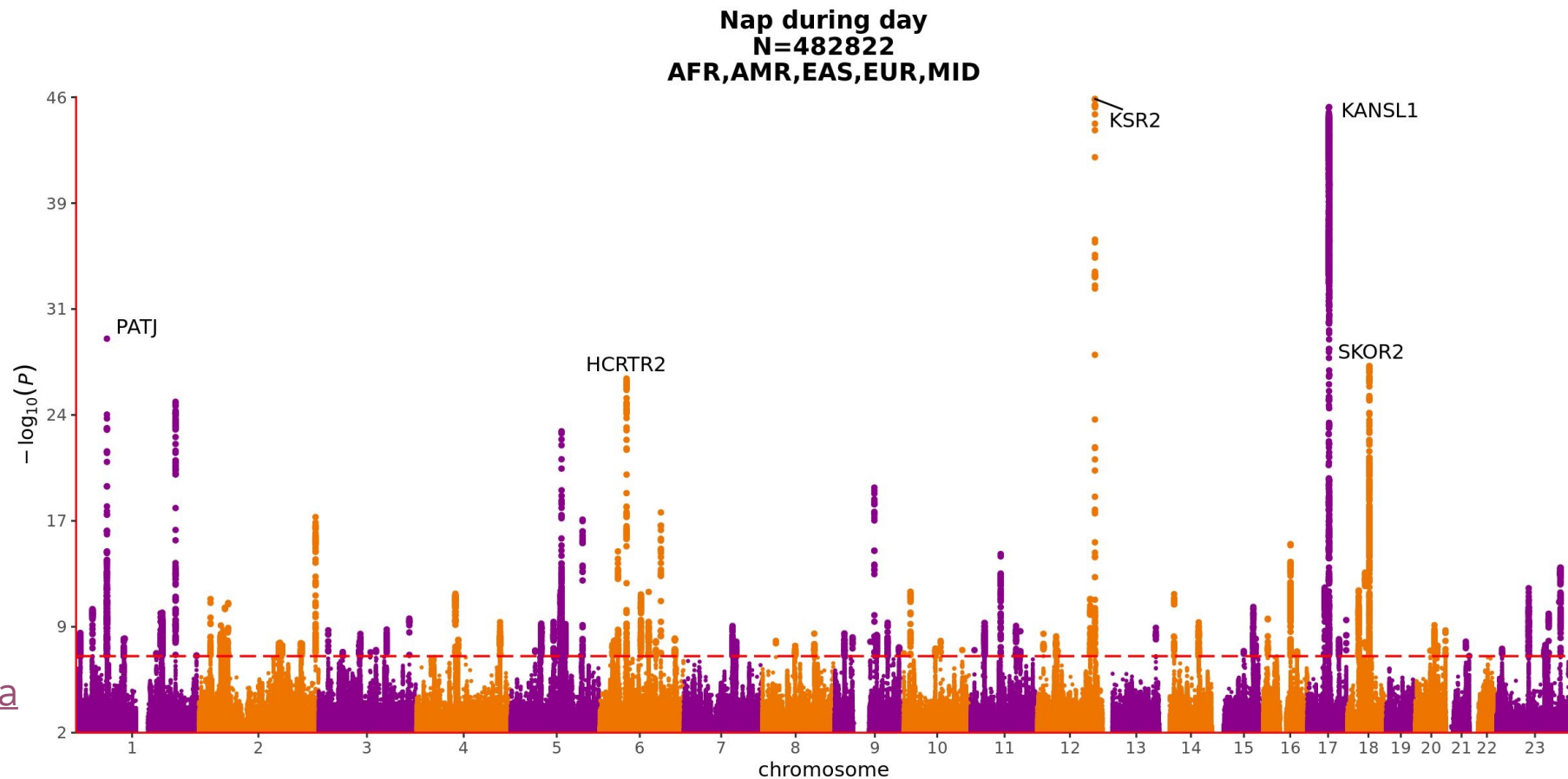


Intro to GWAS

BIOL 435/535: Bioinformatics
Feb 24th, 2022



Gene discovery and polygenic prediction from a genome-wide association study of educational attainment in 1.1 million individuals

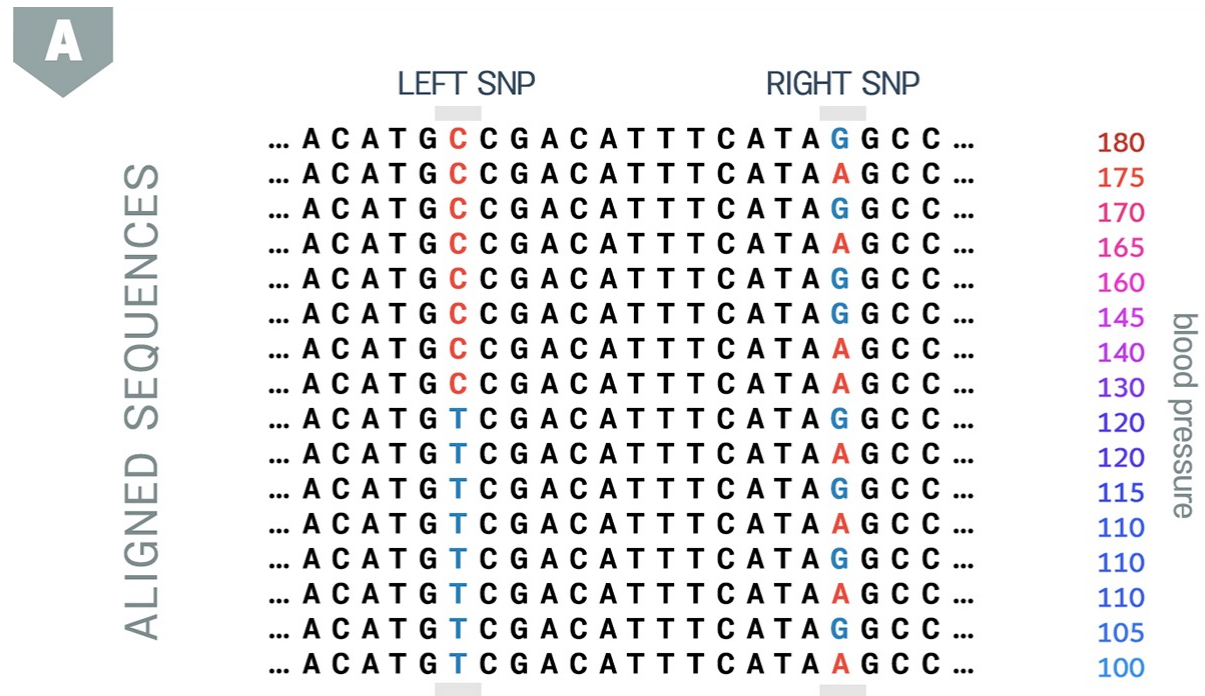
James J. Lee^{1,58}, Robbee Wedow^{2,3,4,58}, Aysu Okbay^{5,6,58*}, Edward Kong⁷, Omeed Maghzian⁷, Meghan Zacher⁸, Tuan Anh Nguyen-Viet⁹, Peter Bowers⁷, Julia Sidorenko^{10,11}, Richard Karlsson Linnér^{5,6,12}, Mark Alan Fontana^{9,12}, Tushar Kundu⁹, Chanwook Lee⁷, Hui Li⁷, Ruoxi Li⁹, Rebecca Royer⁹, Pascal N. Timshel^{14,15}, Raymond K. Walters^{16,17}, Emily A. Willoughby¹, Loïc Yengo^{10, 23} and Me Research Team¹⁸, COGENT (Cognitive Genomics Consortium)¹⁹, Social Science Genetic Association Consortium¹⁸, Maris Alver¹¹, Yanchun Bao²⁰, David W. Clark²¹, Felix R. Day²², Nicholas A. Furlotte²², Peter K. Joshi^{21,24}, Kathryn E. Kemper¹⁰, Aaron Kleinman²², Claudia Langenberg²², Reedik Mägi¹¹, Joey W. Trampush^{25,26}, Shefali Setia Verma²⁷, Yang Wu¹⁰, Max Lam^{28,29}, Jing Hua Zhao²², Zhili Zheng^{10,30}, Jason D. Boardman^{2,3,4}, Harry Campbell²¹, Jeremy Freese²¹, Kathleen Mullan Harris^{22,23}, Caroline Hayward²⁴, Pamela Herd^{20,25}, Meena Kumari²⁰, Todd Lencz^{26,27,28}, Jian'an Luan²², Anil K. Malhotra^{26,27,28}, Andres Metspalu^{11,29}, Lili Milani¹¹, Ken K. Ong²², John R. B. Perry²², David J. Porteous⁴⁰, Marylyn D. Ritchie²⁷, Melissa C. Smart²¹, Blair H. Smith^{41,42}, Joyce Y. Tung²³, Nicholas J. Wareham²², James F. Wilson^{21,24}, Jonathan P. Beauchamp⁴³, Dalton C. Conley⁴⁴, Tõnu Esko¹¹, Steven F. Lehrer^{45,46,47}, Patrik K. E. Magnusson⁴⁸, Sven Oskarsson⁴⁹, Tune H. Pers^{14,15}, Matthew R. Robinson^{10,50}, Kevin Thom⁵¹, Chelsea Watson⁹, Christopher F. Chabris⁵², Michelle N. Meyer⁵³, David I. Laibson⁷, Jian Yang^{10,54}, Magnus Johannesson⁵⁵, Philipp D. Koellinger^{5,6,12}, Patrick Turley^{16,17,59}, Peter M. Visscher^{10,54,59*}, Daniel J. Benjamin^{9,43,56,59*} and David Cesarini^{47,51,57,59}

Genome-wide association studies

- Identify SNPs that are associated with phenotype(s) of interest

H_0 : Phenotype \perp SNP

H_A : Phenotype \sim SNP

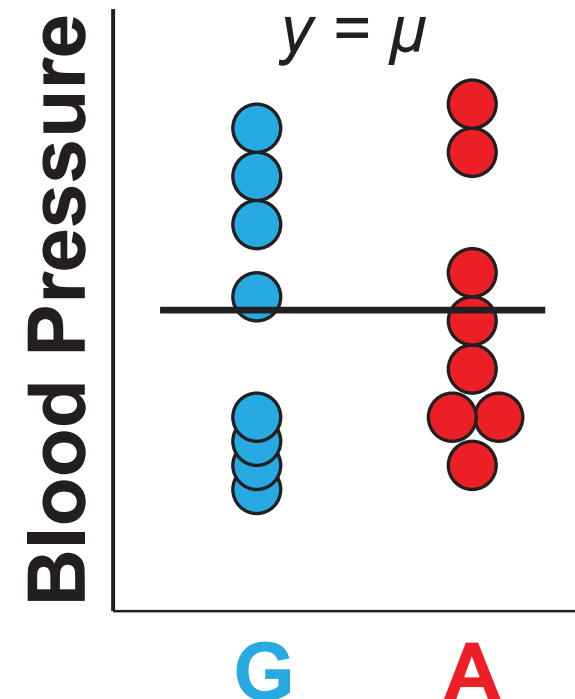


Genome-wide association studies

- Identify SNPs that are associated with phenotype(s) of interest

H_0 : Phenotype \perp SNP

H_A : Phenotype \sim SNP

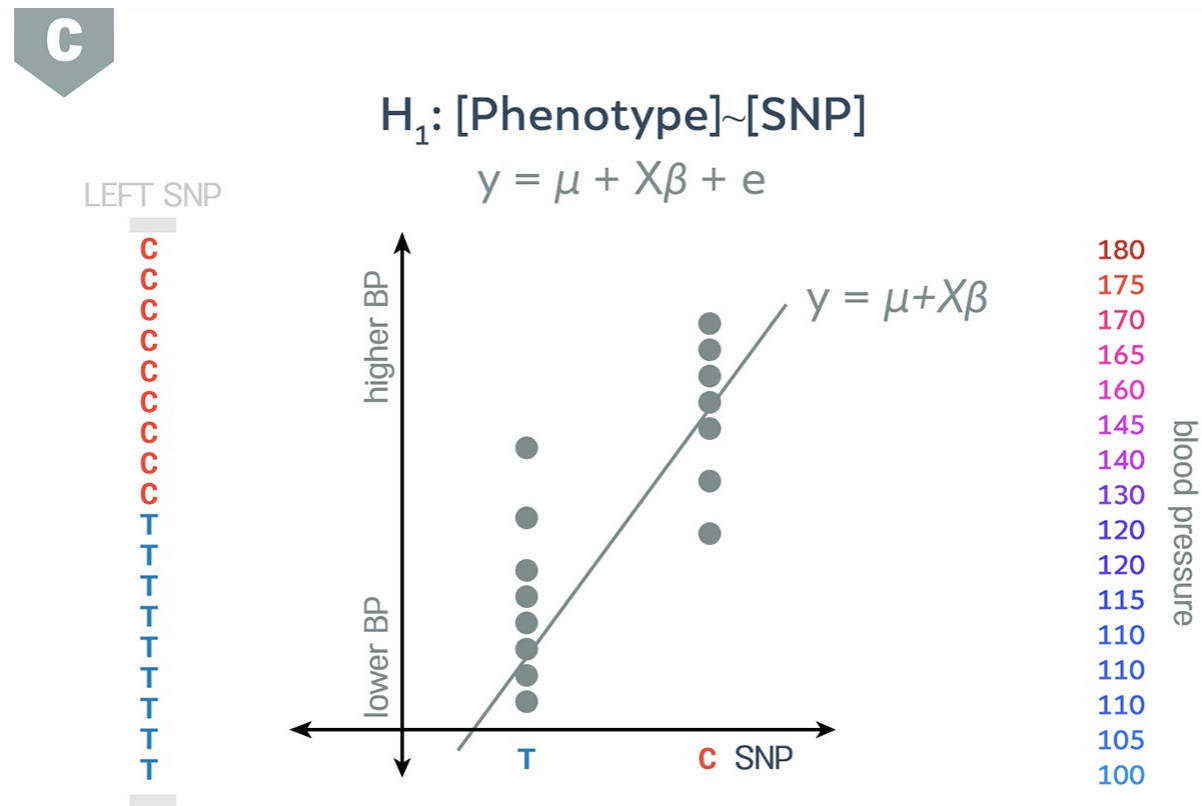


Genome-wide association studies

- Identify SNPs that are associated with phenotype(s) of interest

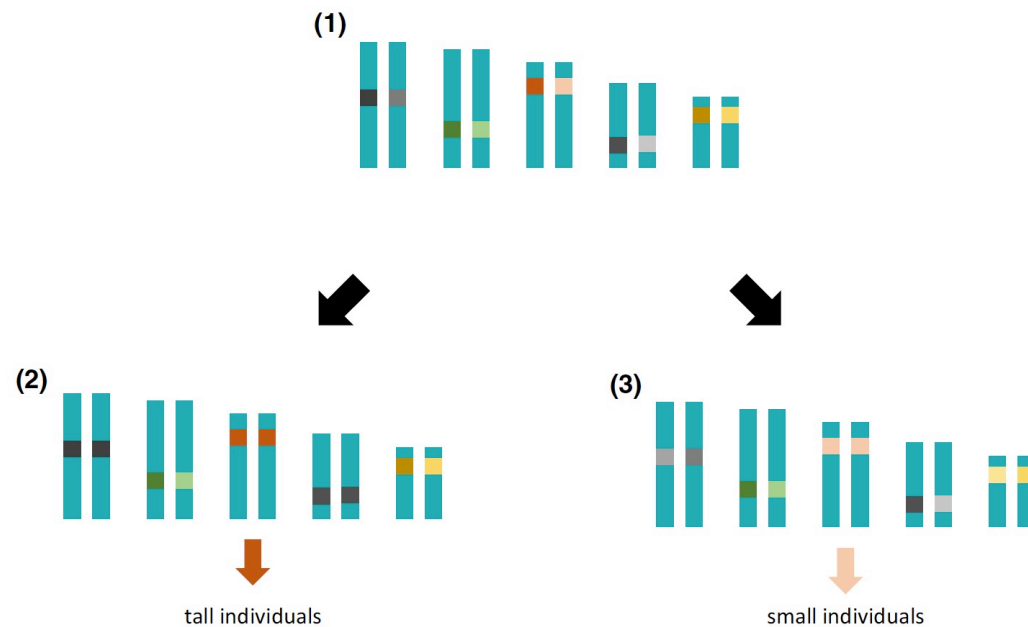
$H_0: \text{Phenotype} \perp \text{SNP}$

H_A : Phenotype \sim SNP



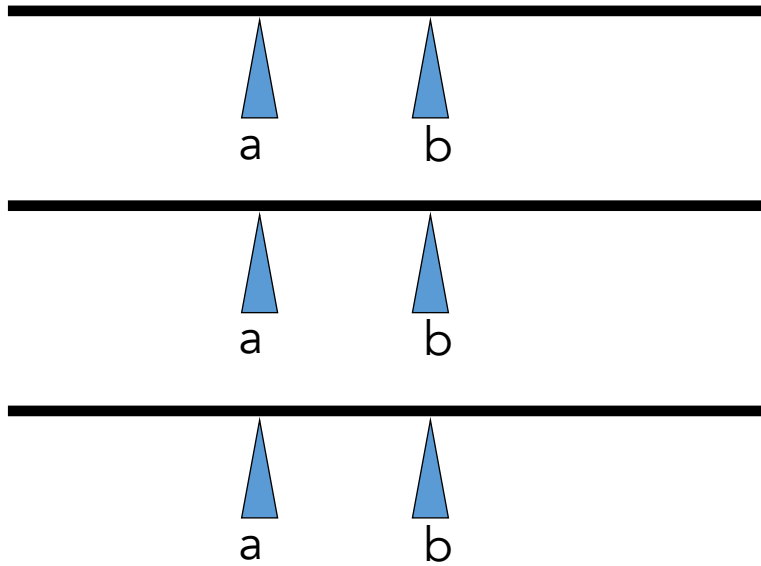
SNP associations may not be due to causal effect

Population structure (differential relatedness among individuals) means that SNPs resulting from **Isolation by Descent** (IBD) will be associated with the phenotype of interest

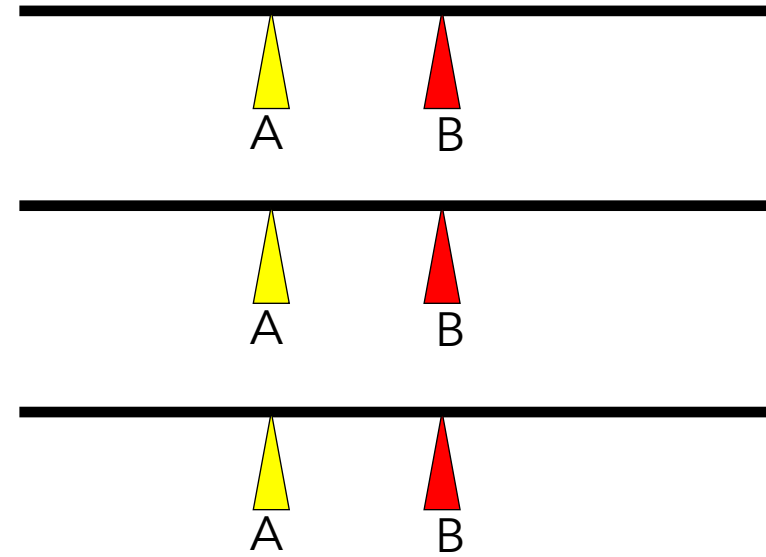


Population structure causes linkage disequilibrium (LD)

Sub-population 1

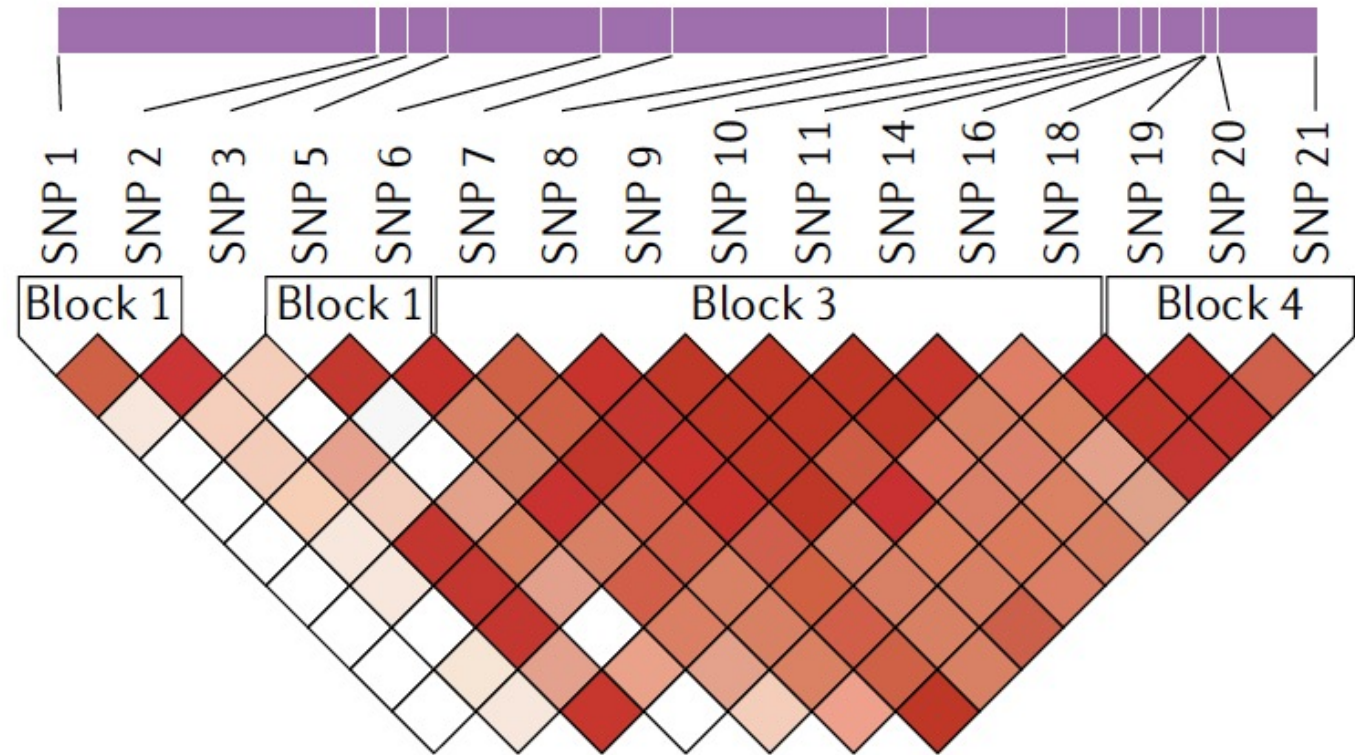


Sub-population 2

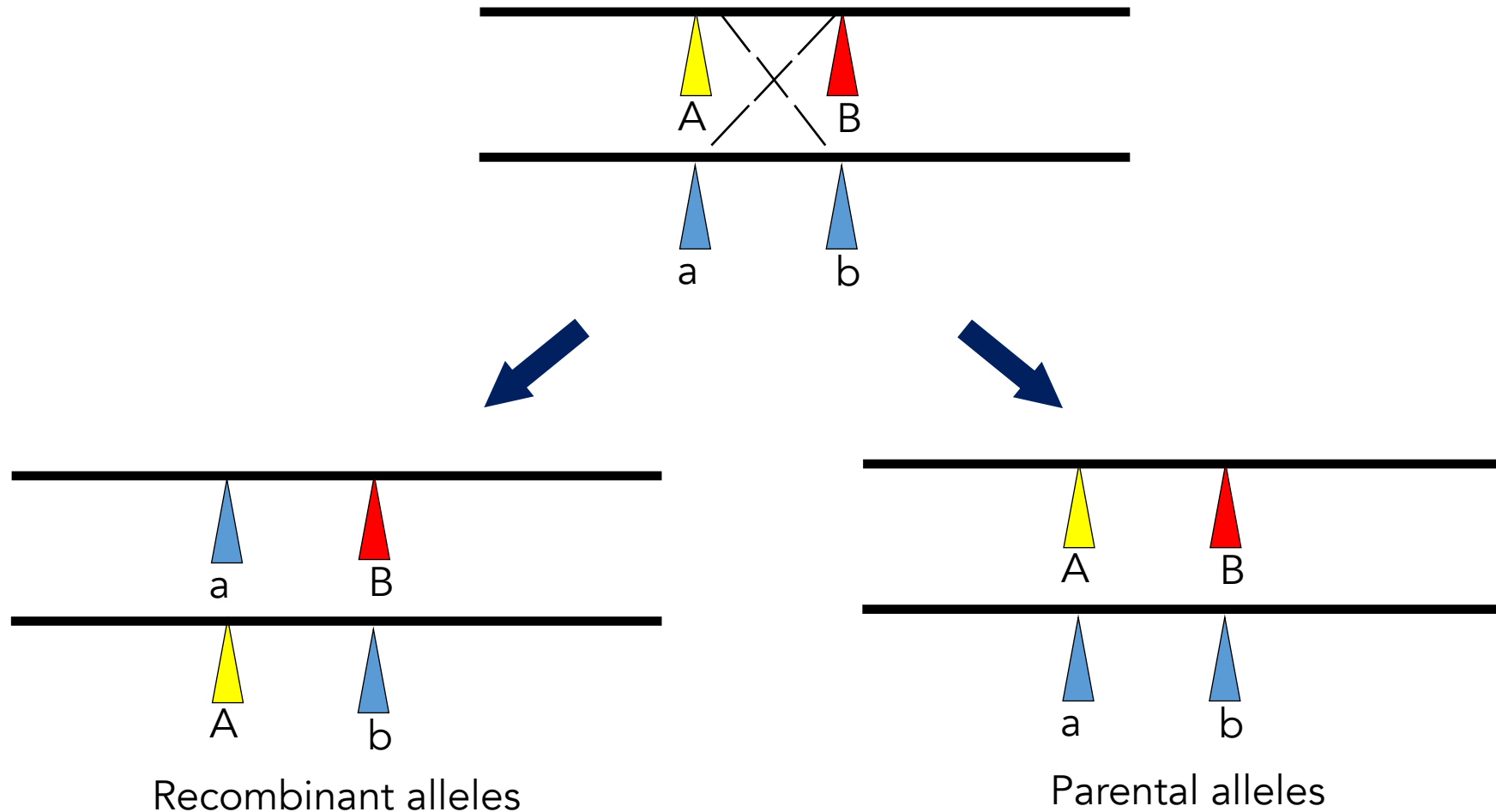


Population structure causes linkage disequilibrium (LD)

Linkage
disequilibrium

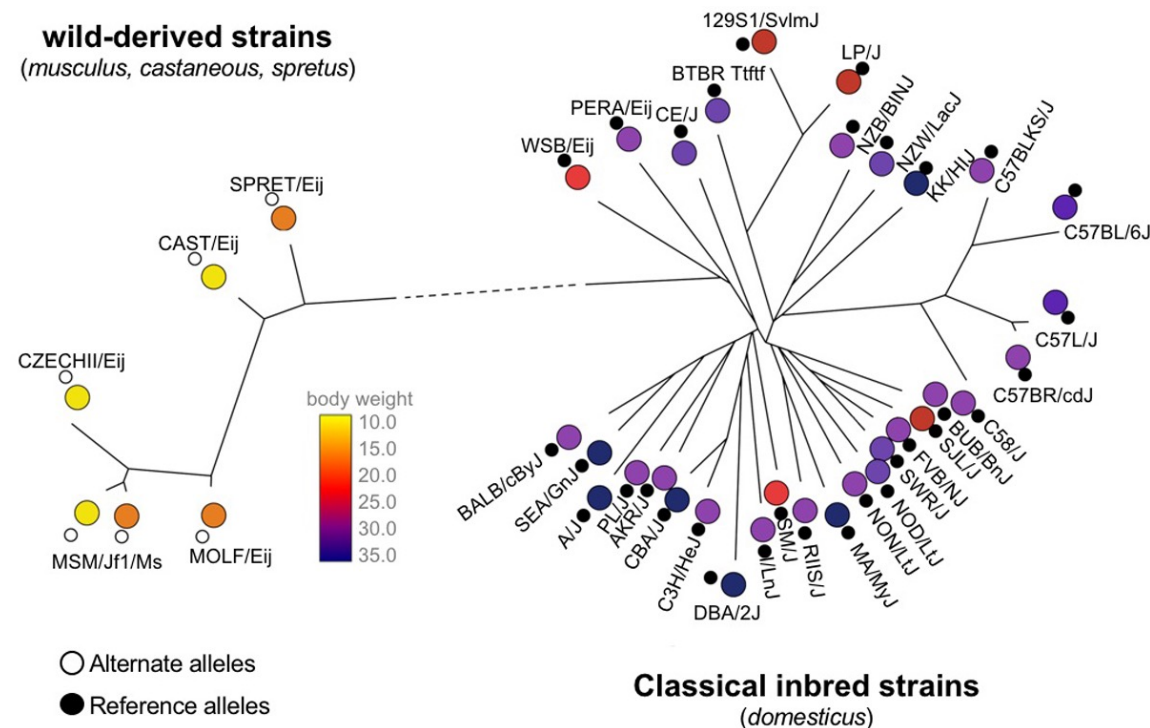


Recombination can break down LD, but it takes time, secondary contact



SNP associations may not be due to causal effect

Population structure (differential relatedness among individuals) means that SNPs resulting from **Isolation by Descent** (IBD) will be associated with the phenotype of interest

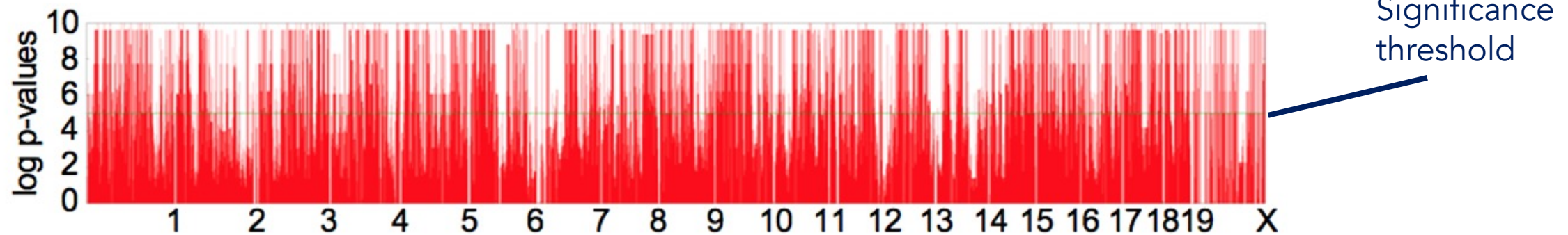


SNP associations may not be due to causal effect

Population structure (differential relatedness among individuals) means that SNPs resulting from **Isolation by Descent** (IBD) will be associated with the phenotype of interest

A

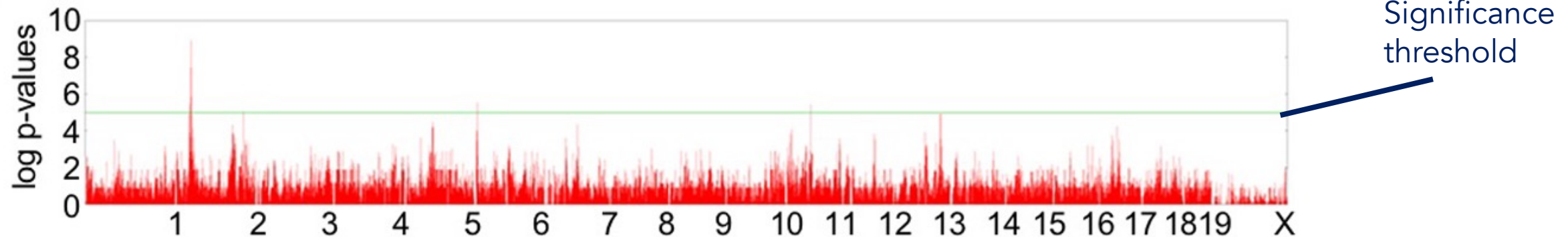
GENOME-WIDE ASSOCIATION MAP



Accounting for population structure can remove spuriously associated SNPs

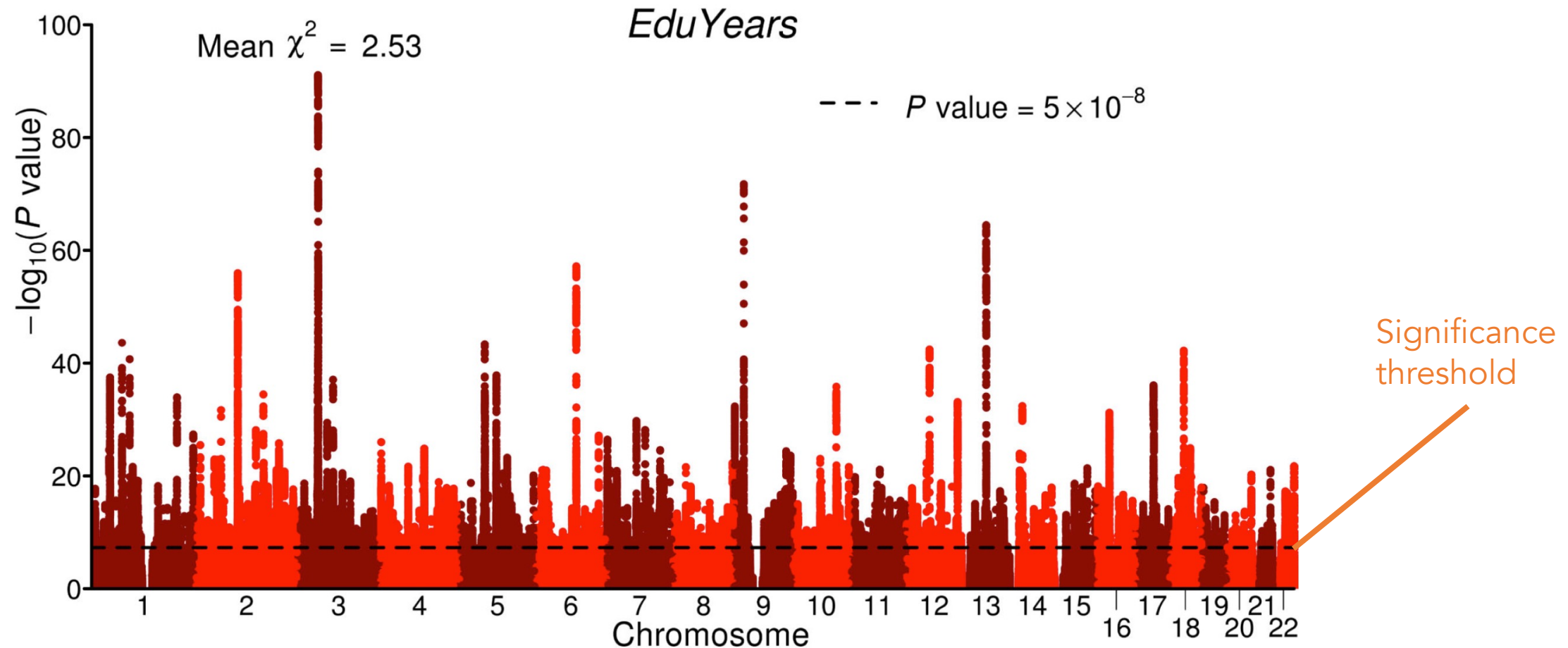
A

GENOME-WIDE ASSOCIATION MAP



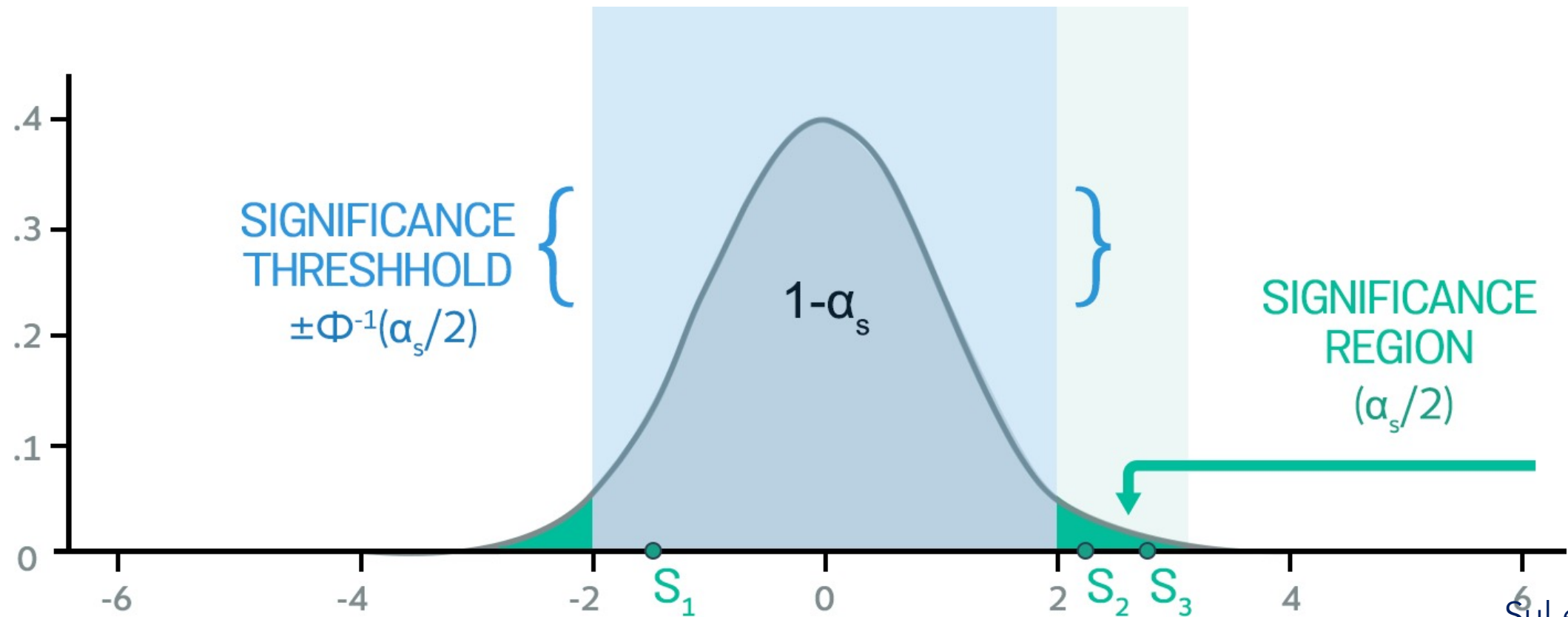
Sul et al 2018

Doing GWAS with only European populations is a problem



Determining Statistical Significance

False Discovery Rate (Benjamini & Hochberg 1996)



SNP databases

dbSNP (NCBI)

OMIM (Online Mendelian Inheritance in Man)

HapMap Project

1000 Genomes Project



Performing GWAS – Linking Genotype to Phenotype

Case-control studies

- Compare SNP frequencies in cases vs. controls

Performing GWAS – Linking Genotype to Phenotype

Case-control studies

Bulk segregant analysis

- Separate phenotypes into bins, compare SNP frequencies across bins

Performing GWAS – Linking Genotype to Phenotype

Case-control studies

Bulk segregant analysis

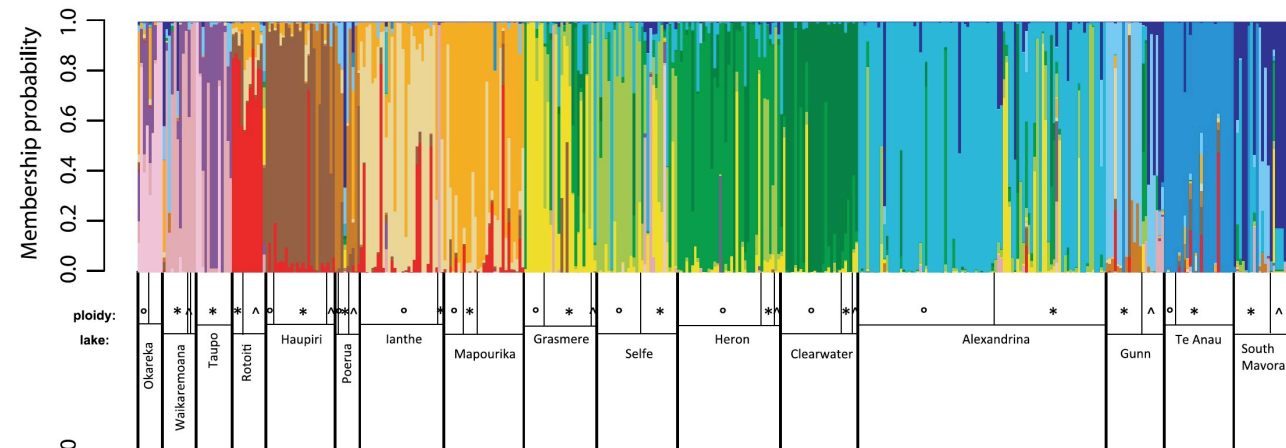
Natural variation

- Multivariate regression of continuous distributions vs. SNP frequencies

Performing GWAS – Key considerations

Problems with the population:

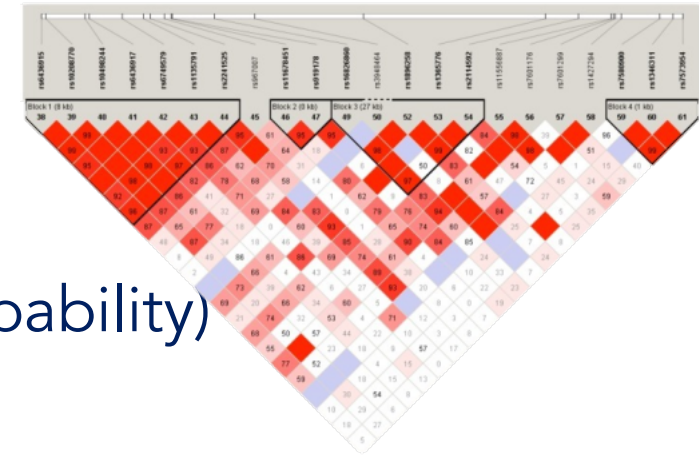
- Hardy-Weinberg Equilibrium
- Population structure
 - individuals non-randomly related to other individuals
- Linkage Disequilibrium (LD)
- Missing data



Performing GWAS – Key considerations

Problems with the population:

- Hardy-Weinberg Equilibrium
- Population structure
- Linkage Disequilibrium (LD) –
Non-random association of two alleles (i.e., > 0.5 probability)
- Missing data



Performing GWAS – Key considerations

Problems with the population:

Problems with the sample:

- Sample size
- Phenotype selection

Performing GWAS – Key considerations

Problems with the population:

Problems with the sample:

Problems with statistical analyses:

- Correcting for multiple comparisons
- False positives
- P-hacking