# Population Stratification

## Human Genomic Epidemiology – Asia
## Virtual Course
## June 13-17, 2022

## Professor Qasim Ayub

qasim.ayub@monash.edu

*Director, Monash University Malaysia Genomics Facility*
*Deputy Head of School (Research)*
*School of Science*

MONASH University | MALAYSIA

Genomics Facility

# Learning Outcomes

➢ Understand modern human population structure.

➢ Interpret principal components and ADMIXTURE analyses.

➢ Describe how population structure confounds genome wide association studies and how to control for it.
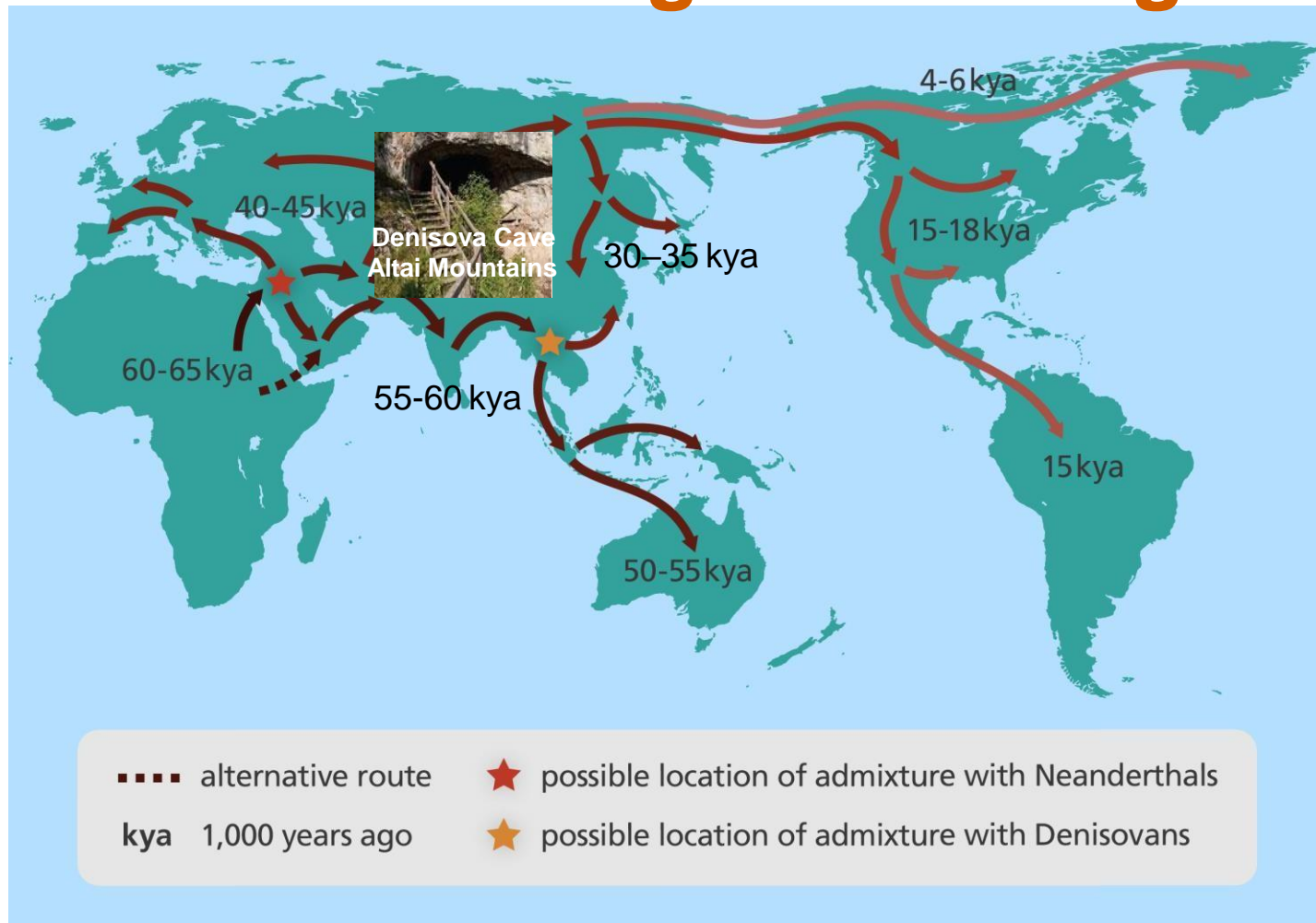
# Outline

➢ Modern human populations and datasets.

➢ Principal Components and ADMIXTURE analysis.

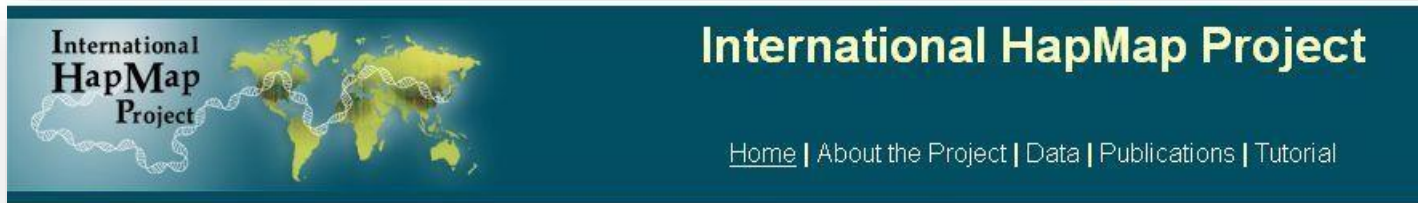➢ Population stratification and GWAS.

# What is a Population?

➢ **Population** is a **spatial-temporal group of interbreeding individuals who share a common gene pool.**

➢ Population genomics aims to understand population structure and relationships.

➢ Population structure is defined by the organization of genetic variation and is driven by the combined effects of evolutionary processes that include recombination, mutation, genetic drift, demographic history (origins, migrations and admixtures) and evolutionary adaptations by natural selection.

# Modern Human Origins and Migrations



> Modern human bony conformation was established in Africa around 330 – 200,000 years ago.

> Genetic evidence supports the fossil and archaeological evidence.

# Population Variation Databases

**http://hapmap.ncbi.nlm.nih.gov/**

A computer security audit revealed security flaws in the legacy HapMap site and NCBI has took it down in June 2016.

**http://www.internationalgenome.org/**

IGSR: The International Genome Sample Resource

Providing ongoing support for the 1000 Genomes Project data

**http://gnomAD.broadinstitute.org/**
**Genome Aggregation Database and gnomAD Browser**

MONASH
University
MALAYSIA

Genomics Facility

# International HapMap Project

# The 1000 Genomes Project Dataset



Admixed

Migrants

CEU n = 99
FIN n = 99
GBR n = 91
CDX n = 93
TSI n = 107
CHB n = 103
JPT n = 104
MXL n = 64
PUR n = 104
IBS n = 107
PJL n = 96
CHS n = 105
ACB n = 96
YRI n = 108
GIH n = 103
CLM n = 94
GWD n = 113
ITU n = 102
KHV n = 99
PEL n = 85
MSL n = 85
ESN n = 99
STU n = 102
BEB n = 86
ASW n = 61
LWK n = 99

- **African (n = 661)**
- **East Asian (n = 504)**
- **South Asian (n = 489)**
- **European (n = 503)**
- **American (n = 347)**

The 1000 Genomes Project Consortium Nature (2015) 526:68-74.

http://www.1000genomes.org/page.php

| Samples | Populations | Mean Coverage | SNPs |
|---------|-------------|---------------|--------|
| 2,504 | 26 | 7.4 X | 84.7 M |

MONASH University
MALAYSIA

Genomics Facility

# Genome Aggregation Database

Whole exomes 125,748
Whole genomes 15,708



http://gnomad.broadinstitute.org

# gnomAD Dataset

| POPULATION | DESCRIPTION | GENOMES | EXOMES | TOTAL |
|---|---|---|---|---|
| AFR | African/African American | 4,368 | 7,652 | 12,020 |
| AMR | Admixed American | 419 | 16,791 | 17,210 |
| ASJ | Ashkenazi Jewish | 151 | 4,925 | 5,076 |
| EAS | East Asian | 811 | 8,624 | 9,435 |
| FIN | Finnish | 1,747 | 11,150 | 12,897 |
| NFE | Non-Finnish European | 7,509 | 55,860 | 63,369 |
| SAS | South Asian | 0 | 15,391 | 15,391 |
| OTH | Other (population not assigned) | 491 | 2,743 | 3,234 |
| | **Total** | 15,496 | 123,136 | 138,632 |

Sample numbers

Legend:
- Other
- Latino
- African
- Ashkenazi Jewish
- European
- South Asian
- East Asian

1000 Genomes   ESP   ExAC   gnomAD

https://gnomad.broadinstitute.org/about

# The HGDP-CEPH Cell Line Panel



Cann *et al.,* 2002. Science **296**:261-262.

# Simons Genome Diversity Project Dataset



**Human genomes(n = 300) from 142 populations**
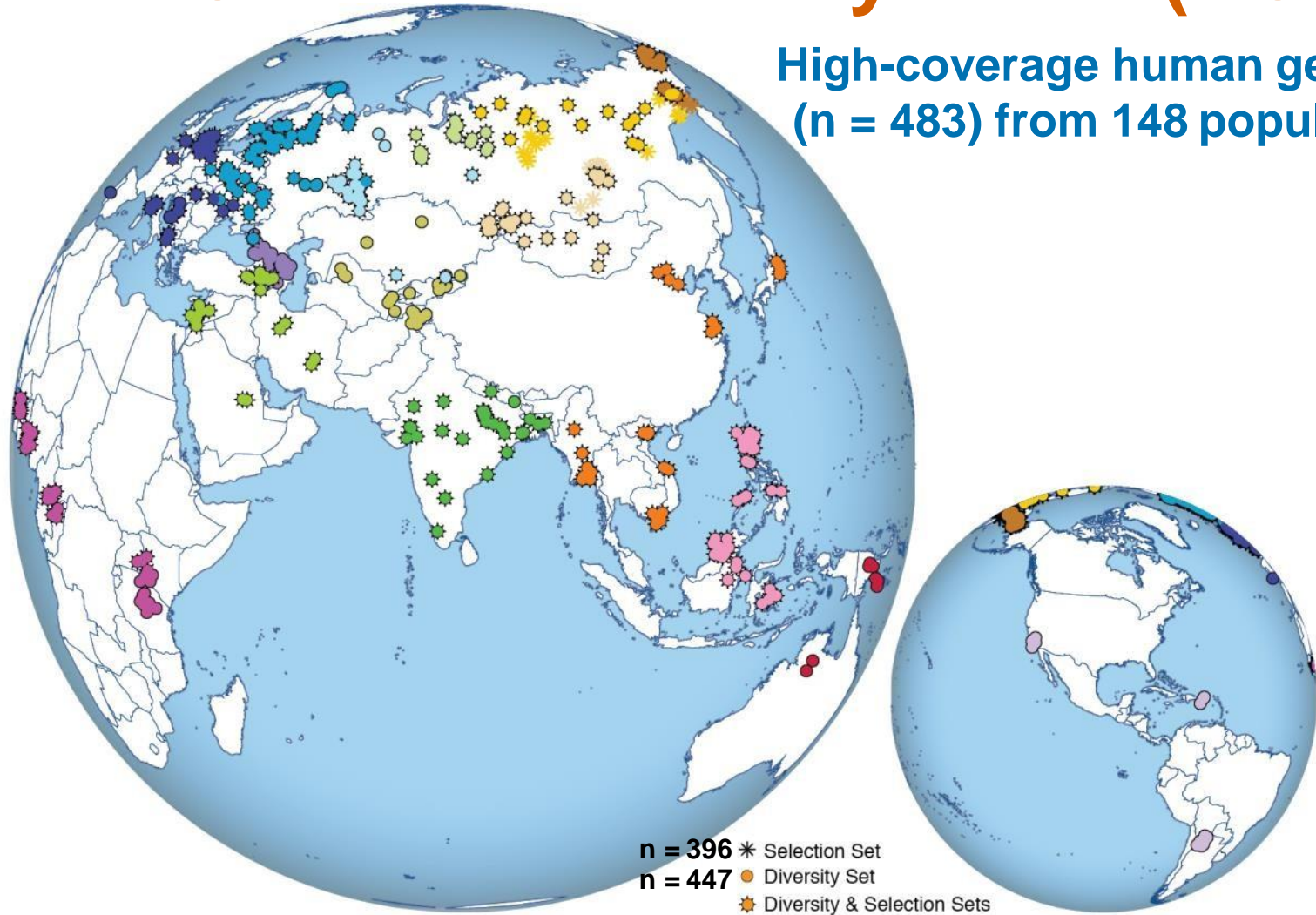https://www.simonsfoundation.org/life-sciences/simons-genome-diversity-project-dataset/

Mallick *et al.*, 2016. Nature 538:201-206.

# Estonian Biocentre
# Human Genome Diversity Panel (EGDP)

**High-coverage human genomes (n = 483) from 148 populations**



**n = 396** ✳ Selection Set
**n = 447** ● Diversity Set
✸ Diversity & Selection Sets

Pagani *et al.*, 2016. Nature 538:238-242.

MONASH University
MALAYSIA

Genomics Facility

# UK Biobank

https://www.ukbiobank.ac.uk/



UK Biobank is a large-scale biomedical database and research resource, containing in-depth genetic and health information from half a million UK participants. The database is regularly augmented with additional data and is globally accessible to approved researchers undertaking vital research into the most common and life-threatening diseases. It is a major contributor to the advancement of modern medicine and treatment and has enabled several scientific discoveries that improve human health.

# TOPMed Program

➢ **Trans-Omics for Precision Medicine (TOPMed) Program funded by NIH.**

➢ The goal of the TOPMed program is to generate scientific resources that will improve the understanding of heart, lung, blood, and sleep disorders and advance precision medicine.

https://www.nhlbiwgs.org/

# Revolution in Personalized Medicine



Stark *et al.* Am J Hum Genet (2019) 104:13-20.

# Worldwide Population Relationships
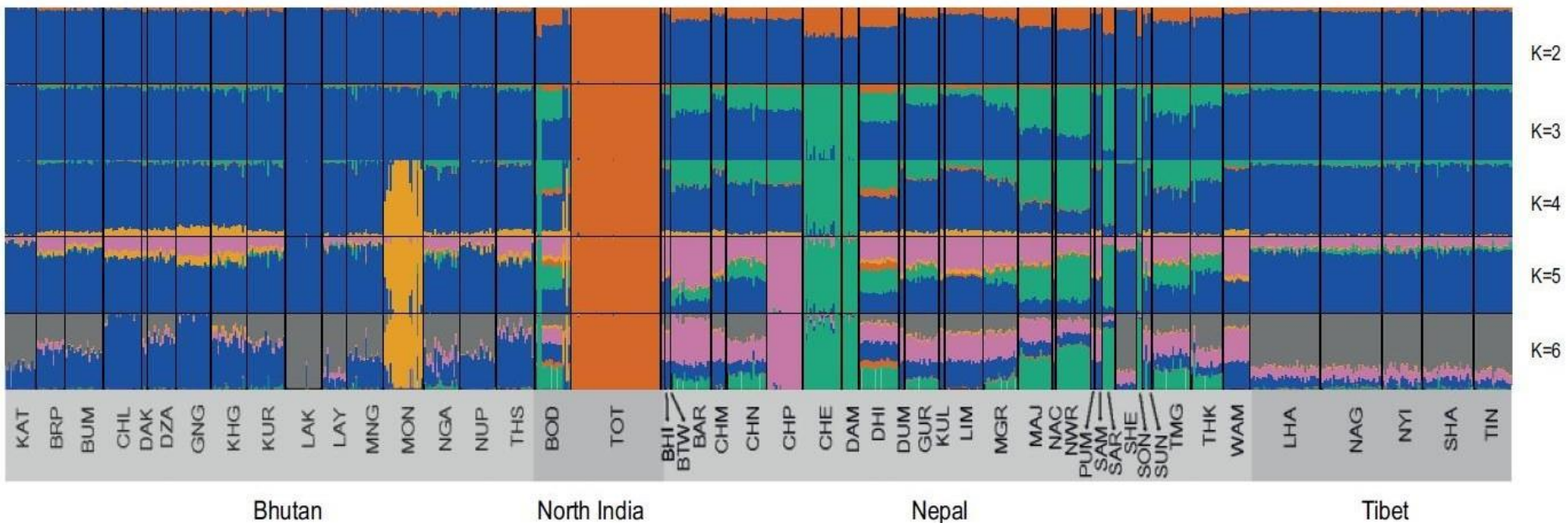## Principal Components Analysis (PCA)

# ADMIXTURE Analysis



- ➢ African populations are genetically more diverse than non-Africans.
- ➢ Genetic diversity outside of Africa tends to be a subset of the diversity within Africa.

# Genetic Drift

➢ Strong effect of random fluctuations in allele frequencies due to population isolation and bottlenecks.



Arciero *et al.* Mol Biol Evol (2018) 35:1916-1933.

# Population Stratification and GWAS

➢ GWAS can be confounded by population stratification—**systematic ancestry differences between cases and controls leading to a spurious association.**

➢ These associations may appear to be significant, but they are driven by the cohort's relatedness rather than variants that truly affect trait or disease risk.

➢ Failure to control for it may lead to confounding, causing a study to fail for lack of significant results or resources to be wasted following false positive signals.

MONASH University
MALAYSIA

Genomics Facility

# How Does it Occur?

➢ Whenever there are substantial variation across ethnicities in the frequency of the variant genotype being considered.

➢ If there is substantial variation across ethnicities in disease rates after adjustment for risk factors, other than the genotype of interest, that were collected in the study.

➢ The allele frequencies track with the disease rates across ethnicities, for reasons other than the effects of the allele of interest. For example, an allele with a clade or gradient of increasing frequency from North to South Asia might track with another factor, such as dietary differences or air pollution, that affects disease risk, thus, introducing bias from population stratification when studying the effect of the allele.

➢ Self-reported ethnic information from study participants does not reduce bias to an acceptable level.

Wacholder *et al. Cancer Epidemiology, Biomarkers and Prevention* (2002) 11:513-520.

MONASH University
MALAYSIA

Genomics Facility

# Population Stratification

|       | Cases | Controls |
|-------|-------|----------|
| **A** | 101   | 20       |
| **a** | 20    | 101      |

**OR = 25.5**

|       | Cases | Controls |
|-------|-------|----------|
| **A** | 100   | 10       |
| **a** | 10    | 1        |

**OR = 1**

Europeans

|       | Cases | Controls |
|-------|-------|----------|
| **A** | 1     | 10       |
| **a** | 10    | 100      |

**OR = 1**

East Asians

MONASH University MALAYSIA

Genomics Facility

# Individual Relatedness

➢ Ancestry differences:
  ➢ Ancestry differences refer to different ancestry among individuals in a study.

  ➢ If an association study contains individuals from different populations or differing degrees of admixture, the individuals will have different degrees of relatedness among them

➢ Cryptic relatedness:
  ➢ Cryptic relatedness exists when some individuals are closely related, but this shared ancestry is unknown to the investigators and the study subjects.

Sul *et al.* PLoS Genetics *(*2018) 14:e1007309.

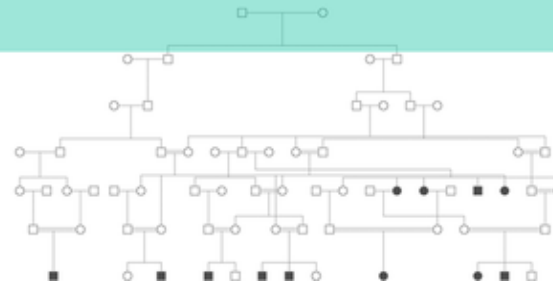https://journals.plos.org/plosgenetics/article?id=10.1371/journal.pgen.1007309

MONASH
University
MALAYSIA

Genomics Facility

# Shared Ancestry



Sul *et al.* PLoS Genetics *(*2018) 14:e1007309.
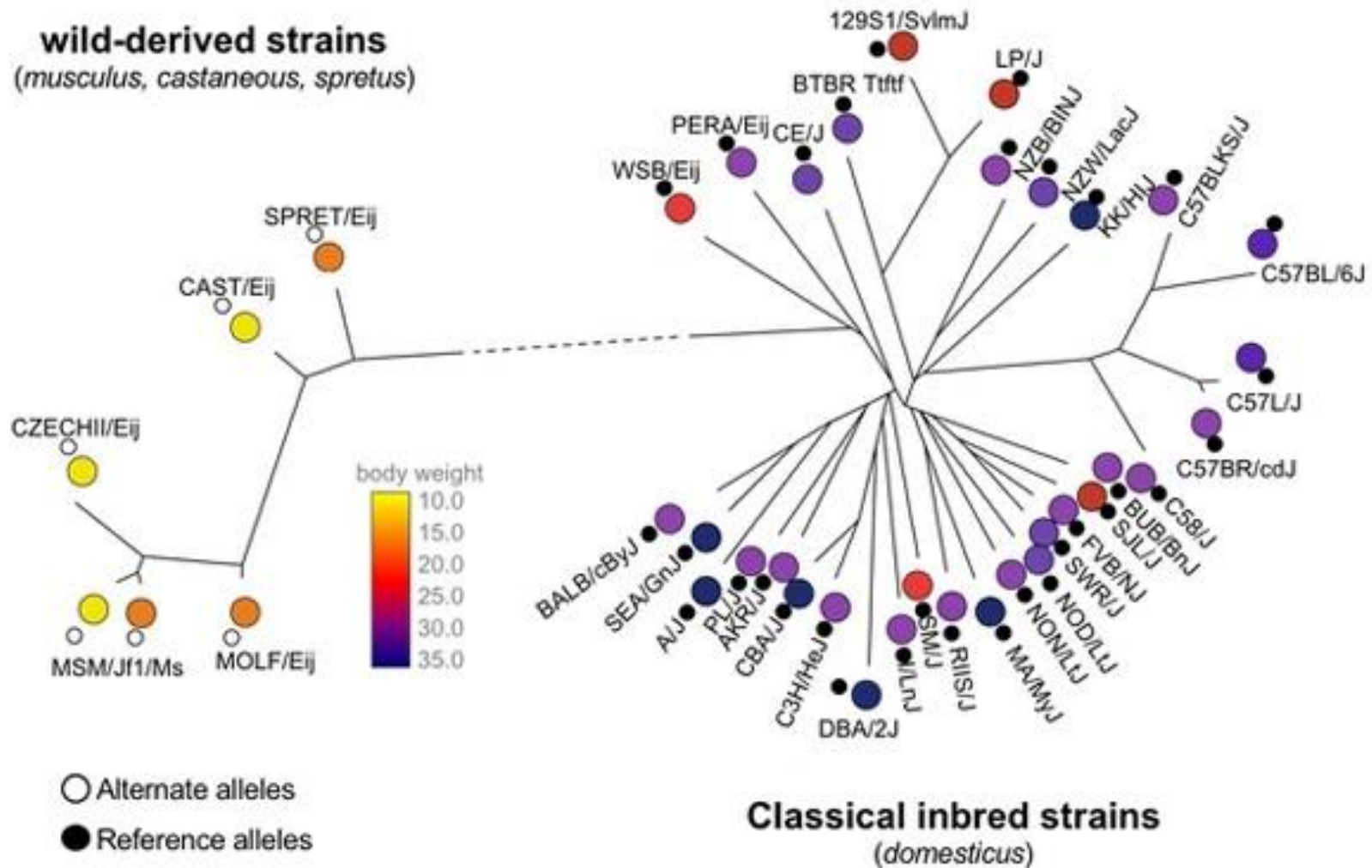
# Population Stratification and GWAS



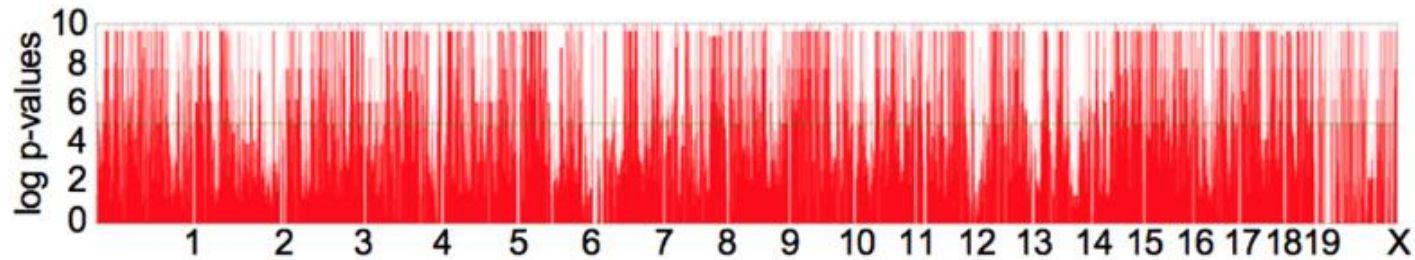Sul *et al.* PLoS Genetics *(*2018) 14:e1007309.
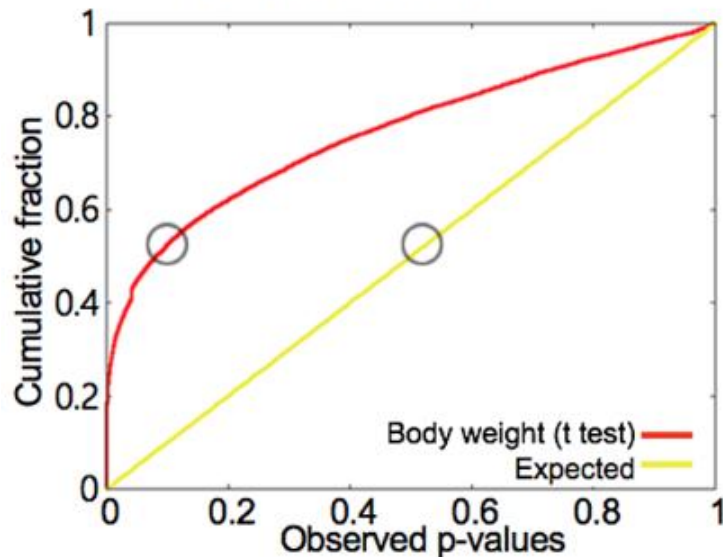
# Individual Relatedness



Sul *et al.* PLoS Genetics *(*2018) 14:e1007309.
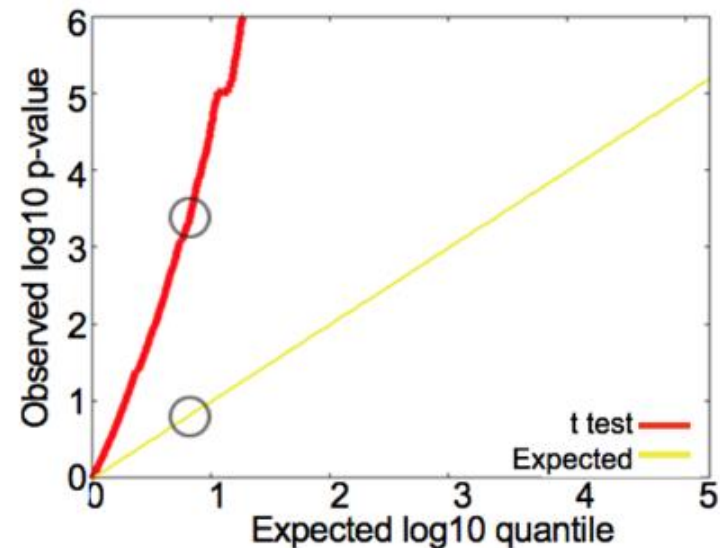
# Effect of Population Structure



**A** GENOME-WIDE ASSOCIATION MAP

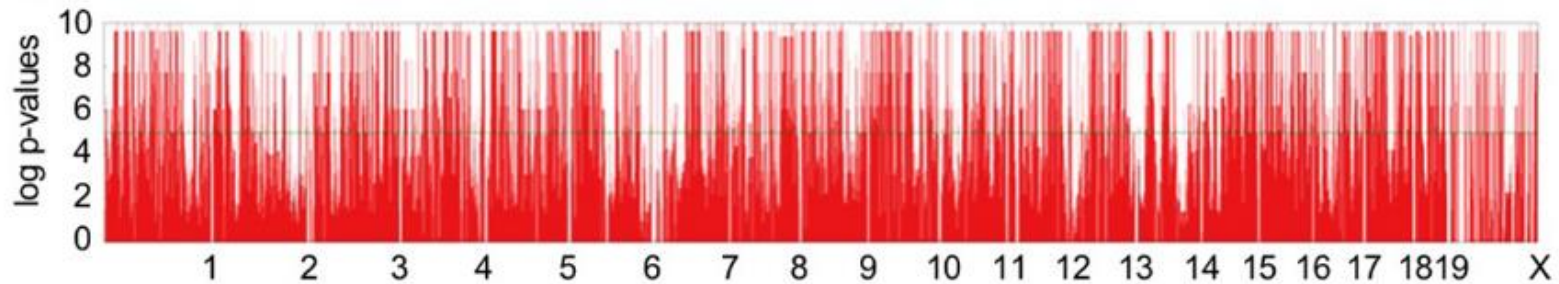**B** CUMULATIVE p-VALUE DISTRIBUTION

**C** Q-Q PLOT

Sul *et al.* PLoS Genetics *(*2018) 14:e1007309.
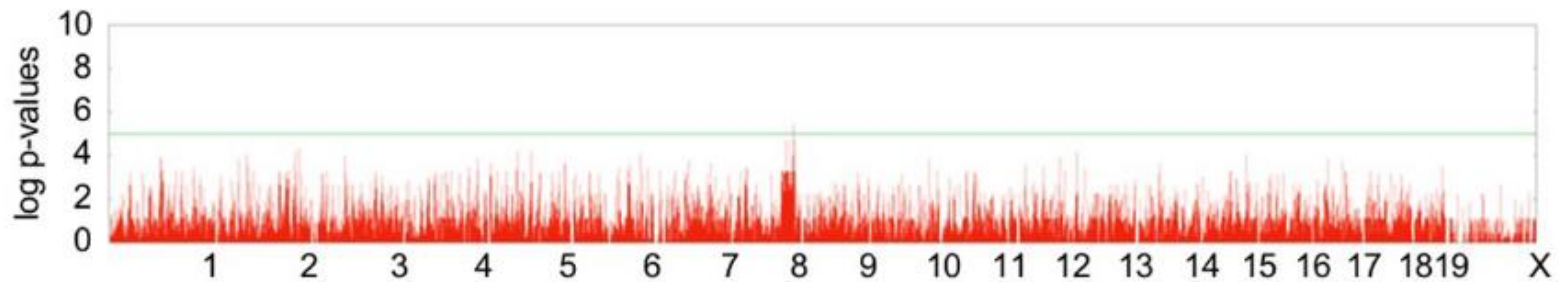
# Correcting for Population Stratification

➢ Replication in different populations:
>  ➢ Major bias in the same direction in populations with substantially different ethnic mixes is very unlikely because the conditions that allowed major bias are unlikely to be repeated.

➢ Genomic control markers:
>  ➢ Genomic control uses markers unrelated to disease to correct the bias.

➢ Correcting for population substructure:
>  ➢ To resolve the problem with population structure with the use of principal components (PC) of the genotyped dataset to model population relationships, which could be interpreted as a proxy for ancestry information, and included in the model as covariates.
>  ➢ Implemented as mixed linear model such as **E**fficient **M**ixed **M**odel **A**ssociation e**X**pedited (EMMAX) in which one SNP is fit in the model as a fixed covariate and, at the same time, a relationship matrix corrects for population structure.

MONASH University MALAYSIA

Genomics Facility

# Effect of Population Structure



Sul *et al.* PLoS Genetics *(*2018) 14:e1007309.

# Questions?

## Post Questions on the Slack Channel

## [qasim.ayub@monash.edu](mailto:qasim.ayub@monash.edu)

# Practical Exercise