# Genome-Wide Association Studies (GWAS)
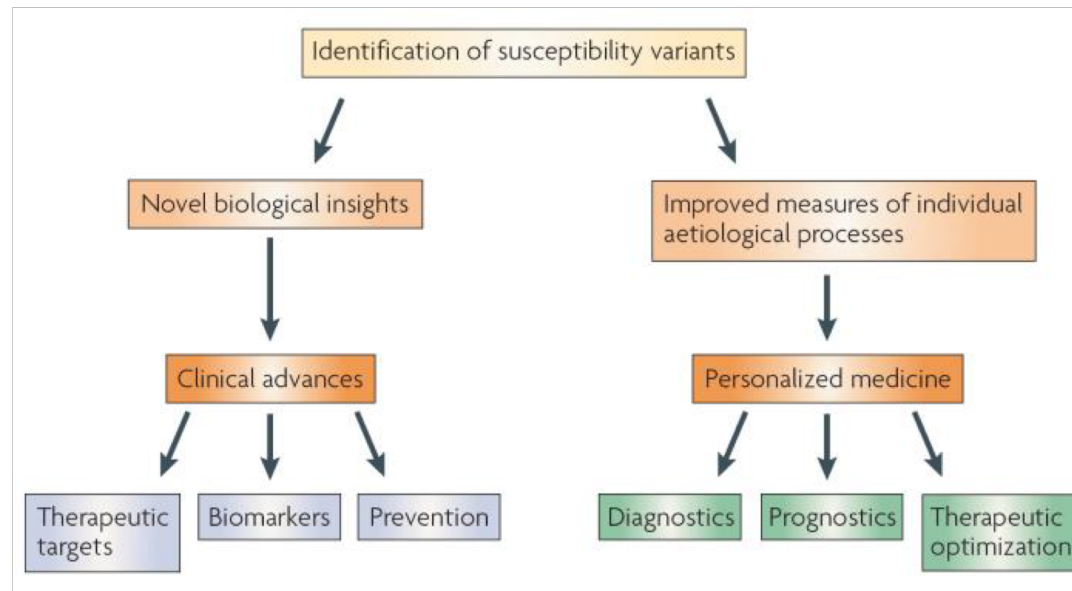
London School of Hygiene and
Tropical Medicine

# Outline

- GWAS overview – Utility & Successes

- GWAS Study Design

- Post GWAS Follow up

- Recap

- Case study

# Purpose of Genetic Association Studies

- Determine if there is a genetic component contributing to phenotype (i.e. disease) under investigation (heritability)

- Identify the genetic region/gene/polymorphism causing the disease

- Determine the effect size of the genetic component

# Genome wide association studies (GWAS)

- High-throughput approach scanning marker across the genome - linking genotype to phenotype

- Relies on dense sets of genetic markers - Usually SNPs and SNP tags for other variation (via LD)

- Usually comparison of variation between affected (cases) and unaffected individuals (controls).

- Goal: Identify markers with significant associations to disease



```
..ACTCGACGATTTACGGTACTTAGGAGCATACGCTAC..
..ACTCTACGATTTACGGTACTTAGGAGCATACGCTAC..
..ACTGTACGATTTACGATACTTAGGAGCATATGCTAC..
..ACTGTACGATTTACGGTACTTAGGAGCATATGCTAC..
..ACTGTACGATTTACGGTACTTAGGAGCATATGCTAC..
..ACTGTACGATTTACGATACTTAGGAGCATAGGCTAC..
..ACTGTACGATTTACGATACTTAGGAGCATAGGCTAC..
..ACTGTACGATTTACGGTACTTAGGAGCATATGCTAC..
..ACTGTACGATTTACGATACTTAGGAGCATAGGCTAC..
```
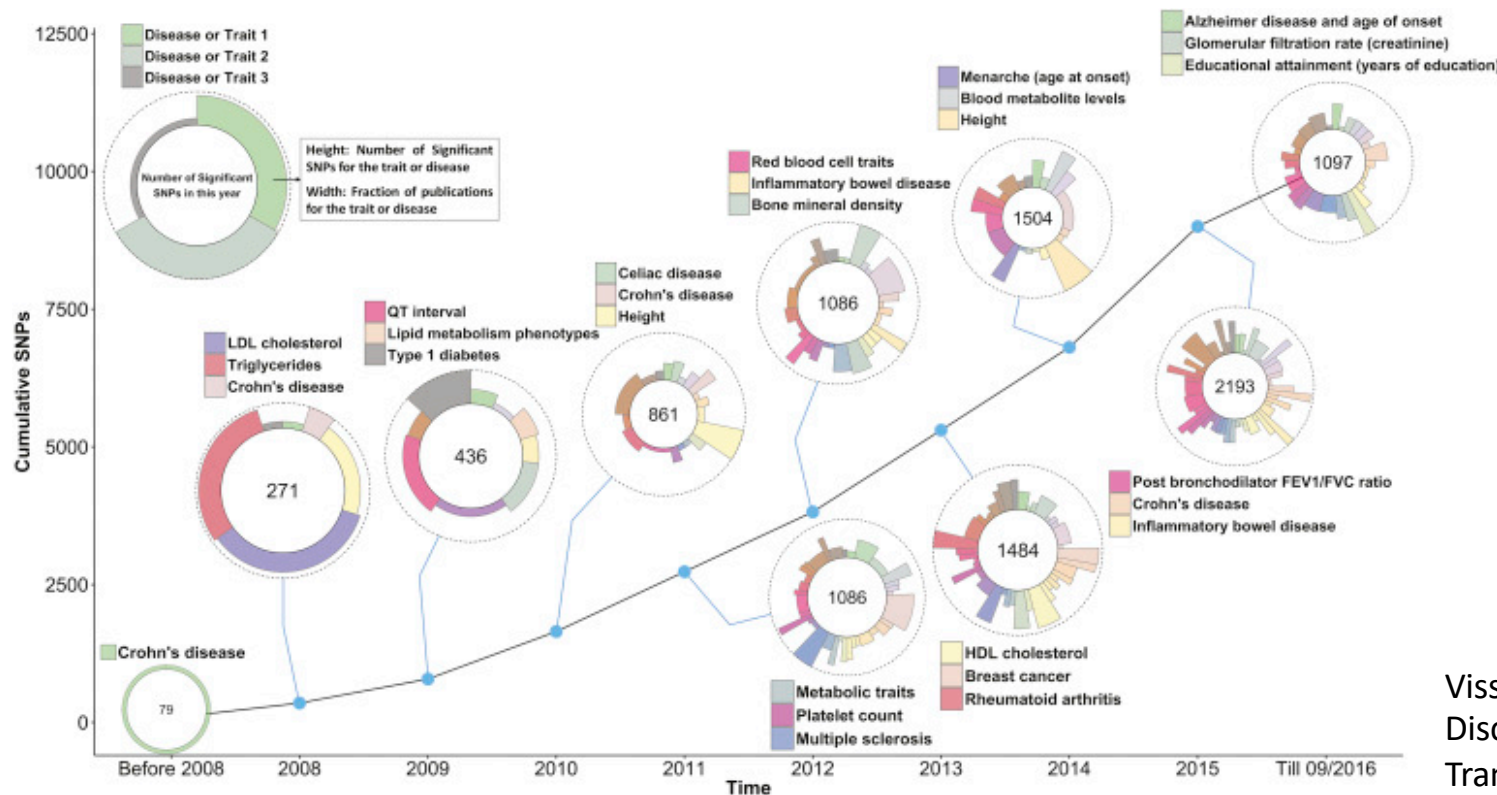
SNPs may have 2, 3 or 4 alleles (most are biallelic)

# Lots of Success

**May 2018 (p≤5X10$^{-8}$)**

- 69,000 trait associations
- >5000 studies
- 3378 publications

NHGRI GWA Catalog
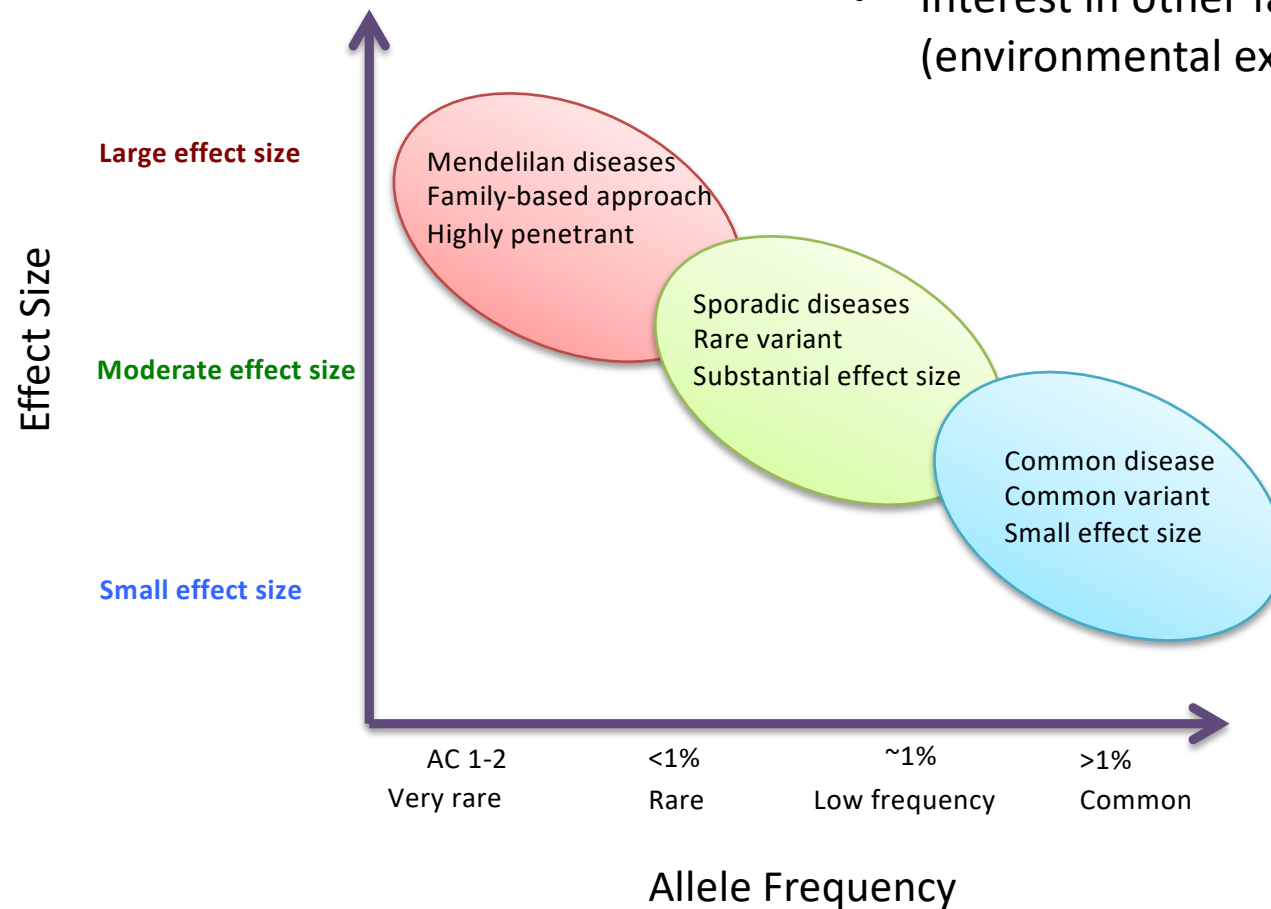www.genome.gov/GWAStudies
www.ebi.ac.uk/fgpt/gwas/

Visscher et al; 2017, 10 Years of GWAS Discovery: Biology, Function, and Translation
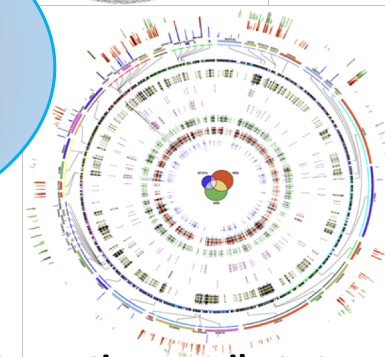
# Variant Identification

Study design is key
- Sample size?
- Interest in other factors of disease (environmental exposures, survival, effect size



Effect Size

Large effect size

Moderate effect size

Small effect size

Mendelilan diseases
Family-based approach
Highly penetrant

Sporadic diseases
Rare variant
Substantial effect size

Common disease
Common variant
Small effect size

AC 1-2
Very rare

<1%
Rare

~1%
Low frequency

>1%
Common

Allele Frequency

# Multifactorial determinants of pathogenesis & clinical outcome



- Co-infection
- Immunosuppression
- Lifestyle & socioeconomic factors
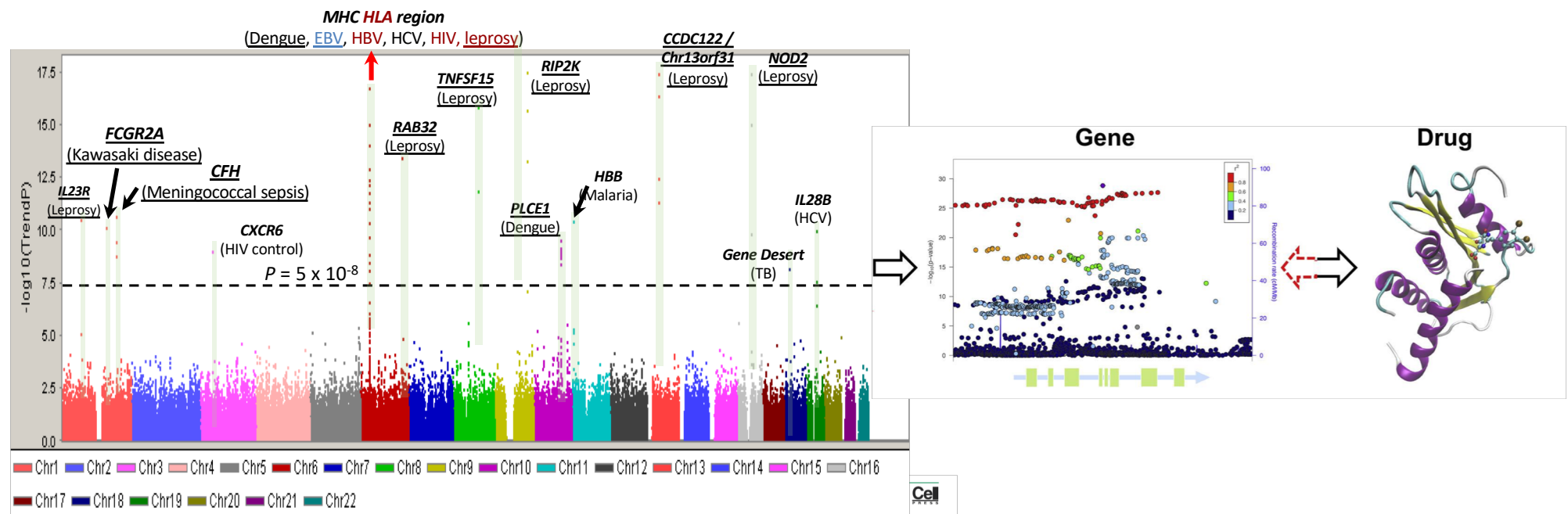
**Environment**

**Pathogen**

**Host**

- What is the extent of pathogen genome diversity?

- Is there a link between genotype and phenotype?

- What are the transmission patterns?

- **Does host genetics contribute to susceptibility to infection/disease outcome?**

- **How do genetic variants influence virus biological function?**

# GWAS of Infectious diseases
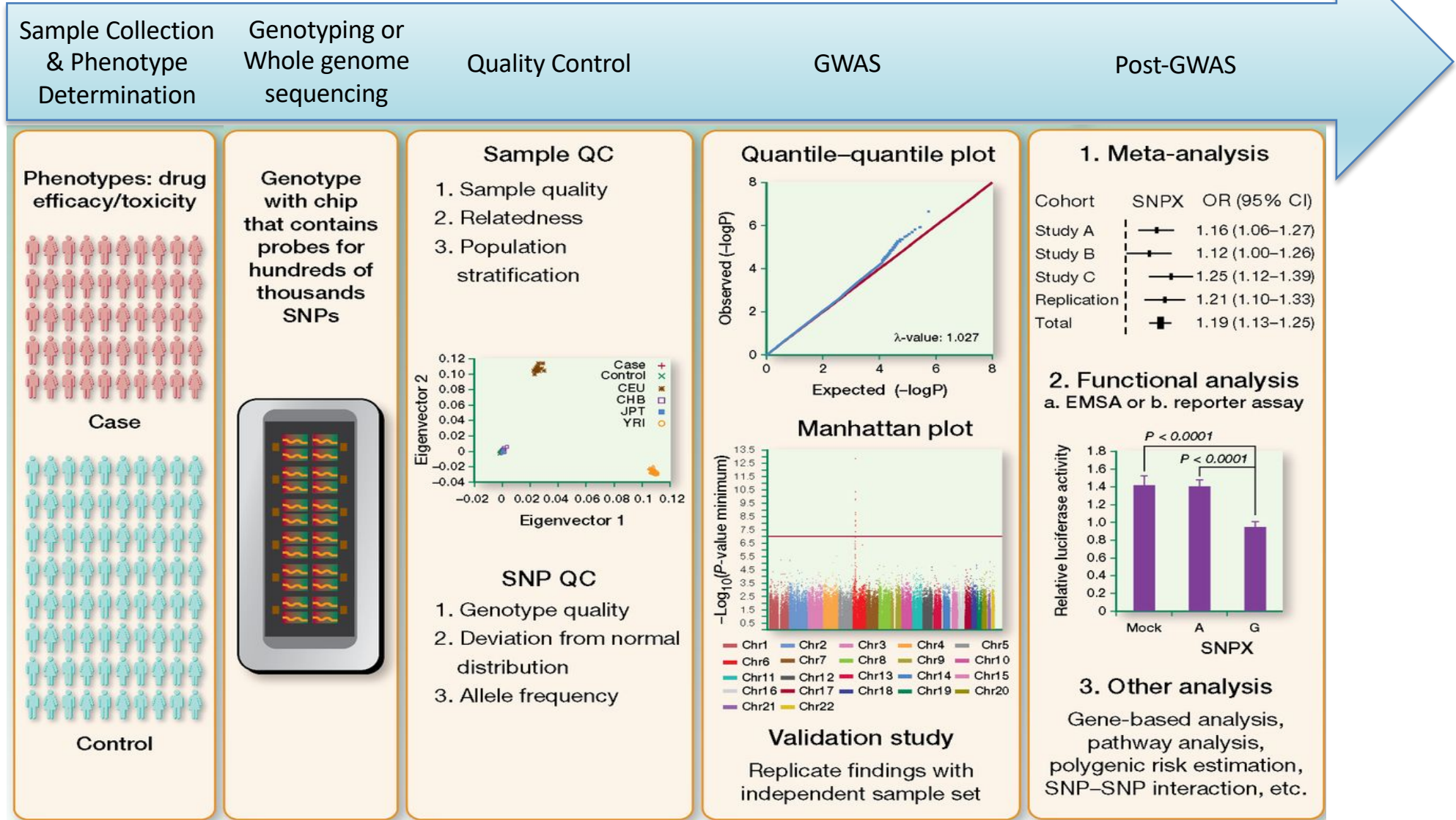
**Phenotypes Studied:**

– Case- Control study: Susceptibility, severity, pathogen clearance, response to vaccination, severe disease

– Quantitative trait: Antibody response, viral load, cell count
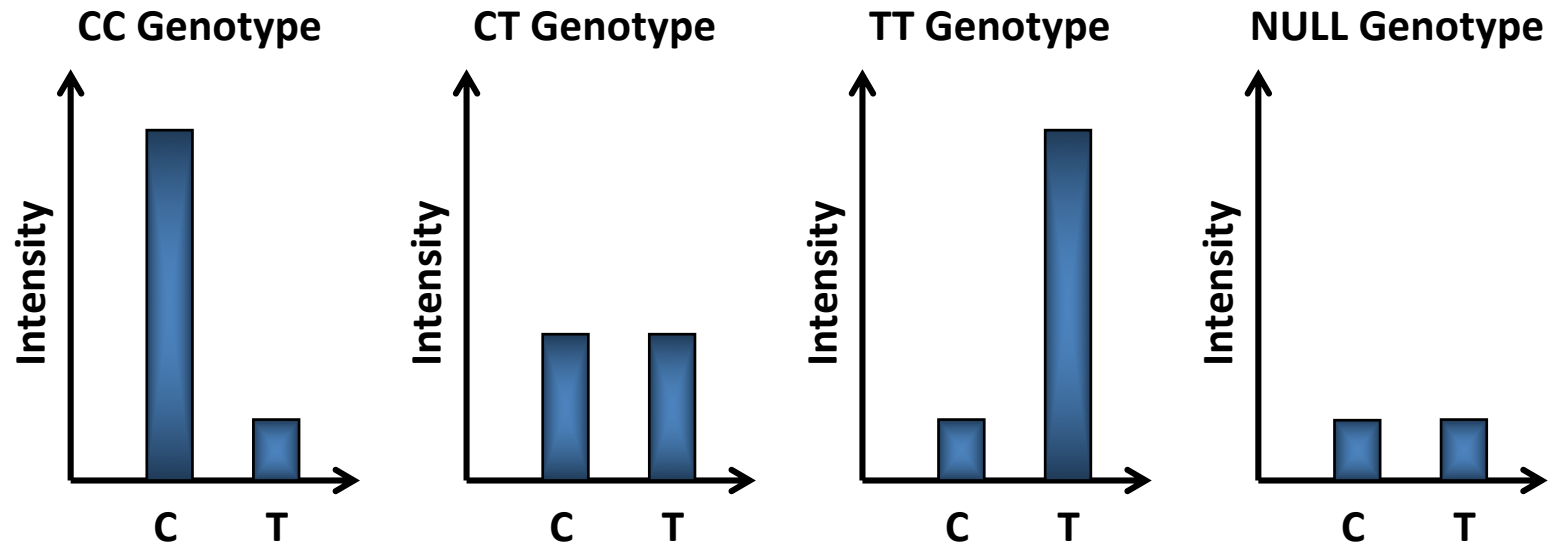


Host–pathogen interactions revealed by human genome-wide surveys

Chiea Chuen Khor[1,2,3] and Martin L. Hibberd[1,2]

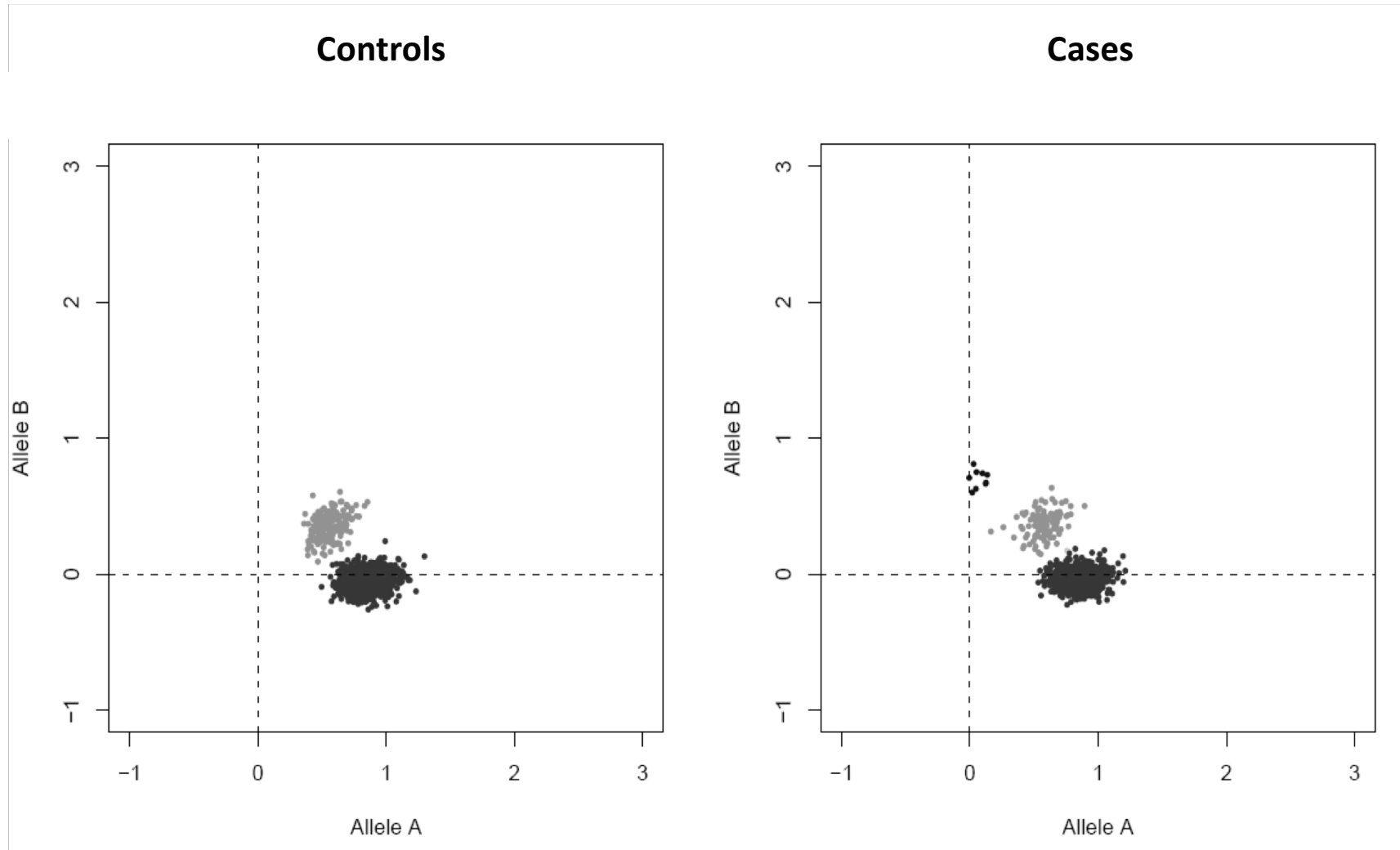[1] Infectious Diseases, Genome Institute of Singapore, Singapore

# GWAS Workflow



**Workflow stages (arrow):** Sample Collection & Phenotype Determination → Genotyping or Whole genome sequencing → Quality Control → GWAS → Post-GWAS

**Sample Collection & Phenotype Determination**

Phenotypes: drug efficacy/toxicity

Case

Control

**Genotyping or Whole genome sequencing**

Genotype with chip that contains probes for hundreds of thousands SNPs

**Quality Control**

Sample QC
1. Sample quality
2. Relatedness
3. Population stratification

SNP QC
1. Genotype quality
2. Deviation from normal distribution
3. Allele frequency

**GWAS**

Quantile–quantile plot
λ-value: 1.027

Manhattan plot

Chr1  Chr2  Chr3  Chr4  Chr5
Chr6  Chr7  Chr8  Chr9  Chr10
Chr11  Chr12  Chr13  Chr14  Chr15
Chr16  Chr17  Chr18  Chr19  Chr20
Chr21  Chr22

Validation study
Replicate findings with independent sample set

**Post-GWAS**

1. Meta-analysis

| Cohort | SNPX | OR (95% CI) |
|---|---|---|
| Study A | | 1.16 (1.06–1.27) |
| Study B | | 1.12 (1.00–1.26) |
| Study C | | 1.25 (1.12–1.39) |
| Replication | | 1.21 (1.10–1.33) |
| Total | | 1.19 (1.13–1.25) |

2. Functional analysis
a. EMSA or b. reporter assay

$P < 0.0001$
$P < 0.0001$

SNPX: Mock, A, G

3. Other analysis
Gene-based analysis, pathway analysis, polygenic risk estimation, SNP–SNP interaction, etc.
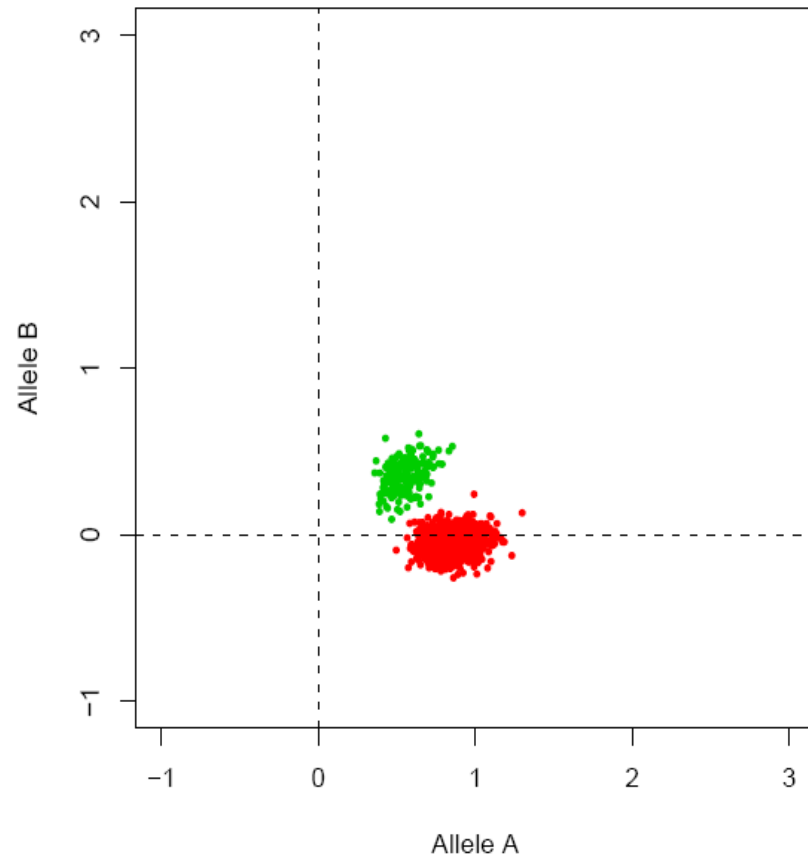
# Raw data is not a genotypes, but Allelic hybridization

# Genotype Calling

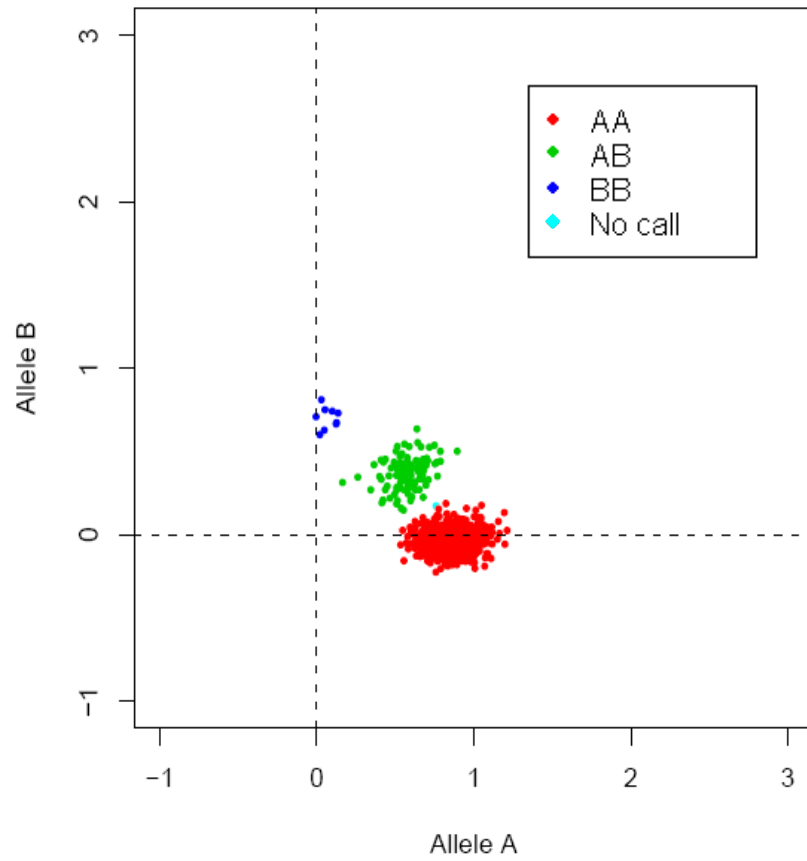# Genotype Calling
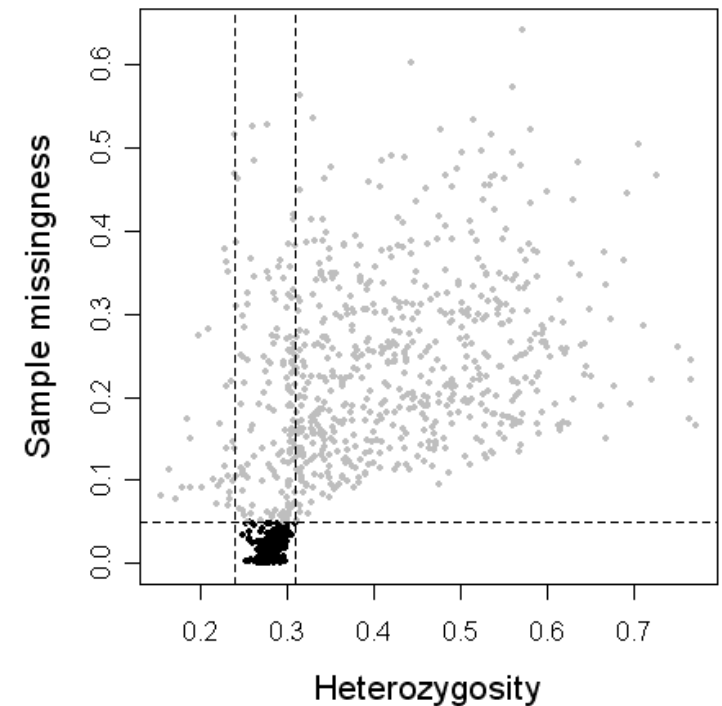
# Need for high quality data

- Number of variants assayed $\Rightarrow$ errors and genotype or sequence miscalls are bound to happen
- If problematic <u>samples</u> not identified and excluded, they can affect the results of the entire experiment
- If <u>SNPs</u> with erroneous genotyping or sequencing not identified and excluded, can produce false signals of associations
- QC samples and SNPs

# Quality Checks

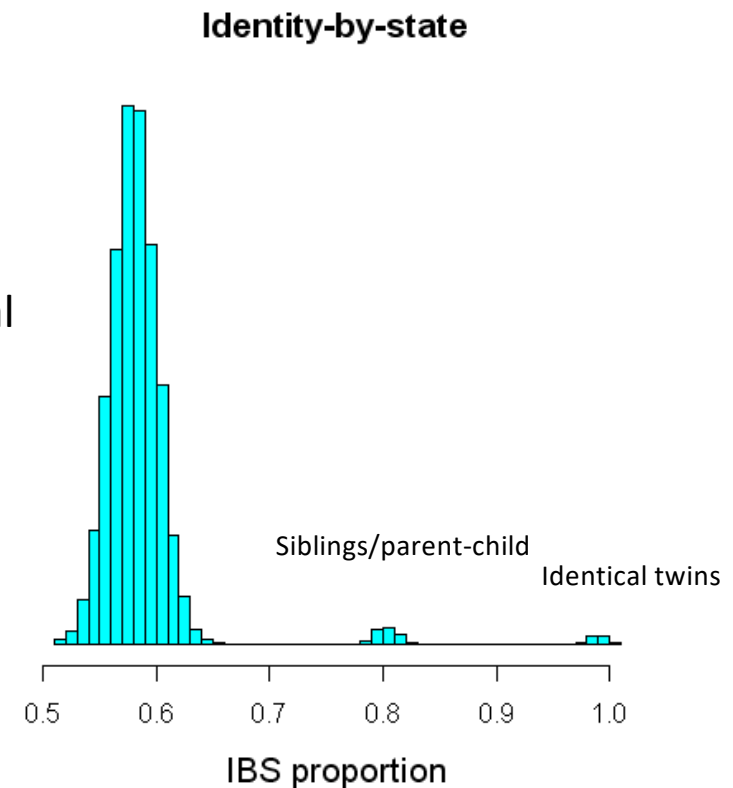| Variable | Comments |
|---|---|
| Genotyping Call Rate | Low call rate often correlates with error. Some low call rate SNPs or samples may still be good. |
| Genotyping Quality | Worse quality score (GenCall) correlates strongly with error rate |
| Sex concordance | Check expectations for X marker heterozygosity and Y marker positive results. Can estimate error rate. |
| Sample Relatedness | Check for related samples (expected or unexpected) |
| Mendelian Inheritance Errors | For trio/family data, can identify problem samples and families. Can estimate error rate. |
| Replicate concordance | Check for consistent genotype calls in duplicate samples |
| Batch effects | Check for genotyping call differences due to plate |
| Hardy-Weinberg Equilibrium | Violation across all sample groups may indicate error, but can also be a good test of association |
| Population Stratification | Check for population substructure using the genome-wide data |

# Sample QC

- Identify SNPs with high rates of missingness and heterozygosity

- Remove samples deviating from average

- Deviations could arise due to several reasons

  – Contamination of samples (high heterozygosity

  – Inbreeding (low heterozygosity)

  – Ancestral differences

  – Data quality / Poor genotype calling

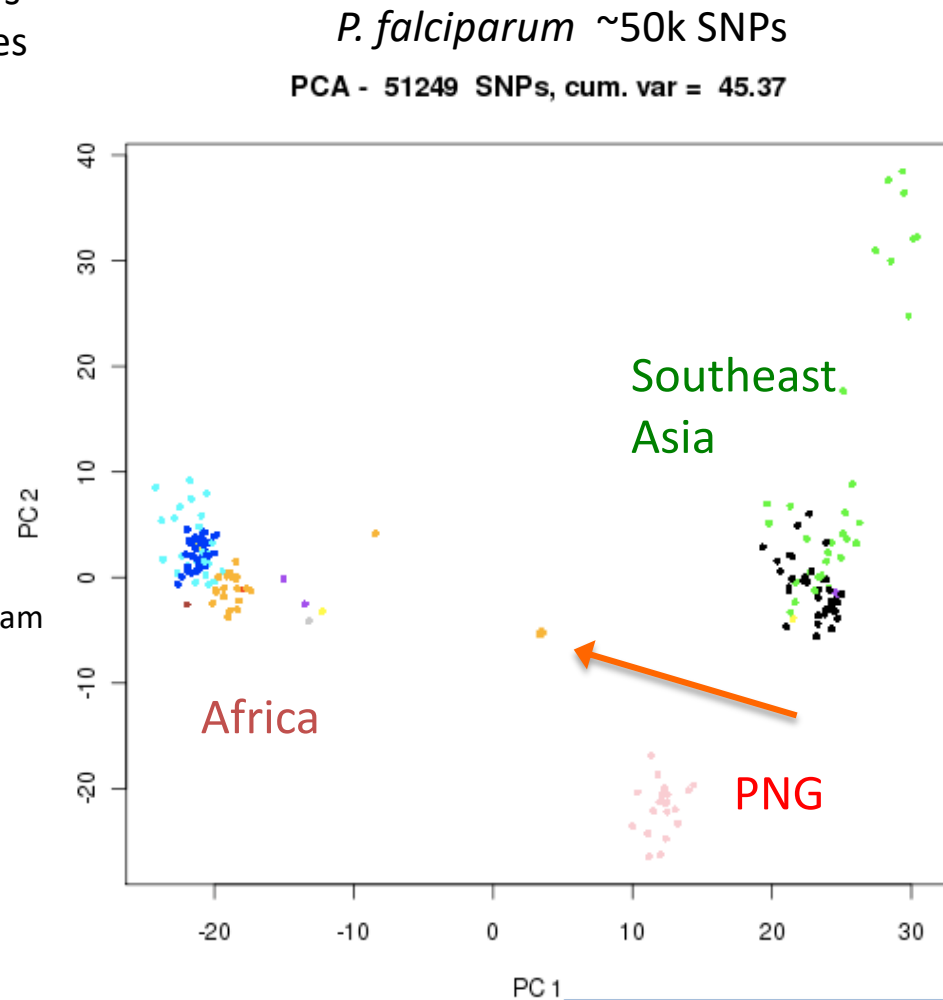- Heterozygotes more likely to be missing

# Sample QC

- Identify related / duplicated samples
- Relatedness is a problem because of overrepresentation of selected alleles, which will bias any multivariate analysis (correlated data!); e.g. PCA or multivariate regression
  - Related samples need to be excluded or taken into account during subsequent analyses
- Related individuals will share more alleles IBS than expected by chance, with the degree of additional sharing proportional to the degree of relatedness.



Identity-by-state

Siblings/parent-child

Identical twins

IBS proportion

# Sample QC

- Population substructure or stratification occurs when samples have different genetic ancestries

- Can lead to spurious associations due to differences in ancestry rather than true associations

- Imperative to check for population structure within samples

- Identify samples that are outliers
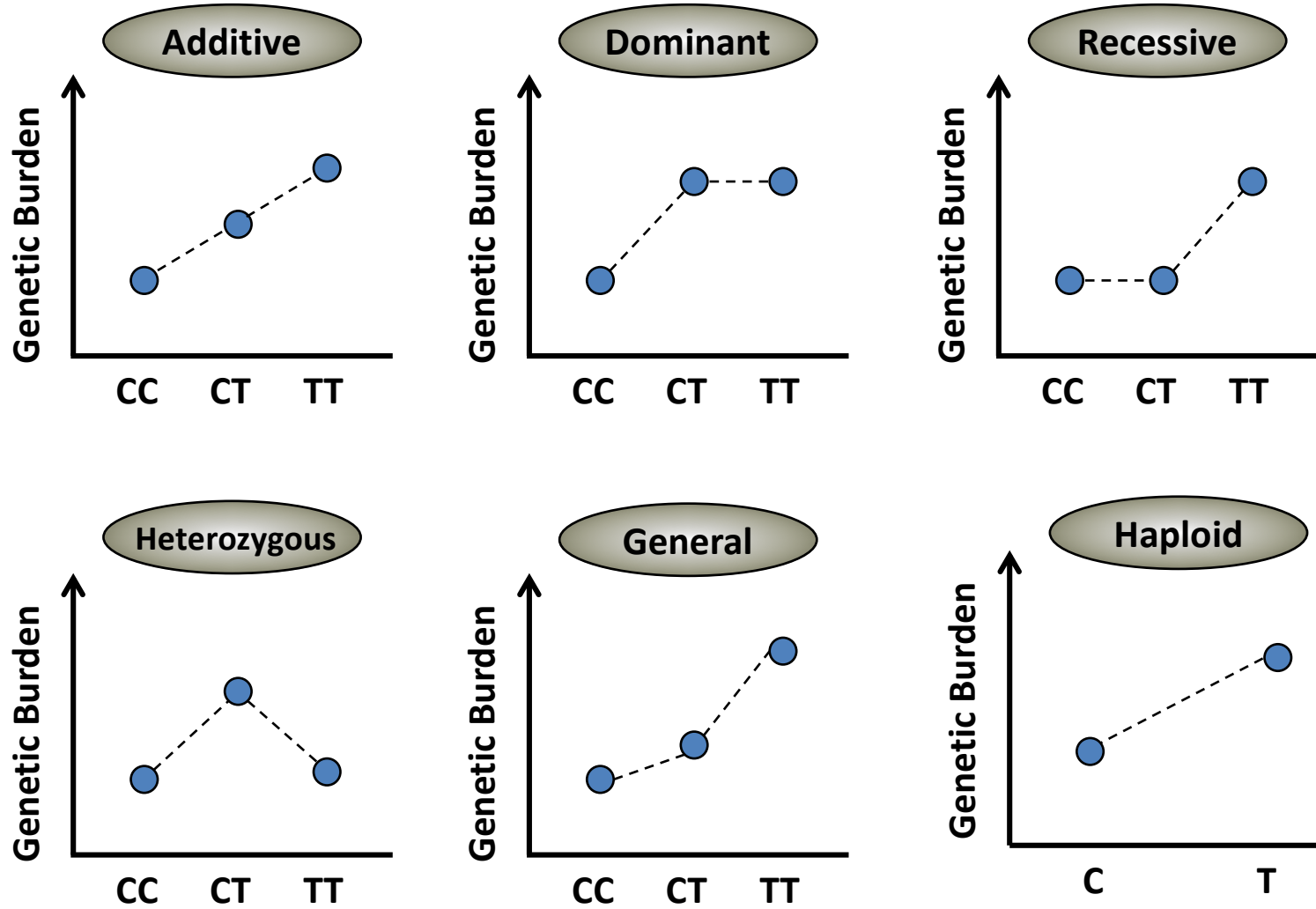  - Can control for structure if identified, in downstream analysis

*P. falciparum* ~50k SNPs

PCA - 51249 SNPs, cum. var = 45.37

Southeast Asia

Africa

PNG

Data from Ocholla et al. 2014

# Testing for associations using regression models

| Outcome | Example | Model |
|---|---|---|
| Continuous | IC50 levels | Linear regression |
| Binary | Malaria status | Logistic regression |

- Using a single SNP in turn, but also can include
  - Interactions
  - Adjustment for confounders
  - Model building and risk prediction strategies
  - Diagnostic tools to assess model fit

- Also non-parametric approaches

Genetic models tested in a regression framework

# Association Studies

**Direct Association**

Tests the genetic variant directly responsible for causing the disease.

# Association Studies



**Direct Association**

Tests the genetic variant directly responsible for causing the disease.

**Indirect Association**

Genetic variant tested is not directly responsible for the disease, but is located near to the disease-causing variant and thus 'correlated', or in linkage disequilibrium (LD).

# Each population has a distinct pattern of genome variation



SNPs

> Most SNPs are correlated with surrounding SNPs. This is known as **linkage disequilibrium (LD)**

> Linkage disequilibrium reflects the common combinations of variants (haplotypes) that exist in the population

# Haplotypes

...ACTC**G**ACGATTTACG**G**TACTTAGGAGCATA**T**GCTAC...
...ACTC**T**ACGATTTACG**G**TACTTAGGAGCATA**T**GCTAC...
...ACTG**T**ACGATTTACG**A**TACTTAGGAGCATA**G**GCTAC...
...ACTGA**A**ACGATTTACG**G**TACTTAGGAGCATA**T**GCTAC...
...ACTG**T**ACGATTTACG**G**TACTTAGGAGCATA**T**GCTAC...
...ACTG**T**ACGATTTACG**A**TACTTAGGAGCATA**G**GCTAC...
...ACTG**G**ACGATTTACG**G**TACTTAGGAGCATA**G**GCTAC...
...ACTG**T**ACGATTTACG**G**TACTTAGGAGCATA**T**GCTAC...

- **A haplotype is an observed sequence of variants**
- Each population has its own pattern of common haplotypes
- By knowing the pattern of haplotypes within a population we may be able to impute genotype at an untyped position

# Why is LD important in humans?

- ~10m genetic variants in the human genome, costly to genotype everything (pre-2012?)
- LD $\Rightarrow$ Reduced amount of genotyping required
- The availability of whole genome sequencing on large numbers of samples makes LD redundant

# Imputation

```
TCCGGACAC C TT C TA A GG  ⎫
TCTGGACAC C TT C TA A GG  ⎪
TCTGTACAC A GG A TT T CG  ⎪  Reference panel:
ACTGGACAC A GG A TT T GG  ⎬  HapMap YRI haplotypes
ACCGTCTTC C TT C TA A CG  ⎪
TCCGGACAC C TT C TA A GG  ⎭
```

```
A..G.....C..C..A..}  Genotyped individual

AC?G....CCTTCTAA..}  Imputed individual
```

Using correlations or 'recurring patterns' in the data to fill in the blanks.

# Imputation



Study 1

Study 2

Study 1 with imputed missing SNPs

- Imputation
    - Requires GWAS genotypes to be used as scaffold
    - Requires reference datasets (e.g. www.hapmap.org; www.1000genomes.org) where the LD (correlation) between SNPs is known and allows imputation of genotypes for variants not typed on a given array. Increasingly these could include reference datasets generated by whole-genome sequencing of subsets of individuals from the populations included in the study
    - There is specialist software to facilitate imputation as well as meta-analysis

# Why impute?

- To predict missing genotypes that haven't been directly typed
    - **Increased power.** The reference panel is more likely to contain the causal variant (or a better tag) than a GWAS array.
    - **Fine-mapping.** Imputation provides a high-resolution overview of an association signal across a locus.
    - **Meta-analysis.** Imputation allows GWAS typed with different arrays to be combined up to variants in the reference panel.

    **What if the LD structure in the imputed population is different to the reference?**

# Association signals across the genome

Manhattan plot: Severe malaria GWAS (n=1,000 cases, 1,500 controls)

# Sickle trait is the strongest known determinant of severe malaria risk



| | Relative risk of severe malaria in children with HbS   AS genotype |
|---|---|
| **Gambia** | **0.11** |
| **Kenya** | **0.17** |
| **Malawi** | **0.08** |

$P = 2 \times 10^{-31}$ ($n = 3630$)

- Genetic factors determine 25% of malaria risk in Kenyan children and sickle trait accounts for only 2% of total variation (Mackinnon et al, 2005)

Signals of malaria association in chromosome 11 in The Gambia

**Signals of malaria association in chromosome 11 in The Gambia**

- Genotyped SNP on Affy500K
- Imputed using HapMap YRI as reference

-log10 P-value

Physical position (Mb)

Signals of malaria association in chromosome 11 in The Gambia

$P \sim 10^{-14}$

- Genotyped SNP on Affy500K
- Imputed using HapMap YRI as reference
- Imputed using sequence data from The Gambia

-log10 P-value

Physical position (Mb)

# Going beyond GWAS

- Need to validate and confirm findings

  - Replication studies and meta analysis

  - In silico annotation tools

- If using genotyping arrays, fine-mapping the causal variant

  - Targeted-resequencing

  - Transethnic mapping

- Functional studies

# Replication

- Assay a small subset of SNPs that arose from GWAS scan
- Ideally within same population, but often unlikely…
- Aim to replicate in other populations for a similarly defined phenotype
- Population structure:
  - Problematic, since we will not have genome-wide data to assess extent of confounding
  - Have to rely on informative surrogates if available (e.g. self-reported ethnicity, language, location)

# Meta-analysis

- Combine multiple genome-wide scans of the same phenotype

- Consistency of phenotypic definition is crucial, given expectation of marginal genetic effects

- Genome-wide pooling, publication bias less of an issue

- Summary stats can be used for analysis

**G6PD 202A and severe malaria**

The Gambia M
The Gambia F
Kenya M
Kenya F
Malawi M
Malawi F
Ghana M
Ghana F

Summary

0.50 0.63    1.00  1.26    2.00

Odds Ratio

# Benefits of GWAS Meta-Analysis

- Increased sample sizes for many disease and continuous trait consortia
    - increased power to detect new loci
    - new pathways and important biological insights gained
    - greater power to detect even smaller effect sizes and greater coverage of allele frequency spectrum
- Power of large collaborations/consortia
    - Design better powered replication and fine-mapping experiments

# Heterogeneity

- Results from meta-analysis of various studies may suggest between study heterogeneity (e.g. especially when combining populations of different ancestry)

- How to interpret heterogeneity?
  - Differences in study design
  - Differences in population structure
  - Differences in environmental exposures
  - False-positive?

# Need to study diverse populations

- Most GWAS have been done in populations of European ancestry

# Hindrance of long LD

- Long LD is valuable at the stage of hunting for associations

- But long LD is a hindrance at fine-mapping – potentially lots of hits



- < LD in African populations lead to > difficulty to detect signals in initial scan, but easier to fine-map causal variants

# Using GWAS To Study Infectious Disease Traits In Africa

## Benefits

- High prevalence of infection
- Identification of functionally relevant loci
- Fine mapping of causal variants

## Challenges

- High genomic diversity
- Pathogen genetic variation
- Lack of African genetic data & resources

# GWAS pipeline recap

```
┌─────────────────────┐      ┌─────────────────────┐      ┌─────────────────────┐
│  Sample collection  │ ───► │    Selection of     │ ───► │   Genotyping and    │
│   and phenotyping   │      │  genotyping array   │      │  genotype calling   │
└─────────────────────┘      └─────────────────────┘      └─────────────────────┘
                                                                      │
                             ┌─────────────────────┐                  ▼
                             │ Transethnic mapping │      ┌─────────────────────┐
                        ───► │  across multiple    │      │   Sample and SNP    │
                             │    populations      │      │   quality control   │
                             └─────────────────────┘      └─────────────────────┘
                                        │                            │
┌─────────────────────┐                 ▼                            ▼
│ Imputation against  │      ┌─────────────────────┐      ┌─────────────────────┐
│   sequence-level    │      │  Identify the causal│      │Population structure │
│   reference panels  │      │      variants       │      │      analysis       │
└─────────────────────┘      └─────────────────────┘      └─────────────────────┘
           ▲                            │                            │
           │                            ▼                            ▼
┌─────────────────────┐      ┌─────────────────────┐      ┌─────────────────────┐
│   Replication and   │ ◄─── │  Clusterplot        │ ◄─── │ Association analysis │
│    meta-analysis    │      │    checking         │      │                     │
└─────────────────────┘      └─────────────────────┘      └─────────────────────┘
                             │ Functional studies  │
```
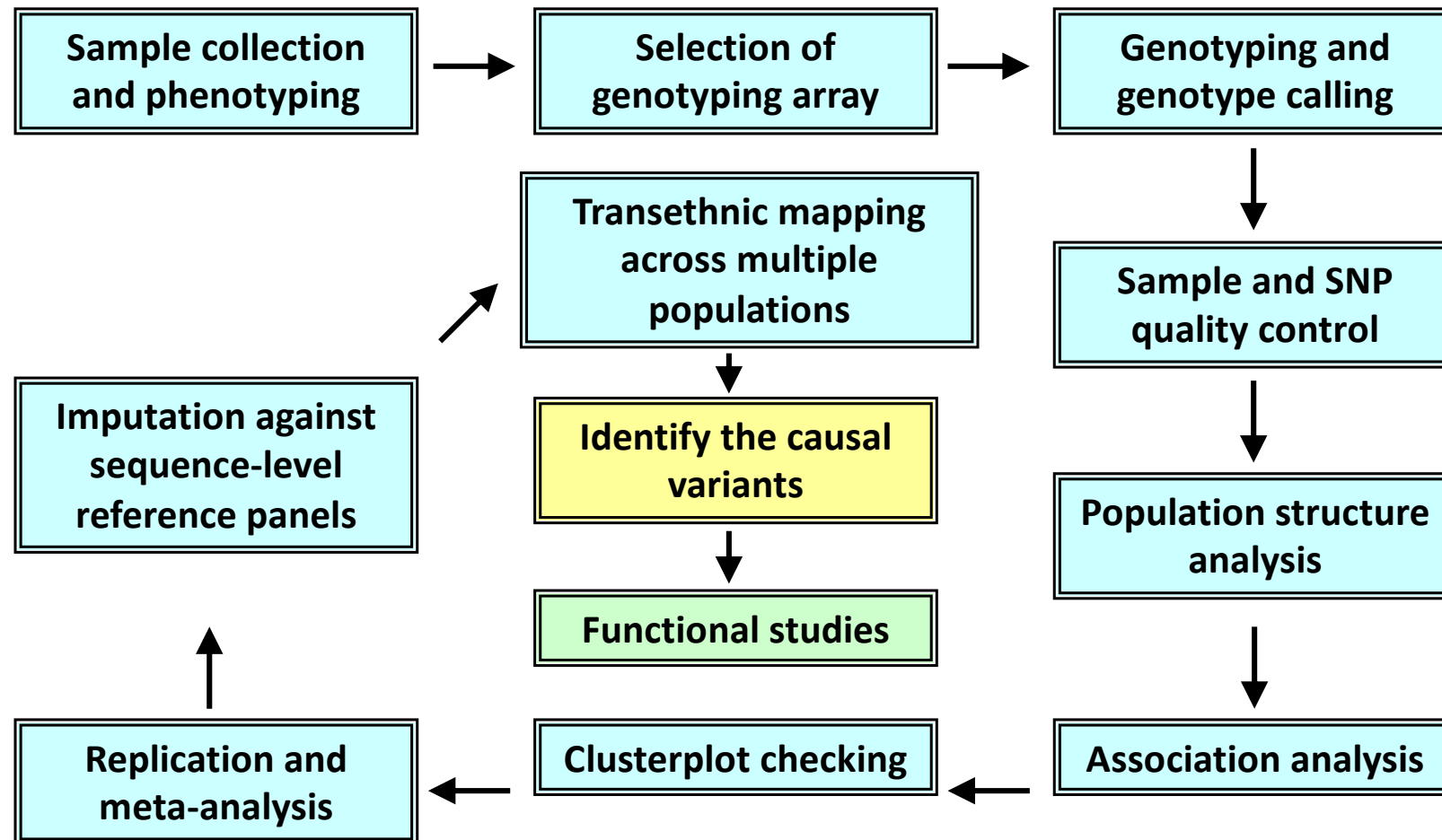
# Useful GWAS analysis tools

- SNP calling
  - Samtools : http://samtools.sourceforge.net
  - GATK : https://software.broadinstitute.org/gatk
  - OptiCall : https://opticall.bitbucket.io

- Data Imputation
  - Impute2: http://mathgen.stats.ox.ac.uk/impute/impute_v2.html
  - Beagle: https://faculty.washington.edu/browning/beagle/beagle.html
  - Sanger Imputation Server: https://imputation.sanger.ac.uk

- Publically available datasets:
  - 1000 Genomes: http://www.internationalgenome.org/data
  - Exac: http://exac.broadinstitute.org
  - UK10K: https://www.uk10k.org
  - HRC: http://www.haplotype-reference-consortium.org
  - African Genome Variation Project: https://www.sanger.ac.uk/science/collaboration/african-genome-variation-project
  - UKBioBank: https://www.ukbiobank.ac.uk

- Analysis:
  - Plink: http://zzz.bwh.harvard.edu/plink/
  - SNPtest: https://mathgen.stats.ox.ac.uk/genetics_software/snptest/snptest.html
  - GEMMA: http://www.xzlab.org/software.html
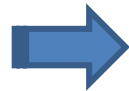  - R Packages: https://cran.r-project.org/web/packages/SNPassoc/SNPassoc.pdf
  - GCTA: http://cnsgenomics.com/software/gcta/#Overview

# Case in Point - GWAS of EBV in an African population

## Whole-genome association study of antibody response to Epstein-Barr virus in an African population: a pilot.

Sallah N[1,2], Carstensen T[1,3], Wakeham K[4,5], Bagni R[6], Labo N[7], Pollard MO[1,3], Gurdasani D[1,3], Ekoru K[1,3], Pomilla C[1,3], Young EH[1,3], Fatumo S[1,3,8], Asiki G[4], Kamali A[4], Sandhu M[1,3], Kellam P[2], Whitby D[7], Barroso I[1], Newton R[4].

# Genome-wide association workflow for EBV serological traits in the Uganda GPC
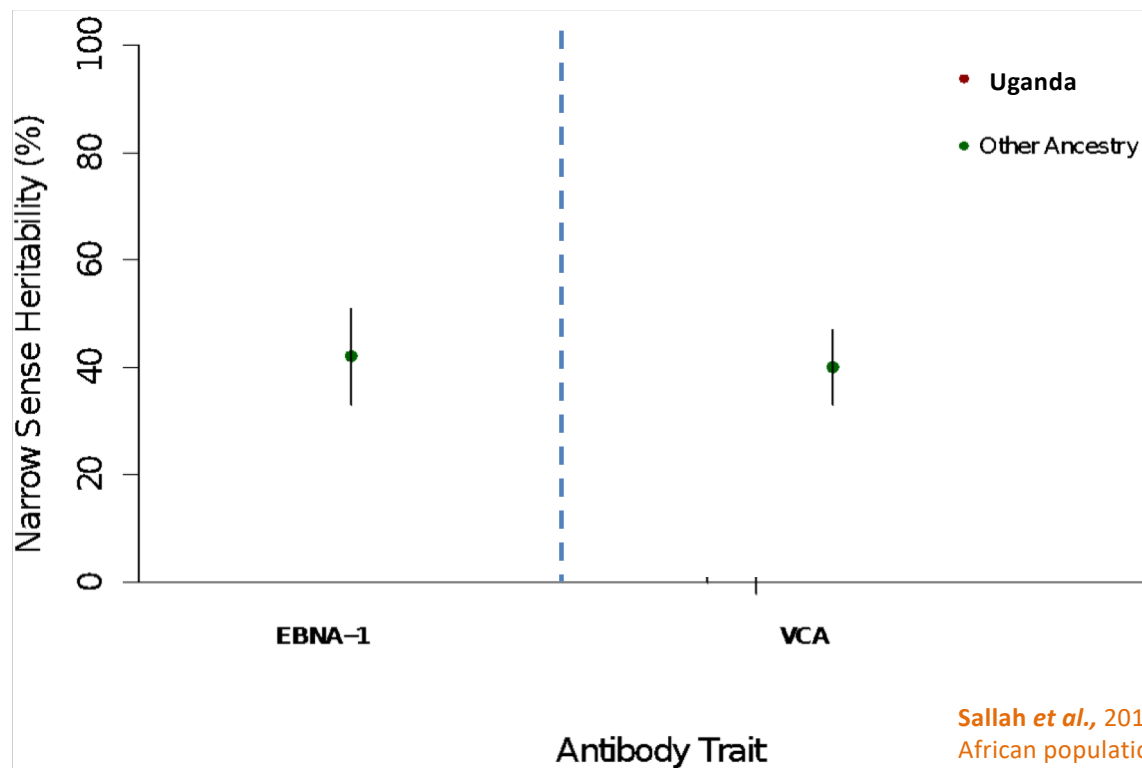
# Heritability of EBV IgG antibody response traits

- Proportion of variation in antibody responses due to host genetics
- $h^2$ based on genotype data using FaST-LMM (Heckerman, et al. 2016)
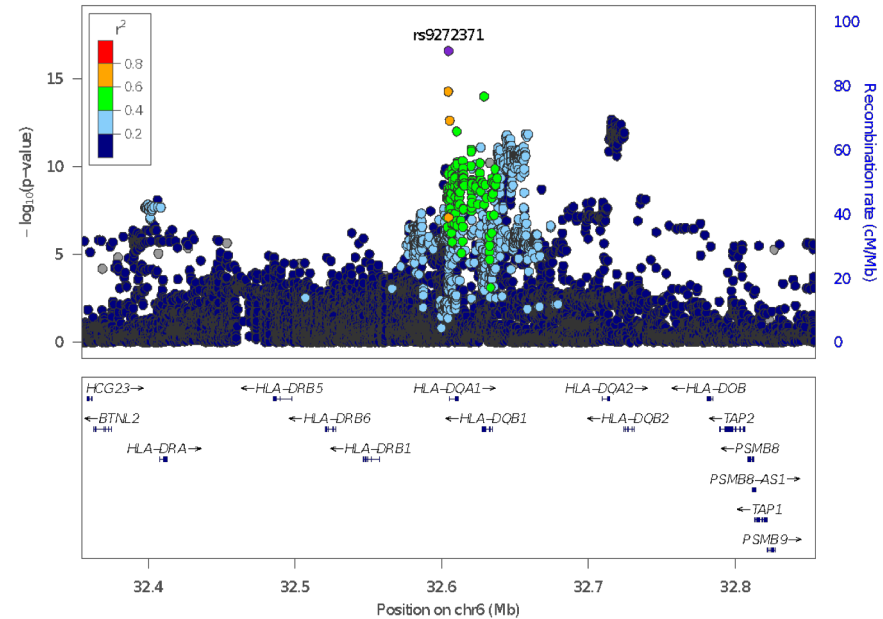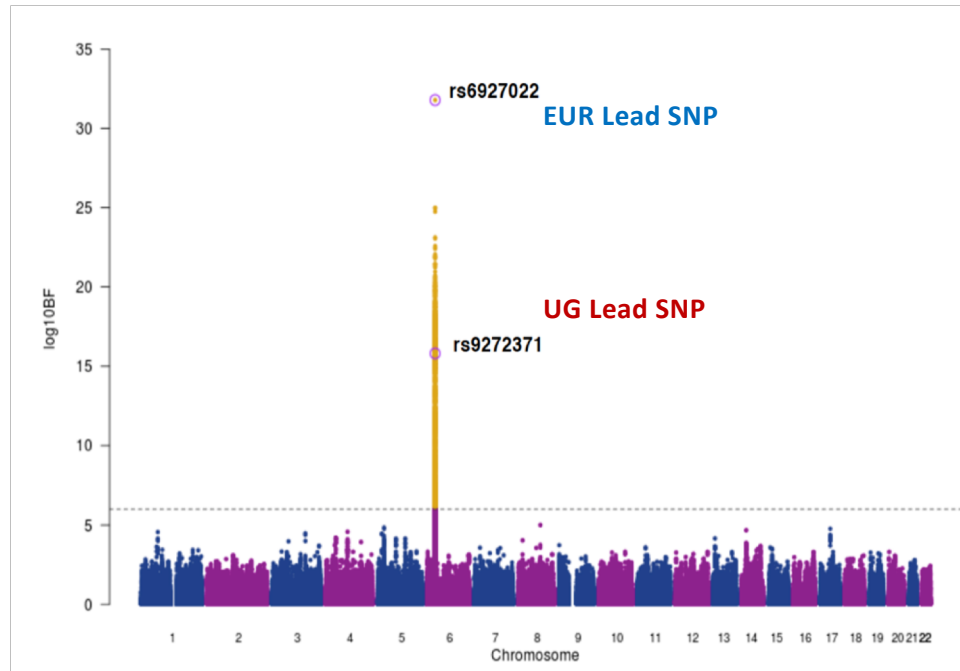- Adjustment for environmental correlation using GPS data

**Lower $h^2$ estimates in Ugandan population**

- Differences in in environmental variation
- Differences in gene-environment interactions
- Differences in variants or allele frequencies/effect sizes contributing to phenotypic variation



Sallah *et al.,* 2017, Whole-genome association study of antibody response to Epstein-Barr virus in an African population: A pilot. *Global Health, Epidemiology and Genomics, 2*. doi:10.1017/gheg.2017.16

# Distinct association signals in the *HLA* class II region for anti-EBNA-1 IgG response



Further analysis shows single signal in Eu and 2 signals in African population

Sallah *et al.,* 2017, Whole-genome association study of antibody response to Epstein-Barr virus in an African population: A pilot. *Global Health, Epidemiology and Genomics, 2.* doi:10.1017/gheg.2017.16
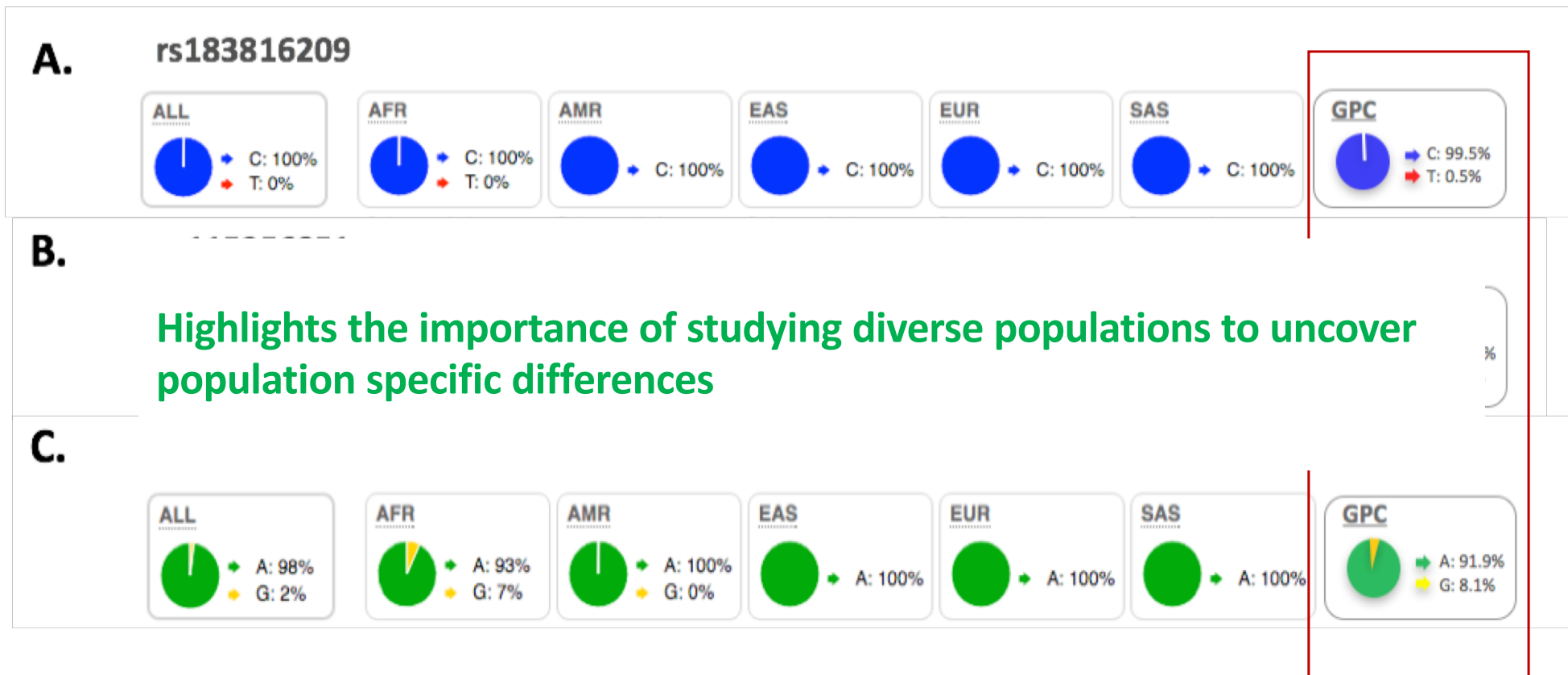
# Other variants identified that are African-Specific



**Highlights the importance of studying diverse populations to uncover population specific differences**

# Future Perspectives

- More data & resources from Africa & other diverse populations needed to leverage GWAS findings to uncover meaningful biological insights

- Large cohorts allow comprehensive analysis of infection – with host and pathogen genomes isolated from the same individuals



Host genome variation

Pathogen sequence variation

Biological function

Neneh Sallah



Tunisia
Morocco
Egypt
Mali
Senegal
Niger
The Gambia
Sudan
Burkina Faso
Nigeria
Guinea
Ethiopia
Sierra Leone
Ghana
Cameroon
Uganda
Kenya
Côte d'Ivoire
Benin
Democratic Republic of Congo
Tanzania
Malawi
Zambia
Mozambique
Zimbabwe
Namibia
Botswana
South Africa
Mauritius

## AN EVOLVING CONSORTIUM

The Human Heredity and Health in Africa (H3Africa) Initiative aims to develop a network of African labs that can do cutting-edge genomics and precision-medicine research with funding from the US National Institutes of Health (NIH) and the Wellcome Trust.

- ■ NIH primary award institution
- ▣ Wellcome primary award institution
- □ Collaborating institution

©nature