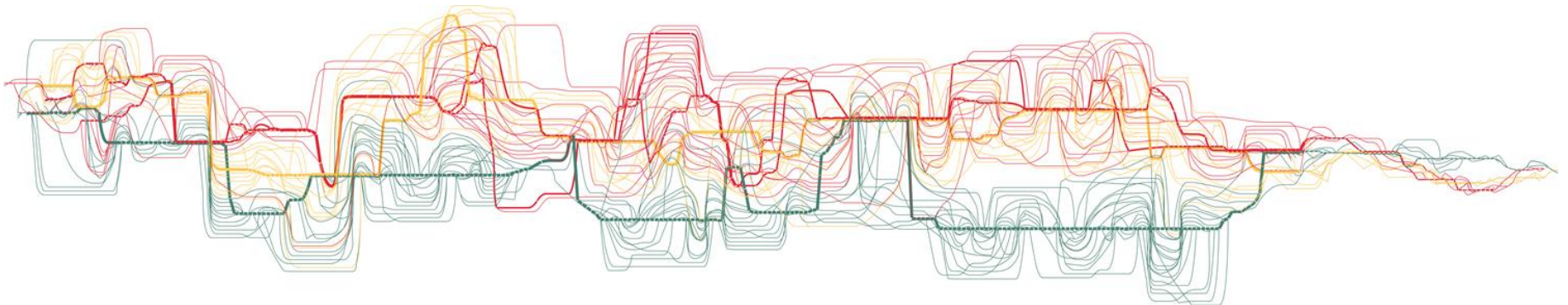
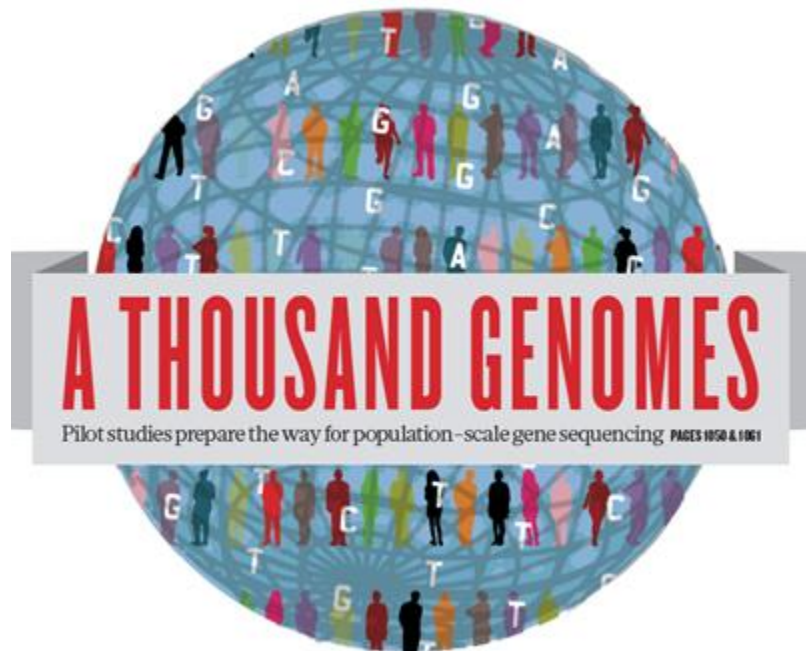


# Making big data work for biomedicine

Gil McVean





# How a DNA first for girl, 4, changed a family's world

A child with a skull abnormality has blazed a trail by having her entire genetic code read. **Mark Henderson writes**

A four-year-old girl has become the first person in Britain to have her entire genetic code read to identify the cause of a disease, in a landmark development that illustrates how personal genetics is changing healthcare.

Katie Warner, from Saffron Walden, Essex, and her parents John and Maria had their genomes sequenced by scientists at the University of Oxford to pin-

diagnosis has been difficult. We might now have a label that makes everything crystal clear. Katie's definitely behind, there are no two ways about it. But we've had problems getting her statemented for school. We know that her condition is going to affect her learning, and we can do something about that immediately: it's going to make the battle we have with education authorities much easier. Starting to understand why, and what she'll be able to do, and not, is a big help."

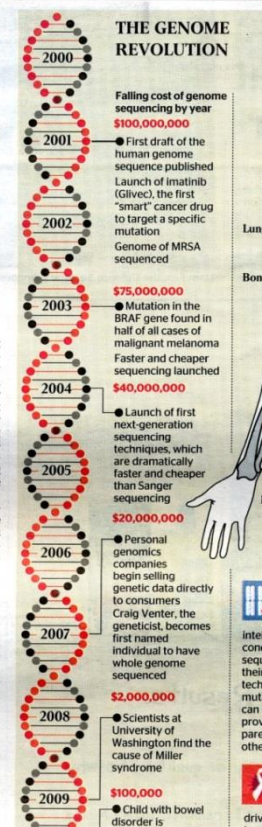
Several children have been diagnosed by genome sequencing in the US, including one who, as a result, was successfully treated for a bowel disorder with a bone marrow transplant. Katie is the first child in Britain to benefit.

Katie has a condition called cranioy-nosis, which causes sections of her skull to fuse early so there is insufficient room for her brain to grow. She has had two operations to relieve pressure on her brain, one when she was just seven months old. The precise cause was unknown, making it difficult for her doctors to give a prognosis.

Though the NHS does not yet provide genome sequencing for unexplained disorders, Katie was referred to Andrew Wilkie, a consultant clinical geneticist at the University of Oxford who specialises in craniofacial disorders. Professor Wilkie is involved in an Oxford research project supported by Illumina, a DNA sequencing company, which is sequencing 500 genomes of people with serious diseases and their



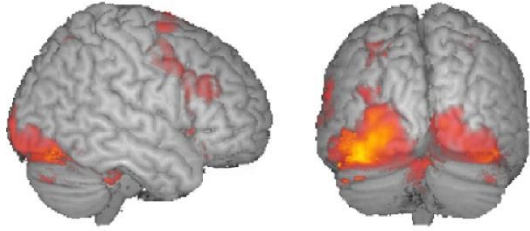
Katie Warner, who has a cranio-facial condition, with her mother Marie



# Big data is...

- Large, population-scale data sets typically made possible by innovations in high throughput technology
  - Genome sequencing
  - Mobile and internet technology
  - Imaging and automated image processing
- Data sets that are large, high dimensional, semi-structured and highly heterogeneous
  - Large, requires distributed and cloud computing
  - Cannot be stored in standard database structures
  - Hard to summarise / visualise
  - Collected in many different ways by many different agents
  - Requiring new statistical and computational methods to analyse

# Medical big data sources



Imaging

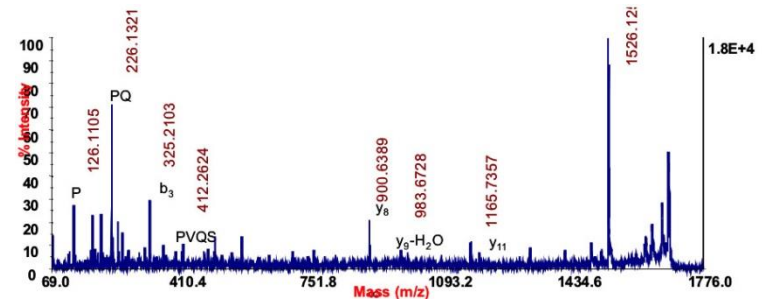
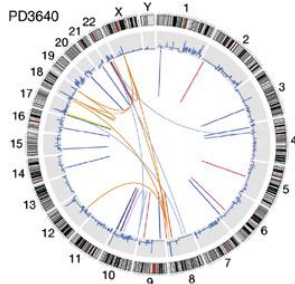


Germ-line genome



Electronic  
medical  
records

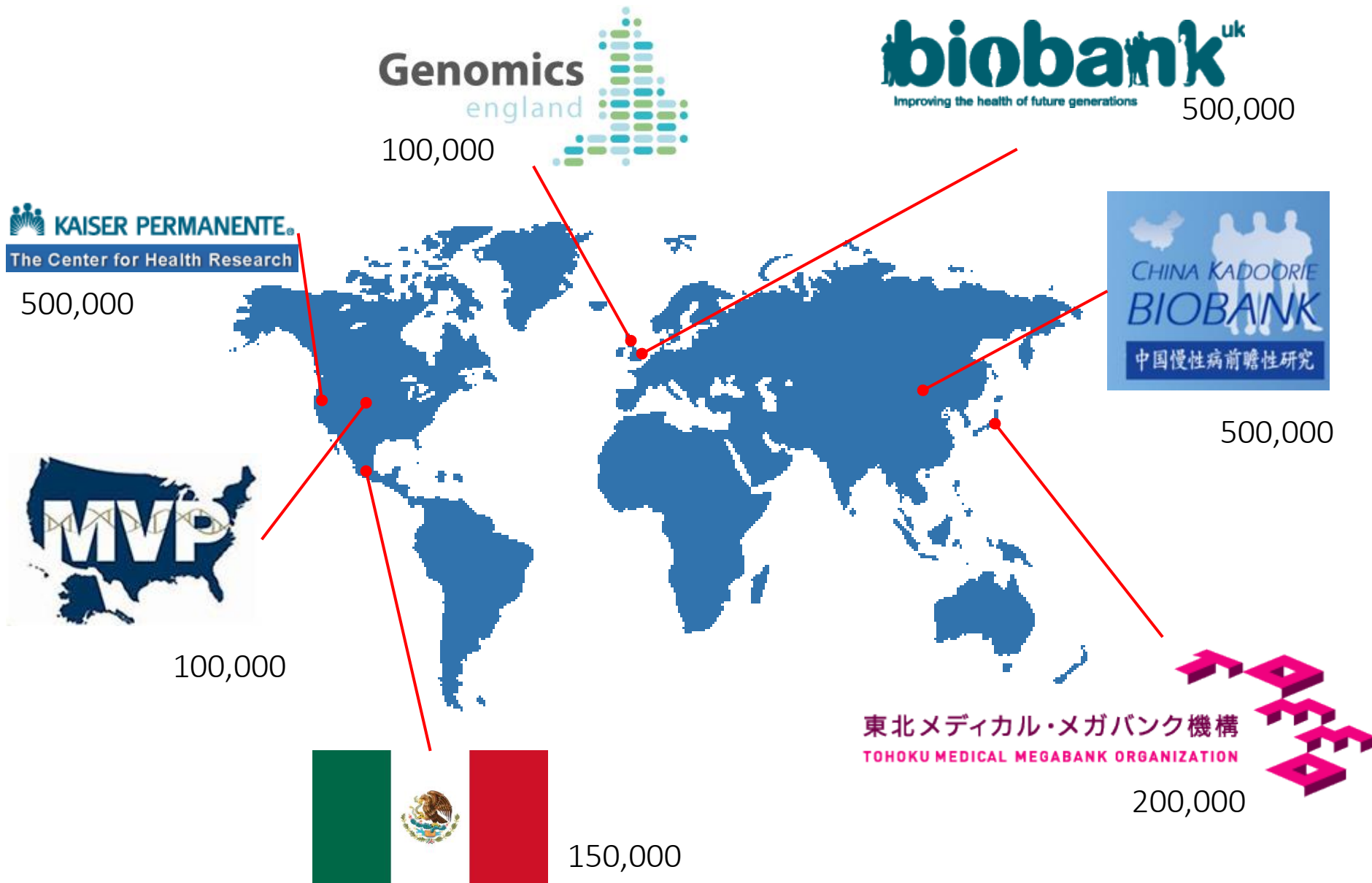
Tumour  
profile



Biochemistry



# Population-scale medical cohorts with genomic data



# The Oxford Big Data Institute



Germ-line genomics



Infectious disease  
surveillance

Electronic Medical  
Record linkage

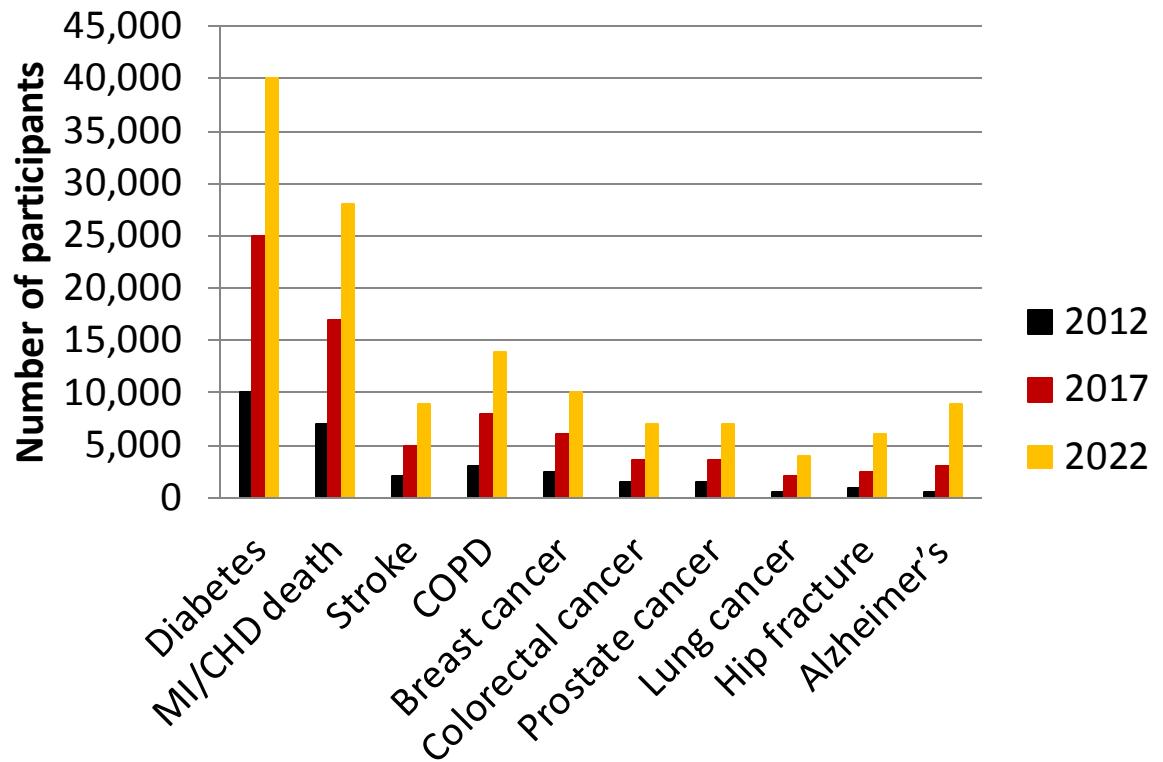
Deep phenotypes in  
Large cohorts



# Big data and epidemiology



**biobank<sup>uk</sup>**  
Improving the health of future generations



wellcome trust

Northwest  
REGIONAL DEVELOPMENT AGENCY

DH Department  
of Health

The Scottish  
Government

Llywodraeth Cymru  
Welsh Assembly Government

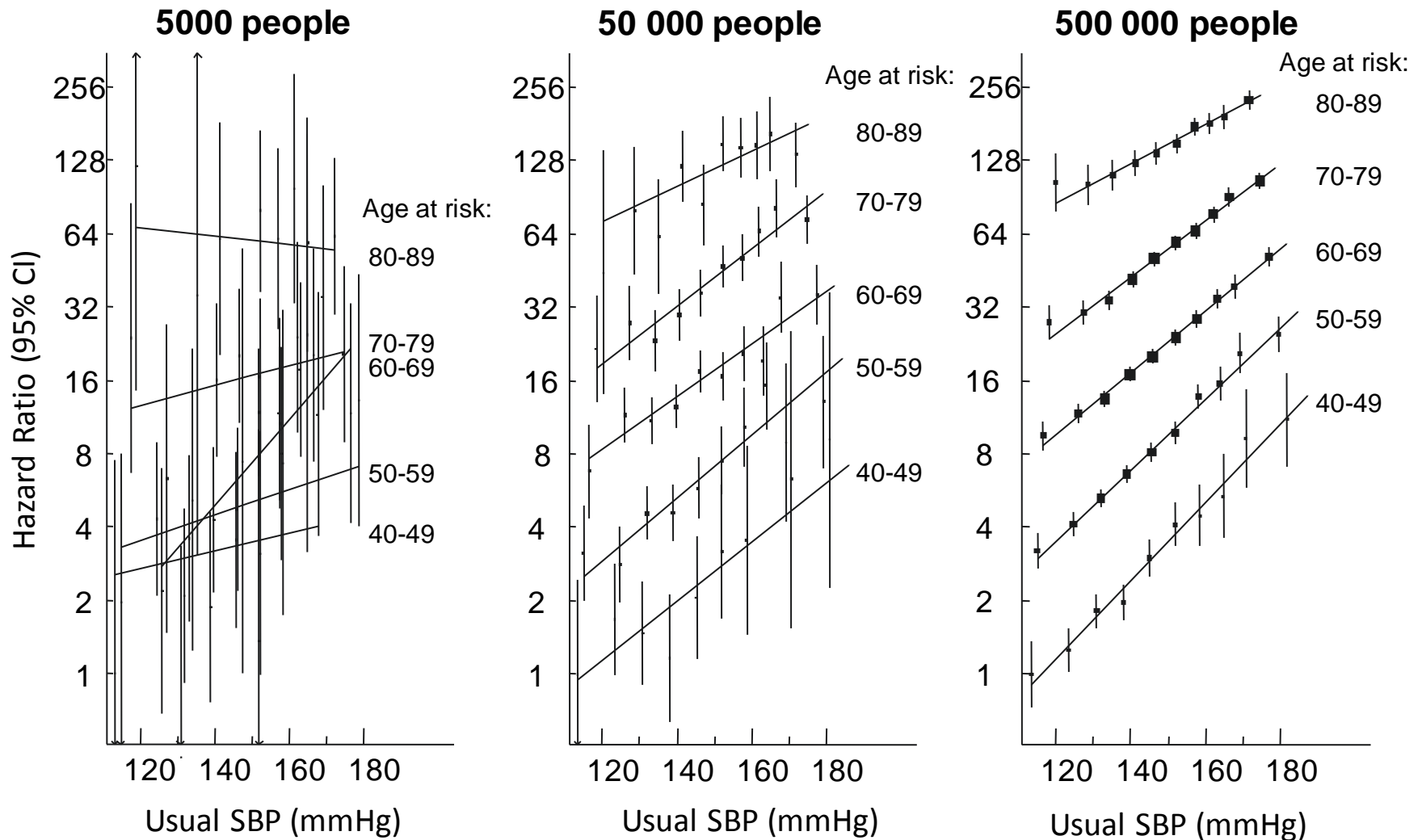
British Heart  
Foundation

# Enhanced phenotyping

- Web-based diet questionnaires on 300,000
  - cognitive assessments planned
- Repeat assessments on 20,000
  - further repeat every few years
- Wrist-worn accelerometers on 100,000
- Standard panel of laboratory assays + genotyping on 500,000
- Imaging visit in 100,000
  - to include whole body and brain MRI, carotid ultrasound, bone (DEXA)
- On-going linkage to additional data sources
  - hospital, primary care, disease register
  - environment



# The value of large numbers: Ischaemic heart disease and systolic blood pressure

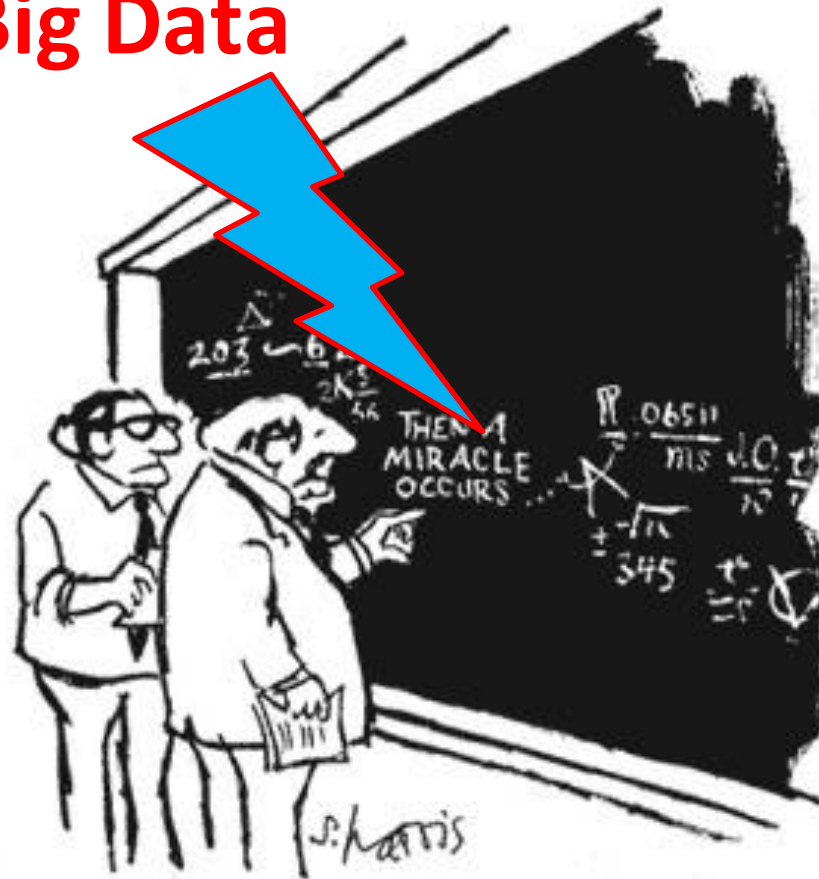


# The promise of biomedical big data

- Better diagnosis
- Better treatment choice
- Improved target discovery
- Improved target validation
- Better outcomes



# Big Data



"I THINK YOU SHOULD BE MORE EXPLICIT HERE IN STEP TWO."

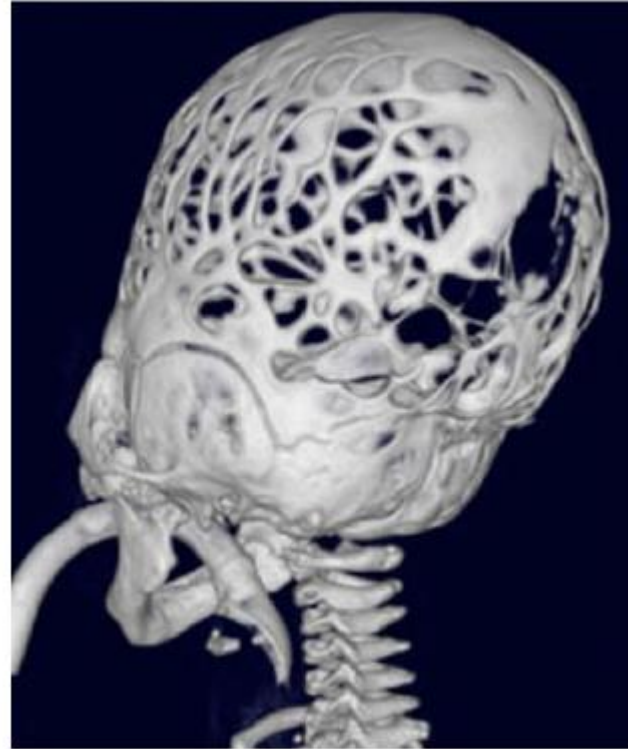
# Putting genomics at the heart of Big Data

- Genetics has a direct and unequivocal relationship to phenotype
  - Germ-line gives exposure from birth
  - Finding the gene immediately informs about the patient
- Genomic data is accurate and easy to collect on a population scale
  - SNP genotype data -> GWAS on 100,000s
  - Sequencing at a population scale
- Genomic data can be used to probe causal relationships between biomarkers and disease
  - Natural variation mimics pharmaceutical interventions
- Genomic data provides a fingerprint for monitoring infectious disease control programmes
  - Can established transmission networks at micro and macro scale

# Big data in the clinic



# Genome sequencing as a clinical tool



# The value of whole genome sequencing

- Whole-genome sequence is the only data type that can detect all types of information relevant to pathology in a single go:

Data type	Large-scale structural changes	Balanced translocations	Distant consanguinity	Uniparental disomy	Novel / known coding variants	Novel /known non-coding variants
Targeted gene sequencing	No	No	No	No	Yes	No
SNP arrays	Yes	No	Yes	Yes	No	No
Array CGH	Yes	No	No	No	No	No
Exome	Partial	No	Partial	Partial	Yes	No
Whole genome	Yes	Yes	Yes	Yes	Yes	Yes

# WGS500 – initiated in 2011

- Collaboration

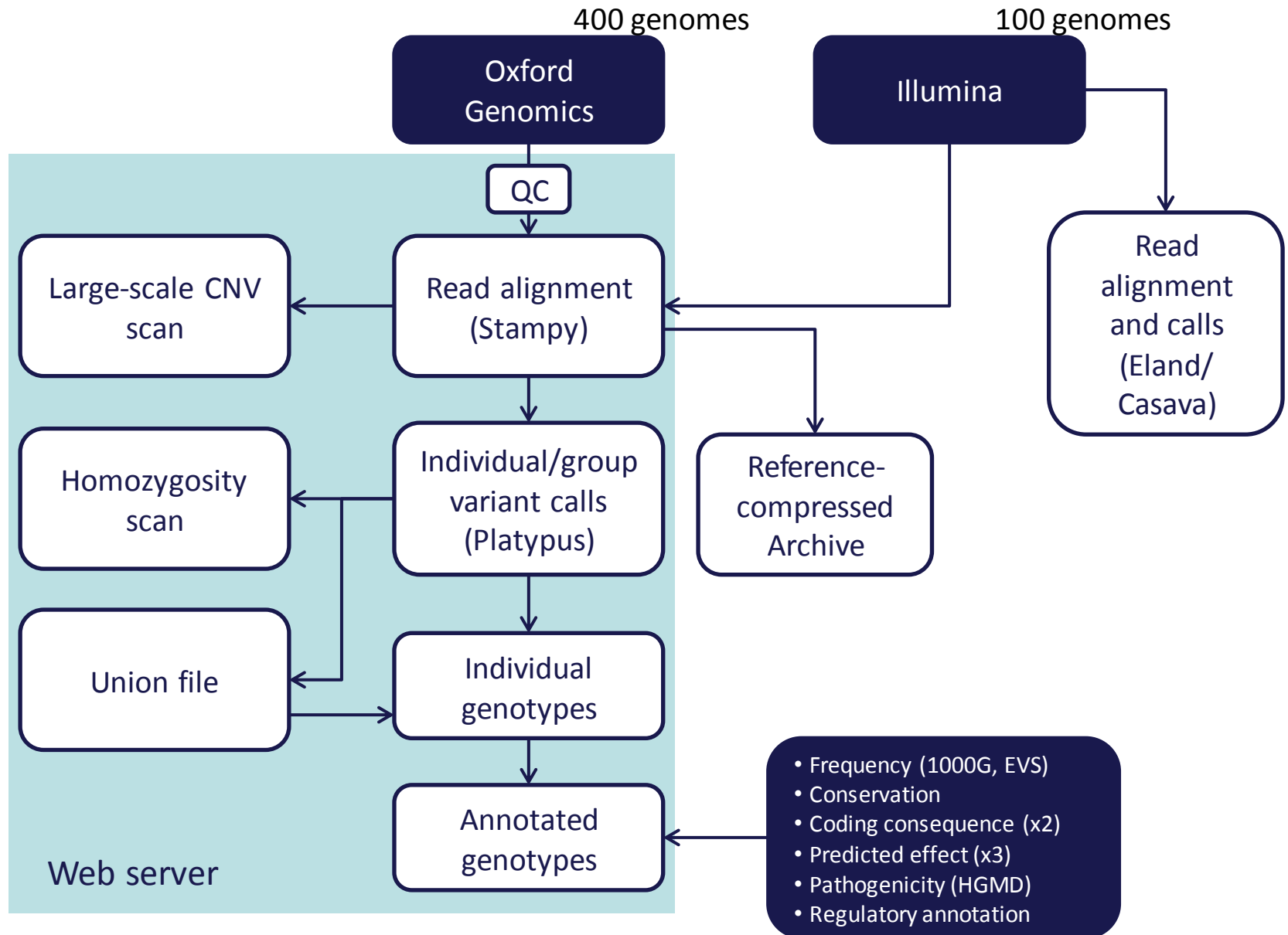


illumina®

- Sequence 500 genomes at 30x
- Diverse set of diseases (>40 phenotypes)
  - Mendelian disorders (without mutation in screened genes)
  - Sporadic (extreme) immune phenotypes
  - Cancers
- Diverse set of experimental designs
  - Familial: Linkage information, trios, quartets
  - Cancer: Tumour-normal, metastases, multiple-mets, ..
- Substantial follow-up (screening and functional) to establish candidacy

What is big data about this study?

# 1. Infrastructure





http://www.ebi.ac.uk/biomed/p110226/

malbec/seq/28303/

pubmed | bioRxiv | HGMD | WGS300 | M2browser | bioinformatics

Analysis of reference that is covered at least once. Estimated heterozygosity (one individual). Sample contamination can increase this estimate. Filters: low quality >= 20, mapping quality >= 30, pairs properly mapped, no indels, maximum 2 high-quality (Q20) variants in read, insert size >= 50 bp. \* Basic coverage is assigned one region that are covered at least once. \* Proportion of reads is equal to coverage in log10 percentage.

### Lane QC statistics and plots

Lane	% GC	% GC <sub>dup</sub>	$\sigma_{GC}$ (%GC)	insert $\pm$ MAD	% exonic	% exon cov'ge	% N	max <sub>cov</sub> %N	%lowQ	%lowQ <sub>red</sub>	avgQ
L1	40.2 $\pm$ 9.5	40.2 $\pm$ 9.5	0.61	401 $\pm$ 25	1.0	57.3	0.0	0.1	5.2	44.8	32.3
L2	39.2 $\pm$ 10.8	39.8 $\pm$ 9.7	0.62	401 $\pm$ 25	1.0	56.7	0.0	0.0	6.7	39.0	31.6

G+C histogram

Insert size histogram

Coverage histogram

Exon/peptide coverage distribution

http://www.ebi.ac.uk/biomed/p110226/

malbec/seq/wgs500/projects/12/

pubmed | bioRxiv | HGMD | WGS300 | M2browser

### WGS 500

Projects Admin Reports

Projects - P110226

### P110226 Project Summary

#### Project Details

Project Number P110226

Disease Spermatocytic seminoma

Disease Type Cancer

PI Andrew Wilkie

Other Contacts Anne Ganiely

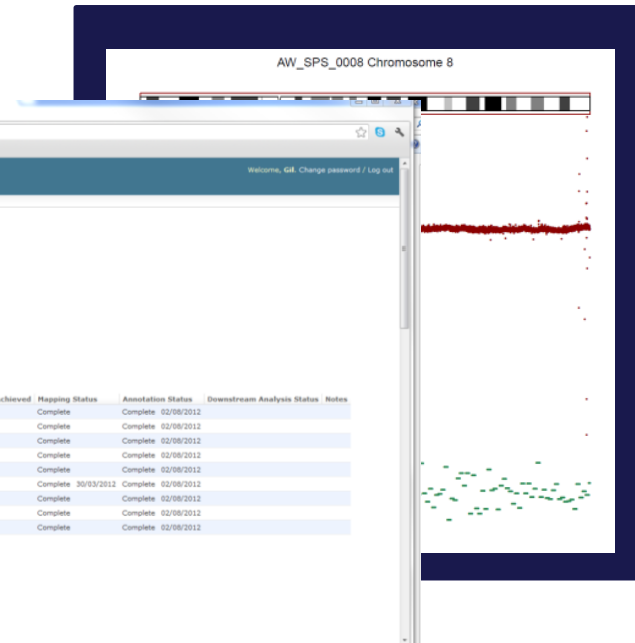
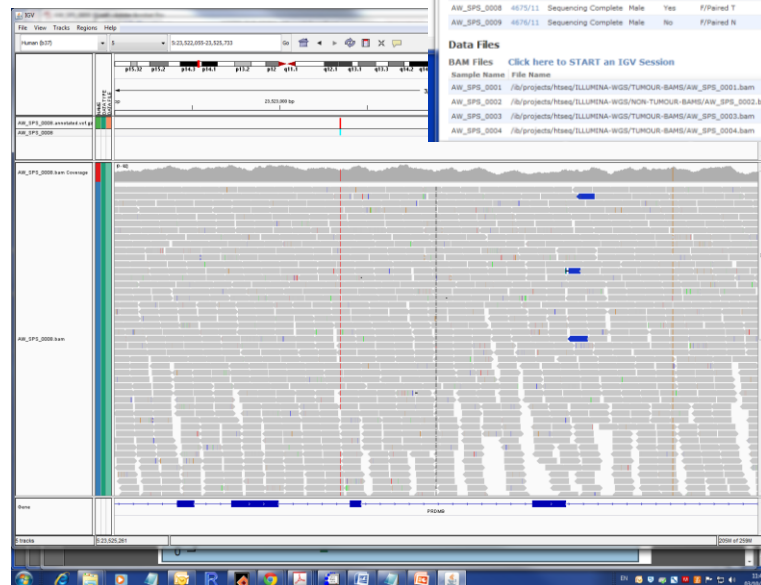
Start Date

Due Date

Analysts Simon McGowan, Anne Ganiely, Eleni Giannoutsou

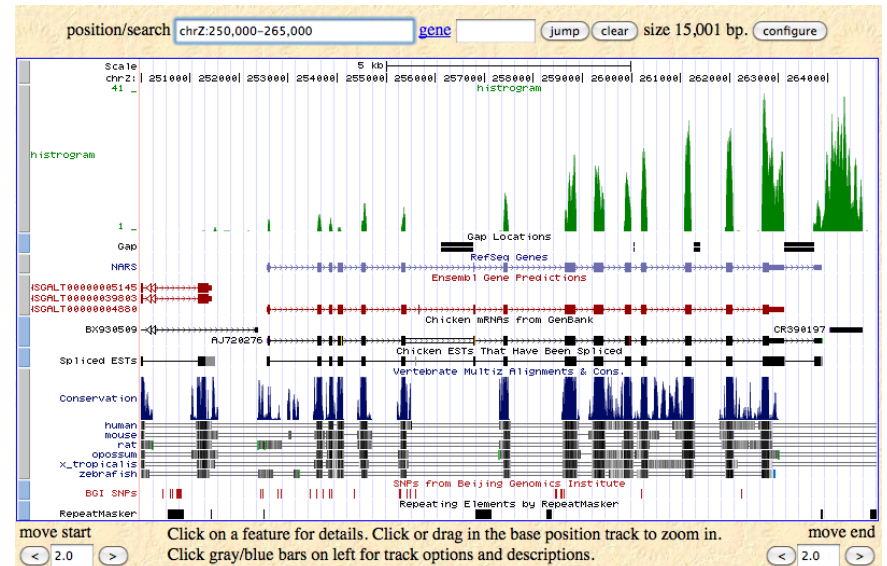
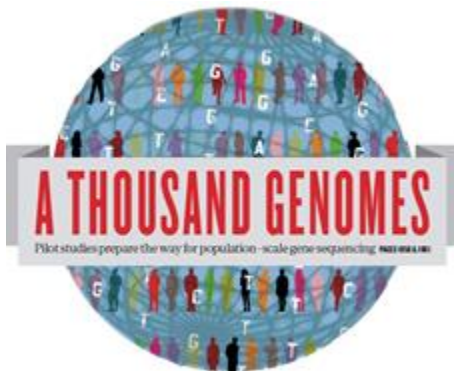
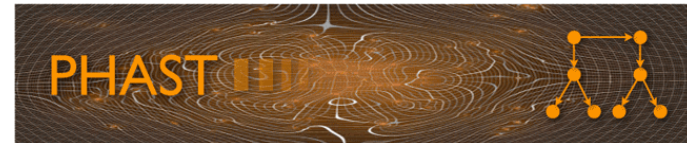
#### Samples

Sample Name	LIMS ID	Sample Status	Gender	Tumour	Relationship	Coverage Req
AWL_SPS_0001	4669/11	Sequencing Complete	Male	Yes	C/Parent T	25.0
AWL_SPS_0002	4669/11	Sequencing Complete	Male	No	C/Parent	25.0



Spermatocytic seminoma (SPS) is a rare germ cell tumour that is slow growing and occurs specifically during adulthood in older males. We originally sequenced 9 samples, including 5 frozen tumour samples (SPS1, SPS6, SPS8) (one of them being a bilateral case (SPS3, SPS4)) and their 4 matched controls (2 from blood genomic DNA (SPS2, SPS5) and 2 from DNA extracted from normal tissue adjacent to the tumours (SPS7, SPS9)). However, it was noted that one of the presumed 'bilateral' sample SPS3 had many more discordant calls than its SPS4 counterpart, suggesting a sampling mismatched. Typing of a few SNPs for SPS3 suggest that the sampling error originated upstream of the WGS process. As a result, the 'trio' has been re-genotyped as a pair (SPS4/SPS5) and SPS3 is therefore an unmatched sample.

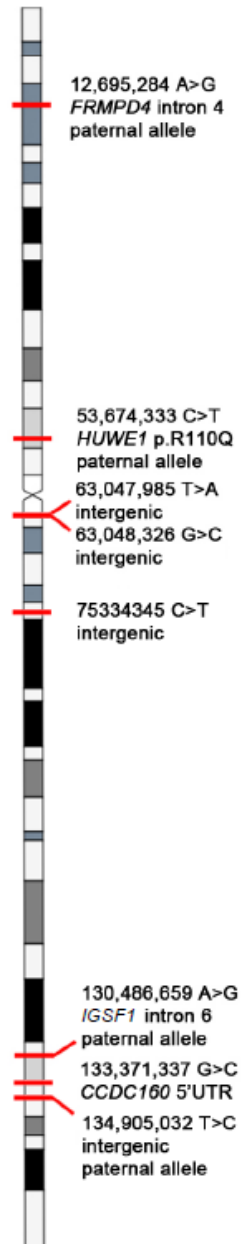
# 3. Use of heterogeneous external data



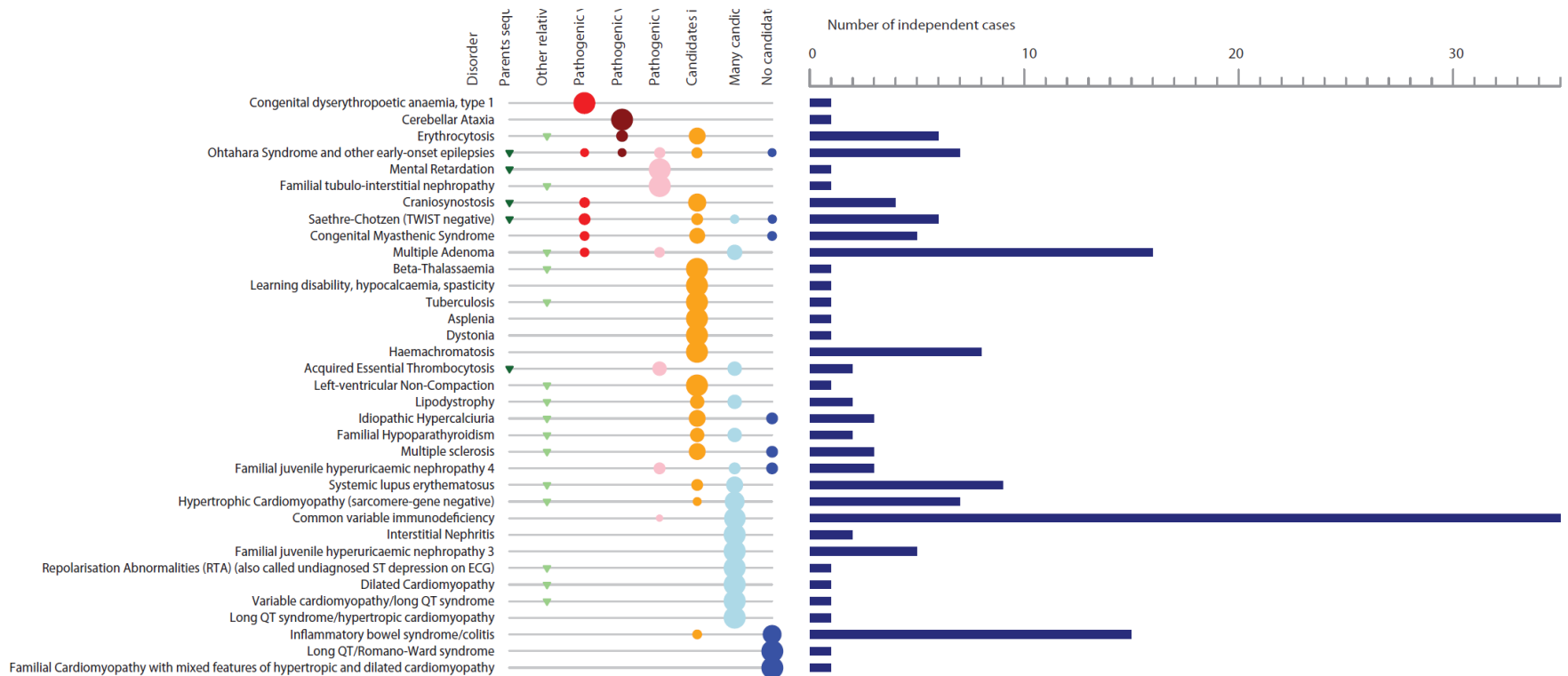
# What did we find?



- Sequenced child and both parents
- De novo mutation in *HUWE1*, a known mental retardation gene
- Skewed inactivation on the X chromosome towards chromosome with mutation
- Multiple additional de novo mutations
- No other *HUWE1* cases in >100 additional cases screened



# Over 25% of cases with clear diagnosis



What limits success?

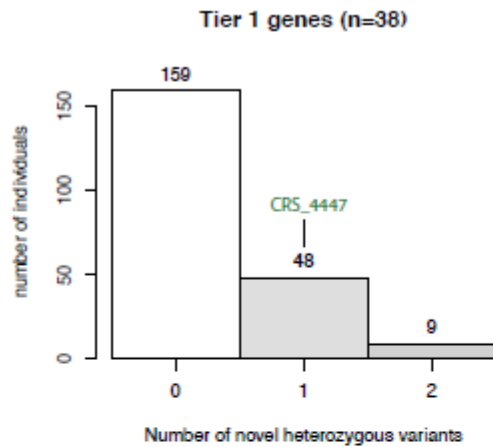


# 1. Ability to predict biological consequence

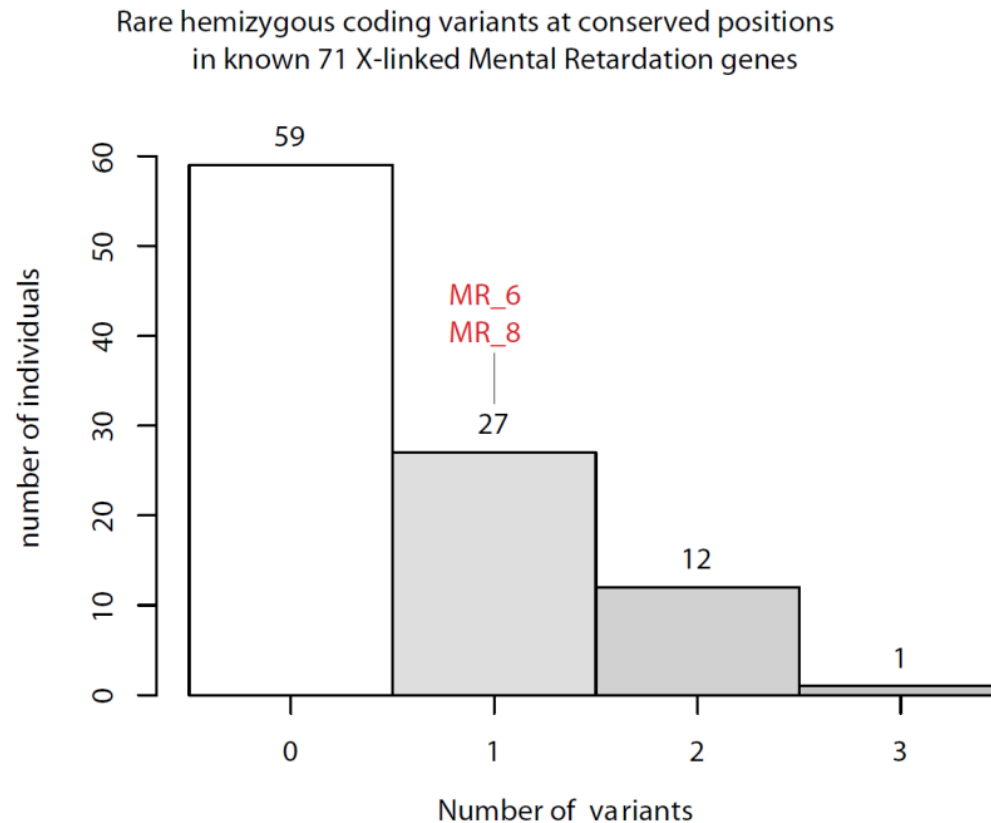
Known disease genes

Protein-protein  
interactions and related  
disorders

Same pathway



# Do 40% of males have mental retardation?



## 2. Sample size: The Genomics England 100k project

- 100,000 genomes sequenced in rare disease and cancer by 2017
- Linkage to medical records



# Global Alliance for Genomics and Health



**Global Alliance**  
for Genomics & Health

[Sign up for updates](#)



[ABOUT GLOBAL ALLIANCE](#)

[OUR WORK](#)

[PARTNERS](#)

[NEWS & EVENTS](#)

[CONTACT US](#)

**Collaborate. Innovate. Accelerate.**

Working together to share knowledge, create networks and accelerate advances in genomics and health.

[Learn More](#)

## What is the Global Alliance?

The Global Alliance for Genomics and Health  
[thehealth.org/about-global-alliance](http://thehealth.org/about-global-alliance)

## What is the Global Alliance doing?

The Global Alliance for Genomics and Health  
has doubled in size since its formation and the

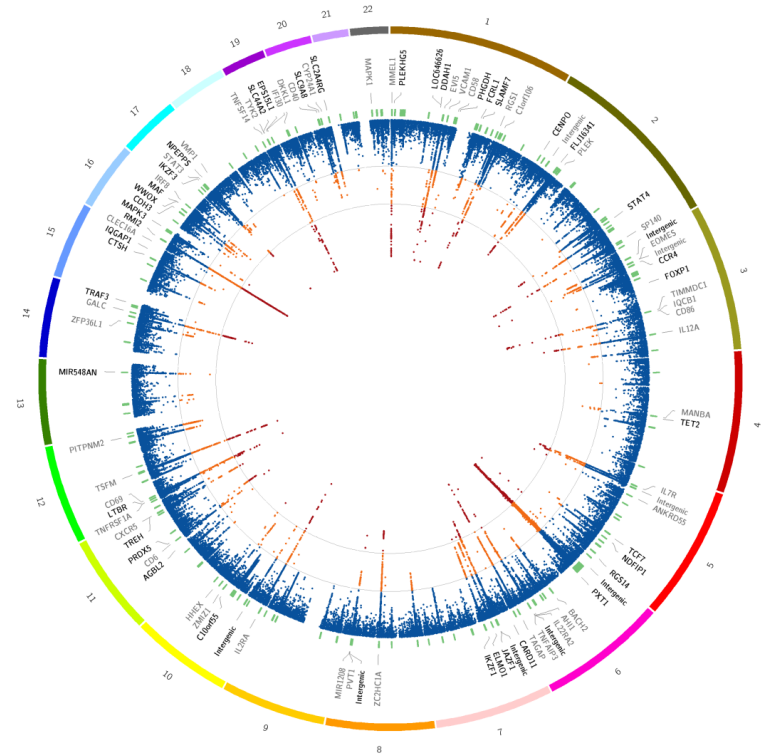
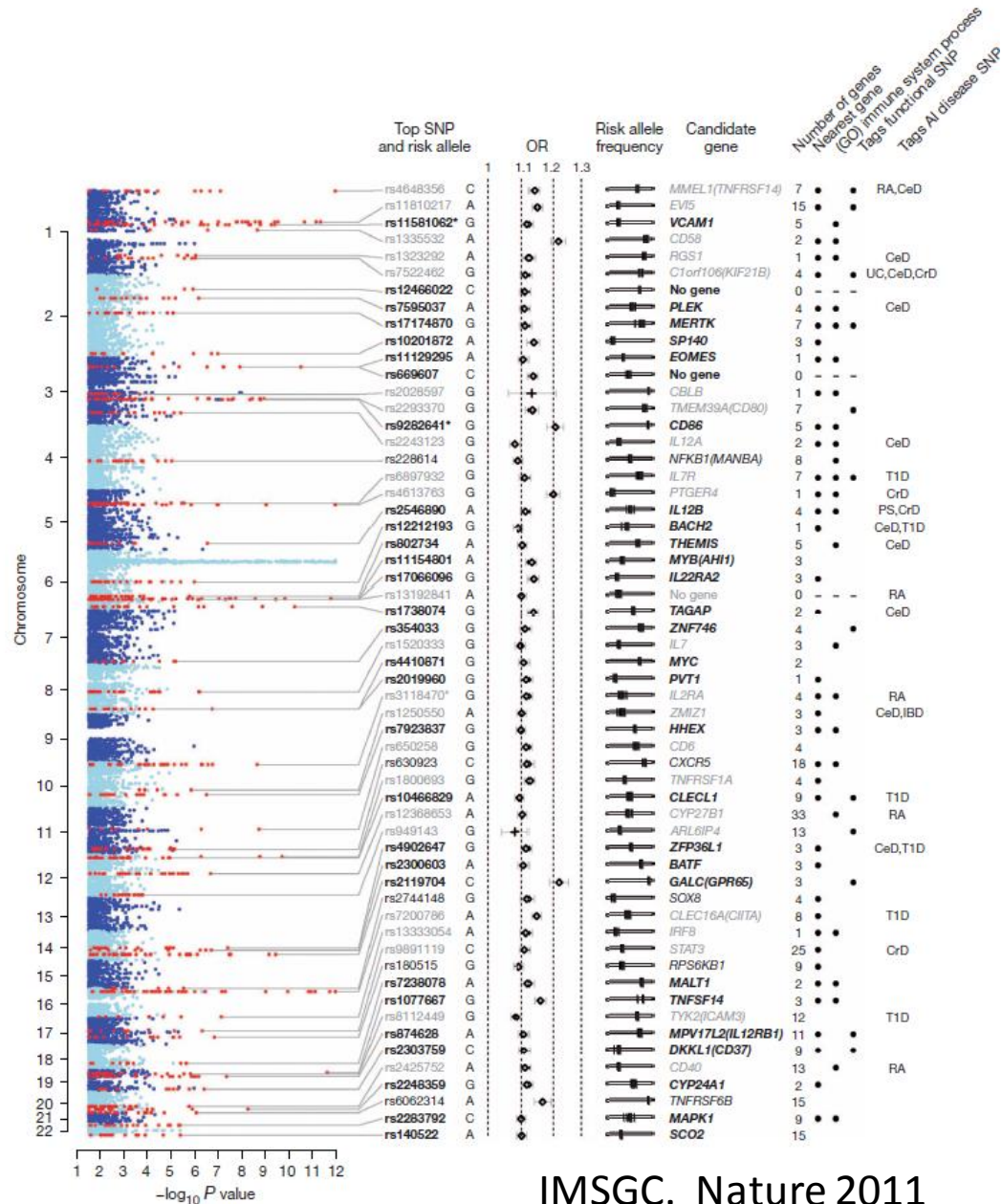
## Who is involved?

The Global Alliance for Genomics and Health  
(Global Alliance) is a broad and inclusive

What will large data sets deliver?



# 1. A (growing) understanding of complex disease mechanisms

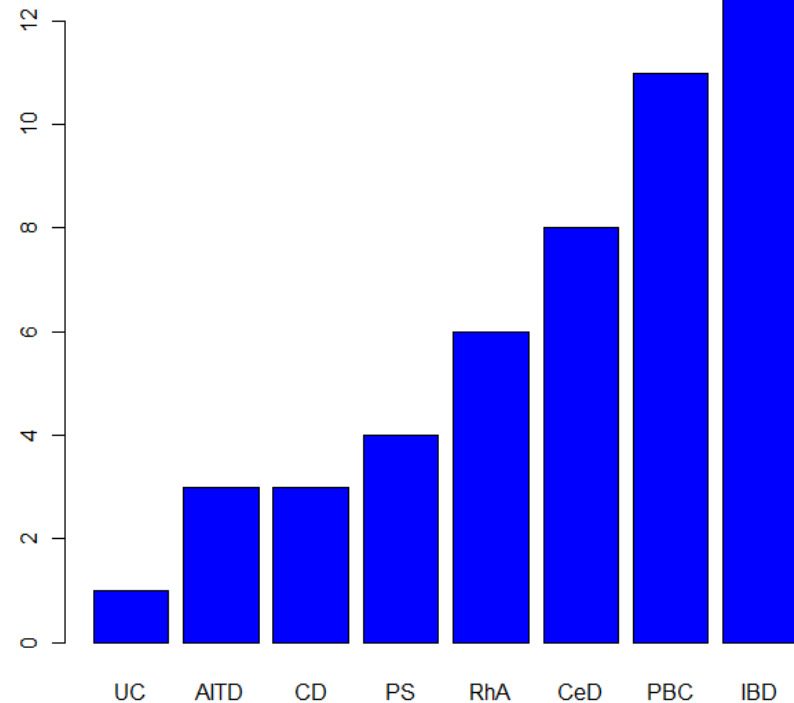


IMSGC. Nature Genet. 2013

IMSGC. Nature 2011

# What have we learned?

- About 113 current risk loci for multiple sclerosis
- Explains c. 30% sibling recurrence risk (10% of that is HLA)
- C. 30% of loci overlap with other autoimmune diseases



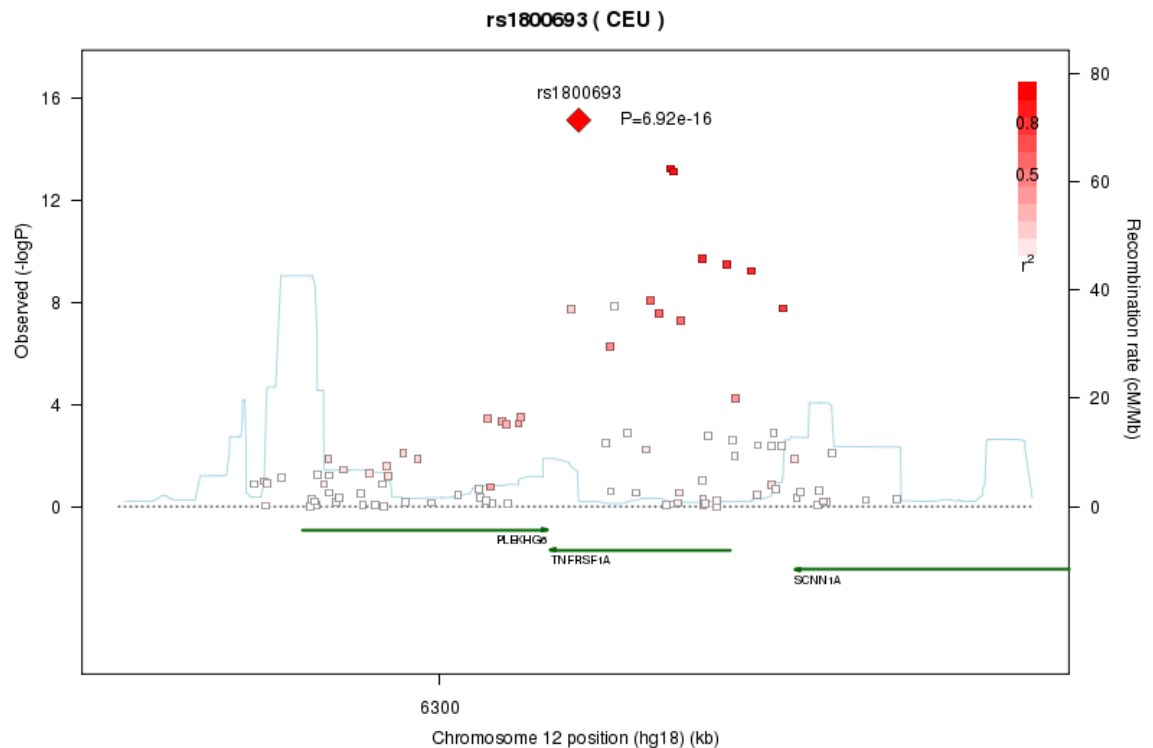
# Can we map causal variants from GWAS?

## ARTICLES

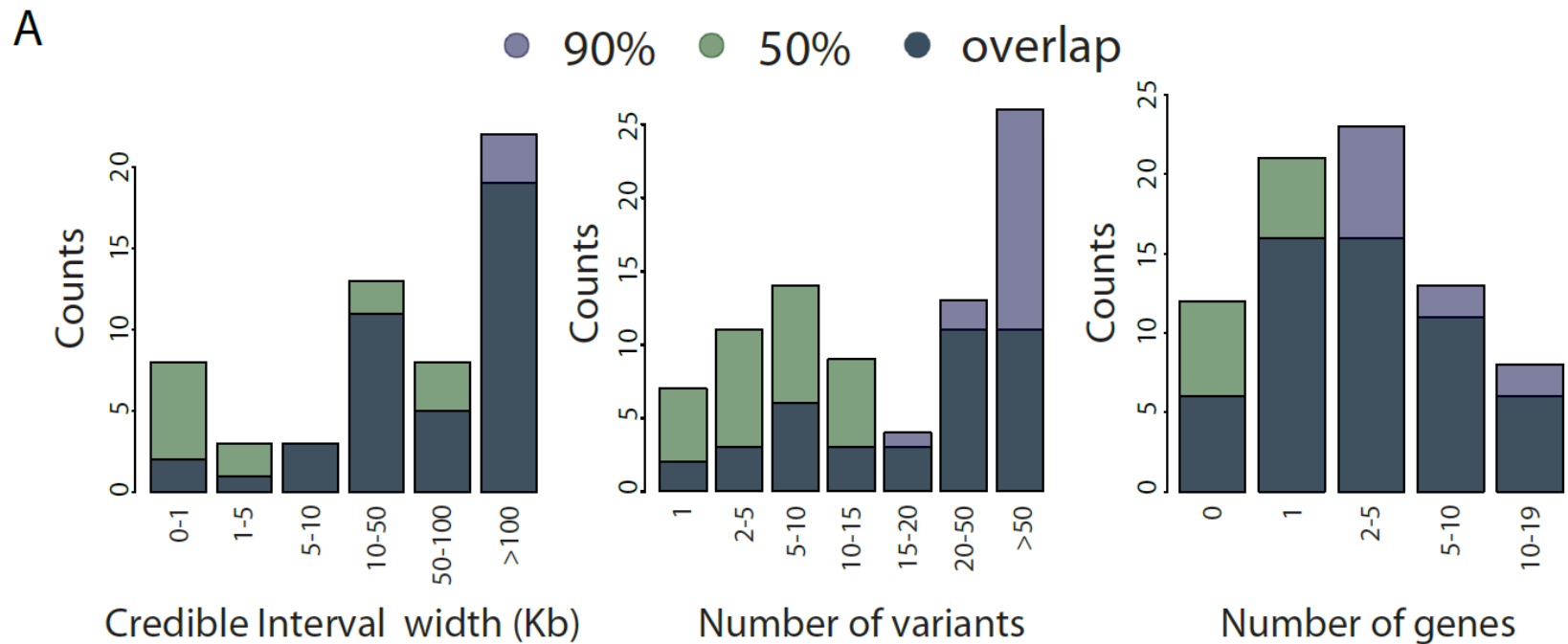
nature  
genetics

### Bayesian refinement of association signals for 14 loci in 3 common diseases

The Wellcome Trust Case Control Consortium<sup>1,2</sup>

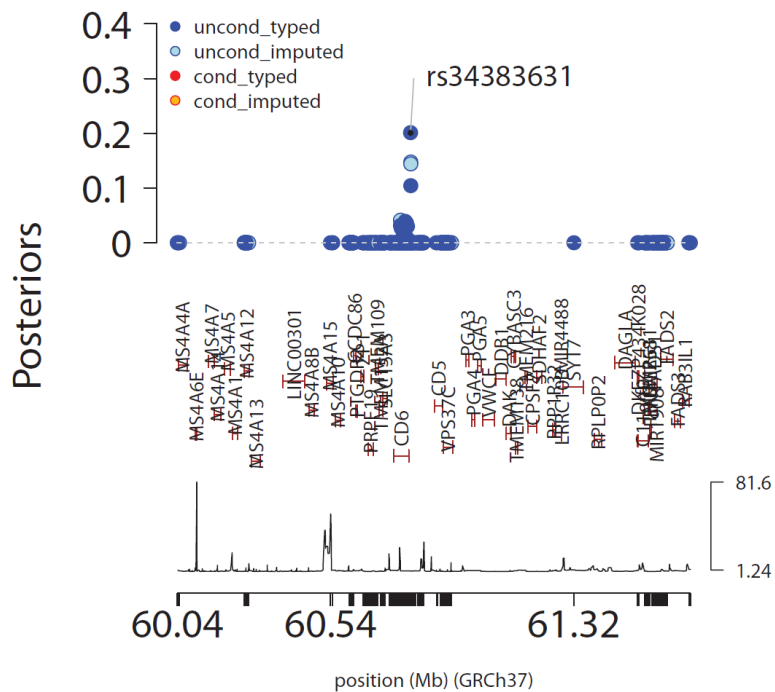


# About 5% of signals fine-map to <5 variants

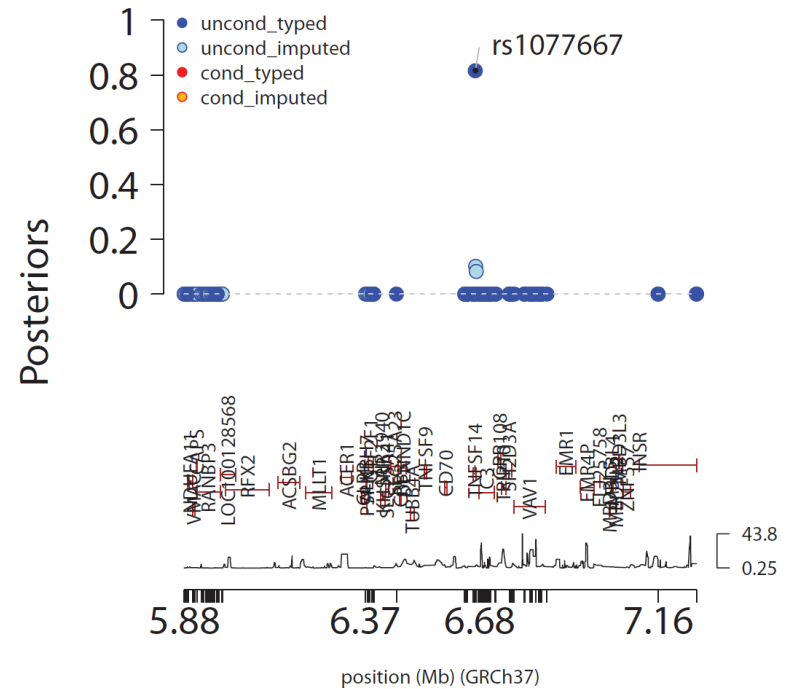


# Occasional success at fine-mapping

meta Chr 11 Gene CD6\_around\_rs34383631 (rs34383631)



meta Chr 19 Gene TNFSF14\_around\_rs1077667 (rs1077667)



# But there is little indication of function....

**Table 3 The 22 variants from the 8 regions with consistent high resolution fine-mapping**

Gene	SNP	Chr	Position <sup>a</sup>	Posterior	GERP	Functional Annotation <sup>b</sup>
<i>TNFSF14</i>	rs1077667	19	6668972	0.81	-3.89	intronic, TFBS / DNase1 peak, correlates with serum levels of TNFSF14
	rs2291668 <sup>c</sup>	19	6669934	0.10	-9.78	intronic / synonymous, TFBS/DNAase1 peak
<i>IL2RA</i>	rs2104286	10	6099045	0.99	-0.47	intronic, correlates with soluble IL-2RA levels
<i>TNFRSF1A</i>	rs1800693	12	6440009	0.70	2.53	intronic, causes splicing defect and truncated soluble TNFRSF1A
	rs4149580 <sup>c</sup>	12	6446990	0.10	2.06	intronic
<i>IL12A</i>	rs1014486	3	159691112	0.79	0.24	-
<i>CD6</i>	rs34383631	11	60793330	0.32	1.66	-
	rs4939490 <sup>c</sup>	11	60793651	0.23	-0.53	-
	rs4939491 <sup>c</sup>	11	60793722	0.23	-0.37	-
	rs4939489	11	60793648	0.16	3.25	-
<i>TNFAIP3</i>	rs632574	6	137959118	0.27	-1.15	-
	rs498549 <sup>c</sup>	6	137984935	0.20	0.52	-
	rs651973	6	137996134	0.17	2.41	downstream of RP11-95M15.1 lincRNA gene
	rs536331	6	137993049	0.15	0.19	upstream of RP11-95M15.1 lincRNA gene
<i>CD58</i>	rs6677309	1	117080166	0.21	-1.18	intronic, TFBS / DNase1 peak
	rs35275493 <sup>c</sup>	1	117095502	0.24	0.75	intronic (insertion)
	rs10754324 <sup>c</sup>	1	117093035	0.22	0.32	intronic
	rs1335532	1	117100957	0.17	-1.32	intronic
<i>STAT4</i>	rs9967792	2	191974435	0.35	-3.96	intronic
	rs10197066 <sup>c</sup>	2	191985459	0.21	0.05	intronic
	rs10804037	2	191991891	0.21	-0.36	intronic
	rs71301540 <sup>c</sup>	2	192001443	0.20	0.08	intronic (deletion)

All listed variants have posterior  $\geq 0.1$  in regions where  $\leq 5$  variants explain the top 50% of the posterior and the top SNP from the frequentist analysis lives in the 90% confidence interval, ordered by maximum posterior.

Posterior denotes the posterior probability of any variant driving association. GERP denotes Genomic Evolutionary Rate Profiling.

<sup>a</sup>Position is based on human genome 19 and dbSNP 137.

<sup>b</sup>Functional data from VEP, eQTL browser, Fairfax et al. (2012), pubmed searches, 1000G. Dash indicates intergenic with no additional annotation.

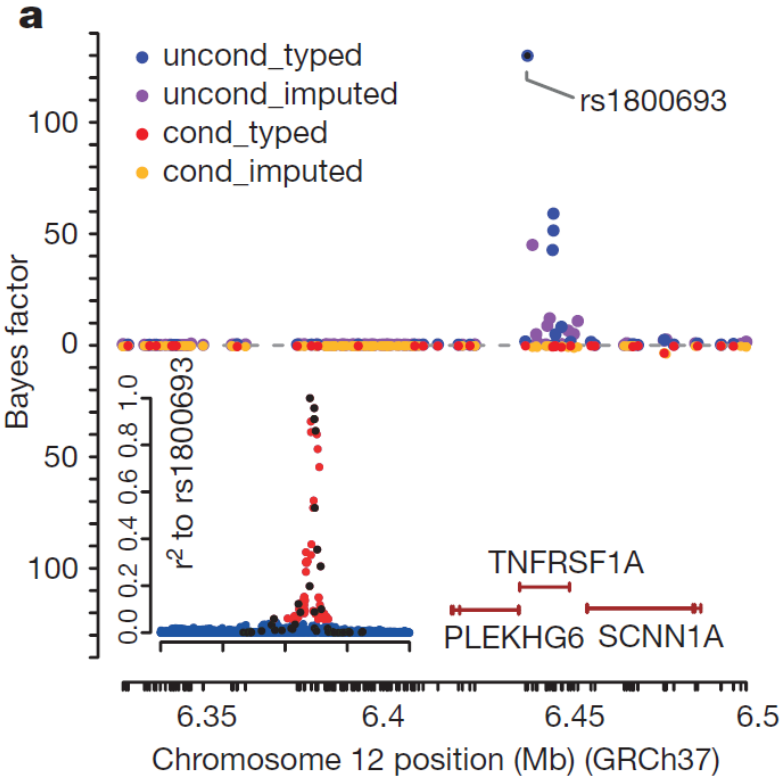
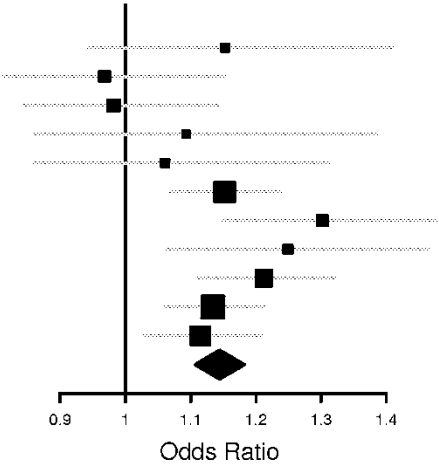
Variants without annotation are intergenic and have no reported regulatory consequence.

<sup>c</sup>Imputed variant.

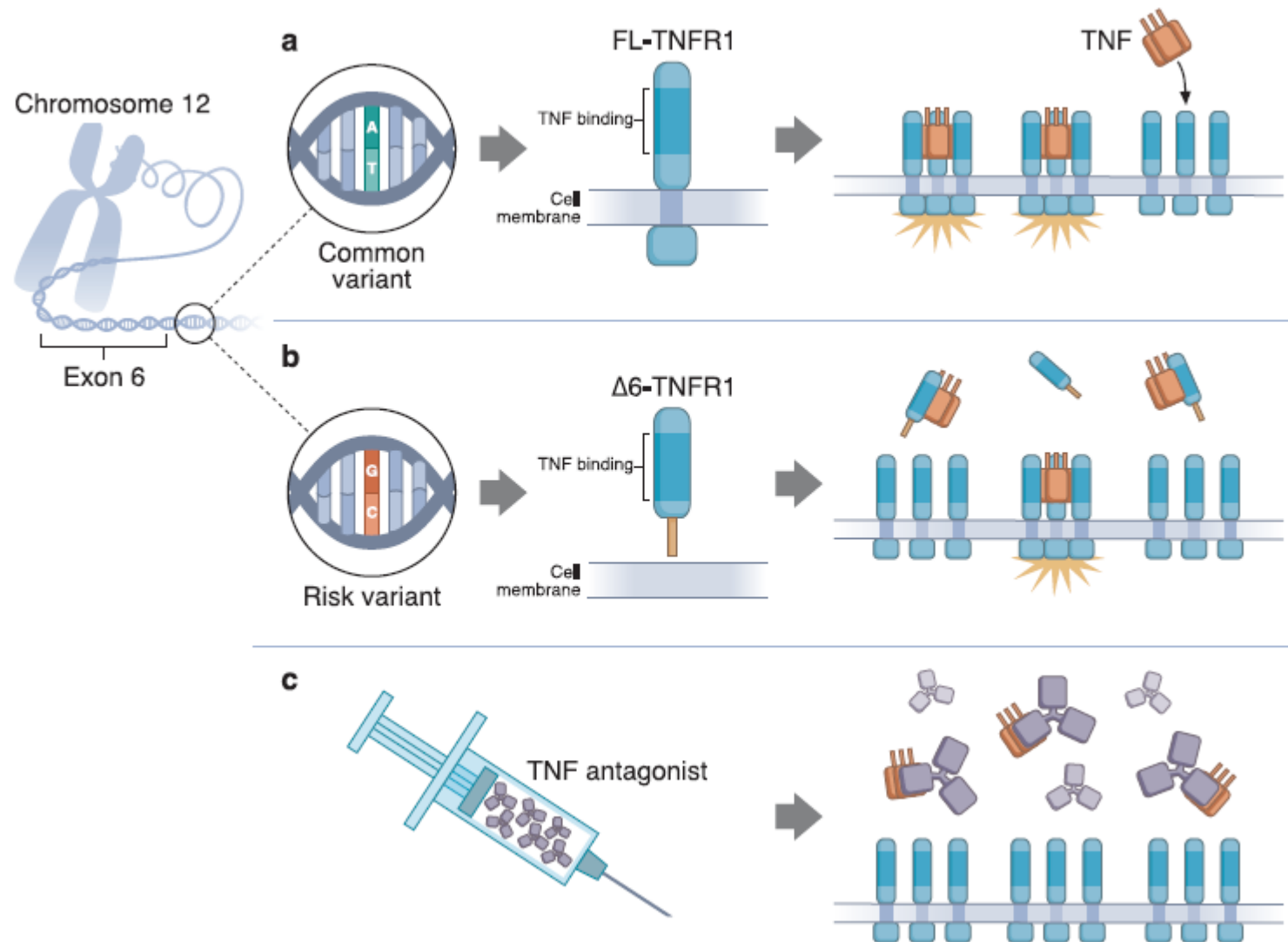
# TNF receptor 1 genetic risk mirrors outcome of anti-TNF therapy in multiple sclerosis

Adam P. Gregory<sup>1\*</sup>, Calliope A. Dendrou<sup>2\*</sup>, Kathrine E. Attfield<sup>2</sup>, Aiden Haghikia<sup>2,3</sup>, Dionysia K. Xifara<sup>4</sup>, Falk Butter<sup>5</sup>, Gereon Poschmann<sup>6</sup>, Gurman Kaur<sup>1</sup>, Lydia Lambert<sup>2</sup>, Oliver A. Leach<sup>2</sup>, Simone Prömel<sup>2</sup>, Divya Punwani<sup>1</sup>, James H. Felce<sup>1</sup>, Simon J. Davis<sup>1</sup>, Ralf Gold<sup>3</sup>, Finn C. Nielsen<sup>7</sup>, Richard M. Siegel<sup>8</sup>, Matthias Mann<sup>5</sup>, John I. Bell<sup>9</sup>, Gil McVean<sup>4</sup> & Lars Fugger<sup>1,2,10</sup>

Stratum	OR	95% CI
AUSNZ	1.15	[0.94–1.41]
Belgium	0.97	[0.81–1.15]
Denmark	0.98	[0.84–1.14]
Finland	1.09	[0.86–1.39]
France	1.06	[0.86–1.31]
Germany	1.15	[1.07–1.24]
Italy	1.30	[1.15–1.48]
Norway	1.25	[1.06–1.46]
Sweden	1.21	[1.11–1.32]
UK	1.13	[1.06–1.21]
US	1.11	[1.03–1.21]
Summary	1.14	[1.11–1.18]



Gregory et al. (2012) Nature.





## 2. Greater ability to validate therapeutic targets



# A case study: Darapladib

- Epidemiological evidence shows that people with lower levels of a particular enzyme (Lp-PLA2) have reduced risk of heart disease.
- GSK developed a drug, darapladib, which inhibits Lp-PLA2
- Common variants around the gene PLA2G7 are modestly associated with Lp-PLA2 levels

...but not with heart disease

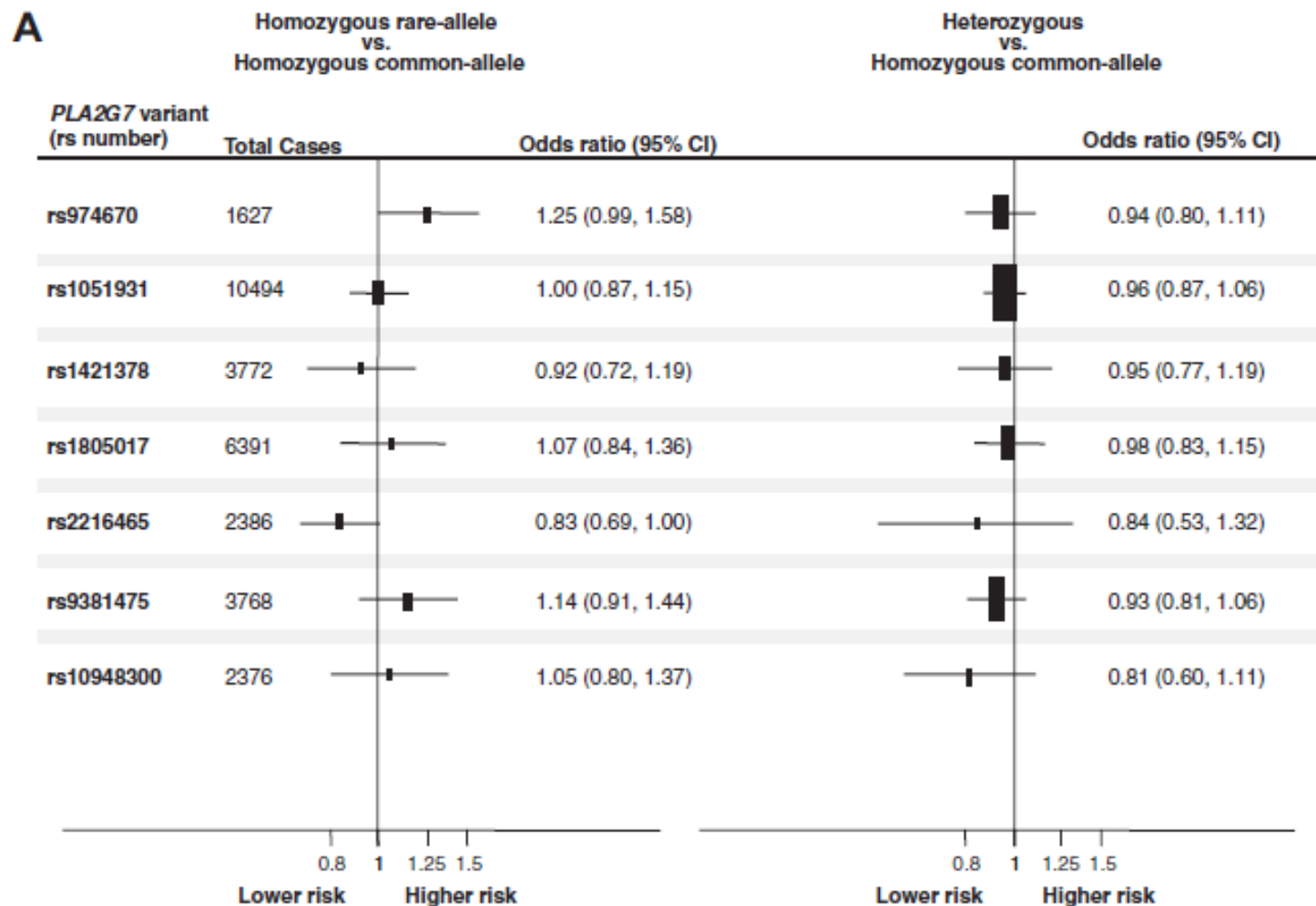
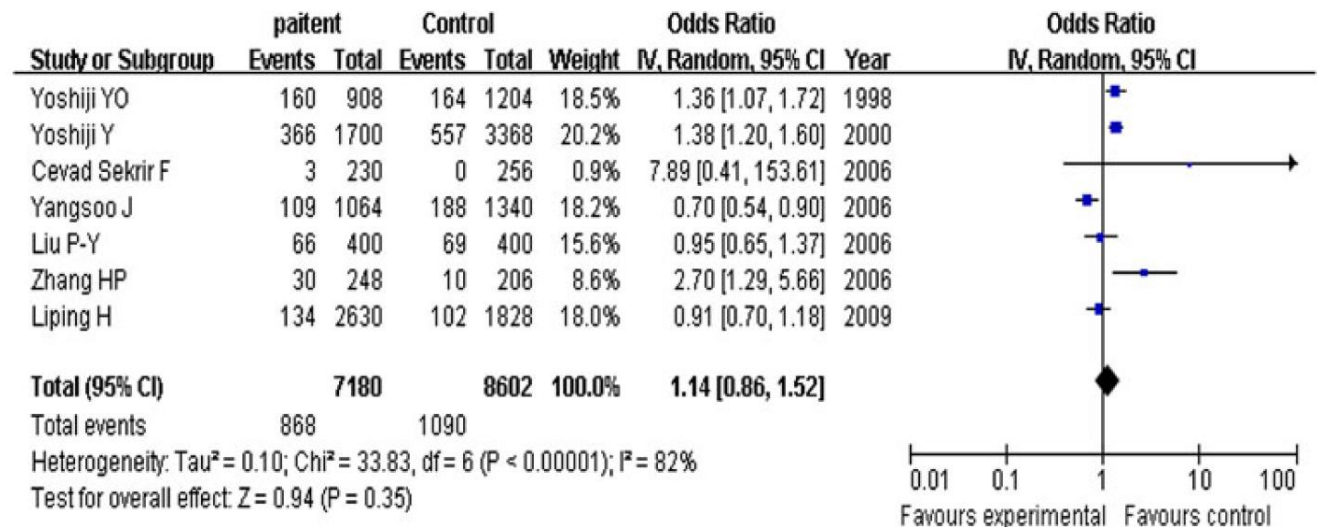


Figure 5. A, Relative odds of CHD associated with *PLA2G7* variants. Data are pooled from up to 10 studies (NPHS-II, EPIC-Norfolk, WH-II, HIFMECH, EAS, AtheroGene, LURIC, Cyprus, SAS, and WTCCC-CHD) including up to 10 494 CHD events. B, Effect of the

# LOF variant in PLA2G7 creates “Lp-PLA2 human knockouts” but is still not associated with risk of CVD

**Fig. 2** Results of association between the V279F polymorphism in PLA2 gene and coronary heart disease under the additive model



# An unhappy ending

- GSK have recently completed “STABILITY” a large clinical trial (c. \$800M) of darapladib.
- On Tuesday, November 12, 2013, GSK announced that the drug had failed to meet Phase III endpoints in a trial of 16,000 patients with acute coronary syndrome. An additional trial of 13,000 patients (SOLID-TIMI 52) is ongoing.

# What does genetics tell us to date about MS treatments?

Compound	Drugs	Mode of action	Relevant genes	Other uses	Comments	eQTLs	GWAS
Teriflunomide	Aubagio	Blocks dihydroorotate dehydrogenase. Inhibits pyrimidine de novo synthesis, hence rapidly dividing cells including activated T cells. Also blocks NF- $\kappa$ B and tyrosine kinases at high dose	DHODH (mutations cause Miller syndrome)	Inactive form (leflunomide) used in severe RA and psoriatic arthritis (also pyrimidine synthesis inhibitor)	Poor efficacy	None	None
Interferon beta-1a	Avonex, Rebif, CinnoVex	Is an interferon type I. Binds to IFN-alpha receptor (IFNAR1/IFNAR2). Cytokine (activate NK cells, macrophages, upregulate antigen presentation). Produced by leukocytes	IFNAR1, IFNAR2 (activate JAK/STAT, Tyk2, etc.). IFN-beta actually 3 products from 3 genes, IFNB1, IFNB3 and IL6 (also called IFNB2). IL6 secreted by CD4 Th cells	Chemotherapy	30% MS patients unresponsive	IFNB1 - no eQTLs reported. IL6 - LPS stimulation specific eQTLs	None
Interferon beta-1b	Betaferon, Extavia	See above	See above	See above	See above	See above	See above
Glatiramer acetate	Copaxone	Random polymer of 4 AA found in MBP	MBP	Dry ARMD (Phase 1)	Doesn't seem to be effective	NA	NA
Fingolimod	Gilenya	Sphingosine 1-phosphate receptor modulator, which sequesters lymphocytes in lymph nodes, preventing them from contributing to an autoimmune reaction	S1PR1 / EDG1	Candidate for heart failure and arrhythmia	Effective treatment	None reported	None
Alemtuzumab	Lemtrada	Monoclonal that binds to CD52 on mature lymphocytes	CD52 (indirect)	Used for CLL	Very serious side effects. Efficacy questioned.	NA	NA
Dimethyl fumarate	Tecfidera	Attached by glutathione, which leads to HO-1 induction (anti-inflammatory). Possible up-regulation of NRF2. HO-1 upregulate IL10 and IL-1R	HMOX1 encode HO-1. IL10, IL1R1, IL1R2	Psoriasis, sarcoidosis, others	Effective treatment	HMOX1 - none reported. IL10 - rs3024490 and rs1554286. None reported for IL1 or receptors.	None. Though IL20 is structurally related to IL10 and is a MS GWAS signal. Note IIs strongly clustered.
Natalizumab	Tysabri	Monoclonal against $\alpha$ 4-integrin	ITGA4 (CD49D)		Low risk of PML caused by reactivation of JC virus	ITGA4 has both eQTLs and multiple exon QTLs reported	None

..at best, confusing

# Big data, big challenges

- Biomedical big data needs many components:
  - High throughput measurement
  - Large cohorts
  - Ease of data access
  - Powerful analysis
  - Appropriate governance
  - Engagement with patients
  - ...
- Medical data has many complexities, but genetics provides a useful instrument to help disentangle causal and indirect associations.
- Much to do to integrate genetic and functional data across many diseases, tissues, cell-types, stimulations, etc.

# With thanks to

- The WGS500 project team
- The IMSGC
- Lars Fugger and his group
- UK Biobank
- Members of the WTCHG