# StatGen Workshop IGSS

Benjamin Neale, Ph.D.
Analytic and Translational Genetics Unit, MGH
Stanley Center for Psychiatric Research & Program in Medical and Population Genetics, Broad Institute

# Analysis of UK Biobank

# GWAS of UK Biobank



**Download and decryption** → **Software development** → **Phenotype wrangling** → **QC and GWAS**

Sam Bryant

Cotton Seed

Andrea Ganna, Duncan Palmer, Caitlin Carey

Liam Abbott
Dan Howrigan

**Also thanks to**:

| Verneri Anttila | Jon Bloom | Mark J. Daly | Jackie Goldstein | Eric Jones | Ruchi Munshi |
| Krishna Aragam | Joanne Cole | Rob Damien | Mary Haas | Sekar Kathiresan | Tim Poterba |
| Alex Baumann | Mark J. Daly | Steven Gazal | Joel Hirschhorn | Dan King | Manuel Rivas |
| | | | | | Sailaja Vedantam |

- Follows health and well-being of 500,000 participants
- Genotyped using the Affymetrix Biobank Array
- Lots of phenotypes collected [needs harmonization]
- Lots of opportunity!

# Data showcase
## http://biobank.ctsu.ox.ac.uk/crystal/

# Sex distribution

502,620 items of data are available, covering 502,620 participants, encoded using Data-Coding 9.



Female [273,455]

Male [229,165]

0   60   120   180   240   300
(thousands)

Counts of participants/items last updated 27 Jul 2017.

# Age distribution at recruitment



Mean = 56.5286
Std.dev = 8.09516

# Example self-report

# PHESANT!



Copious thanks to Millard LAC, Davies NM, Gaunt TR, Davey Smith G, Tilling K. PHESANT: a tool for performing automated phenome scans in UK Biobank. bioRxiv (2017)
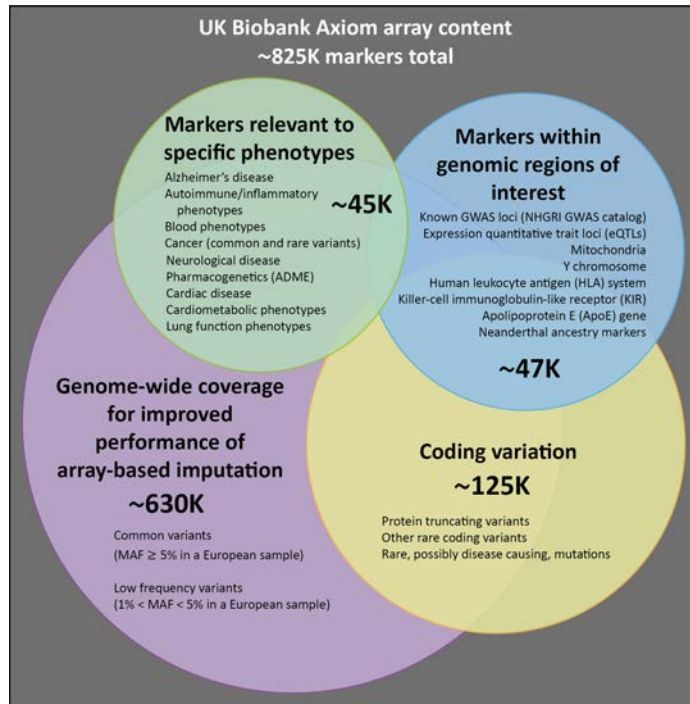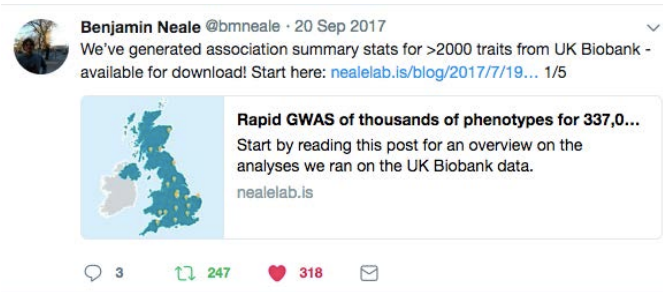
# What's on the array?



Imputed to HRC

# Round 1 GWAS

- Last fall, the Neale lab…
  - GWASed 2,419 phenotypes
    - Blogged about it
    - Put them on dropbox
      - And people made browsers
  - Estimated h² for all of them
  - Made an h² browser
    - Blogged about that too

Nealelab.is/blog

# Scalability



10hr

{ 10 compute hours

# Scalability

# Scalability

# Association results for many things! Taking cholesterol lowering meds

**LDL Cholesterol**

Genetic correlation of 0.47 with LDL, 0.58 with triglycerides and 0.51 with total cholesterol

GLGC Nat Genet

6 months later, we did it all again

# Why Round 2 of UKB GWAS?

- Missing a batch of imputed SNPs
  - Corrected data released in March

- Hadn't gotten permissions for *all* the phenotypes
  - Expanded UKB application

- Feedback on improvements for the GWAS
  - Age, sex, stratification

# Round 2: QC Updates

- Variant QC:
  - Added the new imputed data
  - Added chrX variants
  - Added VEP missense and PTVs with MAF > 1e-6
  - Net: 3 million more variants
    - 13.8 million total

- Sample QC:
  - Relaxed restriction to "white British" samples

# How "White British" is defined

- What is your ethnic group?
  - **White**
  - Mixed
  - Asian or Asian British
  - Black or Black British
  - Chinese
  - Other ethnic group
  - Do not know
  - Prefer not to answer

- What is your ethnic background?
  - **British**
  - Irish
  - Any other white background
  - Prefer not to answer

- Don't be defined as a PCA outlier
  - Bayesian outlier detection algorithm on PCs 1&2, 3&4, and 5&6

# How "White British" is defined



× In white British ancestry subset
× Self-reported ethnic background British
○ Other ethnic background

# Widening out definition of Europeans

- Get mean and SD of top 6 PCs among the "white British"

- Draw ellipse in PCA space with radius of 7 SDs along each PC axis
  - Provides good predictive accuracy for self-reporting "White" vs. other ethnicities

- Discard any self-reported as non-white

- Final N (after QC): 361,194
  - Previously 337,199

# Round 2: GWAS Changes

- Add age, age$^2$, sex*age, and sex*age$^2$ as covariates

- Increase number of PC covariates from 10 to 20

- Compute PCs within the GWAS sample rather than using the PCs computed by UKB on the full sample

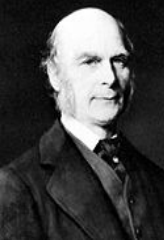- In addition to main GWAS, run sex-specific GWAS [withou sex covariates

# Let's go to the code!

https://github.com/Nealelab/UK_Biobank_GWAS



We'll start with the readme

# Francis Galton
# Twin and family studies

RATE OF REGRESSION IN HEREDITARY STATURE.
Fig. (a)

- Relatives are more similar

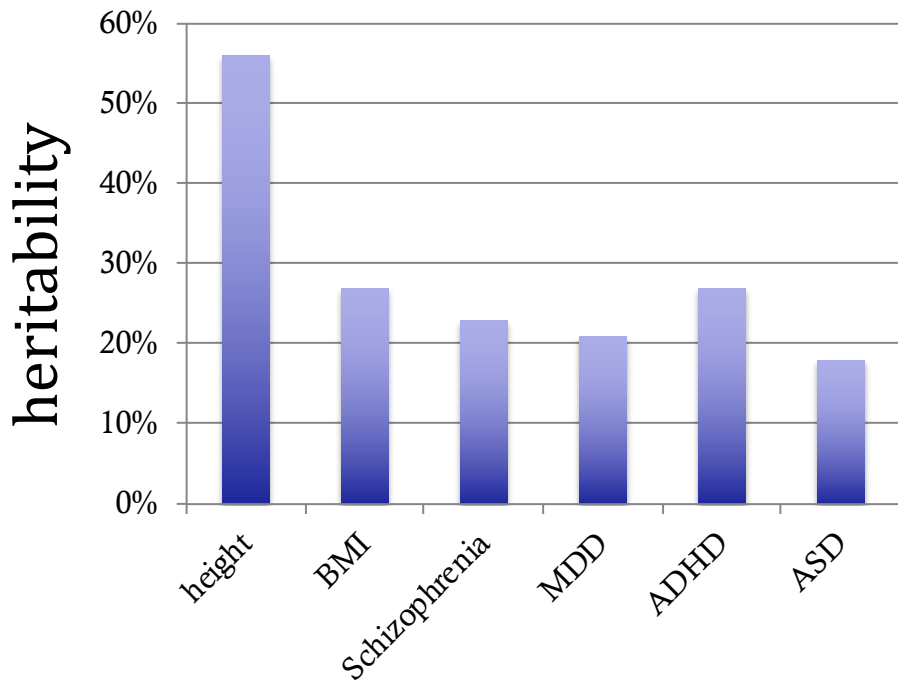## Meta-analysis of the heritability of human traits based on fifty years of twin studies

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

Average estimate of heritability 49%
69% of twin studies support a purely additive genetic model

# GREML/GCTA



- Use estimated genetic similarity

# LD Score regression

With thanks



Brendan Bulik-Sullivan     Hilary Finucane     Po-Ru Loh     Mark Daly     Alkes Price

# How does LD shape association?

# How does LD shape association?

Lonely SNPs [no LD]

LD blocks

# How does LD shape association?

| Lonely SNPs [no LD]

LD blocks

\* Causal variants

Association

All markers correlated with a causal variant show association

# How does LD shape association?

| Lonely SNPs [no LD]

▣ LD blocks

\* Causal variants

Association

Lonely SNPs only show association if they are causal

# What happens under polygenicity?

| Lonely SNPs [no LD]
▢ LD blocks
✳ Causal variants

Assuming a uniform prior, we see SNPs with more LD friends showing more association

The more you tag, the more likely you are to tag a causal variant

# Simulated polygenic architecture
## Lambda = 1.30 LD score intercept = 1.02

# What happens under stratification?

| Lonely SNPs [no LD]
▨ LD blocks
✳ Causal variants

Under pure drift we expect LD to have no relationship to differences in allele frequencies between populations

# UK controls versus Sweden controls
## Lambda = 1.30 LD score intercept = 1.32

# PGC Schizophrenia

Lambda = 1.48

Intercept = 1.06

Slope $p$-value $< 10^{-300}$

Overwhelming majority of inflation is consistent with polygenic architecture

# LD Score regression



Draw polygenic effects from
$N(0, n/m^2)$, var =

What is the $E[\chi^2]$ for variant $j$?

$$E[\chi_j^2] = 1 + Na + \boxed{\frac{h_g^2 N}{M}} l_j$$

New estimator of heritability

where N=sample size, M=# of SNPs, a=inflation due to confounding,
$h^2 g$ is heritability (total obs.) and $l_j$ is the *LD Score*

Bulik-Sullivan et al. Nature Genetics 2015
Yang et al. EJHG 2011

$$l_j = \sum_{k \neq j} r_{jk}^2$$

# 9,928 GWAS later... let's talk $h^2$ using LD score regression

$$E\left[\chi_j^2\right] = 1 + Na + \frac{h_g^2 N}{M} l_j$$

Estimating heritability from GWAS summary statistics

# How do round 2 ldsc results compare?

Raymond Walters

- Intercept less significant
- h2 more significant with stable estimates



Intercept -$\log_{10}$(p) of old



h$^2$ -$\log_{10}$(p) of old

# Let's look at heritability

Raymond Walters



Lymphocyte count
Reticulocyte count
Reticulocyte %
High light scatter reticulocyte %



Reticulocyte count

# What about sex-specific effects?

Raymond Walters

- Sex-specific GWAS allow us to scan for:
  - Differences in female vs. male $h^2$
    - E.g. could indicate differences in variance of environmental effects, measurement differences
  - female vs. male $r_g < 1$
    - E.g. relative effects of different SNPs differ by sex

- Can also test for SNP-level differences
  - Slower and labor intensive, so $h^2$, $r_g$ can help prioritize

- To start: look at 448 phenotypes with Neff > 10000 in both sexes and z-score of h2 > 4 is at least 1 sex

# Strong h² observed in both sexes

- >70% of traits at least nominally heritable in each sex
  - P < .05

- Mean h² ~ .09

- Consistent with joint analysis of both sexes

# Is h² equal across sexes?

## h² strongly correlated across sex



## ~10% of traits have nominally different h2 between sexes

| description | Fem. h2 | Male h2 | P diff |
|---|---|---|---|
| Average weekly beer plus cider intake | 0.0416 | 0.1152 | 3.11E-10 |
| Diastolic blood pressure, automated | 0.1799 | 0.1160 | 1.13E-06 |
| Systolic blood pressure, automated | 0.1768 | 0.1208 | 1.03E-05 |
| Number of operations, self-reported | 0.0845 | 0.0491 | 2.53E-05 |
| Duration of vigorous activity | 0.0037 | 0.0555 | 3.91E-05 |

# Female (1) vs male (0) GWAS

Michel Nivard  Mattijs van der Zee



50_pheno_sex GWAS

$h^2$ (ldsc) = 0.012 (0.002)

Whaaaa?

# Differential ascertainment bias



negative rg (absence of low males or high females)

positive rg (absence of high males or low females)

-1.0   -0.8   -0.6   -0.4   -0.2   0.0    0.0    0.2    0.4    0.6    0.8    1.0

- p<7.69e-04  - p<.05  - p>=.05

GPC_A

GPC_C

Alzheimer

Neuroticism

TV time (UKB)

MDD PGC1

Neuroticism (UKB)

BMI (UKB)

Body fat

Health rating (UKB)

DS_Full

Hip circumference (UKB)

Weight (UKB)

Waist circumference (UKB)

N sexpartners (UKB)

Risktaking (UKB)

Cannabis

Cannabis use (UKB)

Alcohol

StrenSport (UKB)

PC time (UKB)

Smoking cigs per day

Alcohol (UKB, 20414)

IQ

Education years

IQ (UKB, 20016)

IQ (UKB, 20191)

Household income (UKB)

# Male/Female genetic correlation



Raymond Walters

- Next step is to look at genetic correlation between female and male results for each trait
  - Again using LD score regression

- Focus on 448 traits with significant $h^2$ in at least one sex
  - After Bonferroni correction for 865 traits

# Genetic correlation estimate between females and males



Cerebellum, Pancreas, Brain cortex, Disease: Female height

Female:Male Genetic Correlation

# Phenotypes with male/female rg significantly < 1 (p < 1e-5)



Trunk fat percentage* (P=8e-28)
Age first had sexual intercourse (P=1e-22)
Hip circumference (P=3e-16)
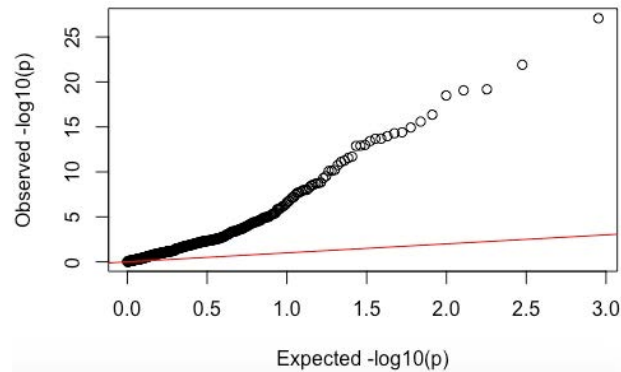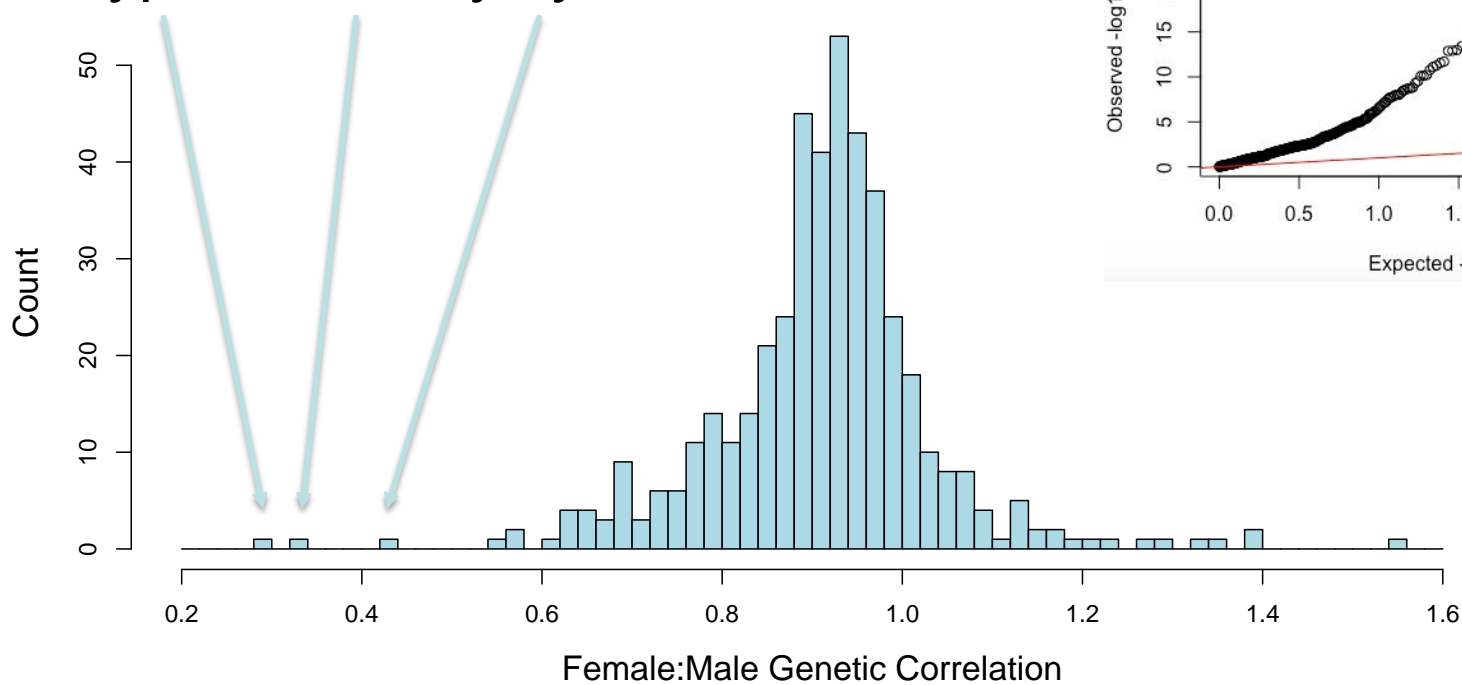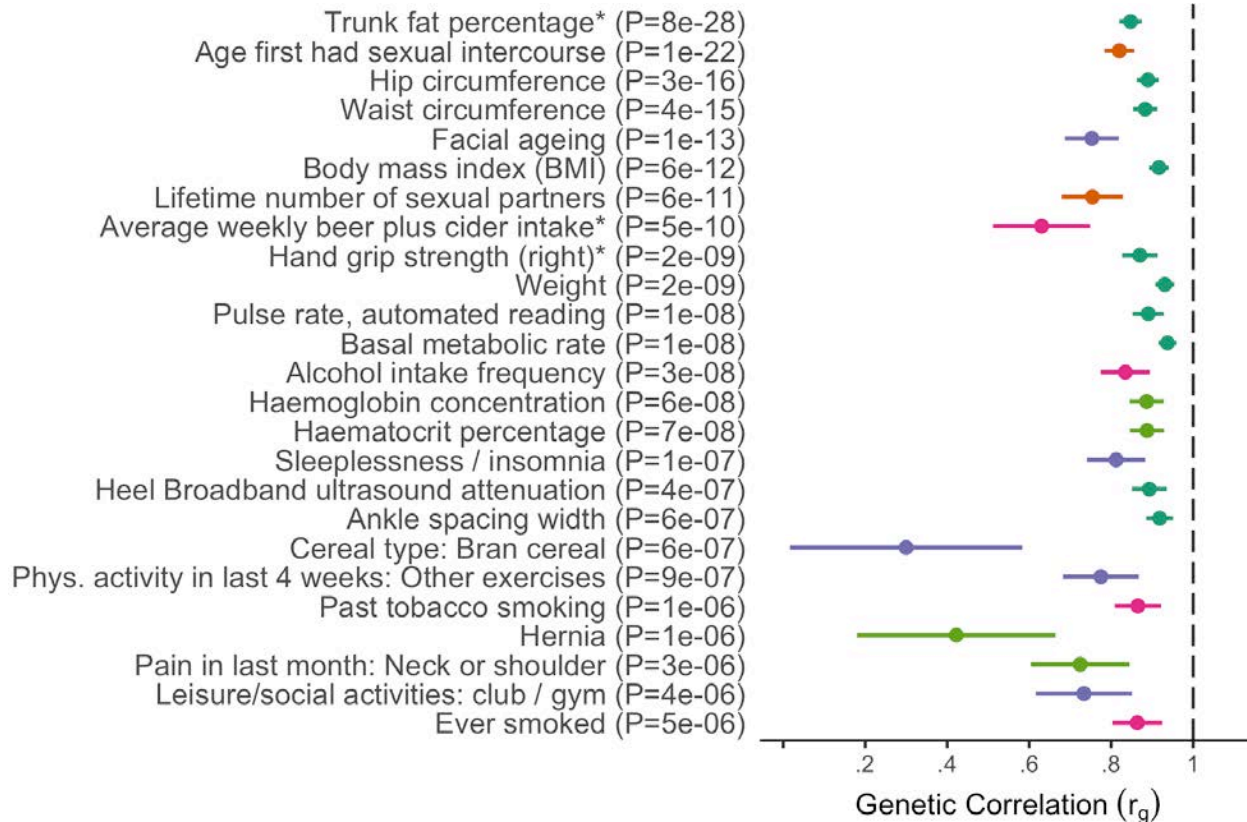Waist circumference (P=4e-15)
Facial ageing (P=1e-13)
Body mass index (BMI) (P=6e-12)
Lifetime number of sexual partners (P=6e-11)
Average weekly beer plus cider intake* (P=5e-10)
Hand grip strength (right)* (P=2e-09)
Weight (P=2e-09)
Pulse rate, automated reading (P=1e-08)
Basal metabolic rate (P=1e-08)
Alcohol intake frequency (P=3e-08)
Haemoglobin concentration (P=6e-08)
Haematocrit percentage (P=7e-08)
Sleeplessness / insomnia (P=1e-07)
Heel Broadband ultrasound attenuation (P=4e-07)
Ankle spacing width (P=6e-07)
Cereal type: Bran cereal (P=6e-07)
Phys. activity in last 4 weeks: Other exercises (P=9e-07)
Past tobacco smoking (P=1e-06)
Hernia (P=1e-06)
Pain in last month: Neck or shoulder (P=3e-06)
Leisure/social activities: club / gym (P=4e-06)
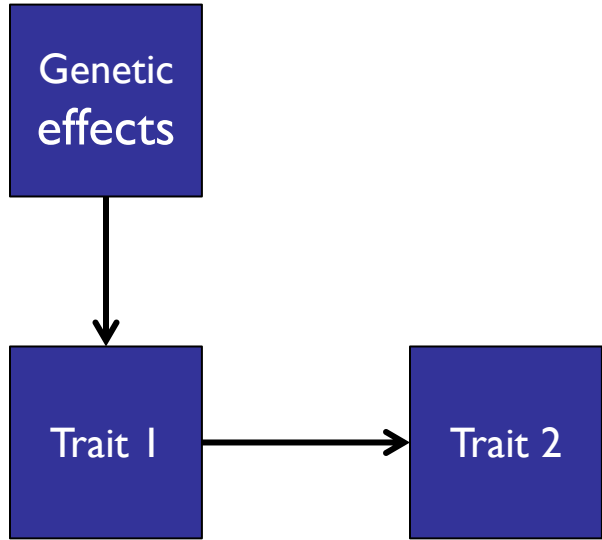Ever smoked (P=5e-06)

Genetic Correlation (r$_g$)

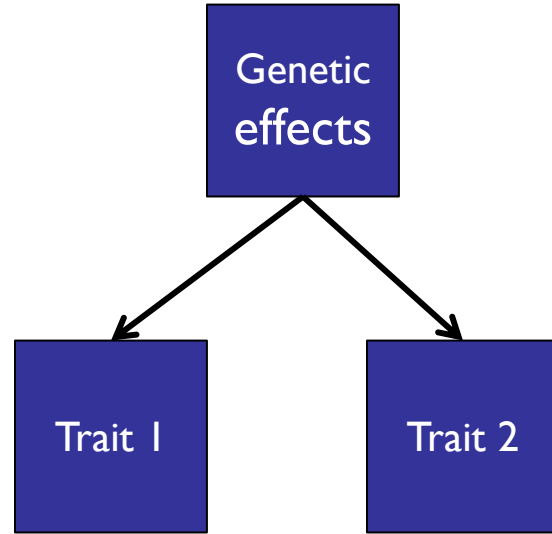# Genetic Correlation Method in:

An atlas of genetic correlations across human diseases and traits

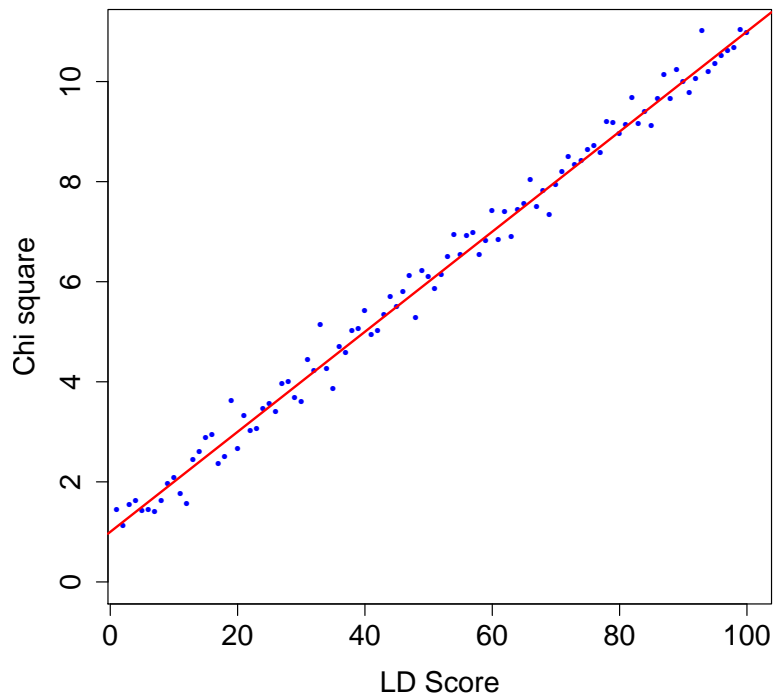# Potential sources of genetic correlation



Trait 1 exerts causal effect on Trait 2

Genetic effects influence
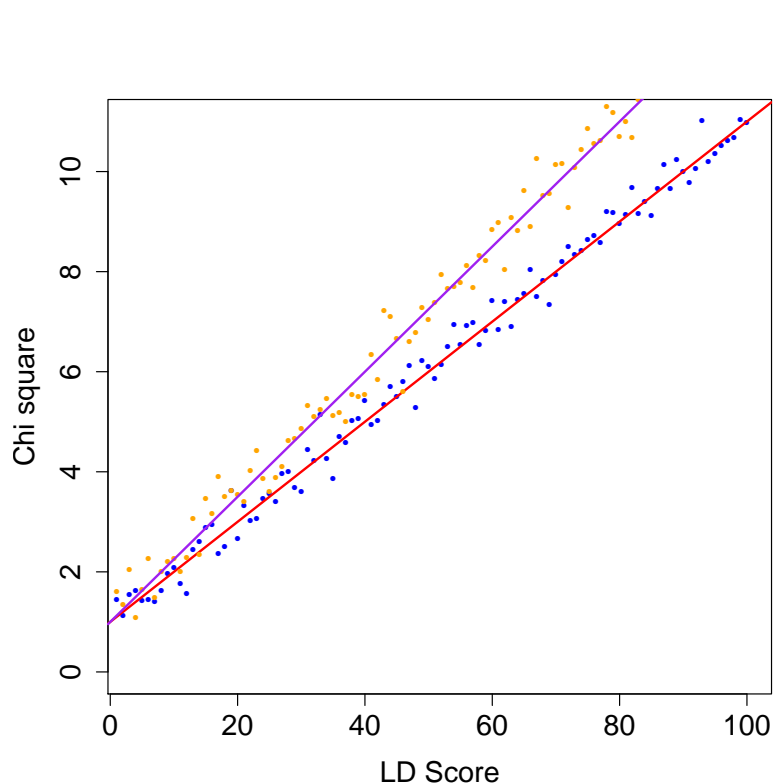Trait 1 and Trait 2

# LD Score regression
# Genetic correlation



| Trait 1

Slope estimates heritability

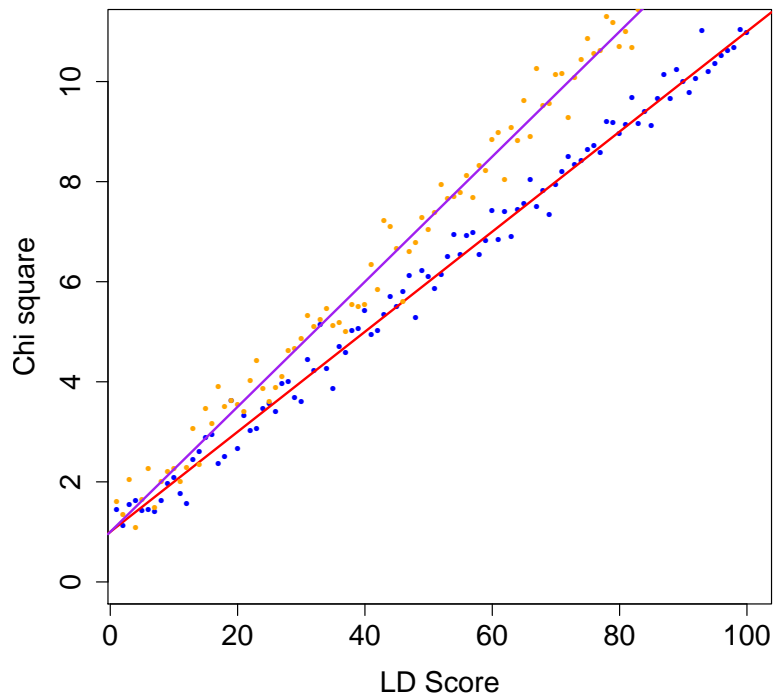# LD Score regression
# Genetic correlation



| Trait 1
| Trait 2

We can a second trait and obtain two heritability estimates
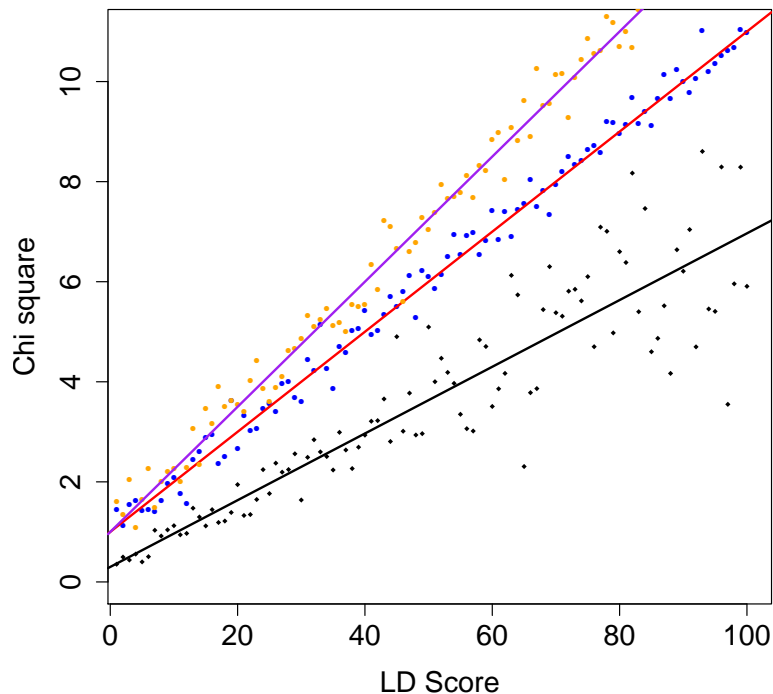
# LD Score regression
# Genetic correlation



| Trait 1
| Trait 2

$Z*Z = \chi^2$

So we can estimate genetic covariance from the product of the Z-scores
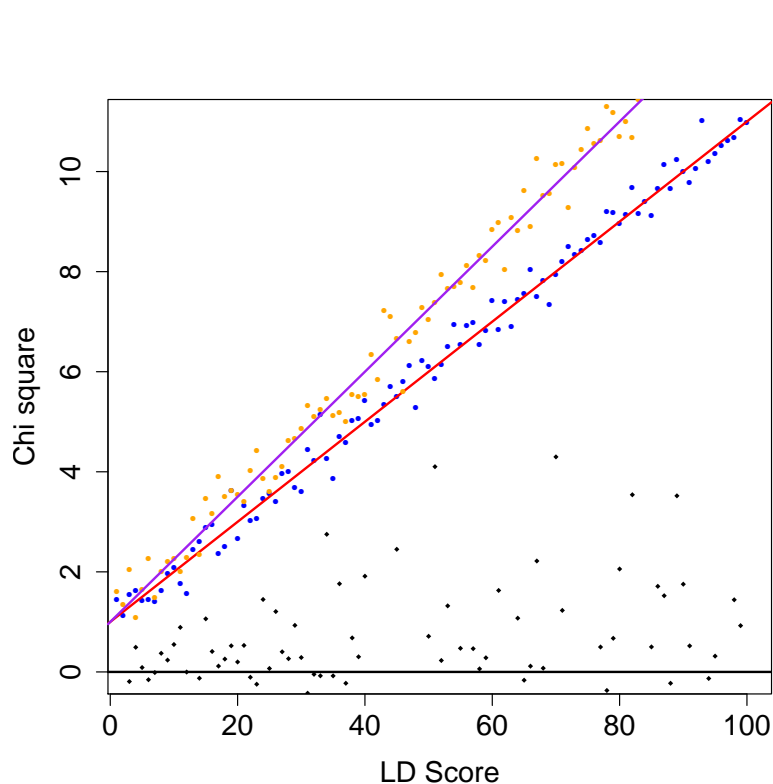
# LD Score regression
# Genetic correlation



Trait 1
Trait 2
$R_G$

$Z*Z = \chi^2$

So we can estimate genetic covariance from the product of the Z-scores for the two traits
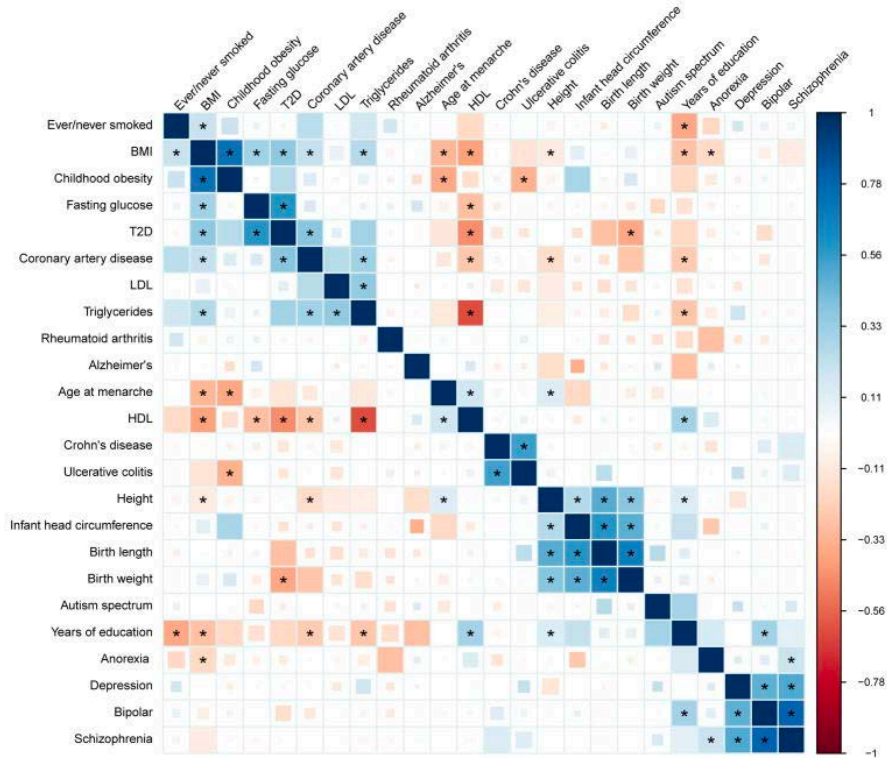
$R_G = 0.5$

# LD Score regression
# Genetic correlation



| Trait 1
| Trait 2
| $R_G$

Here $R_G = 0$

This approach is robust to sample overlap as all variants are equally inflated

# Genetic correlations



- Genetic correlation is a widespread phenomenon

# Brainstorm Project

Verneri Anttila

Aiden Corvin

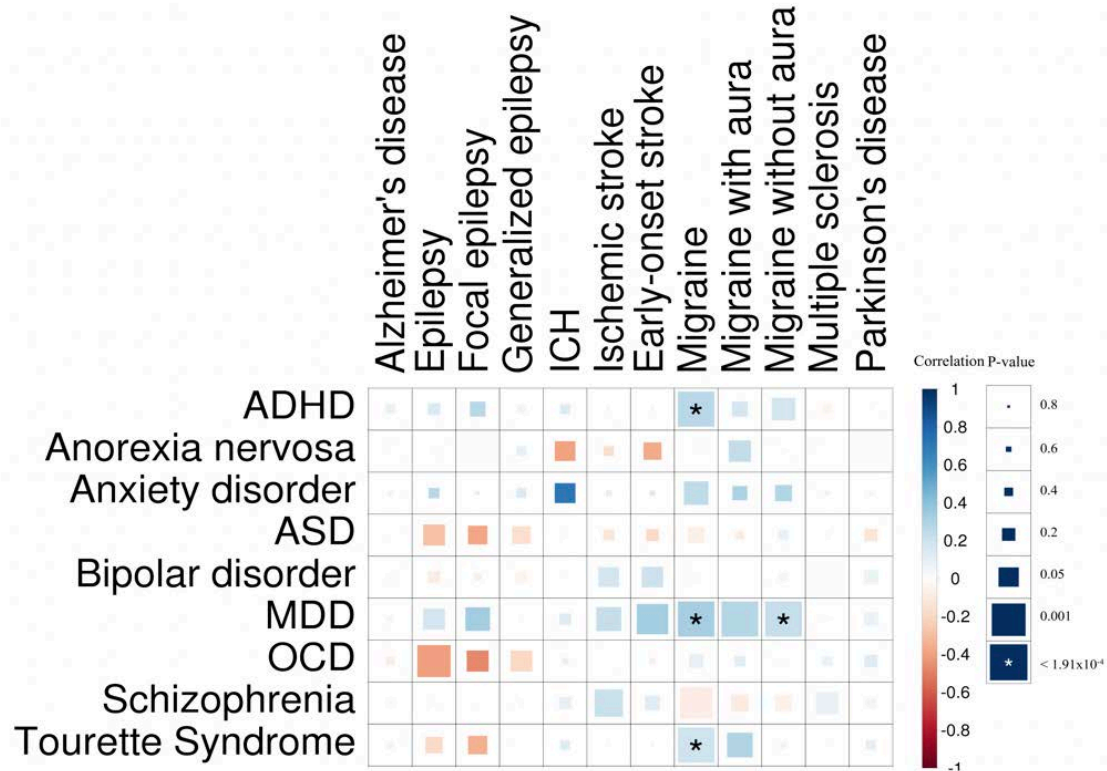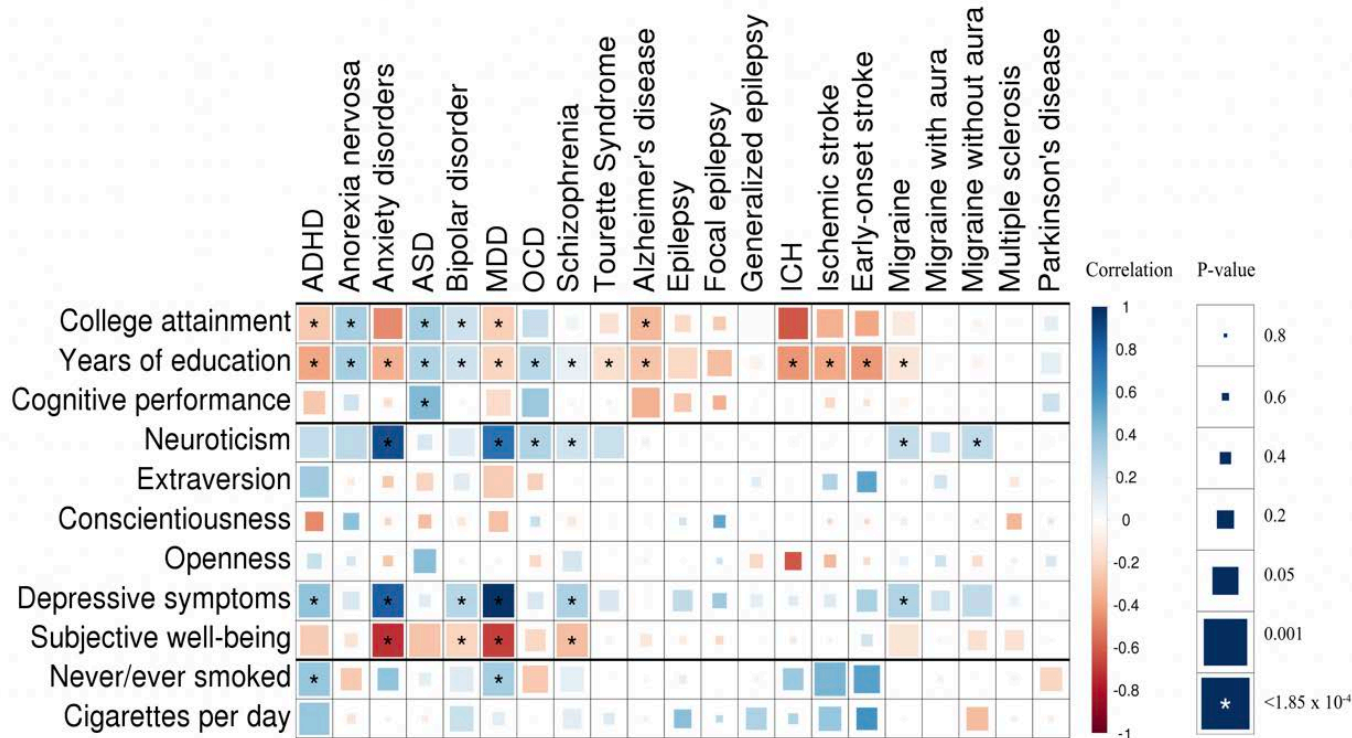| | | |
|---|---|---|
| **Brendan Bulik-Sullivan** | Alessandro Biffi | Hailiang Huang |
| **Hilary Finucane** | Jeremiah Scharf | Andrea Byrnes |
| Jonathan Rosand | Kenneth Kendler | Dongmei Yu |
| Aarno Palotie | Stephan Ripke | Laramie Duncan |
| Mark Daly | Alkes Price | Kai-How Farh |
| Patrick Sullivan | Chris Cotsapas | Namrata Gupta |
| Bobby Koeleman | Padhraig Gormley | Miriam Raffeld |
| Nick Wood | Zhi Wei | …and many, many others |
| Julie Williams | Rainer Malik | in their respective study groups |

# Brainstorm within psychiatry

# Brainstorm within neurology

# Brainstorm – across neurology and psychiatry

# Brainstorm – take it further?

# Comprehensive evaluation of genetic correlation

Duncan Pal[...]



https://ukbb-rg.hail.is/        https://github.com/astheeggeggs/UKBB_ldsc_r2