Genome-wide association studies II: Identifying genetic associations with complex traits

Gavin Band gavin.band@well.ox.ac.uk MSc Global Health Science and Epidemiology Genetic Epidemiology Module Wednesday 5th Mar 2025



Learning objectives

Understand a genome-wide association study (GWAS) and the concept of a hypothesisfree approach to studying genetic associations.

Have a working knowledge of the different steps involved in the conduct of GWAS, including study design, quality control and basic analyses.

Be able to interpret and critically appraise evidence from genome-wide association studies.

Understand the relevance of replication, meta-analysis and consortia, and multiancestry approaches, in genome-wide association studies.

Appreciate the use of post-GWAS analyses including fine mapping, gene and pathway analyses, and the concept of causal variants.

Main points in this lecture

How polygenic do traits get, anyway?

Polygenicity – Genetic architecture - Consortia & Meta-analysis – GWAS trajectory



Extracting biological information from GWAS

Pathway analysis – fine mapping - pleiotropy



Actual results from the Wellcome Trust Case-Control Consortium study:





"Common variant, common trait" hypothesis

Are you convinced?

Additive model, N=2,000 cases + 3000 controls per phenotype

Actual results from the Wellcome Trust Case-Control Consortium study:

Number





"Common variant, common trait" hypothesis

Are you convinced?

Maybe we haven't found them all how could we find more?

Additive model, N=2,000 cases + 3000 controls per phenotype

Remember the formula



Genotype frequency

GWAS revolution



Mills & Rahal, "A scientometric review of genome-wide association studies", Communications Biology 2019

NHGRI GWAS Catalog: https://www.ebi.ac.uk/gwas/

GWAS revolution



Mills & Rahal, "A scientometric review of genome-wide association studies", Communications Biology 2019

NHGRI GWAS Catalog: https://www.ebi.ac.uk/gwas/

GWAS revolution



Mills & Rahal, "A scientometric review of genome-wide association studies", Communications Biology 2019

NHGRI GWAS Catalog: https://www.ebi.ac.uk/gwas/

Inflammatory bowel disease (Crohn's disease and ulcerative colitis)





N = 125,992 IBD cases 1.2 million controls

> 600 association signals.

Abstract citation ID: jjad212.0008 OP08

Multi-ancestry genome-wide association study of inflammatory bowel disease identifies 125 novel loci and directly implicates new genes in disease susceptibility

L. Fachal¹, on behalf of the International IBD Genetics Consortium ¹Wellcome Sanger Institute, Human Genetics, Hinxton- Saffron Walden, United Kingdom

Type 2 diabetes



N = 74,000 T2D cases And 824,000 controls



Fig. 5 | **The relationship between effect size and MAF.** Conditional- and joint-analysis effect size (*y* axis) and MAF (*x* axis) for 403 conditionally independent SNPs. Previously reported T2D-associated variants are shown in green, and novel variants are shown in purple. Stars and circles represent the 'strongest regional lead at a locus' and 'lead variants for secondary signals', respectively.

403 signals

"conditionally independent" meaning some of them overlap the same regions

genetics

ARTICLES https://doi.org/10.1038/s41588-018-0241-6

Fine-mapping type 2 diabetes loci to single-variant resolution using high-density imputation and islet-specific epigenome maps

GWAS of human height

In 5.4 million individuals





~90% heritability



? signals

Article

A saturated map of common genetic variants associated with human height

https://doi.org/10.1038/s41586-022-05275-y Received: 19 December 2021 Accepted: 24 August 2022 Common single-nucleotide polymorphisms (SNPs) are predicted to collectively explain 40–50% of phenotypic variation in human height, but identifying the specific variants and associated regions requires huge sample sizes'. Here, using data from a genome-wide association study of 5.4 million individuals of diverse ancestries, we

N = 5.4 million

GWAS of human height

In 5.4 million individuals





~90% heritability



12,111 independent signals

Collectively explaining 50% of heritability

Article

A saturated map of common genetic variants associated with human height

https://doi.org/10.1038/s41586-022-05275-y	Common single-nucleotide polymorphisms (SNPs) are predicted to collectively						
Received: 19 December 2021	explain 40-50% of phenotypic variation in human height, but identifying the specific						
Accepted: 24 August 2022	variants and associated regions requires huge sample sizes ¹ . Here, using data from a						
Published online: 12 October 2022	show that 12.111 independent SNPs that are significantly associated with height						
Open access	account for nearly all of the common SNP-based heritability. These SNPs are clustered						
Check for updates	within 7,209 non-overlapping genomic segments with a mean size of around 90 kb, covering about 21% of the genome. The density of independent associations varies across the genome and the regions of increased density are enriched for biologically relevant genes. In out-of-sample estimation and prediction, the 12,111 SNPs (or all SNPs in the HapMap 3 panel ³) account for 40% (45%) of phenotypic variance in populations of European ancestry but only around 10–20% (14–24%) in populations of other ancestries. Effect size-associated regions and gene notiritization are similar						

They collectively explain ~ 50% of heritability In European ancestry people

N = 5.4 million

Comparing across traits



ARTICLE

DOI: 10.1038/s41467-018-06805-x OPEN

Large-scale genome-wide enrichment analyses identify new trait-associated genes and pathways across 31 human phenotypes

Xiang Zhu ^{1,2} & Matthew Stephens ^{2,3}

With all this data it's possible to fit more sophisticated models that estimate the amount of polygenicity across traits.

Common variant, common disease hypothesis





A complex trait.

Caused by many factors, each having a small overall effect. Including

- Many genetic variants, including common ones
- Environmental factors
- Gene-environment or gene-gene interactions

- ...

How are these studies possible? Consortia and meta-analysis

Consolidation question from last lecture

WTCCC2 GWAS of multiple sclerosis (9,772 cases and 7,376 controls).

	MMEL1(TNFRSF1
124	EVI5
and de circuit	VCAM1
	CD58
12 C	RGS1
	C1orf106(KIF21B)
	No gene
	PLEK
2.4	MERTK
Sec.	SP140
1 A A	EOMES
Man	No gene
	CBLB
	TMEM39A(CD80)
	CD86
	IL12A
	NFKB1(MANBA)
1.28	/ IL7R
22 C	PTGER4
and .	IL12B
	BACH2
-St. :	THEMIS
and the second second	MYB(AHI1)
5 	IL22RA2
	No gene
100	TAGAP
	ZNF746
A.C.	/ IL7
Care -	MYC
Sec. 1	PVT1
	IL2RA
27°.	ZMIZ1
	HHEX
	CD6
· · · · · · · · · · · · · · · · · · ·	CXCR5
0	TNFRSF1A
	CLECL1
	CYP27B1
Bar .	ARL6IP4
	ZFP36L1

For further information about terms used below, how over the red question marks.							
	Region						
dbSNP id.?	rs11581062						

status:² novel association physical position:² 01:101,180,107 association region:² 01:100,983,315-101,455,310 functional tag:² N/A nearest gene:² SLC30A7 candidate gene:² VCAMI*

Signal

 p-value discovery:²
 3.7e-10

 OR discovery (95% CI):³
 1.13 (1.09-1.18)

 p-value replication:²
 4.20e-02 (one-sided)

 OR replication (95% CI):³
 1.07 (0.99-1.15)

 p-value combined:²
 2.50e-10

 OR combined (95% CI):³
 1.12 (1.1-1.13)

 Risk (non-risk) allele:
 G(A)

Allele frequencies?

Country	controls / cases	control / case frequency
Australia	- / 647	- / 0.32
Belgium	- / 544	- / 0.33
Denmark	- / 332	- / 0.32
Finland	2165 / 581	0.23 / 0.24
France	347 / 479	0.31 / 0.34
Germany	1699 / 1100	0.29 / 0.31
Ireland	- / 61	- / 0.34
Italy	571 / 745	0.30 / 0.33
Norway	121 / 953	0.26 / 0.28
Poland	- / 58	- / 0.27
Spain	- / 205	- / 0.36
Sweden	1928 / 685	0.27 / 0.28
UK	5175 / 1854	0.29 / 0.32
USA	5370 / 1382	0.29 / 0.32

Proximal genes[?]



Can you explain?

DPH5, EXTL2, S1PR1, SLC30A7, VCAM1*

Consortia and meta-analysis

To generate such large sample sizes for "common" (but still relatively rare) diseases, requires setting up large multi-centre collaborations. This is fun to be involved in but comes with its own analysis challenges....

Dealing with population structure



This study suffered from a key problem. Can you see what it is?

Dealing with population structure



This study suffered from a key problem. Can you see what it is?



First two "principal components" obtained purely from the genotypes Case-control sampling is correlated with genome-wide genetic variation.

"Confounding by population structure"



Population structure: solutions

Instead of simple 2x2 table

1. Regression including principal components

outcome ~ genotype + *PCs*



Plot of first two principal components obtained from the genome-wide genotypes

Uses just the strongest directions of variation in relatedness (population structure)

2. Linear mixed model

outcome ~ genotype +



Include a **genetic relatedness matrix computed from genome-wide genotypes** in the association test

Uses the entire matrix of relationships



MS study

Most p-values are now not inflated

Anatomy of an association analysis

All GWAS should report data in a way that can be re-used by future studies. This study used several previous GWAS to conduct replication. All the details are given in a supplementary table:

	WAS + replicatio	GWAS	UK only GWAS	on-UK only GWA	combined replication	eneMSA NL replicatio	ene MSA US replicati	eneMSA CH replicatio	ANZ replication	BWH replication
Risk	OR (95%)	log1 OR (Bay (95% esFa	0 0R (95%	OR (95%		OR pval* (95%	0R pval* (95%	OR pval* (95%	OR pval* (95%	OR (95%
🛛 Gene 🔽 lle	▼ pval T C1 ▼	pval 🔍 Cl 🔍 cto	🚺 pval 🔽 Cl 💌	pval 💌 C1 💌	pval* 💌 OR (95% C 💌	*▼ C1▼ inf ▼	*▼ C1▼ inf ▼	* C 1 inf	* C 1 i nf	pval Cl 💌 inf
MMEL1 C	1.00E-14 1.14 (1	3.10E-14 1.16 (1 11.3	9 0.0073 1.12 (7.10E-13 1.17 (1	0.0085 1.08 (1.01-1.15)	0.26 1.1 (0. 0.94	0.18 1.1 (0. 1.01	0.24 1.11 (0 1.03	0.006 1.15 (1 1	0.41 1.02 (0
5 EVI5 A	5.80E-15 1.15 (1	6.50E-12 1.15 (1 9.1	5 2.90E-05 1.2 (1	2.70E-08 1.14 (1	1.00E-04 1.14 (1.06-1.22)	0.088 1.23 (0 1.05	0.59 0.97 (0 0.91	0.71 0.92 (0 0.94	0.023 1.12 (1 0.97	0.0059 1.18 (1 1
SLC30A7 G	2.50E-10 1.12 (1	3.70E-10 1.13 (1 7.4	3 0.00047 1.16 (1.70E-07 1.13 (1	0.042 1.07 (0.99-1.15)	0.57 0.99 (0 1.01	0.095 1.09 (0 0.99	0.013 1.18 (0 0.91	0.57 0.99 (0 1.01	0.095 1.09 (0 0
EXTL2 A	4.00E-08 1.09 (1	3.70E-07 1.1 (1. 4.5	2 0.00096 1.14 (6.00E-05 1.08 (1	0.017 1.08 (1.01-1.15)	0.025 1.11 (1 1	0.088 1.09 (0 1.01	0.45 1.01 (0 0.88	0.025 1.11 (1 1	0.088 1.09 (0 1.



This is a common analysis approach: to gain sample size, use meta-analysis to combine results across several component studies. Then look for consistency between the studies.

$$v_{meta} = 1/\left(\sum_{i} \frac{1}{v_i}\right) \qquad \beta_{meta} = \left(\sum_{i} \frac{\beta_i}{v_i}\right) \times v_{meta} \qquad (Where v denotes squared standard error)$$

"Inverse variance weighted fixed-effect meta-analysis", gives results approximately equal to joint analysis of genotype data.

We now have thousands of GWAS signals across thousands of traits. What do they teach us about the underlying biology?





 \checkmark

DNA gets physically packaged up into chromosomes...





 $\sqrt{}$

DNA gets physically packaged up into chromosomes...





...inside cells, where it is transcribed to form proteins and other molecules...





 $\sqrt{}$

DNA gets physically packaged up into chromosomes...



...inside cells, where it is

transcribed to form proteins and other molecules...





...that combine to make individuals...

...that affect how the cells behave, forming different organs...

1



 $\sqrt{}$

DNA gets physically packaged up into chromosomes...



...inside cells, where it is transcribed to form proteins and other molecules...





...whose success is affected by the traits they have...



...that combine to make individuals...

...that affect how the cells behave, forming different organs...



 $\sqrt{}$

...that gets physically packaged up into chromosomes...



...passing on DNA, with mutations and recombination, to new generations...



There is complex biology at all stages



...whose success is affected by the traits they have...



...that combine to make individuals...

...that affect how the cells behave, forming different organs...

7

...inside cells, where it is transcribed to form proteins and other molecules...



...passing on DNA, with mutations and recombination, to new generations...



...that gets physically packaged up into chromosomes...



microarrays. genome sequencing

There is complex biology at all stages

And we can measure it.

Biomarker measurements

Л

Clinical phenotype

measurements





...that combine to make individuals...

...that affect how the cells behave, forming different organs...

RNA-seq, spectroscopy, antibody binding



...inside cells, where it is transcribed to form proteins and other molecules...

Chromatin state marker assays, ChIP-seg, ...

Gaining biological knowledge from GWAS

There are several ways we can try to translate knowledge of associations into new biological insights. I will try to describe a few of these.

- Fine-mapping can we identify the actual causal variants underlying these associations, and hence discover specific proteins and disease pathways?
- Pathway analysis even if we can't fine-map, we can still try to assess whether associations group into particular biological pathways that might shed light on biology
- Pleiotropy how are associations shared between traits?



 \checkmark





Evenue 1. e nothway analysis

Example 1: a pathway analysis





...that combine to make individuals...

7

Pathway analysis

Pathway analyses and gene enrichment analysis seek to determine whether there is a statistical tendency for association signals to fall into known groups of related genes. These can be

- Known biological pathways (functional networks of proteins and molecules, performing known specific biological functions) – such as those available from the KEGG and Reactome databases
- More general classifications of genes by function, such as those from the Gene Ontology Project

A slightly different direction is to try to group signals by genome function – for example, do they lie in exons? Or gene promoters? Or in regulatory regions active in particular cells?

https://www.genome.jp/kegg/

https://reactome.org

http://geneontology.org

Pathway analysis example

The primary cause of MS has typically been thought to be inflammation causing downstream neurodegeneration – with some debate about this. Can the GWAS of MS we discussed shed light on this?

100 PT 10 PT	-rs4648356	C	0	1 1	—	MMEL1(TNFRSF14)	7	•	
	>rs11810217	A	: 0	1	-	EVI5	15		
1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	rs11581062*	G	0	1 1		VCAM1	5	1	
	rs1335532	A		0		CD58	2		
Barn and the	-rs1323292	A	0			RGS1	1		
and a second sec	rs7522462	G		1 1	-	C1orf106(KIE21B)	4		
and the second s	ro10466000	~	Ľ	1 1		No gono	*	•	
	1512400022	~		1 1				2.1	
*	15/595037	A	P	i i		PLER	4	•	•
and the state	rs1/1/48/0	G	0	1 1		MERIK	7	•	•••
and the second sec	rs102018/2	A		1 1		SP140	3	•	
the second s	rs11129295	A	•	1 1	-	EOMES	1	•	•
Protest	rs669607	C	•	1		No gene	0		
	rs2028597	G	+++-	1		CBLB	1	•	•
A A A A A A A A A A A A A A A A A A A	rs2293370	G	0	1 1		TMEM39A(CD80)	7		
	rs9282641*	G	1	0		CD86	5	•	•
	rs2243123	G	0	1	-	IL12A	2		•
	rs228614	G	0	1 1		NFKB1(MANBA)	8		
1	rs6897932	G	0	1 1		IL7R	7		
and the second se	-rs4613763	G			-	PTGFR4	1		
	rs2546890	Δ		T i		II 12B			
The second s	re12212193	G		1 1	_	BACH2	4		
· ·	1012212100	~	1			THEMIS	6	۰.	-
TRAL A	15002734	21	ľ.	1 1			0		
	1311134001	~	•				3		
	rs17066096	G	•	1		IL22RA2	3	•	
	rs13192841	A	?	1 1		No gene	0		
	rs1738074	G	0			TAGAP	2	•	
and a second sec	rs354033	G	•	1 1		ZNF746	4		•
And the second s	//rs1520333	G	•		-	IL7	3	3	•
224-	rs4410871	G	0	1		MYC	2		
	rs2019960	G	0	1 1	4	PVT1	1	•	
	rs3118470*	G	0	1 1		IL2RA	4	•	•
	/rs1250550	A	٩	1 1		ZMIZ1	3	•	
	rs7923837	G	۵	1 1	-	HHEX	3	•	•
Ales	//rs650258	G	0	1	-	CD6	4		
	rs630923	C	0	1 1		CXCR5	18	•	•
The state of the s	///rs1800693	G	0	1 1		TNFRSF1A	4		
	rs10466829	A	ó	1 1		CLECL1	9		
Martin *	urs12368653	A	6			CYP27B1	33	1	
	rs949143	G	Later 1	1 1		ABI 6IP4	13		
	rs4902647	G	6	1 1		ZEP36I 1	3		
	re2300603	Δ		1 1	_	BATE	2		
	re2119704	0				GALC(GPR65)	2	÷.	
11, 1000 000	132113104	0		1		COVO			
and a second sec	152744140	G		i i		OLEO ICA(OUTA)	4		
and a	157200780	^		1 1		(DEC TOA(CITTA)	8		-
	IIIIII 13333054	A	0	1 1		IRF8	1	• •	•
There Train to the	///rs9891119	C	•	1 1		STAT3	25	•	
and the second sec	rs180515	G	0		-	RPS6KB1	9	•	
A fin and date of a state of the	rs7238078	A	•	1		MALT1	2	•	•
	rs1077667	G	0	1 1		TNFSF14	3	•	•
A	///rs8112449	G	0	1 1		TYK2(ICAM3)	12		
A Starty M. C.	rs874628	Α	Ð		-	MPV17L2(IL12RB1)	11	•	
	rs2303759	С	Þ			DKKL1(CD37)	9	•	
and all it.	//rs2425752	A	0		-	CD40	13	,	•
	rs2248359	G	0	1		CYP24A1	2	•	
	rs6062314	A	10	1		TNFRSF6B	15		
	rs2283792	C				MAPK1	9		
	rs140522	A		1	-	SC02	15	-	1010
	13140022	· · · ·				0001	10		



As the main figure shows, many of the association signals looked like they were near immune-system related genes.

www.well.ox.ac.uk/wtccc2/ms/

Pathway analysis example

We:

- Assigned SNPs to their nearest gene using the available annotation
- Used the Gene Ontology Project to classify genes into functionally related groups
- Conducted a statistical test (Fisher's exact test) to identify whether the nearest genes were enriched in each group.



T-helper-cell differentiation pathway (from Ingenuity Pathway Analysis software) Particularly strong enrichment was observed for immune system pathways – notably in "T cell activation and proliferation" (P=1.9x10⁻⁹)

"Although GO immune system genes only account for 7% of human genes, in 30% of our association regions the nearest gene to the lead SNP is an immune system gene"

Published: 10 August 2011

Genetic risk and a primary role for cell-mediated immune mechanisms in multiple sclerosis

The International Multiple Sclerosis Genetics Consortium & The Wellcome Trust Case Control Consortium

Fine-mapping

"Fine-mapping" is the general term used for attempts to narrow down association signals to the underlying causal variants. A typical process involves:

- Gathering complete information on genetic variation in the region of interest

 for example by deep-sequencing a large number of individuals. (Large databases such as gnomAD / TopMed now make this easier.)
- Gathering information on genome function including gene structure and regulatory regions.
- Potentially leveraging data from different ancestral backgrounds, hoping that differences in LD patterns will help narrow down signals.
- Fitting models that attempt to parse apart multiple associations in the same region

Possible underlying mechanisms are pretty diverse and a healthy dose of genomic detective work is often needed.



V





Fine-mapping example 1 Complex genetic variation





...that combine to make individuals...

7


GWAS of susceptibility to severe malaria

Study sample	S				cm~2	Whole-genom	e sec	quenc	es	
Group	Cases	Controls	TOTAL	A brown	هم چ	Group	Trios	Duos	Other	TOTAL
Africa					N.	Gambia				
Gambia	2567	2605	5172		{}	FULA	31	1	5	100
Mali	274	183	457		, ∠	JOLA	32	1	2	100
Burkina Faso	733	596	1329	L'IL & Any on	6.	MANDINKA	33	0	1	100
Ghana	399	320	719	3 Contrad		WOLLOF	32	1	3	98
Nigeria	113	22	135			Burkina Faso				
Cameroon	592	685	1277			MOSSI	0	0	57	57
Malawi	1182	1317	2499) - The port		Cameroon				
📕 Tanzania	416	403	819			BANTU	5	3	11	31
📕 Kenya	1681	1615	3296	1 JUL		SEMIBANTU	8	0	7	32
Asia				land de		👂 Tanzania				
Vietnam	718	546	1264	$\nabla \circ /$	19	CHAGGA	21	2	13	80
Oceania						PARE	22	2	7	77
PNG	402	374	776			WASAAMBA	23	6	9	90

GWAS in 17,000 severe malaria cases and population controls From 12 sites in Africa, Oceania, and SE Asia. Genotyped on the Illumina Omni 2.5M array + whole-genome sequences for imputation

Malaria Genomic Epidemiology Network. "Insights into malaria susceptibility using genome-wide data on 17,000 individuals from Africa, Asia and Oceania". Nature Communications (2019). <u>https://doi.org/10.1038/s41467-019-13480-z</u>



www.malariagen.net

Natural resistance is driven by red blood cell variation



Natural resistance is driven by red blood cell variation



SNPs on chromosome 4 are associated with proection against severe malaria



Chromosome 4

The association has quite large effect



> 30% protective effect per copy of the derived allele

Standard error(log OR) \approx

 $\overline{N} \times f(1-f) \times \phi(\overline{1-\phi})$

We had an exciting association. But fine-mapping has proven to be difficult for many GWAS loci.

To hope for success we might need:

- Good candidates for the functional gene?

SNPs on chromosome 4 are associated with proection against severe malaria



Glycophorins encode the 'MNS' blood group (antigenic molecules on RBC surface)

Glycophorins



Grimes and Slater, The Inherited Metabolic Diseases, 1994

Glycophorins are receptors for *P.falciparum* during red blood cell invasion

P. Falciparum parasite



We had an exciting association. But fine-mapping has proven to be difficult for many GWAS loci.

To hope for success we might need:

Good candidates for the functional gene?

We had an exciting association. But fine-mapping has proven to be difficult for many GWAS loci.

To hope for success we might need:

Good candidates for the functional gene?

We had an exciting association. But fine-mapping has proven to be difficult for many GWAS loci.

To hope for success we might need:

Good candidates for the functional gene?

Structural variants create deletions, duplications, and hybrid genes

The MNS blood group is highly diverse, with over 45 known antigens.

Encoded by single nucleotide polymorphisms and structural variants



Deleted / duplicated / hybrid genes

We had an exciting association. But fine-mapping has proven to be difficult for many GWAS loci.

To hope for success we might need:

- Good candidates for the functional gene?
- Good candidates for the causal mutation(s)?

Steps to fine-map

Step 1: type or sequence as much of the genetic variation in the region as possible – hope to catch the causal mutation.

Step 2: re-analyse the association.

Step 3: look for functional mutations

A regional reference panel capturing structural variation

We used the >3,600 samples including

- 1000 Genomes Project Phase III reference panel
- plus our newly-sequenced samples



...to call SNPs and indels <u>and</u> structural variation.

Illustration of structural variant calling:



(this sample has a deletion in this region)

A regional reference panel capturing structural variation

We used the >3,600 samples including

- 1000 Genomes Project Phase III reference panel
- plus our newly-sequenced samples



...to call SNPs and indels <u>and</u> structural variation.

Illustration of structural variant calling:



(this sample has a deletion in this region)

...our method infers the copy number

The region turned out to have a lot of structural variation

Deletions

Duplications



14% of Africans carry a CNV affecting these genes

The region turned out to have a lot of structural variation

Deletions

Duplications



14% of Africans carry a CNV affecting these genes

Before fine-mapping



Original GWAS result

After fine-mapping



Result after incorporating genetic variation discovered in sequenced samples.

After fine-mapping





This is how a microarray cluster plot should look: 3 clusters for AA / AB / BB genotypes

Actually this signal was evident in our cluster plots







What we saw in this region

Still true that nothing seemed to be functional. What next? Protective: relative risk ~ 0.6



This is how a microarray cluster plot should look: 3 clusters for AA / AB / BB genotypes



What we saw in this region



We were able to use cluster plots to confirm individuals in our GWAS really do carry the complicated structural variant "DUP4".

DUP4 is pretty complicated – what could it be?

What is DUP4?





Leffler et al, "Resistance to malaria through structural variation of red blood cell invasion receptors", Science (2017) <u>https://doi.org/10.1126/science.aam6393</u>

What is DUP4?



Leffler et al, "Resistance to malaria through structural variation of red blood cell invasion receptors", Science (2017) <u>https://doi.org/10.1126/science.aam6393</u>



Article Red blood cell tension protects against severe malaria in the Dantu blood group

https://doi.org/10.1038/s41586-020-2726-6 Received: 20 November 2018 Accepted: 19 June 2020 Published online: 16 September 2020

 Silvia N. Karluki¹⁰⁰, Alejandro Marin-Menendez³⁰⁰, Viola Introini³⁰⁰, Benjamin J. Ravenhill⁴, Yen-Chun Lin³, Alex Macharia¹, Johnstone Makale¹, Metrine Tendwa¹, Wilfred Nyamu¹, Jurij Kotar³, Manuela Carrasquilla³, J. Alexandra Rowe⁸, Kirk Rockett⁴, Dominic Kwiat/Kowski^{24,7}, Michael P. Weekes⁴, Pietro Cicuta³¹⁰², Thomas N. Williams^{(MABER} & Julian C. Rayne^{-2,ABER}

https://doi.org/10.1038/s41586-020-2726-6

Dantu is globally rare...

The Dantu blood group has been found in:

1 in 44,112	Londoners*
0 in 1,000	Germans [†]
1 in 320	African Americans [†]
0 in 2870	Gambians [‡]

...but found at high frequency in east Africa

The Dantu blood group has been found in:

1 in 44,112	Londoners [*]
0 in 1,000	Germans [†]
1 in 320	African Americans [†]
0 in 2870	Gambians [‡]
1 in 12	Malawians [‡]
1 in 6	Kenyans (from the Kilifi region)

Allele frequency:



The circle of genetic causation



 \checkmark





Fine-mapping example 2 Cell-specific gene regulation





...that combine to make individuals...

7



Natural resistance is driven by red blood cell variation



Association near 2nd exon of *ATP2B4*



The associated SNPs cover a region around the second exon. None of these SNPs make changes to the protein. What could be going on? "Canonical" gene model for ATP2B4

ATP2B4 = a red cell "calcium pump"

Cartoon of a gene



Introns (these get spliced out of the RNA)

Cartoon of a gene



If the DNA is accessible here, transcription factors will bind and help to 'turn on' transcription
Two ways to look at transcription









ATP2B4 is widely expressed...





Data from ENCODE / Roadmap

Malaria-associated region

...but shows chromatin differences in RBCs



2nd exon



Data from Xu et al Dev Cell (2012)

ATP2B4 is widely expressed...

Measured RNA transcription (RNA-seq)

1st exon

2nd exon

GENCODE v19 transcripts	direction of transcription				
ESC1		-			the star inter-
ES-deriv ¹					
Epithelial ¹					
HSC & B-cell1				<u> </u>	bi i luiti (())
Blood & T-cell ¹		distant.			Mhaili i saini hi
Neurosph ¹		all states	a distant and the distant of the	LE L 3	Martin Lake
Heart'				1112	أستأ أكتدانك
Other ¹					and a strength of the
Brain ¹		then:		مريبه من السينية المريد ال	and the first state of the stat
Digestive '				11.1.1	in all a line
Muscle ¹					
Thymus ¹					
ENCODE2012 (except K562)1					
K5621			The second se		Antik

Non-erythroid cells (i.e. no red blood cells)

ATP2B4 has an erythroid-specific transcript

Measured RNA transcription (RNA-seq)

1st exon 2nd exon GENCODE v19 transcripts direction of transcription ESC1 111111 ES-deriv¹ 1.111 Epithelial¹ 1 111 | 1 Ш HSC & B-cell¹ 1 111-1 11 Blood & T-cell¹ 1 I tool of 11 Neurosph¹ 1.111.1.1 11 1 1 1 1 1 1 11.1 Heart¹ - Children 1 11 11 1 1 1 1 1 Other1 Brain¹ 11 Diaestive¹ 1.11.1.1 1.1

Muscle¹ Thvmus¹ ENCODE2012 (except K562)1 11 1 K5621 11.1 proervthroblast² 1.111.1.1 111111 11.1.1 header and a early basophilic² 11.1 late basophilic² 1 11 1 1 1 11 والأشباط أترهيا orthochromatic² polychromatic² Bone marrow erythroblast3 Fetal liver erythroblast³

Erythroid cells show a different expression pattern.

Red cells do not have nuclei, so to capture mRNA expression in red cells, these studies experimentally differentiated stem cells into the erythroid lineage, and measured transcription before enucleation.

niiiniiiiiiiiiiiii Nationalii

Martin Lake

الالاست فتشدانة

ATP2B4 has an erythroid-specific transcript

Measured RNA transcription (RNA-seq)

	1 st exon	2 nd exon	
ESC' ES-deriv' Epithelial' HSC & B-cell' Blood & T-cell' Neurosph' Heart' Other' Brain' Digestive' Muscle' Thymus' ENCODE2012 (except K562)' K562'			
proerythroblast ² early basophilic ² late basophilic ² orthochromatic ² polychromatic ² Bone marrow erythroblast ³ Fetal liver erythroblast ³			
FANTOM5 transcripts			
GWAS posterior (SM)			GWAS SNPs

Putting together data from a variety of sources suggests the existence of an *alternative transcription start site* near the GWAS signal, but only active in erythrocytes. How can this be?



The transcription of genes in red blood cells is controlled by a particular set of transcription factors – a key one is GATA1.

GATA1 is named after the DNA motif it recognises:



v1.factorbook.org

GATA1 binds just upstream of 2nd exon

Measured GATA1 binding

	1 st exon	2 nd exon		
ESC ¹ ES-deriv ¹ Epithelial ¹ HSC & B-cell ¹ Blood & T-cell ¹ Neurosph ¹ Heart ¹ Other ¹ Brain ¹ Digestive ¹ Muscle ¹ Thymus ¹ ENCODE2012 (except K562) ¹ K562 ¹				
proerythroblast ² early basophilic ² late basophilic ² orthochromatic ² polychromatic ² Bone marrow erythroblast ³ Fetal liver erythroblast ³				bei eller a bela bei eller a bela
FANTOM5 transcripts				
GWAS posterior (SM)			GWAS SNPs	
GATA1 peaks		_	•	

ChIP-seq experiments show GATA1 binds just upstream of our new exon. Moreover, one of the associated SNPs disrupts the GATA1 motif.

One of the malaria-associated SNPs disrupts the GATA site



Leads to a prediction:

- The risk allele *creates* GATA motif and is associated with increased *ATP2B4* expression in RBCs.
- The protective allele removes the GATA motif and the gene is not expressed.

Does this really hold up?

Leads to a prediction:

- The risk allele *creates* GATA motif and is associated with increased *ATP2B4* expression in RBCs.
- The protective allele removes the GATA motif and the gene is not expressed.

N = 24 experimentallydifferentiatederythrocyte precursorcells

Erythrocyte-specific calcium control at ATP2B4



Erythrocyte-specific calcium control at ATP2B4



Learning biology from GWAS - summary

Have highlighted two of the complexities that could occur when trying to fine-map genetic association signals.

They are pretty fascinating and luckily there is lots more of tis type of thing to find!

Anything that can happen, does happen. ...and there is lots of data!

Learning biology from GWAS - summary Long-distance interactions in the genome Non-coding variants Changes to gene expression Polygenic effects (lots of variants involved) Cell-type / tissue heterogeneity **Pleiotropy** (a variant affects lots of phenotypes at once) Genetic interactions Host-pathogen interactions Repetitive DNA / repeat expansions Genome structural variation Genome evolution

Anything that can happen, it does happen.

Prospective cohort studies

A new crop of studies aims to create a database of deep genotype, phenotype, and exposure data across large cohorts of individuals sampled from the population or from health services. Examples:



Precision Medicine Initiative, All of Us (US)





CartaGene (Canada)



FinnGen (Finland)





China Kadoorie Biobank



The 100,000 genomes project (UK)



To all residents,

An opportunity to take part in research and learn new information about your blood pressure and future risk of disease.

You are invited to take part in Our Future Health, the UK's largest ever health research programme. If you take part, you will have the chance to find out more about your health now, and your risk of developing some diseases in the future.

Today, too many people spend many years of their life in poor health. Our Future Health aims to help prevent, detect and treat diseases earlier. Diseases like dementia, cancer, diabetes, heart disease and stroke.

Our Future Health needs up to five million people. Everyone aged 18 and over living in the UK is eligible to take part.

Taking part includes answering some online questions about yourself, providing a blood sample, and having your blood pressure measured at a local clinic.

In the future you will have the option to receive information on your risk of some diseases including diabetes, heart disease and some cancers. This will be calculated using the information you provide and analysis of the DNA in your blood sample.



Scan this QR code for more info and to sign up

Or visit ourfuturehealth.org.uk/join/0518

£10 voucher

Sign up using the QR code or website link above and you will be eligible for a £10 voucher to recognise the time and effort of volunteering. You can find more information on the back of this letter.

You can also share this invitation with other members of your household. If you have any questions, please call 0808 501 5634 or email support@ourfuturehealth.org.uk

Yours sincerely,

Raghib Ali

Raghib Ali OBE MD FRCP(UK) Chief Medical Officer, Our Future Health NHS Consultant in Acute Medicine

She kur

Professor Sir John Bell GBE, FRS Chairman, Our Future Health

.....



https://ourfuturehealth.org.uk

Recruiting now

Learning objectives

Understand a genome-wide association study (GWAS) and the concept of a hypothesisfree approach to studying genetic associations.

Have a working knowledge of the different steps involved in the conduct of GWAS, including study design, quality control and basic analyses.

Be able to interpret and critically appraise evidence from genome-wide association studies.

Understand the relevance of replication, meta-analysis and consortia, and multiancestry approaches, in genome-wide association studies.

Appreciate the use of post-GWAS analyses including fine mapping, gene and pathway analyses, and the concept of causal variants.

Conclusions and summary

- Most human traits are highly heritable
- For 'complex' traits, the effects are made up of many genetic variants often with modest effects - polygenicity
- GWAS study designs can find these variants. They rely on large samples and dense genotyping, and patterns of linkage disequilbirum to detect signals.
- A major frontier is to understand the biology and translate these findings into clinically useful insights and predictions.

(We need people like you to do this.)











V

Thanks for listening!



