# Understanding the genetics of complex traits I

Gavin Band gavin.band@well.ox.ac.uk

**BA Human Sciences** 

Thursday 6th Mar 2025



#### Main lecture messages

1. Most human phenotypes are highly heritable (a large proportion of variation is due to genetics)

2. But many 'complex' traits are *not* mendelian - they are polygenic

 The discovery of this fact is due to *genome-wide association studies* (GWAS), the first of which was conducted in the mid 2000s.

We will go into this in some detail – methodology, population genetics, GWAS in practice

#### 4. Biology is hard

The human genome is ~3.2 billion base pairs long. About 1 in 100 – 1000 of those bases vary between people.



# What proportion of phenotypic variation is due to genetic variation?

#### Human traits are highly heritable

We don't have to guess!

Idea: if genetics determines a trait, then *more genetically similar individuals should have more similar phenotypes.* 

We can estimate how much genetics determines trait variation by comparing trait similarity in more genetically similar and less genetically similar individuals, such as monozygotic and dizygotic twins.

Meta-analysis of the heritability of human traits based on fifty years of twin studies

Tinca J C Polderman<sup>1,10</sup>, Beben Benyamin<sup>2,10</sup>, Christiaan A de Leeuw<sup>1,3</sup>, Patrick F Sullivan<sup>4–6</sup>, Arjen van Bochoven<sup>7</sup>, Peter M Visscher<sup>2,8,11</sup> & Danielle Posthuma<sup>1,9,11</sup>

(2015)

Large meta-analysis of > 2000 twin studies (Browse the results at: <u>https://match.ctglab.nl</u>)

#### Human traits are highly heritable

Idea: if genetics determines a trait, then *more genetically similar individuals should have more similar phenotypes.* 



All studied traits Compare trait correlations between twins.

(Adult) height is *much* more similar between monozygotic than dizygotic twins. <u>The heritability</u> is about 90%.

*Heritability* is the proportion of trait variation explained by inherited factors (including genetics). Can be estimated as  $h^2 \approx 2 \times (r_{MZ} - r_{DZ})$ .

## Human traits are highly heritable

If genetics determines a trait, then *more* genetically similar individuals should have more similar phenotypes.

## Meta-analysis of the heritability of human traits based on fifty years of twin studies

(2015)

Tinca J C Polderman<sup>1,10</sup>, Beben Benyamin<sup>2,10</sup>, Christiaan A de Leeuw<sup>1,3</sup>, Patrick F Sullivan<sup>4–6</sup>, Arjen van Bochoven<sup>7</sup>, Peter M Visscher<sup>2,8,11</sup> & Danielle Posthuma<sup>1,9,11</sup>

Dizygotic Monozygotic 1.0 Age 18-64 years 0.8 0.6 0.4 0.2 n -0.2 Ment. Blood Endocr. General High-L Imm. Ment. Other Spec. Structure Temp. Weight Depr. Conduct Heart Structure Funct. of Hyperkingland metab. Height cognitive beh. dis beh. dis of the pressure Food system anxiety personal pers. maint. dis. episode brain funct. etic dis. of mouth funct. funct. funct. funct. funct. alc. tob. dis. dis. eyeball funct. funct. TM7 0.59 0.67 0.39 0.53 0.42 0.65 0.65 0.52 0.92 0.68 0.58 0.56 0.55 0.69 0.41 0.41 0.68 0.89 0.42 0.76  $\Box r_{DZ}$ 0.29 0.20 0.26 0.22 0.52 0.34 0.18 0.37 0.19 0.36 0.24 0.47 0.28 0.30 0.30 0.41 0.17 0.33 0.21 0.34 Structure Adult height Blood pressure "Higher Depression of the  $h^2 \approx 90\%$  $h^2 \approx 60\%$ level  $h^2 \approx 42\%$ eveball cognitive  $h^2 \approx 70\%$ function"  $h^2 \approx 80\%$ 

Lots of theoretical caveats might apply here – see Lecture 1. But in general it is true that a large proportion of variation in most human phenotypes is caused by genetics.

# Two possible extreme genetic architectures



#### Example: Huntingdon's

Cell, Vol. 72, 971-983, March 26, 1993, Copyright © 1993 by Cell Press

#### A Novel Gene Containing a Trinucleotide Repeat That Is Expanded and Unstable on Huntington's Disease Chromosomes

The Huntington's Disease Collaborative Research Group\*

Introduction

Affects ~1 in 20,000 people of European ancestry (less in Africa and Asia)

#### Discovered by looking in families



#### A "Mendelian" trait

# Two possible extreme genetic architectures



#### Example: Huntingdon's

Cell, Vol. 72, 971-983, March 26, 1993, Copyright © 1993 by Cell Press

#### A Novel Gene Containing a Trinucleotide Repeat That Is Expanded and Unstable on Huntington's Disease Chromosomes

The Huntington's Disease Collaborative Research Group\*

Introduction

Affects ~1 in 20,000 people of European ancestry (less in Africa and Asia)

#### Discovered by looking in families



A "Mendelian" trait

# End of an era



"Linkage Mapping was successful in identifying the genetic basis of many human diseases in which the disease penetrance resembles a simple Mendelian model e.g. Huntington's disease, Cystic Fibrosis, some forms of breast cancer, Alzheimers, ..."

"...but the literature is now replete with linkage screens for an array of common 'complex' disorders such as schizophrenia, manic depression, autism, asthma, type I and type II diabetes, Multiple Sclerosis, Lupus. Although many of these studies have reported significant linkage findings, none has lead to convincing replication"

- Risch "Searching for genetic determinants in the new millennium" Nature (2000)

# Common variant, common disease hypothesis



# Common variant, common disease hypothesis





# 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

# Summary

- Most human phenotypes are highly heritable a large proportion of phenotype variation seems to be caused by genetics. ~60% on average!
- In principle this heritability could occur in different ways for example through single variants with strong effects, or through multiple variants with small effects.
- By the 2000s family studies had identified the causes of several mendelian traits, but had failed to solve the genetics of multiple complex diseases.

Was the "common variant, common disease" hypothesis true?

# End of the linkage era



# The birth of GWAS



## **GWAS** roadmap

1. Collect as many samples as possible

How many samples?

2. Genotype the at as many variants across the genome as possible

How many variants? Which ones?

An ad hoc but widely used threshold

3. Run a statistical test for genotype-phenotype association How to test? Can we deal with confounders?

To produce this:



Lots of statistical tests so to get excited we need strong evidence e.g.  $P < 5 \times 10^{-8}$ 

## GWAS roadmap

- → Testing for association
  - Confounding and the importance of quality control
  - What variants to genotype, and how? LD and the HapMap study
  - A real GWAS study WTCCC

## Testing for association

#### Imagine a genetic variant that affects risk of disease



## Testing for association



If genotype *G* causes disease, then carrying *G* will make you more likely to have disease.



## Testing for association



If genotype G causes disease, then carrying G will make you more likely to have disease.

$$\underline{Relative \ risk} = \frac{P(\text{disease}|G)}{P(\text{disease}|g)} > 1 \qquad \qquad \text{Using probability} \\ \underbrace{\text{notation}}$$

If the genotype causes disease, then the relative risk will be different from 1

How to estimate relative risk?

$$RR = \frac{P(\text{disease} \mid G)}{P(\text{disease} \mid g)}$$
  
Disease frequencies  
given genotype

(in population)

#### How to estimate relative risk?



(in population)

To estimate the relative risk, we just need to **measure the genotypes** in some disease cases and population controls.

#### How to estimate relative risk?



## Key fact



#### The odds ratio in a sample of cases and controls\* estimates the population relative risk.

Stricyly this applies to 'population controls', but also approximately true for 'true' disease controls, as long as the disease is not too common.

#### Example: O blood group and severe malaria

Cases were ascertained as children arriving in hospital with severe symptoms compatible with malaria & parasitaemia, in a hospital in Kilifi, eastern Kenya. Controls were ascertained from new births in the same hospitals.

		non-
	0	0
Severe malaria cases	686	843
Controls:	839	700

*N*=3,068 samples MalariaGEN 2019 doi: 10.1038/s41467-019-13480-z Can you compute the odds ratio?

#### Example: O blood group and severe malaria

Cases were ascertained as children arriving in hospital with severe symptoms compatible with malaria & parasitaemia, in a hospital in Kilifi, eastern Kenya. Controls were ascertained from new births in the same hospitals.

		non-
	0	0
Severe malaria cases	686	843
Controls:	839	700

N=3,068 samples MalariaGEN 2019 doi: 10.1038/s41467-019-13480-z  $OR = \frac{686}{843} \times \frac{700}{839} = 0.68$ 

Suggests people with O blood group get severe malaria at ~70% of the rate of people without

$$OR = \frac{686}{843} \times \frac{700}{839} = 0.68$$

Could say: "O blood group is associated with ~30% reduced risk of severe malaria."

But how much statistical evidence is there that this is a real effect?

The key association test summary statistics

#### Effect size estimate $\hat{\beta} = \log(OR)$

#### Standard error

se



How strong is the estimated effect? Often described on log(OR) scale

How much noise is there in the estimate, because we only have a finite sample?

P-value

Informally, a small p-value means the effect is unlikely to be zero How unlikely was such a big estimate, if actually there was no effect?

In practice computed from the beta and standard error:

 $P = \Phi^{-1}\left(\frac{\log(\mathrm{OR})}{\mathrm{se}}\right)$ 

Normal distribution function

### Incredibly useful formula

Fact: the standard error is largely determined by the study design.

Here is a very useful formula which approximates it in the 2x2 table example:



The standard error depends on sample size, frequency, and case/control ratio. It gets smaller (at rate  $\frac{1}{\sqrt{N}}$ ) as the sample size increases.

#### How many samples did we need anyway?

E.g. suppose the variant we're looking for has frequency f = 20%and the effect size is RR = 1.5. How many samples do we need?

 $P = 5 \times 10^{-8}$  corresponds to an effect about 5.5 standard errors from zero, so very roughly we need a standard error at least as small

$$\frac{\log(1.5)}{5.5} \approx 0.07$$
$$\operatorname{se}(\log OR) \approx \frac{1}{\sqrt{2N \times f(1-f) \times 0.5^2}}$$



as

Answer: we need thousands!

#### Example: O blood group is associated with malaria protection



Estimate is about 5 standard errors from zero  $P = 9.6 \times 10^{-8}$ 

## Major possible confounders



Before testing, it is imperative to look carefully at genotyping and perhaps remove samples or variants that have genotyped poorly

Because both genotypes and environments vary with geography, you should also expect to have to deal with any issues of population structure – can be either by removing samples or 'controlling' for structure.

#### Association testing in practice

In practice you would use a 'regression' method\*, rather than this simple 2x2 table approach to make these estimates:

- More flexible, e.g. allows modelling additive, dominance or recessive effects
- Can include other covariates which help explain the phenotype including *confounders*

\*E.g logistic regression (for case/control traits) or linear regression for continuous traits.

## **GWAS** roadmap

1. Collect as many samples as possible

How many samples?

2. Genotype the at as many variants across the genome as possible and do careful QC

How many variants? Which ones?

An ad hoc but widely used threshold

3. Run a statistical test for genotype-phenotype association How to test? Can we deal with confounders?

To produce this:



Lots of statistical tests so to get excited we need strong evidence e.g.  $P < 5 \times 10^{-8}$ 

# The birth of GWAS



Microarrays developed in the late 90's / early 2000's.

For the first time was possible to rapidly type hundreds of thousands or millions of SNPs

#### Patterns of inheritance generate linkage disequilbrium



patterns

#### Patterns of inheritance generate linkage disequilbrium



Idea: maybe we can just genotype a dense set of marker genotypes E.g. if we genotyped  $\triangle$ , we might pick up the true signal at 3

# The HapMap project estimated LD

The extent of LD depends on the amount of recombination.

#### A haplotype map of the human genome

The International HapMap Consortium

Inherited genetic variation has a critical but as yet largely uncharacterized role in human disease. Here we report a public database of common variation in the human genome: more than one million single nucleotide polymorphisms (SNPs) for which accurate and complete genotypes have been obtained in 269 DNA samples from four populations, including ten 500-kilobase regions in which essentially all information about common DNA variation has been extracted. These data document the generality of recombination hotspots, a block-like structure of linkage disequilibrium and low haplotype diversity, leading to substantial correlations of SNPs with many of their neighbours. We show how the HapMap resource can guide the design and analysis of genetic association studies, shed light on structural variation and recombination, and identify loci that may have been subject to natural selection during human evolution.

#### International HapMap Project doi:10.1038/nature0422 (2005)

A database of > 1M SNPs found in European, African, and Asian ancestry individuals (A subset of the samples later used in the 1000 Genomes Project)



Recombination turns out to be highly nonuniform. It is concentrated in *recombination hotspots*. So mutations are carried on longer haplotypes than had been expected.

Shared haplotype lengths

Map of recombination rate



Block-like structure of LD (correlations between SNPs in two different regions)

Tag SNP set size	Common SNPs captured (%)				
	YRI	CEU	CHB + JPT		
10,000	12.3	20.4	21.9		
20,000	19.1	30.9	33.2		
50,000	32.7	50.4	53.6		
100,000	47.2	68.5	72.2		
250,000	70.1	94.1	98.5		
As in Table 7, tag SNPs w Haploview, selecting SNP;	ere picked to capture or s in order of the fraction	ommon SNPs in HapMa of sites captured. Com	p release 16c1 using mon SNPs were		

captured by fixed-size sets of pairwise tags at  $r^2 \ge 0.8$ .

HapMap estimated how many SNPs genome-wide would need to be typed to capture (by LD) most common genetic variants. E.g. 250,000 would capture ~95% of SNPs in European populations.

# The birth of GWAS



Microarrays developed in the late 90's / early 2000's.

For the first time was possible to rapidly type hundreds of thousands or millions of SNPs

#### How a microarray works



Wash the DNA over and let it hybridise to millions of probes – one for each SNP

Flourescent markers are then attached. A picture is taken of the array.

## A microarray gives you intensities, not genotypes

#### For each (well-genotyped) SNP, you get back this:



Each dot represents DNA from one individual. X axis = image intensity for  $1^{st}$  allele Y axis = image intensity for  $2^{nd}$  allele A clustering algorithm has been used to turn the intensity values (x/y axis values) into genotype calls (colours). A microarray gives you intensities, not genotypes

#### For each SNP, you get back this:



Each dot represents DNA from one individual. X axis = image intensity for 1<sup>st</sup> allele Y axis = image intensity for 2<sup>nd</sup> allele

#### Or this if you're less lucky:



Small genotyping errors in cases or controls could easily confound the study

Careful quality control needed with these technologies

# The birth of GWAS



Microarrays developed in the late 90's / early 2000's.

For the first time was possible to rapidly type hundreds of thousands or millions of SNPs

#### Anatomy of a GWAS – what to look for

1. Collect as many cases and controls as possible

2. Genotype (or impute) them at as many variants across the genome as possible

3. Deal with potential confounders – careful data quality control and handle population structure.

4. Estimate relative risks, and look for statistical evidence that of  $RR \neq 1$ 

5. If estimate is many standard deviations from zero, bingo! We may have found a true causal effect.

6.Replicate in other studies, or find other corroborating evidence?

7. (Now try to understand the underlying biology.)



# A real GWAS study - WTCCC

## Genome-wide association study of 14,000 cases of seven common diseases and 3,000 shared controls

The Wellcome Trust Case Control Consortium\*

Nature (2007)

Studied seven common diseases in the UK

Bipolar disorder, Coronary Artery Disease, Crohn's disease, Hypertension, Rheumatoid arthritis, Type 1 and Type 2 Diabetes

Genotyped at 500,000 SNPs across the genome

doi:10.1038/nature05911

# A real study - WTCCC



www.wtccc.org.uk

#### Anatomy of a GWAS – what to look for

1. Collect as many cases and controls as possible

2. Genotype (or impute) them at as many variants across the genome as possible

3. Deal with potential confounders – careful data quality control and handle population structure.

4. Estimate relative risks, and look for statistical evidence that of  $RR \neq 1$ 

5. If estimate is many standard deviations from zero, bingo! We may have found a true causal effect.

6. Does it replicate in other studies, or have other corroborating evidence?

7. (Now try to understand the underlying biology.)

N=2,000 cases and 3,000 controls

Genotyped at 500k SNPs

Have they done adequate data quality control? Have they dealt with possible confounders?

Did they find anything with strong evidence?

Is it convincing?

What about biology?

Collection	Missingness	Heterozygosity	External discordance	Non-European ancestry	Duplicate	Relative	Total
58C	9	0	4	6	4	1	24
UKBS	8	0	5	14	0	15	42
BD	30	0	0	9	77	13	129
CAD	41	1	0	13	2	5	62
CD	43	4	6	54	131	18	256
HT	29	0	0	2	6	11	48
RA	47	1	0	26	53	9	136
T1D	7	2	1	18	6	3	37
T2D	36	1	0	11	16	11	75
Total	250	9	16	153	295	86	809

**Supplementary Table 4 | Exclusion summary by collection.** Six filters were applied for sample exclusion: 1. SNP call rate < 97% (missingness). 2. Heterozygosity > 30% or < 23% across all SNPs. 3. External discordance with genotype or phenotype data. 4. Individuals identified as having recent non-European ancestry by the Multidimensional Scaling analysis (see Methods). 5. Duplicates (the copy with more missing data was removed) 6. Individuals with too much IBS sharing (>86%); likely relatives. Where individuals could be excluded for more than one reason, they appear in the leftmost such column.

#### They then threw away 809 samples!

#### Due to:

- Poor genotyping rates
- Evidence of contamination (too many heterozygous genotypes)
- Evidence of being not of European ancestry
- A duplicate, or close relative of another sample





Some of the poor quality data was apparently due to batch effects.



PCA computes genome-wide relationships between samples and then looks for directions of greatest variation. Since relatedness typically decreases with geographic distance, principal components typically reflect geography. To avoid confounding by population structure, the samples were all supposed to be from the United Kingdom, and with European ancestry.

They used a method called *principal components analysis* to detect ancestry against the HapMap project samples. Some non-European ancestry individuals had been typed.

153 individuals were excluded on this basis.



#### Using quantile-quantile plots to assess residual confounding



They also excluded 25,567 SNPs from the study for

- High missing data rates
- Deviation from Hardy-Weinberg equilibrium (lecture 1) in controls
- Frequency differences between the two control groups
- And they visually inspected cluster plots for remaining SNPs

If there are few true signals, and if we have removed confounders – then P-values should largely come from a uniform distribution - they should lie on the diagonal.

#### Phew!



#### The main result of the study



Number of associations with strong evidence

The study found 25 associations at their nominal P-value threshold.

Twelve of these provided
replication of previously
implicated variants.
Thirteen were new
associations.

The traits clearly differ in their genetic architecture

Some SNPs were associated
 with some evidence with
 multiple traits (mainly for
 the autoimmune diseases).



Effect sizes were generally modest

E.g. across the 9 associations with Crohn's disease, the maximum estimated odds ratio was 1.54, (similar to the O blood group example)

(A strong effect with Type 1 Diabetes was also observed in the MHC locus)



#### Zooming in to a GWAS 'hit' plot

Sometimes called a 'locus zoom' plot. Here are some things to look for:



osition of SNPs in the referenc genome assembly

#### Summary

- GWAS is a very simple study design in principle just genotyping a lot of cases and controls, and test for association. The hard parts are in the implementation details
- In the early 2000's, The HapMap and other projects enabled the first GWAS by mapping SNPs genome-wide, and describing human haplotype variation.and patterns of LD. High-throughput genotyping microarray technology was developed to type these SNPs.
- The WTCCC was one of the first large GWAS studies. It provided compelling evidence that the 'common variant, common disease' hypothesis really holds.
- Although the overall design is simple, we are looking for small differences in risk between cases and controls (often RR = 1.5 or smaller). Consequently a lot of careful work is needed to ensure there is no subtle confounding – e.g. from sample collection, genotyping and data quality issues, or environmental covariates.



We have clearly learned something about the biology of these traits.

...so what?

#### Where next?

We have clearly learned something about the biology of these diseases the 'common variant, common disease' hypothesis is really true – at least for some traits, to some extent.

Raises several questions which we will get into in the next lecture, such as:

- So how polygenic do traits get?
- What about the biology underlying these associations?



#### Let's zoom in

# **Biology is hard**



# Biology is hard



No genes under the main association signal!



## **Biology is hard**

Association observed with CAD over a ~100kb region of chromosome 9. This is unquestionably a real association (it has been replicated in several independent studies).

The functional mechanism of this association is not fully solved; it probably involves regulation of expression of the two nearby genes *CDKN2A/B*.

Neither gene was an obvious candidate beforehand - thus, this association does point to novel biology.







This association with Type 2 Diabetes turned out to be through a second, related trait (obesity), again unquestionably a real effect. But as of 2018 the functional mechanism remains unclear. Expression of *FTO* is known to affect obesity, but the SNPs may also affect expression of another gene, IRX3, 200kb away.

Smemo et al, Nature 2014



This pattern has turned out to be typical. It has generally proven extremely hard to narrow down GWAS associations to underlying 'causal' variants.

LD is a double-edged sword.

Next lecture: we will look at this.

#### Anatomy of a GWAS – what to look for

1. Collect as many cases and controls as possible

2. Genotype (or impute) them at as many variants across the genome as possible

3. Deal with potential confounders – careful data quality control and handle population structure.

4. Estimate relative risks, and look for statistical evidence that of  $RR \neq 1$ 

5. If estimate is many standard deviations from zero, bingo! We may have found a true causal effect.

6.Replicate in other studies, or find other corroborating evidence?

7. (Now try to understand the underlying biology.)



## **Consolidation question**

Region

N/A

PLEK

Signal



#### Multiple Sclerosis GWAS Browser

This site accompanies the article "Genetic risk and a primary role for cell-mediated immune mechanisms in multiple sclerosis", The International Multiple Sclerosis Genetics Consortium (IMSGC) and the Wellcome Trust Case Control Consortium 2 (WTCCC2), Nature (2011). The data sources for this page are those described in the above article, and were current at the time of article preparation. Show more details



## GWAS of multiple sclerosis (2011) 9772 cases, 17,376 controls from across Europe

www.chg.ox.ac.uk/wtccc2/ms/ (I think this url requires the trailing /)

Visit the above site and make sure you understand what is shown. Pick a signal and try to work out

- What is the estimated effect size?
- How strong was the evidence?
- Did it replicate?
- Does the association signal look sensible does it follow LD patterns, and do the cluster plots look sensible?
- Can you figure out what the nearby genes do? (Warning: this can be a time sink!)

**Bonus question**: read the paper and try to figure out the questions on the checklist.

Next lecture: Friday 7th Mar @14:30

# Understanding the genetics of complex traits II

Gavin Band gavin.band@well.ox.ac.uk

**BA Human Sciences** 

Friday 7th March 2025

