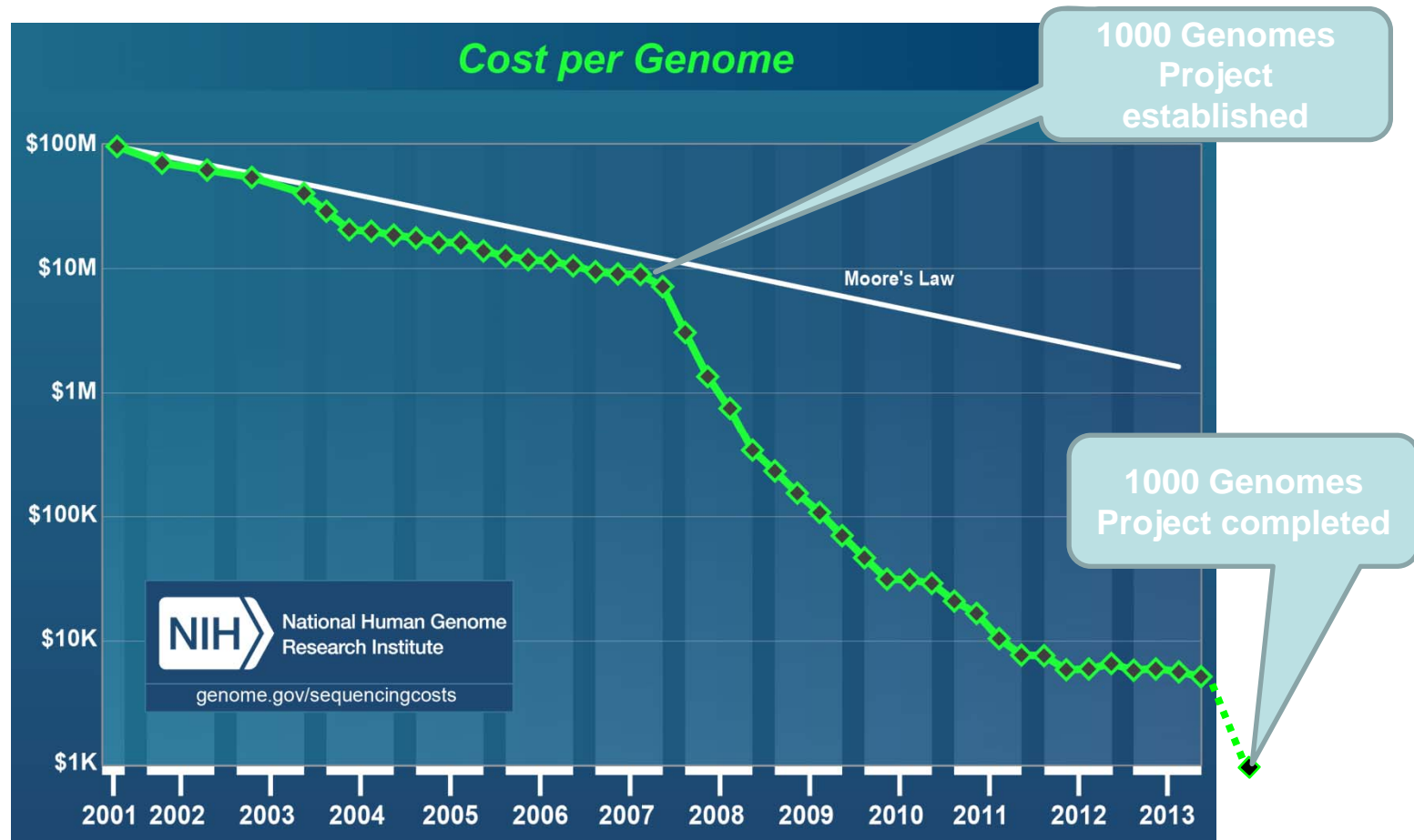# The 1000 Genomes Project

Chris Tyler-Smith

The Wellcome Trust Sanger Institute

Hinxton, CB10 1SA

# Topics

- Why was the 1000 Genomes Project established?

- What has the project achieved?

- What is its importance and legacy?

# Human genome sequencing costs 2001-2014

# The view in 2007

- Major developments in sequencing technology on the horizon

- Becoming possible to sequence multiple whole human genomes

- No single group thought they could do this alone at scale

- Need to establish an international collaborative project
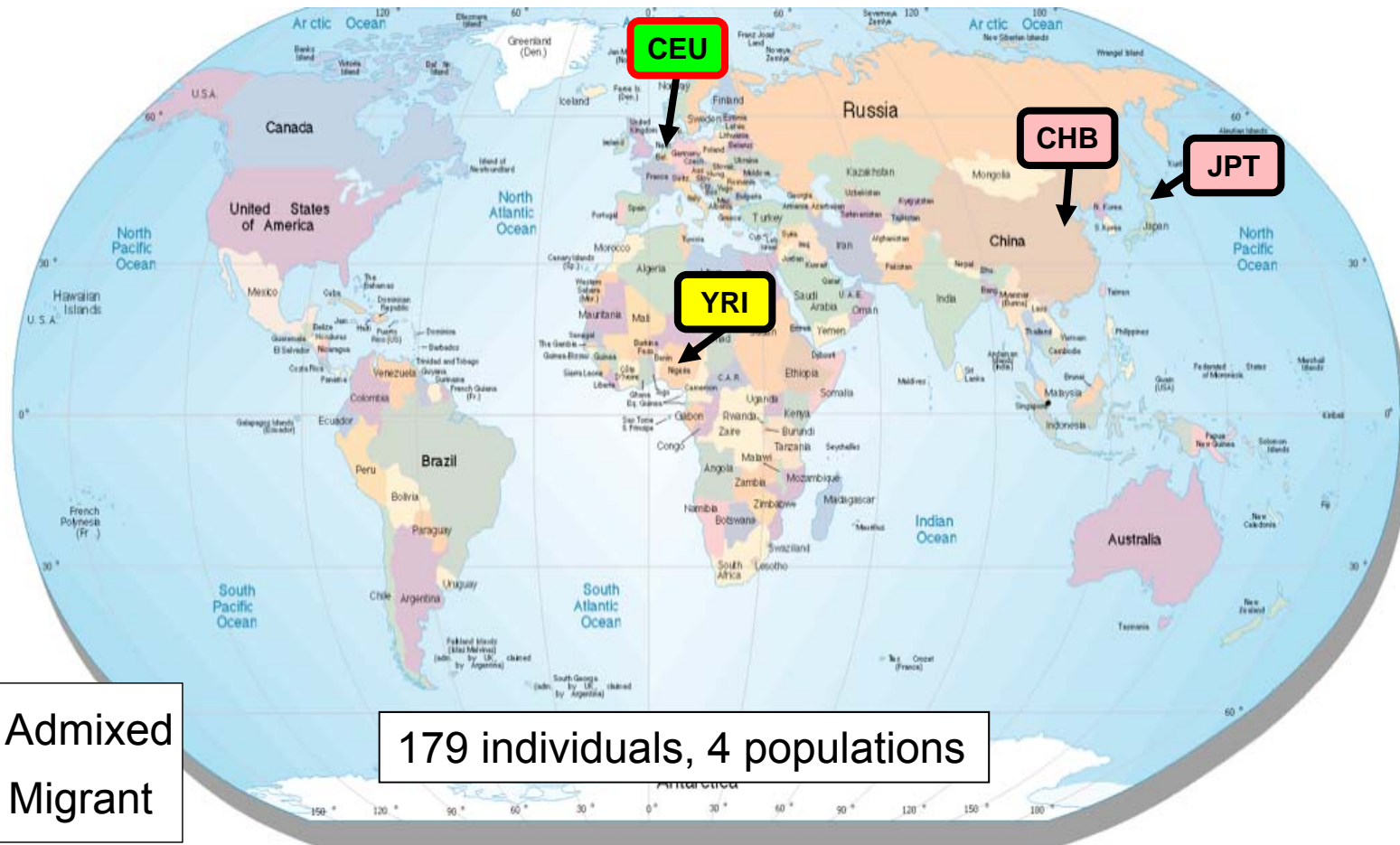
# Aims of the 1000 Genomes Project



- Primary goal: to develop a public resource of genetic variation to support the next generation of medical association studies

- Find all accessible variants ≥1% across the genome and 0.1-0.5% in gene regions

- Estimate allele frequencies, identify haplotype backgrounds, etc.

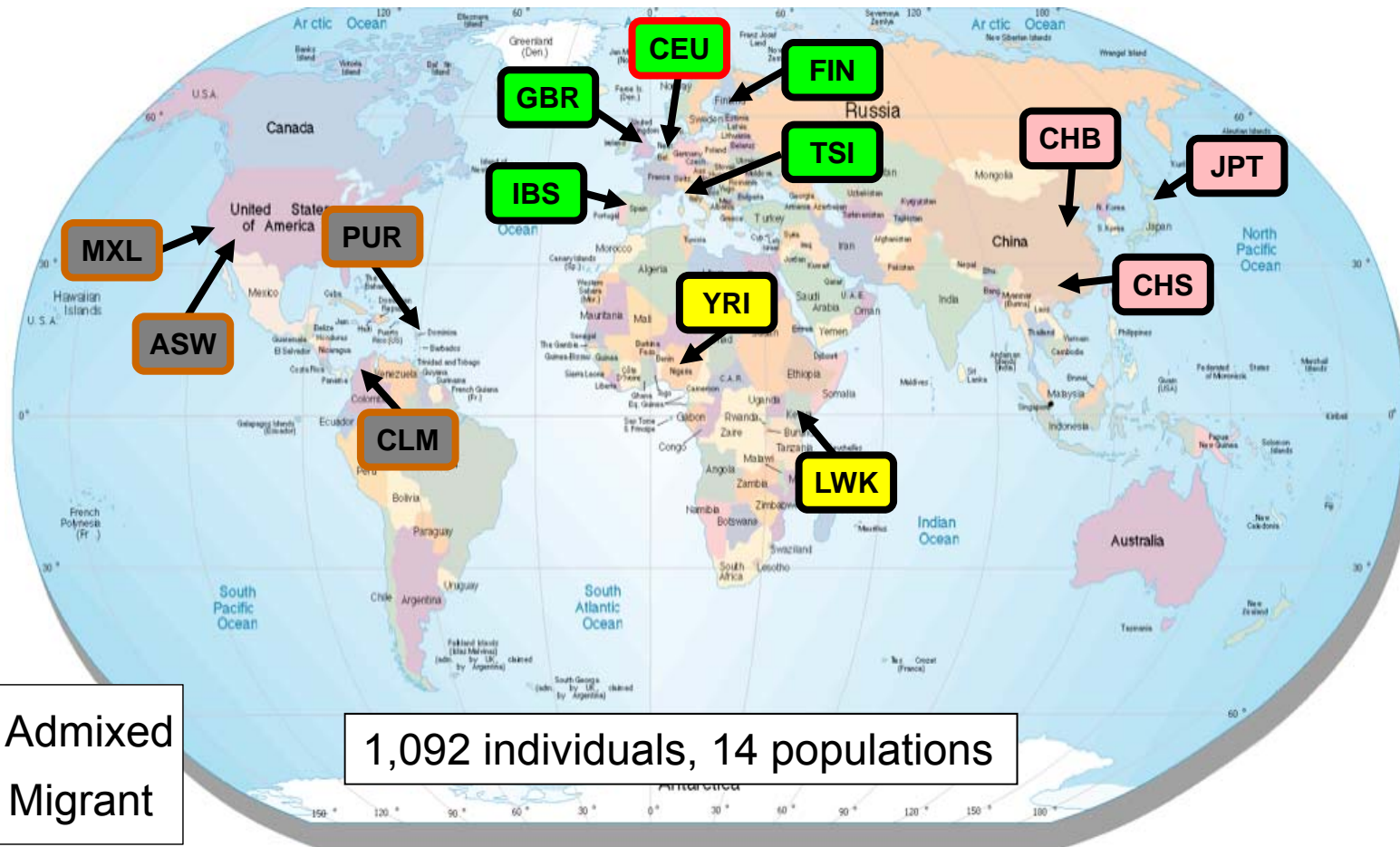wellcome trust
**sanger**
institute

# Sample choice

- Consent for full free web release of sequence data needed

- When the project began, only HapMap samples met these criteria

- Consent process developed, additional samples recruited

- All individuals are anonymous adults, able to consent, with no phenotype information

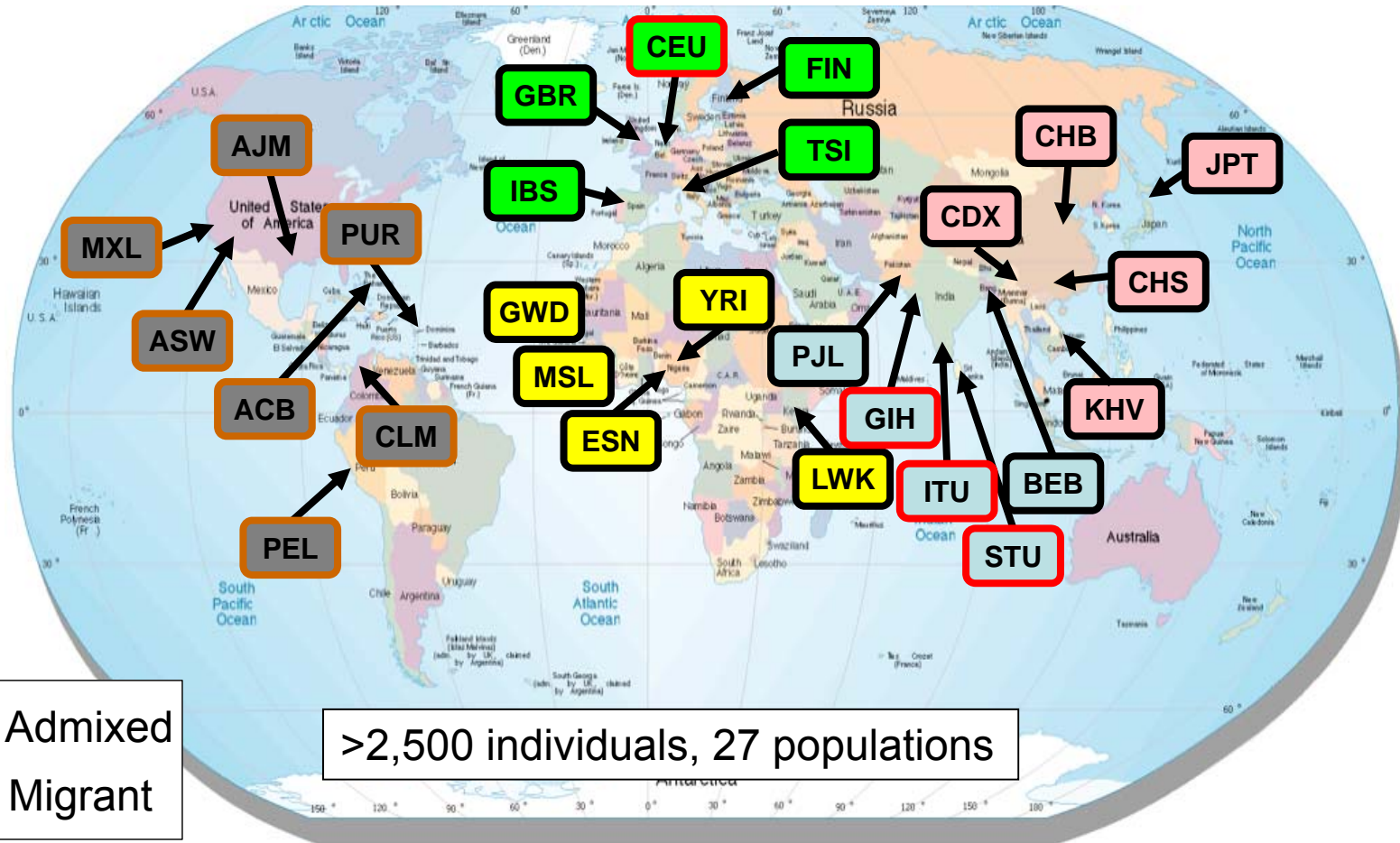# Pilot project samples (2010)



179 individuals, 4 populations

Admixed
Migrant

# Phase 1 samples (2012)



Admixed
Migrant

1,092 individuals, 14 populations

wellcome trust
**sanger**
institute

# Full project samples (2014)



Admixed
Migrant

>2,500 individuals, 27 populations

# Project design considerations

- In early 2008, still several million $ to sequence a human genome at high coverage (30x)

- But shared genetic variants can be effectively discovered by sequencing at low coverage (2x) and combining information from multiple individuals

wellcome trust
**sanger**
institute

# Stages of the 1000 Genomes Project

- **Pilot, published 2010:**
  - 179 genomes at 2x-3x, 2 trios + ~1000 genes at high coverage

- **Phase I of main project, published 2012:**
  - 1,092 genomes at 3x-4x + high-coverage whole exomes

- **Phase 3**
  - >2,500 genomes at 4x-6x + exomes, available spring 2013
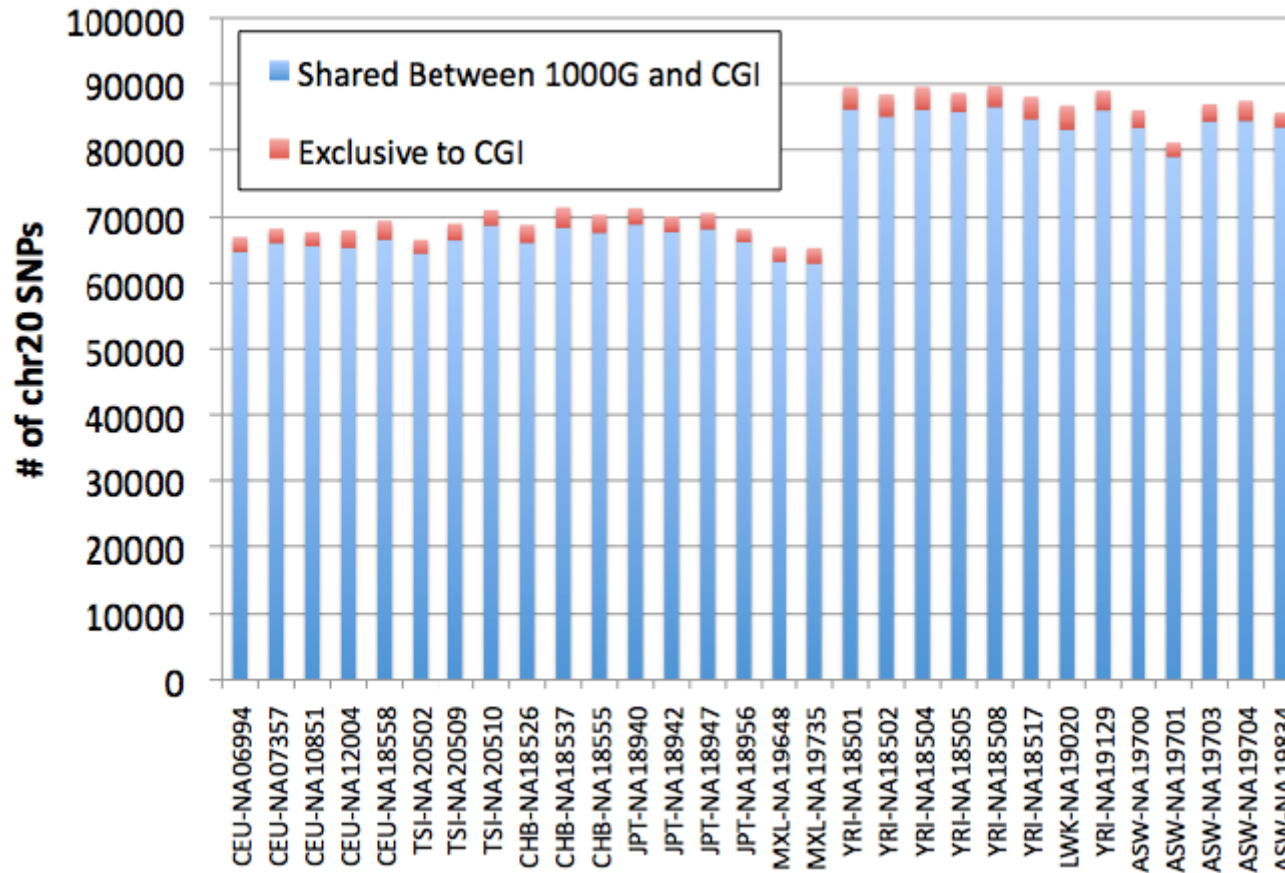  - Analysis during 2013-14

# Many variants discovered

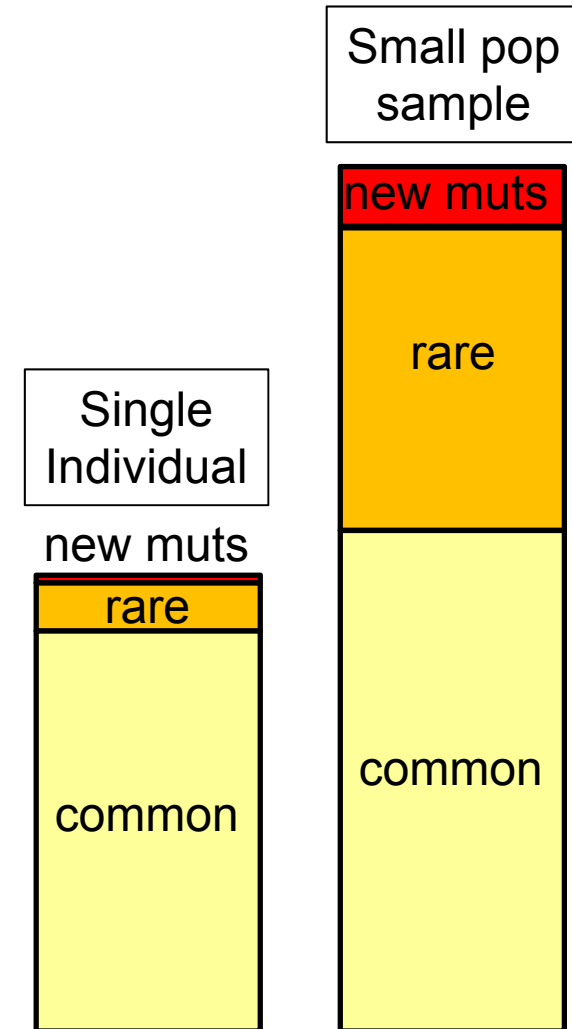| Variant type | Pilot | Phase 1 |
| --- | --- | --- |
| **Total SNPs** | 15.2 M | 37.9 M |
| Known SNPs | 6.8 M | 8.2 M |
| Novel SNPs | 8.4 M | 29.7 M |
| **Short indels** | 1.5 M | 3.8 M |
| **Large deletions** | 14 K | 14 K |

**New generation of SNP chips**
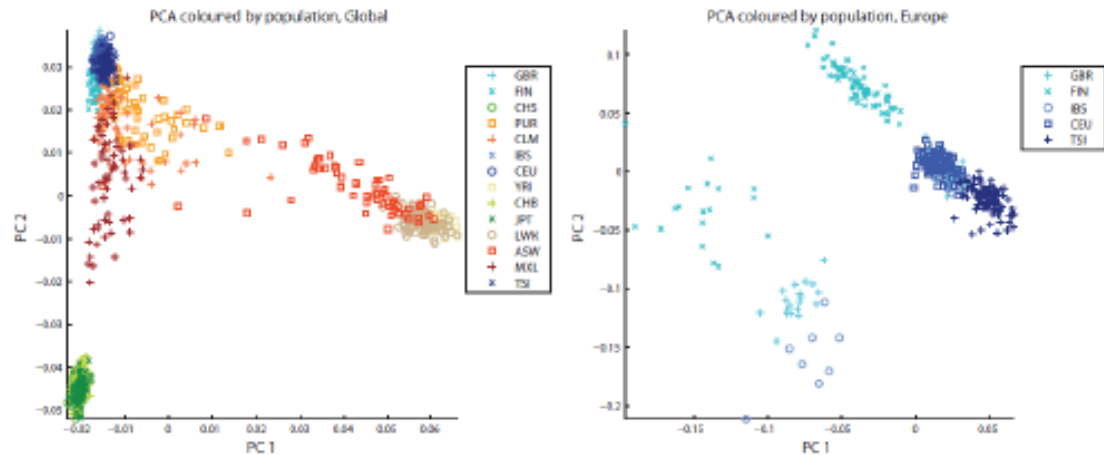
# Discovered most of the variants present in a genome

# Features of human genetic variation

- 3-4 million variants per individual

- Most have no functional consequences

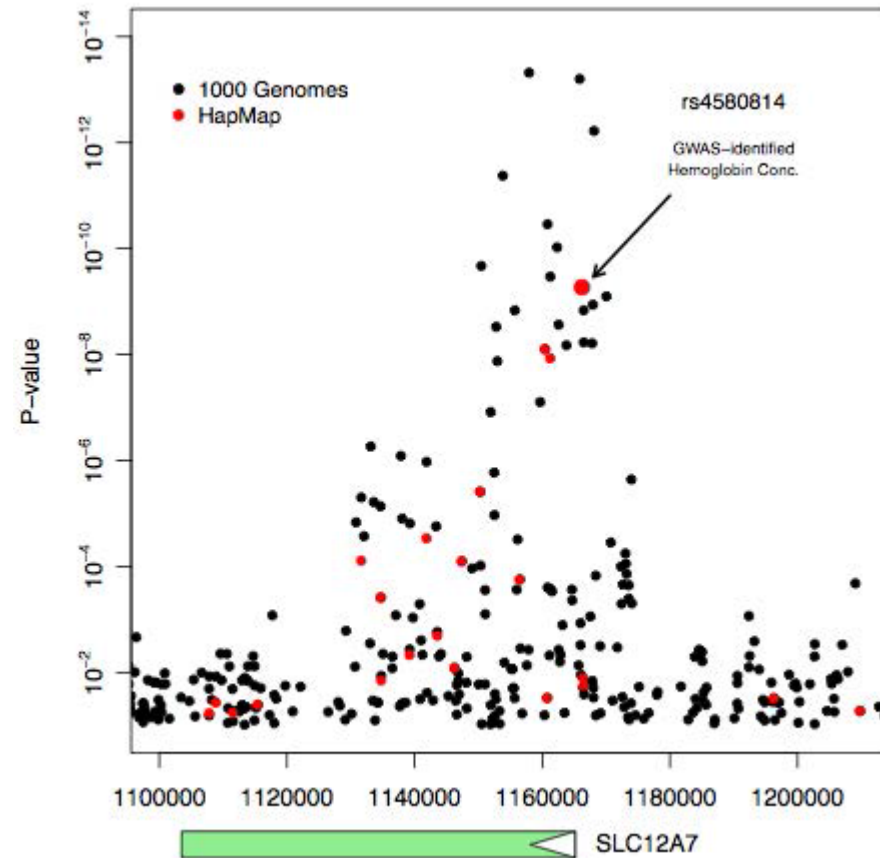- Some are common in the population, some rare

Small pop sample

new muts

rare

Single Individual

new muts

rare

common

common

wellcome trust
**sanger**
institute

# Insights from common variants (1)

- Geographical origin is predicted by genotype

# Insights from common variants (2)

- GWASs are usually carried out using SNP chips
- The best hit on the chip is often not the causal variant
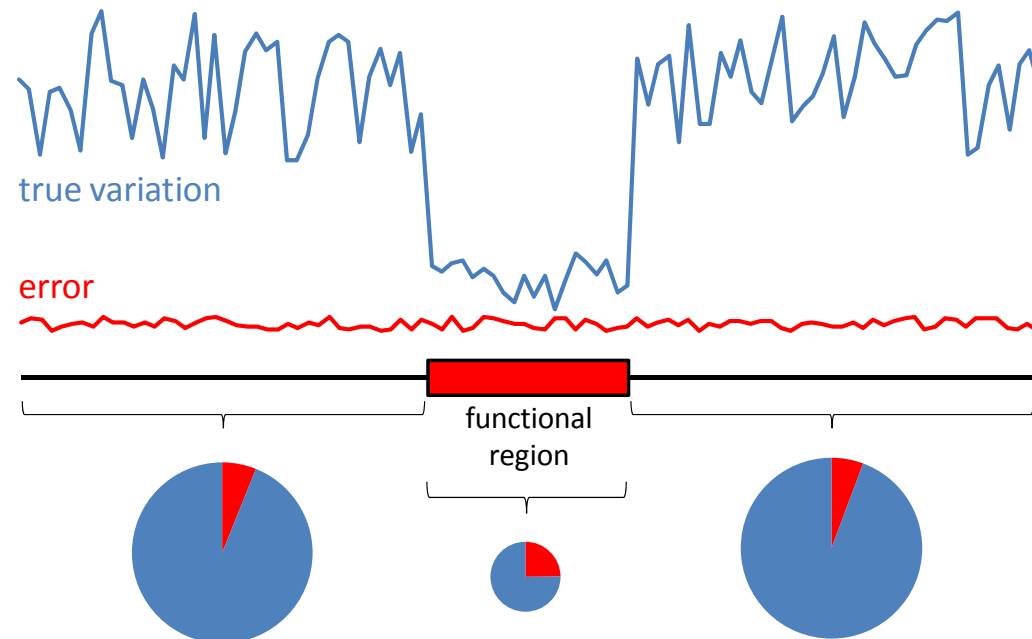- Better candidates for the causal variant can be obtained by imputation using sequence data

# Strong functional variants are mostly rare in the population

# Insights into functional variants (1)

- Enriched in errors of many kinds, extra QC and validation required



true variation

error

functional region

wellcome trust
**sanger**
institute

# Insights into functional variants (2)

- Each individual carries ~100 genes in an inactive form, ~20 with both copies inactive

2,951 raw > 1,269 filtered

| Category | Filtered number/individual (CEU) | |
|---|---|---|
| | All | Homozygous |
| **nonsense SNP** | 26.2 | 5.2 |
| **splice SNP** | 11.2 | 1.9 |
| **frameshift indel** | 38.2 | 9.2 |
| **large LoF deletion** | 28.3 | 6.2 |
| **total** | 103.9 | 22.3 |

MacArthur *et al*. (2012) *Science* **335**, 823-828

wellcome trust
**sanger**
institute

# Insights into functional variants (3)

- Each individual carries ~2 (0-7) known disease-causing variants, and these are expected to impact health in ~10% of carriers
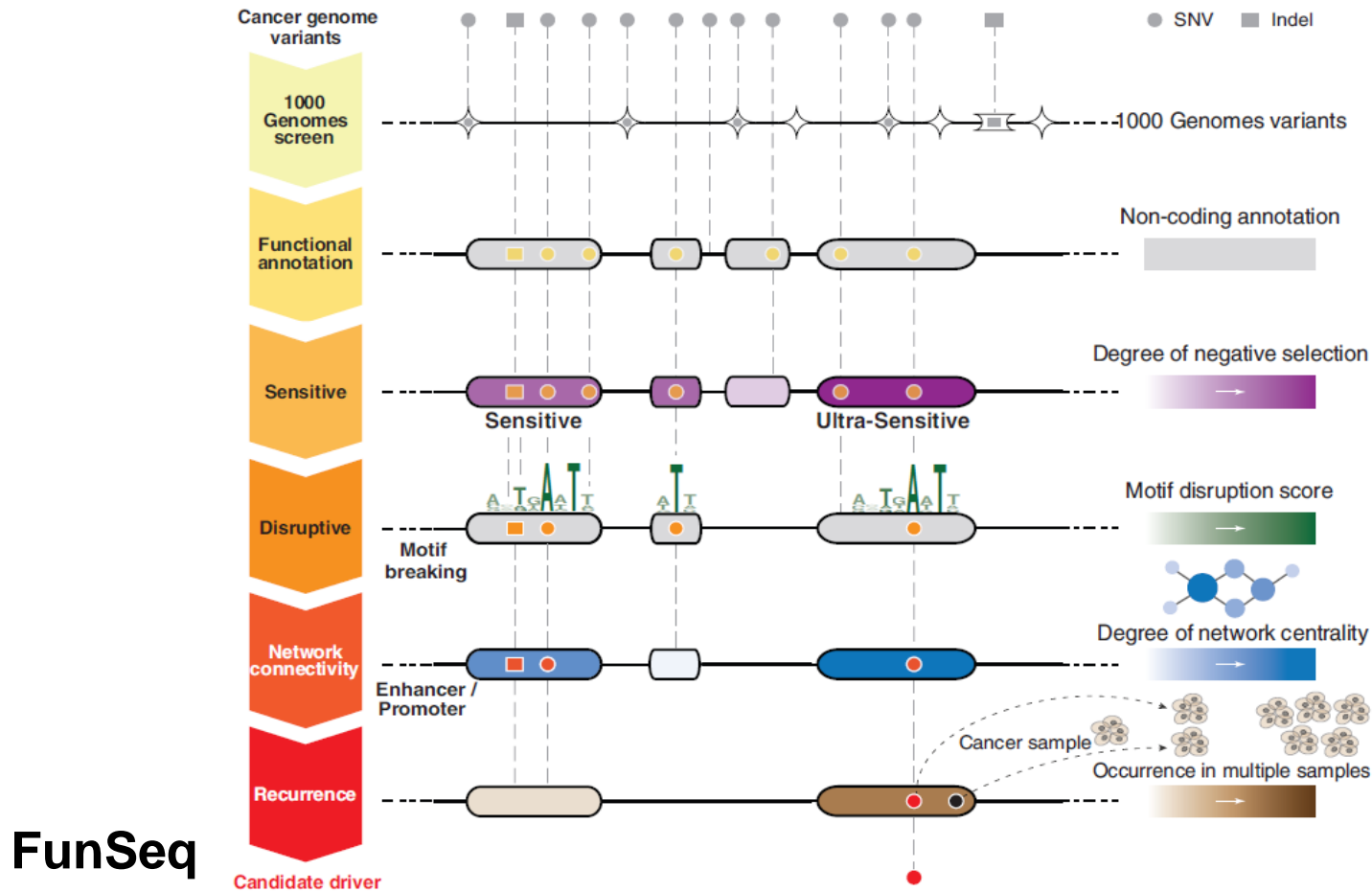


578 raw > 45 filtered

---

Xue *et al.* (2012) *Am. J. Hum. Genet.* **91**, 1022-1032

**SME Bioinformatics, Hinxton**

**6th March 2014**

# Insights into functional variants (4)

- Population genetics can help to discover new *non-coding* functional variants

Khurana *et al.* (2013) *Science* **342**, 84

# Insights into functional variants (4)



Khurana *et al.* (2013) *Science* **342**, 84

# Insights into functional variants (4)



Khurana *et al*. (2013) *Science* **342**, 84

**New Non-coding cancer driver mutations**

**FunSeq**

# Final thoughts

- One of the very few large-scale sources of open access genomic data

- Samples (cell lines) available

- No phenotypes

- A standard and lasting resource for the human genomics field

- A legacy model for additional populations

wellcome trust
**sanger**
institute

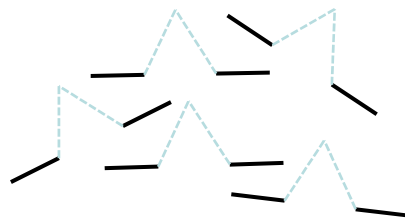# Acknowledgements

# http://www.1000genomes.org/

# Sequencing technology

**What we want:**

TCACAATGTA

**What we get:**

**Base-calling errors**

**Mapping errors**

Map to reference sequence

**Duplicate reads**

Short reads derived from single molecules
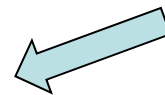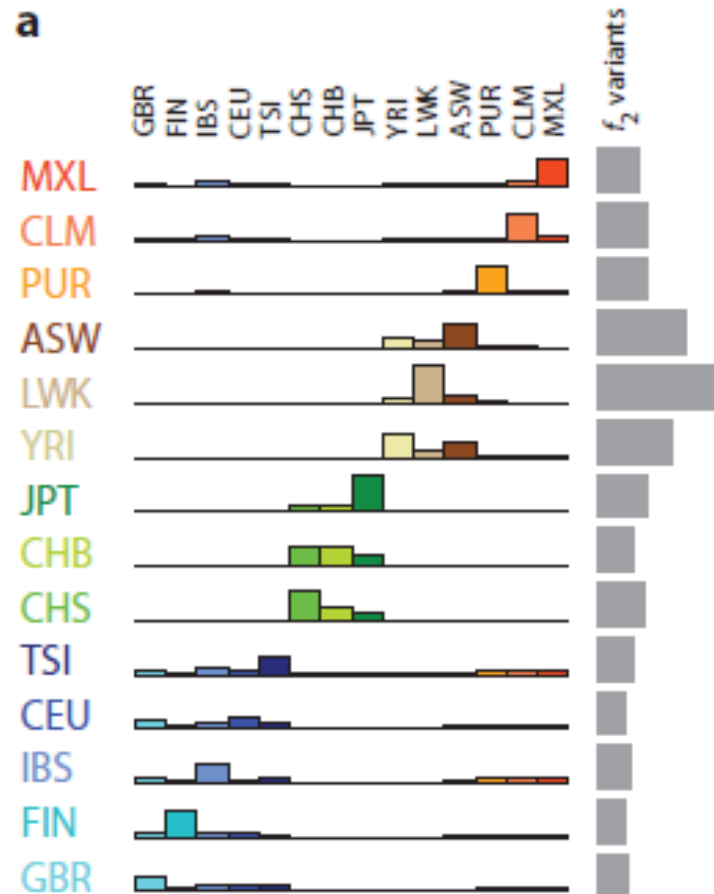
Call variants

**Variant-calling errors**

Large but error-prone datasets
Need to filter and validate

# Rare variants tend to be population-specific



Gil McVean, Adam Auton

Phase 1

**SME Bioinformatics, Hinxton**

**6th March 2014**