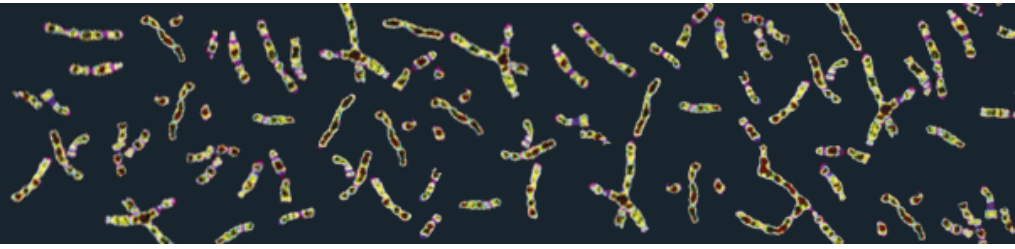# 1000 Genomes
## A Deep Catalog of Human Genetic Variation

# The 1000 Genomes Project

Gil McVean

Department of Statistics, Oxford

# How can we achieve large-scale GWAS with genomic sequence data now?

well come

# The Daily Tel

Tuesday, April 21, 2009

## Misery for Slumdog star

# 'False hope' in hunt for genetic cures

By Richard Alleyne and Kate Devlin

A LEADING scientist has claimed that the hope that genetic research could provide a cure for a host of common illnesses has proved a "false dawn".

Prof Steve Jones, a geneticist, said the belief that a few genes held the key to ridding the world of conditions such as cancer and diabetes had proved to be "plain wrong",

into genes and that there was a danger it had become "largely unfounded". "Just a couple of years ago, there was real optimism, that a new era of understanding was around the corner," he said. "That did not last long, for hubris has been replaced with concern."

Prof Jones added: "Of course there have been some successes, but it is the 'cure all' aspect of the work that has proved unfounded.

"It is the nature of the busi-
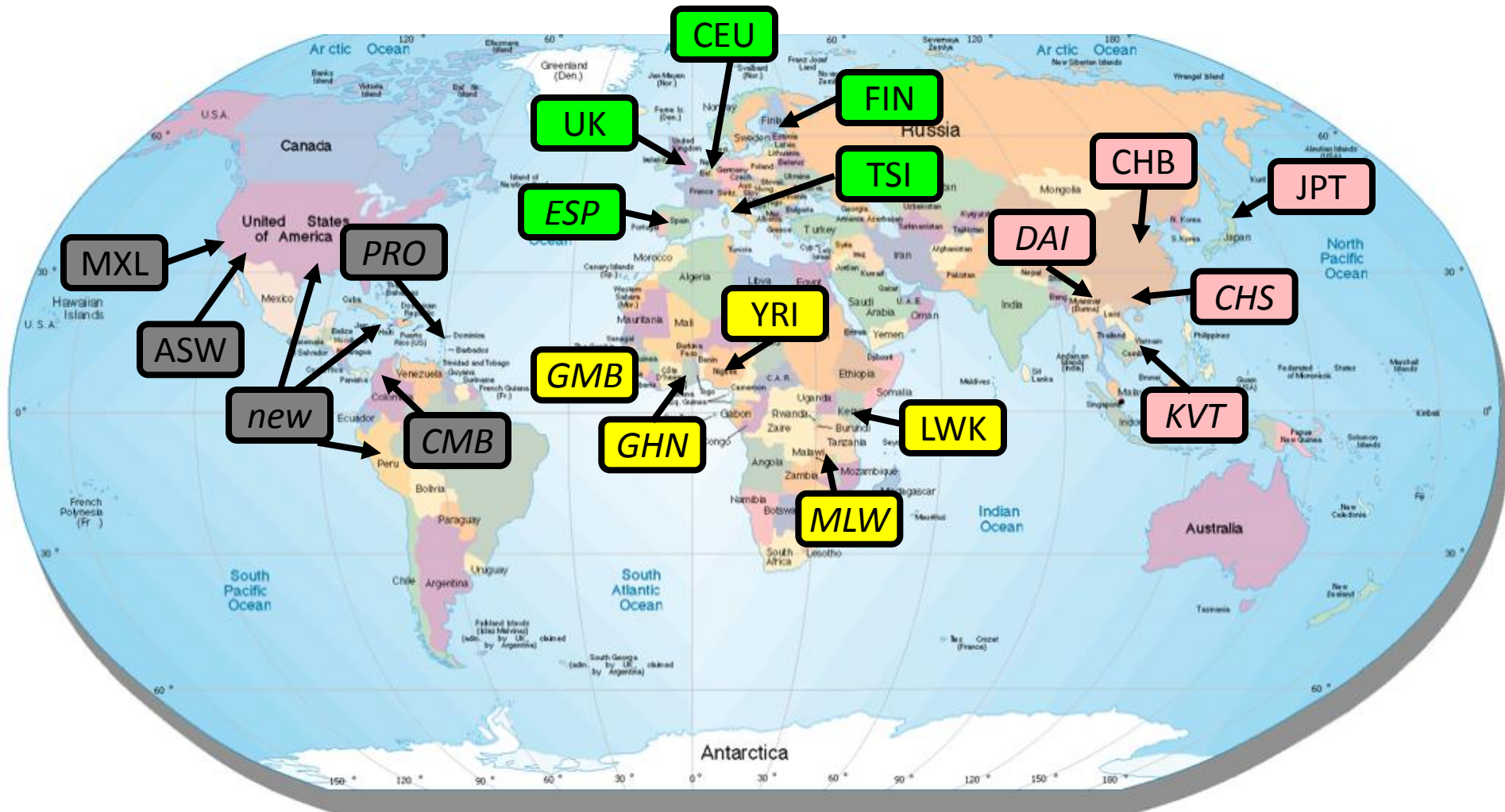
# The dark matter of genetics

# Why can we explain only a fraction of the genetic risk?

- For most complex disease/phenotypes, the proportion of the variance explained by GWAS hits is less than 5%

- What explains the missing heritability?
    - **Common, but untagged SNPs?**
    - **Structural variation?**
    - **Rare variants?**
    - **GxG interactions?**
    - **GxE interactions?**

**The 1000 Genomes Project**
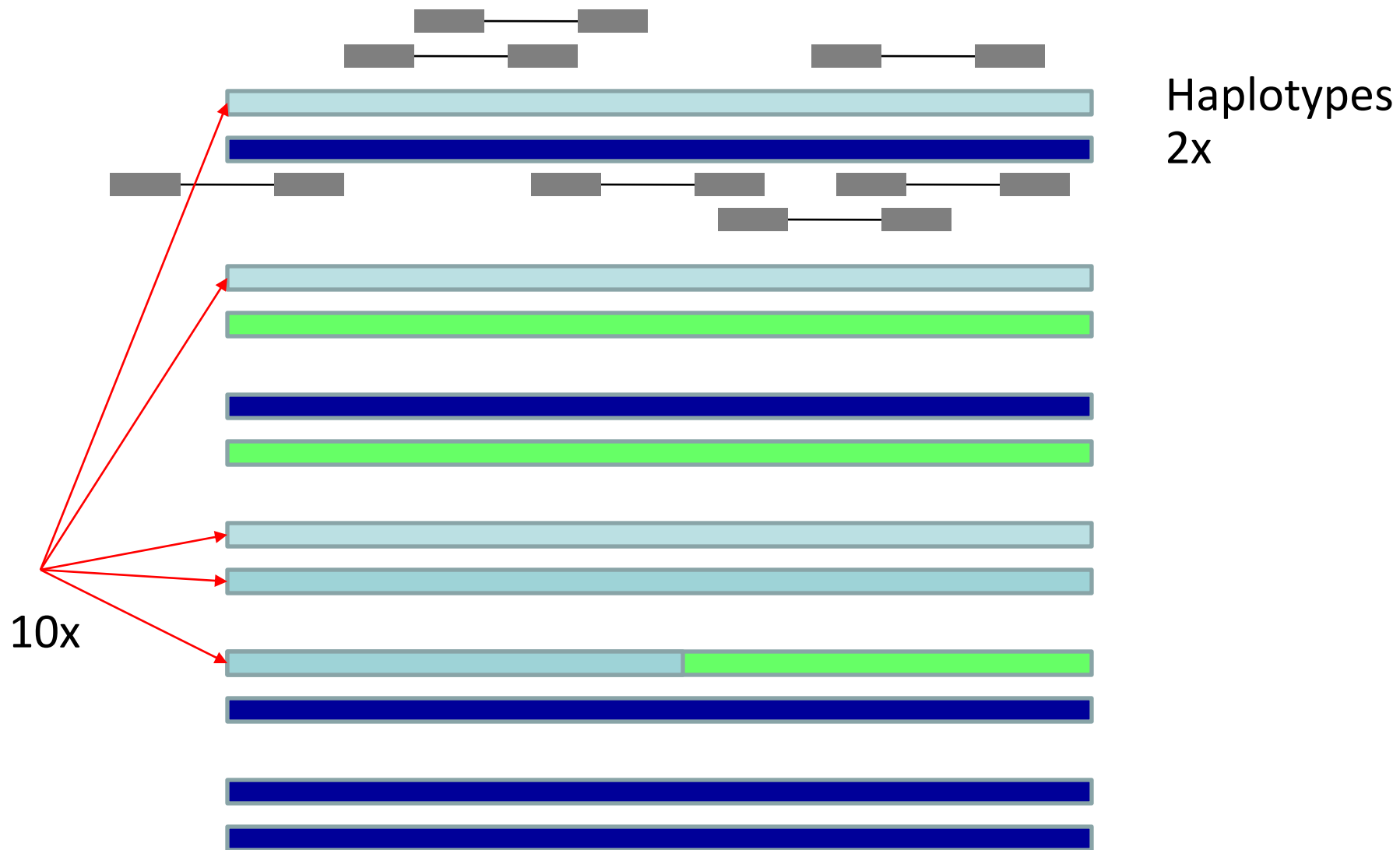
# What is the 1000 Genomes Project?

- A catalogue of **all** types of genetic variation, including **rare** variants (c. 1% frequency) obtained by sequencing at least 1000 individuals from geographic centres of major medical genetics interest

- A large international collaboration
  - UK, USA, China, Germany

- An exploration of the use of next-generation technologies for population-scale genome sequencing

- A resource for accelerating the rate of identifying disease mechanisms in the follow-up to disease-association studies

# Samples for the main project



**Major population groups comprised of subpopulations of c. 100 each**

# Population-scale sequencing



Haplotypes
2x

10x

# Pilot experiments

- Pilot 1
  - Low-coverage (4x-8x) on 60 unrelated individuals from each of CEU, YRI and CHB+JPT

- Pilot 2
  - High-coverage (20x diploid) on 2 trios (one from CEU, one from YRI)

- Pilot 3
  - Exons from 1000 genes to 20x in c. 1000 samples (largely European)

## Complete!

From The Times

May 19, 2009

# Discovery of DNA variations promises bespoke treatment for disease

Mark Henderson, Science Editor

The prospect of personalised medical care based on the genetic profiles of patients has moved closer with the discovery of millions of fresh ways in which DNA can vary from person to person.

An initiative to create a comprehensive atlas of human genetic differences has delivered spectacular early results that are already advancing the search for the genetic origins of conditions such as heart disease, diabetes and cancer.

The first phase of the international 1,000 Genomes Project has identified about 11 million new places where the human genome varies, doubling the tally known to science. Researchers have now begun to sift these variants for links to disease.

Insights from the work will accelerate development of drugs and diagnostic techniques, and pave the way for an era of bespoke medicine in which the treatment and prevention of disease are tailored to individuals' genes.

# Data processing innovation and standards



| Process | Data | Unit (one file per …) | Who? (italics if not done yet) |
|---|---|---|---|
| 1. Submit | Primary data **SRF** | *lane* | production centres |
| 2. Extract | Primary data **fastq** | *lane* | DCC |
| 3. Map sample | Sample alignment **SAM** | *lane* | Sanger *to DCC* |
| 4. Recalibrate | Mismatch table / QC data | *lane* | Sanger *to DCC* |
| | Recalibrated data **fastq** | *lane* | Sanger *to DCC* |
| 5. Map | Lane alignment **SAM** | *lane* | Data Processing |
| 6. Merge and remove dups | Library alignment **SAM** | *library* | Data Processing |
| 7. Merge | Platform alignment **SAM** | *platform/individual* | Data Processing |
| 8. Calc likelihoods | Platform likelihoods **GLF** | *platform/individual* | Data Processing |
| 9. Combine l'hoods | Individual likelihoods **GLF** | *individual* | Data Processing |
| 10. Apply priors | Posterior probabilities **GLF** | *individual* | Data Processing |
| 11. Call SNPs/indels | Candidate SNPs/indels | *experiment/population* | Data Processing |
| 12. Call genotypes | Genotypes/haplotypes | *individual* | Data Processing |
| 13. Collect read pair info | Anomalous read pairs | *library* | Structural Variation |
| 13. Collect depth info | Depth information | *library* | Structural Variation |
| 14. Call SVs | Structural variants | *experiment and individual* | Structural Variation |

# Sequence AlignMent format

```
(a) coor    12345678901234    567890123456789012345678901234 5
    ref     AGCATGTTAGATAA**GATAGCTGTGCTAGTAGGCAGTCAGCGCCAT

    r001+           TTAGATAAAGGATA*CTG
    r002+         aaaAGATAA*GGATA
    r003+      gcctaAGCTAA
    r004+                   ATAGCT..............TCAGC
    r003-                            ttagctTAGGC
    r001-                                      CAGCGCCAT


(b) @SQ SN:ref LN:45
    r001 163 ref  7 30 8M2I4M1D3M = 37   39 TTAGATAAAGGATACTA *
    r002   0 ref  9 30 3S6M1P1I4M *  0    0 AAAAGATAAGGATA     *
    r003   0 ref  9 30 5H6M        *  0    0 AGCTAA    *    NM:i:1
    r004   0 ref 16 30 6M14N5M     *  0    0 ATAGCTTCAGC       *
    r003  16 ref 29 30 6H5M        *  0    0 TAGGC     *    NM:i:0
    r001  83 ref 37 30 9M          =  7  -39 CAGCGCCAT         *


(c) ref  7 T 1 .   |ref 12 T 3 ...  |ref 17 T 3 ...
    ref  8 T 1 .   |ref 13 A 3 ...  |ref 18 A 3 .-1G..
    ref  9 A 3 ... |ref 14 A 2 .+2AG.+1G |ref 19 G 2 *.
    ref 10 G 3 ... |ref 15 G 2 ..   |ref 20 C 2 ..
    ref 11 A 3 ..C |ref 16 A 3 ...  |...
```

Bioinformatics (2009) http://samtools.sourceforge.net

# Variant Call Format



**www.1000genomes.org/wiki/doku.php?id=1000_genomes:analysis:vcfv3.2**

# www.1000genomes.org

# ftp.1000genomes.ebi.ac.uk

# Read-scale view

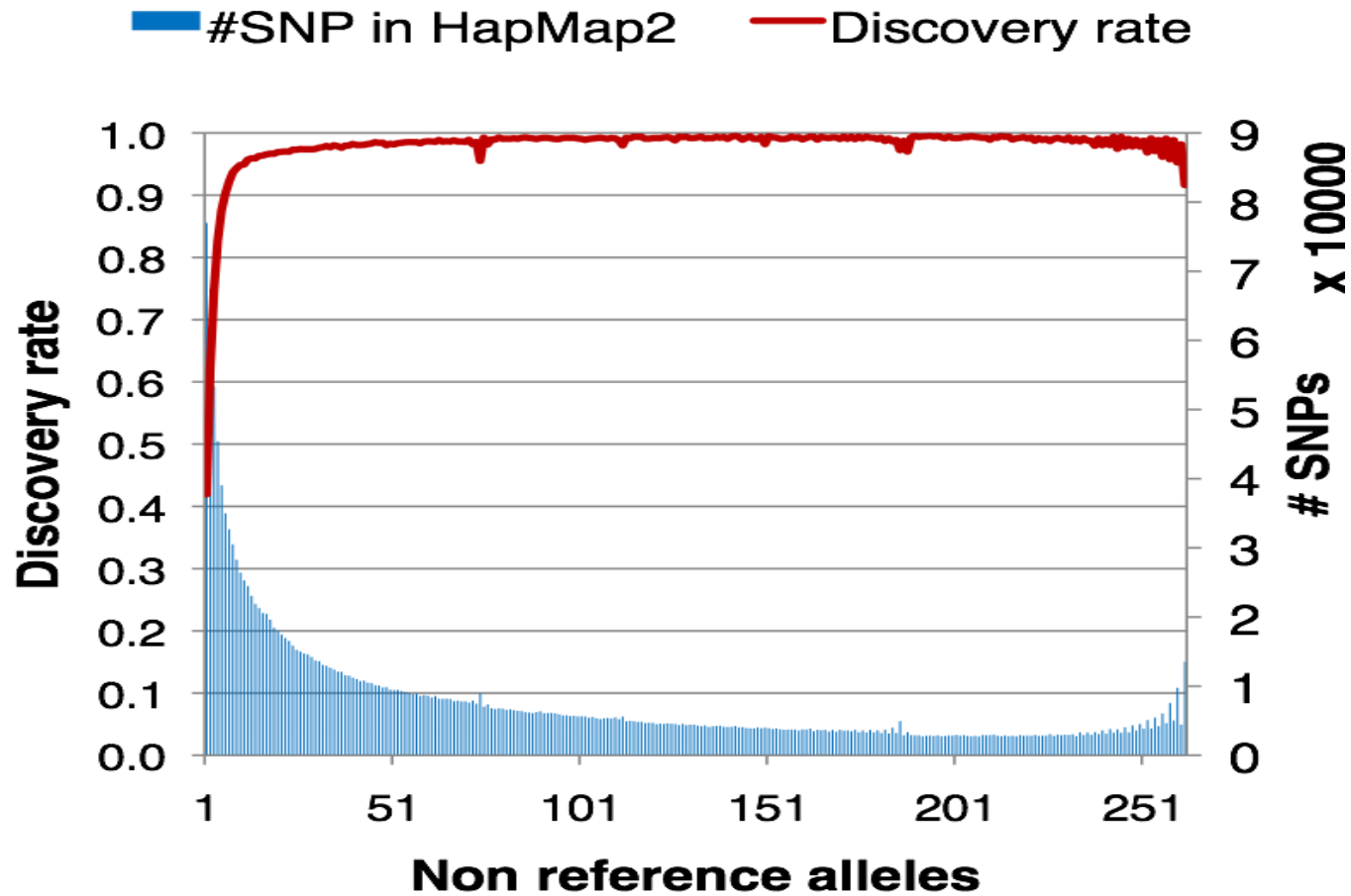**Eric Banks (Broad)**

# Genome wide SNP discovery

- Total 17.2 M SNPs called

- Previously ~12M SNPs "known" (dbSNP 129)
  - 7.9M confirmed
  - 9.2M novel



**Total SNPs**

CEU    YRI

2.80    1.09    5.65

4.84

0.78    0.48

1.54

CHB+JPT

**Novel SNPs**

CEU    YRI

2.20    0.38    4.38

0.50

0.29    0.26

1.35

CHB+JPT

Le Quang

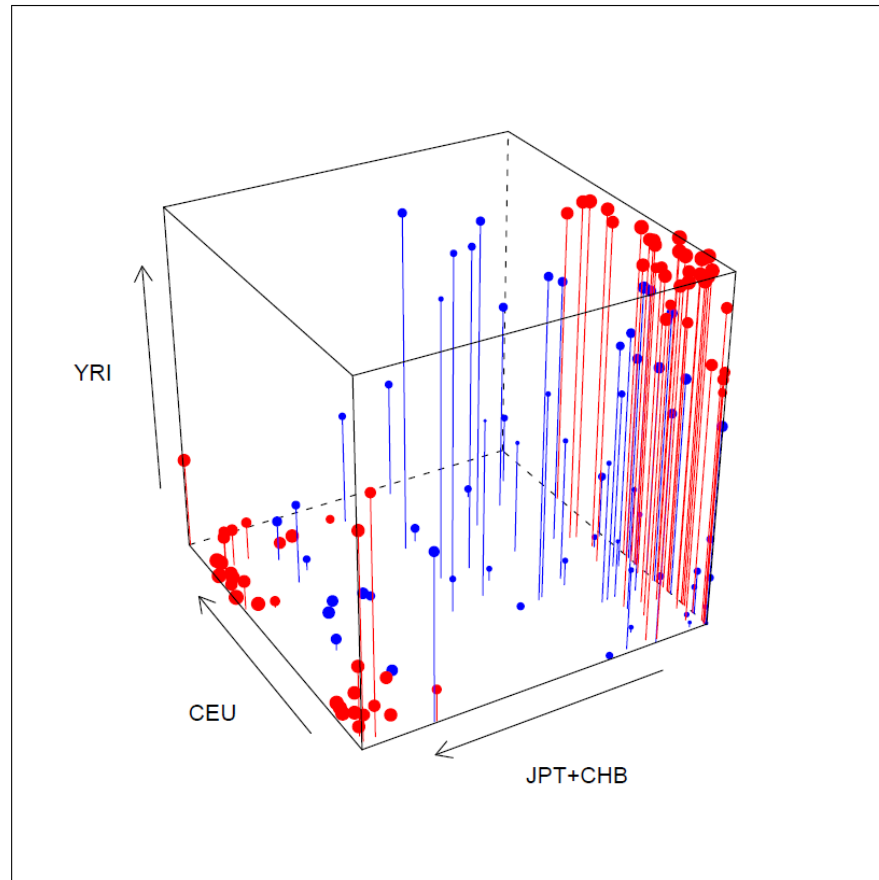# Completeness



**Durbin, Le Quang**

# Genotype accuracy on HapMap2



- This is about where simulations suggest we should be with 2-4x on 60 samples

- Much higher than independent calls

# Some surprises – high Fst SNPs



Ryan Hernandez, Adam Auton

# Using the 1000G data **now**

# Imputation

… 1110101010101011 …
… 0011111000111 …
… 1111000011101 …
… 0010101110010 1 …

Reference panel
(1000G)

… 1.2..1.0.0..22…

Genotypes in
additional
samples from
standard product

IMPUTE

… 112201102001 22 …

Imputed
genotypes

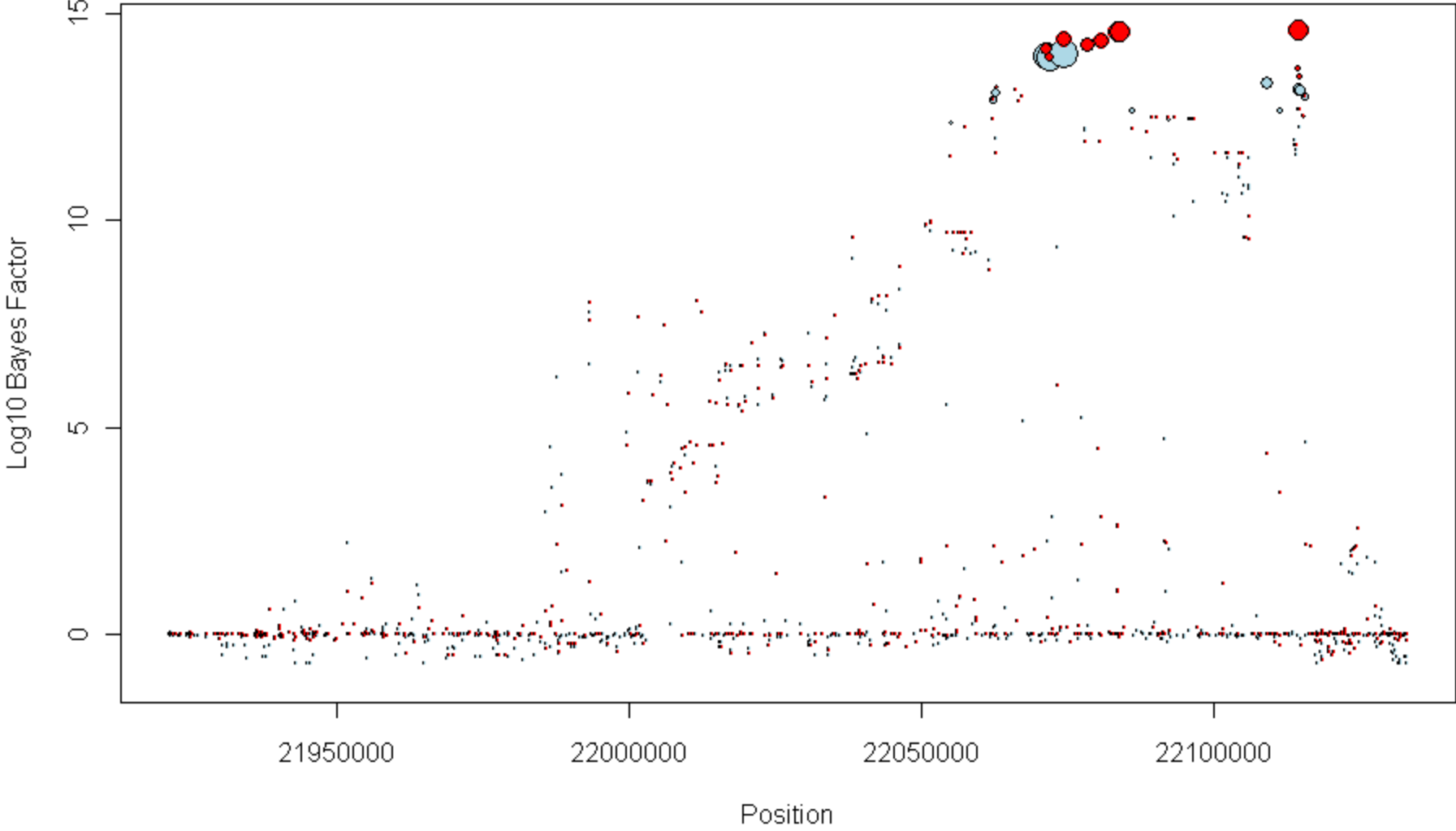# What can we use imputation for in GWAS/fine-mapping?

- To
  - Help define genomic regions likely to contain causal variants
  - Define a small set of SNPs to take through to additional genotyping?
  - Fine-map?

- Accuracy depends on completeness of imputation source and accuracy of imputed genotypes
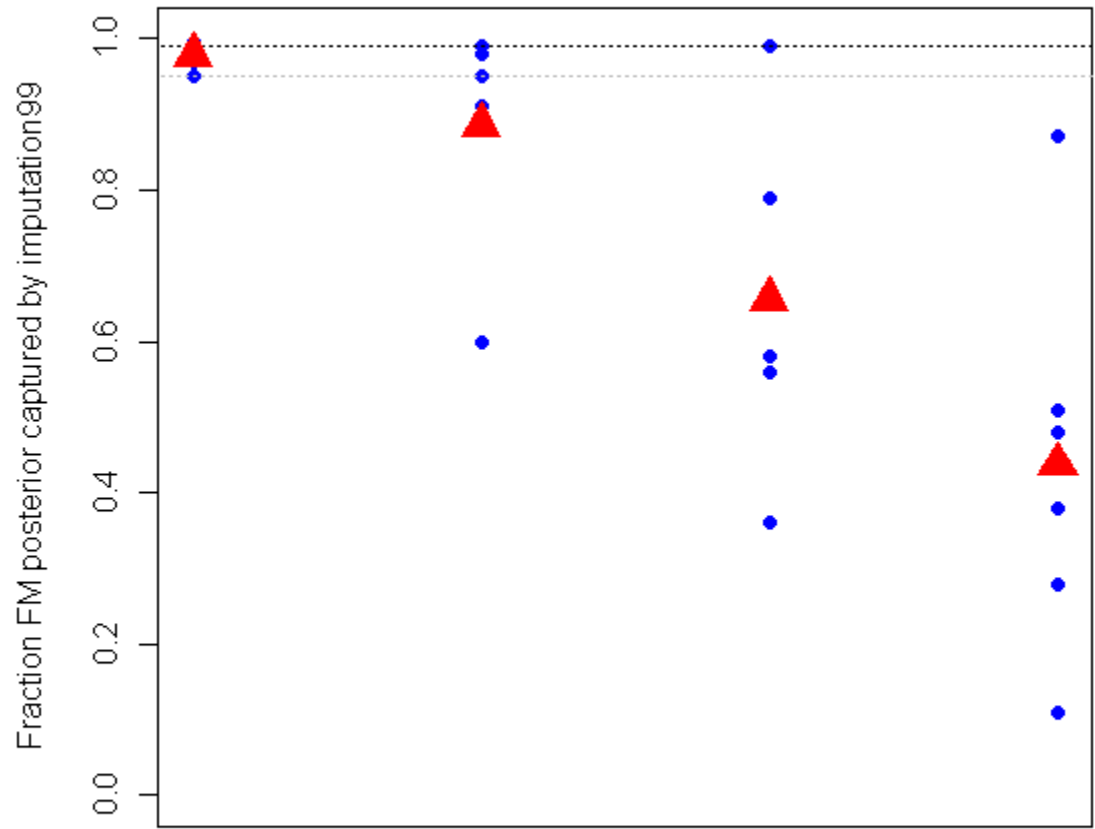
# Imputation resources

| Data source | Type | Number haplotypes | Number SNPs | Includes |
|---|---|---|---|---|
| HapMap2 | GT | 120 | 671 | Mainly common SNPs, phased from trios |
| Reseq | RS | 64 | 1543 | SNPs found from resequencing. Genotypes called independently |
| FM panel | GT | 64 | 1699 | All SNPs from RS pilot and dbSNP for which design possible |
| 1000 Genomes pilot (CEU) | RS | 114 | 2561 | SNPs found from RS, integrated into haplotype structure of HM3 SNPs |

# How can we measure imputation?

CAD: CDKN2A/2B

# How good is imputation?

- Imputation is never going to be as convincing as genotyping

- BUT it is sufficiently accurate, at least for common variants, to
  - Define sets of SNPs of interest
  - Exclude SNPs
  - Indicate whether the signal is likely to be localised through additional genotyping (a function of power and haplotype structure)

- The completeness of the 1000G project data is extremely valuable
  - The WTCCC decided to stop any further sequencing of controls or cases

# Open questions

- How reliably can you detect structural variation from low-coverage data?

- How do you combine information across populations?

- Can we generate reliable de novo assemblies from the data?

# Acknowledgements

- Oxford
  - Adam Auton, Zam Iqbal, Gerton Lunter, Jules Maller, Simon Myers, Jonathan Marchini, Peter Donnelly

- The Wellcome Trust Case Control Consortium

- The 1000 Genomes Project