

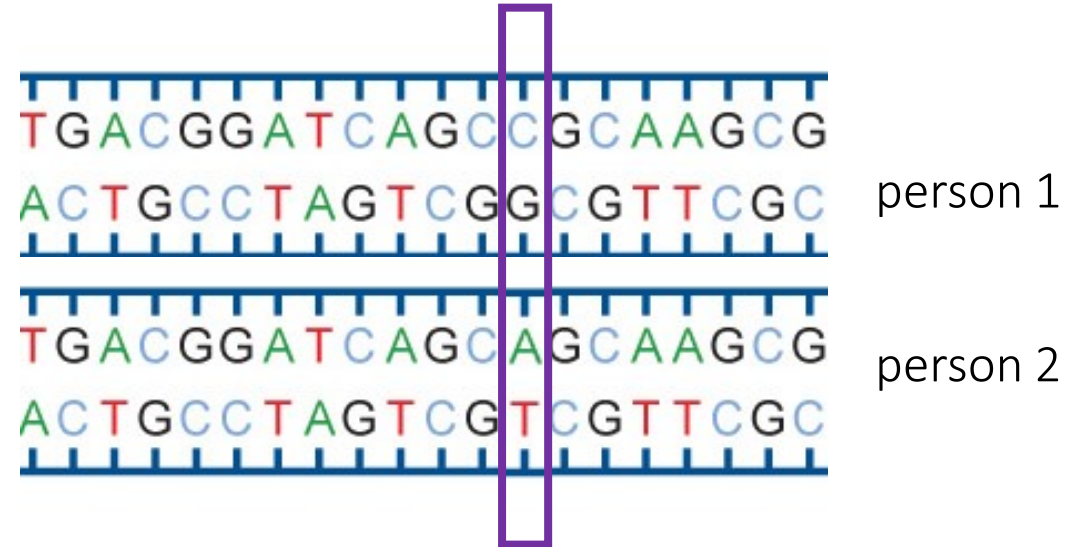
Introduction to common variation, quality control, GWAS, and PLINK (Part I)

Lucia Colodro Conde and Katrina Grasby

Introduction to common variation



adenine (A), thymine (T), cytosine (C), guanine (G)



Genetic variation: differences in the sequence of DNA among individuals.

Mutation: a newly arisen variant

Genetic variant: any specific region of the genome which differs between two genomes.

Allele: version of a variant

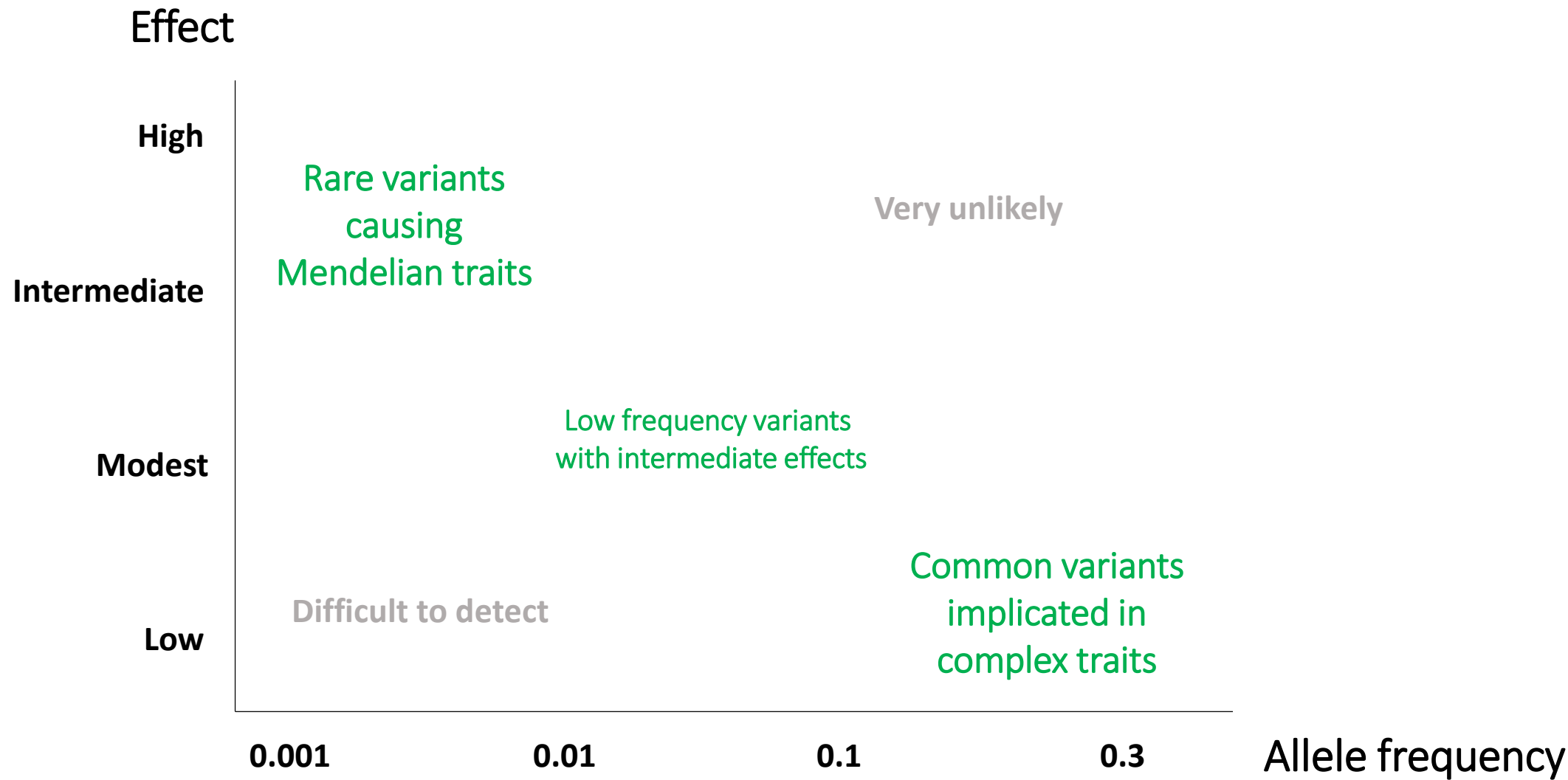
Allele frequency: incidence of an allele in a population.

Minor allele frequency (MAF): frequency at which the less common allele occurs in a given population.

Rare variant: a genetic variant present in $< 1\%$ of the alleles in the population

Common variant: a genetic variant present in $\geq 1\%$ of the alleles in the population

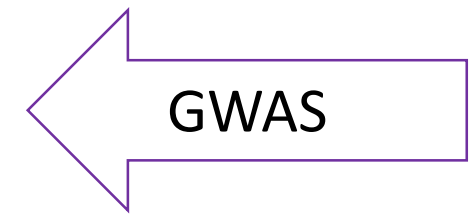
Note 1% is arbitrary



Examples of genetic variation

Sequence variation

- **Single nucleotide**
- substitutions
- insertions | 'indels'
- deletions



Structural variation

- **2bp to 1,000bp**
- VNTRs: microsatellites, minisatellites
- indels
- inversions
- di-, tri-, tetranucleotide repeats

- **1kb to submicroscopic**
- copy number variants
- segmental duplications
- inversions, translocations
- copy number variant regions
- microdeletions, microduplications

- **Microscopic to subchromosomal**
- segmental aneusomy
- chromosomal deletions (losses)
- chromosomal insertions (gains)
- chromosomal inversions
- intrachromosomal translocations
- chromosomal abnormality
- heteromorphisms
- fragile sites

- **Whole chromosomal to whole genome**
- interchromosomal translocations
- ring chromosomes, isochromosomes
- marker chromosomes
- aneuploidy
- aneusomy

Knight JC (2009). Genetics and the general physician: insights, applications and future challenge. *QJM*.

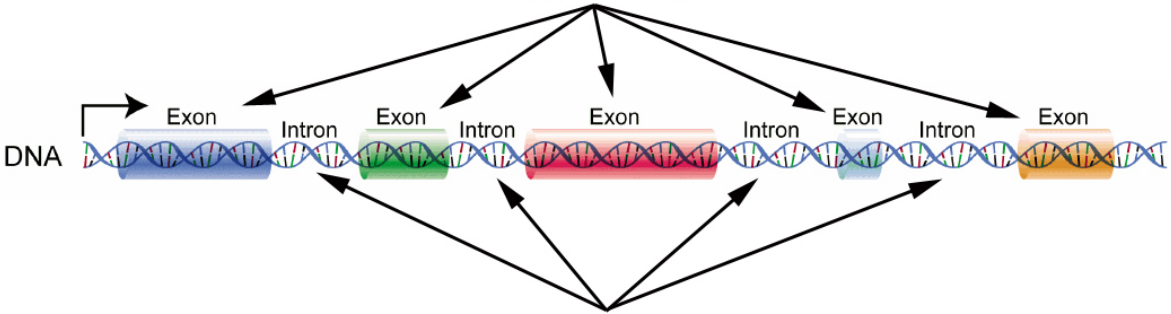
Sherer SW et al (2007). Challenges and standards in integrating surveys of structural variation. *Nat Genet*

SNP (single nucleotide polymorphism):

variation at a single base pair in a DNA sequence among individuals.

	Chrom.	DNA sequence	Genotype
Person 1	Mat	GTA ACTTGGGATCT A GACCAATAGAT	A A
	Pat	GTA ACTTGGGATCT A GACCAATAGAT	
Person 2	Mat	GTA ACTTGGGATCT A GACCAATAGAT	A C
	Pat	GTA ACTTGGGATCT C GACCAATAGAT	
Person 3	Mat	GTA ACTTGGGATCT C GACCAATAGAT	C C
	Pat	GTA ACTTGGGATCT C GACCAATAGAT	

Coding Regions



Non-coding Regions

TYPES OF SNPs

Non-coding region

They may affect regulation of genes in coding regions

Coding region

Synonymous

They do not change the amino acid sequence

Non-Synonymous

They change the amino acid sequence

Missense

They change one amino acid in a protein (which may have an effect in the protein)

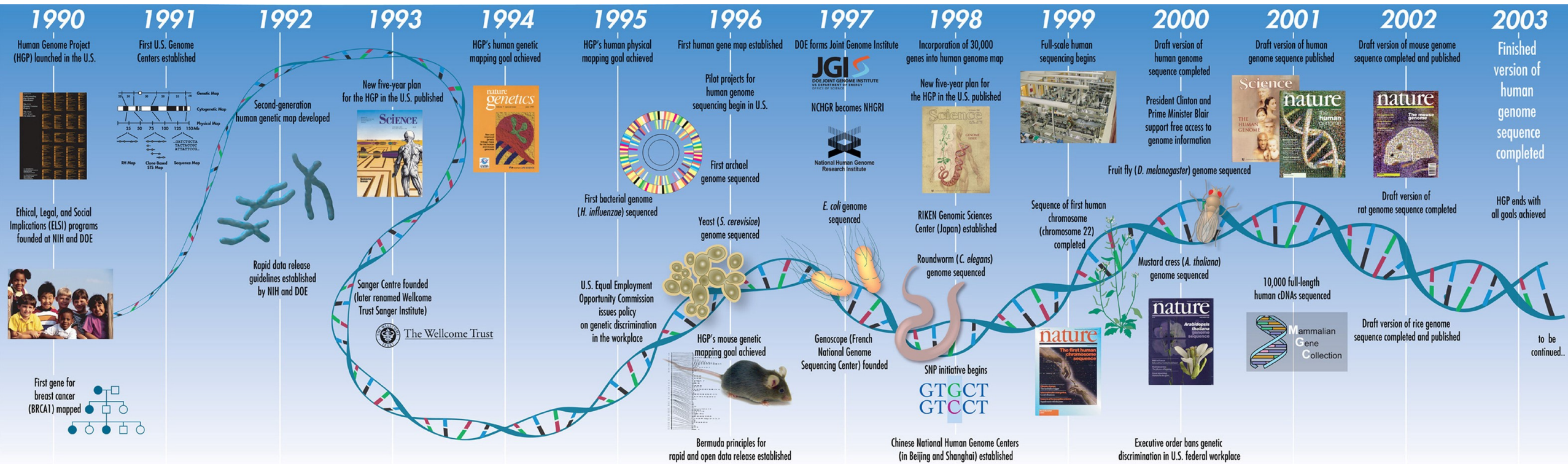
Nonsense

They produce a premature stop codon, that results in an incomplete and usually non-functional protein product

Insertion–deletion variants (indels):

one or more base pairs are present in some genomes but absent in others in relation to the reference

	Chrom.	DNA sequence	Genotype
Person 1	Mat	GTA ACTTGGGATCT GAT GACCAGATAG	R D
	Pat	GTA ACTTGGGATCT --- GACCAGATAG	
Person 2	Mat	GTA ACTTGGGATCT GAT GACCAGATAG	R R
	Pat	GTA ACTTGGGATCT GAT GACCAGATAG	
Person 3	Mat	GTA ACTTGGGATCT --- GACCAGATAG	D D
	Pat	GTA ACTTGGGATCT --- GACCAGATAG	



Collins et al 2003, A vision for the future of genomics research, *Nature*

27 October 2005 | www.nature.com/nature

THE INTERNATIONAL WEEKLY JOURNAL OF SCIENCE

nature



HapMap (haplotype map) Project

270 samples:

30 parent-offspring trios of the Yoruba from Ibadan, Nigeria (YRI)

30 trios of Utah residents with European ancestry (CEU)

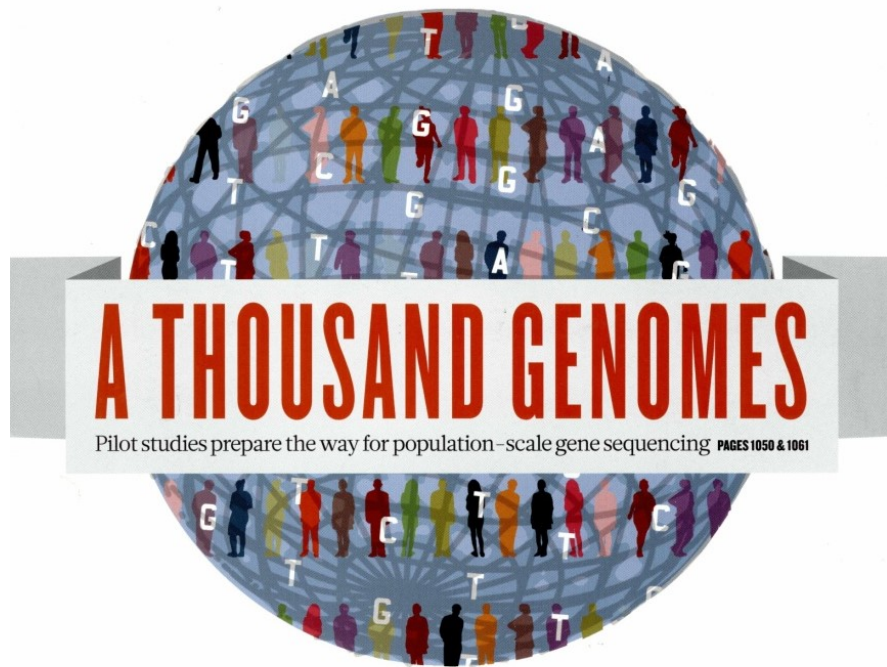
45 individuals from Beijing, China (CHB)

45 individuals from Tokyo, Japan (JPT)

The International HapMap Consortium (2005). A haplotype map of the human genome. *Nature*.

nature

THE INTERNATIONAL WEEKLY JOURNAL OF SCIENCE



1000 Genomes Project

Phase 1: 1,092 individuals from 14 populations..

Phase 3: 2,504 individuals from 26 populations (~500 samples form each 5 continental ancestry groups, with ~5 populations for each group)

Population	Code	Population Color	Continental Group Color	Analysis Panel	Phase 1	Phase 3
African ancestry						
Esan in Nigeria	ESN			AFR		99
Gambian in Western Division, Mandinka	GWD			AFR		113
Luhya in Webuye, Kenya	LWK			AFR	97	99
Mende in Sierra Leone	MSL			AFR		85
Yoruba in Ibadan, Nigeria	YRI			AFR	88	108
African Caribbean in Barbados	ACB			AFR/AMR		96
People with African Ancestry in Southwest USA	ASW			AFR/AMR	61	61
Americas						
Colombians in Medellin, Colombia	CLM			AMR	60	84
People with Mexican Ancestry in Los Angeles, CA, USA	MXL			AMR	66	64
Peruvians in Lima, Peru	PEL			AMR		85
Puerto Ricans in Puerto Rico	PUR			AMR	55	104
East Asian ancestry						
Chinese Dai in Xishuangbanna, China	CDX			EAS		93
Han Chinese in Beijing, China	CHB			EAS	97	103
Southern Han Chinese	CHS			EAS	100	105
Japanese in Tokyo, Japan	JPT			EAS	89	104
Kinh in Ho Chi Minh City, Vietnam	KHV			EAS		99
European ancestry						
Utah residents (CEPH) with Northern and Western European ancestry	CEPH			EUR	85	99
British in England and Scotland	GBR			EUR	89	91
Finnish in Finland	FIN			EUR	93	99
Iberian Populations in Spain	IBS			EUR	14	107
Toscans in Italia	TSI			EUR	98	107
South Asian ancestry						
Bengali in Bangladesh	BEB			SAS		86
Gujarati Indians in Houston, TX, USA	GIH			SAS		103
Indian Telugu in the UK	ITU			SAS		102
Punjabi in Lahore, Pakistan	PJL			SAS		96
Sri Lankan Tamil in the UK	STU			SAS		102
Total					1092	2504

HUMAN STEM CELLS

BEYOND THE COURT CASE
Implications for the law, industry and ethics
PAGE 1031

OCEAN PRODUCTIVITY

PHOSPHATE DOWN THE AGES
Key nutrient plentiful after 'snowball' Earth
PAGES 1052 & 1088

AUTUMN BOOKS

THE RECURRING UNIVERSE
Lee Smolin on Roger Penrose's grand idea
PAGE 1034

NATURE.COM/NATURE

28 October 2010 £10
Vol. 467, No. 7319



The 1000 Genomes Project Consortium (2012). An integrated map of genetic variation from 1,092 human genomes. *Nature*.

The 1000 Genomes Project Consortium (2015). A global reference for human genetic variation. *Nature*.

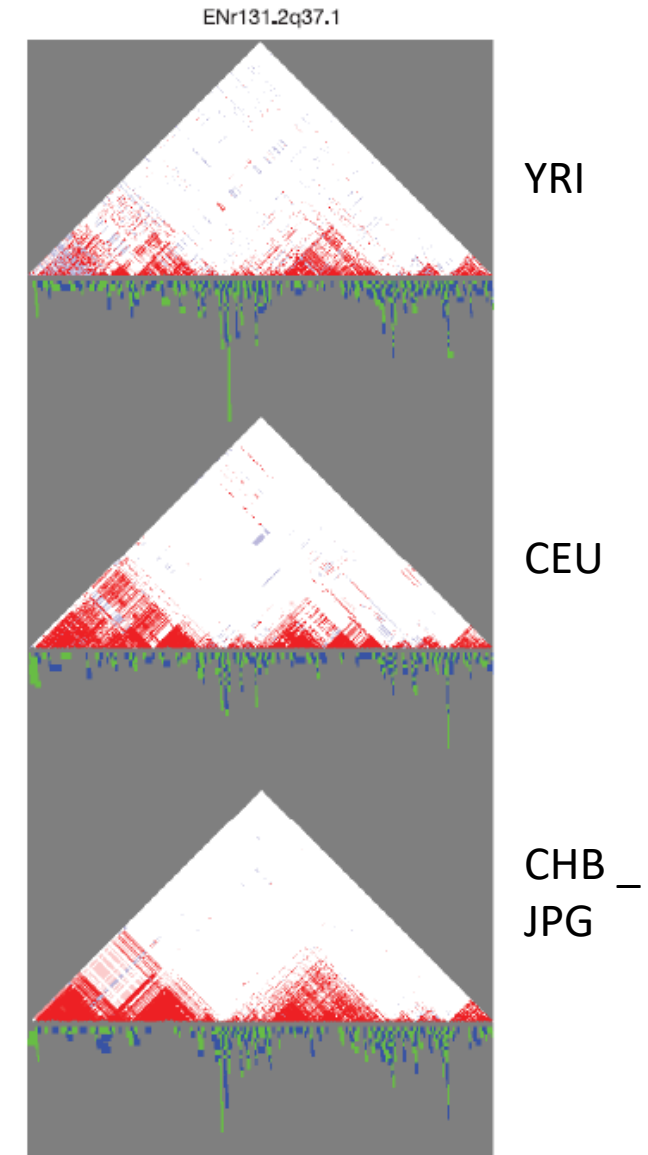
The Haplotype Reference Consortium (HRC)

nature
genetics

A reference panel of 64,976 haplotypes for genotype imputation

These projects have provided information on:

- Patterns of human common genetic variation
- Linkage disequilibrium (LD) and allele frequencies differences in populations
- Tag SNPs for the design of SNP arrays to facilitate imputation and GWAS



Practical. Where to start with.

GWAS have been facilitated by the development of relatively inexpensive SNP arrays.

How do we make the information provided by SNP arrays usable?