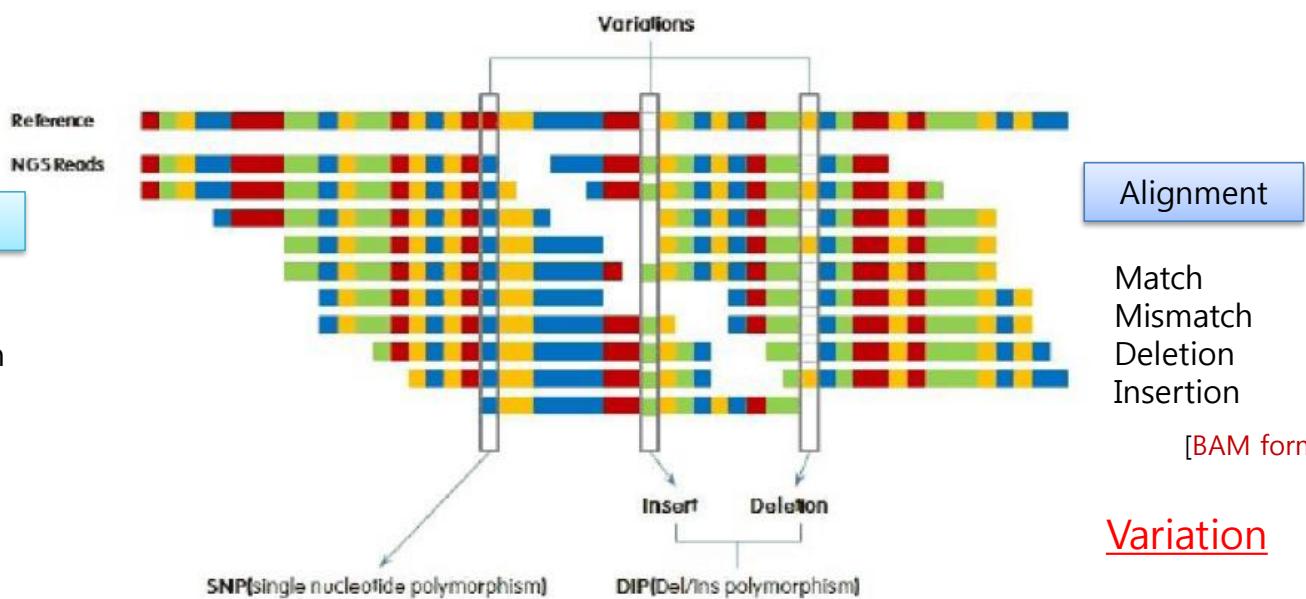
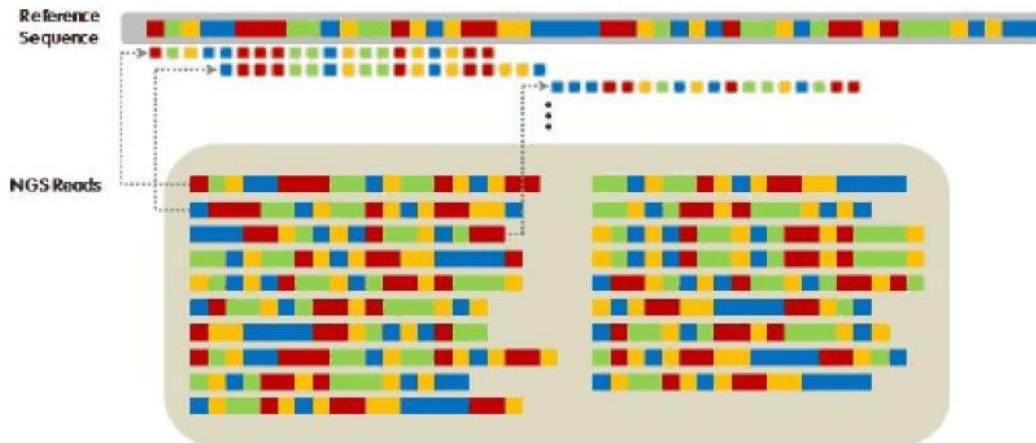
- Transform genome sequences into data structure (condensed, indexed, fast search)
- Find billions of short reads sequence there (near perfect match)
- Location of matched reads (Genomic coordination, Alignment)

Mapping billions of short reads on genome

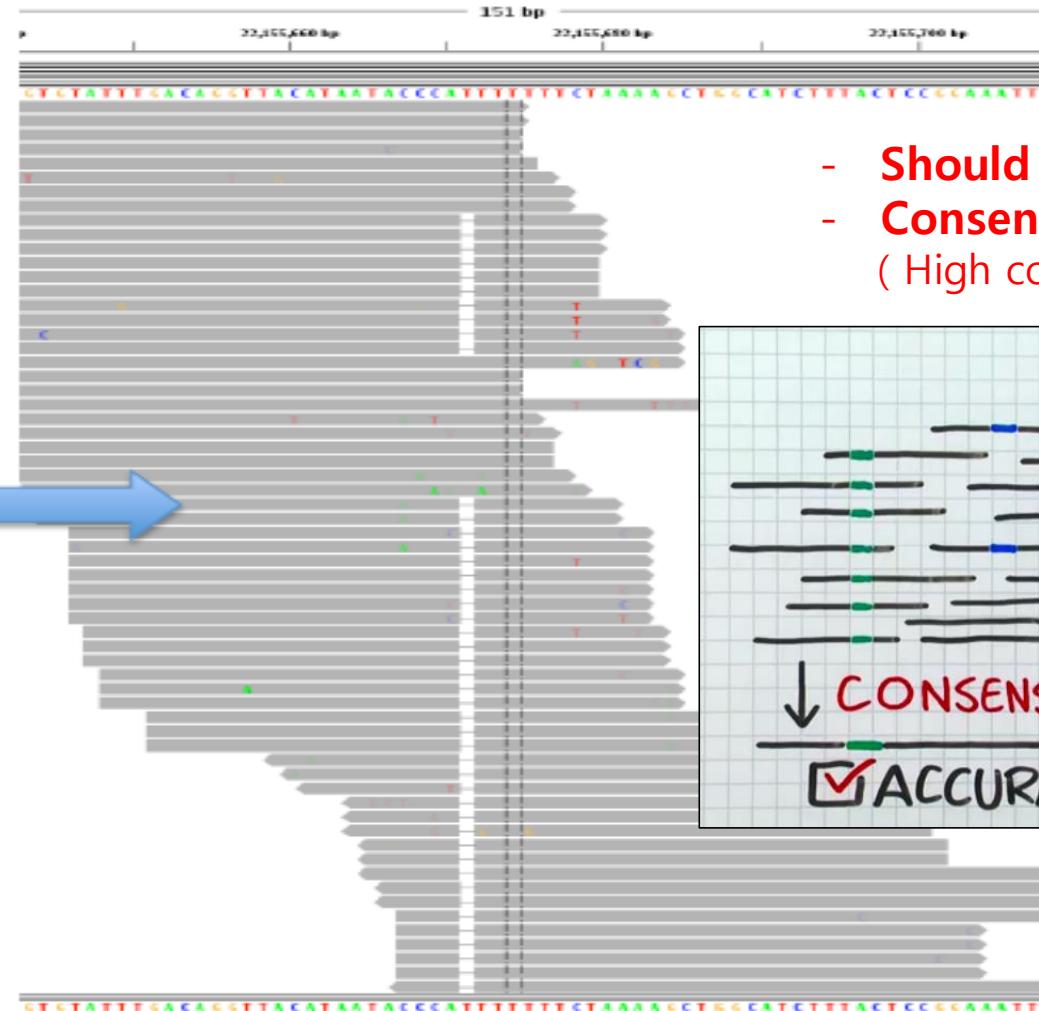


BWT (Burrow-Wheeler transform) based mapping algorithm

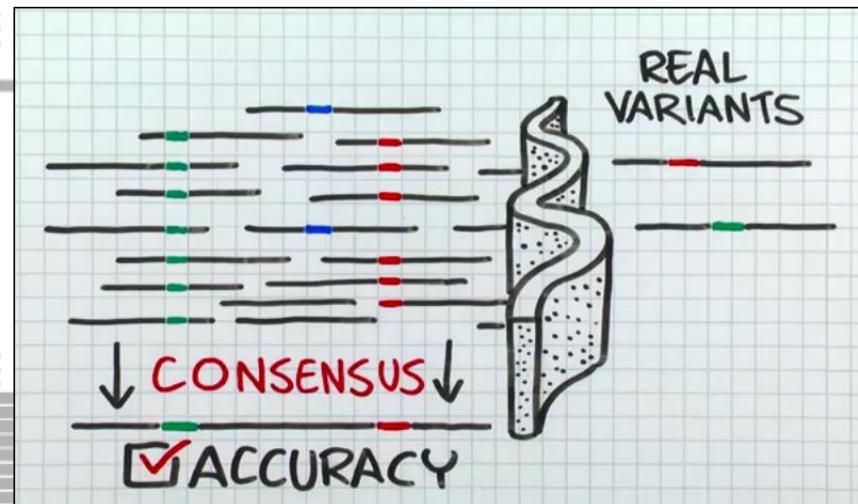
Sequence Alignment (Mapping)



Variation in NGS reads : Real or not ?

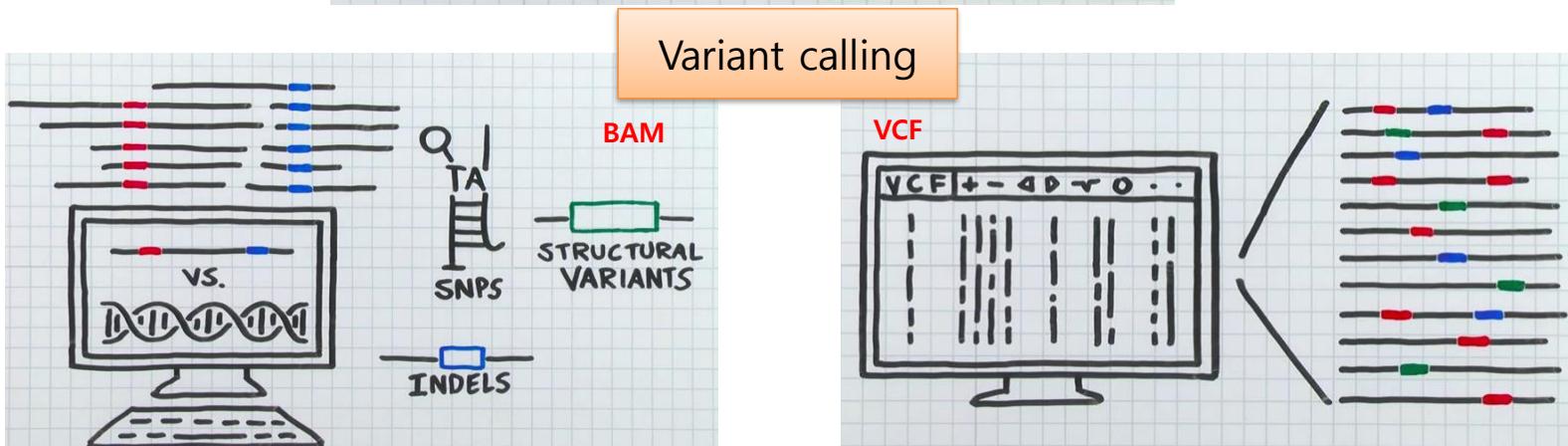
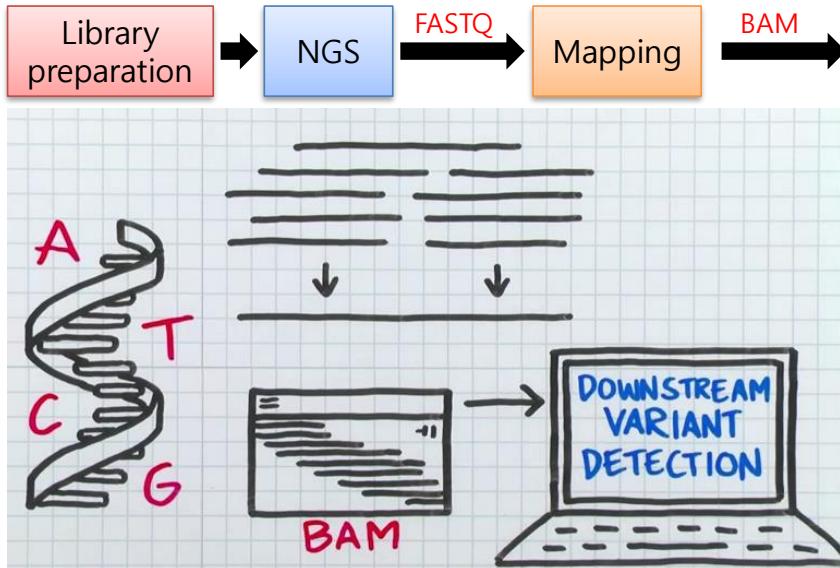


- Should consider Q-score
- Consensus
(High coverage > Statistics)



“Variant calling”

NGS WGS analysis for variation detection

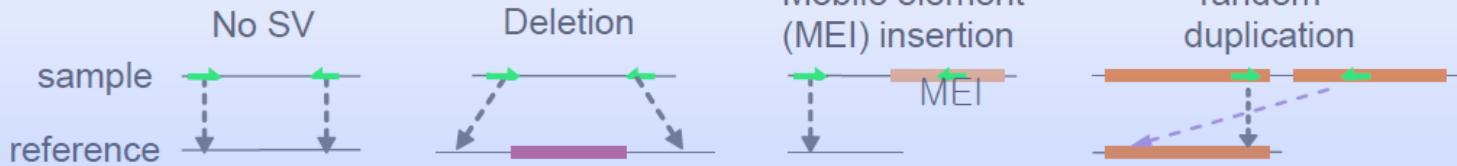


SV Discovery with diverse resources

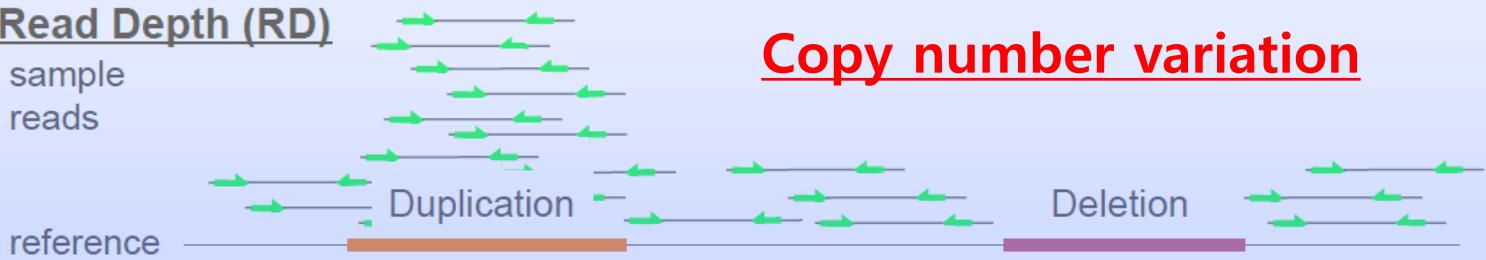
Discovery of Structural Variation

- Challenging in mapping
- Difficult for interpretation

Read Pairs (RP)



Read Depth (RD)

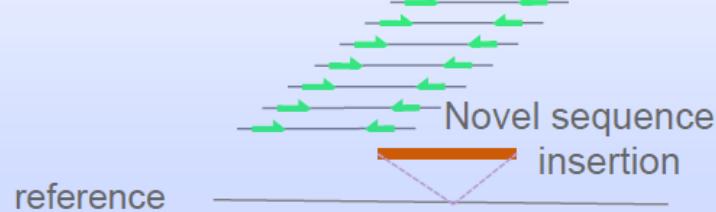


Copy number variation

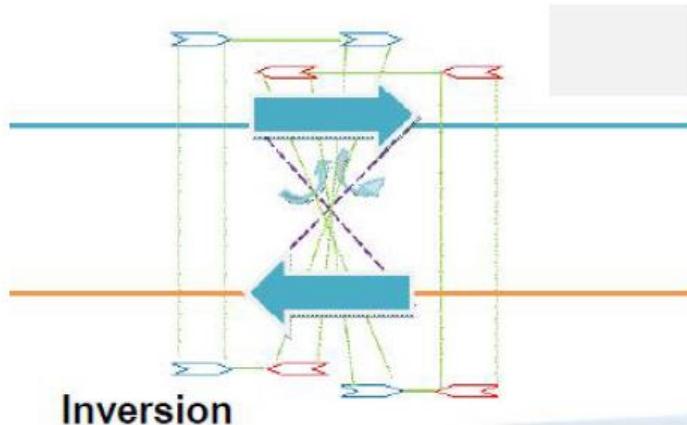
Split Reads (SR)



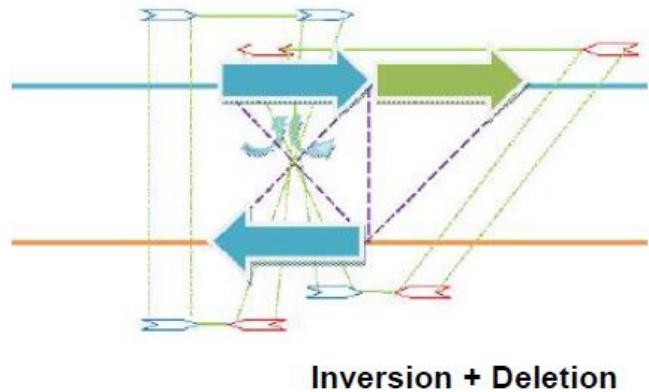
Assembly (AS)



SV Discovery with diverse resources



Inversion



Inversion + Deletion

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

