



Genome-wide Association Studies

Yongjin Park



a place of mind
THE UNIVERSITY OF BRITISH COLUMBIA

Today's lecture

- Fundamentals in Population Genetics
- Genome-wide association studies
- What are the limitations of GWAS?

Genetics: It all started from pea plants

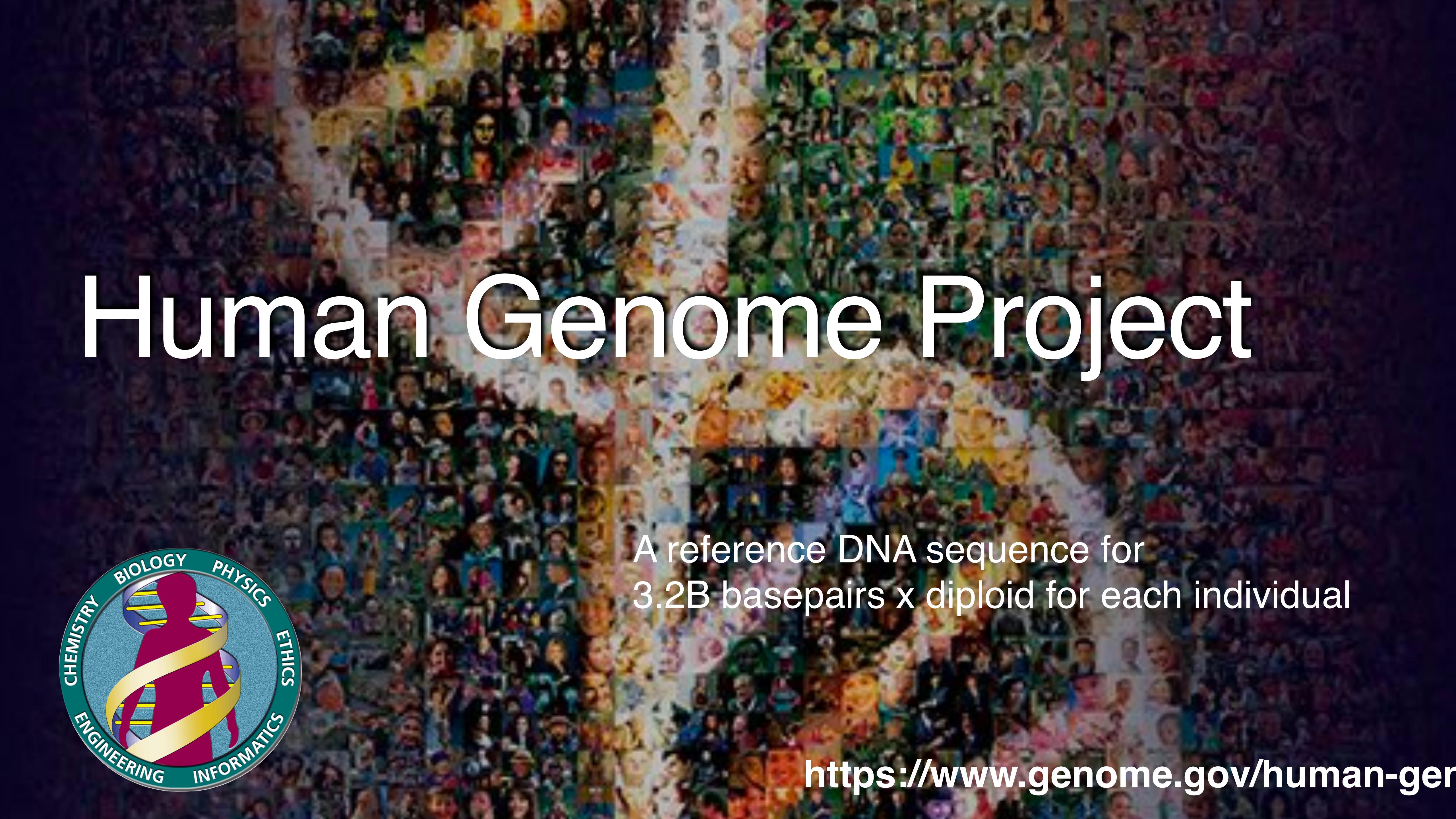
		pollen ♂	
		B	b
pistil ♀	B	BB	Bb
	b	Bb	bb



Gregor Mendel
(1822-1884)

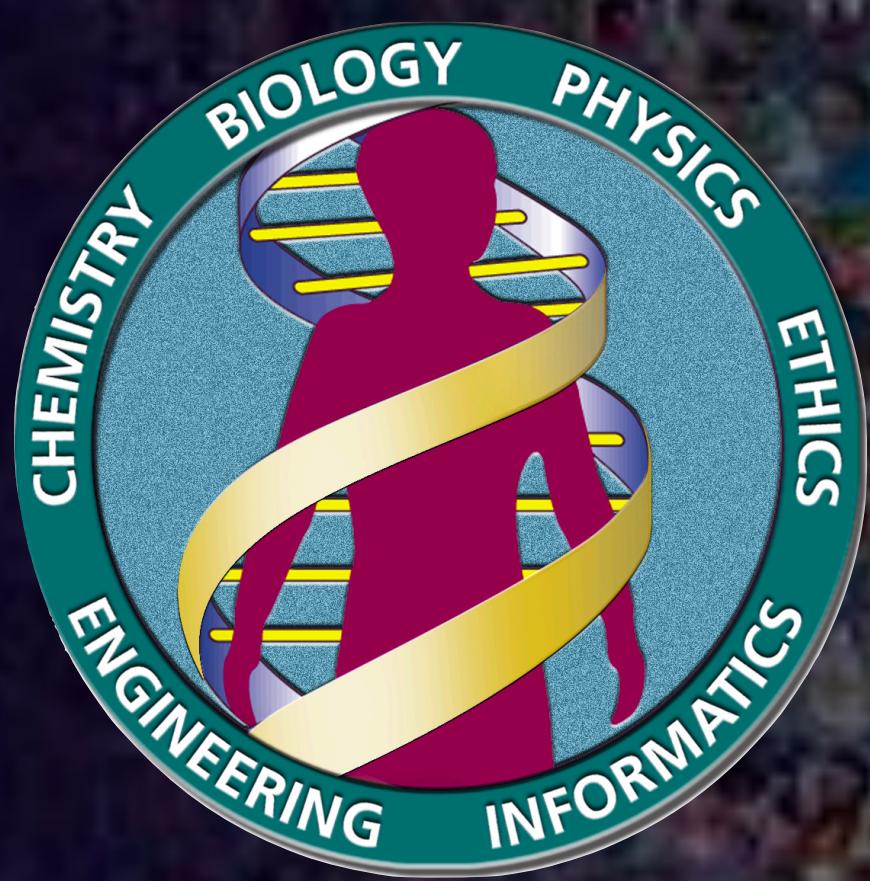
Key concepts emerged:

- Gene = a unit of heredity that transfers from a parent to offspring
- Allele = a different form of a gene [from a Greek word, αλληλο, αλλος, "allos", other]



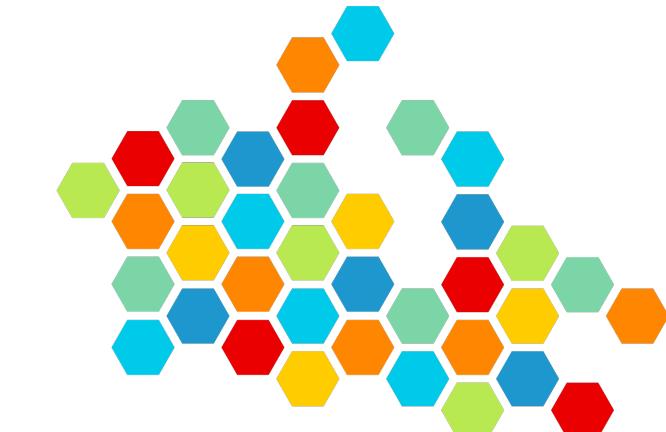
Human Genome Project

A reference DNA sequence for
3.2B basepairs x diploid for each individual



<https://www.genome.gov/human-genome-project>

Human genetics revolution



Canadian Partnership
for Tomorrow's Health



To Solve 3 Cold Cases, This Small County Got a DNA Crash Course

Forensic genealogy helped nab the Golden State Killer in 2018. Now investigators across the country are using it to revisit hundreds of unsolved crimes.



Human Genomes help keep track of human evolution

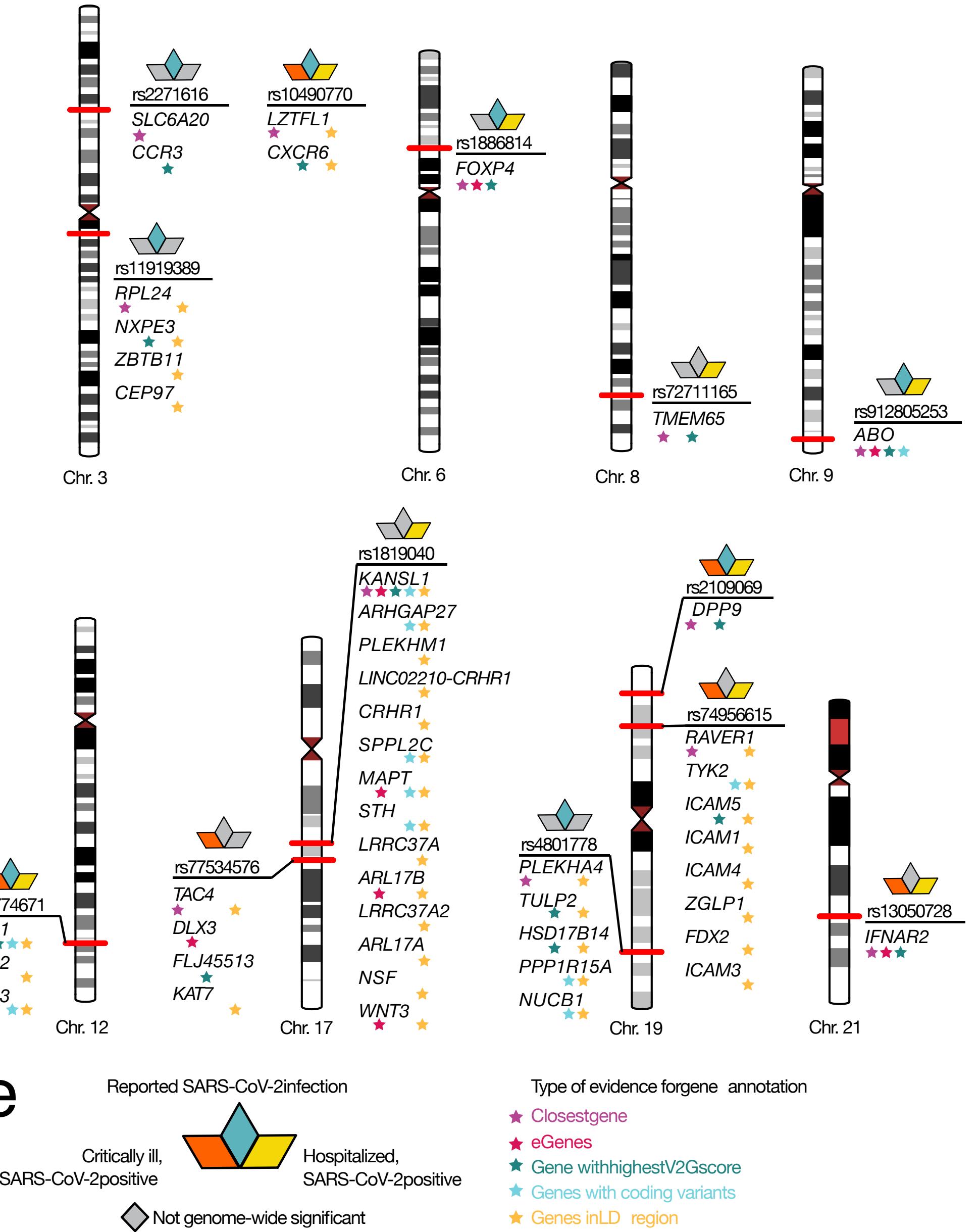


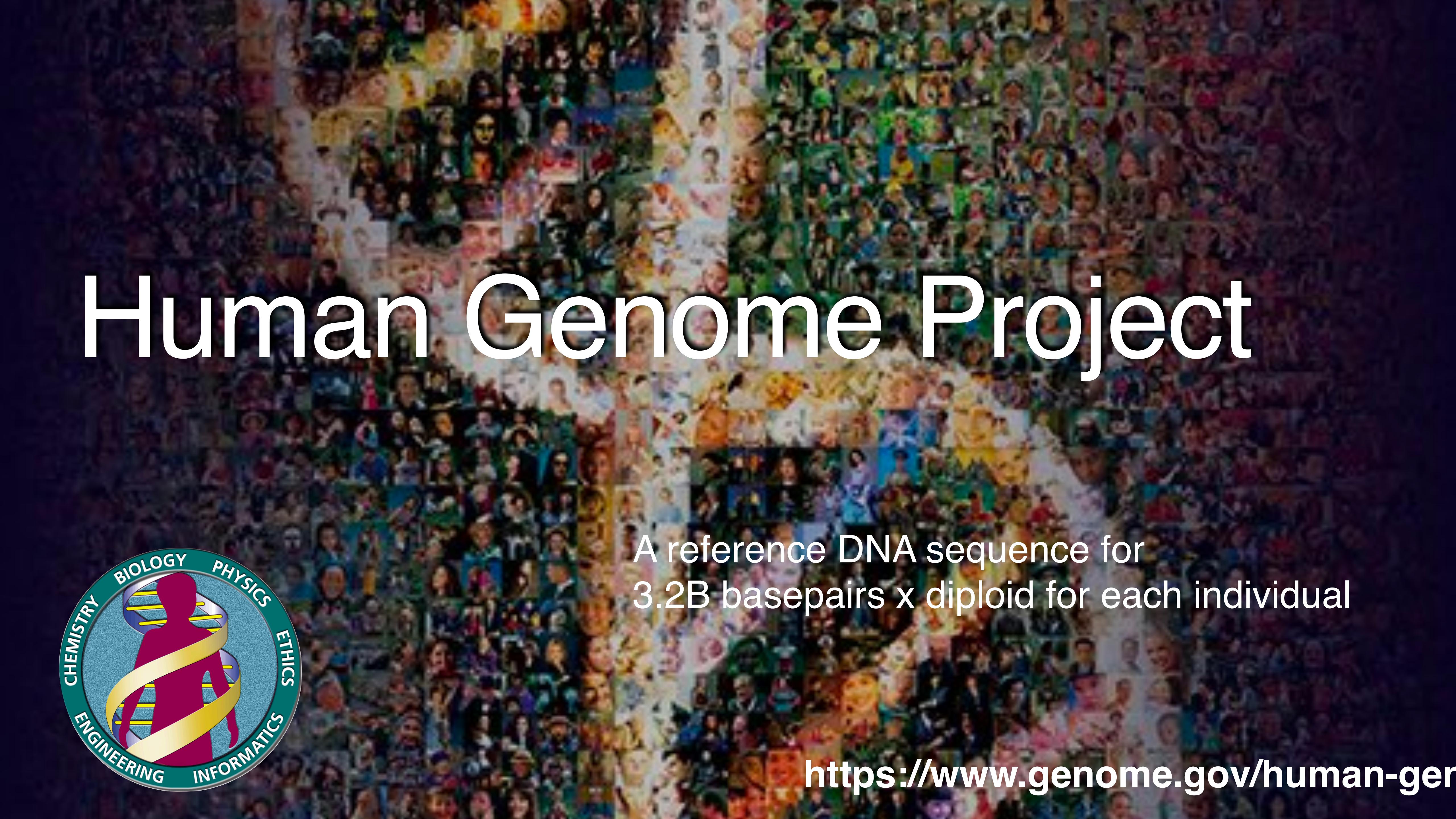
COVID-19 Host Genetics



March 12, 2021,
GWAS paper
in medRxiv

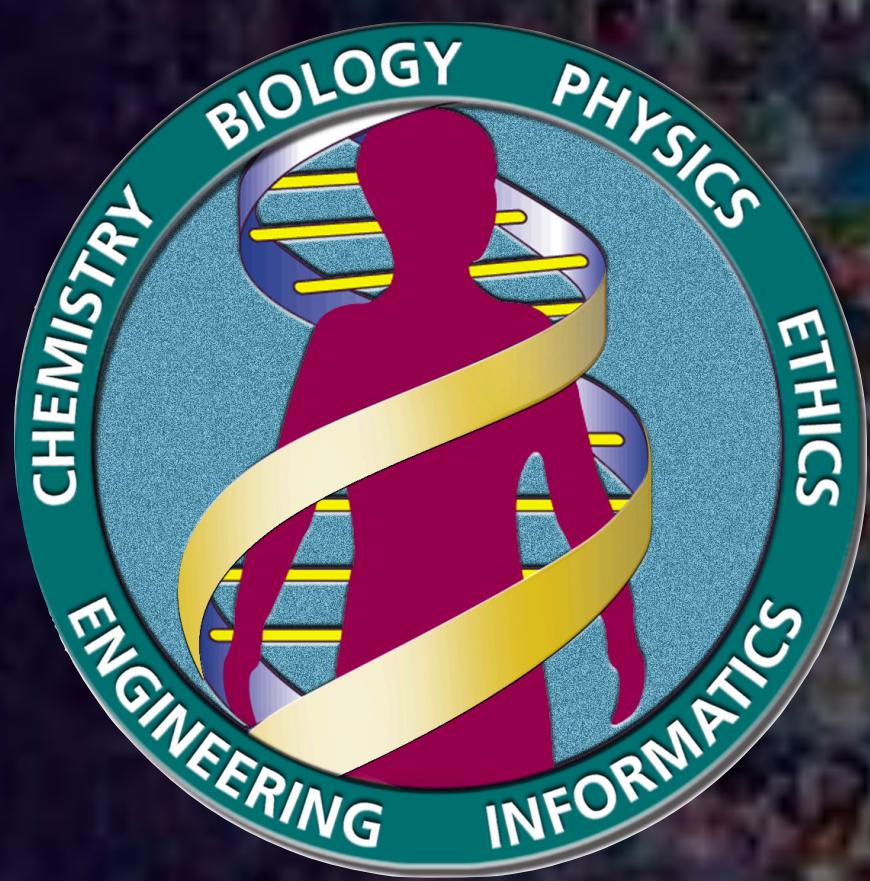
An unprecedented pace
of GWAS profiling





Human Genome Project

A reference DNA sequence for
3.2B basepairs x diploid for each individual



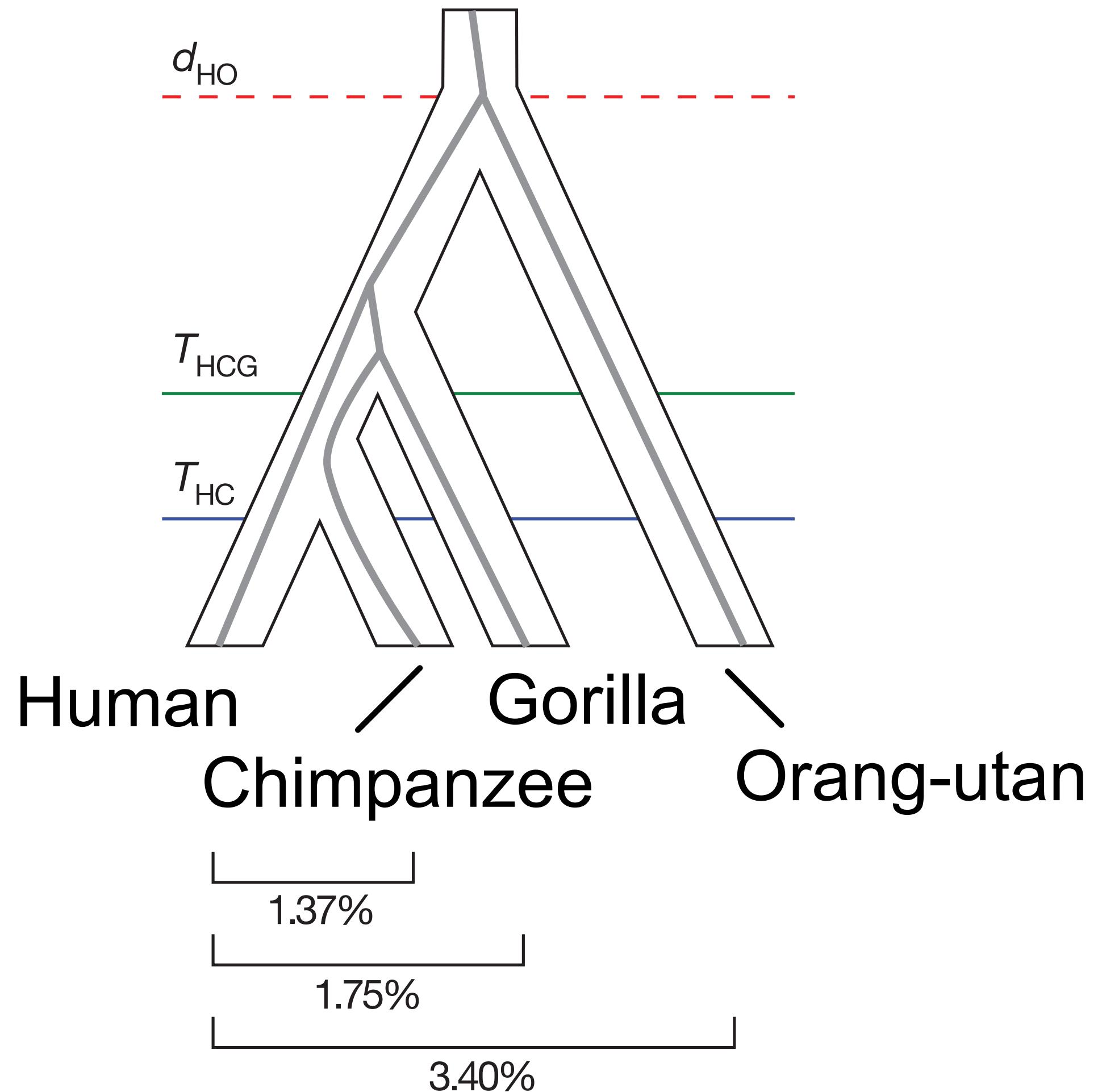
<https://www.genome.gov/human-genome-project>

How did Human Genome Project start the revolution?

We have a reference panel
of genomic sequence
information...

Why is it important to me?

99% genetic information
shared across humans



Scally .. Durbin, *Natur*

If we are 99% identical, then
what is the 1% difference?

Which part of the human genome is variable?

Published: 18 December 2003

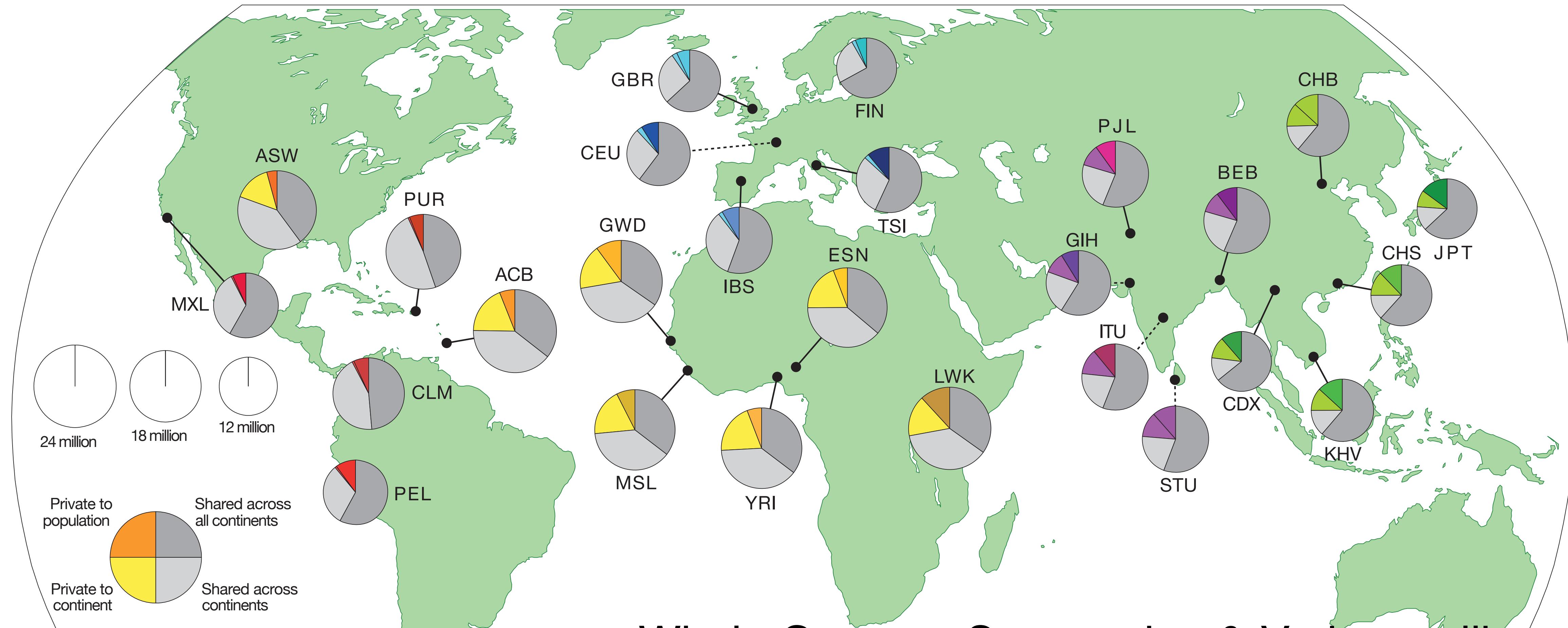
The International HapMap Project

†The International HapMap Consortium

Nature **426**, 789–796 (2003) | Cite this article

80k Accesses | **4231** Citations | **59** Altmetric | Metrics

The 1000 genomes project



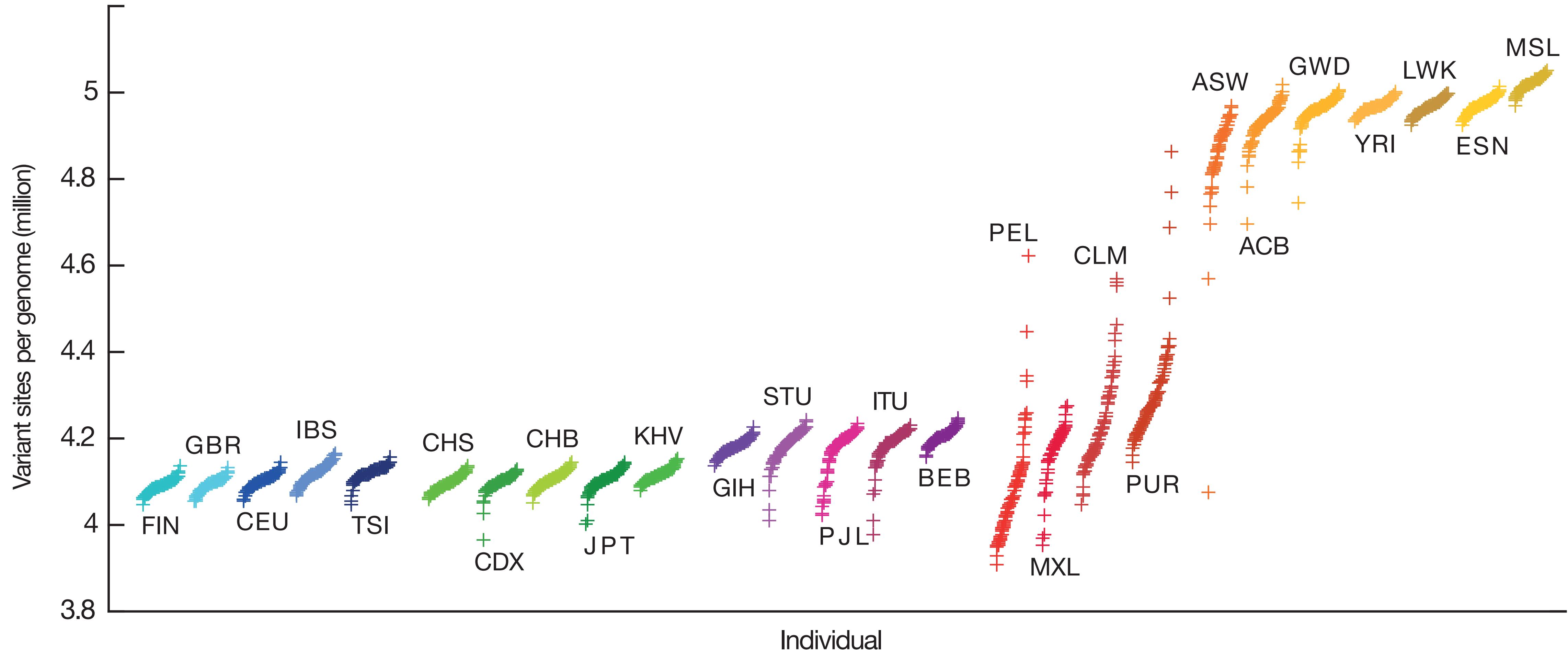
Whole Genome Sequencing & Variant calling
across 1000 individuals in many different groups

Genetic variation across human population

A typical genome

We find that a typical genome differs from the reference human genome at 4.1 million to 5.0 million sites (Fig. 1b and Table 1). Although >99.9% of variants consist of SNPs and short indels, structural variants affect more bases: the typical genome contains an estimated 2,100 to 2,500 structural variants (\sim 1,000 large deletions, \sim 160 copy-number variants, \sim 915 Alu insertions, \sim 128 L1 insertions, \sim 51 SVA insertions, \sim 4 NUMTs, and \sim 10 inversions), affecting \sim 20 million bases of sequence.

A typical genome differs from the reference at 4.1 to 5 million sites



Allele

A different form of a gene
[from a Greek word,
ἀλληλο, ἀλλος, "allos",
other]

A different version of a gene

A different version of the
same variant

Mostly used for a gene

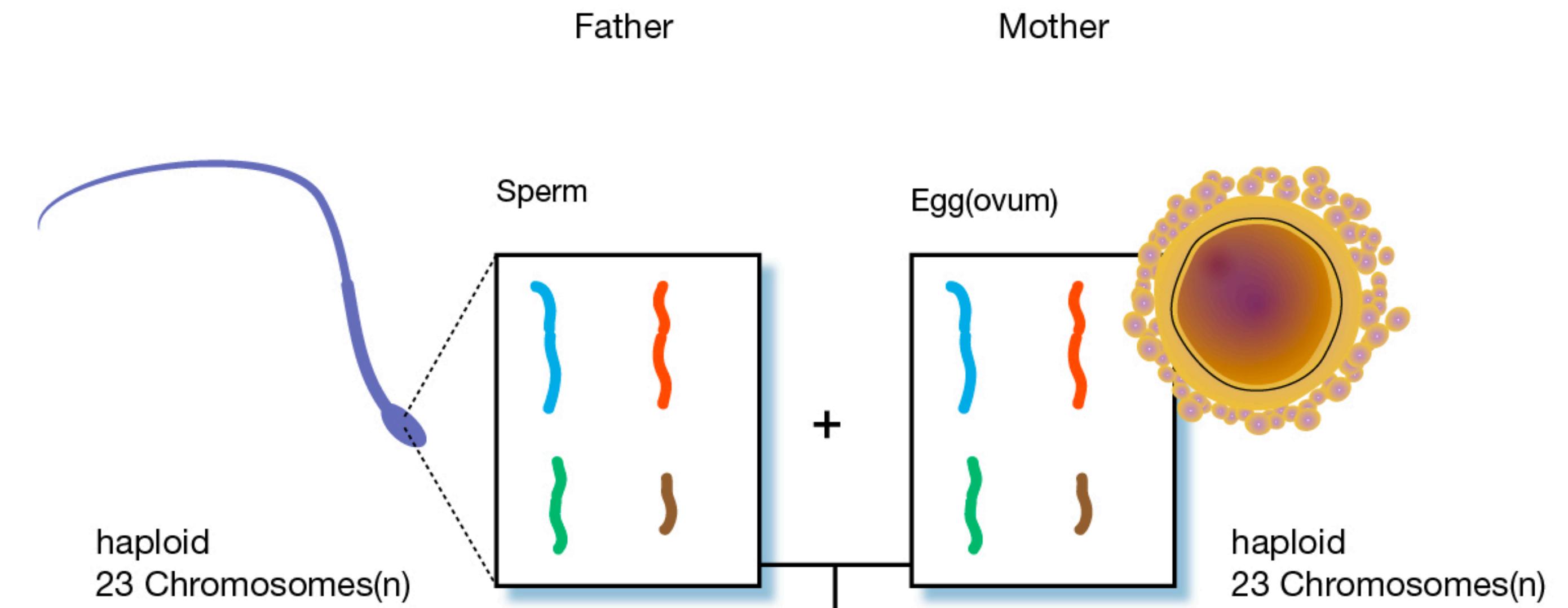
Variant

A specific region of the
genome differs between
two genomes.

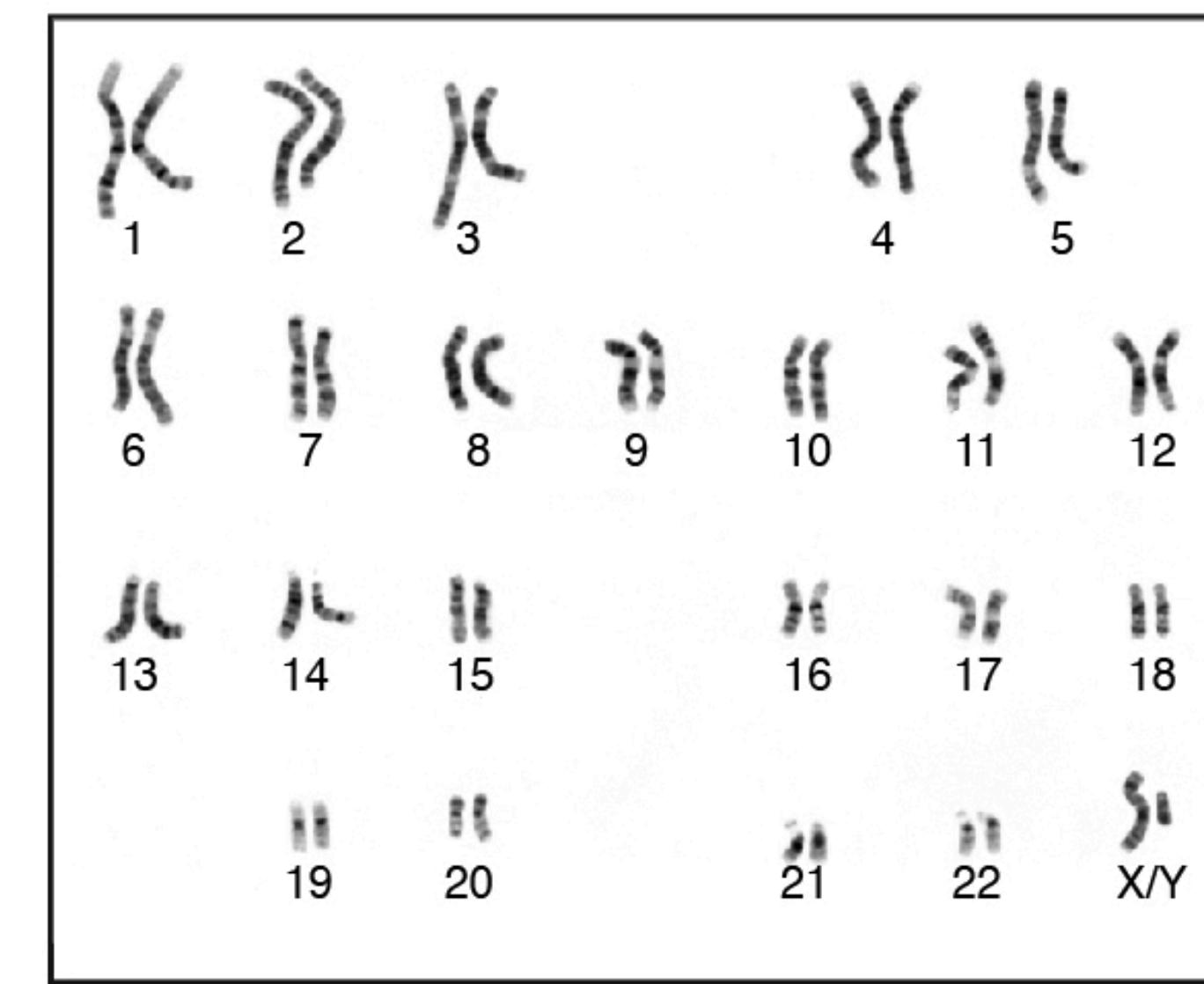
Across two or more
genomes, there could have
genetic variants due to a
mutational process

Ploidy (diploid)

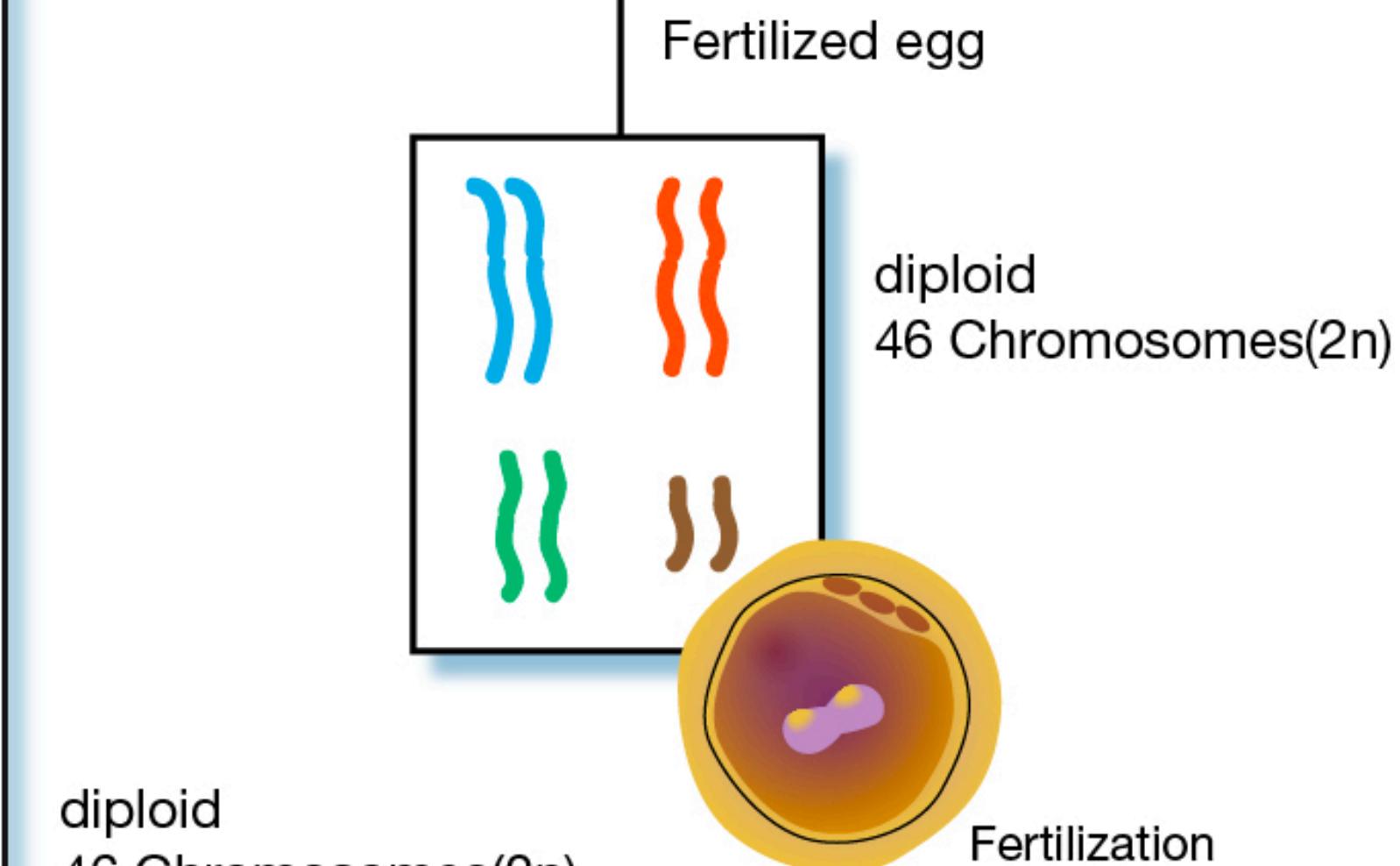
One copy from
maternal genome
another copy from
paternal genome



When do we have
polyploidy or
aneuploidy?



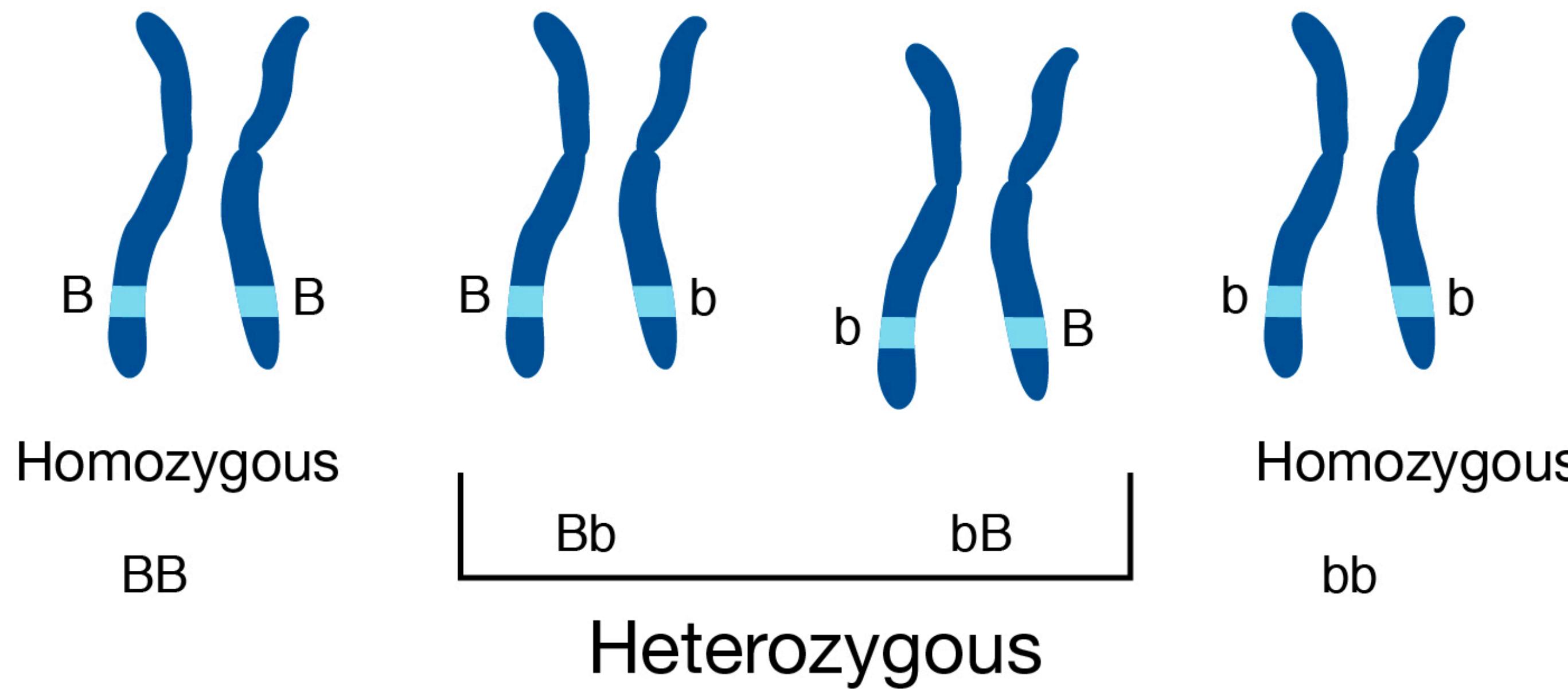
diploid
46 Chromosomes(2n)



Biallelic vs. triallelic vs. multiallelic

"Bi"-allelic: two different forms exist for this heritable unit (most of the human variants)

"Tri"-allelic: three different forms exist (much rare than biallelic variants)



For biallelic variant:

- reference allele
- alternative allele

Many types of genetic variants

ACTCGTGACCGCATGCATCTTCATTGATGC

ACTCGTGACCGCATGCATCGTCAATTGATGC

Reference

Insertion

Deletion

Reference

Alternative

ACTGACGCATGCATCATGCATGC

ACTGACGCATG**GT**A CATCATGCATGC

ACTGACG--TGCATCATGCATGC

Indel

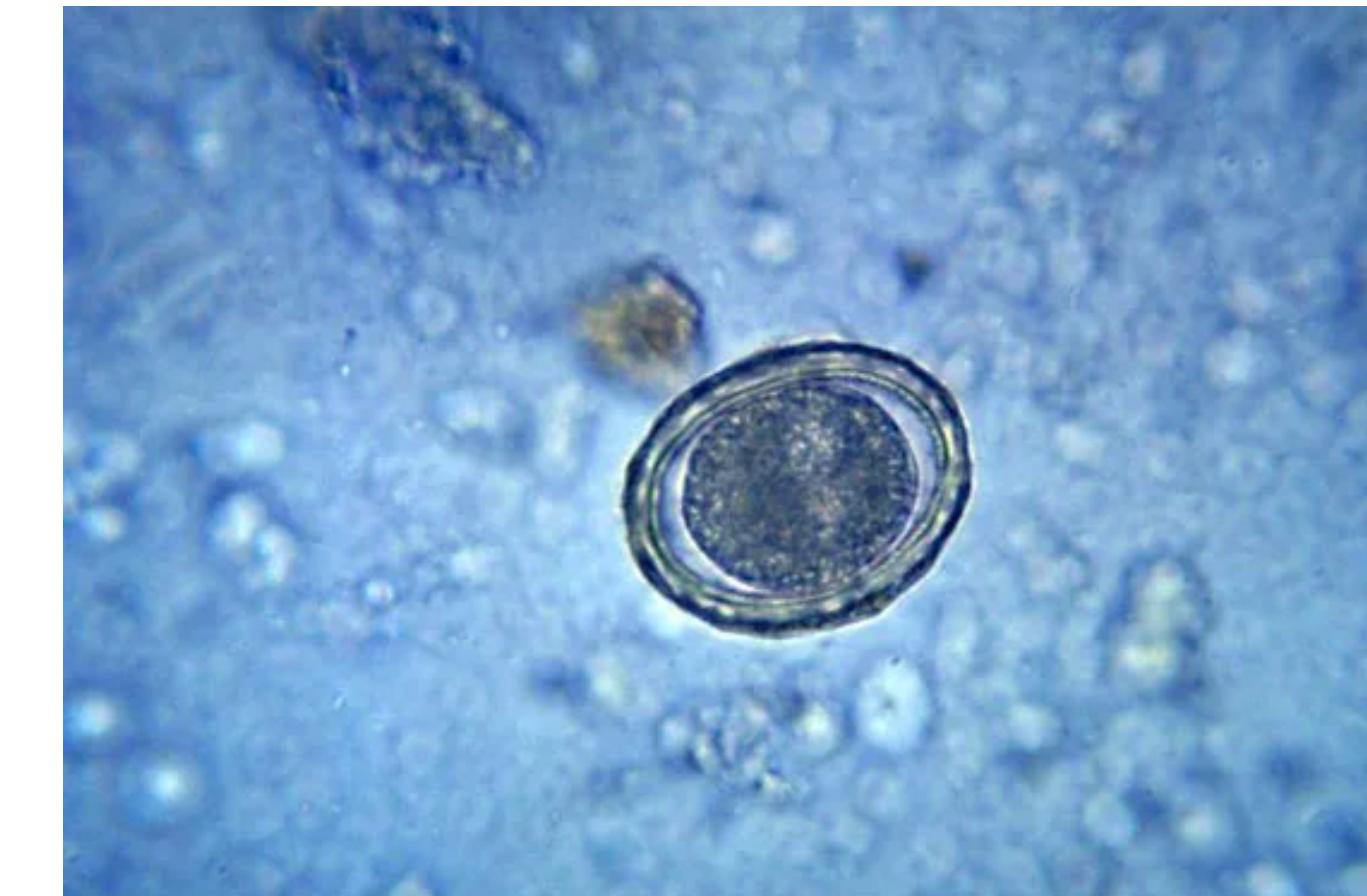
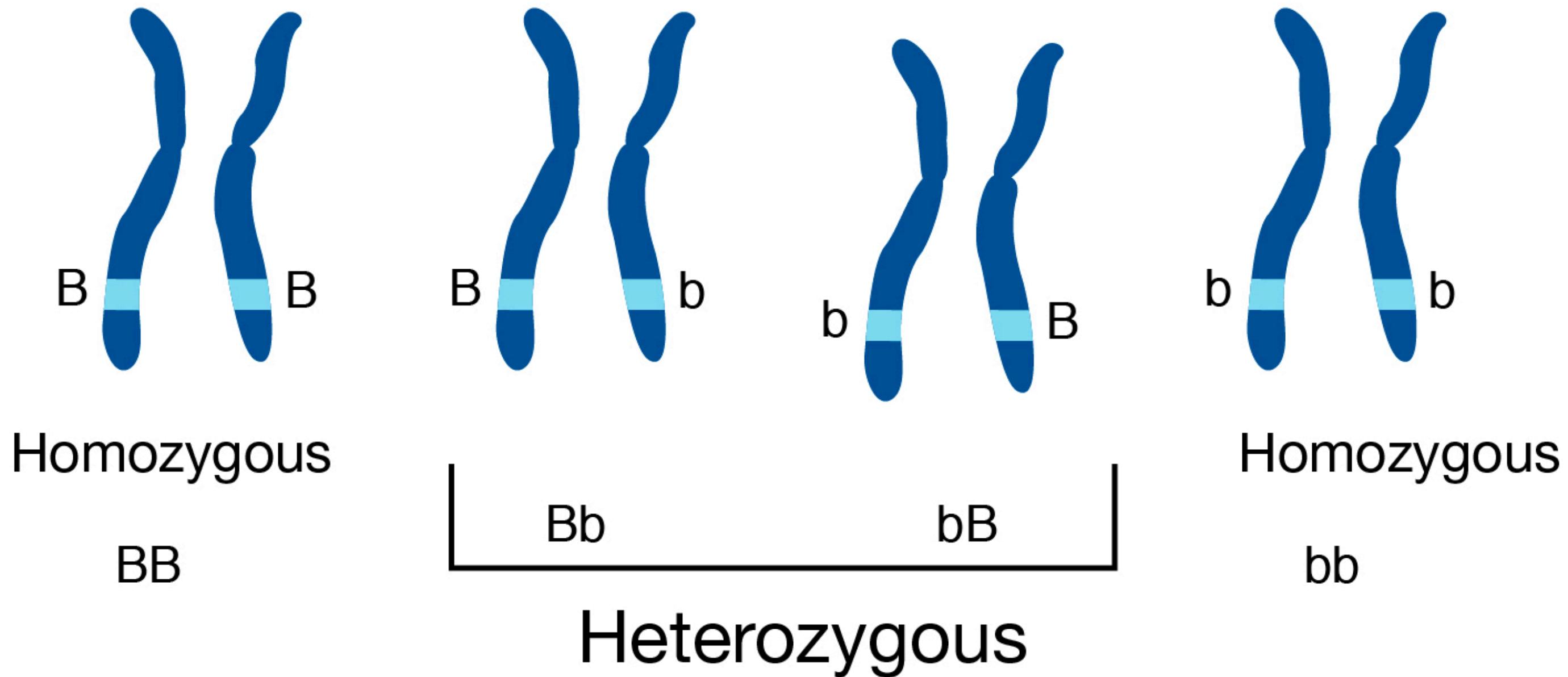
Homozygote vs. heterozygote

One copy from maternal genome
another copy from paternal genome

Homozygous: the maternal == paternal copy

Heterozygous: the maternal != paternal copy

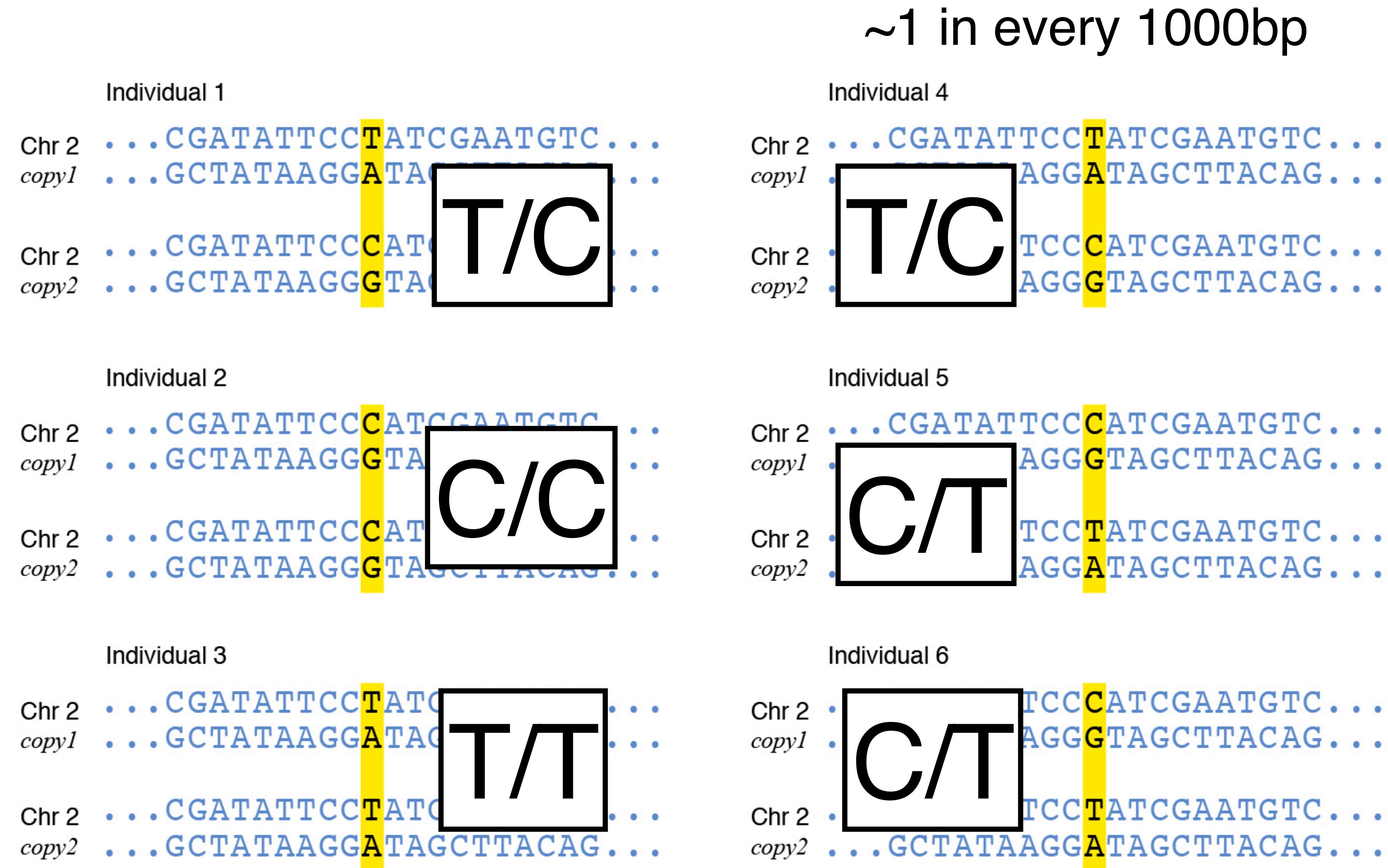
- "Homo" vs. "hetero" \iff the same vs. different
- Zygote (fertilized egg cell)



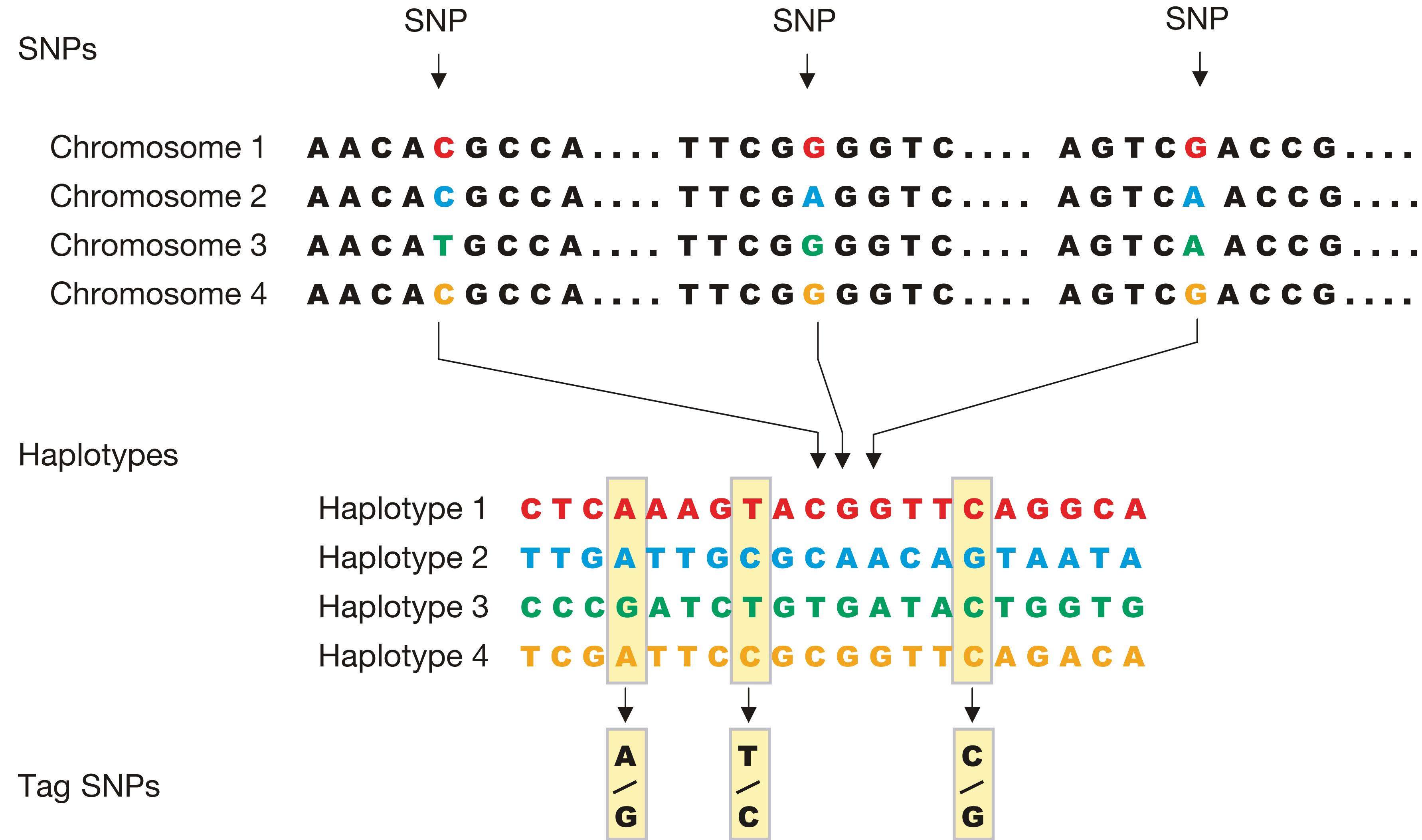
Single Nucleotide Polymorphism

Single nucleotide polymorphism is way too many syllables, so you can understand why we just say "snip". And this is really a simple concept.

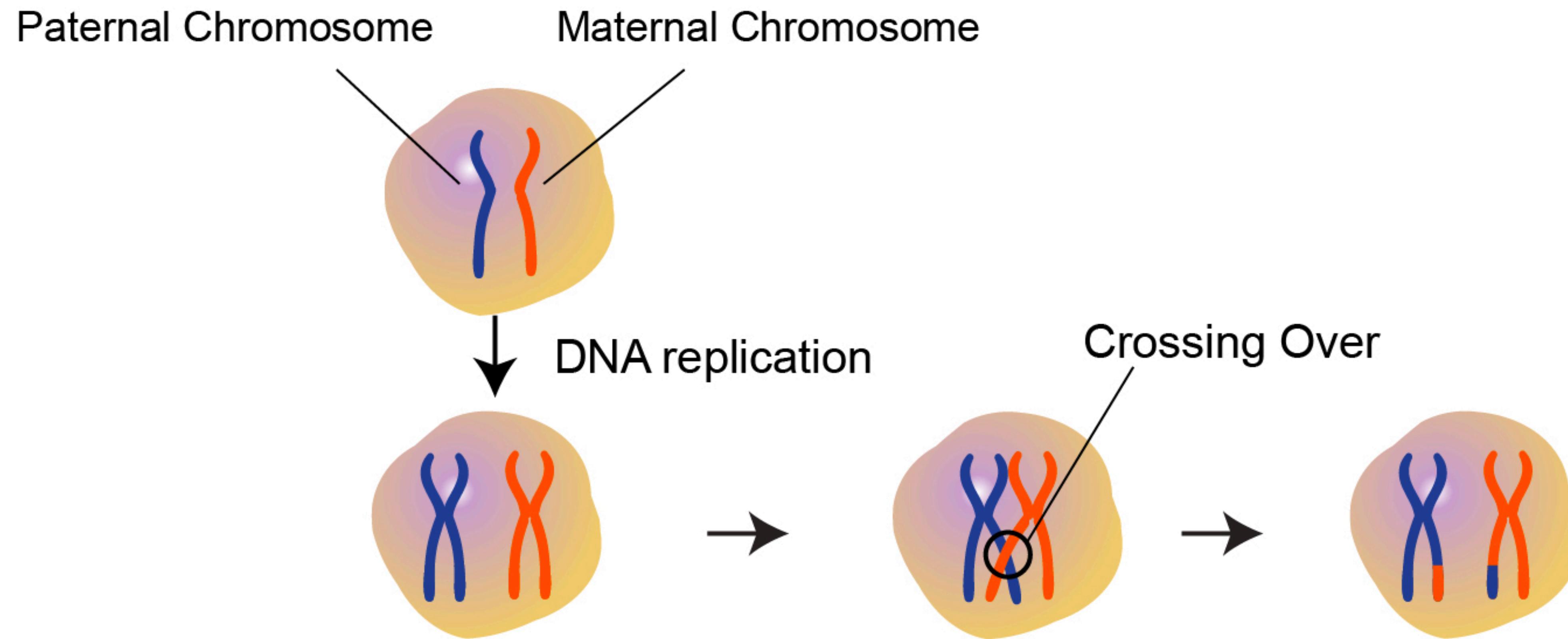
These are the places in the genome where people differ.



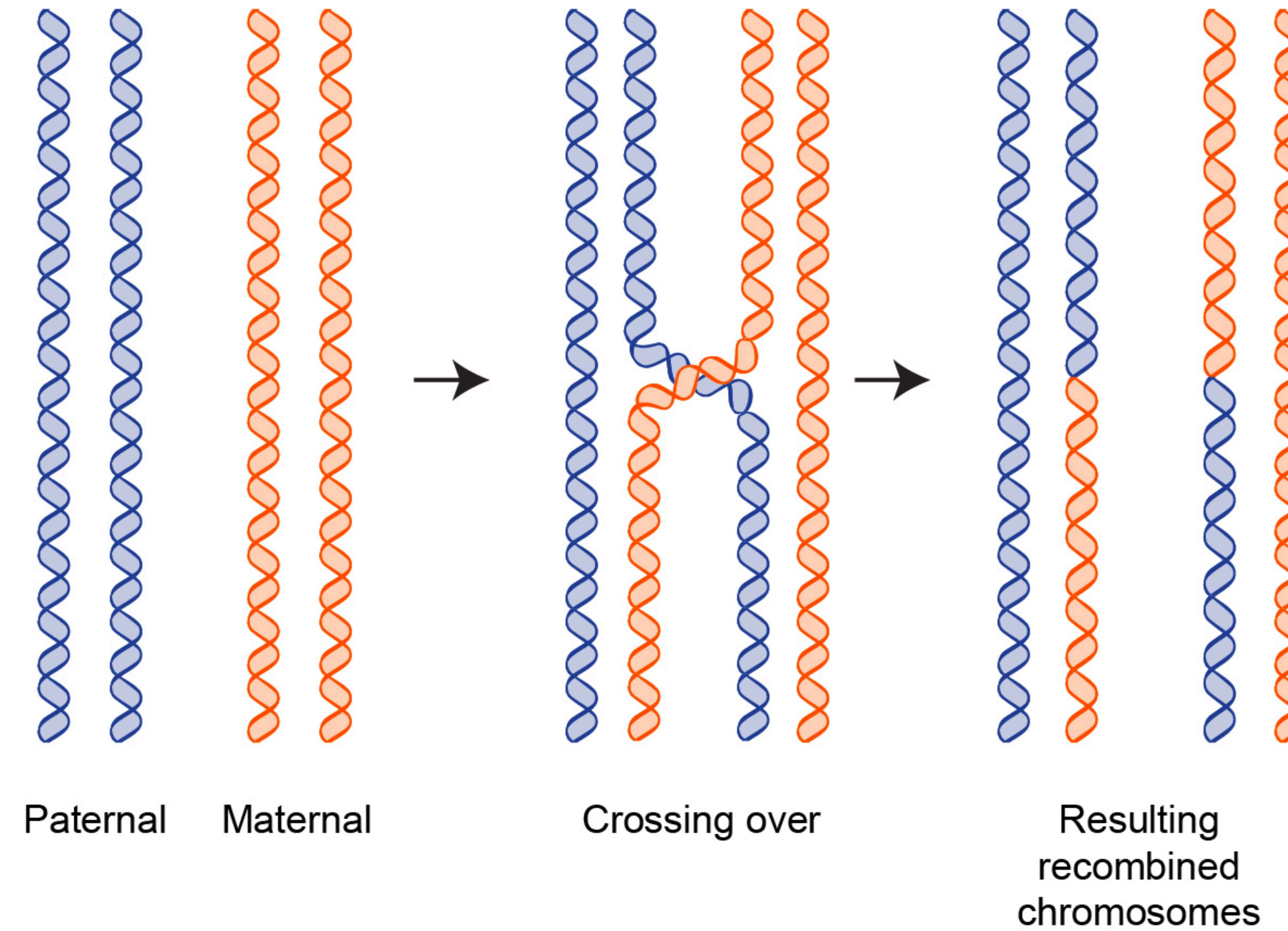
~ 1M common SNPs



Recombination: Mixing the maternal and paternal copies

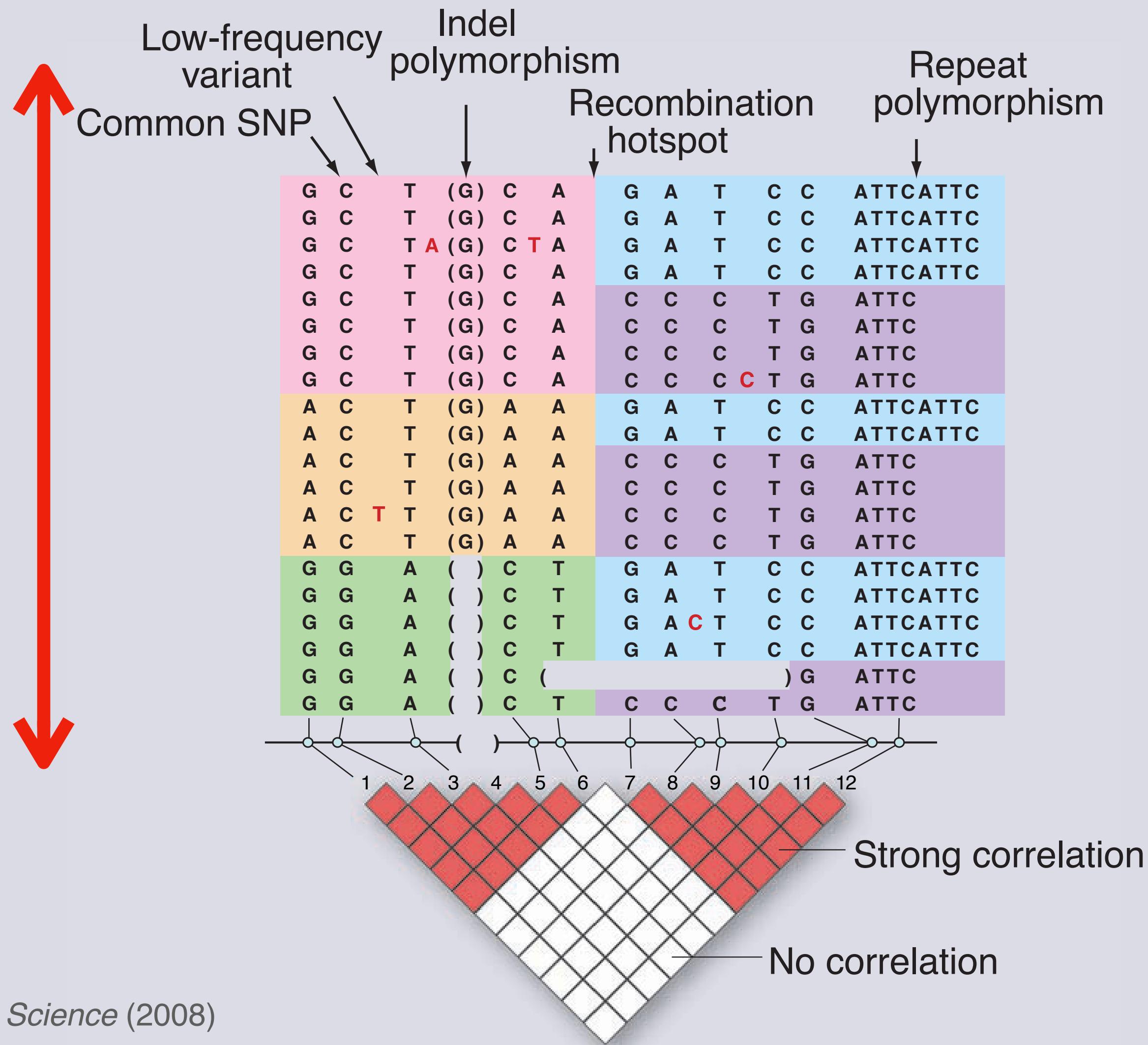


Recombination: Mixing the maternal and paternal copies



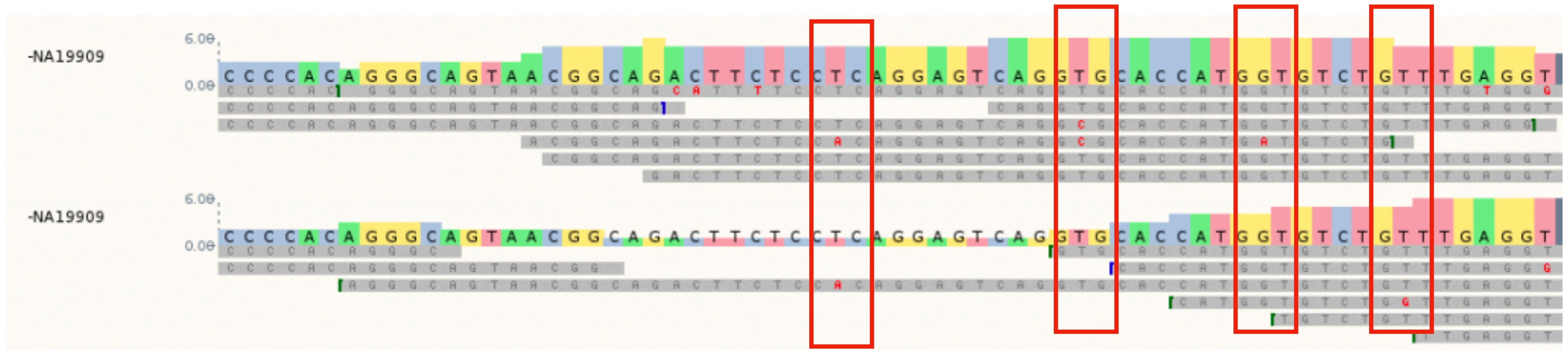
Genetic variation in one figure

across
many
individuals
(diploid
genomes)

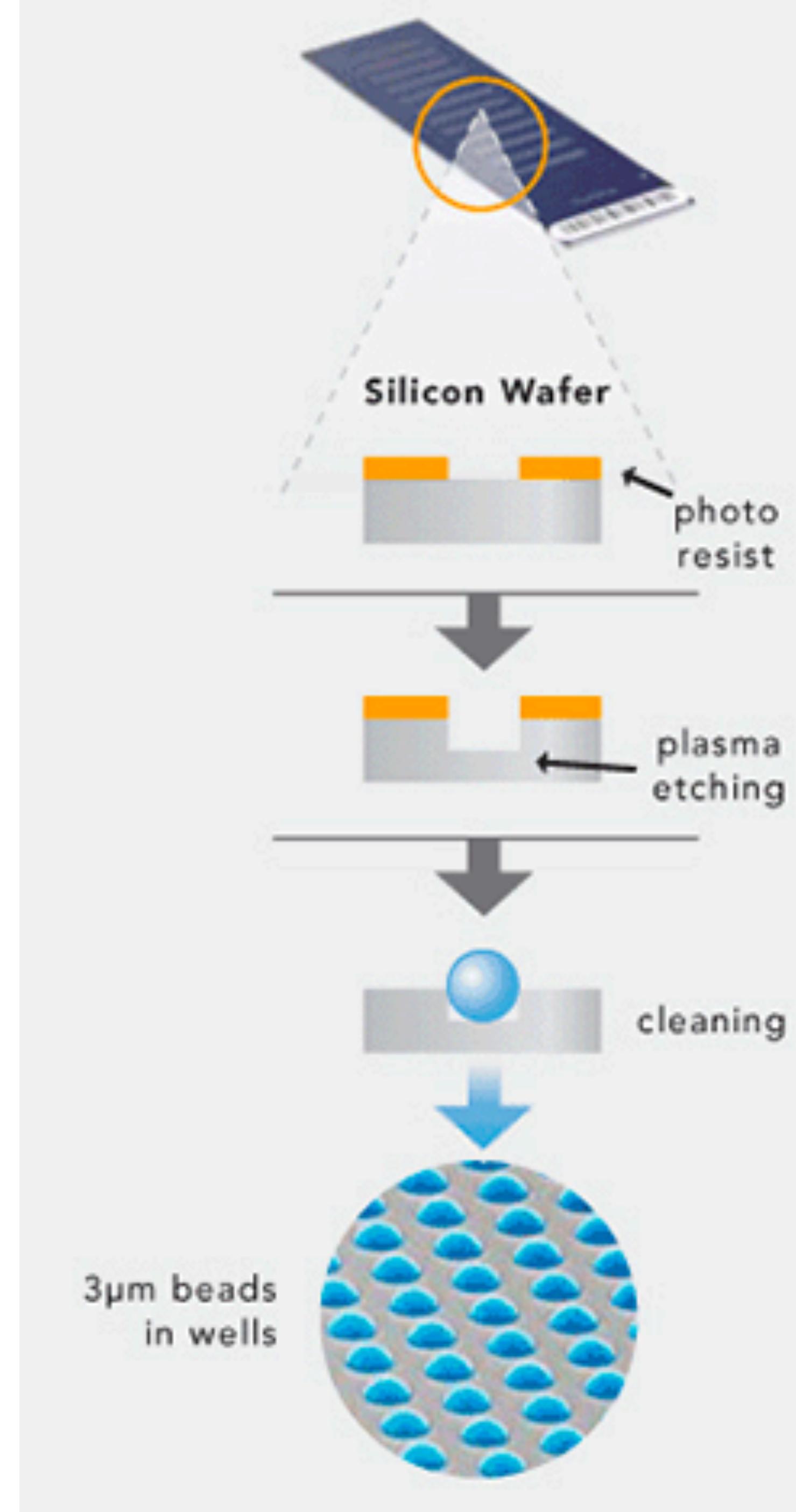
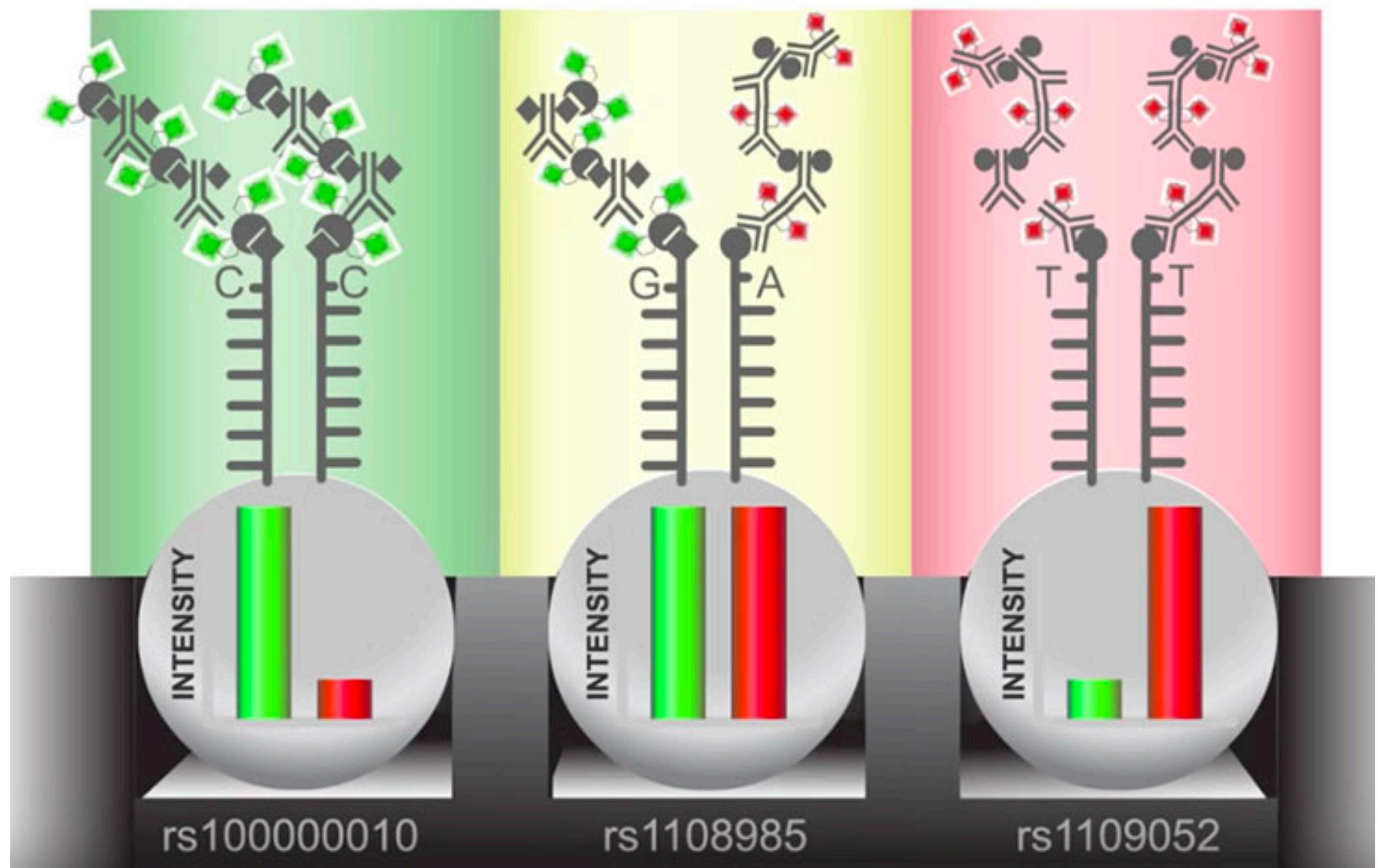


Common SNPs
Insertion/deletion
Other low-freq. variants
Other structural variants
Recombination hotspot

How do we call/quantify variants?



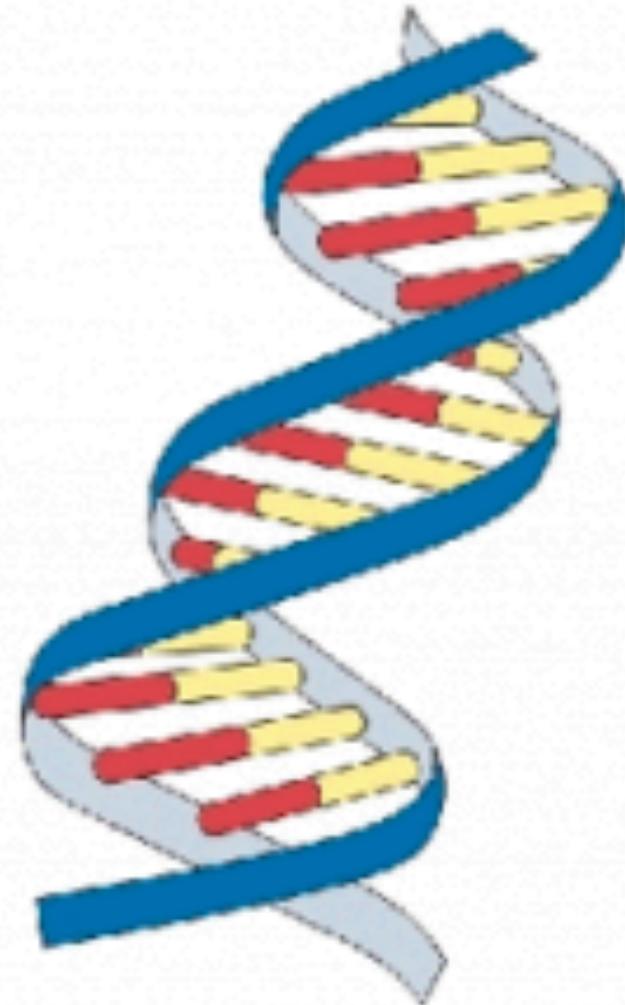
How do we measure variants?



Today's lecture

- Fundamentals in Population Genetics
- Genome-wide association studies
- What are the limitations of GWAS?

Genotype vs. phenotype



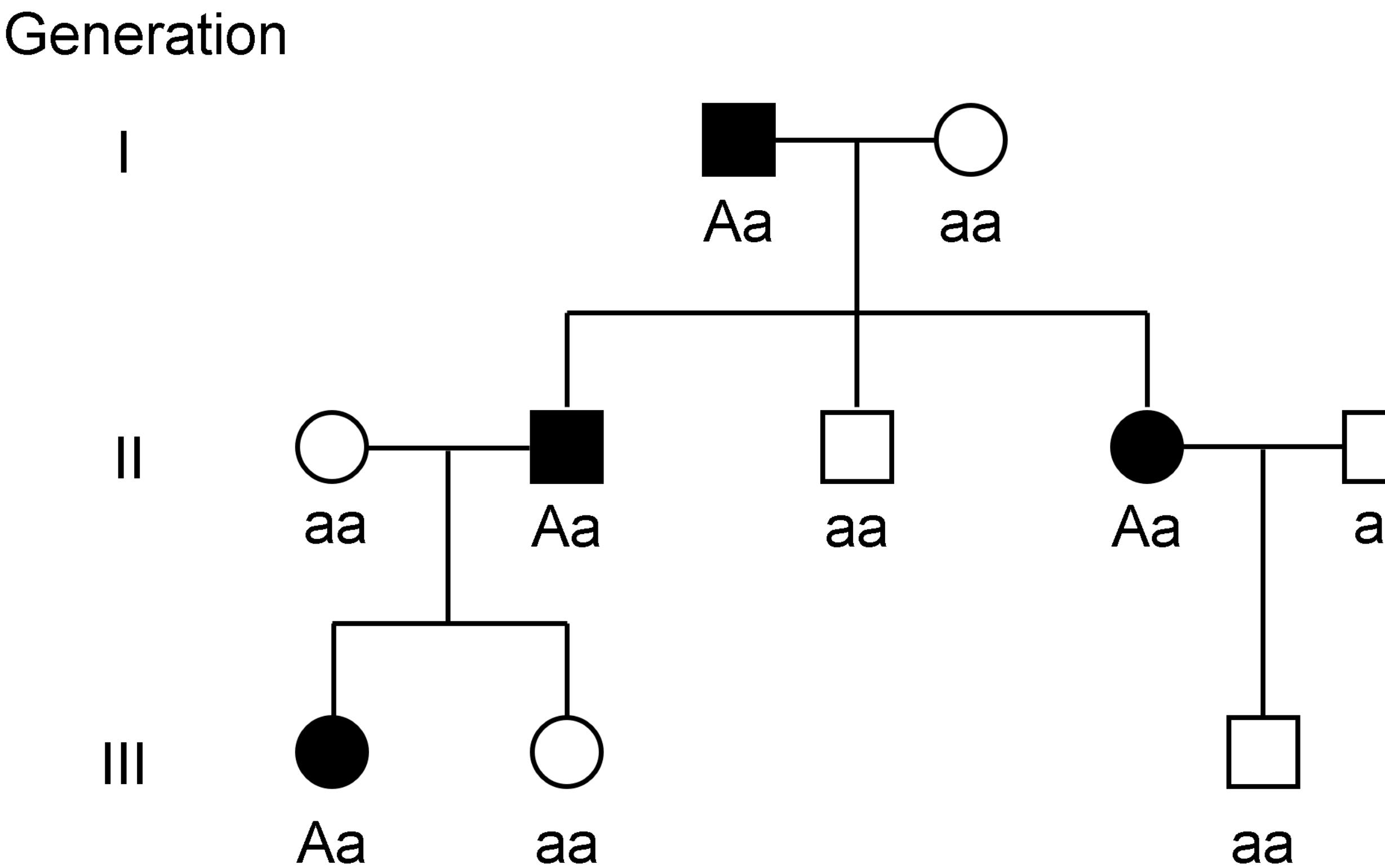
How do we test this?



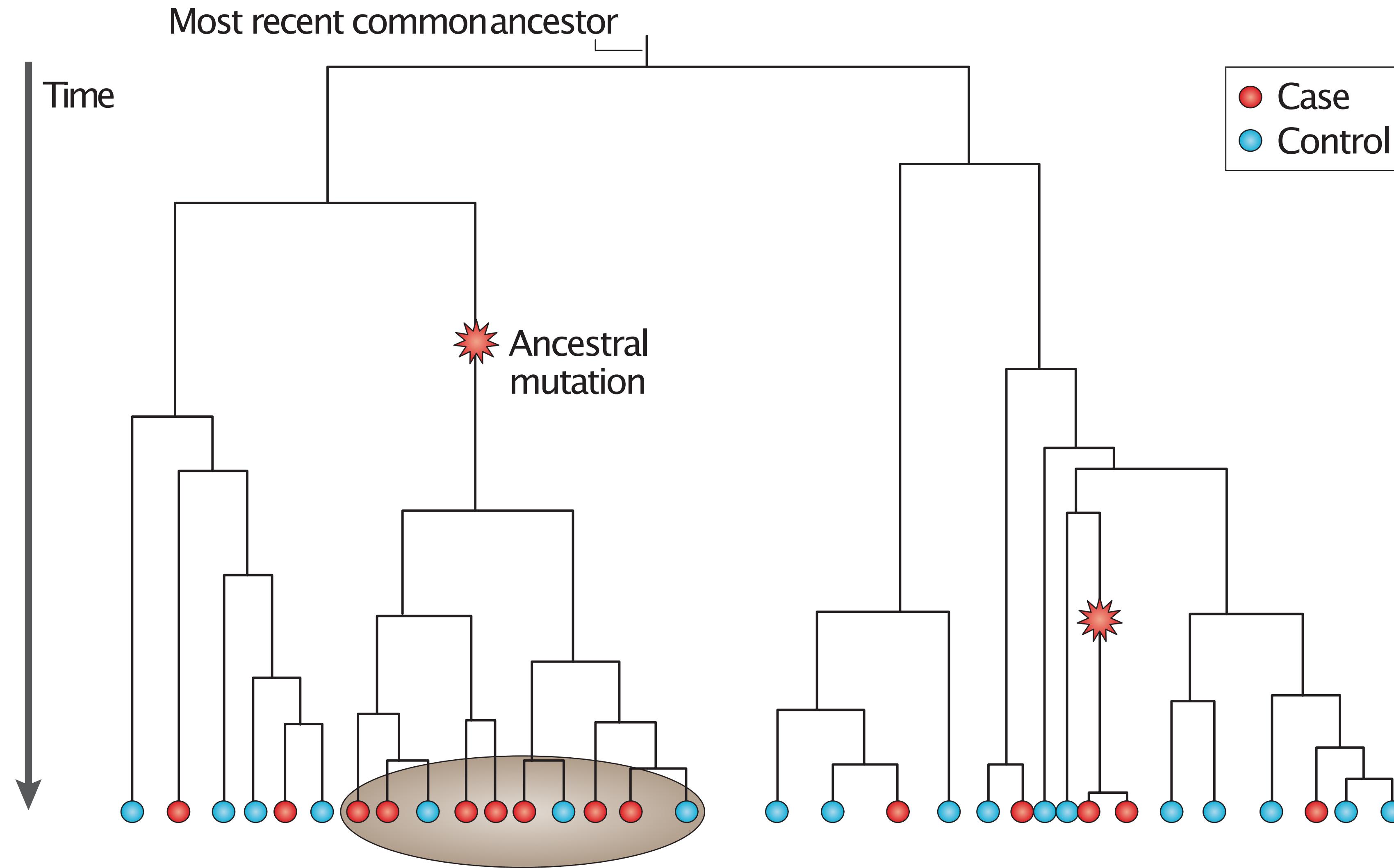
Genotypes are the genetic make-up of an individual.

Phenotypes are the physical traits and characteristics of an individual and are influenced by their genotype and the environment.

How do we associate genetic variants to traits?

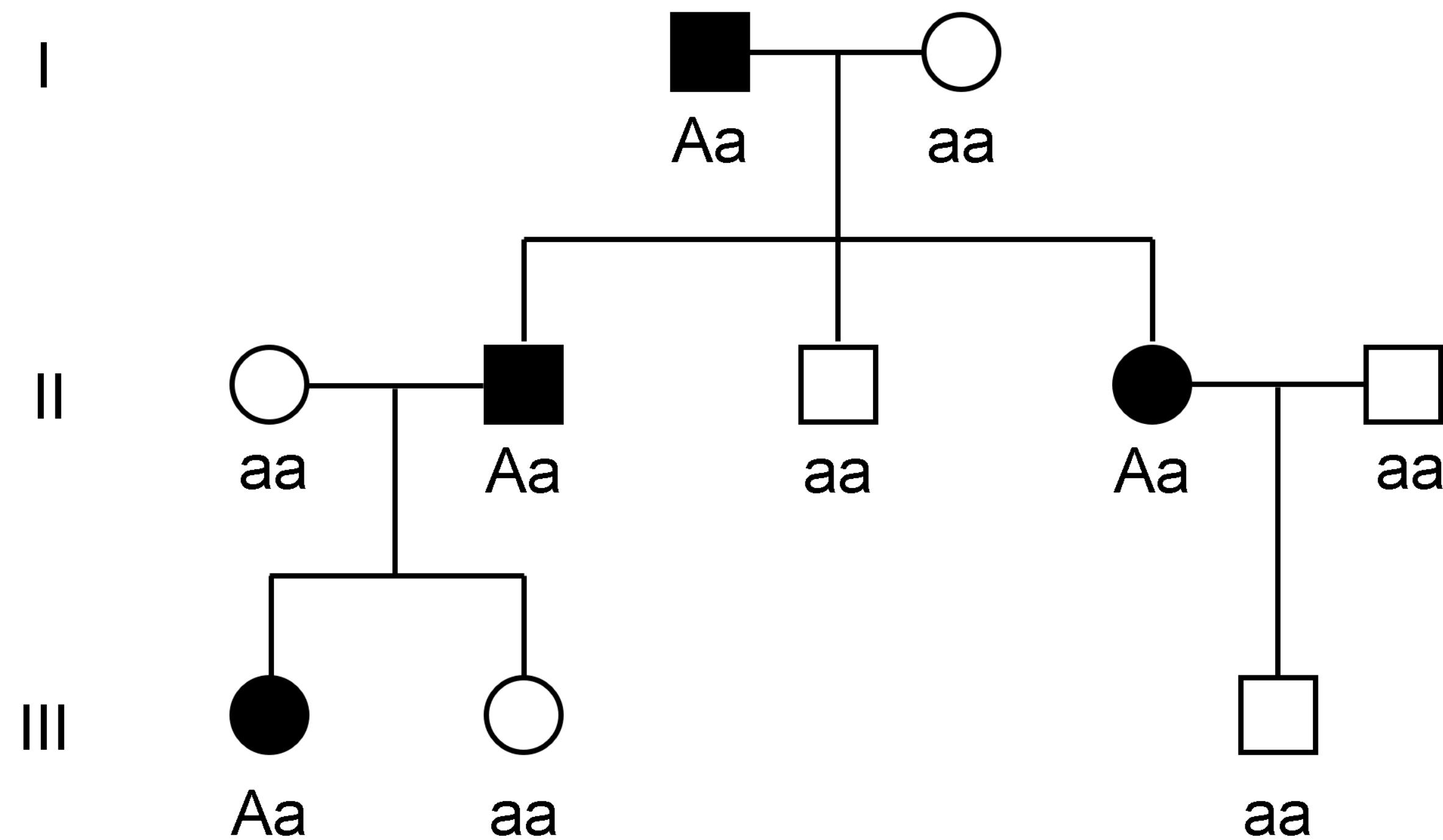


How can a genetic disease occur?



Mendelian disorders ≈ a single disease gene

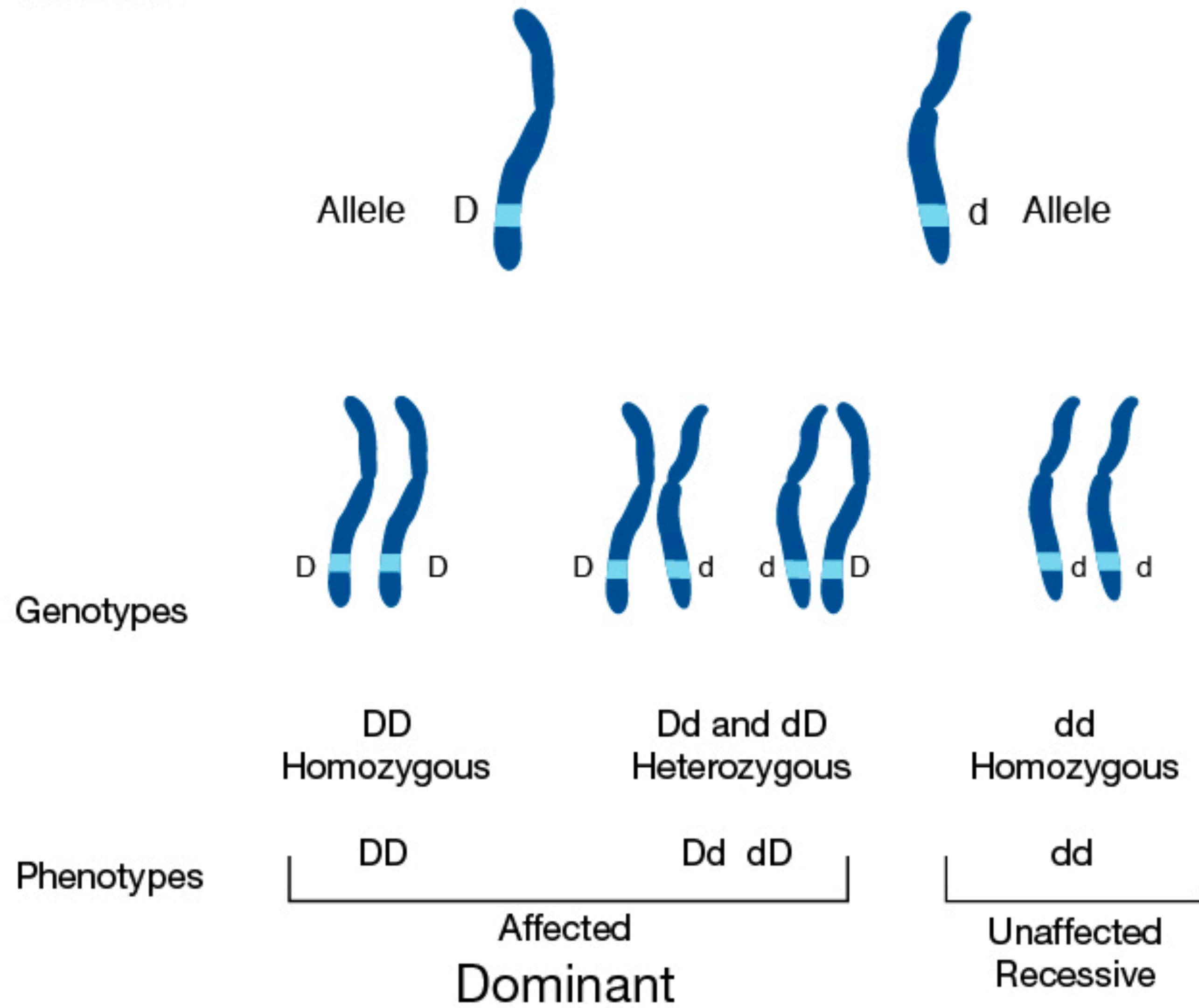
Generation



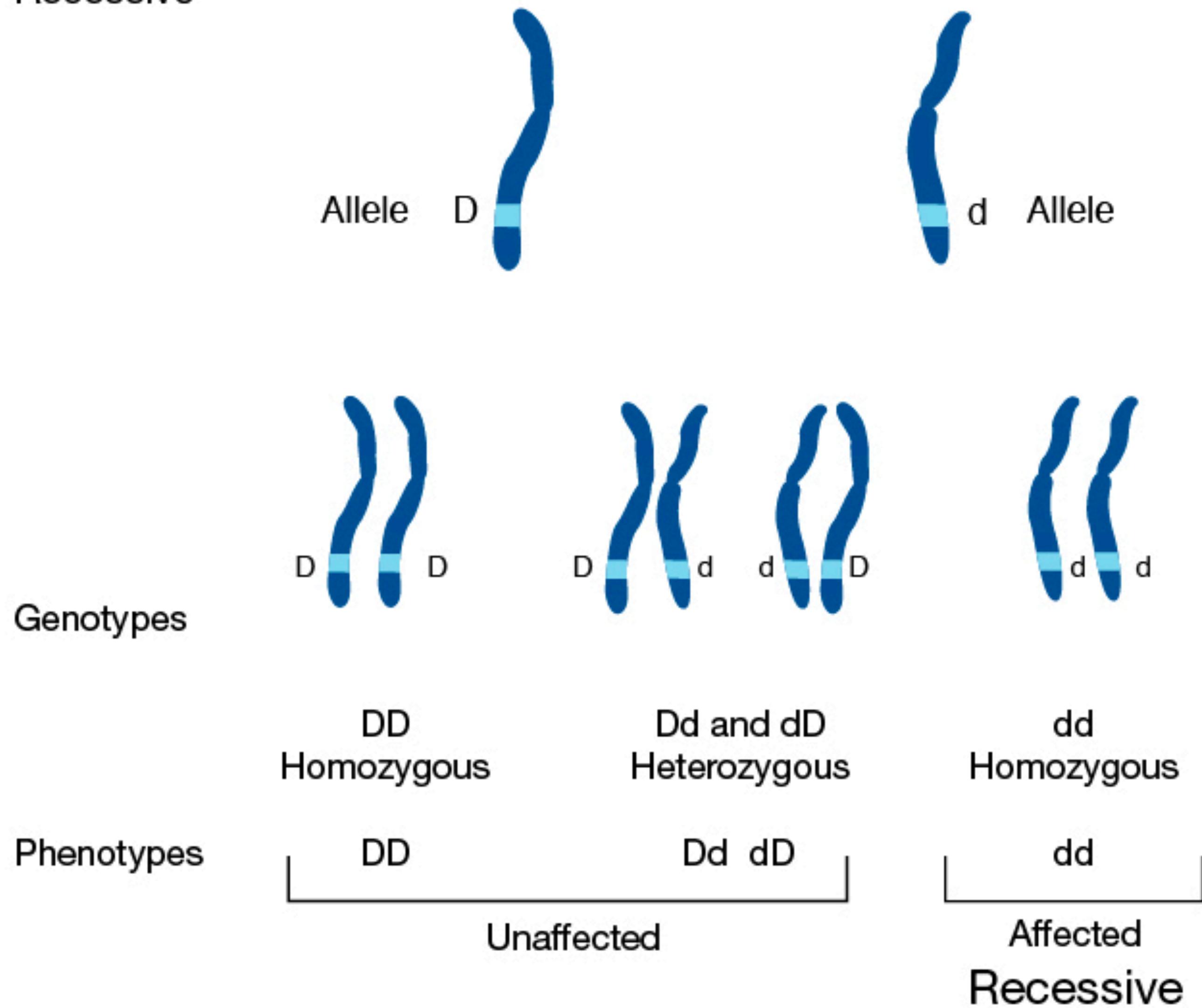
- Cystic Fibrosis
- Sickle-cell anemia
- Phenylketonuria
- Huntington's disease
- ...

Very high penetrance, monogenic

Huntington's Disease Dominant



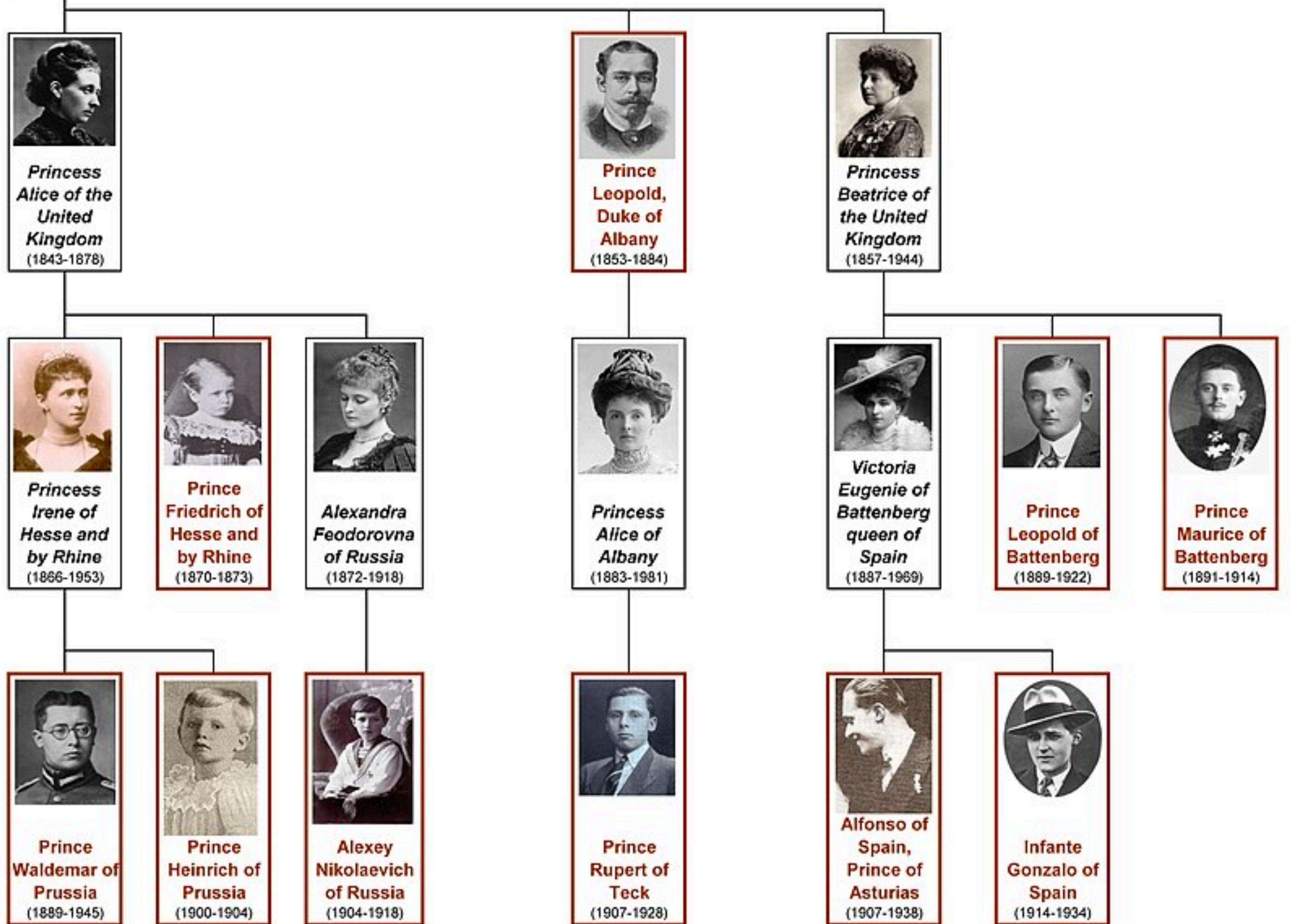
Sickle Cell Anemia or Cystic Fibrosis
Recessive



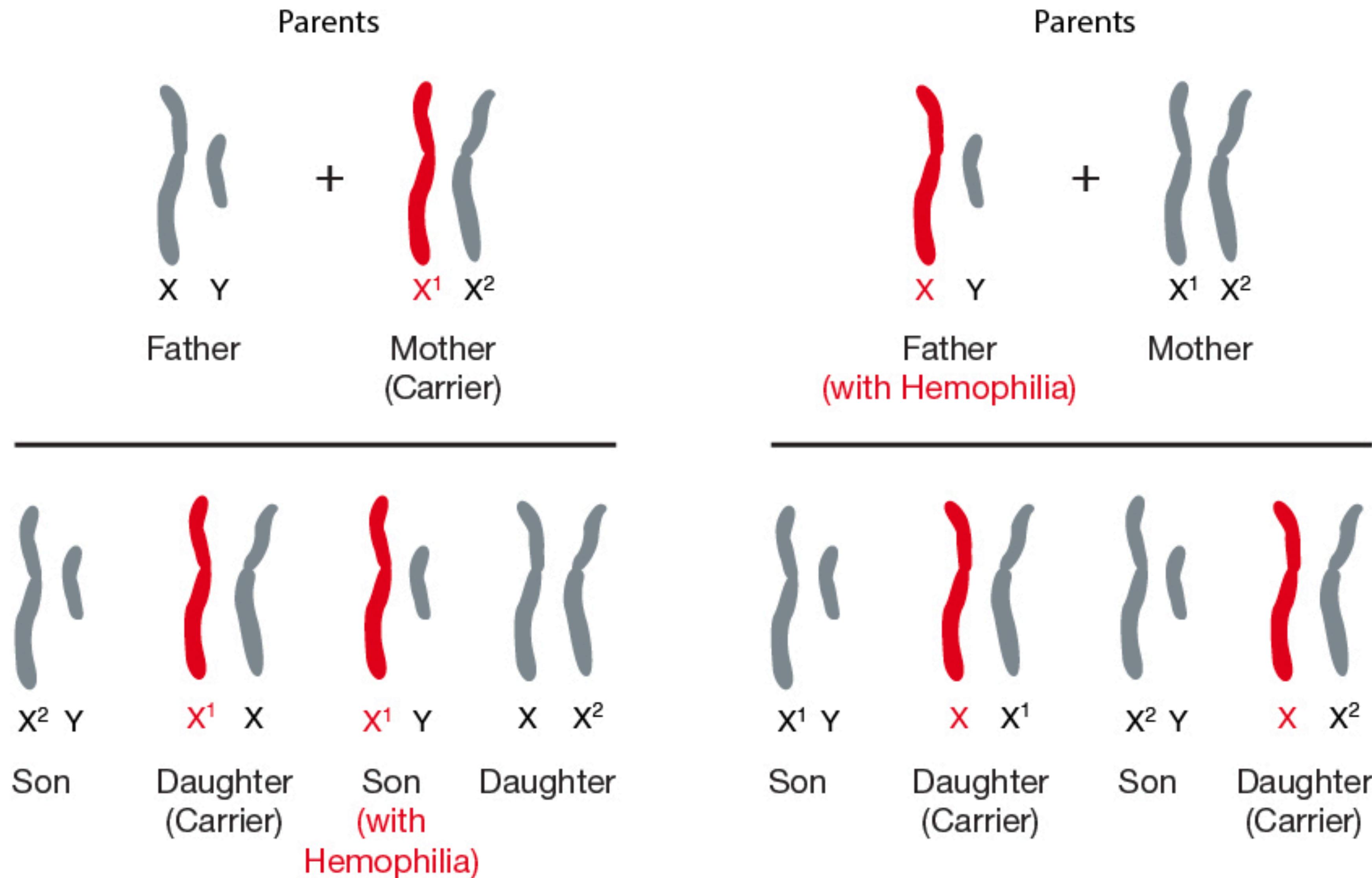


Queen
Victoria
(1819-1901)

Hemophilia in the Queen Victoria's family



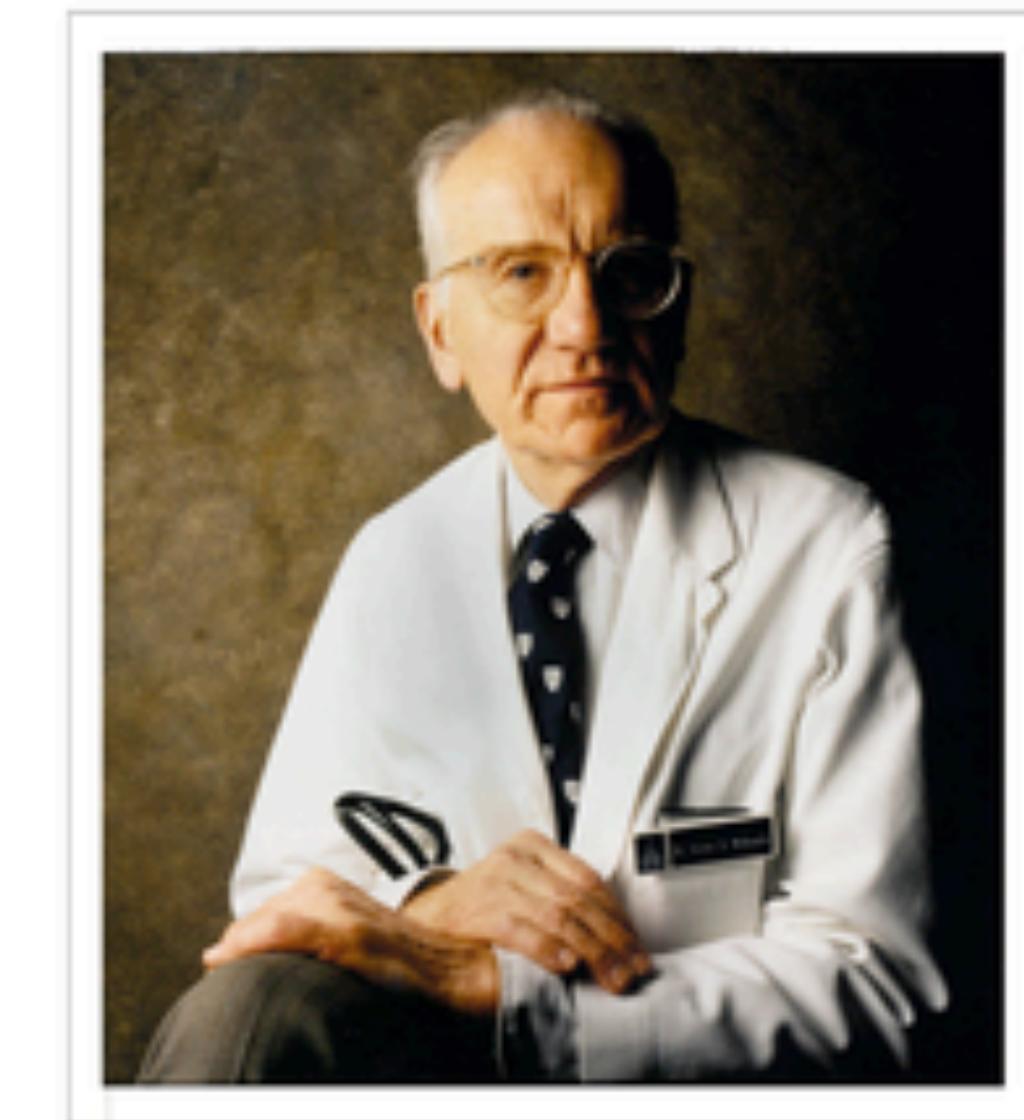
Hemophilia



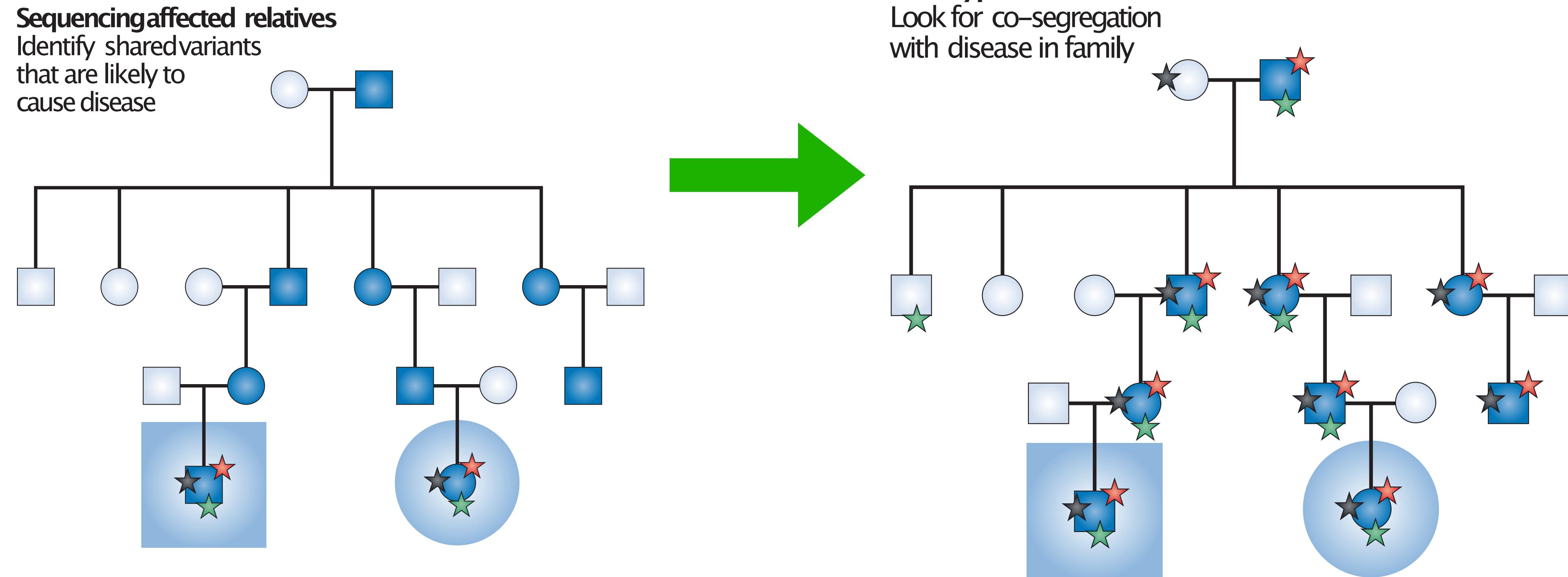
OMIM® – Online Mendelian Inheritance in Man®

Welcome to OMIM®, Online Mendelian Inheritance in Man®. OMIM is a comprehensive, authoritative compendium of human genes and genetic phenotypes that is freely available and updated daily. The full-text, referenced overviews in OMIM contain information on all known mendelian disorders and over 16,000 genes. OMIM focuses on the relationship between phenotype and genotype. It is updated daily, and the entries contain copious links to other genetics resources.

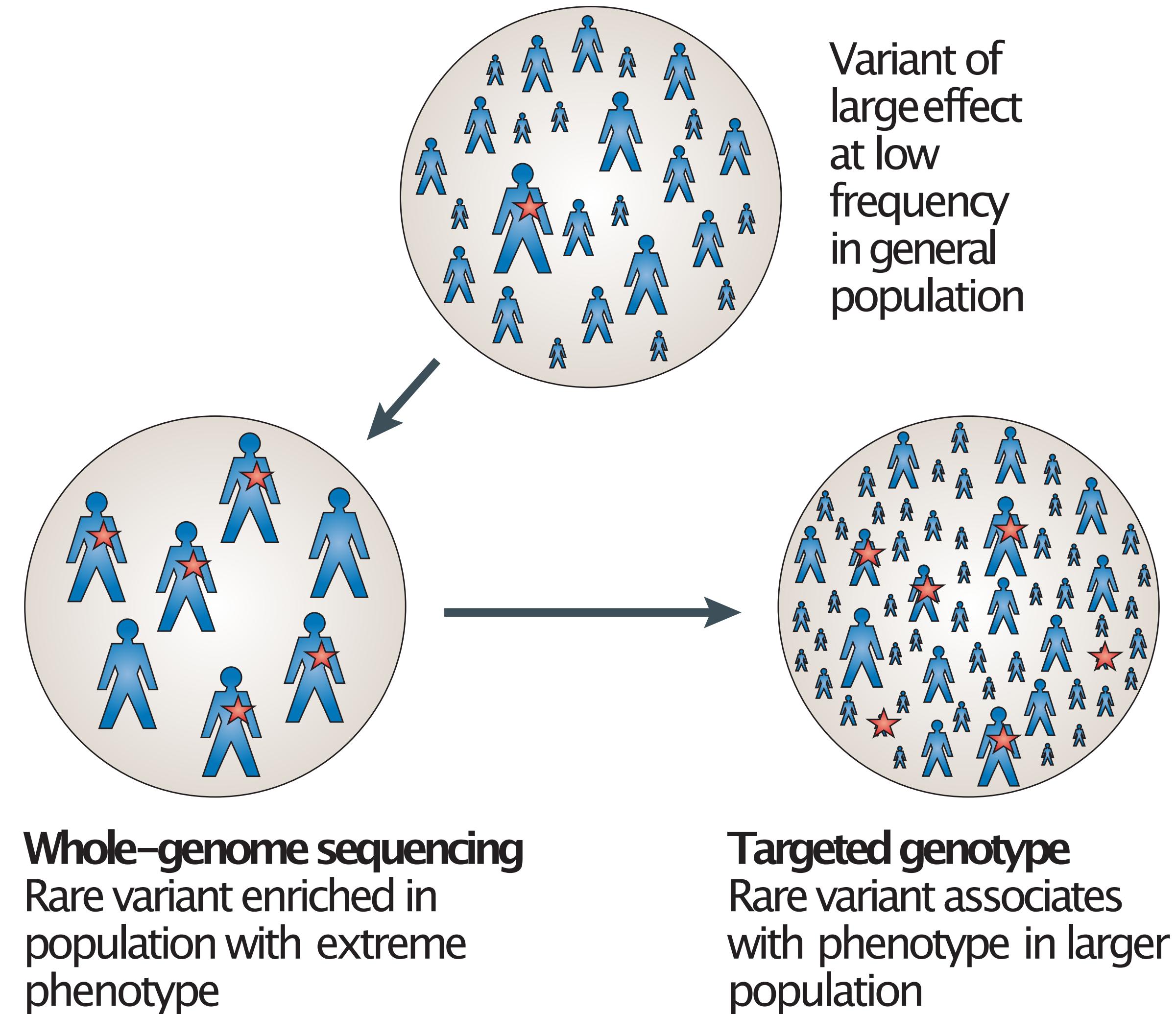
This database was initiated in the early 1960s by Dr. Victor A. McKusick as a catalog of mendelian traits and disorders, entitled Mendelian Inheritance in Man (MIM). Twelve book editions of MIM were published between 1966 and 1998. The online version, OMIM, was created in 1985 by a collaboration between the National Library of Medicine and the William H. Welch Medical Library at Johns Hopkins. It was made generally available on the internet starting in 1987. In 1995, OMIM was developed for the World Wide Web by NCBI, the National Center for Biotechnology Information.



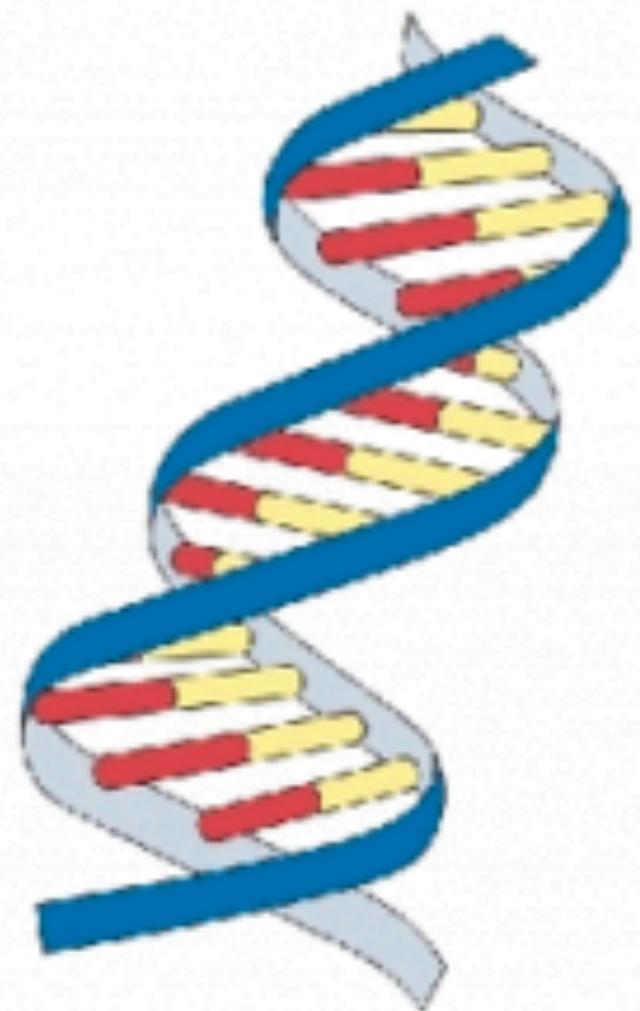
Targetted genotyping and sequencing



Population-level enrichment?



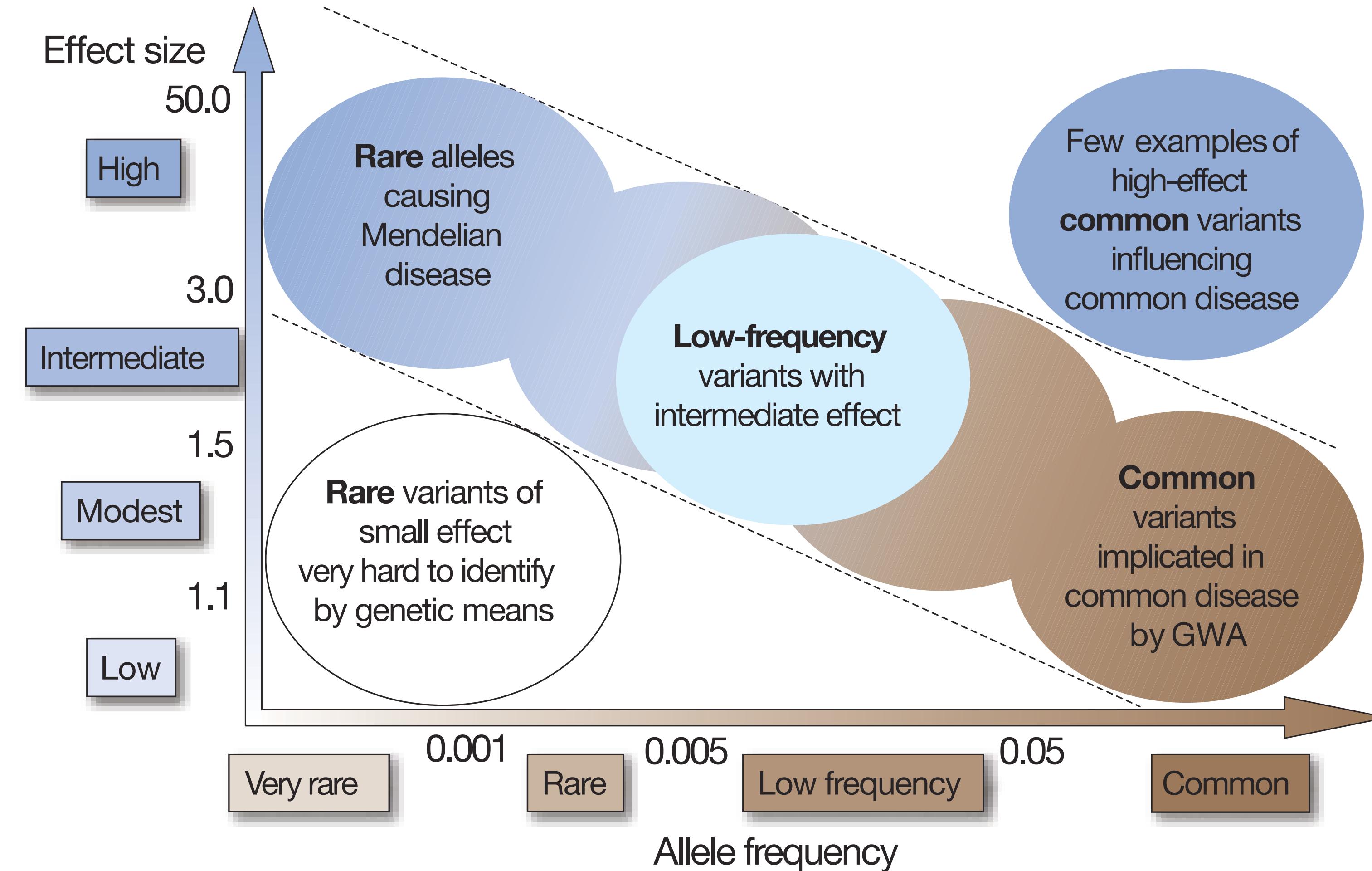
Can we identify variants associated with traits?



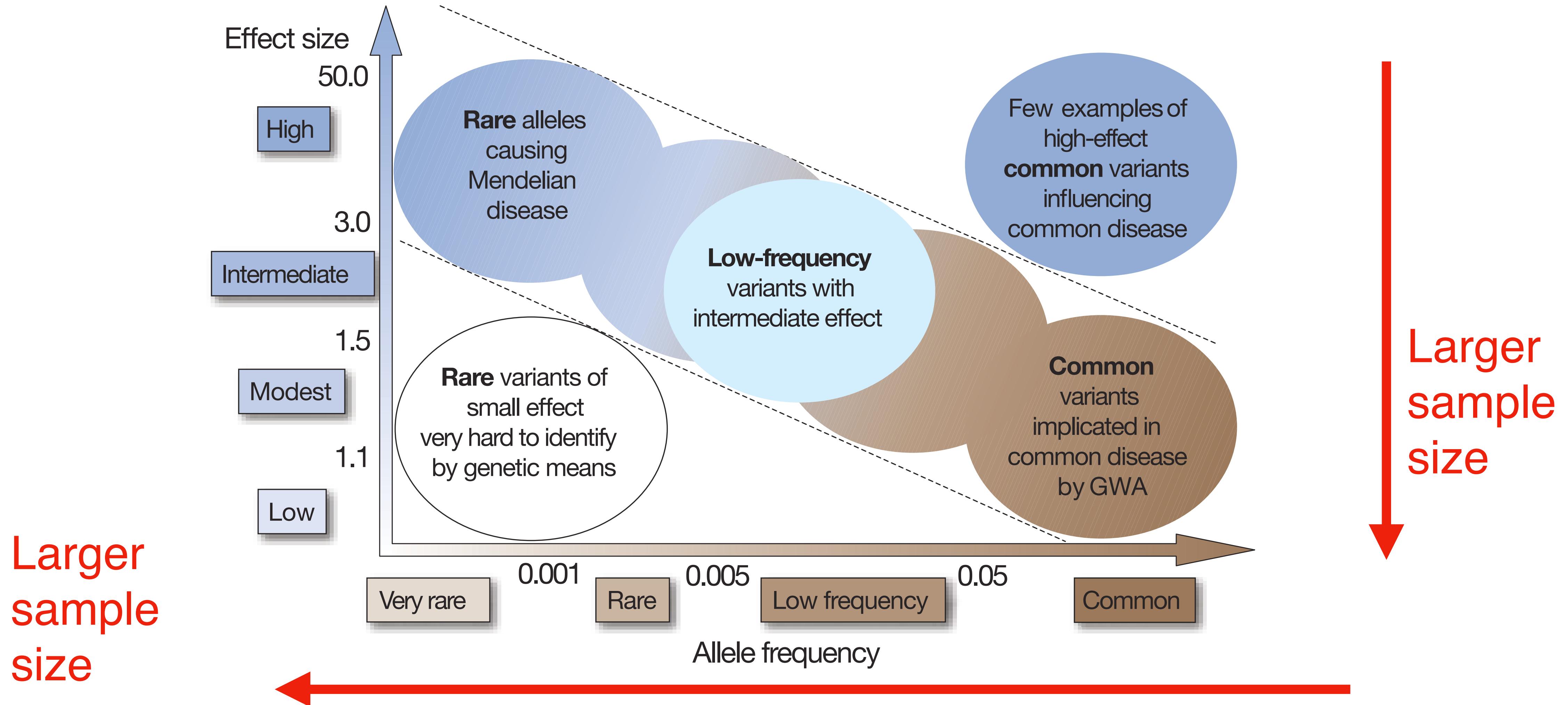
Genotypes are the genetic make-up of an individual.

Phenotypes are the physical traits and characteristics of an individual and are influenced by their genotype and the environment.

Will trait-associated variants emerge?



Will population-level genetic studies be of worth?



We need a large sample size!

Sample size matters (e.g., Schizophrenia GWAS)

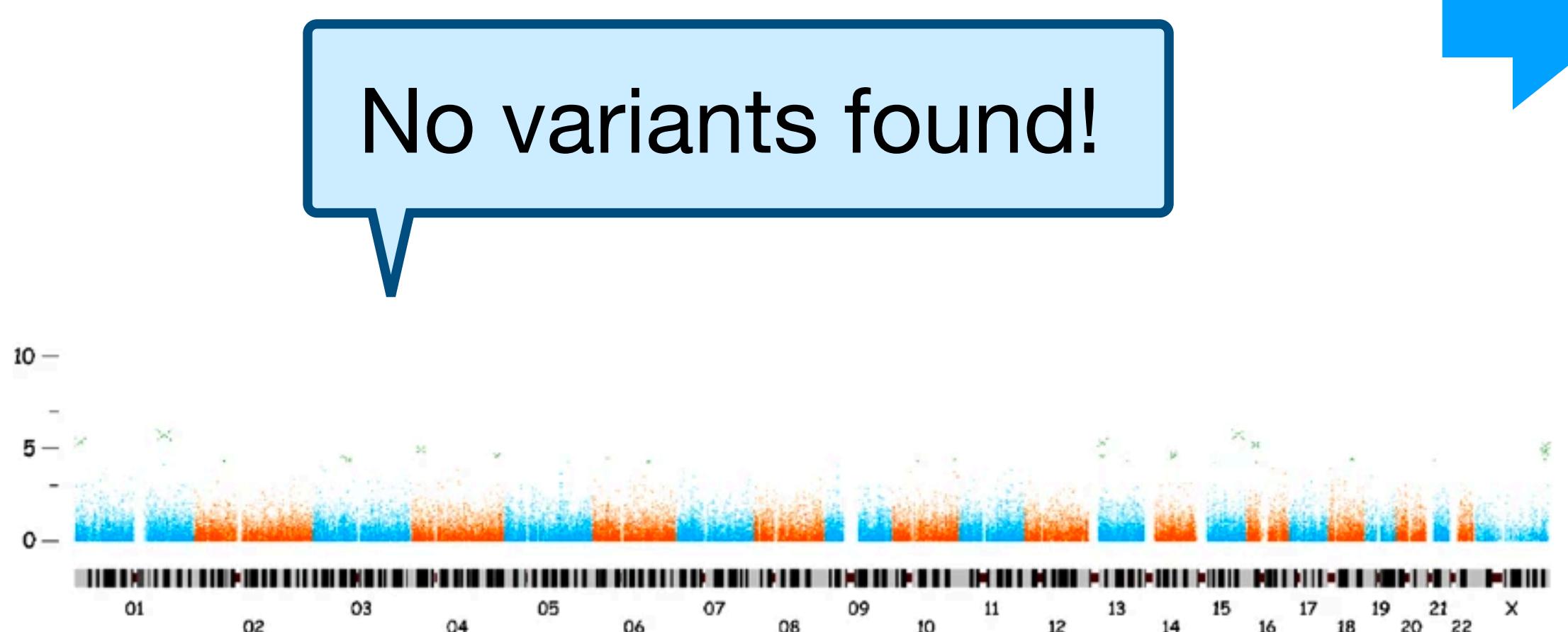
IMMEDIATE COMMUNICATION

Genomewide association for schizophrenia in the CATIE study: results of stage 1

PF Sullivan^{1,2}, D Lin³, J-Y Tzeng⁴, E van den Oord⁵, D Perkins⁶, TS Stroup⁶, M Wagner⁷, S Lee³, FA Wright³, F Zou³, W Liu⁸, AM Downing⁹, J Lieberman¹⁰ and SL Close⁹

N=733 cases vs. 733 controls

No variants found!

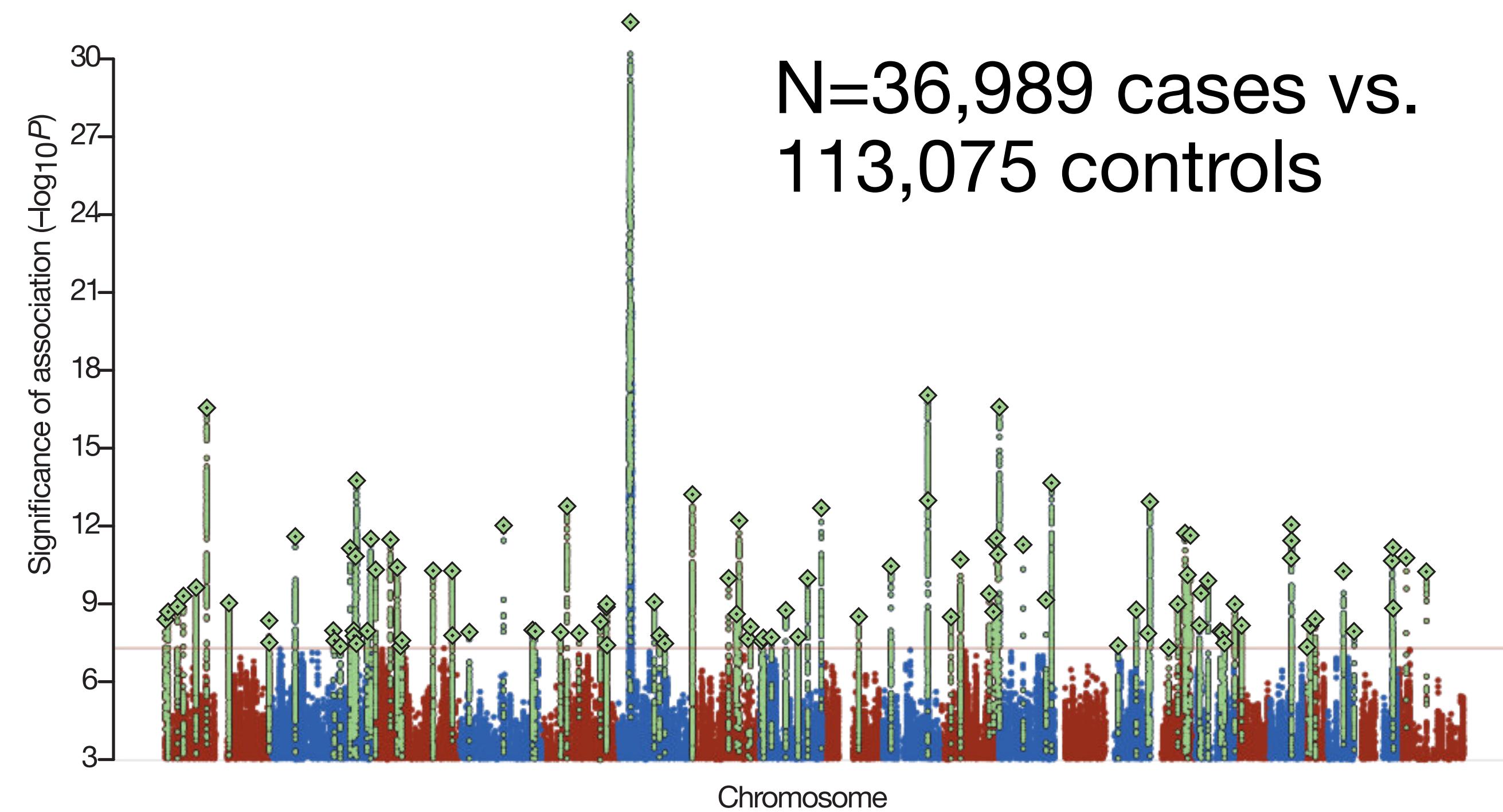


Sullivan *et al.* (2008)

Biological insights from 108 schizophrenia-associated genetic loci

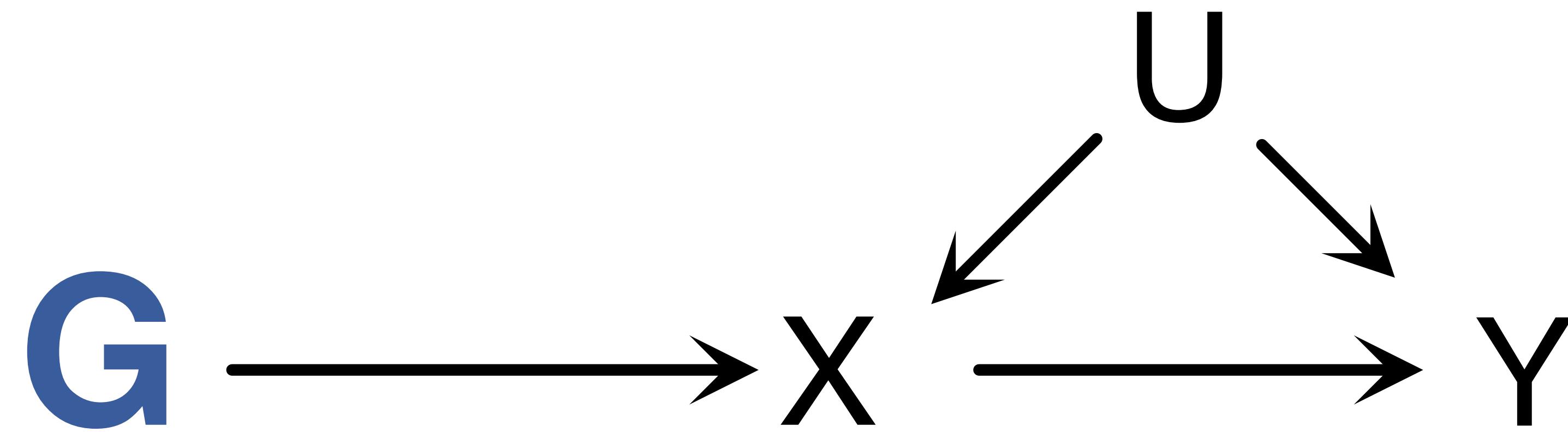
Schizophrenia Working Group of the Psychiatric Genomics Consortium*

N=36,989 cases vs.
113,075 controls



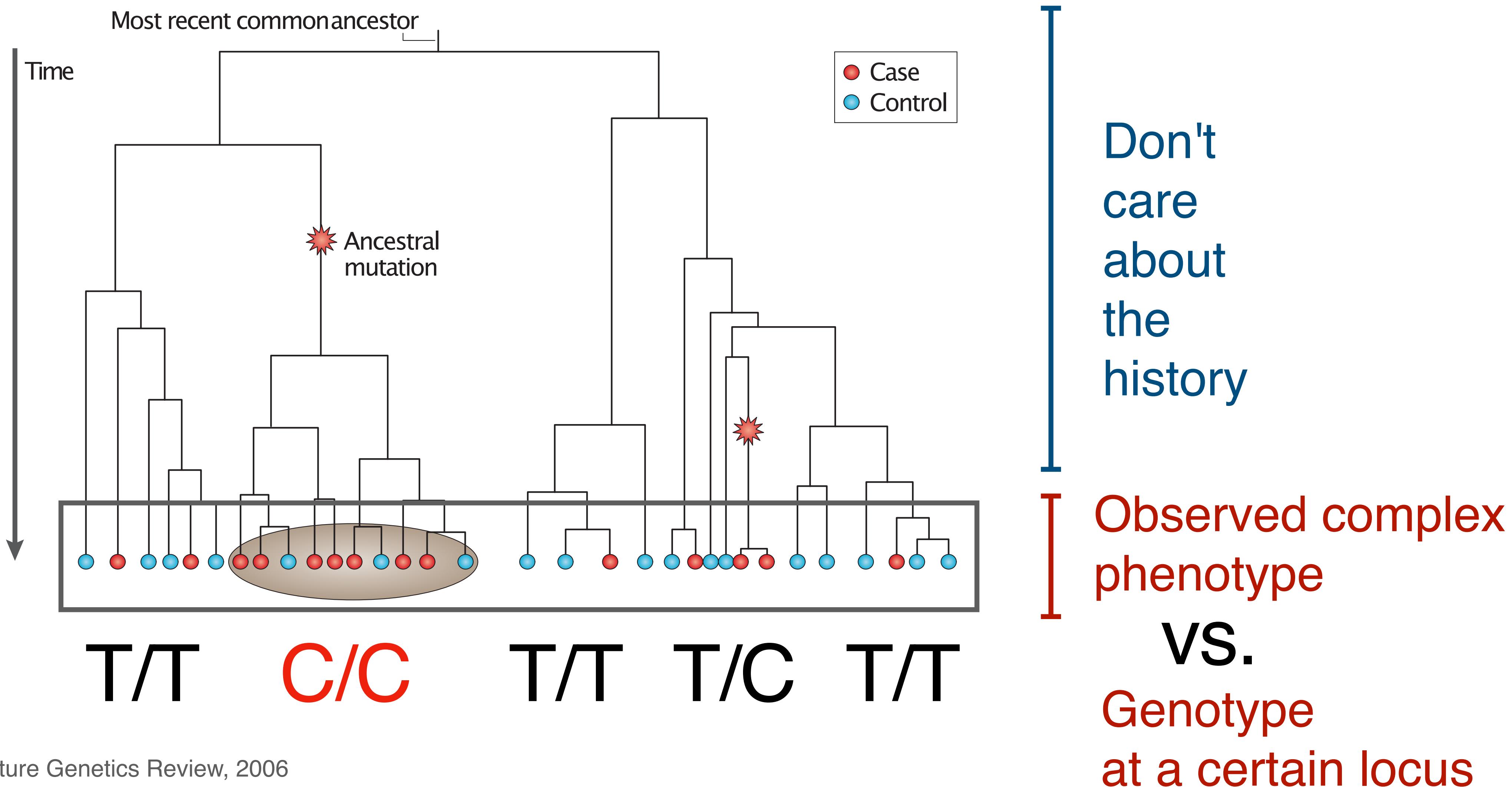
Schizophrenia Working Group of PGC, *Nature* (2014)

Genetic associations can lead to other discoveries



E.g., Mendelian Randomization (in the previous lecture)

Genetic association tests compare genotype vs. phenotype variation across individuals



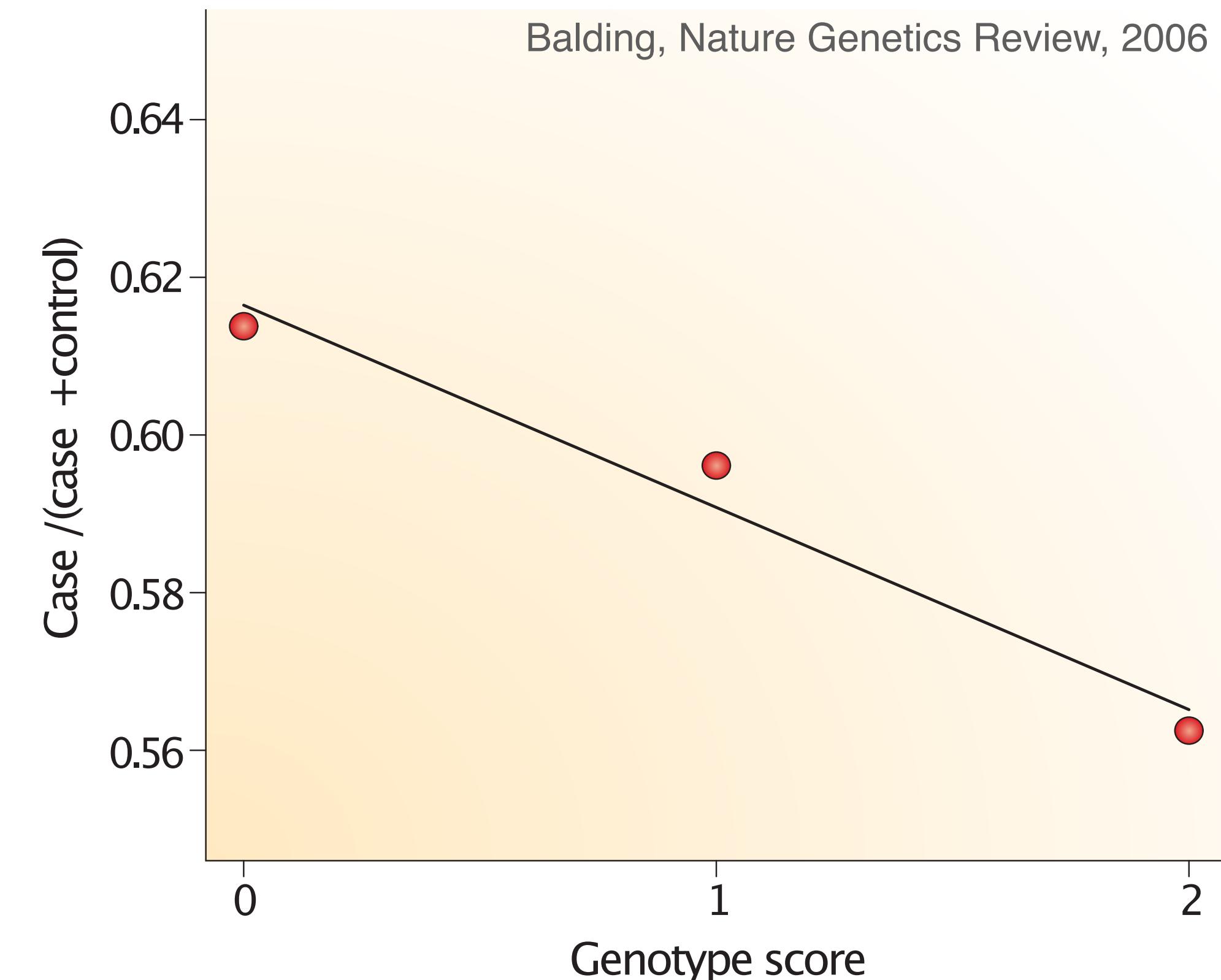
A typical genetics study design

Resistant
to a disease

More
susceptible
to a disease
(e.g., COVID)

	T/T	0
	C/C	2
	T/T	0
	T/C	1
	T/T	0

T = a major allele
C = a minor allele



A typical genetics study design

Resistant
to a disease

More
susceptible
to a disease
(e.g., COVID)

		T/T	0
		C/C	2
		T/T	0
		T/C	1
		T/T	0

Genotype (dosage) of a variant j

$$X_{ij} \in \{0,1,2\}$$

$$Y_i \in \{0,1\}$$

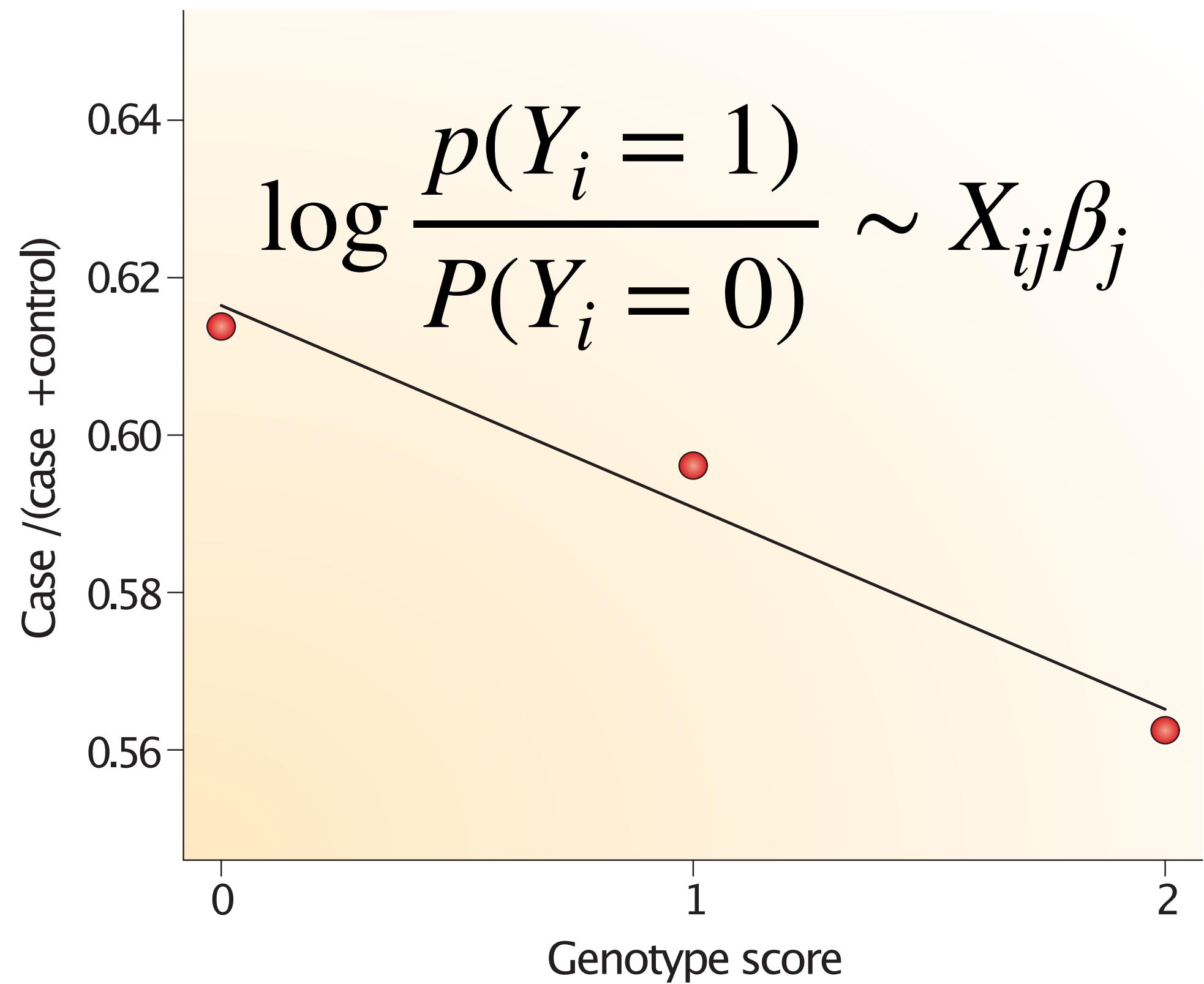
Genotype (dosage) of a variant j

$$\log \frac{p(Y_i = 1)}{P(Y_i = 0)} \sim X_{ij}\beta_j$$

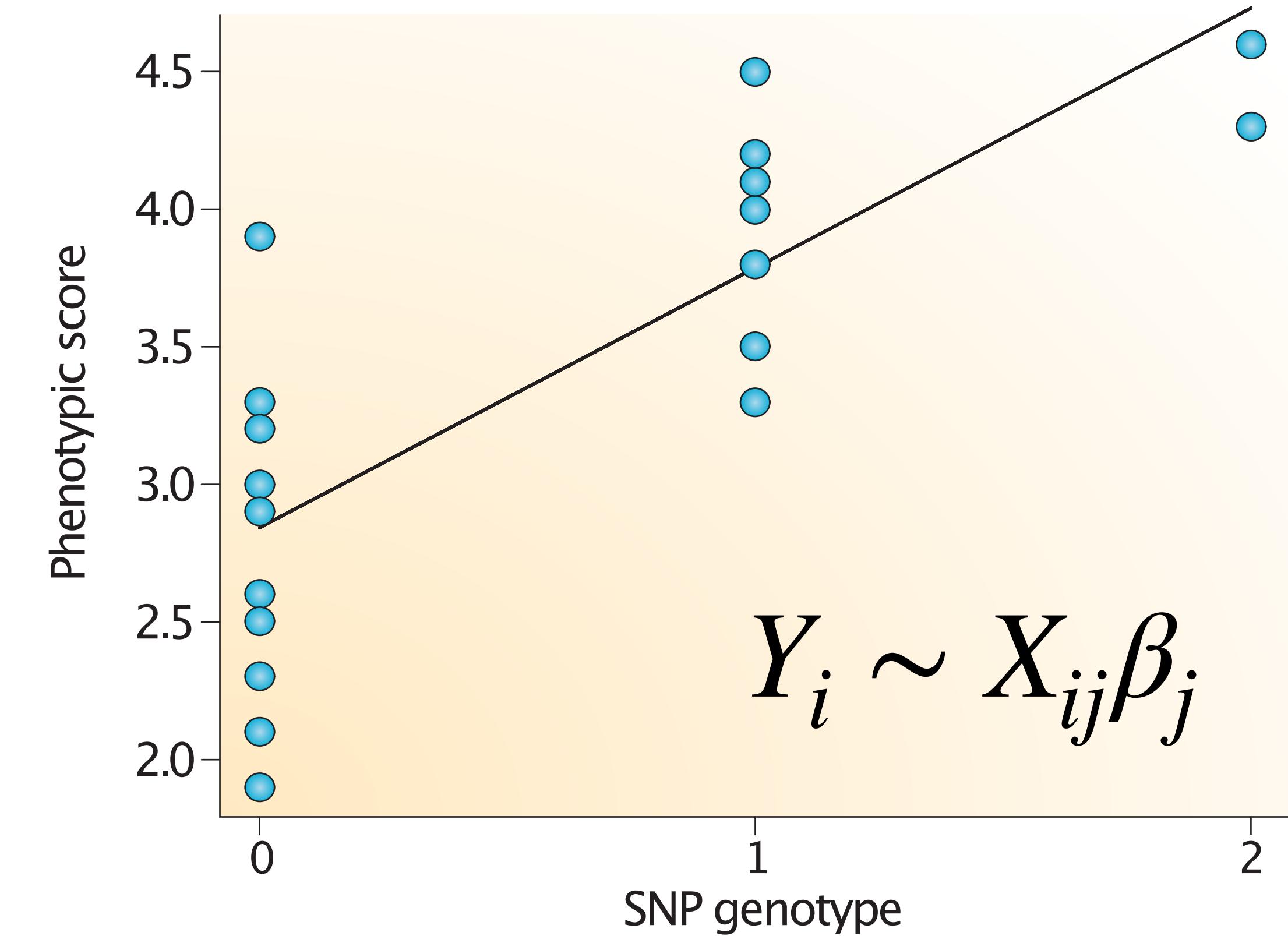
T = a major allele C = a minor allele

A traditional genetics study design

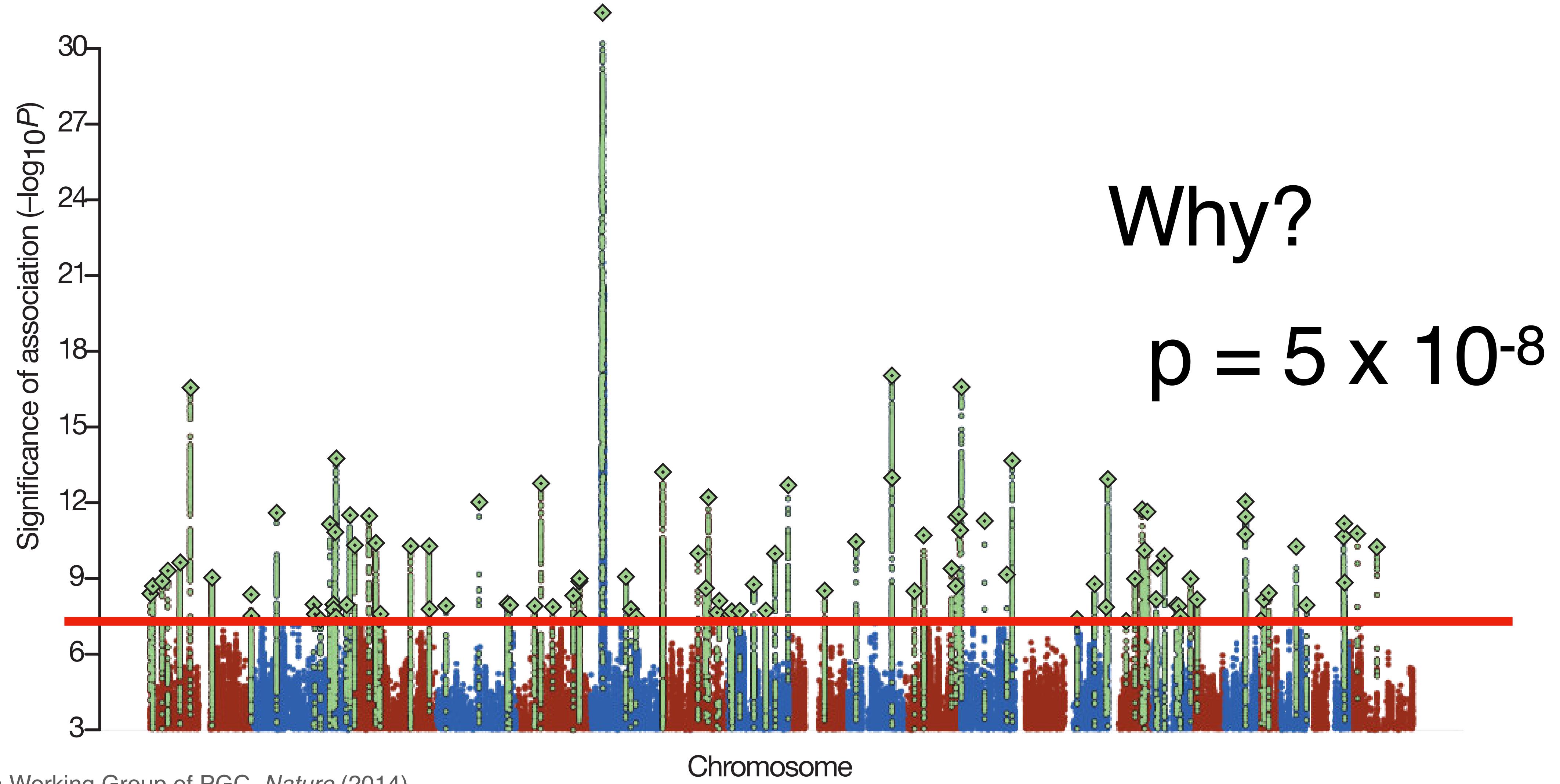
Case-control GWAS



Quantitative trait GWAS

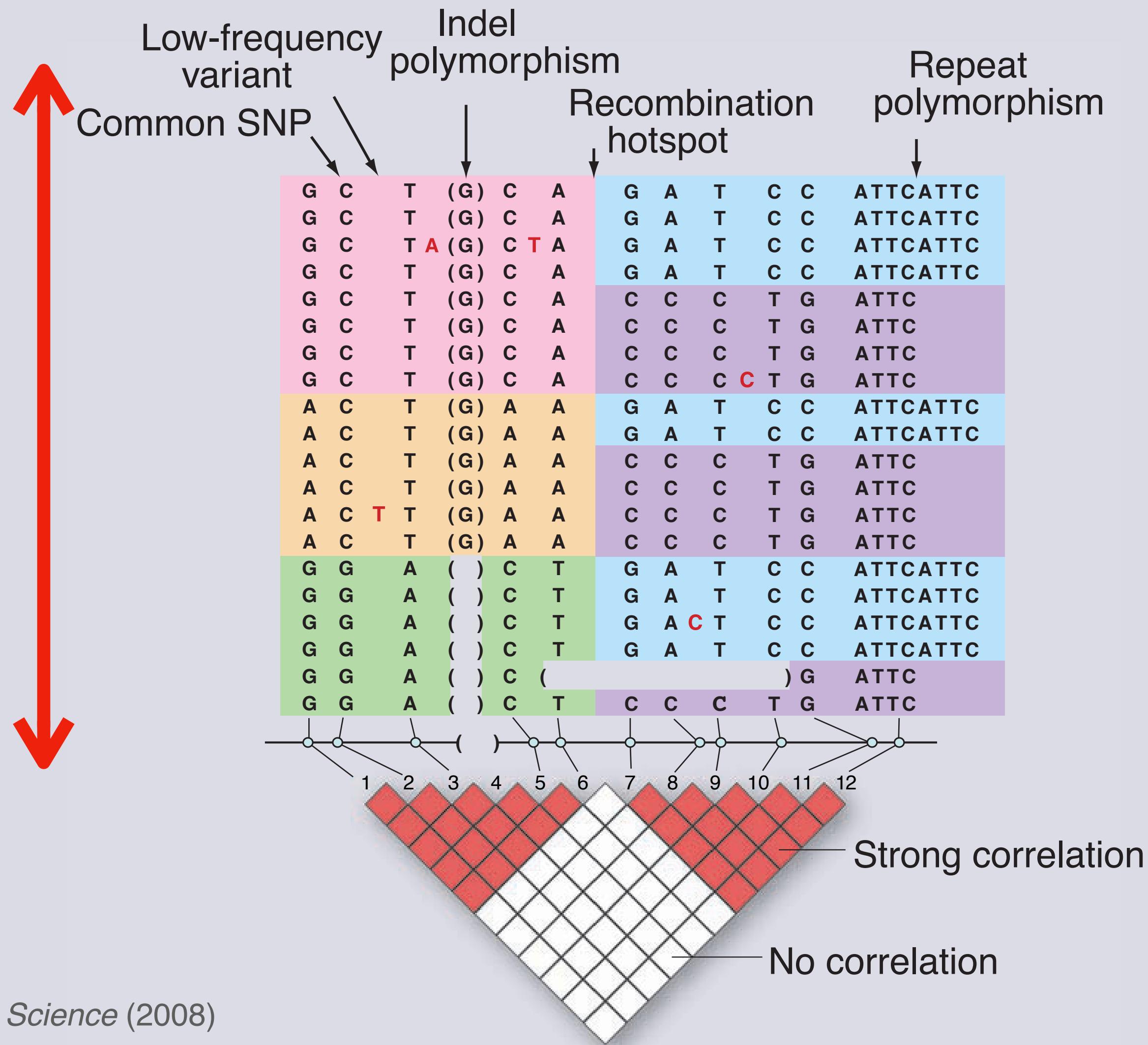


Genome-wide association study (many loci)

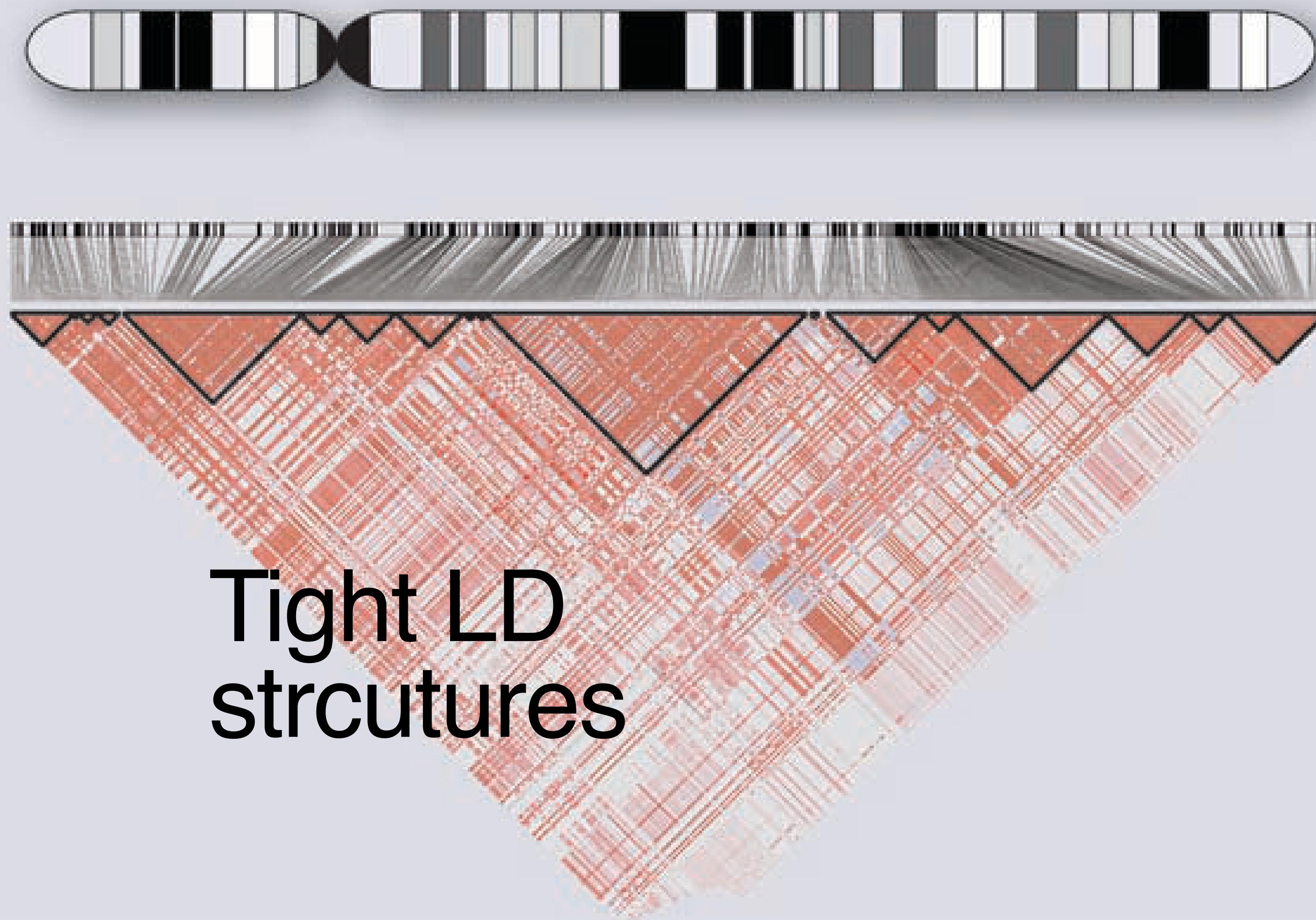


Genetic variation in one figure

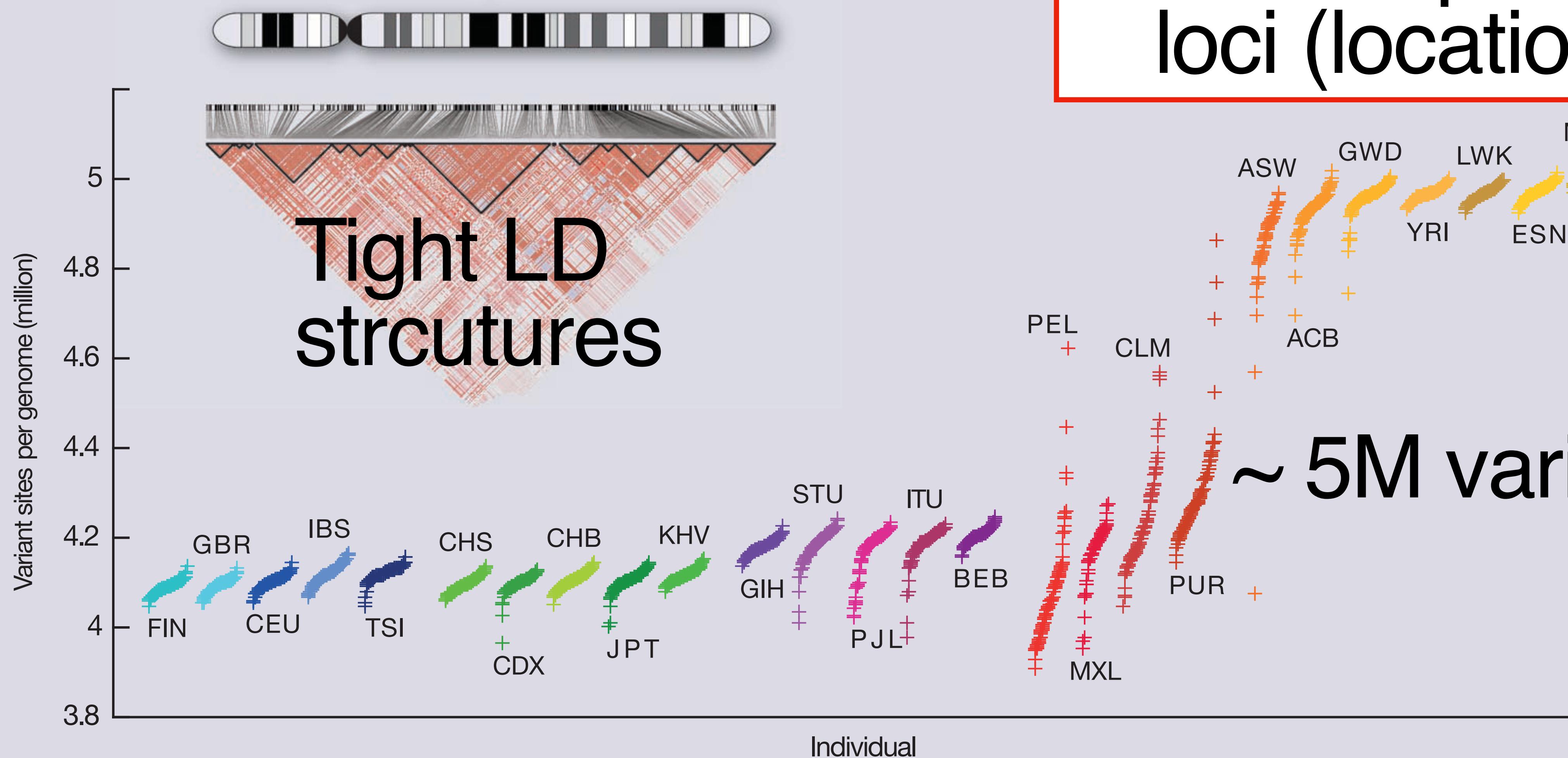
across
many
individuals
(diploid
genomes)



Common SNPs
Insertion/deletion
Other low-freq. variants
Other structural variants
Recombination hotspot



1M "tagging" SNPs across the Human Genome



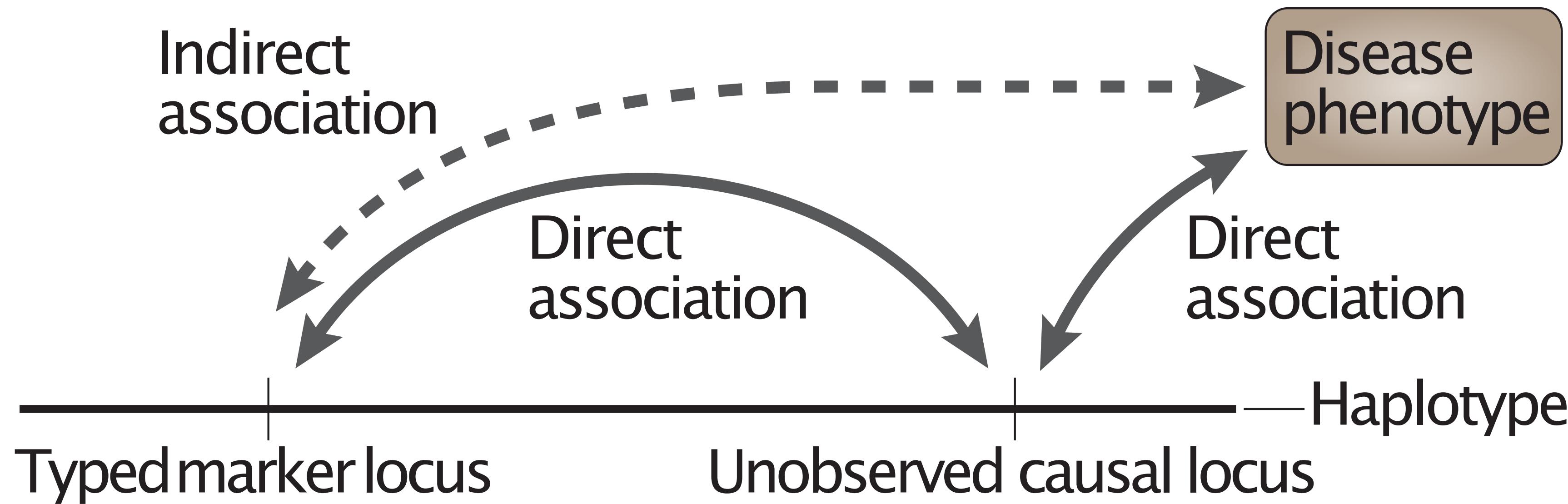
Genome-wide significance level

$$p_j < \frac{0.05}{\text{number of (independent) tests}} = \frac{0.05}{10^6} = 5 \times 10^{-8}$$

$$\text{FWER} = P \left(\cup_{j=1}^{10^6} \{p_j < 0.05/10^6\} \mid H_0 \right)$$

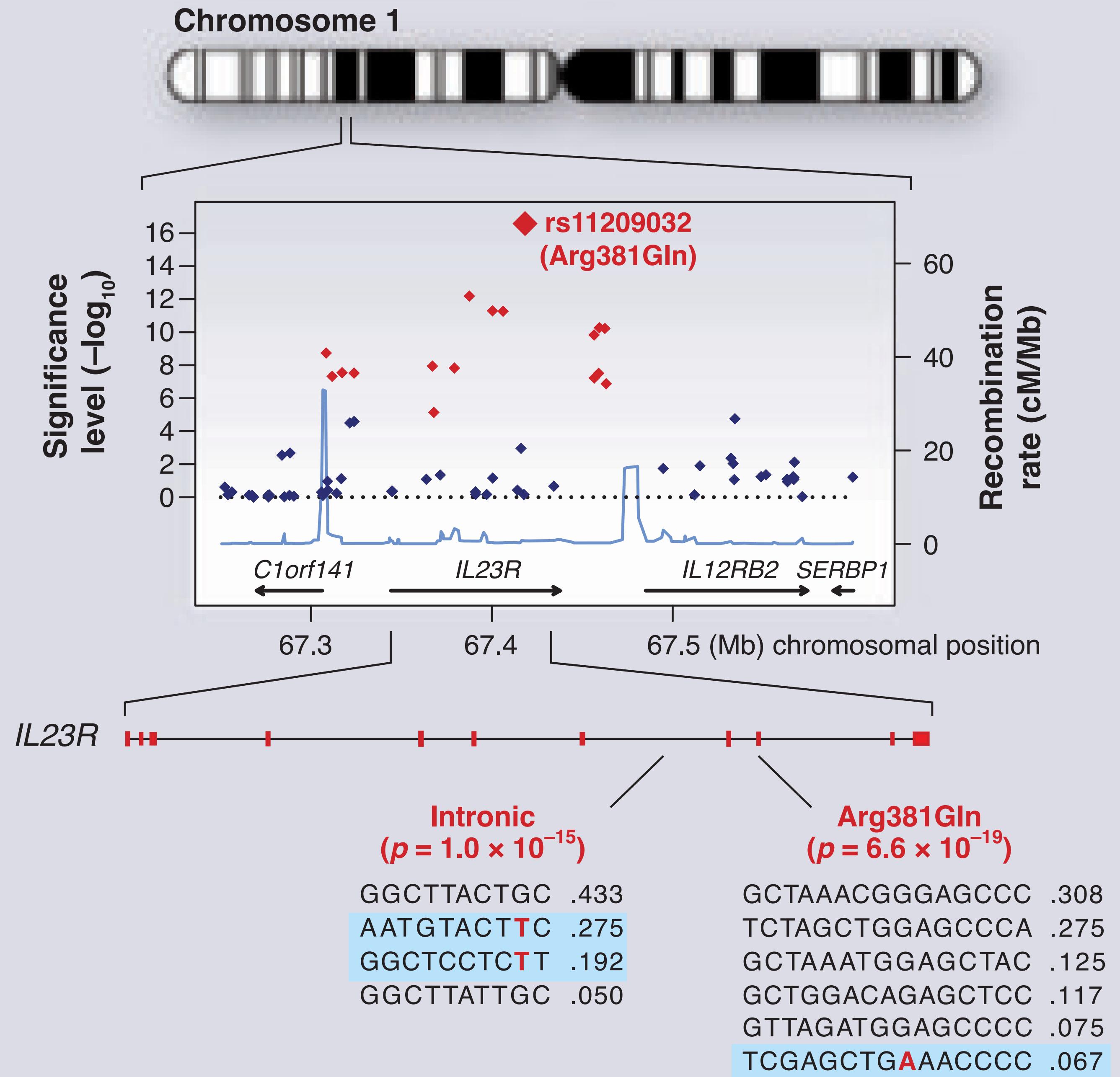
Union bound $< \sum_{j=1}^{10^6} P(p_j < 0.05/10^6 \mid H_0) = \sum_{j=1}^{10^6} \frac{0.05}{10^6} = 0.05$

Types of association (typed \approx tagging SNP)



GWAS only teach
approximate locations of
causal variants in the
genome

A close-up view of Genome-wide- significant loci



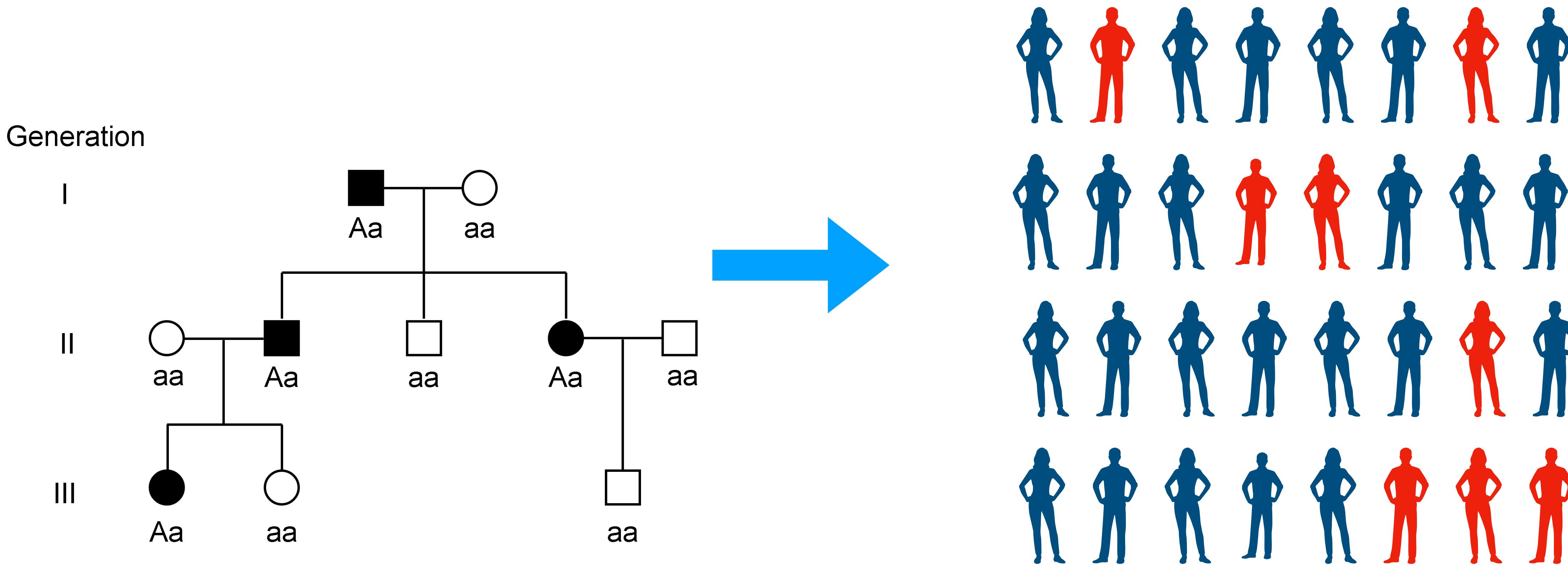
Summary for GWAS

- Massive multiple hypothesis testing associations between genetic variants and phenotypic changes
- Most studies focus on biallelic variants, coding them as {0, 1, 2}
- Most association statistics concern an additive effect of genotypes (linearity)
- Genetic variants closely located in the genome are strongly correlated with each other due to recombination
- Implicitly assume ~ 1M genetic variants (of tagging, or representative SNPs)
- GWAS only teach an approximate genomic locations associated with phenotypes

Today's lecture

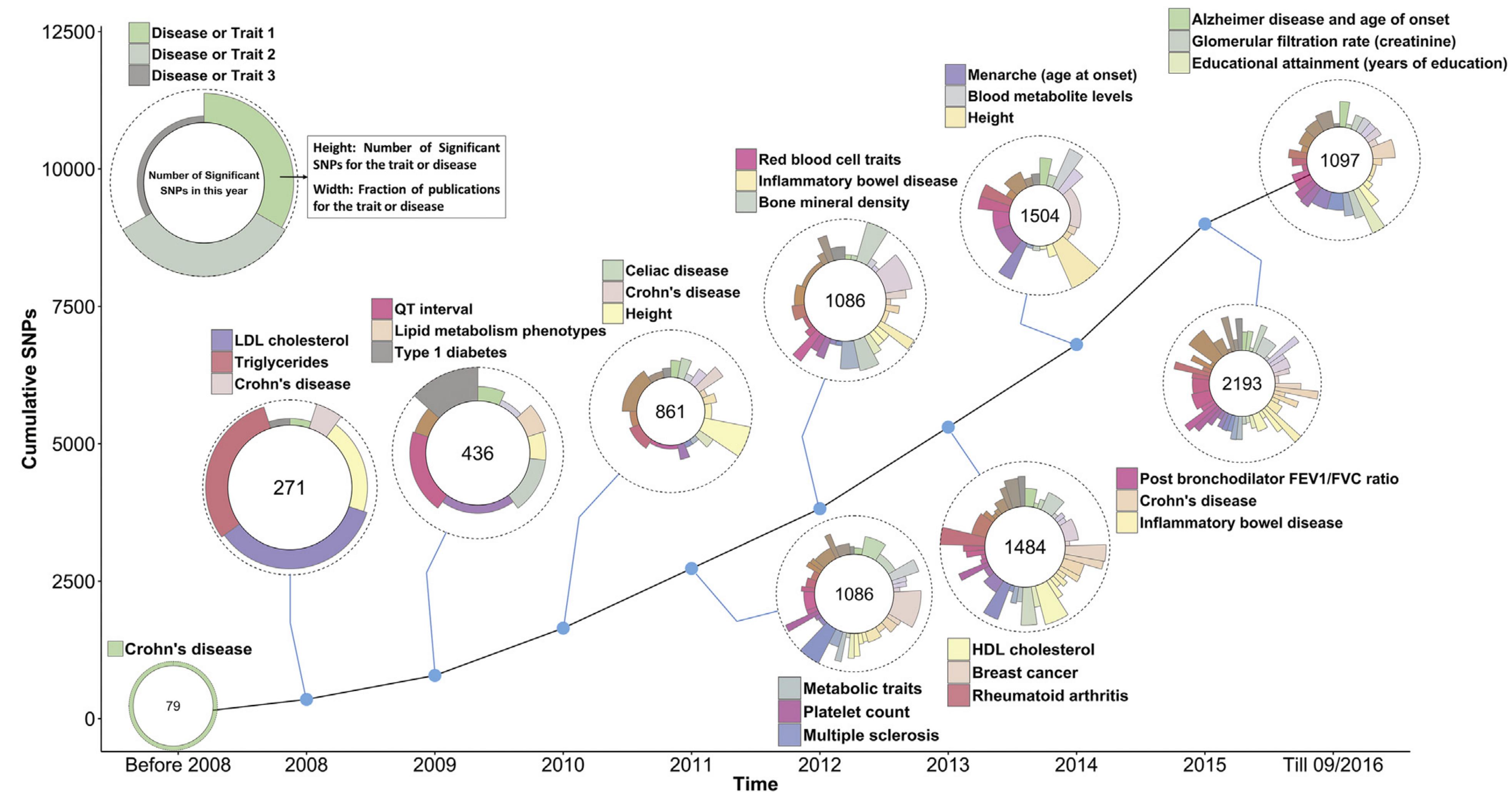
- Fundamentals in Population Genetics
- Genome-wide association studies
- What are the limitations of GWAS?

The success of GWAS changed the paradigm of human genetics studies

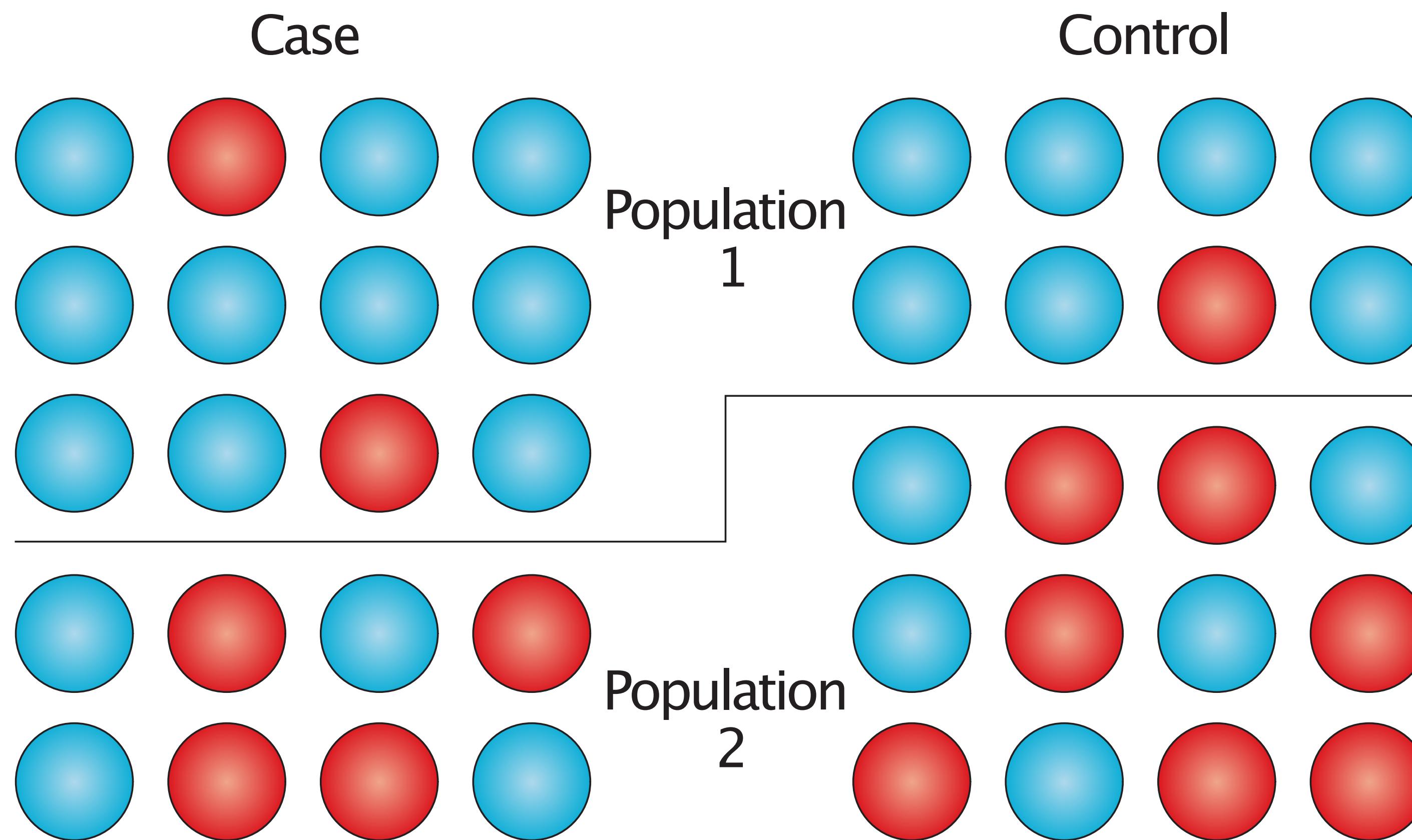


Are we assuming a
homogeneous population?

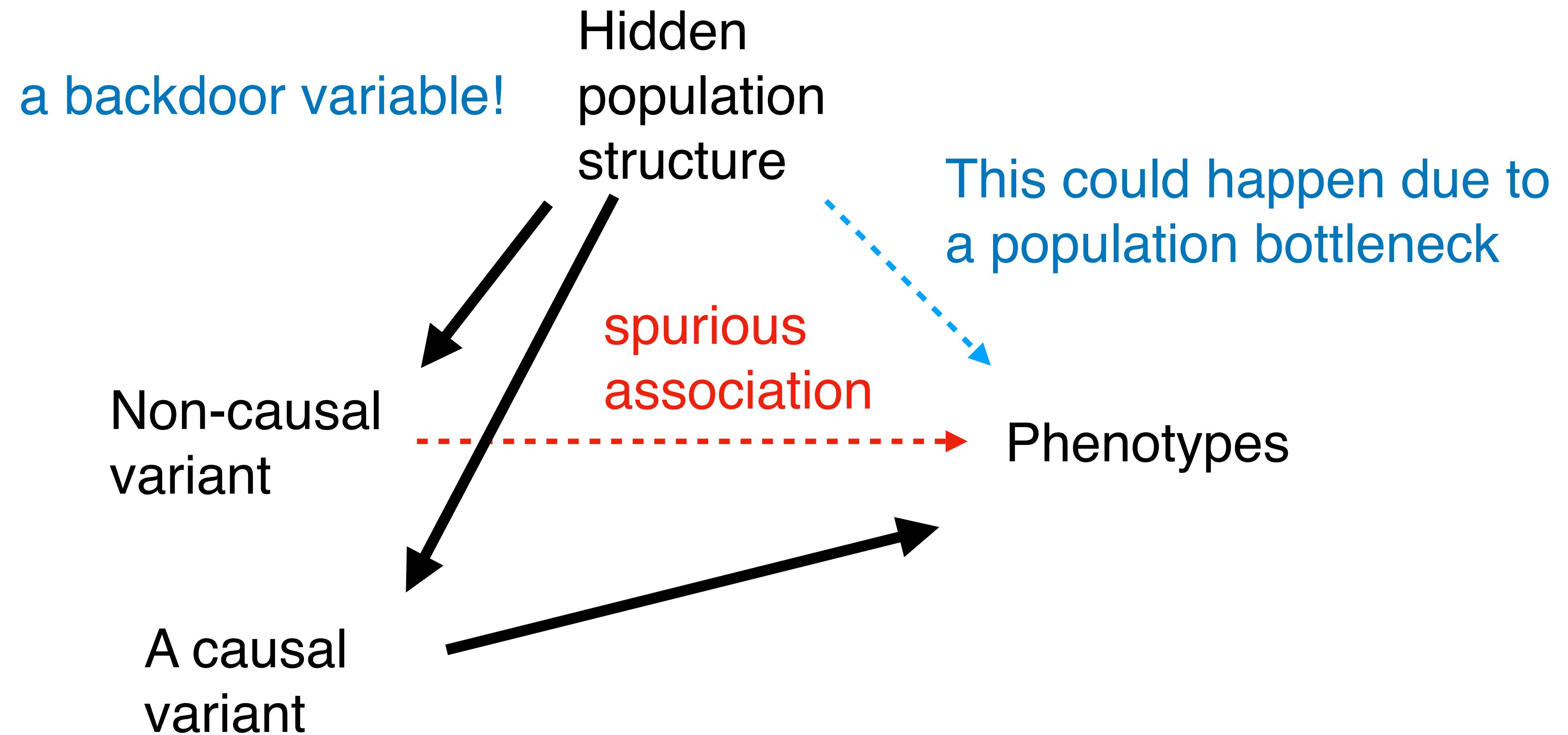
GWAS data/results increased rapidly



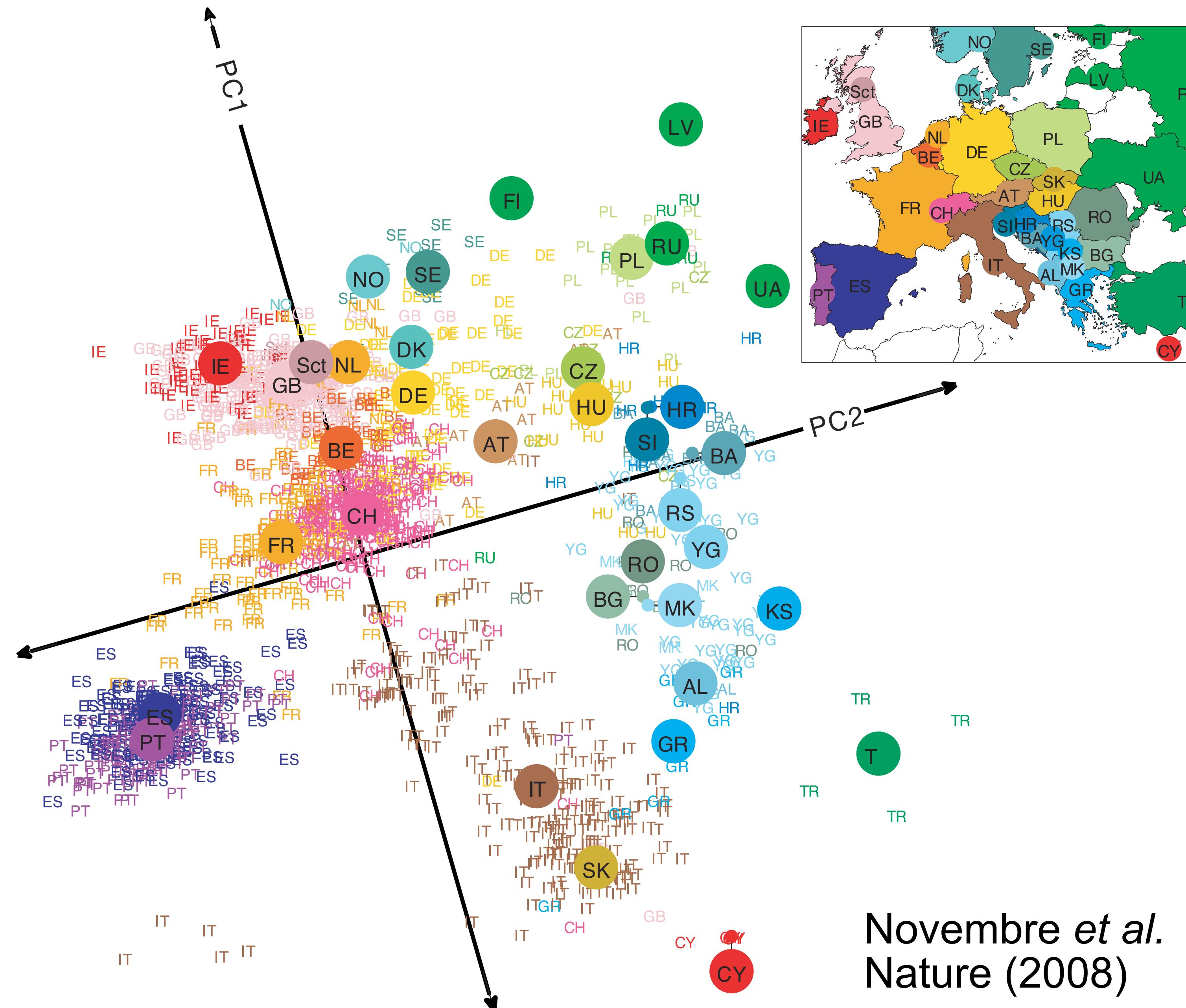
What if there were a hidden population structure?



Hidden population structure in a causal graph

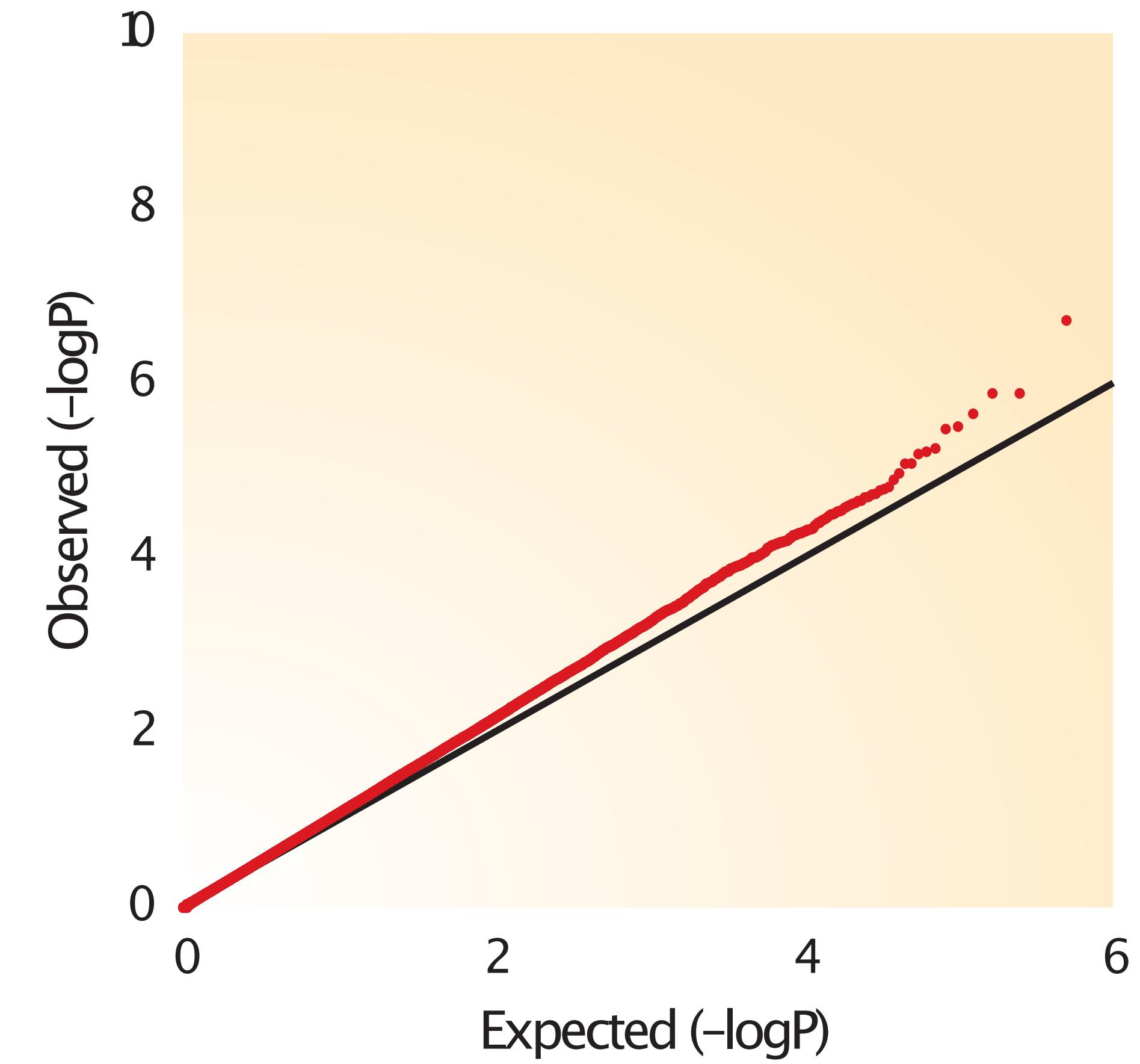
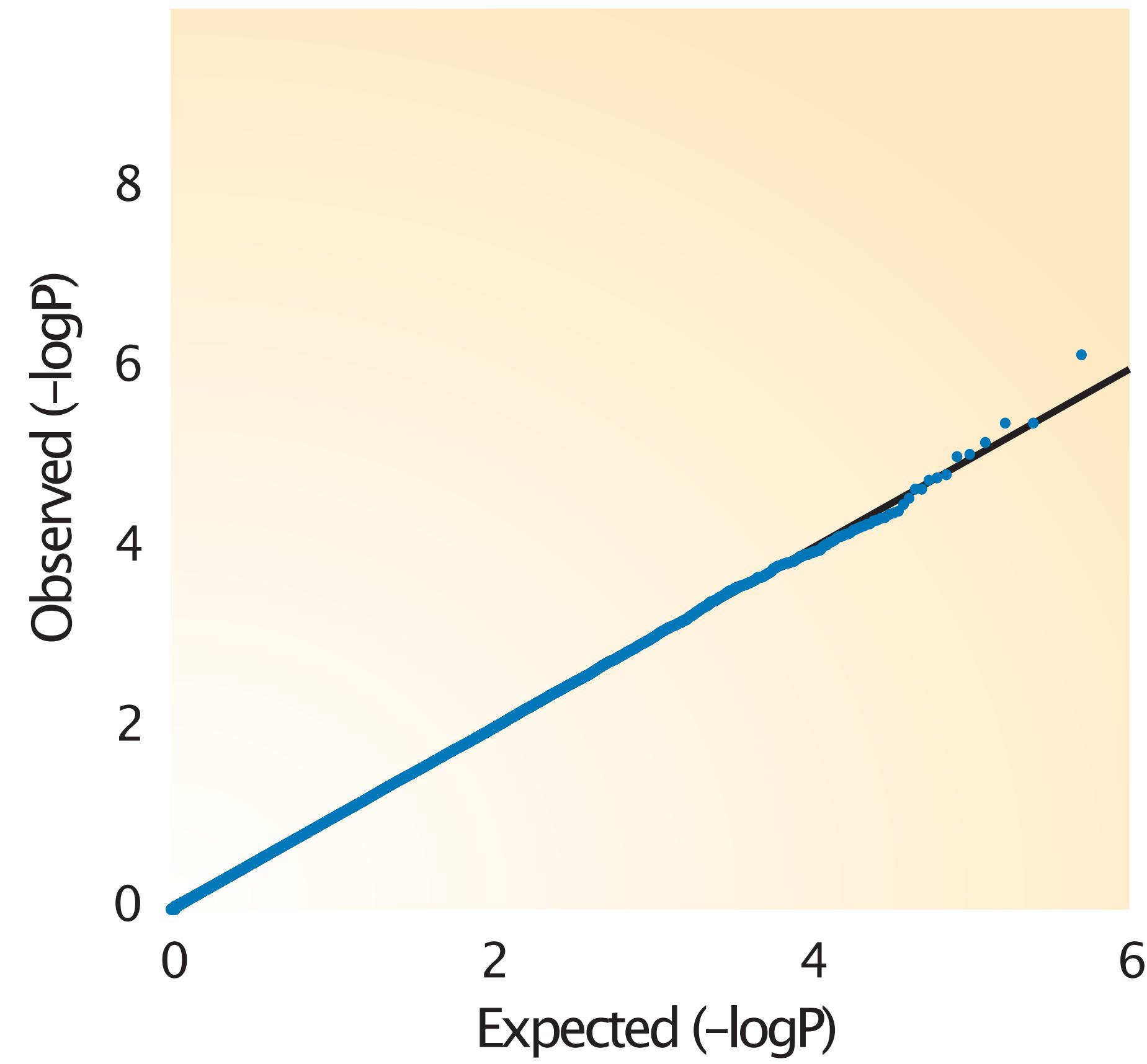


The human population as a result of the human migration history



Novembre et al.
Nature (2008)

Population structures may inflate GWAS stats



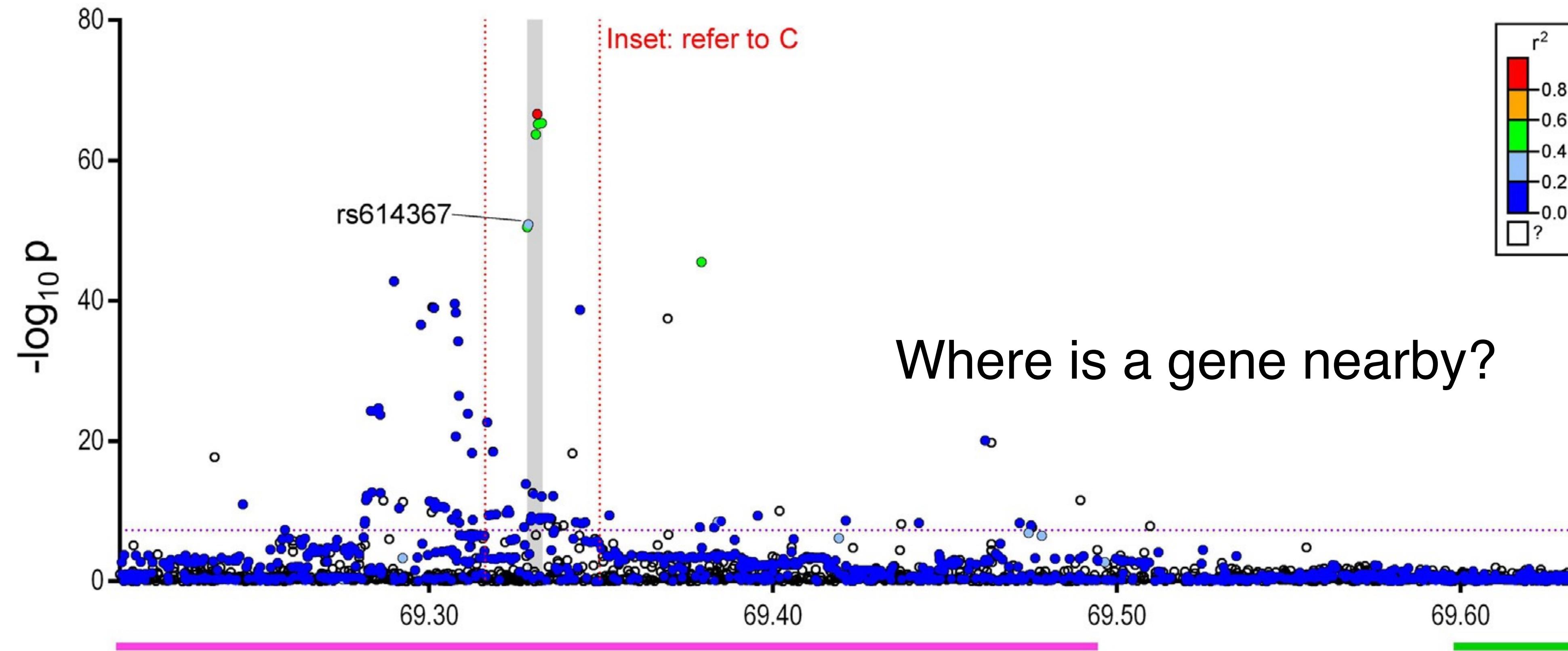
What could be a generative model?

$$Y_i = \sum_j X_{ij} \beta_j + \sum_k U_{ik} \gamma_k + \epsilon_i$$

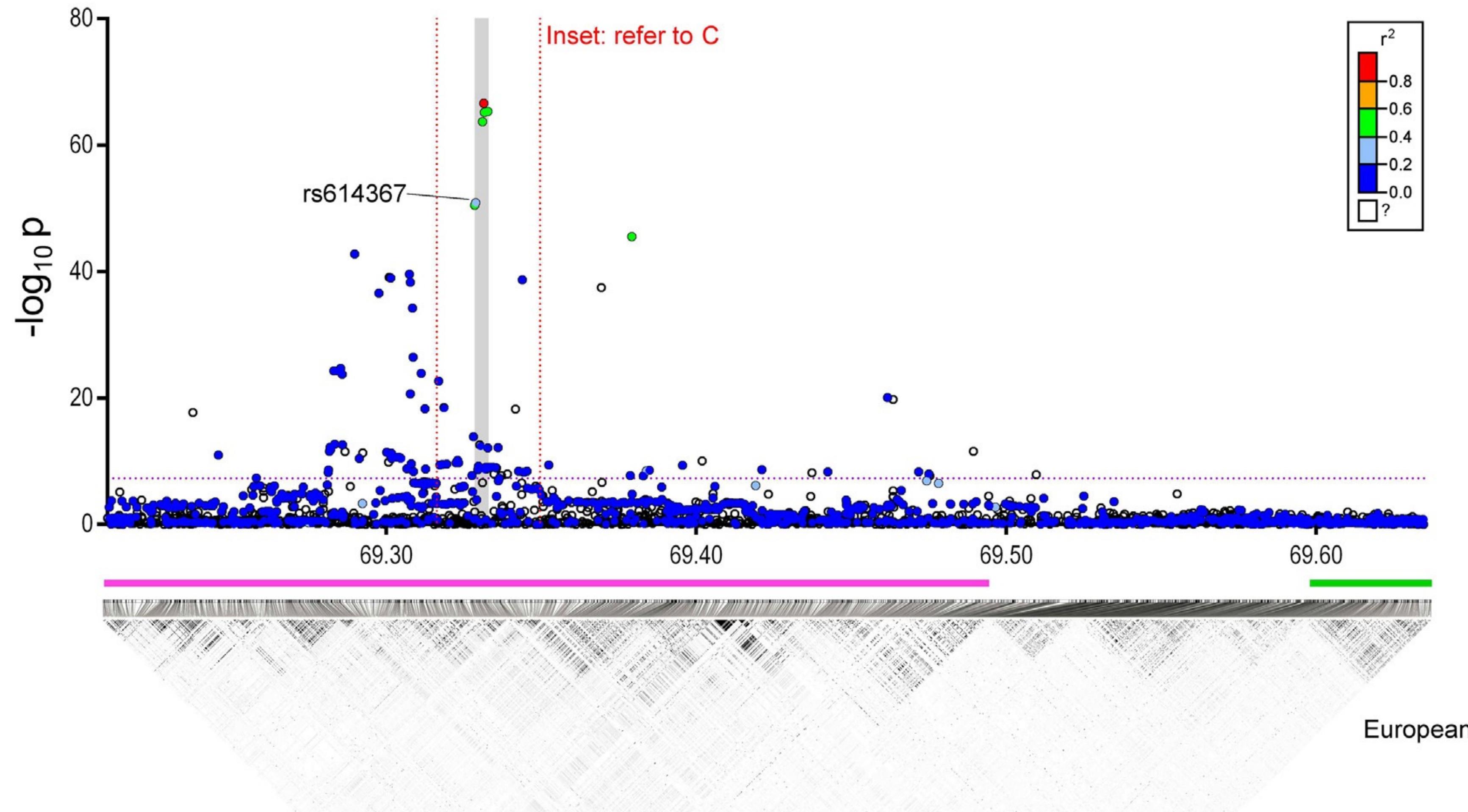
genotype
phenotype **popuation**

$$X_{ij} = \sum_k U_{ik} \alpha_k + \delta_{ij}$$

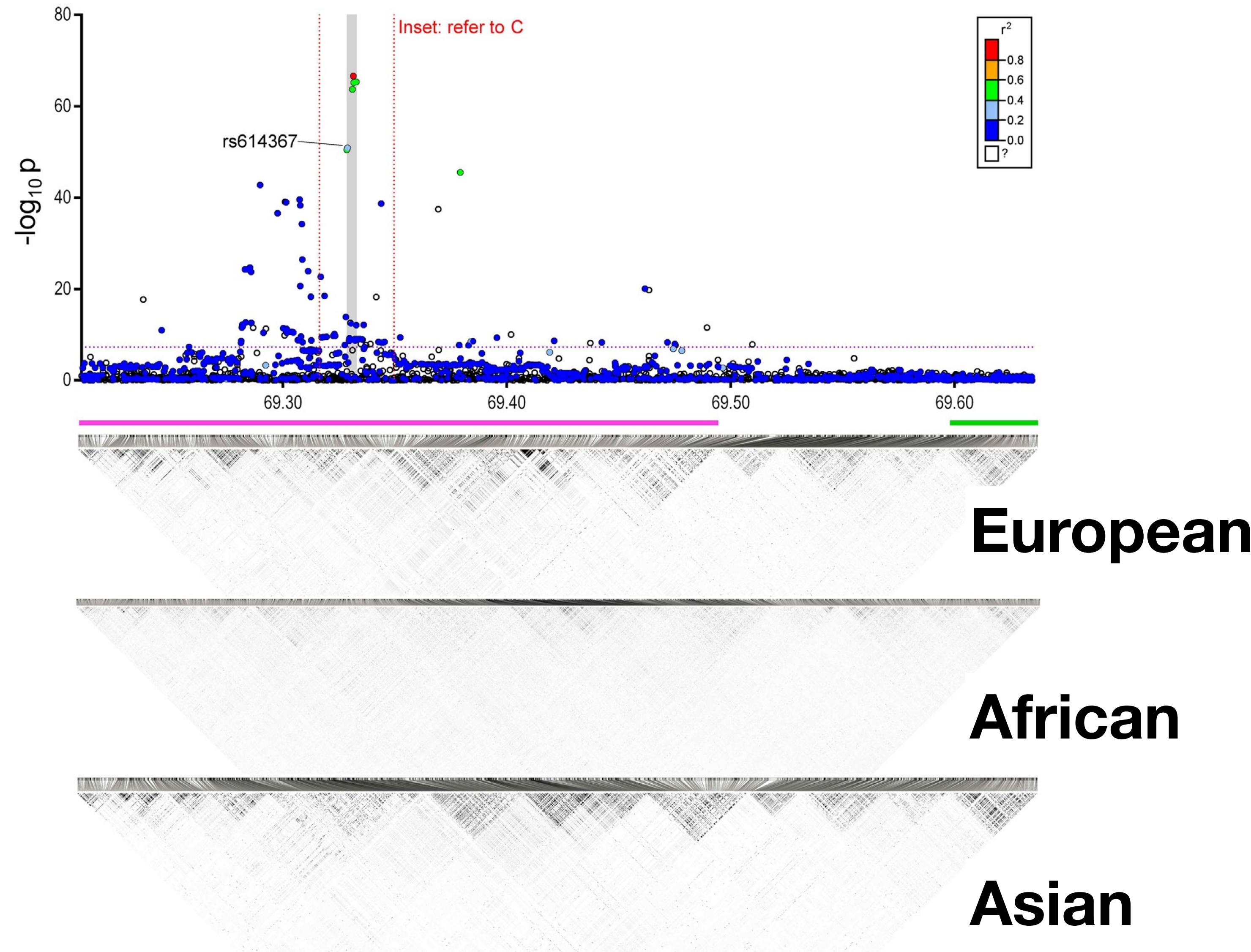
90% of GWAS hits on the non-coding regions



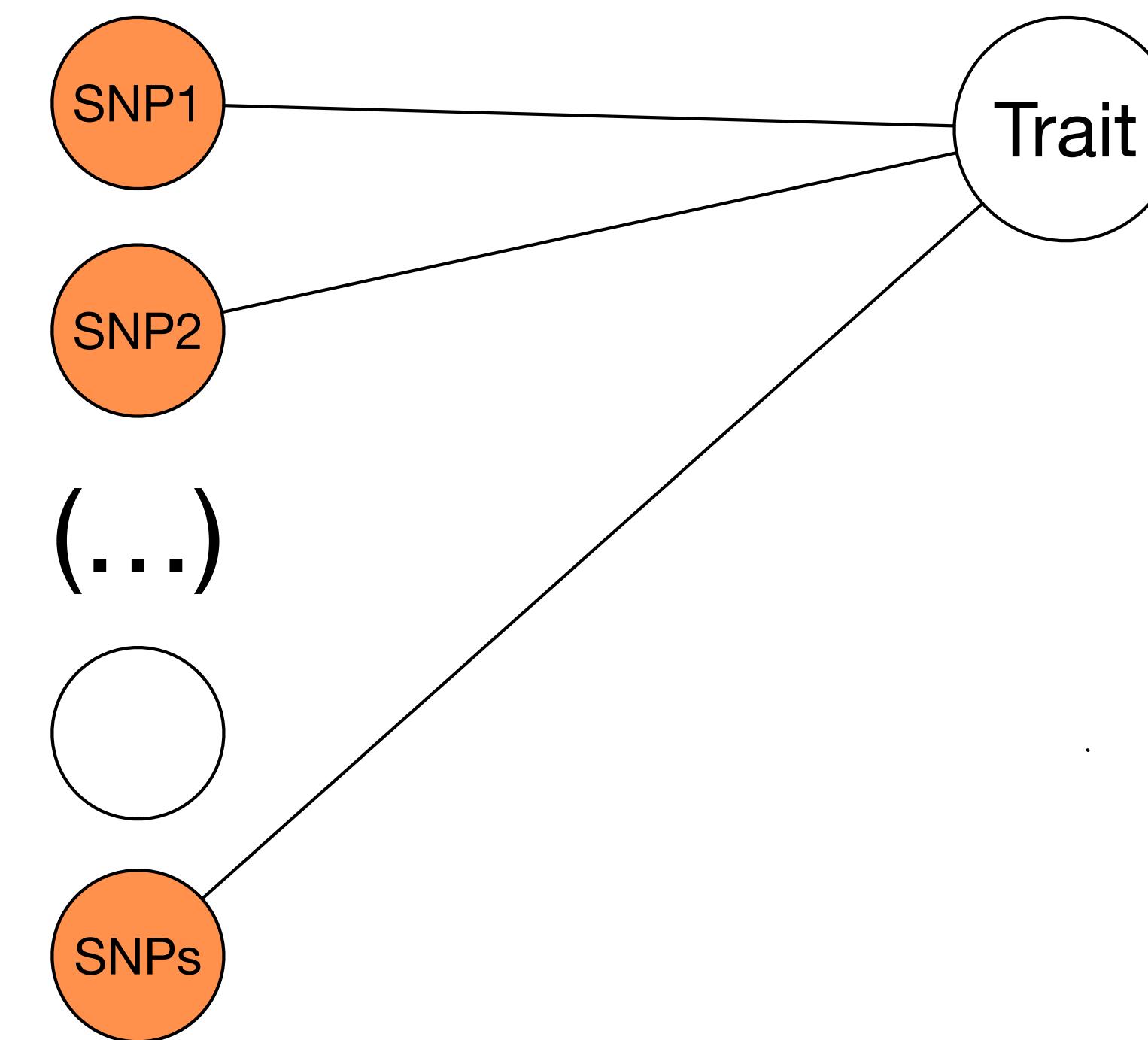
LD structure, many significantly associated variants



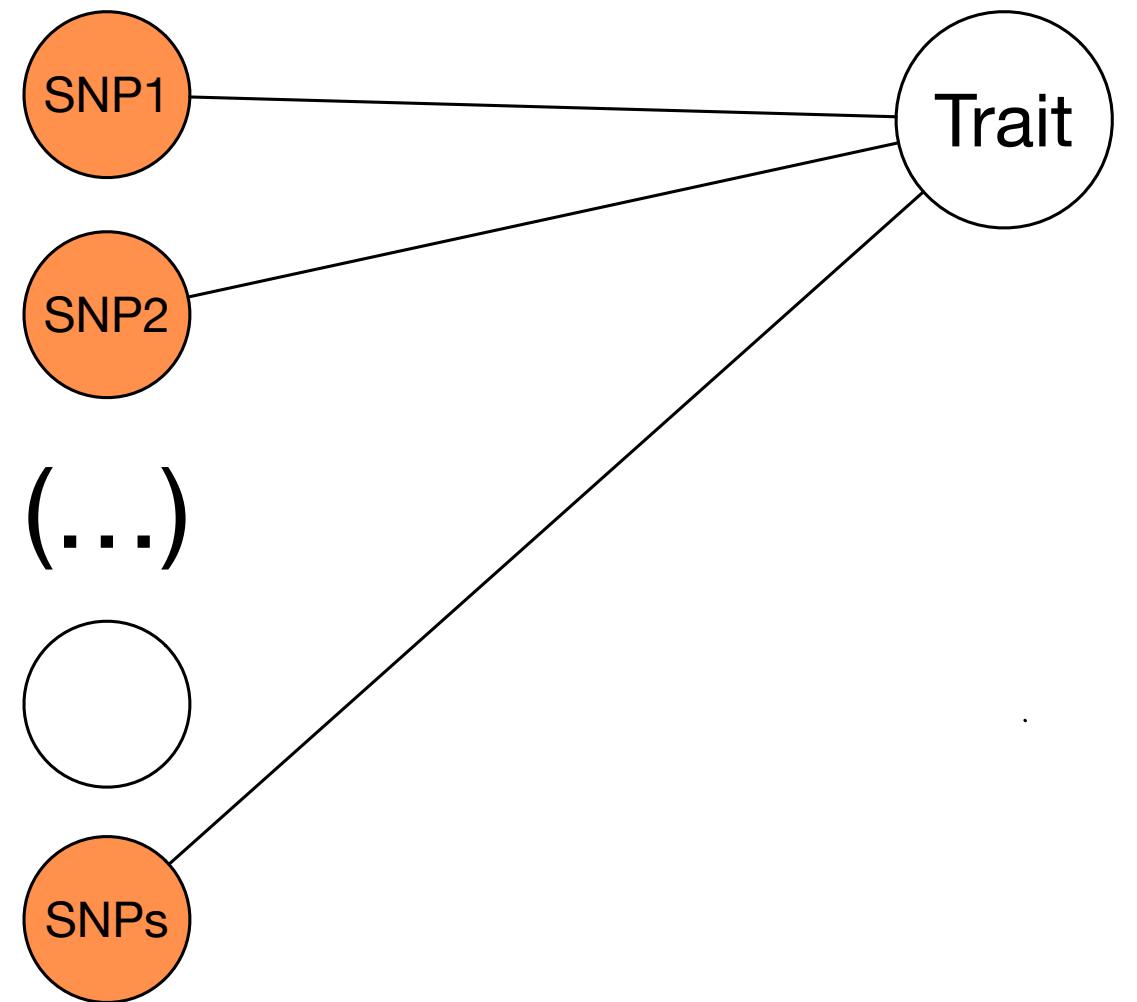
Different population structures, different LD patterns



Many complex traits are polygenic



A missing heritability problem



A thousand, if not thousands of independent weak effect variants explain total heritability

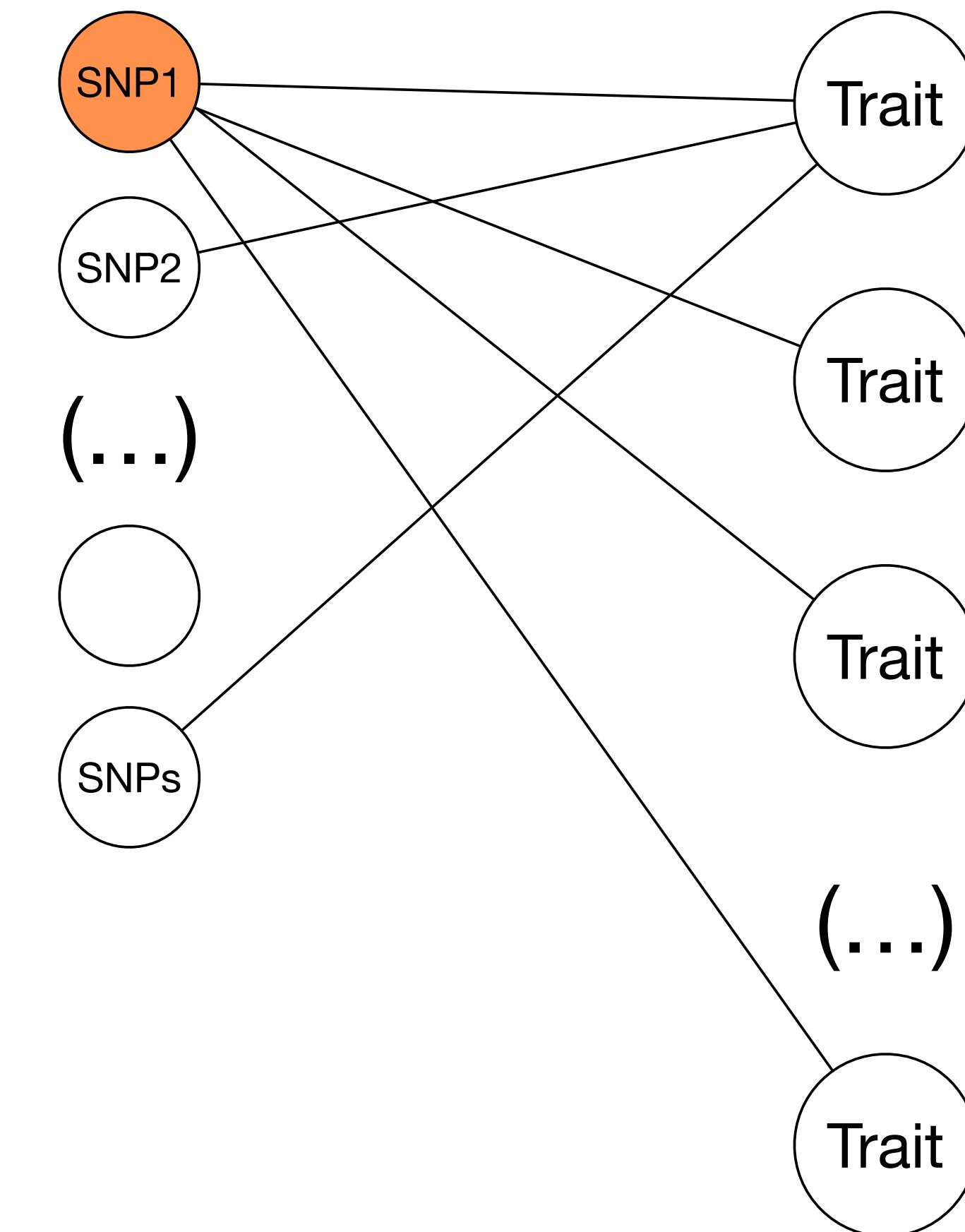


The case of the missing heritability

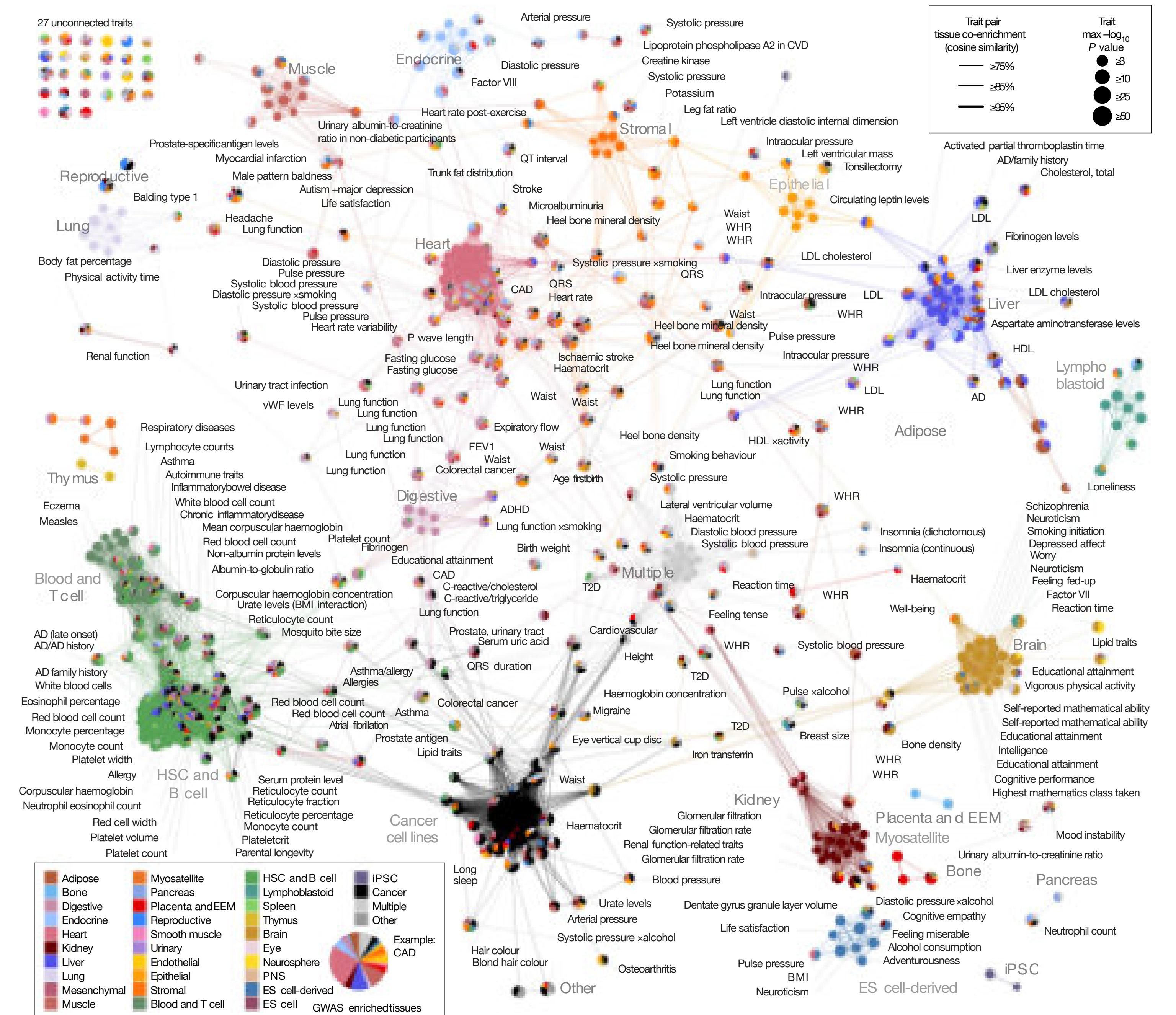
A single variant may act on many traits

Pleiotropy

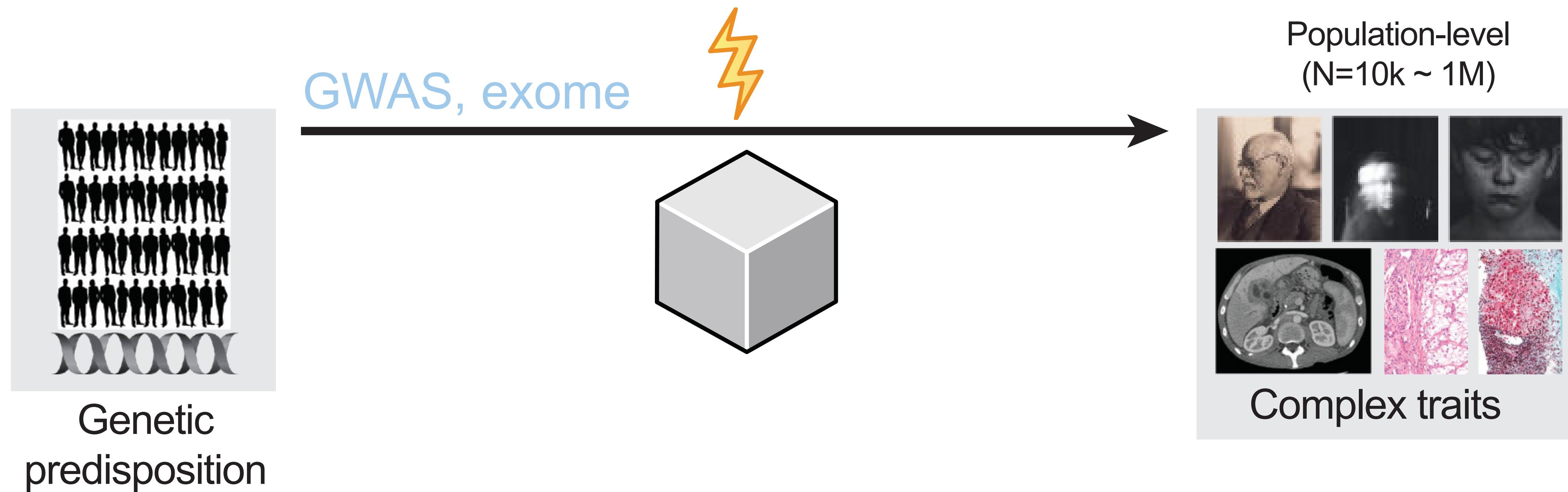
A single variant is associated, over and over with many different human traits!



Most traits are also connected with other traits because of pleiotropy



There is no contextual information in GWAS



A genetic variant \neq a gene \neq a whole mechanism

Today's lecture

- Fundamentals in Population Genetics
- Genome-wide association studies
- What are the limitations of GWAS?

