

# Genome-Wide Association Studies (GWAS)

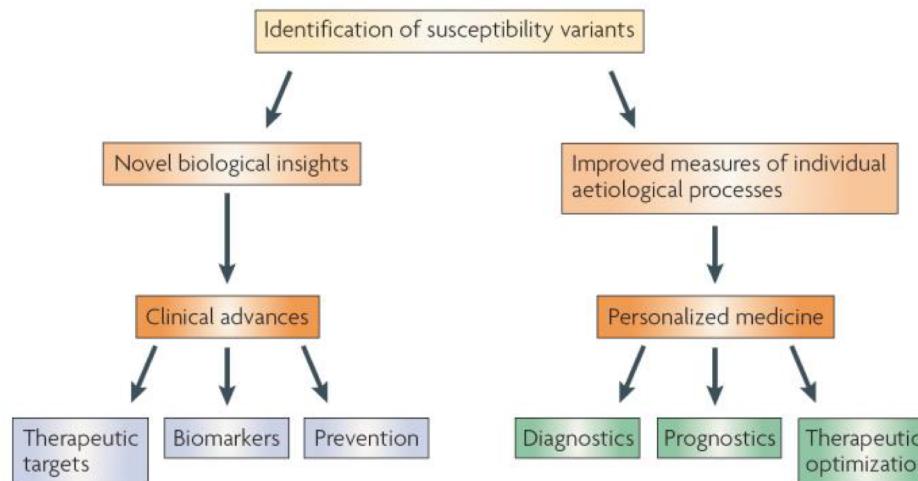
London School of Hygiene and  
Tropical Medicine

# Outline

- GWAS overview – Utility & Successes
- GWAS Study Design
- Post GWAS Follow up
- Studying host genomics in diverse populations e.g. Africa
- Recap

# 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



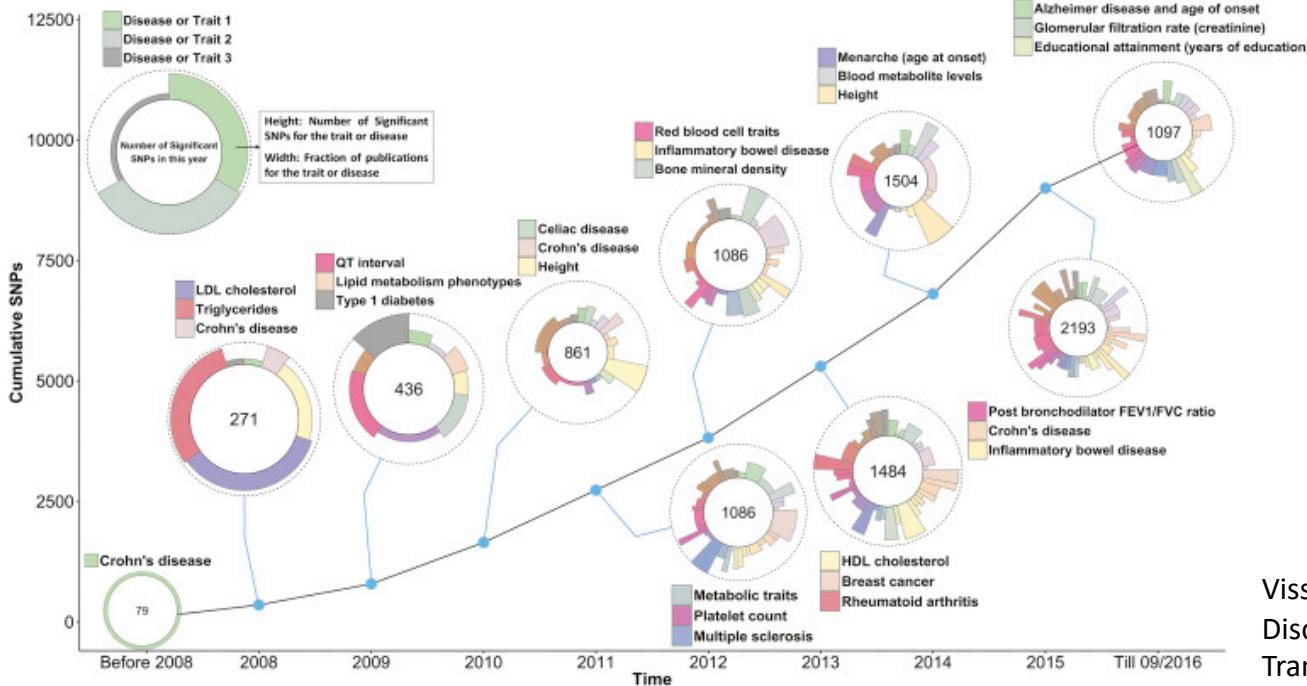
..ACTC**G**ACGATTACG**G**TACTTAGGAGCATA**CG**CTAC ..  
..ACTC**T**ACGATTACG**G**TACTTAGGAGCATA**CG**CTAC ..  
..ACTG**T**ACGATTACG**A**TACTTAGGAGCATA**T**GCTAC ..  
..ACTG**T**ACGATTACG**G**TACTTAGGAGCATA**T**GCTAC ..  
..ACTG**T**ACGATTACG**G**TACTTAGGAGCATA**T**GCTAC ..  
..ACTG**T**ACGATTACG**A**TACTTAGGAGCATA**GG**CTAC ..  
..ACTG**T**ACGATTACG**A**TACTTAGGAGCATA**GG**CTAC ..  
..ACTG**T**ACGATTACG**G**TACTTAGGAGCATA**T**GCTAC ..  
..ACTG**T**ACGATTACG**A**TACTTAGGAGCATA**GG**CTAC ..

SNPs may have 2, 3 or 4 alleles (most are biallelic)

# Lots of Success

May 2018 ( $p \leq 5 \times 10^{-8}$ )

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

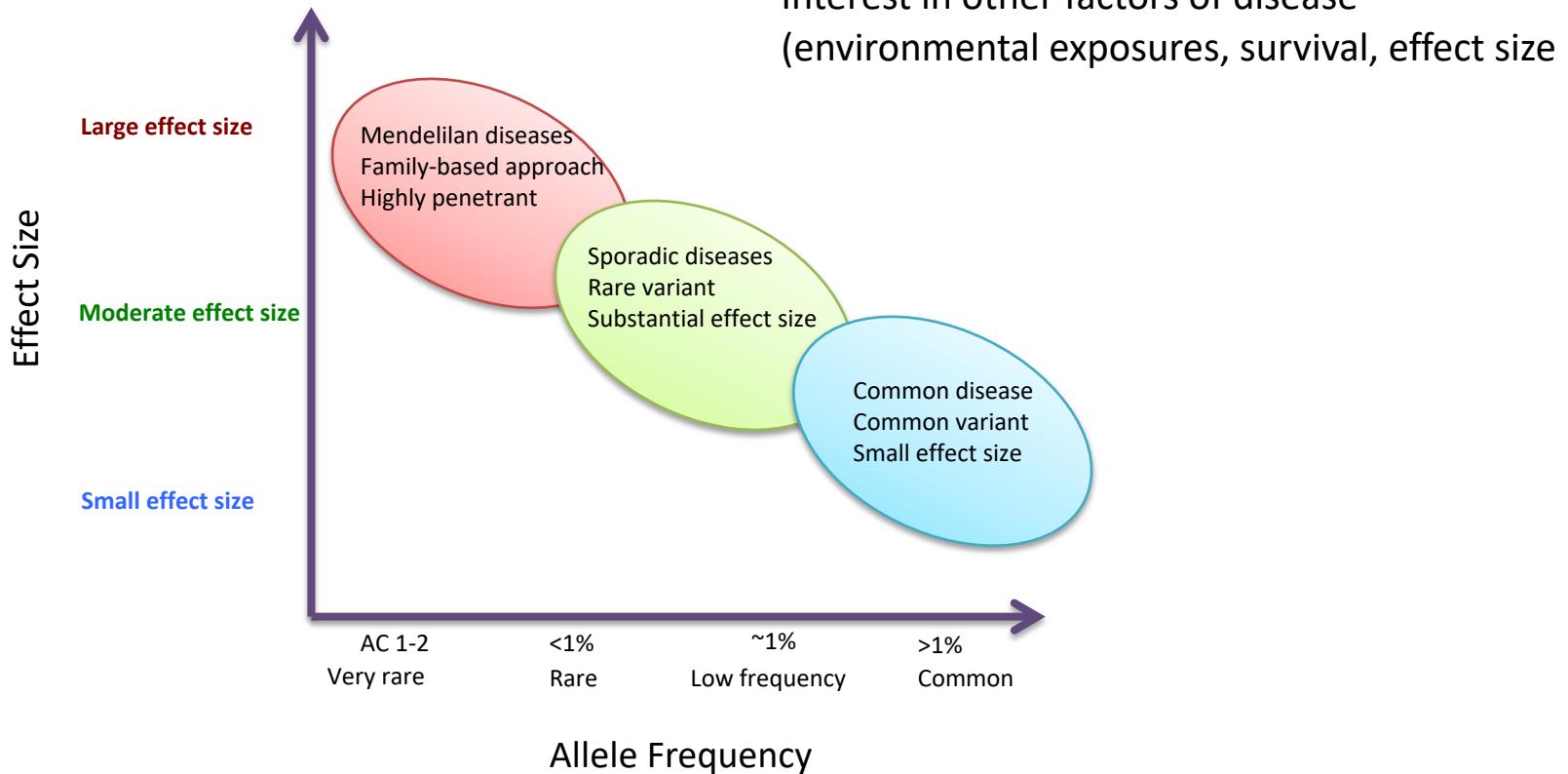


NHGRI GWA Catalog  
[www.genome.gov/GWASStudies](http://www.genome.gov/GWASStudies)  
[www.ebi.ac.uk/fgpt/gwas/](http://www.ebi.ac.uk/fgpt/gwas/)

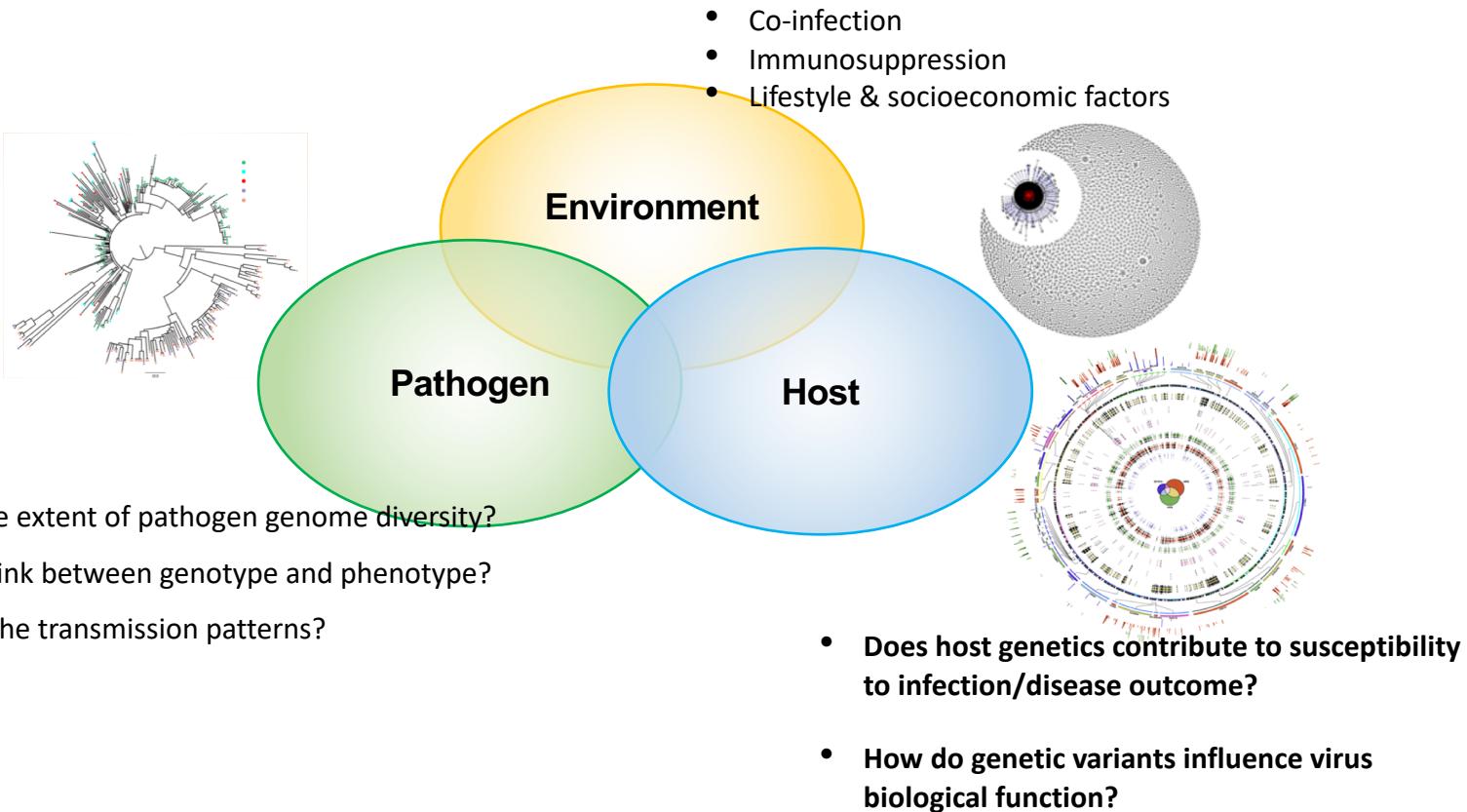


Visscher et al; 2017, 10 Years of GWAS Discovery: Biology, Function, and Translation

# Variant Identification



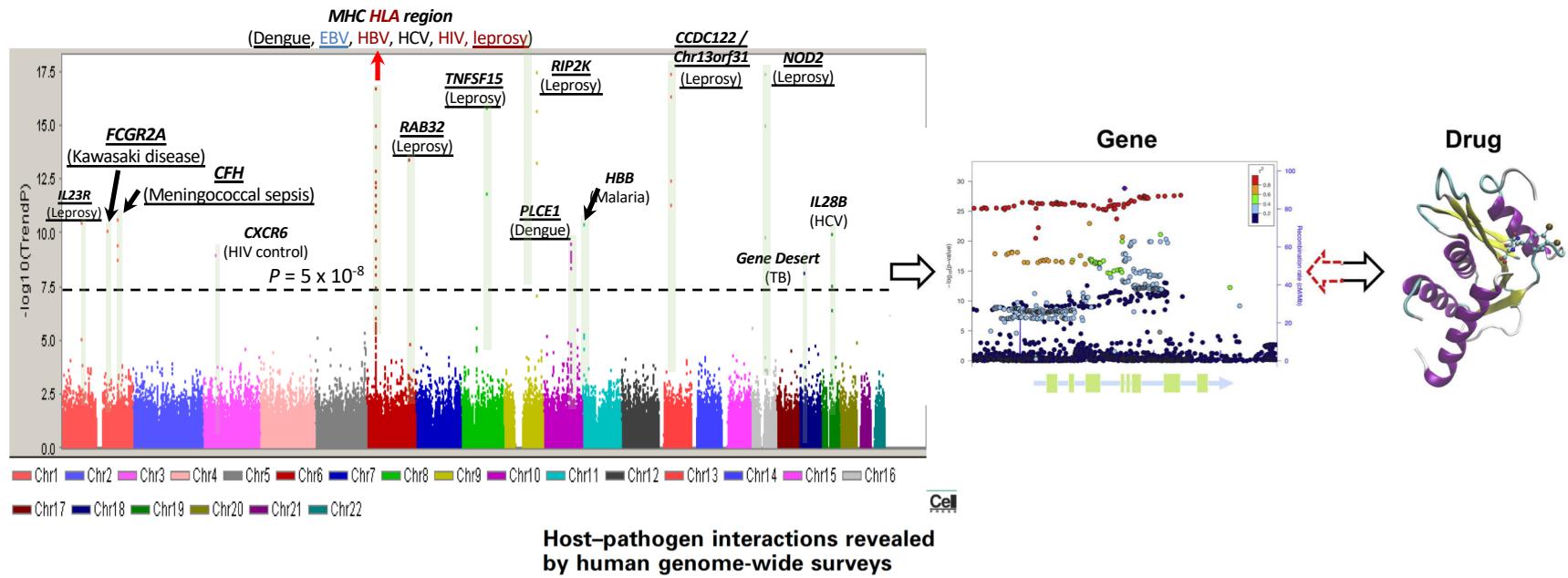
# Multifactorial determinants of pathogenesis & clinical outcome



# 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



# Infectious Diseases - Insights

Nat Genet. 2009 Jun;41(6):657-65. Epub 2009 May 24.

## Genome-wide and fine-resolution association analysis of malaria in West Africa.

Jallow M, Teo YY, Small KS, Rockett KA, Deloukas P, Clark TG, Kivinen K, Bojang KA, Conway DJ, Pinder M, Sirugo G, Sisay-Joof F, Usen S, Auburn S, Bumpstead SJ, Campino S, Coffey A, Dunham A, Fry AE, Green A, Gwilliam R, Hunt SE, Inouye M, Jeffreys AE, Mendy A, Palotie A, Potter S, Ragoussis J, Rogers J, Rowlands K, Somaskantharajah E, Whittaker P, Widden C, Donnelly P, Howie B, Marchini J, Morris A, SanJoaquin M, Achidi EA, Aqbenvega T, Allen A, Amodu O, Corran P, Dijimde A, Dolo A, Doumbo OK, Drakeley C, Dunstan S, Evans J, Farrar J, Fernando D, Hien TT, Horstmann RD, Ibrahim M, Karunaweera N, Kokwaro G, Koram KA, Lemnge M, Makani J, Marsh K, Michon P, Modiano D, Molyneux ME, Mueller I, Parker M, Pesu N, Plowe CV, Pujalal O, Reeder J, Reyburn H, Riley EM, Sakuntabhai A, Singhvilanont P, Sirima S, Tall A, Taylor TE, Thera M, Trovo-Blomberg M, Williams TN, Wilson M, Kwiatkowski DP; Wellcome Trust Case Control Consortium; Malaria Genomic Epidemiology Network.

Nat Genet. 2010 Sep;42(9):739-41. Epub 2010 Aug 8.

## Genome-wide association analyses identifies a susceptibility locus for tuberculosis on chromosome 18q11.2.

Thye T, Vanberg FO, Wong SH, Owusu-Dabo E, Osei I, Gyapong J, Sirugo G, Sisay-Joof F, Enimil A, Chinbuah MA, Floyd S, Warndorff DK, Sichali L, Malema S, Crampin AC, Ngwira B, Teo YY, Small K, Rockett K, Kwiatkowski D, Fine PE, Hill PC, Newport M, Lienhardt C, Adegbola RA, Corrah T, Ziegler A; African TB Genetics Consortium; Wellcome Trust Case Control Consortium, Morris AP, Meyer CG, Horstmann RD, Hill AV.

Nat Genet. 2011 Oct 16;43(11):1139-41. doi: 10.1038/ng.960.

## Genome-wide association study identifies susceptibility loci for dengue shock syndrome at MICB and PLCE1.

Khor CC, Chau TN, Pang J, Davila S, Long HT, Ong RT, Dunstan SJ, Wills B, Farrar J, Van Tram T, Gan TT, Binh NT, Tri le T, Lien le B, Tuan NM, Tham NT, Lanh MN, Nguyet NM, Hieu NT, Van N Vinh Chau N, Thuy TT, Tan DE, Sakuntabhai A, Teo YY, Hibberd ML, Simmons CP.

Nature. 2012 Aug 15. doi: 10.1038/nature11334. [Epub ahead of print]

## Genome-wide association study indicates two novel resistance loci for severe malaria.

Timmann C, Thye T, Vens M, Evans J, May J, Ehmen C, Sievertsen J, Muntau B, Ruge G, Loaq W, Ansong D, Antwi S, Asafo-Adjei E, Nguah SB, Kwakye KO, Akoto AO, Sylverken J, Brendel M, Schuldert K, Lolei C, Franke A, Meyer CG, Aqbenyeqa T, Ziegler A, Horstmann RD.

Proc Natl Acad Sci U S A. 2012 Aug 7;109(32):13052-7. Epub 2012 Jul 23.

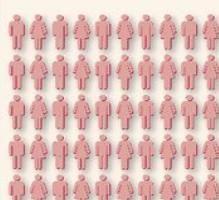
## Sequence-based association and selection scans identify drug resistance loci in the Plasmodium falciparum malaria parasite.

Park DJ, Lukens AK, Neafsey DE, Schaffner SF, Chang HH, Valim C, Ribacke U, Van Tyne D, Galinsky K, Galligan M, Becker JS, Ndiaye D, Mboup S, Wiegand RC, Hartl DL, Sabeti PC, Wirth DF, Volkman SK.

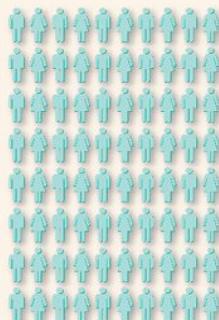
# GWAS Workflow

## Sample Collection & Phenotype Determination

Phenotypes: drug efficacy/toxicity



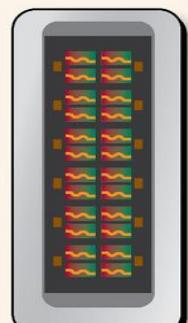
Case



Control

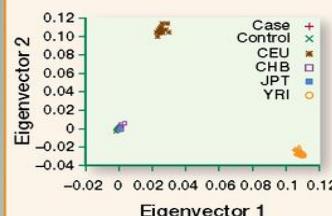
## Genotyping or Whole genome sequencing

Genotype with chip that contains probes for hundreds of thousands SNPs



## Quality Control

1. Sample quality
2. Relatedness
3. Population stratification

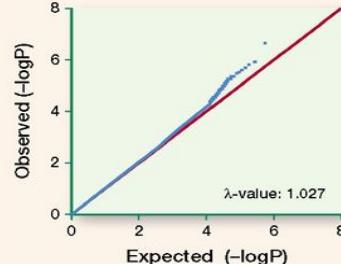


### SNP QC

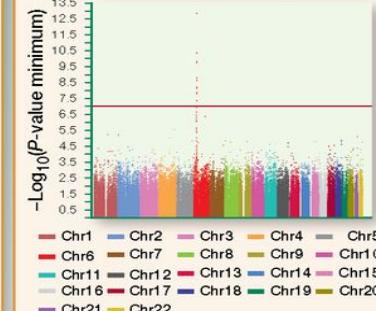
1. Genotype quality
2. Deviation from normal distribution
3. Allele frequency

## GWAS

### Quantile-quantile plot



### Manhattan plot



### Validation study

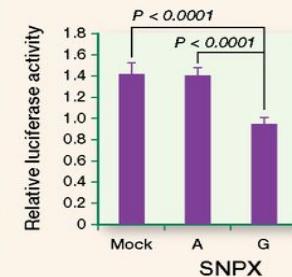
Replicate findings with independent sample set

## 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)

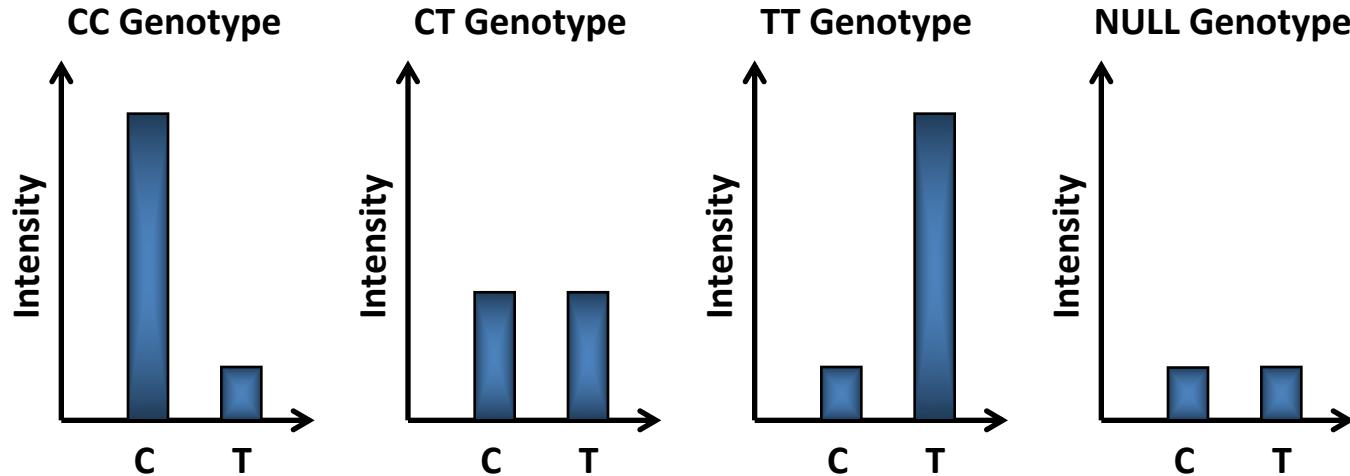
### 2. Functional analysis a. EMSA or b. reporter assay



### 3. Other analysis

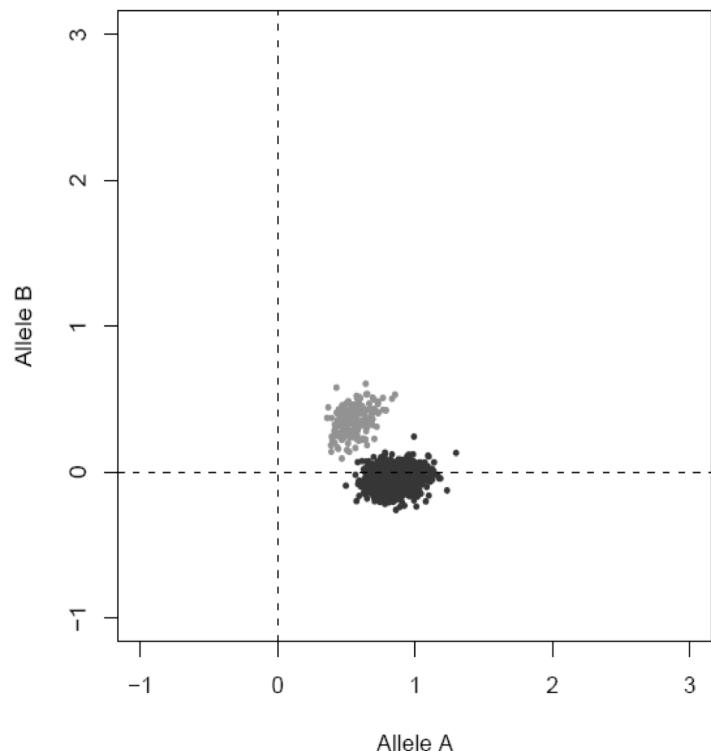
Gene-based analysis, pathway analysis, polygenic risk estimation, SNP-SNP interaction, etc.

# Raw data is not a genotypes, but Allelic hybridization

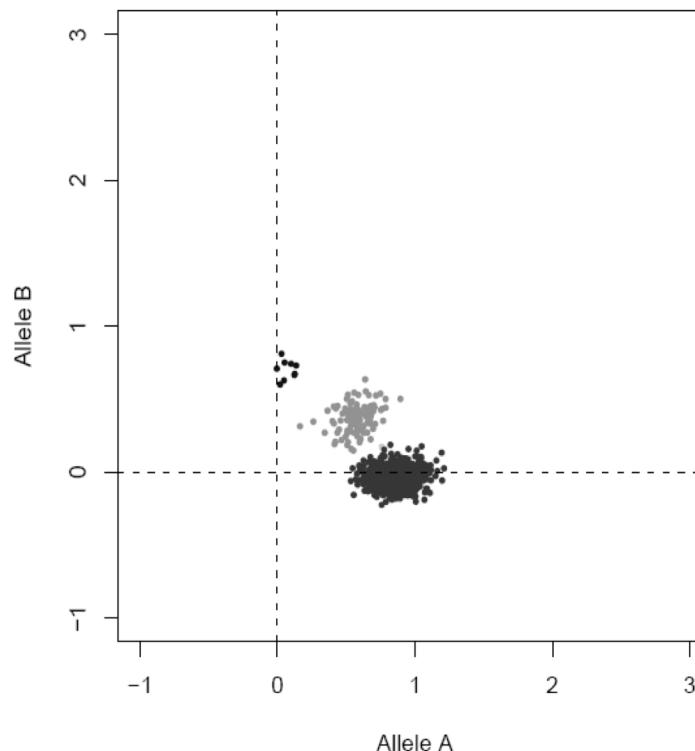


# Genotype Calling

Controls

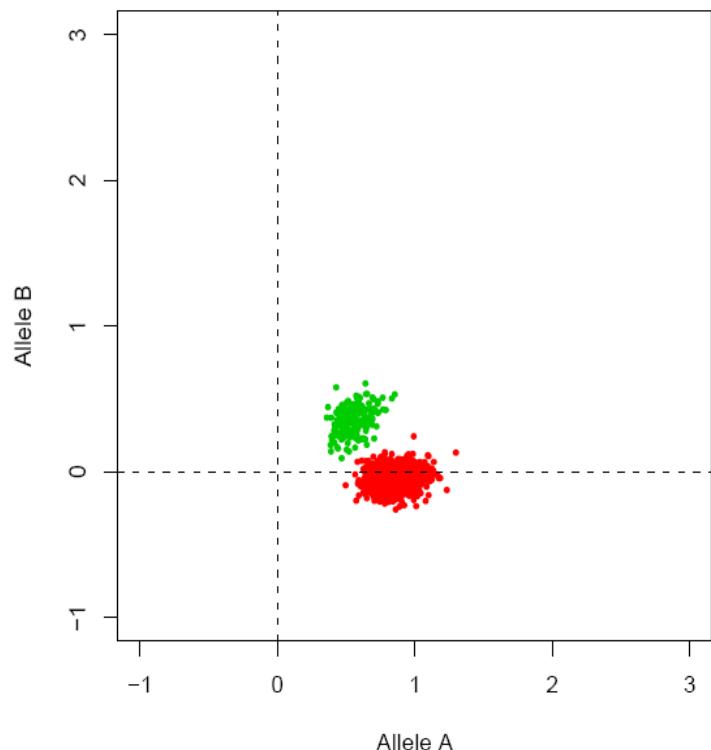


Cases

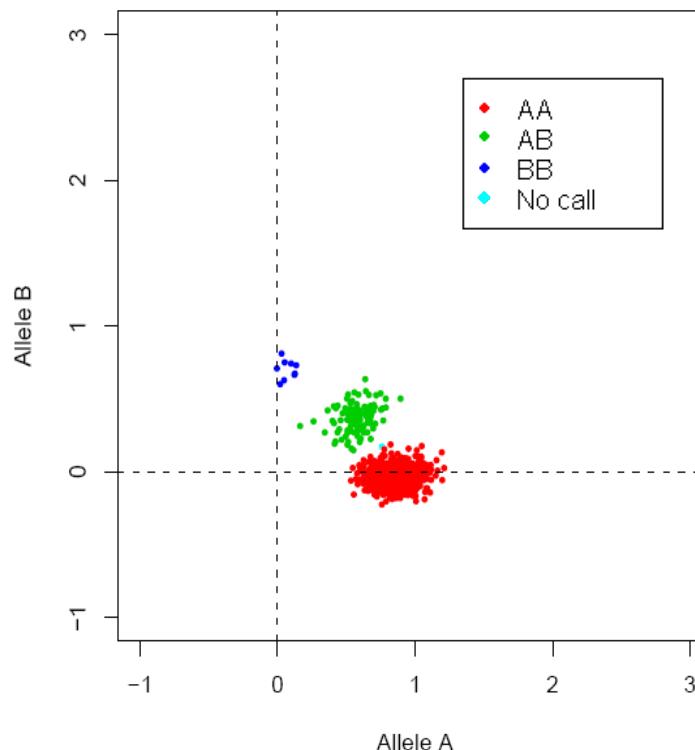


# Genotype Calling

Controls



Cases



# Need for high quality data

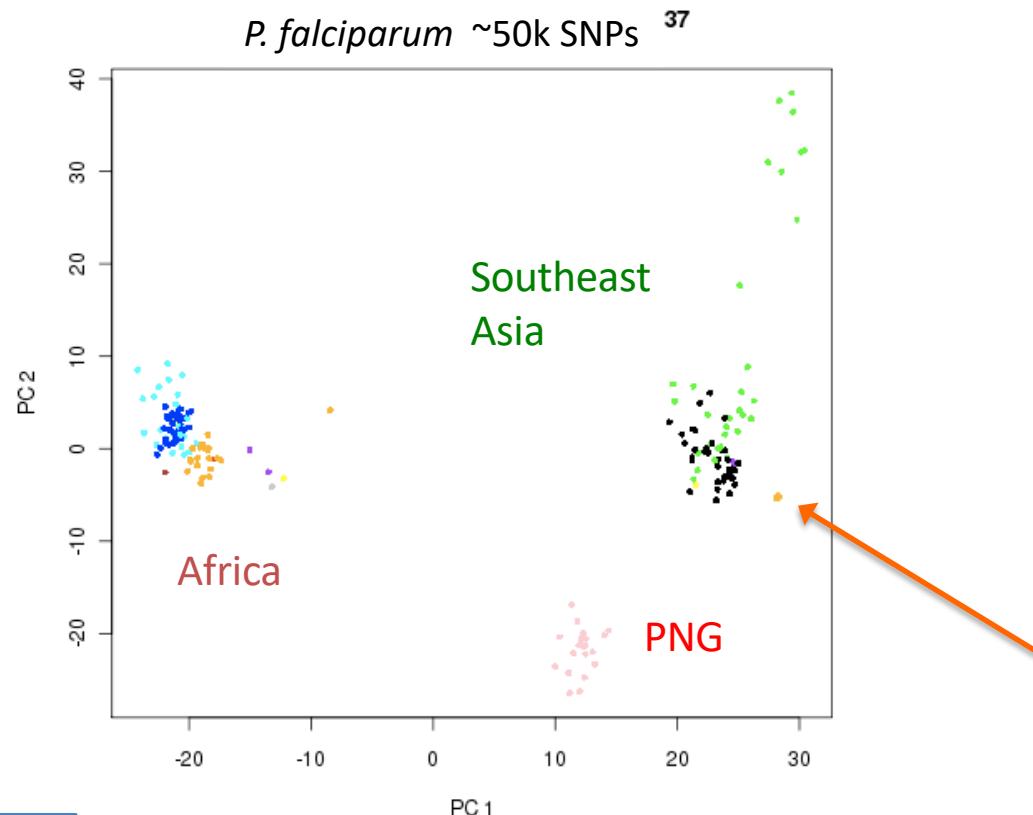
- Number of variants assayed  $\Rightarrow$  errors and genotype or sequence miscalls are bound to happen
- If problematic samples not identified and excluded, they can affect the results of the entire experiment
- If SNPs with erroneous genotyping or sequencing not identified and excluded, can produce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

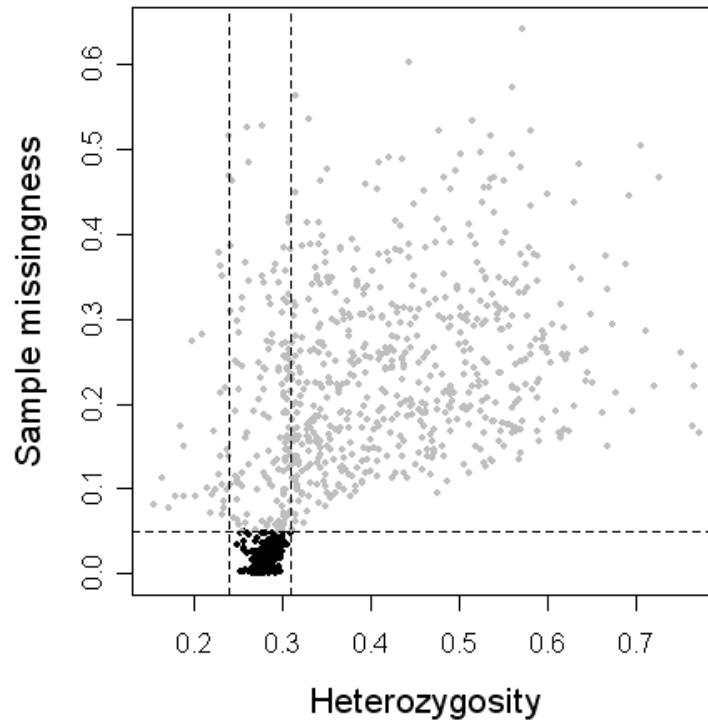
# Sample QC

- Identify samples that are outliers



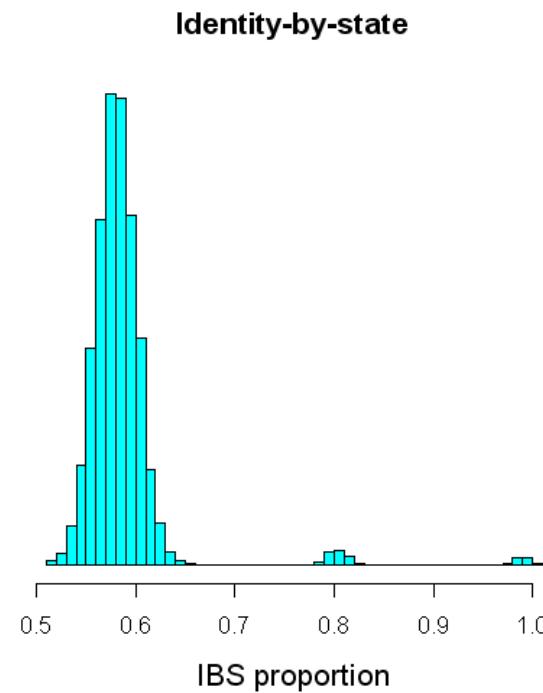
# Sample QC

- Identify SNPs with high rates of missingness and heterozygosity

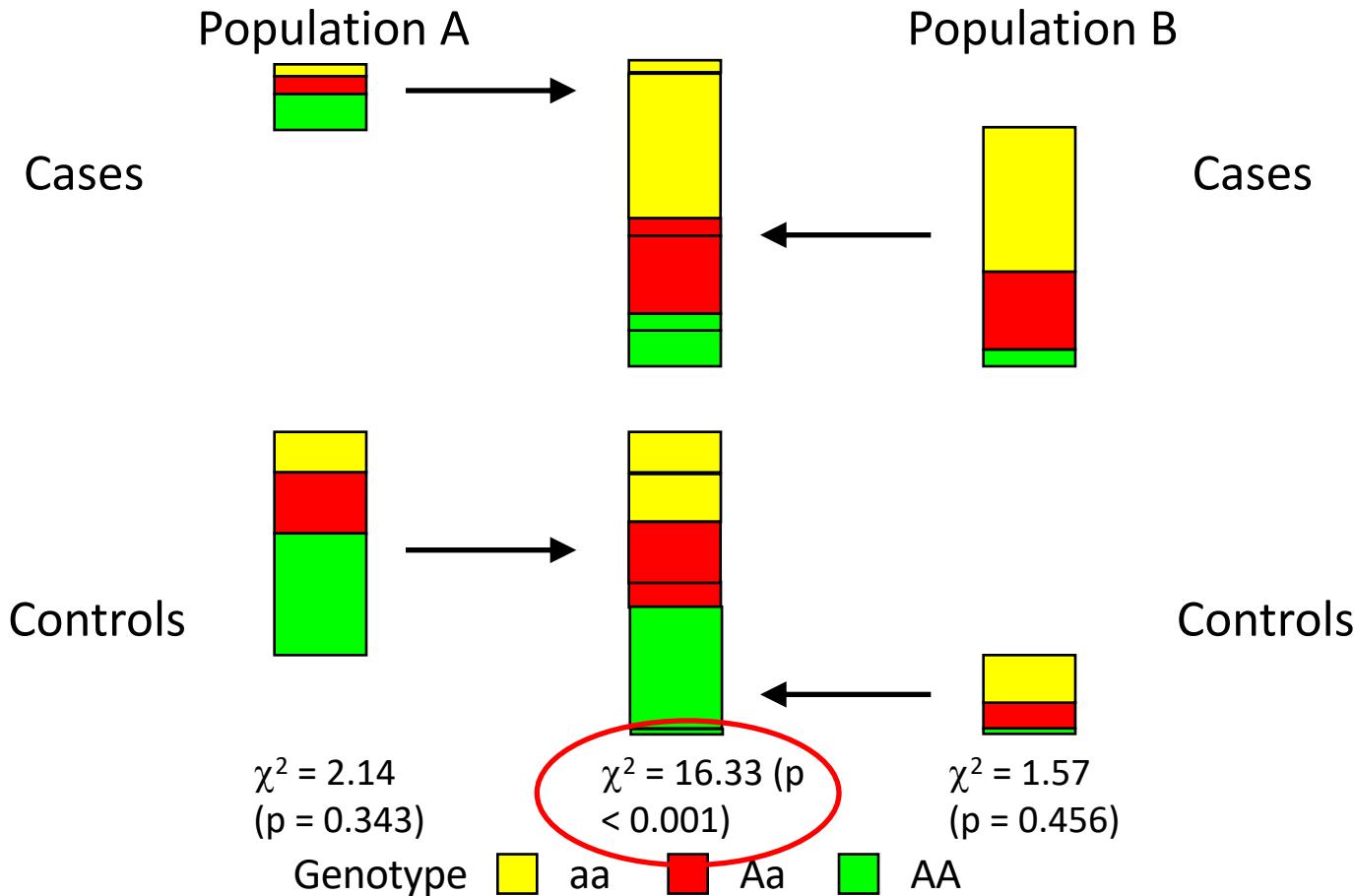


# Sample QC

- Identify samples with high rates of missingness and heterozygosity
- Identify related / duplicated samples



# Effects of population structure in an association study

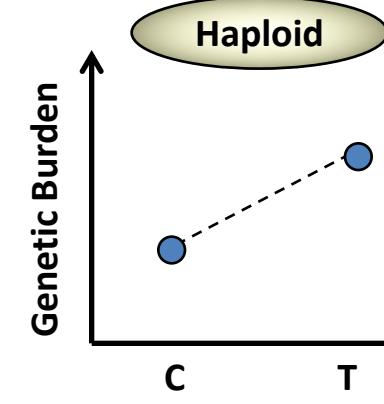
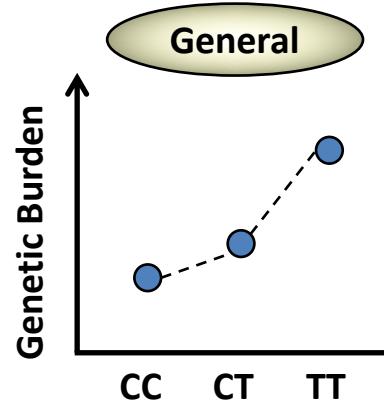
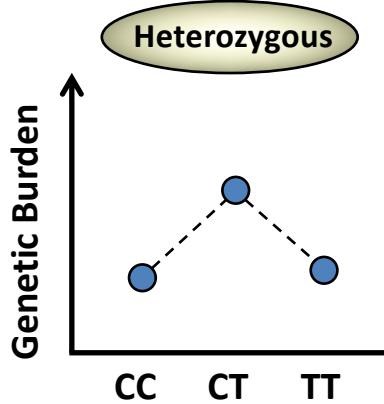
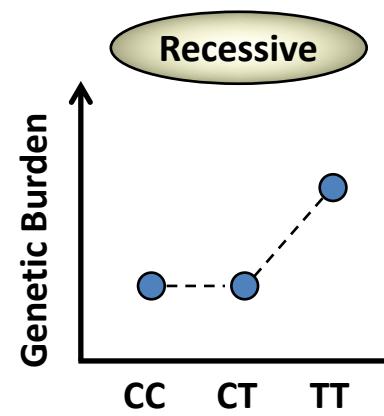
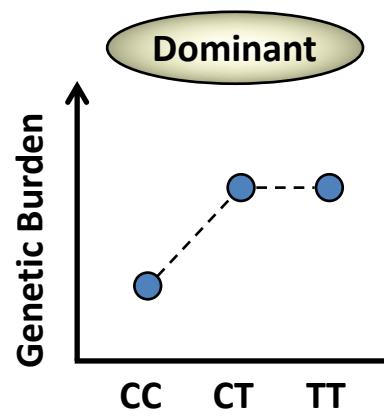
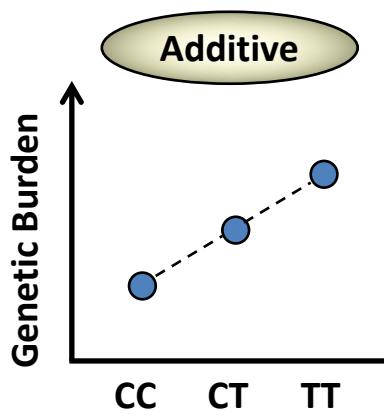


# Testing for associations using regression models

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

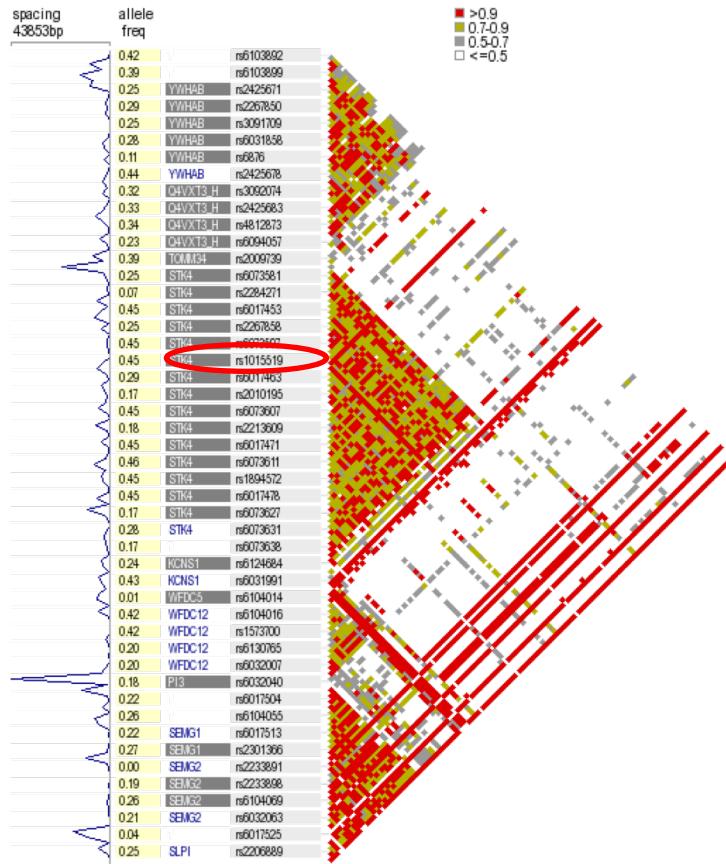
# Genetic models tested in a regression framework



# Association Studies

## Direct Association

Tests the genetic variant directly responsible for causing the disease.



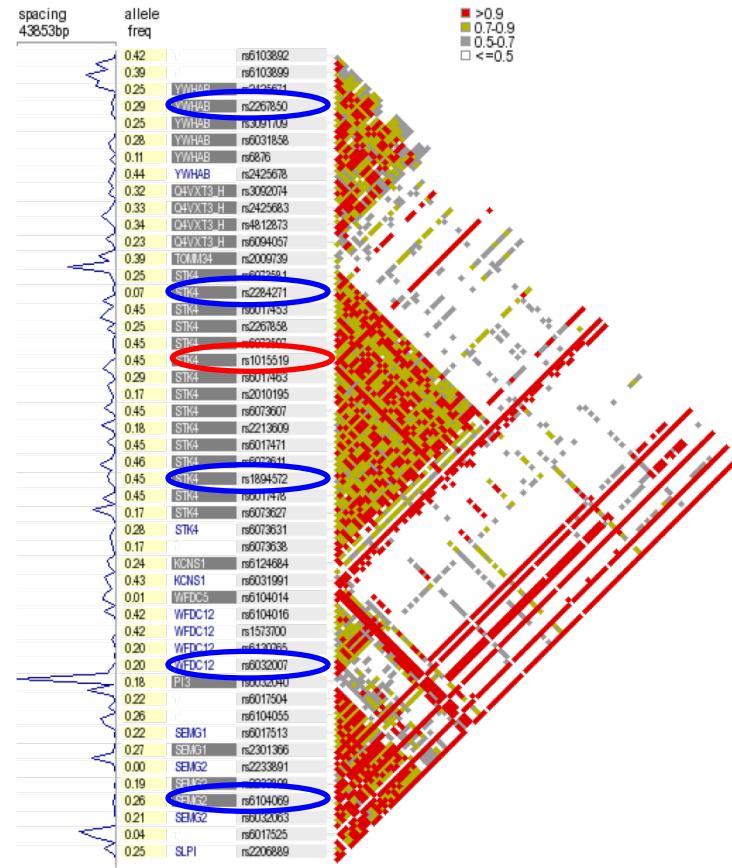
# Association Studies

## Direct Association

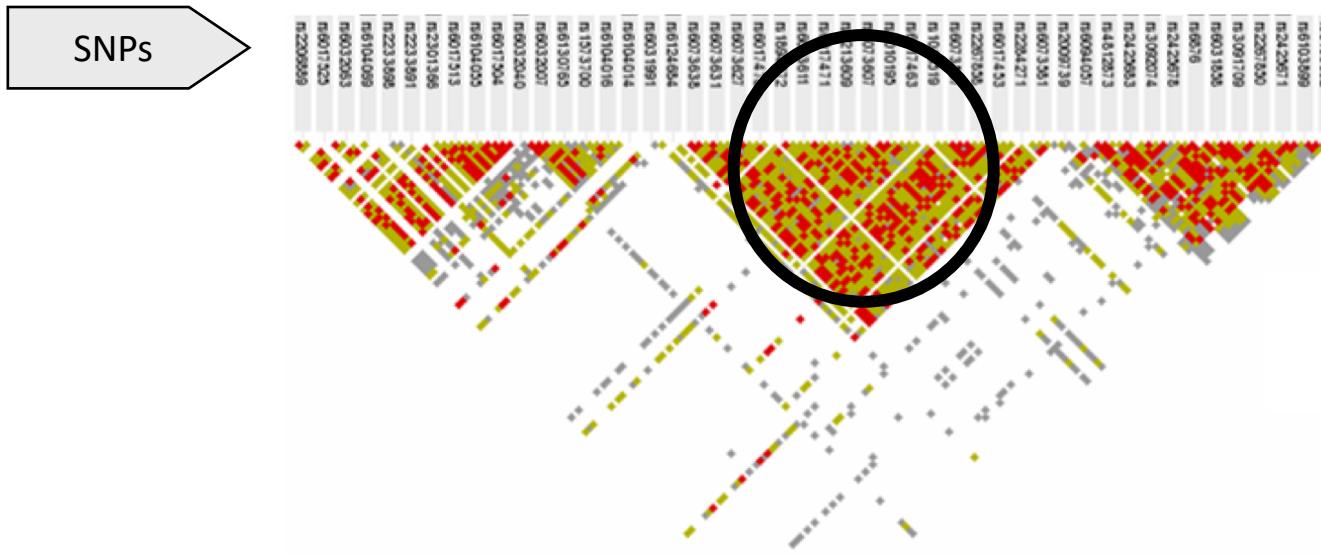
Tests the genetic variant directly responsible for causing the disease.

## 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



- > 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

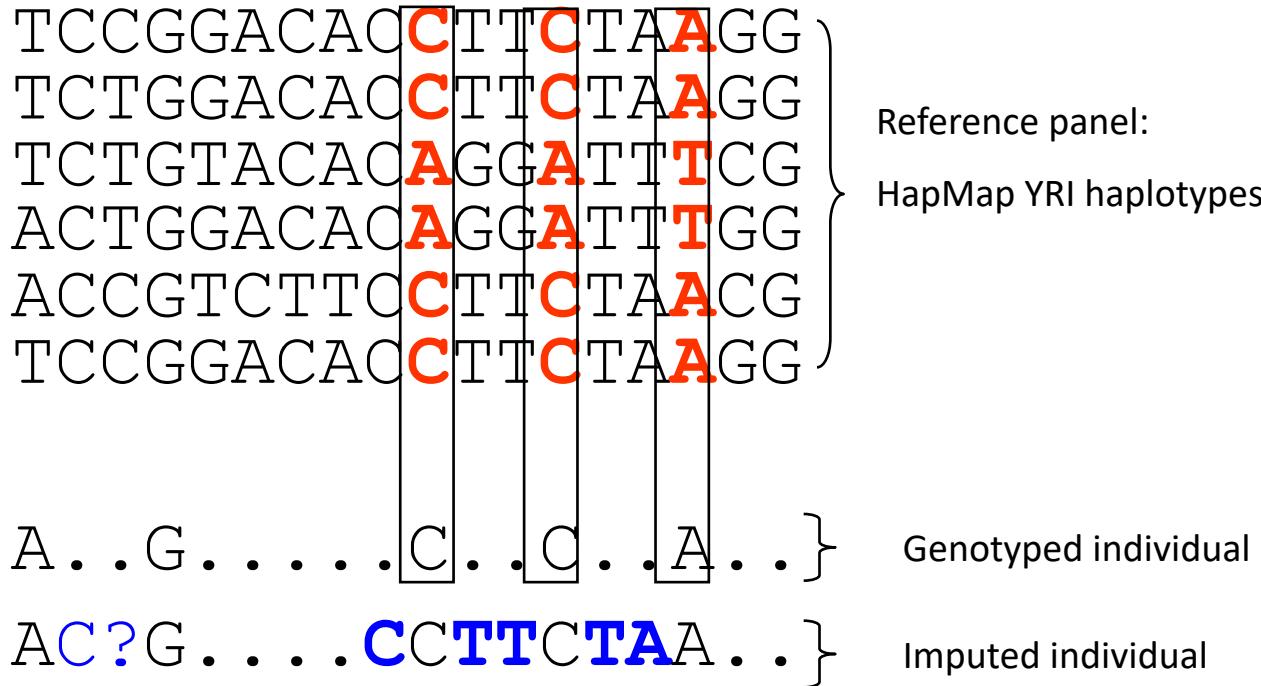
```
.. .ACTCGACGATTACGGTACTTAGGAGCATATGCTAC ..  
.. .ACTCTACGATTACGGTACTTAGGAGCATATGCTAC ..  
.. .ACTGTACGATTACGATACTTAGGAGCATAGGCTAC ..  
.. .ACTGAACGATTACGGTACTTAGGAGCATATGCTAC ..  
.. .ACTGTACGATTACGGTACTTAGGAGCATATGCTAC ..  
.. .ACTGTACGATTACGATACTTAGGAGCATAGGCTAC ..  
.. .ACTGGACGATTACGGTACTTAGGAGCATAGGCTAC ..  
.. .ACTGTACGATTACGGTACTTAGGAGCATATGCTAC ..
```

- 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?

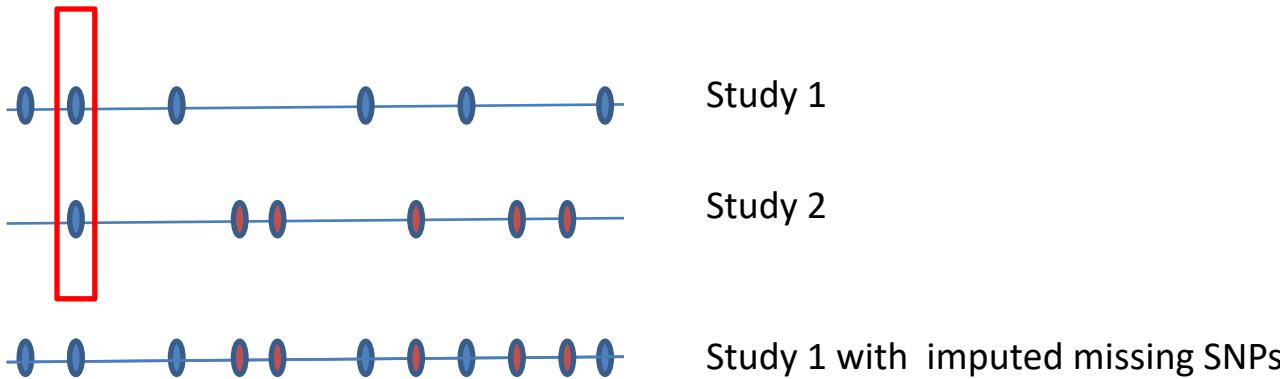
- ~10m genetic variants in the human genome, costly to genotype everything (pre-2012?)
- LD  $\Rightarrow$  Reduced amount of genotyping required
- The availability of whole genome sequencing on large numbers of samples makes LD redundant

# Imputation



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

# Imputation



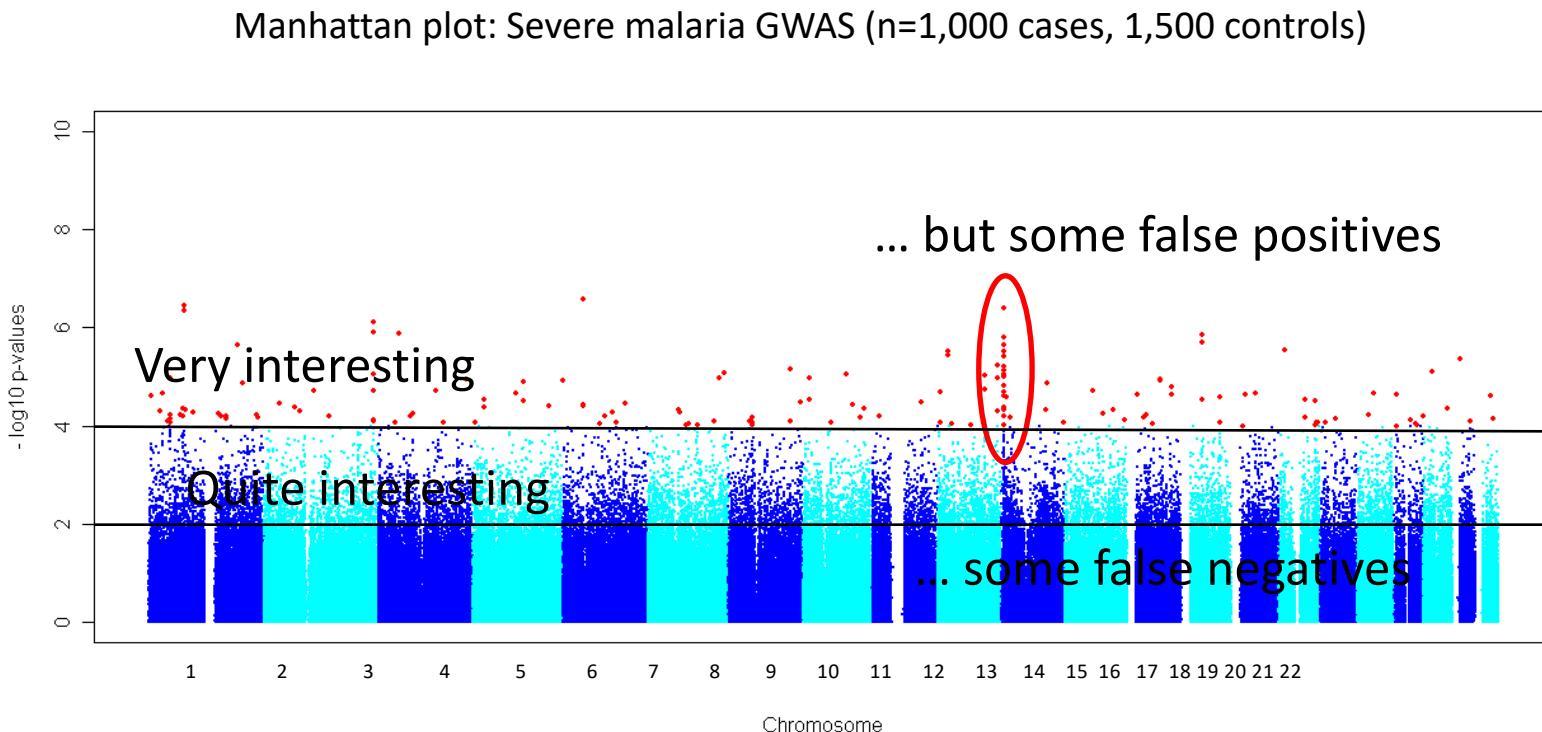
- Imputation
  - Requires GWAS genotypes to be used as scaffold
  - Requires reference datasets (e.g. [www.hapmap.org](http://www.hapmap.org); [www.1000genomes.org](http://www.1000genomes.org)) where the LD (correlation) between SNPs is known and allows imputation of genotypes for variants not typed 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?

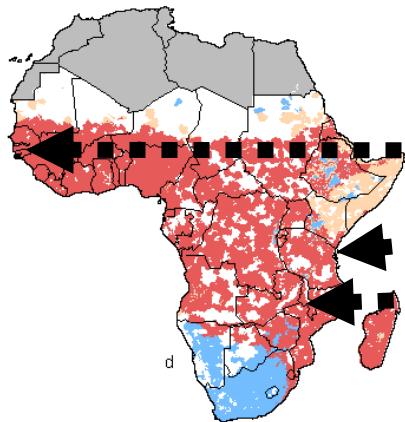
- To predict missing genotypes that haven't been directly typed
  - **Increased power.** The reference panel is more likely to contain the causal variant (or a better tag) than a GWAS array.
  - **Fine-mapping.** Imputation provides a high-resolution overview of an association signal across a locus.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



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

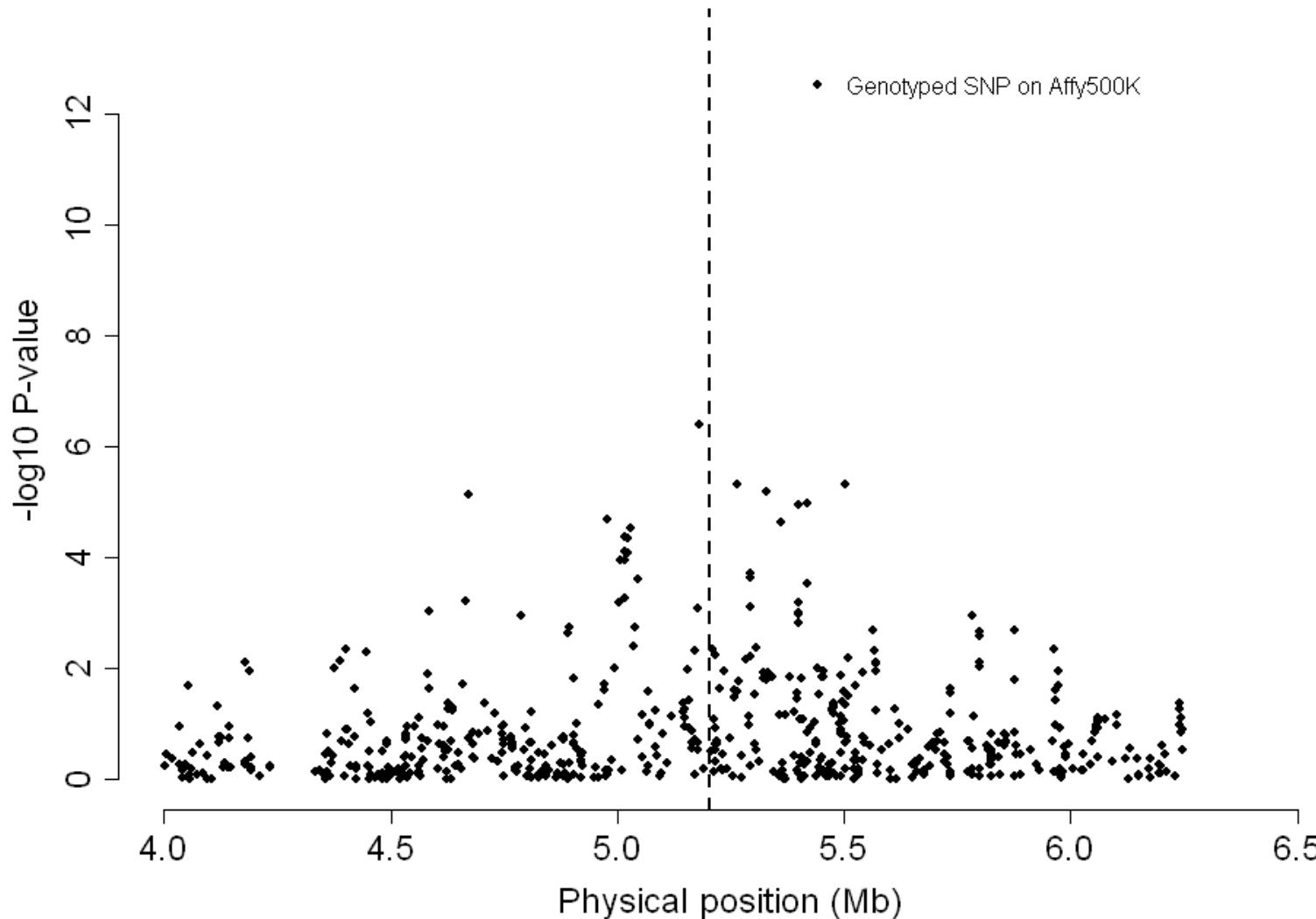


Relative risk of severe malaria in children with HbS AS genotype	
Gambia	0.11
Kenya	0.17
Malawi	0.08

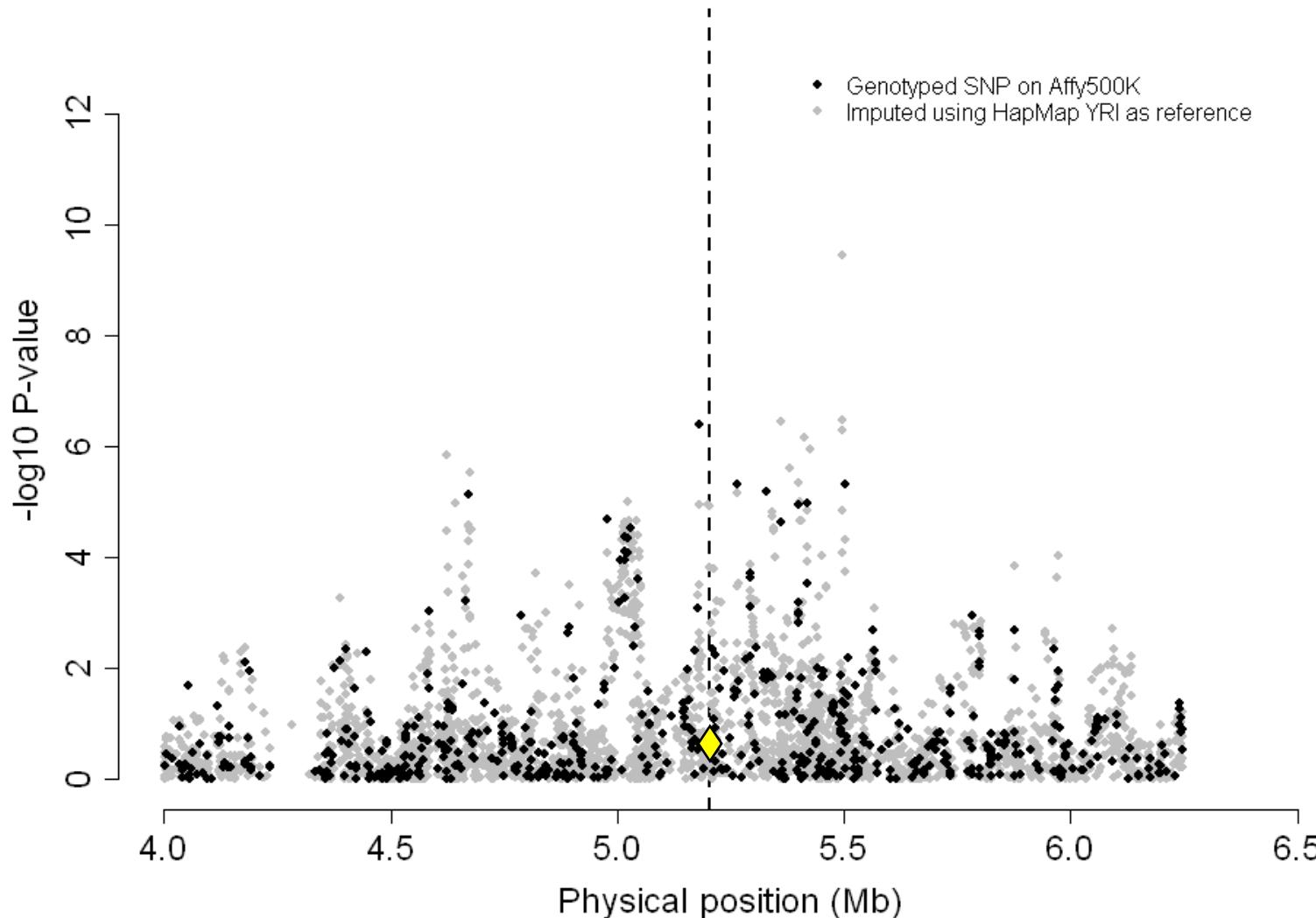
$$P = 2 \times 10^{-31} (n = 3630)$$

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

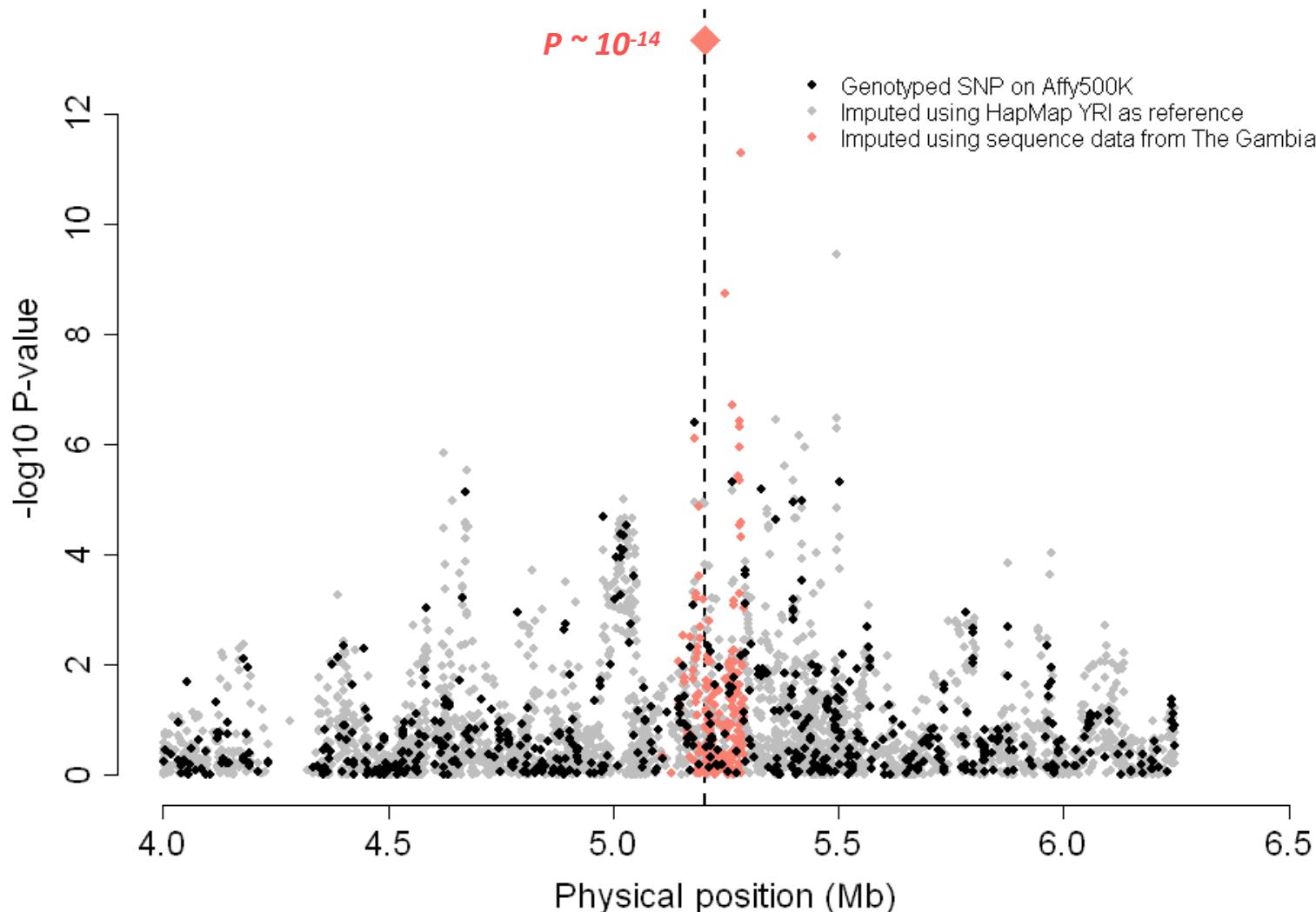
## Signals of malaria association in chromosome 11 in The Gambia



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



# Signals of malaria association in chromosome 11 in The Gambia



# Going beyond GWAS

- Need to validate and confirm findings
  - Replication studies and meta analysis
- If using genotyping arrays, fine-mapping the causal variant
  - Targeted-resequencing
  - Transethnic mapping
- Functional studies

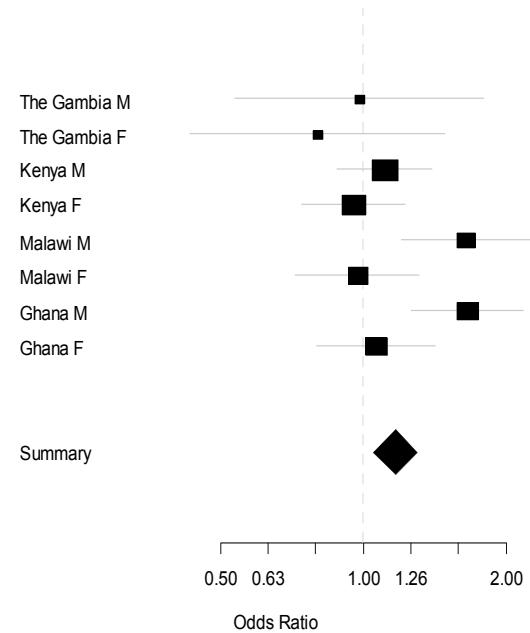
# Replication

- 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



# 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

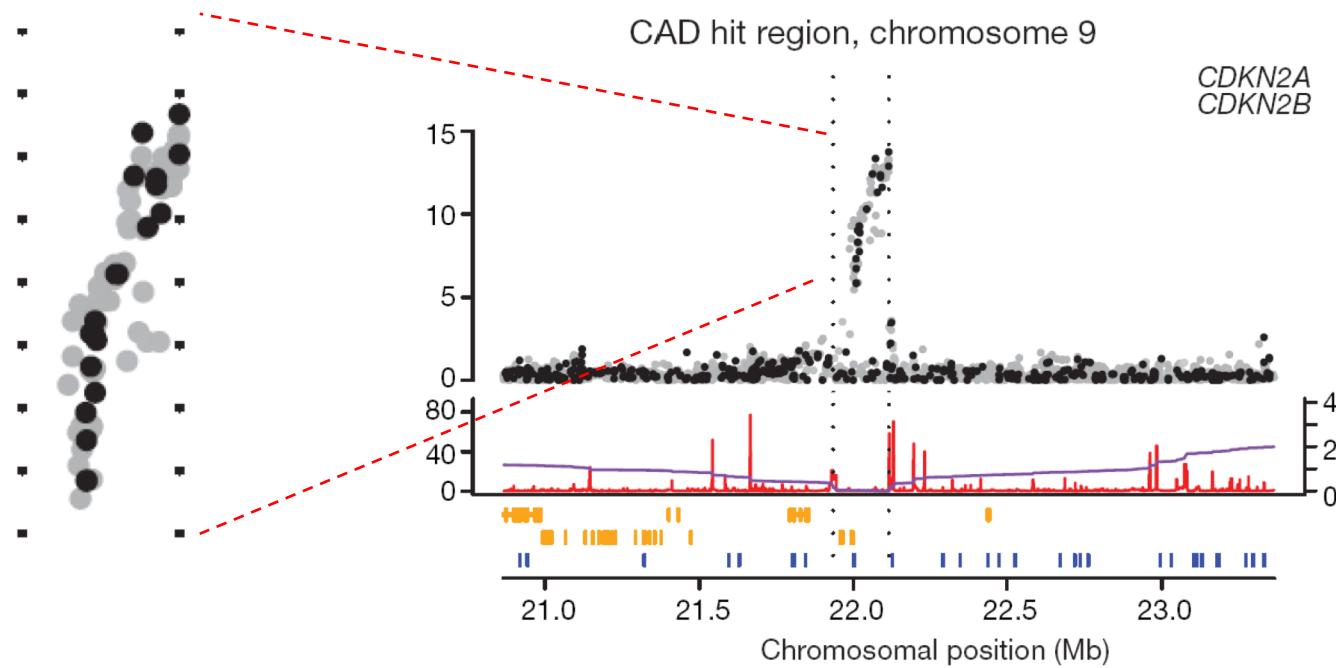
- Results from meta-analysis of various studies may suggest between study heterogeneity (e.g. especially when combining populations of different ancestry)
- How to interpret heterogeneity?
  - Differences in study design
  - Differences in population structure
  - Differences in environmental exposures
  - 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

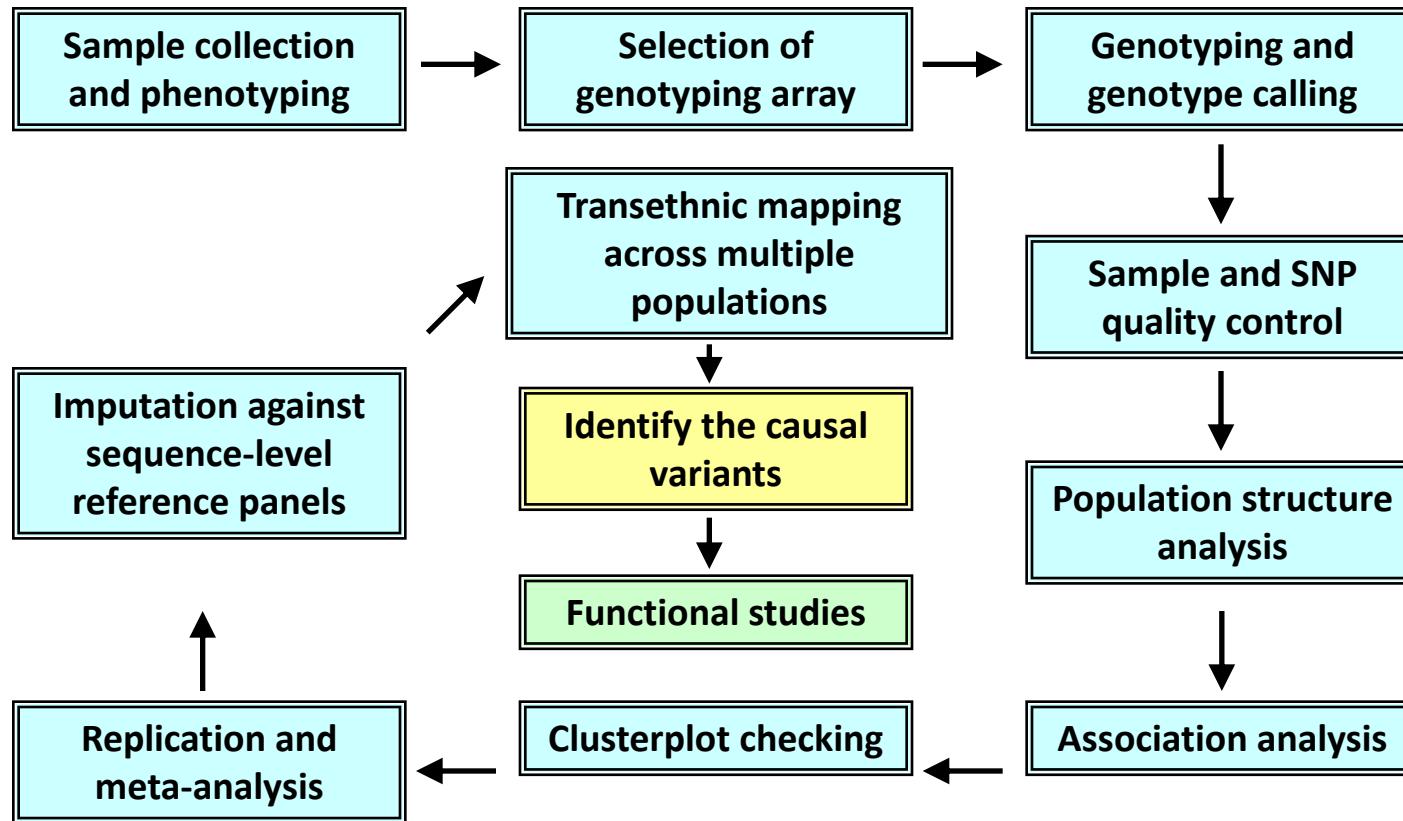
## Benefits

- 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



# 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](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](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>

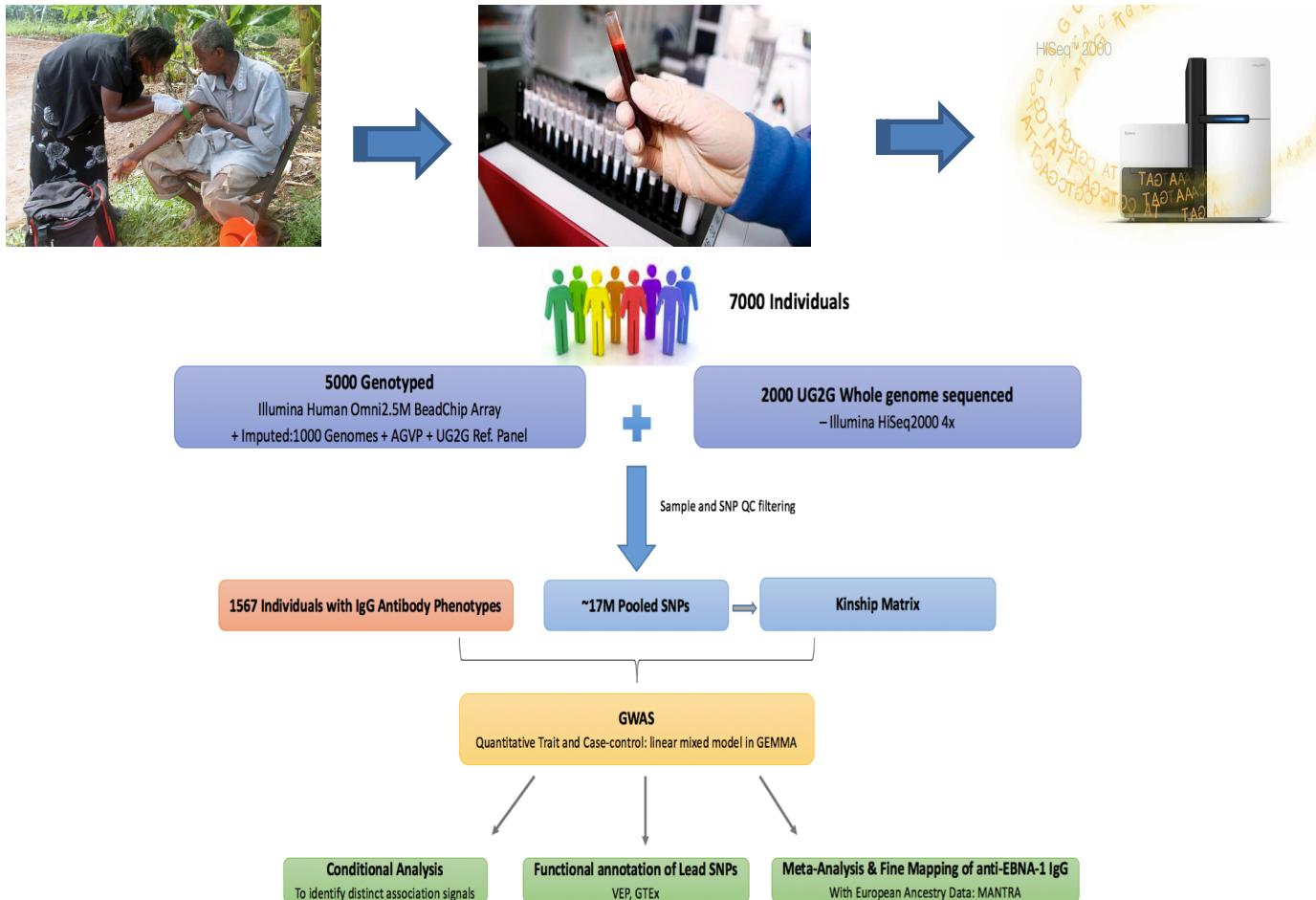
# Case in Point - GWAS of EBV in an African population

Glob Health Epidemiol Genom. 2017 Nov 27;2:e18. doi: 10.1017/gheg.2017.16. eCollection 2017.

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

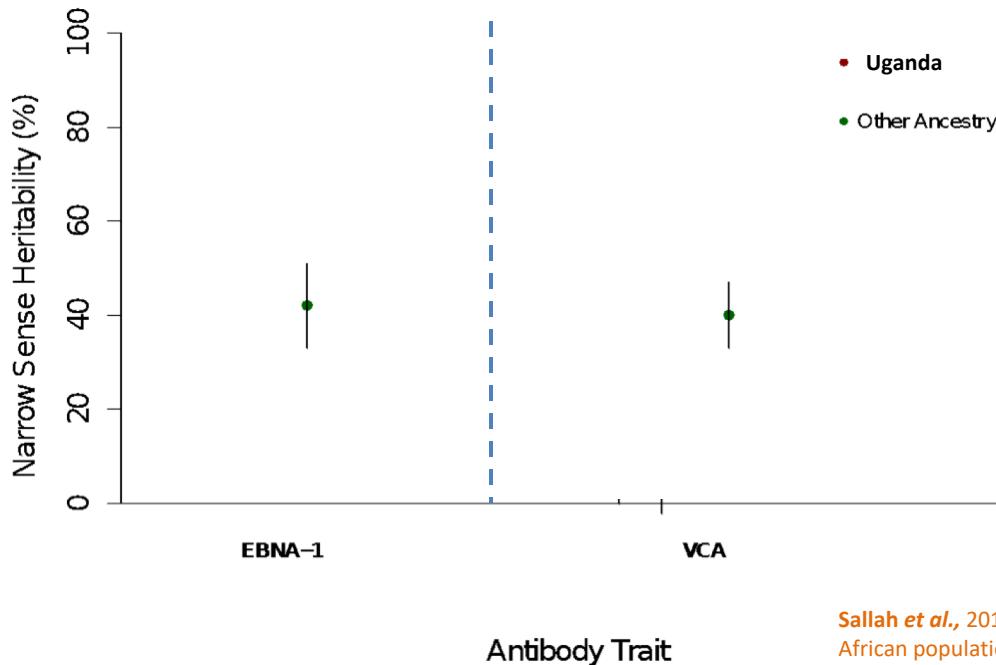
Sallah N<sup>1,2</sup>, Carstensen T<sup>1,3</sup>, Wakeham K<sup>4,5</sup>, Bagni R<sup>6</sup>, Labo N<sup>7</sup>, Pollard MO<sup>1,3</sup>, Gurdasani D<sup>1,3</sup>, Ekoru K<sup>1,3</sup>, Pomilla C<sup>1,3</sup>, Young EH<sup>1,3</sup>, Fatumo S<sup>1,3,8</sup>, Asiki G<sup>4</sup>, Kamali A<sup>4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

# 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

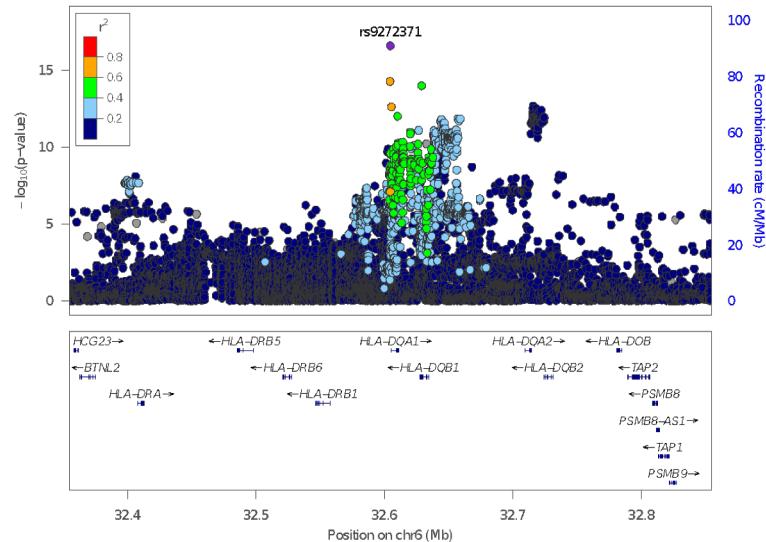
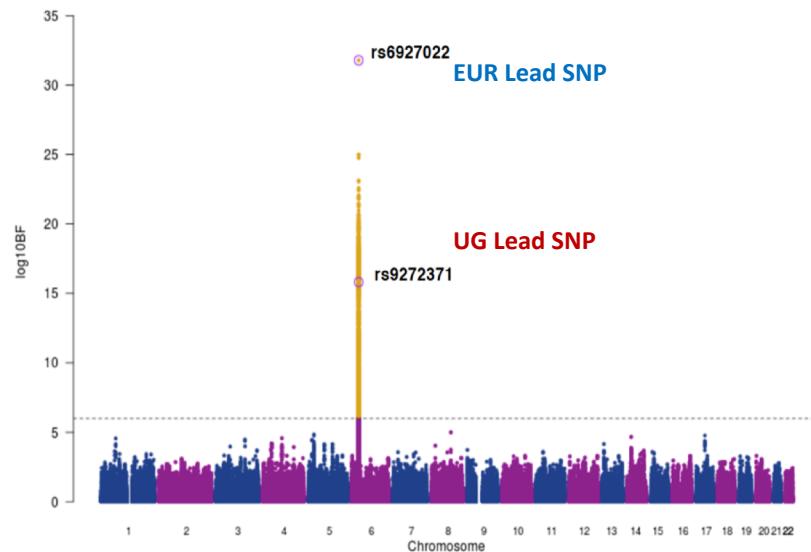


## Lower $h^2$ estimates in Ugandan population

- Differences in 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

# Other variants identified that are African-Specific

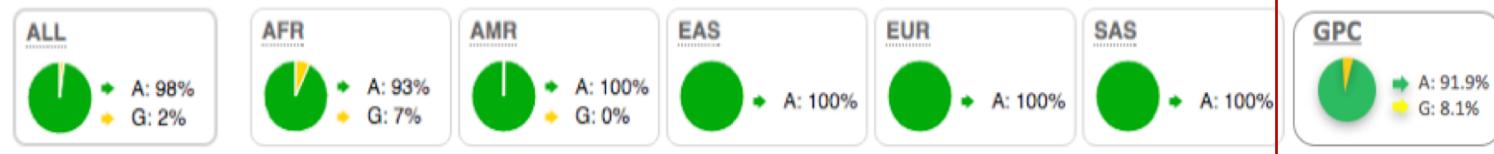
A. rs183816209



B.

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

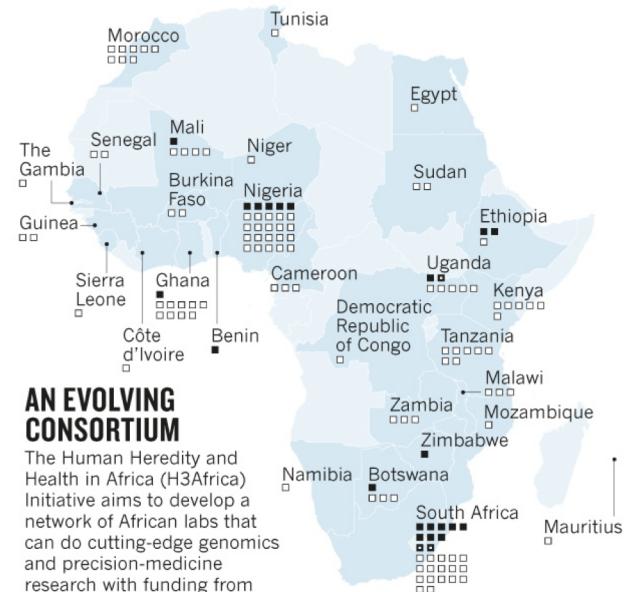
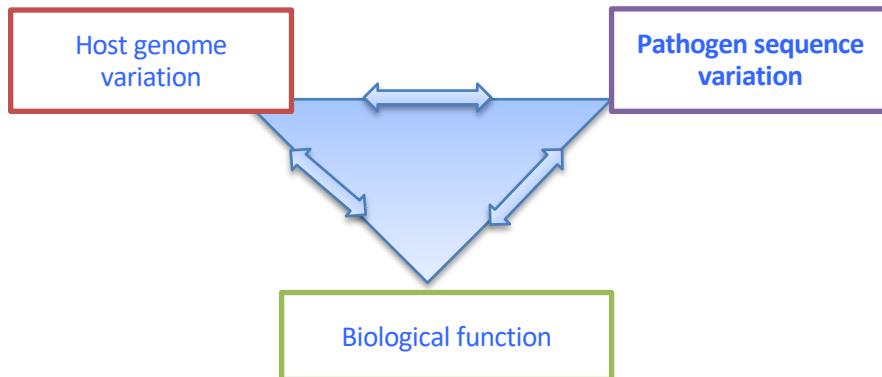
C.



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

# 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



## 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