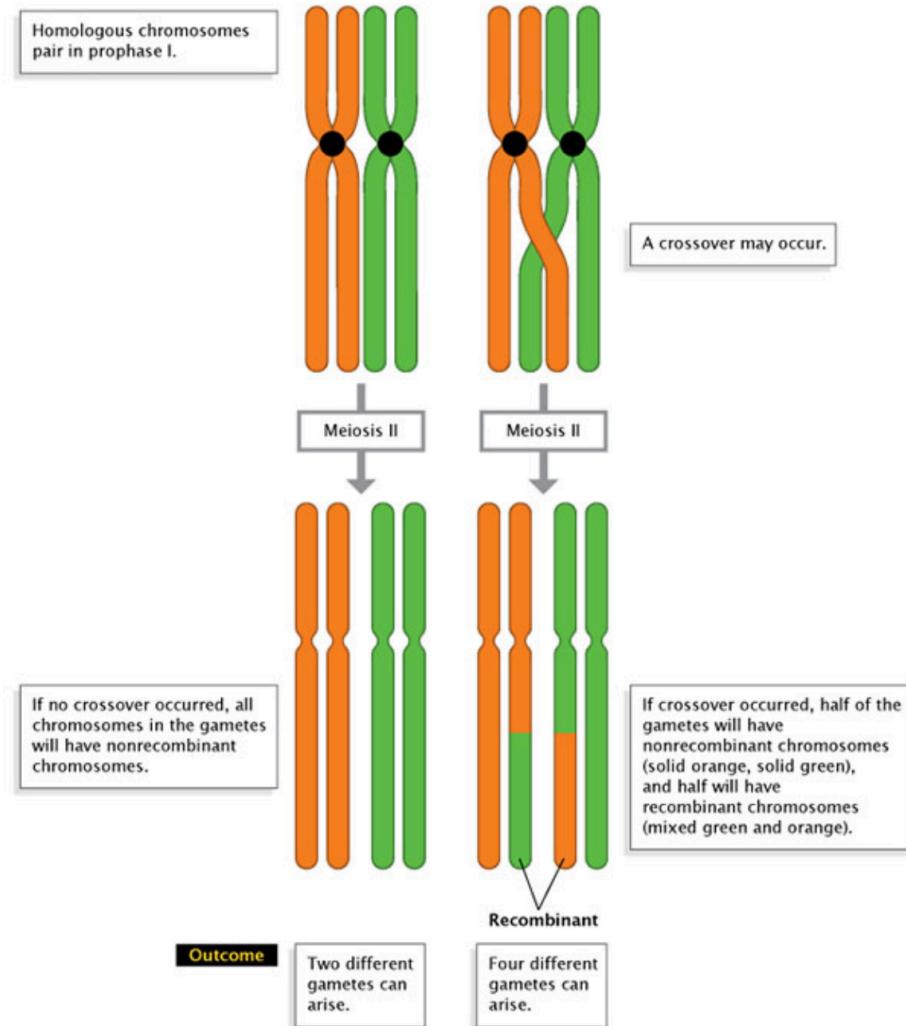


# **Module 10: Inheritance of Polygenic Diseases**

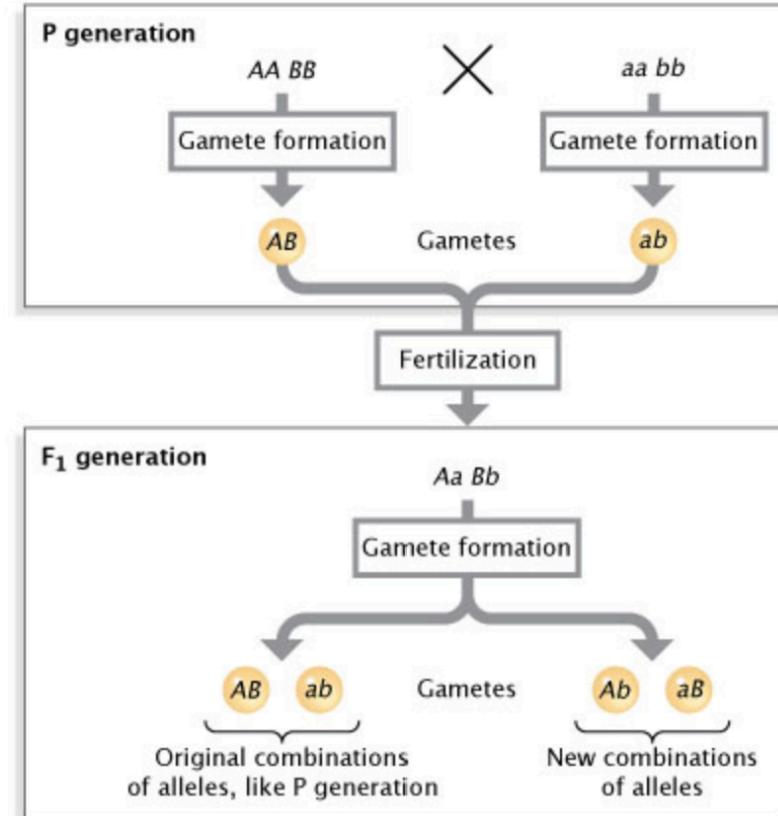
# Key Concepts

- Linkage Analysis
- HWE and Allele Frequency
- Genome Wide Association Studies
- Applications

# Recombination<sup>1</sup>

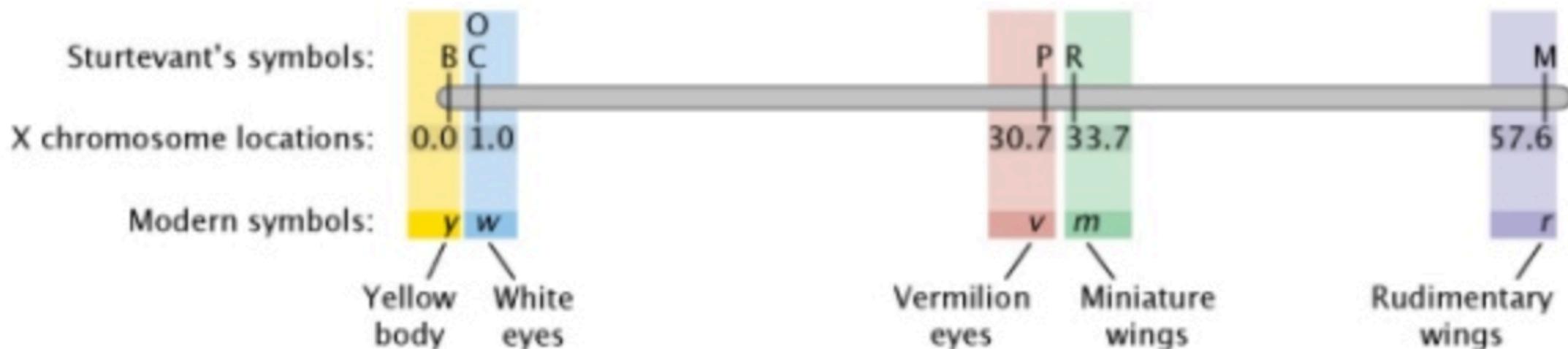


- Alleles recombine during crossover events in meiosis



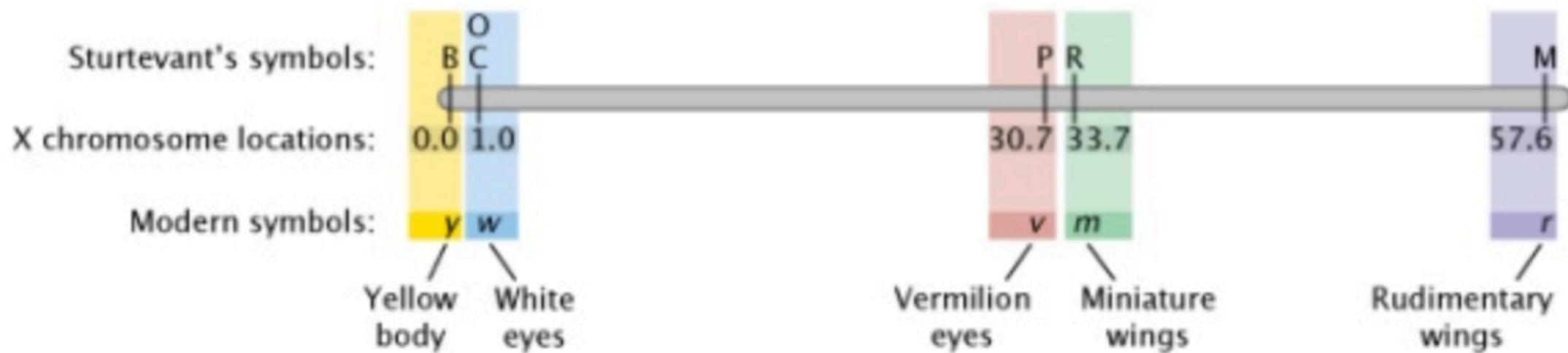
# Genetic Linkage

- The frequency of recombination is related to the distance between the genes on a chromosome.<sup>1</sup>
- Violation of Mendel's law of independent assortment<sup>1</sup>



# Genetic Linkage

- The closer two genes were to one another on a chromosome, the greater their chance of being inherited together.<sup>1</sup>
- Genes located farther away from one another on the same chromosome were more likely to be separated during recombination.<sup>1</sup>



# Linkage Map

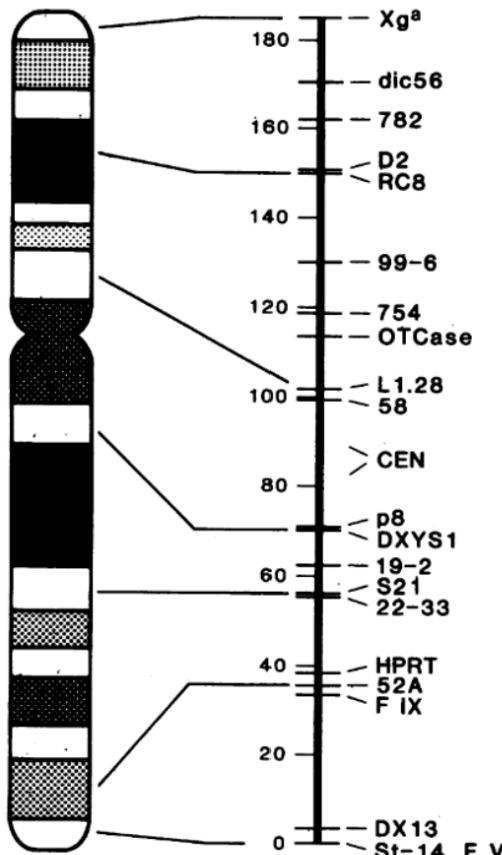
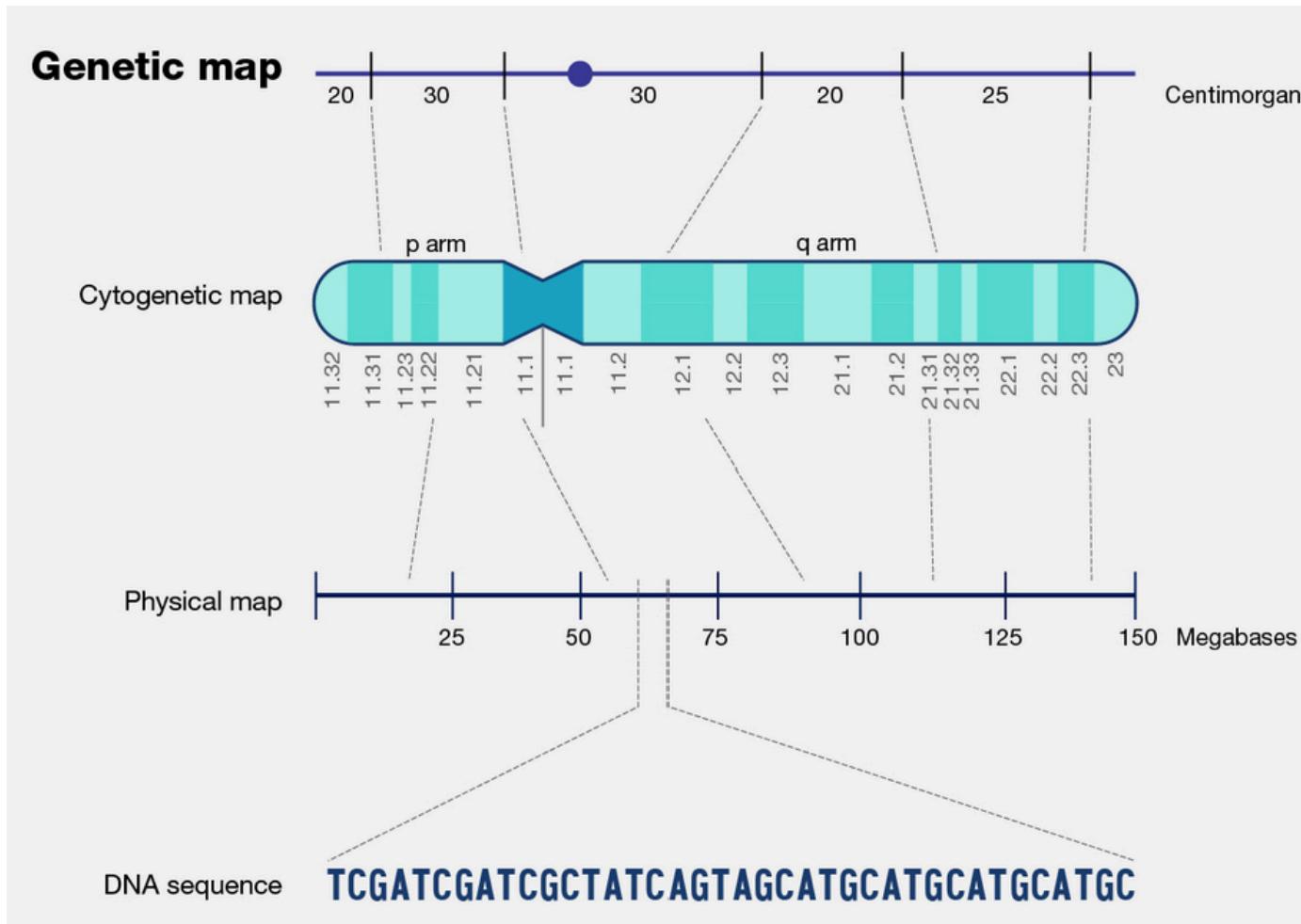


Fig. 1. Linkage map of the human X chromosome. Genetic distances are represented in recombination units, and physical locations of selected markers are indicated on the ideogram.

- **Linkage map:** (also called genetic map) the order and the linear distances between the linked genes<sup>1</sup>
- **Map unit (mu):** arbitrary unit of distance representing a recombination frequency of 1%<sup>1</sup>
  - renamed as **centimorgan (cM)**
- Completely linked genes = 0 mu
- Unlinked genes > 50 mu
- Linked genes < 50 mu

Drayna & White (1985)

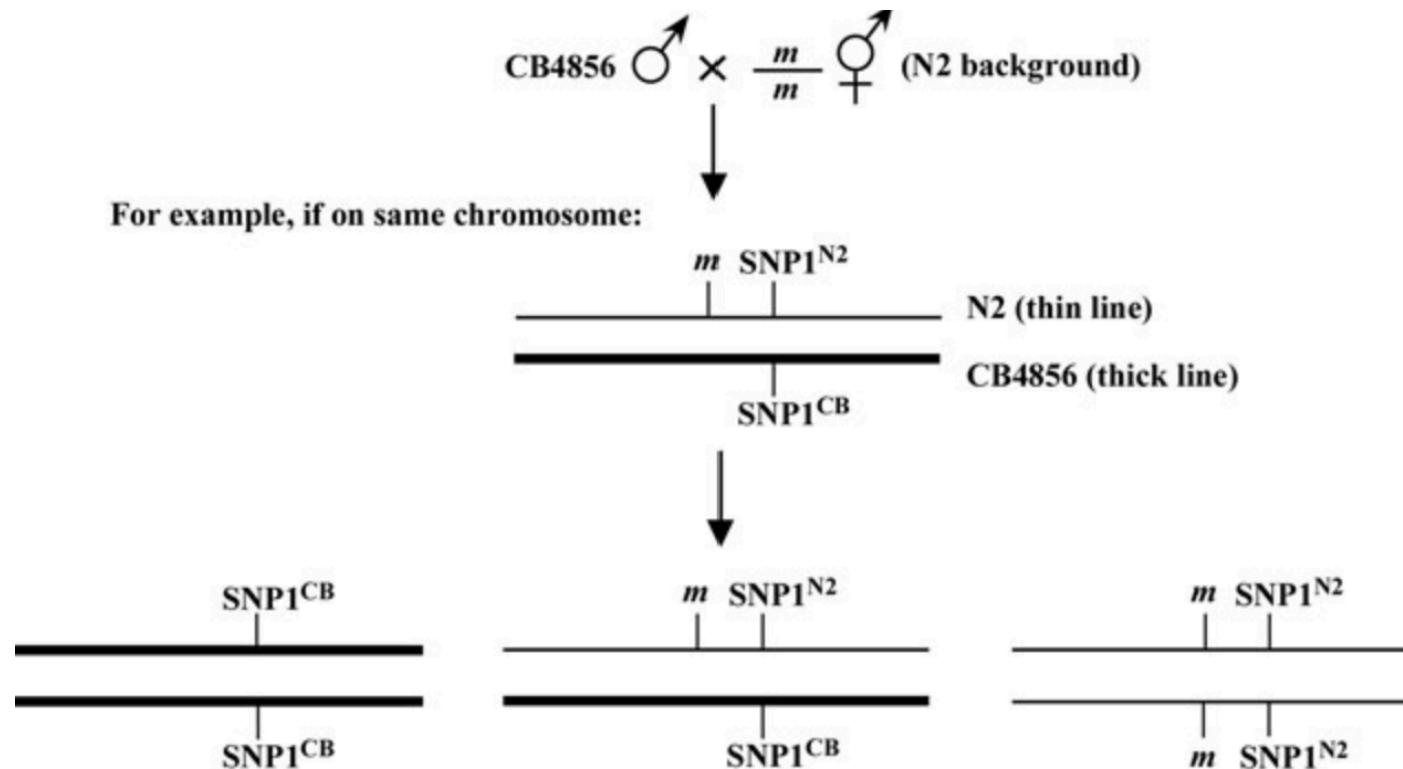
# Linkage Map



- On average 1 cM corresponds to 1 million base pairs.<sup>2</sup>
- Linkage map is NOT the same as a cytogenetic map (karyotype)<sup>3</sup>

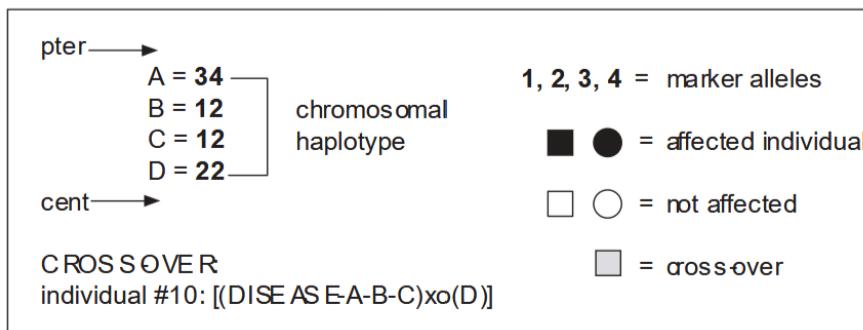
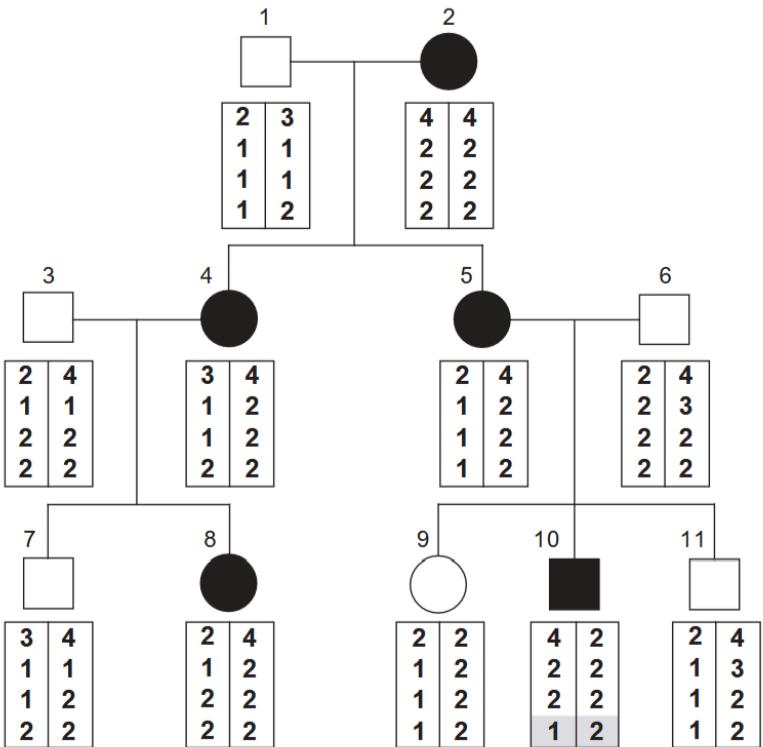
# Two-point linkage

- Calculate the linkage between a putative disease locus and a single marker locus



## Linkage Analysis: Next Steps

- To identify the smallest region of the genome that should contain the disease gene, the minimum-candidate region (MCR), we can perform:
  - i. Haplotype analysis or
  - ii. Multipoint linkage analysis

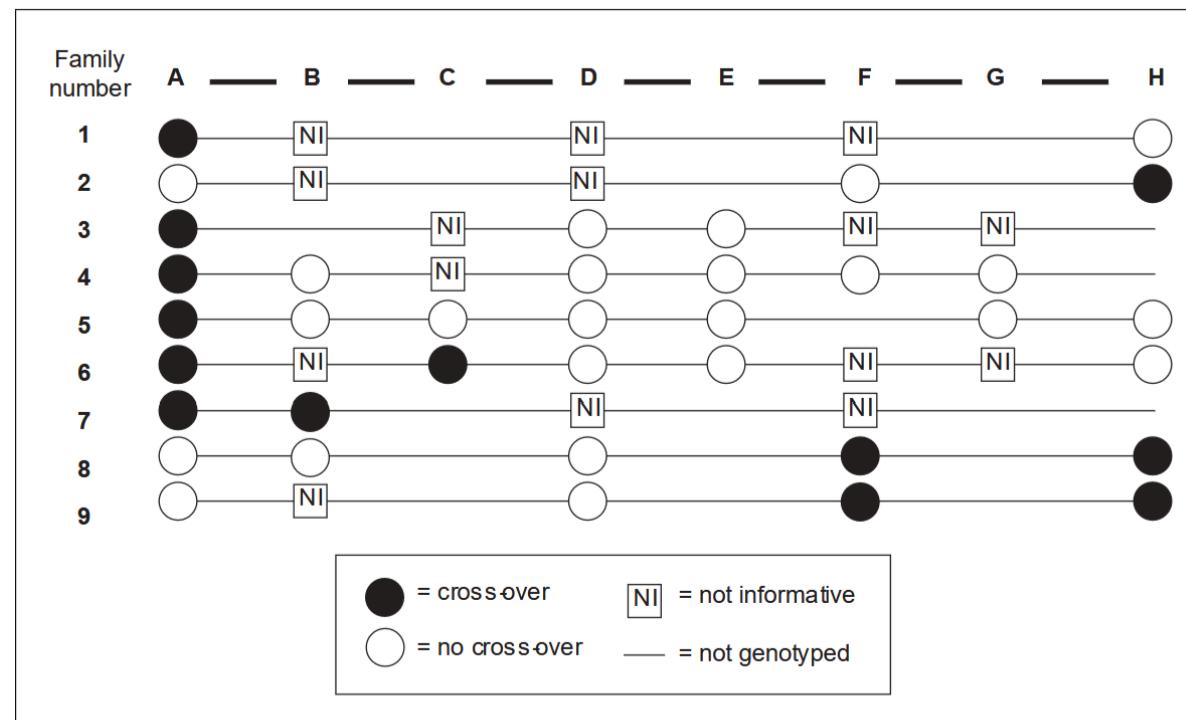


## Haplotype Analysis<sup>7</sup>

- **Haplotype:** the arrangement of alleles on sister chromatids
- Construct haplotypes by examining the ordered markers genotyped in an individual and arranging them by chromatid, on the basis of parental types

# Haplotype Analysis<sup>7</sup>

- After summarizing all cross-over individuals in the study, the most likely region for the disease gene is the region where no cross-overs are found.



**Figure 1.4.7** Method for summarizing all cross-over individuals in all family data used in a study. The most likely region for the disease gene is between markers C and F, the region where no cross-overs are found.

## Multipoint analysis - Parametric<sup>4</sup>

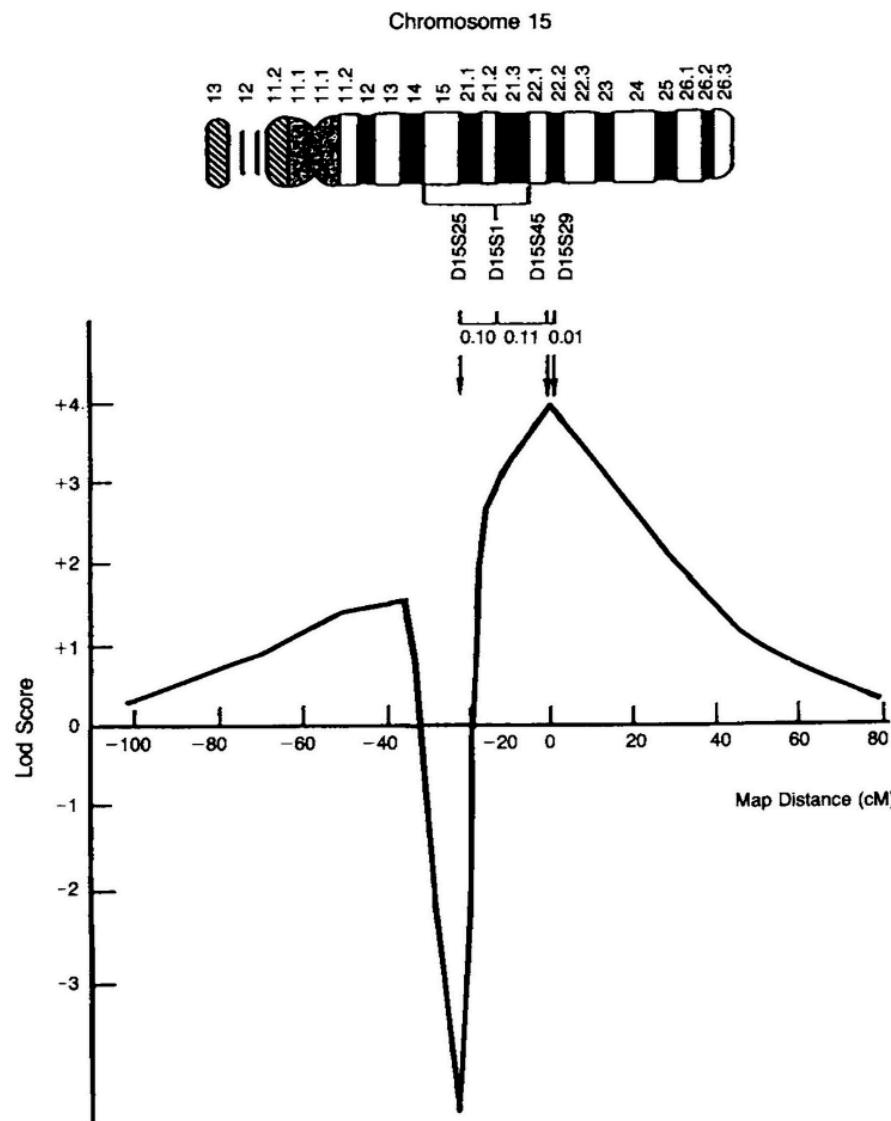
- Assume a model describing the probability of phenotype given genotype at the disease locus
- Calculate the likelihood ratio under the hypothesis that a disease gene is at  $x$ , versus the hypothesis that it is unlinked to  $x$ .

$$\text{LOD score: } Z(x) = \log_{10} \left( \frac{L(x)}{L(\infty)} \right) \text{ where}$$

- $L(x)$  = Likelihood of the disease locus at position ( $x$ ) on the marker map  
= the probability of **linkage** at a given recombination fraction ( $\theta$ )
- $L(\infty)$  = Likelihood of the disease locus being infinitely far from a set of markers  
= the probability of **no linkage** ( $\theta = 0.5$ )

## How to use LOD score<sup>5</sup>

- LOD score of 3 means the odds are 1,000:1 that the two genes are linked and therefore inherited together.
- **LOD score>3.3** → Strong evidence for linkage (genome-wide significance  $p<0.05$ )
  - Will still require replication with independent dataset
- $3.3 > \text{LOD score} > 1.86$  → Suggestive linkage
- $\text{LOD score} < 1.86$  → Nominal linkage



## Multipoint analysis

- Multipoint LOD scores generate a curve across a chromosome to identify the **likely position of a disease gene** (maximum of the curve) given that the max. LOD score  $\geq 3.3$ .<sup>4</sup>
- (Figure) Shows proof of linkage for Marfan Syndrome (autosomal dominant) in 5 families (LOD score = 3.92,  $\theta = 0.0 \pm 0.11$ ). The most probable location of the gene for the disease is D15S45 (LOD score = 3.32,  $\theta = 0.0 \pm 0.12$ ).<sup>6</sup>

# LOD scoring algorithms<sup>7</sup>

## Lander-Green algorithm

- Scales linearly with  $m$  markers , but exponentially with  $n$  individuals in the pedigree.
- More suitable for smaller pedigrees with many markers.
- Tool: MERLIN

## Elston-Stewart algorithm

- Scales linearly with  $n$  individuals, but exponentially with  $m$  markers.
- More suitable for large pedigrees.
- Tool: FASTLINK (inbred pedigrees; multiple founder pairs), VITESSE (without loops, one founder pair)

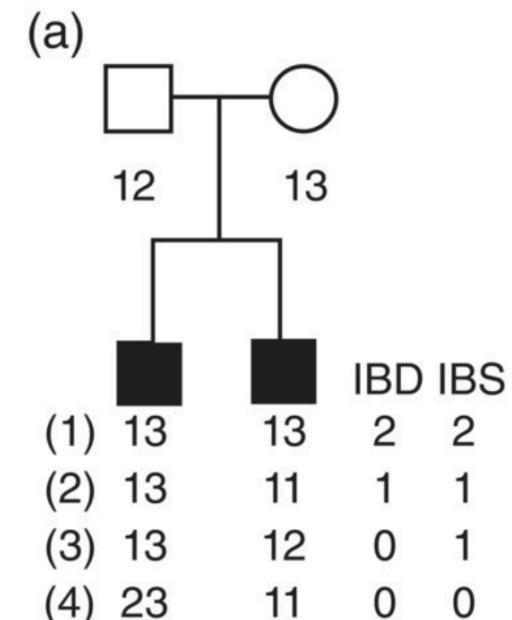
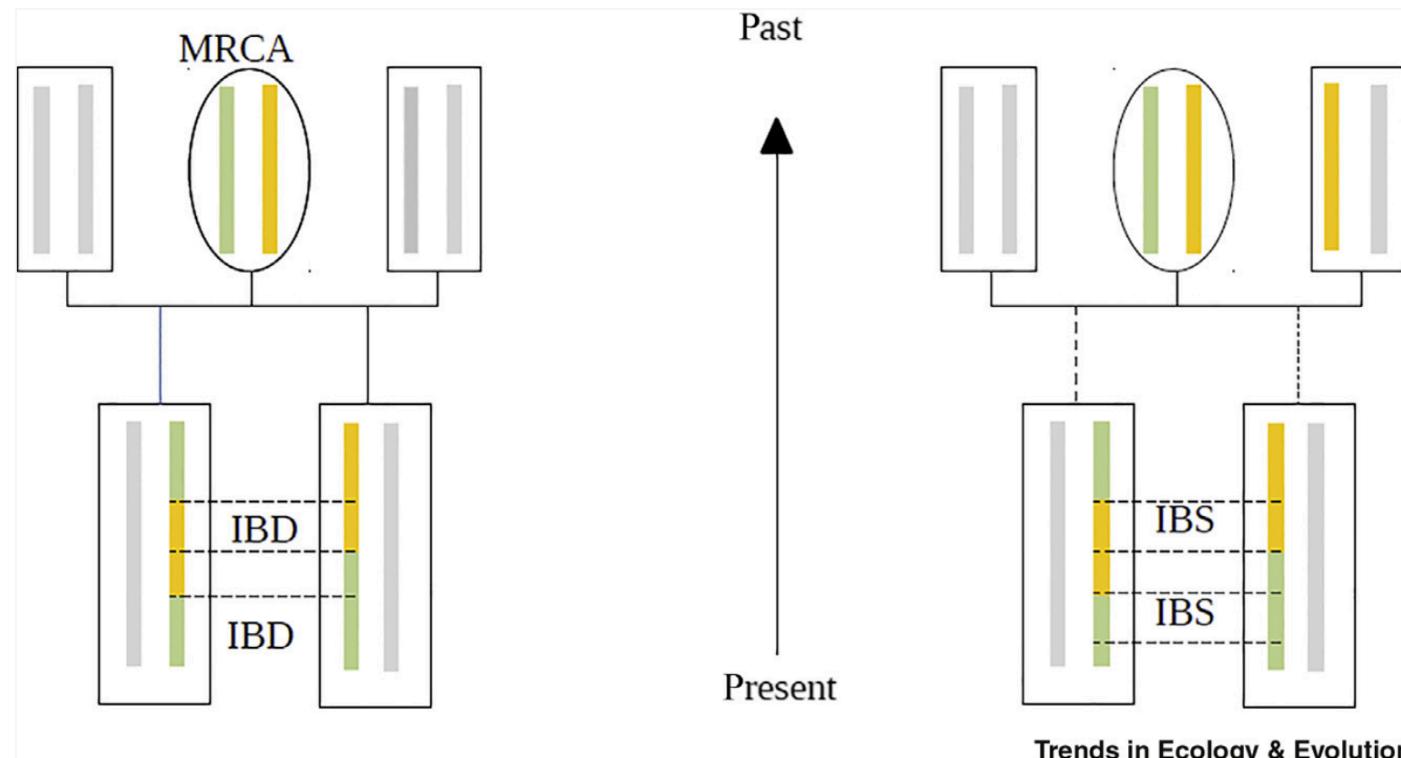
## Multipoint analysis - Non-parametric<sup>7</sup>

- Doesn't assume the mode of inheritance so it is suitable for complex traits
- **Assumptions:**
  - The trait has some genetic component.
  - The genes involved follow Mendelian laws of inheritance
- Significantly less than the power of parametric analysis
- **Use case:** Complex diseases with completely unknown inheritance model

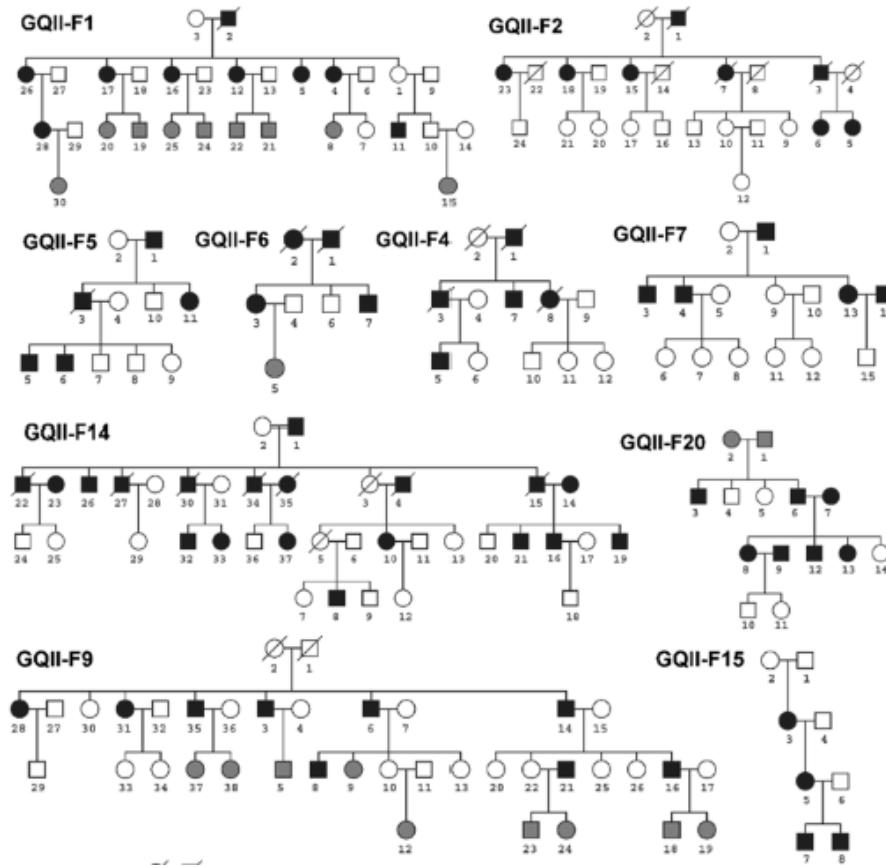
# Multipoint analysis - Non-parametric<sup>7,8</sup>

Quantifies the degree to which related individuals "share" alleles at the marker loci

1. Identity by state (IBS)
2. Identity by descent (IBD)



# Genome-Wide Linkage Analysis



## SCIENTIFIC REPORTS

OPEN

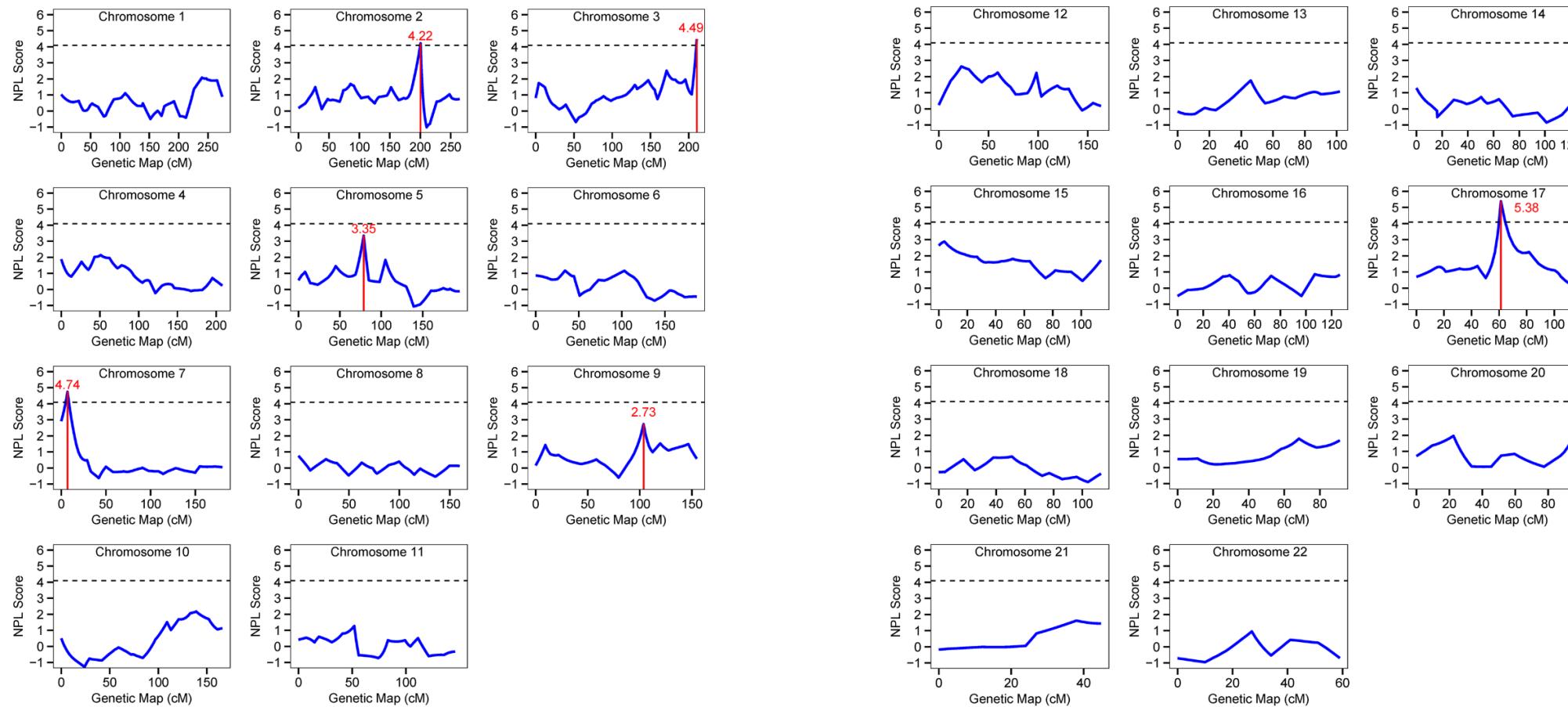
### Genome-Wide Linkage Analysis of Large Multiple Multigenerational Families Identifies Novel Genetic Loci for Coronary Artery Disease

Yang Guo<sup>1,2</sup>, Fan Wang<sup>1,2</sup>, Lin Li<sup>1,2</sup>, Hanxiang Gao<sup>1,3</sup>, Stephen Arkacki<sup>1,2</sup>, Isabel Z. Wang<sup>1,4</sup>, John Barnard<sup>1,5</sup>, Stephen Ellis<sup>6</sup>, Carlos Hubbard<sup>6</sup>, Eric J. Topol<sup>7</sup>, Qiuyun Chen<sup>1,2</sup> & Qing K. Wang<sup>1,2,6,8</sup>

- 24 large and multigenerational families with 433 family members
- All family members were genotyped with markers spaced by every 10 cM

# Genome-Wide Linkage Analysis

- The threshold of significant linkage is  $NPL \geq 4.08$



# Genome-Wide Linkage Analysis

CAD Locus (Genetic Map <sup>a</sup> )	Genomic Region <sup>b</sup>	RefSeq Genes <sup>c</sup>	Genes Related to Cardiovascular diseases <sup>d</sup>
17q21.2 (56.9–83.1 cM)	34.40–57.50 Mb	514	<i>CCL3, CCL4, CCL3L3, CCL4L1, CCL4L2, CCL3L1, TADA2A, PLXDC1, TCAP, PNMT, PGAP3, ERBB2, IKZF3, CSF3, MED24, THRA, NR1D1, CCR7, KRT12, KRT20, GAST, HAP1, JUP, FKBP10, CNP, KCNH4, HCRT, STAT5B, STAT5A, STAT3, ATP6V0A1, MLX, RAMP2, WNK4, BECN1, AOC3, BRCA1, SOST, PYY, G6PC3, HDAC5, GRN, ITGA2B, FZD2, ADAM11, GJC1, CCDC103, GFAP, HEXIM1, MAP3K14, CRHR1, MAPT, WNT3, GOSR2, MYL4, ITGB3, MRPL10, PNPO, MIR10A, UBE2Z, GIP, IGF2BP1, B4GALNT2, ZNF652, NGFR, ITGA3, PDK2, SGCA, COL1A1, XYLT2, CACNA1G, LUC7L3, NME1, MMD, AKAP1, MPO, MIR142</i>
7p22.2 (1.4 –13.0 cM)	0.88–7.25 Mb	87	<i>GPER1, MAFK, NUDT1, GNA12, SDK1, ACTB, AIMP2, EIF2AK1, RAC1</i>
2q33.3 191.9–202.4 cM	177.33–192.47 Mb	88	<i>NFE2L2, PDE11A, RBM45, TTN, CCDC141, ZNF385B, ITGA4, NEUROD1, PDE1A, FRZB, ITGAV, ZSWIM2, CALCRL, TFPI, COL3A1, COL5A2, SLC40A1, PMS1, MSTN, STAT1, STAT4</i>
3q29 (206.5–216.0 cM)	188.96–193.86 Mb	31	<i>TP63, CLDN16, UTS2B, HRASLS, OPA1</i>
5q13.2 (74.9–80.0 cM)	66.68–71.63 Mb	39	<i>PIK3R1, CCNB1, OCLN, SMN2, SMN1, NAIP, MCCC2, CARTPT</i>
9q22.33(103.6–105.3 cM)	101.72–104.22 Mb	20	<i>TGFBR1, NR4A3, INVS</i>

## Applications of Linkage Analysis in Clinical Genomics

- Historically used in family-based studies to identify Mendelian disease genes.
- Previously required sequencing candidate regions, but now NGS allows rapid follow-up of suggestive linkage signals.

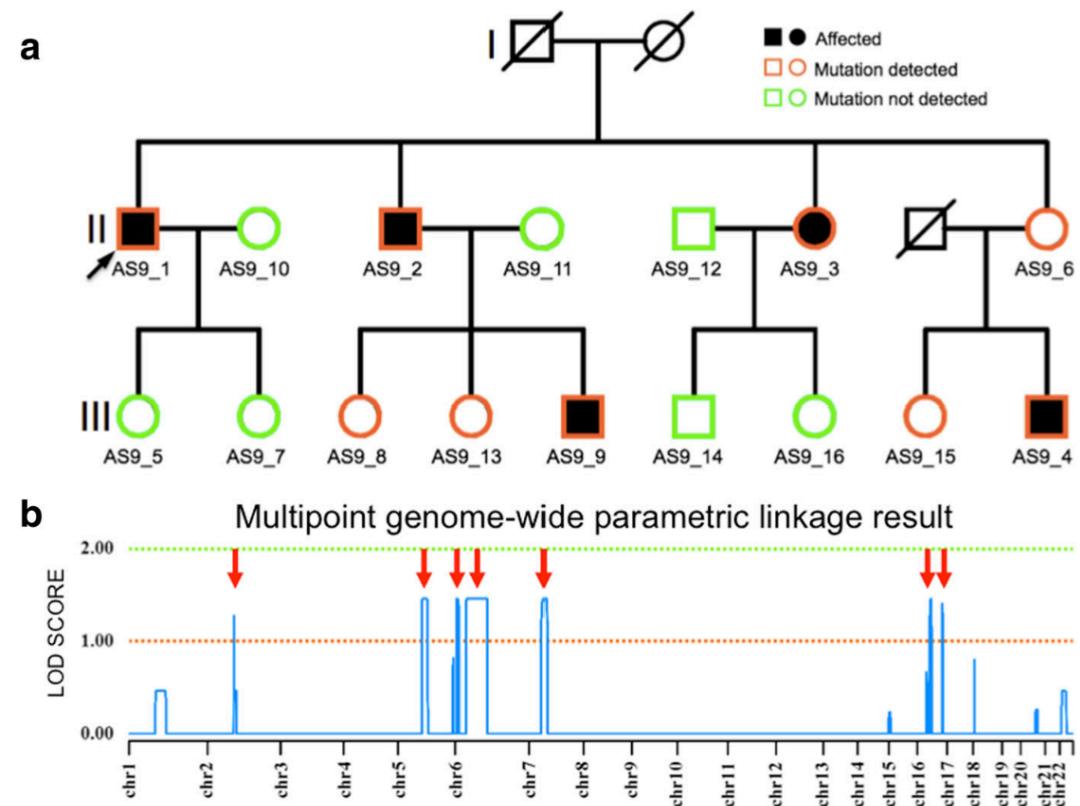
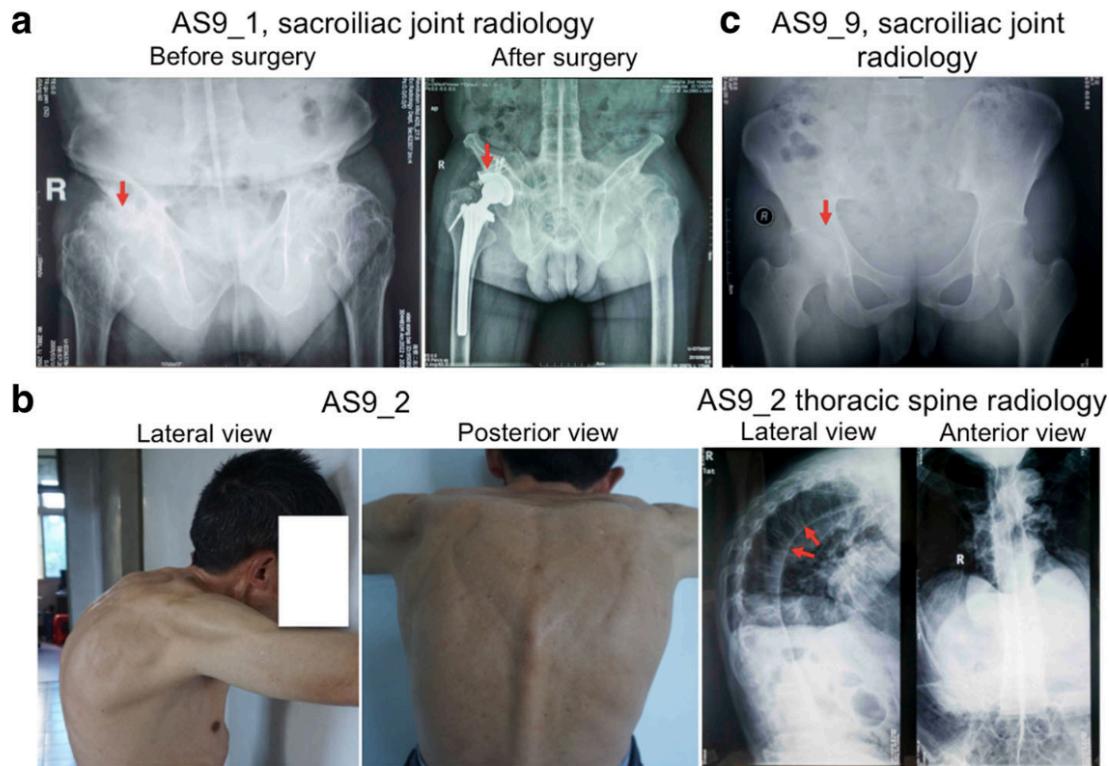


# Identification of *ANKDD1B* variants in an ankylosing spondylitis pedigree and a sporadic patient

Zhiping Tan<sup>1,2\*</sup>, Hui Zeng<sup>1,2</sup>, Zhaofa Xu<sup>3</sup>, Qi Tian<sup>3</sup>, Xiaoyang Gao<sup>3</sup>, Chuanman Zhou<sup>3</sup>, Yu Zheng<sup>3</sup>, Jian Wang<sup>1,2</sup>, Guanghui Ling<sup>4</sup>, Bing Wang<sup>5</sup>, Yifeng Yang<sup>1,2</sup> and Long Ma<sup>3\*</sup> 

**Background:** Ankylosing spondylitis (AS) is a debilitating autoimmune disease affecting tens of millions of people in the world. The genetics of AS is unclear. Analysis of rare AS pedigrees might facilitate our understanding of AS pathogenesis.

**Methods:** We used genome-wide linkage analysis and whole-exome sequencing in combination with variant co-segregation verification and haplotype analysis to study an AS pedigree and a sporadic AS patient.



# Population Genetics

The quantitative study of the distribution of genetic variation in populations and of how the frequencies of genes and genotypes are maintained or change over time both within and between populations.

- HWE and Allele Frequency
- Genome Wide Association Studies

# Hardy-Weinberg Equilibrium

The genetic variation in a population will remain constant from one generation to the next in the absence of disturbing factors.

Assumes:

- Large population
- Random mating
- Allele frequencies remain constant over time because there are no mutations, natural selection, nonrandom mating, genetic drift, and gene flow.

# Hardy-Weinberg Equilibrium

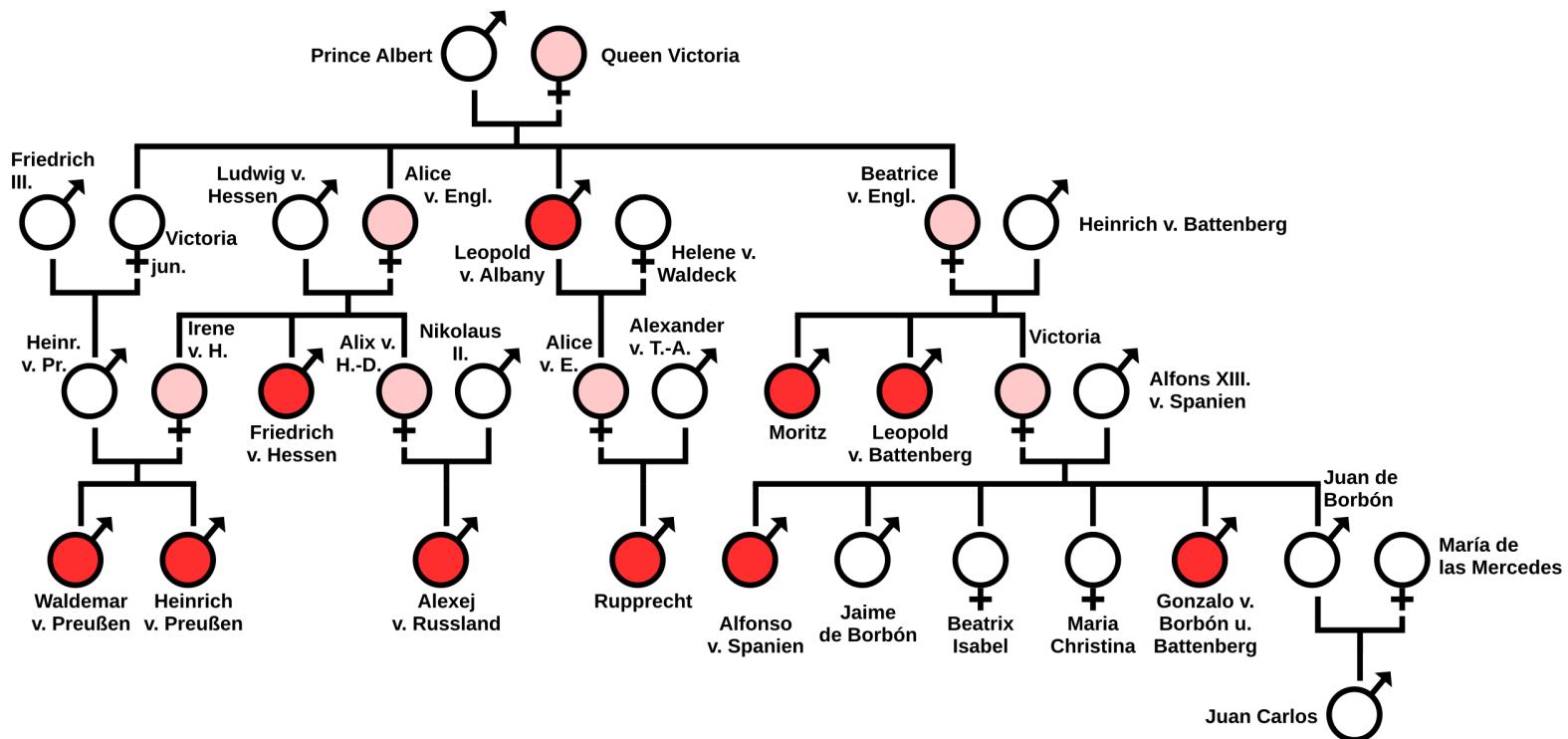
$$p^2 + 2pq + q^2 = 1$$

where:

- $p$  and  $q$  are **allele frequencies**.
- $p^2$ ,  $2pq$ , and  $q^2$  are **genotype frequencies**.

# Factors That Disturb HWE: Nonrandom Mating

In human populations, nonrandom mating may occur due to: stratification, assortative mating, consanguinity.

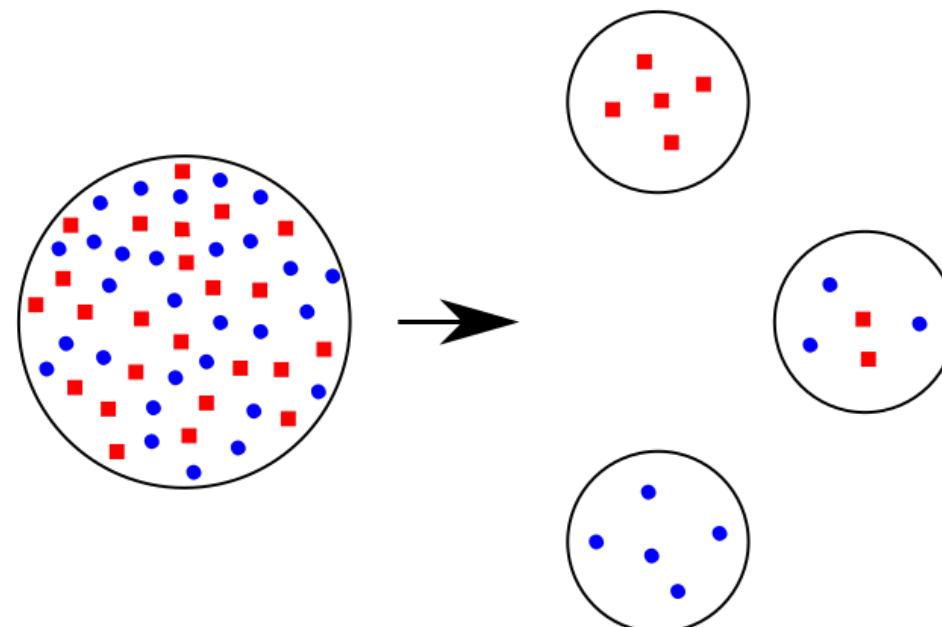


## Factors That Disturb HWE: Mutation

- Selection pressure on mutations determines the allele frequencies in the population.
- Fitness is the outcome of the joint effects of survival and fertility.
- When a genetic disorder limits reproduction so severely that the fitness is zero (i.e.,  $s = 1$ ), it is thus referred to as a **genetic lethal**.

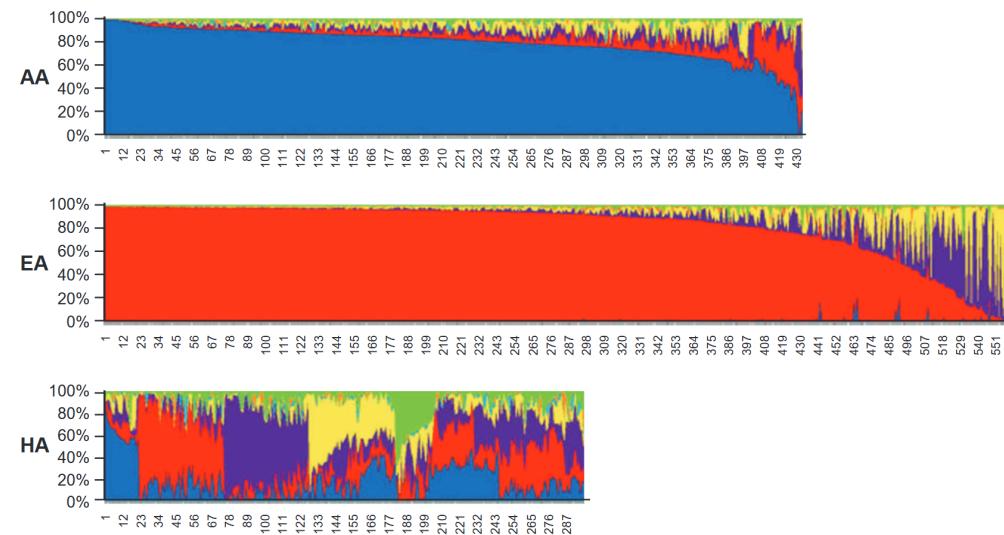
# Factors That Disturb HWE: Genetic Drift

- Random effects of environment or other chance occurrences can change the frequency of the disease allele when the population is small.
- **Founder effect:** the loss of genetic variation that occurs when a new population is established by a very small number of individuals from a larger population



# Factors That Disturb HWE: Migration and Gene Flow

- Migration can change allele frequency by the process of **gene flow** (the slow diffusion of genes across a barrier).
- **Genetic admixture:** the genes of migrant populations with their own allele frequencies are gradually merged into the gene pool of the population into which they have migrated.





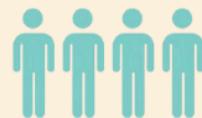
# Multi-ancestry genome-wide association meta-analysis of Parkinson's disease

## Study participants

Goal: Collate the largest and most diverse set of participants in Parkinson's disease genomics



European  
39,275 cases  
18,618 proxy cases  
1.5M controls



East Asian  
7,046 cases  
176,756 controls



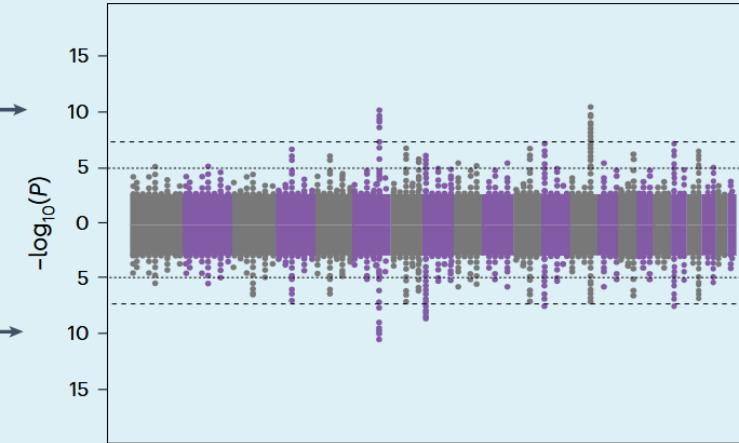
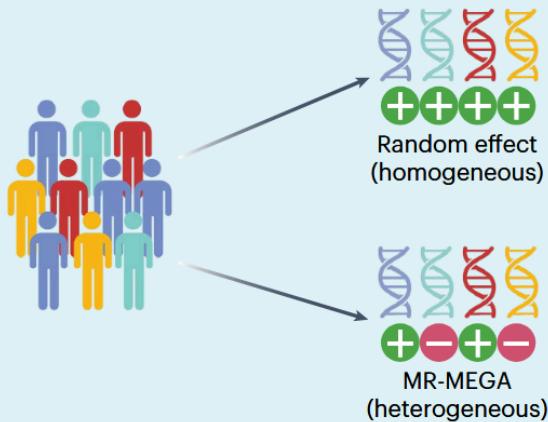
Latino  
2,440 cases  
582,220 controls



African  
288 cases  
193,985 controls

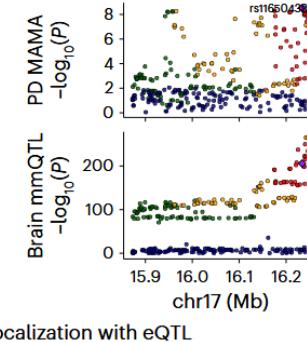
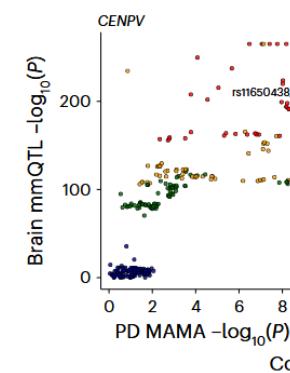
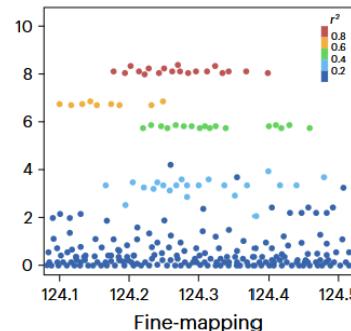
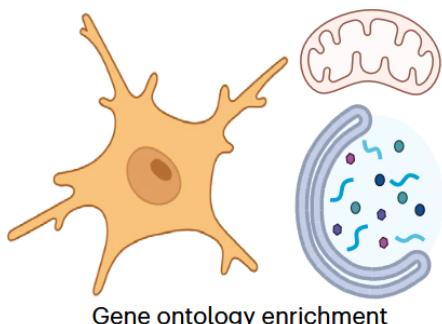
## Multiancestry genome-wide meta-analysis

Goal: Identify common SNPs that are associated with Parkinson's disease risk that are applicable across different ancestries

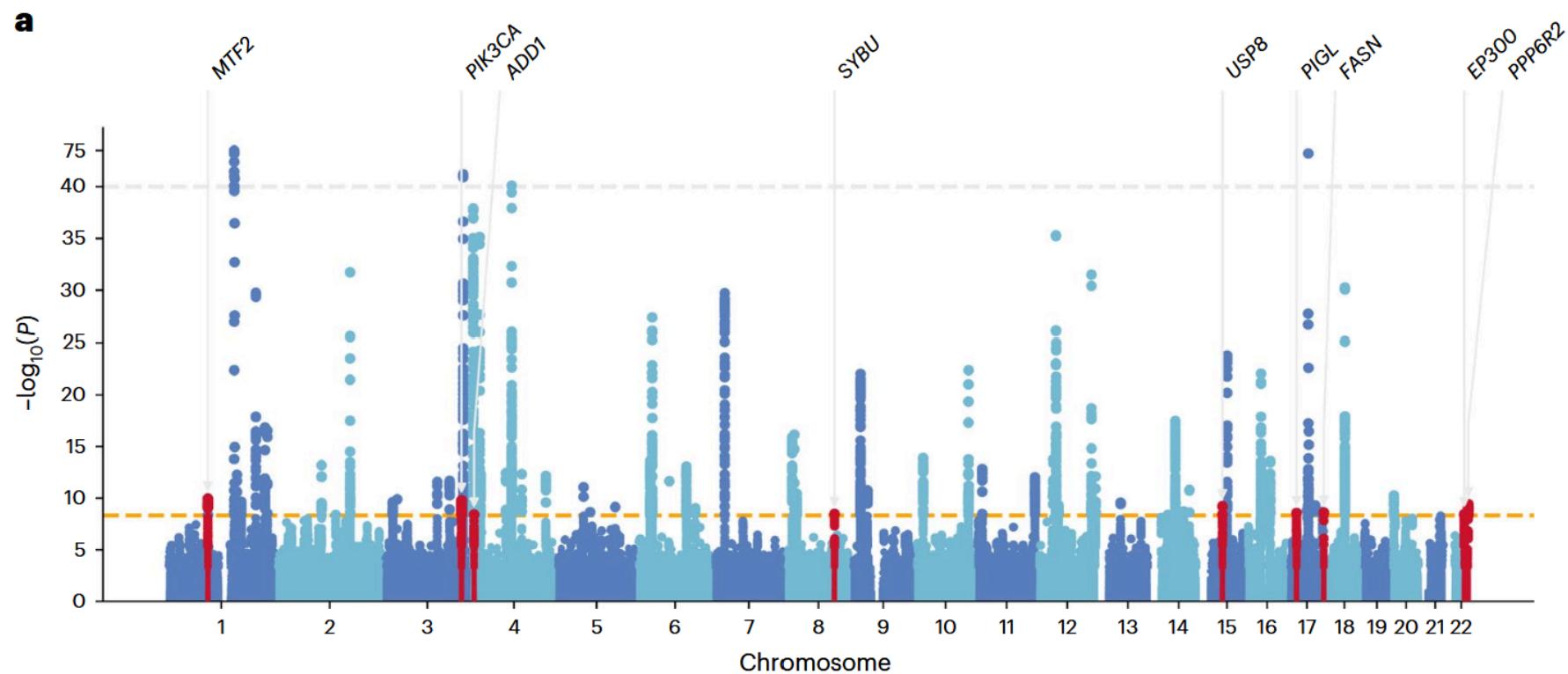


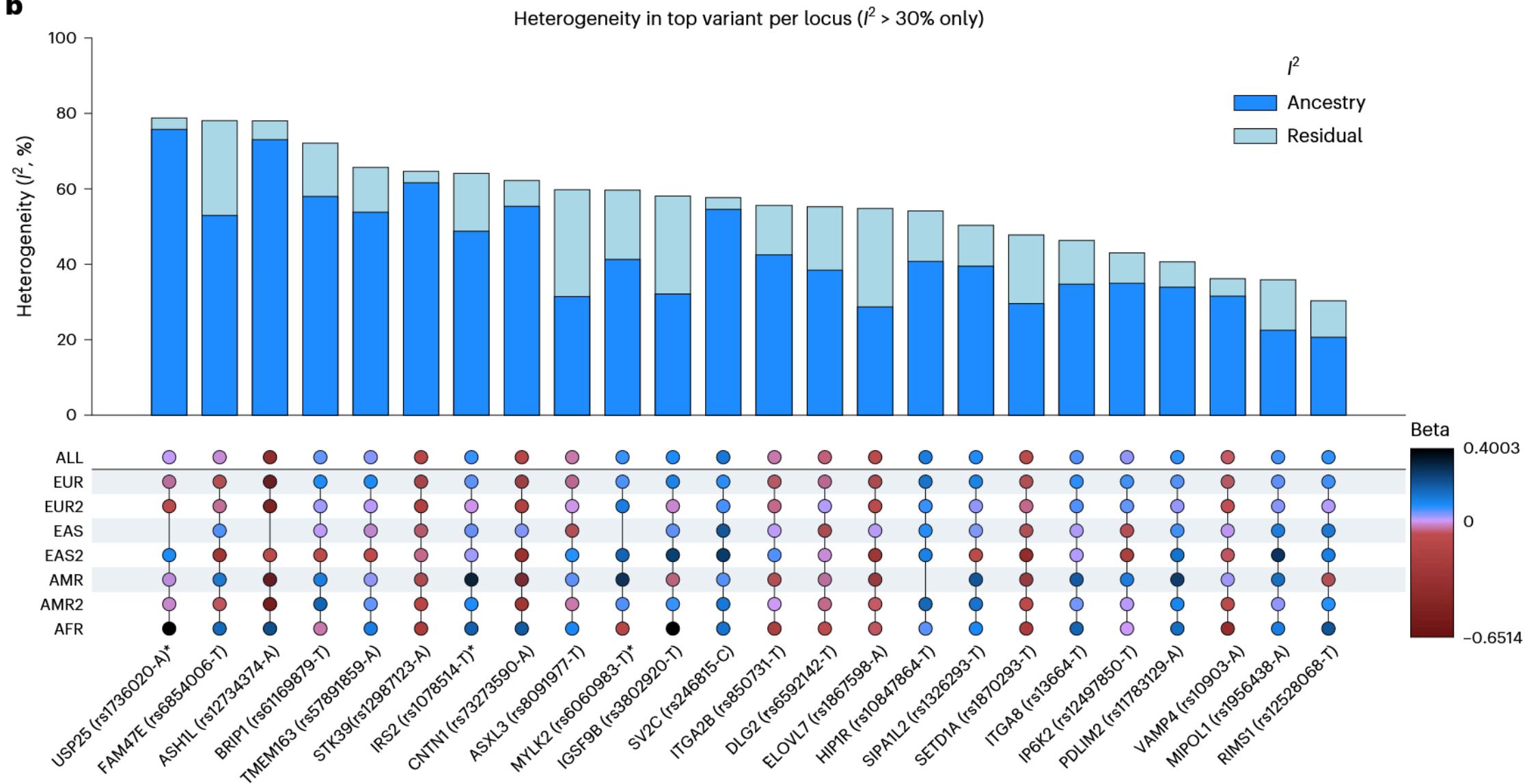
## Downstream analyses

Goal: Interpret the meta-analysis results and identify potential targets and biological mechanisms



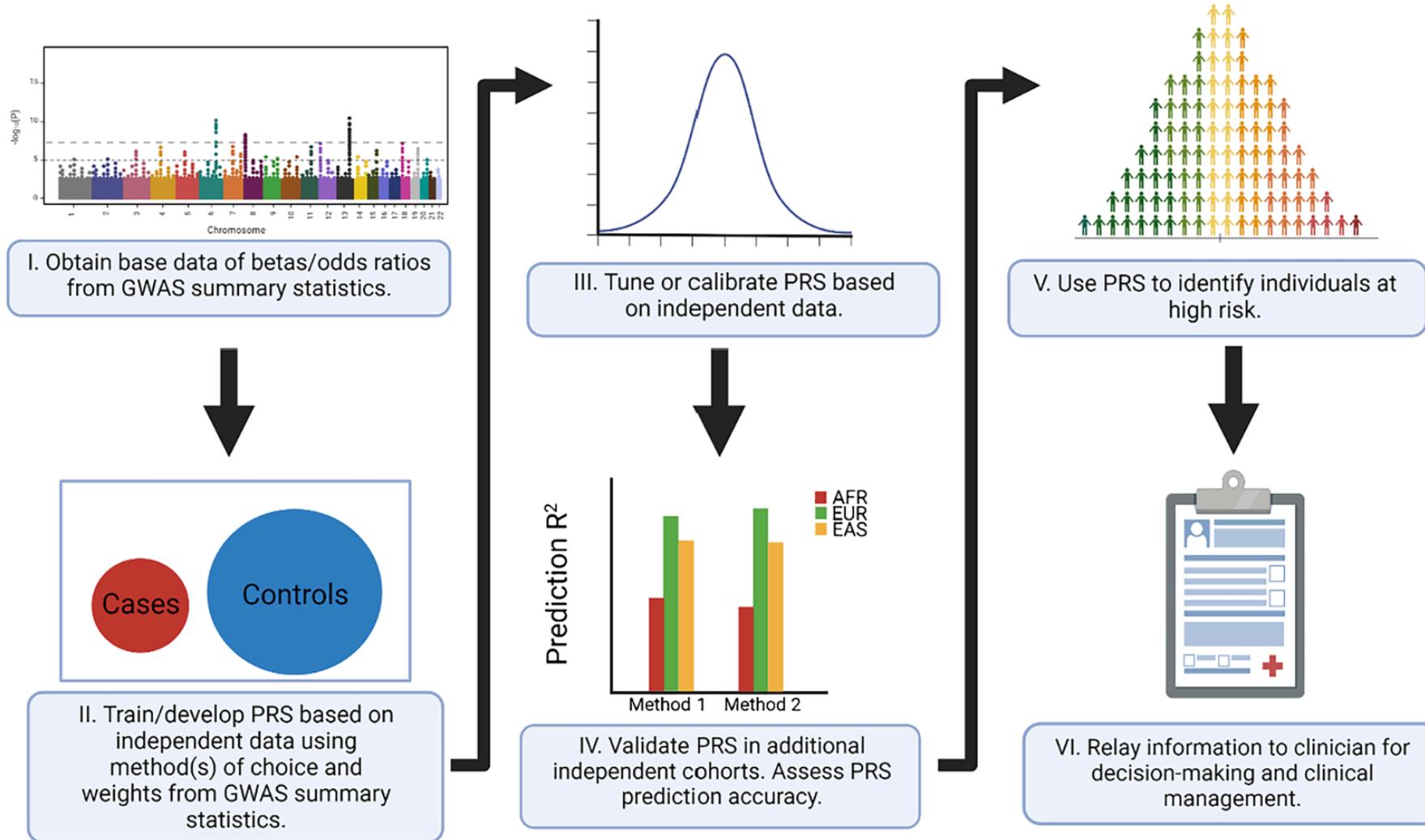
# Manhattan Plot



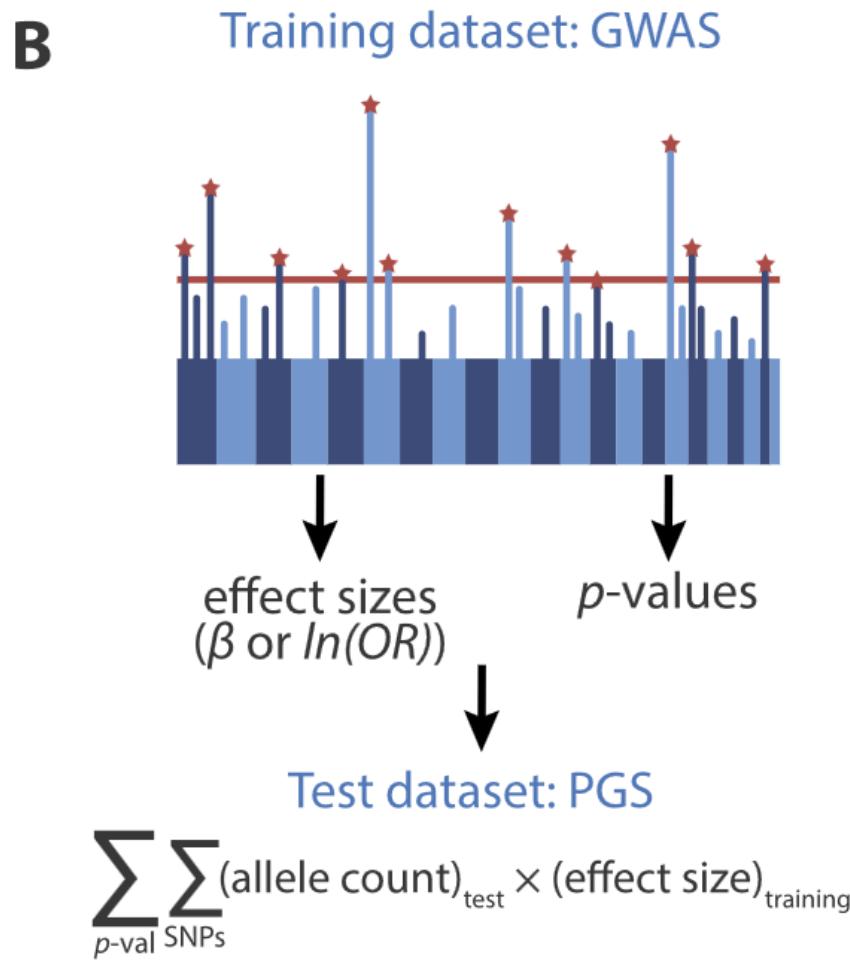
**b**

Locus	Number of significant SNPs	Nominated variant	CHR:BP:A1:A2	Nearest gene	Known PD gene ±1 MB	Functional consequence
11	6	<a href="#">rs57891859</a>	2:135464616:A:G	<i>TMEM163</i>	<i>TMEM163</i>	intronic
19	926	<a href="#">rs34311866</a>	4:951947:C:T	<i>TMEM175</i>	<i>TMEM175</i>	exonic
23	1483	<a href="#">rs356182</a>	4:90626111:A:G	<i>SNCA</i>	<i>SNCA</i>	ncRNA intronic
24	121	<a href="#">rs13117519</a>	4:114369065:T:C	<i>CAMK2D</i>	<i>CAMK2D</i>	intergenic
45	1371	<a href="#">rs10847864</a>	12:123326598:G:T	<i>HIP1R</i>	<i>HIP1R</i>	intronic
60	1	<a href="#">rs55818311</a>	19:2341047:C:T	<i>SPPL2B</i>	<i>LSM7</i>	ncRNA exonic

# Clinical Application of GWAS



# Polygenic Risk Score



- Sum of the weighted allele counts of independent SNPs found to be associated with a trait or disorder in a GWAS
- For quantitative traits, the weights are the linear regression effect sizes, i. e., the  $\beta$  coefficients.
- For case/control phenotypes, the weights are the natural logarithm of the odds ratio,  $\ln(OR)$ .