

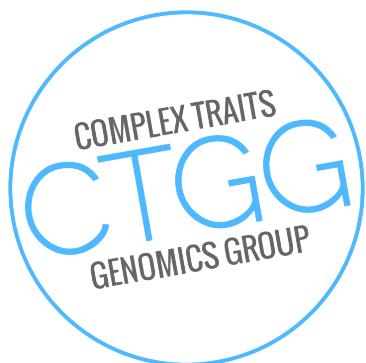


THE UNIVERSITY
OF QUEENSLAND
AUSTRALIA

Mixed linear model analyses of human complex traits using SNP data

Jian Yang

Queensland Brain Institute
The University of Queensland



Why do we need a mixed linear model?

Linear model

- $y = b_0 + x_1 b_1 + x_2 b_2 + \dots + x_p b_p + e$

y = phenotype

x_i = independent variable

$y \sim N(b_0 + x_1 b_1 + x_2 b_2 + \dots + x_p b_p, \sigma^2_e)$

b_0 = mean term

$b_1 \dots b_p$ = effect sizes (regression coefficients)

e = residual, $e \sim N(0, \sigma^2_e)$

Linear model

- In matrix form

$$\mathbf{y} = \mathbf{X}\mathbf{b} + \mathbf{e}$$

$$\mathbf{y} = \{y_j\}_{n \times 1}; \mathbf{X} = \{X_{ij}\}_{n \times p}; \mathbf{b} = \{b_i\}_{p \times 1}; \mathbf{e} = \{e_j\}_{n \times 1}$$

- Estimation

$$\hat{\mathbf{b}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$$\text{var}(\hat{\mathbf{b}}) = \sigma_e^2 (\mathbf{X}^T \mathbf{X})^{-1}$$

Special cases

- Simple regression

$$y = b_0 + x_1 b_1 + e$$

$$\hat{b}_1 = \hat{\mathbf{b}} = \mathbf{X}_1^\top \mathbf{y} / (\mathbf{X}_1^\top \mathbf{X}_1)$$

$$E(\hat{b}_1) = b_1 = \text{cov}(x_1, y) / \text{var}(x_1)$$

$$\text{var}(\hat{b}_1) = \sigma_e^2 / [n * \text{var}(x_1)]$$

- Conditional analysis

$$y | b_2 \dots b_p = b_0 + x_1 b_1 + e$$

Limitations

- $n > p$: sample size needs to be $>$ than the number of parameters
- All the effect sizes are treated as fixed (we have no idea about the variation in effect sizes)
- What if $n \ll p$?

What is a mixed linear model (MLM)?

- $\mathbf{y} = \mathbf{X}\mathbf{b} + \mathbf{Z}\mathbf{u} + \mathbf{e}$

Fixed effects: \mathbf{b} (special case: $\mathbf{X} = \mathbf{1}$ and $\mathbf{b} = \mathbf{b}_0$)

Random effects: $\mathbf{u} = \{u_i\}$, $\mathbf{u} \sim N(\mathbf{0}, \sigma^2_u \mathbf{A})$

\mathbf{A} = correlation matrix between u_i and u_j

$$E(\mathbf{y}) = \mathbf{X}\mathbf{b}$$

$$\text{var}(\mathbf{y}) = \mathbf{V} = \mathbf{Z}\mathbf{A}\mathbf{Z}^T \sigma^2_u + \mathbf{I}\sigma^2_e$$

Parameter estimation

- Estimation of variance components (σ^2_u)

$$\log L = -\frac{1}{2}(\log |\mathbf{V}| + \log |\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X}| + \mathbf{y}^T \mathbf{P} \mathbf{y})$$

$$\mathbf{P} = \mathbf{V}^{-1} - \mathbf{V}^{-1} \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1}$$

- Prediction of random effects (\mathbf{u})

$$\hat{\mathbf{u}} = \sigma^2_u \hat{\mathbf{Z}}^T \mathbf{P} \mathbf{y}$$

- Estimation of fixed effects (\mathbf{b})

$$\hat{\mathbf{b}} = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$$\text{Linear model: } \hat{\mathbf{b}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

MLM analysis of human complex traits

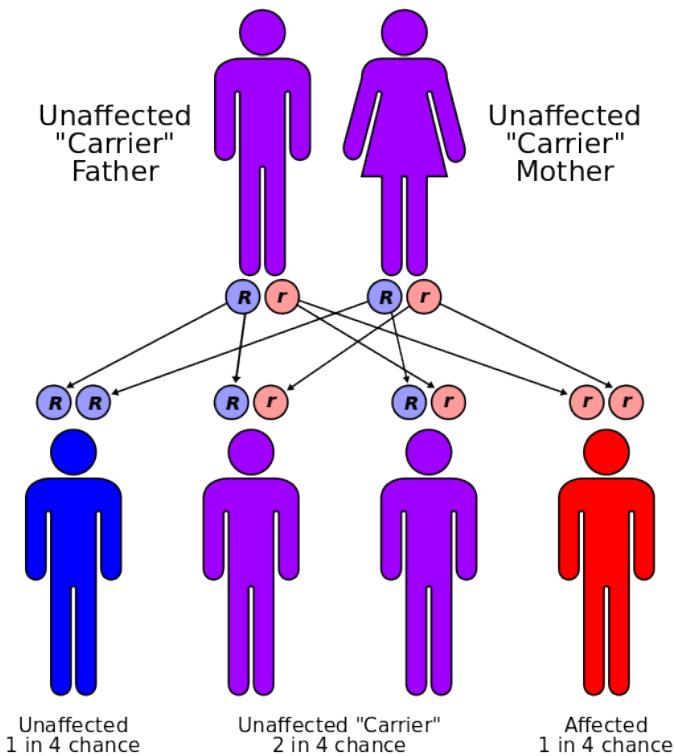
- Animal and plant breeding
 - predicting breeding values
 - linkage mapping (QTL mapping)
- Human genetics (before 2007)
 - pedigree based analysis of variance (heritability)
 - linkage mapping
- Human genetics (after 2007)
 - **estimating SNP-based heritability**
 - association analysis
 - genetic risk prediction

Background

Mendelian traits

Complex traits

Cystic fibrosis



Human height



Schizophrenia



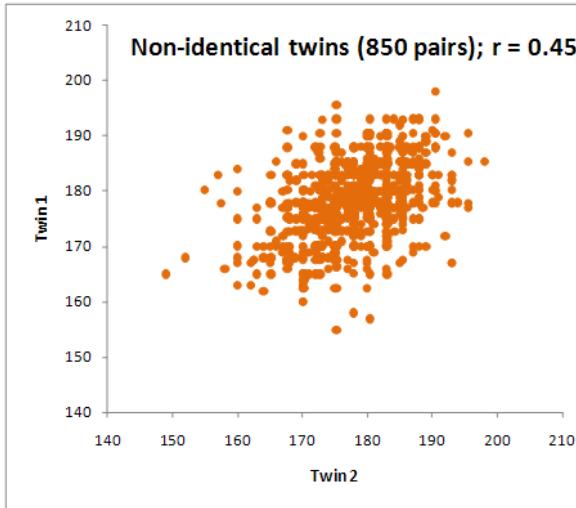
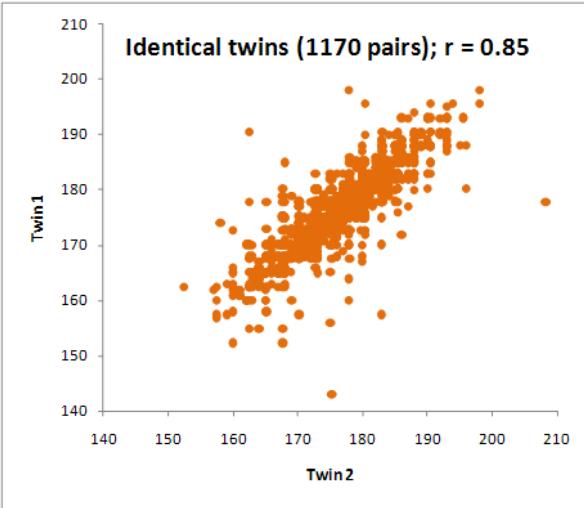
Obesity



Major questions

- Are these traits heritable?
- If so, what is the heritability?
- How many genes involved and where are they located?

Resemblance between twins for human height



Heritability = ~80%

Resemblance between relatives for body mass index (BMI)

Relatedness	Correlation
Full-sibs	0.36
Father-son	0.28

Heritability = 40%~60%

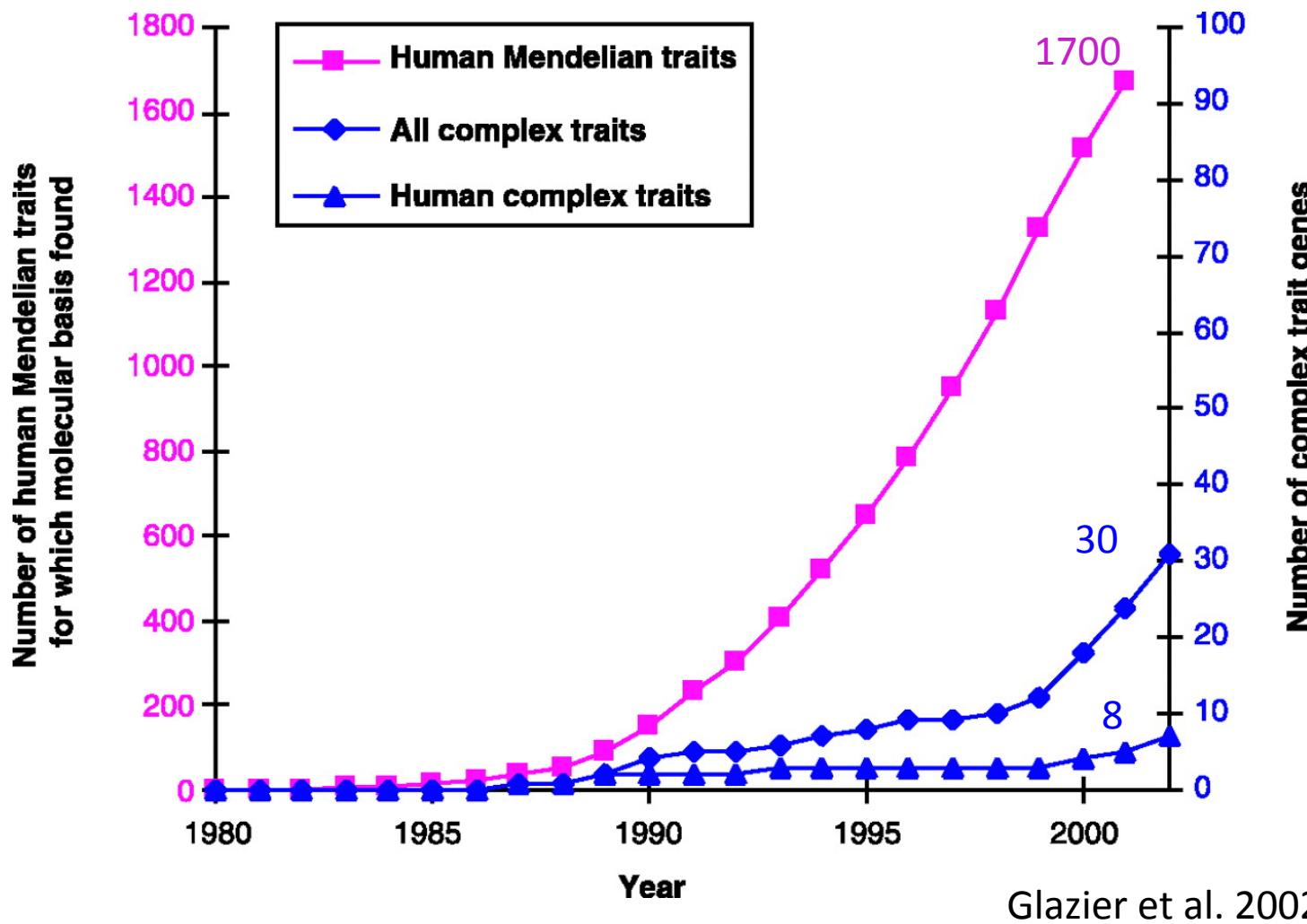
Risk of schizophrenia (%)

Heritability = ~80%

	Population	1st degree relative	MZ twin
Schizophrenia	1	10	50

Complex traits such as height, BMI and SCZ are highly heritable.

Identifying genes underlying complex traits



8 genes for human complex traits before 2002

Genome-wide Association Study (GWAS)

Linear model (simple regression)

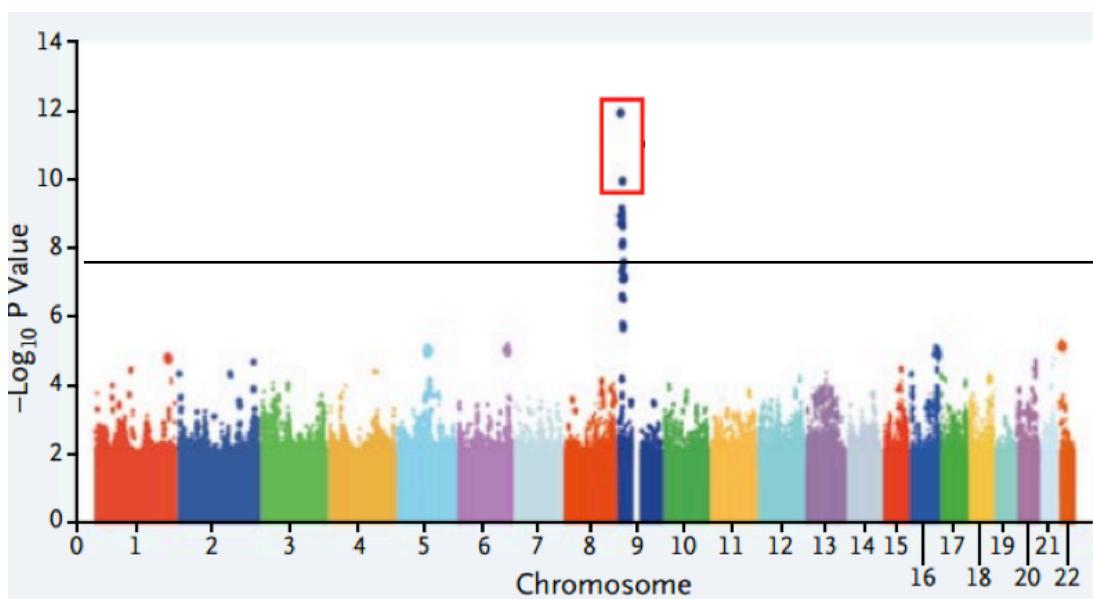
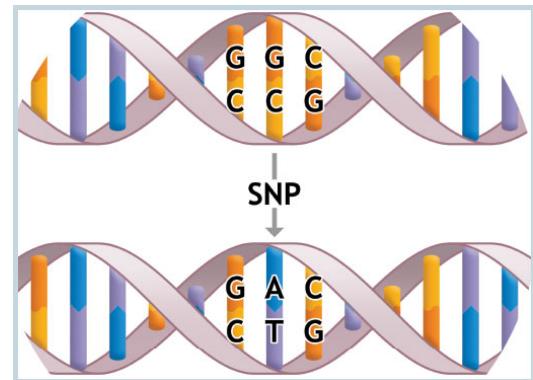
$$y = b_0 + x_1 b_1 + e$$

y = trait value

x_1 = SNP genotype (0, 1 or 2)

$$\hat{b}_1 = \mathbf{X}_1^T \mathbf{y} / (\mathbf{X}_1^T \mathbf{X}_1) = \text{cov}(x_1, y) / \text{var}(x_1)$$

$$\text{SE}^2(\hat{b}_1) = \sigma_e^2 / [n \text{ var}(x_1)]$$

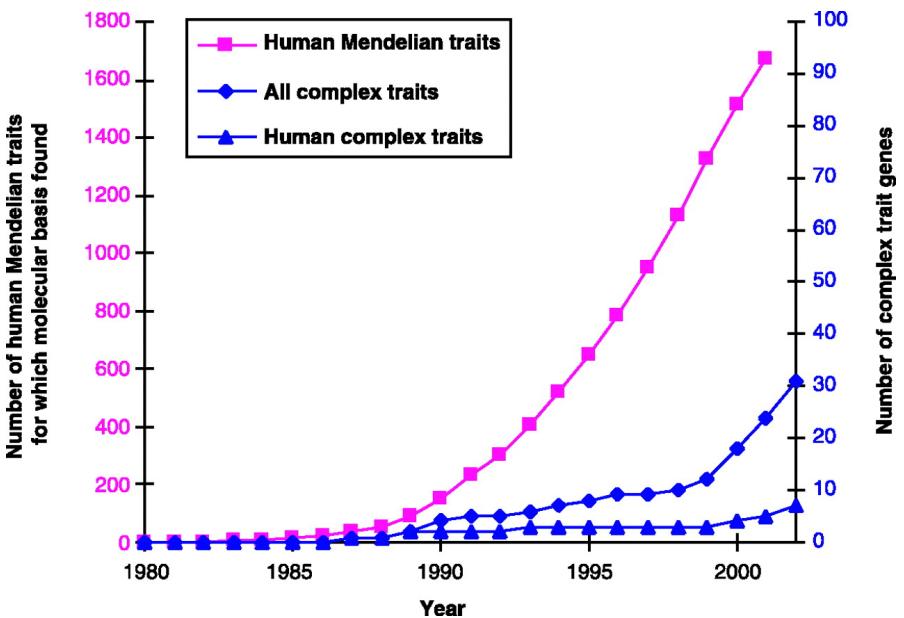


Genome-wide threshold $P = 5 \times 10^{-8}$

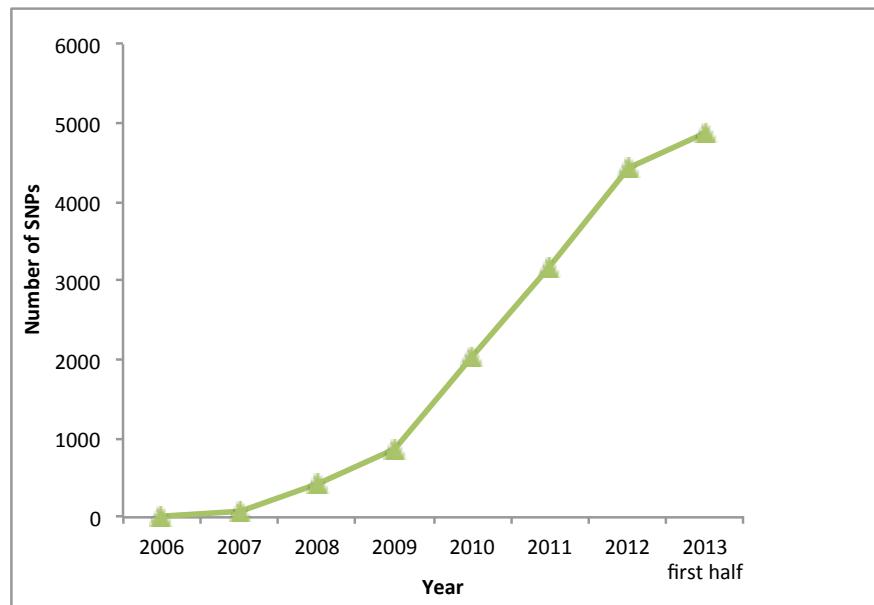
Manolio 2010 NEJM

An explosion of gene discoveries

Prior to GWAS



GWAS



Glazier et al. 2002 Science

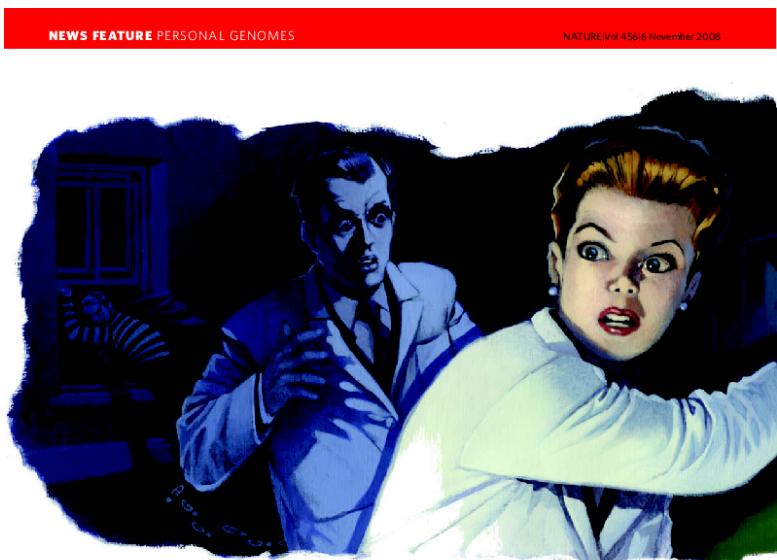
~5000 genetic variants associated with ~650 traits / diseases

The missing heritability problem

Height:

- 180 loci
- ~180K samples
- < 10% of variance explained
- heritability = ~80%

Lango Allen et al. 2010 Nature



The case of the missing heritability

When scientists opened up the human genome, they expected to find the genetic components of common traits and diseases. But they were nowhere to be seen. **Brendan Maher** shines a light on six places where the missing loot could be stashed away.

Schizophrenia:

- 22 loci
- ~21K cases / ~38K controls
- < 3% of variance explained
- heritability = ~80%

Ripke et al. 2013 Nat Genet

BMI:

- 32 loci
- ~250K samples
- ~1% of variance explained
- Heritability = 40% ~ 60%

Speliotis et al. 2010 Nat Genet

Fitting all SNPs in a MLM

- $\mathbf{y} = \mathbf{W}\mathbf{u} + \mathbf{e}$
 $\mathbf{W} = \{w_{ij}\}_{n \times m}$, w_{ij} = standardised SNP genotype
 $\mathbf{u} \sim N(\mathbf{0}, I\sigma^2_u)$
 $\text{var}(\mathbf{y}) = \mathbf{Z}\mathbf{Z}^T\sigma^2_u + I\sigma^2_e$
variance explained = $m\sigma^2_u / (m\sigma^2_u + \sigma^2_e)$
- Let $\mathbf{g} = \mathbf{Z}\mathbf{u}$
 $\mathbf{y} = \mathbf{g} + \mathbf{e}$
 $\mathbf{g} \sim N(\mathbf{0}, \mathbf{A}\sigma^2_g)$, \mathbf{A} = genetic relationship matrix
 $\text{var}(\mathbf{y}) = \mathbf{A}\sigma^2_g + I\sigma^2_e$
variance explained = $\sigma^2_g / (\sigma^2_g + \sigma^2_e)$
- $\text{var}(\mathbf{y}) = (1/m)\mathbf{Z}\mathbf{Z}^T(m\sigma^2_u) + I\sigma^2_e$
 $\mathbf{A} = \mathbf{Z}\mathbf{Z}^T / m$

Reconciling family studies and GWAS

Family studies: comparing phenotypic similarity to family relatedness

– ***Our method: comparing phenotypic similarity to genetic similarity (estimated from SNPs) in unrelated individuals***

GWAS: testing a SNP at a time in unrelated samples

– ***Our method: Estimating the contribution from all SNPs together***

ANALYSIS

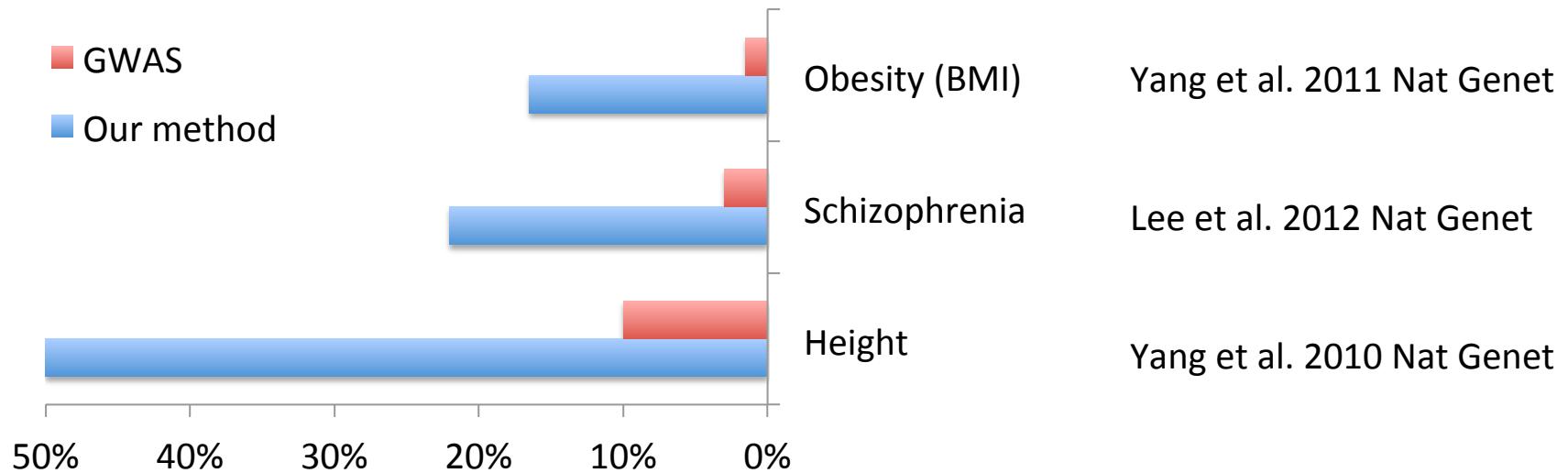
nature
genetics

~50% of variation explained by all SNPs for height vs. ~10% from GWAS

Common SNPs explain a large proportion of the heritability for human height

Jian Yang¹, Beben Benyamin¹, Brian P McEvoy¹, Scott Gordon¹, Anjali K Henders¹, Dale R Nyholt¹, Pamela A Madden², Andrew C Heath², Nicholas G Martin¹, Grant W Montgomery¹, Michael E Goddard³ & Peter M Visscher¹

GWAS vs All-SNP estimation



Many genetic variants each with a small effect contributing to the trait variation

Genome partitioning

- Single component MLM

$$\mathbf{y} = \mathbf{g} + \mathbf{e} \text{ (or } \mathbf{y} = \mathbf{W}\mathbf{u} + \mathbf{e})$$

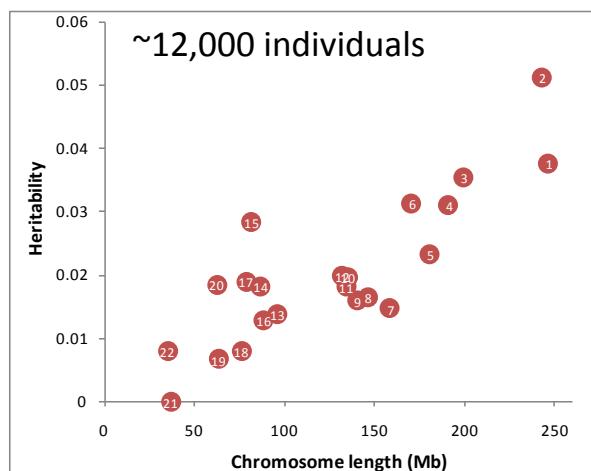
- Multi-component MLM

$$\mathbf{y} = \mathbf{g}_1 + \mathbf{g}_2 + \dots + \mathbf{g}_{22} + \mathbf{e}$$

$$\text{var}(\mathbf{y}) = \mathbf{A}_1\sigma^2_{g1} + \mathbf{A}_2\sigma^2_{g2} + \dots + \mathbf{A}_{22}\sigma^2_{g22} + \mathbf{I}\sigma^2_e$$

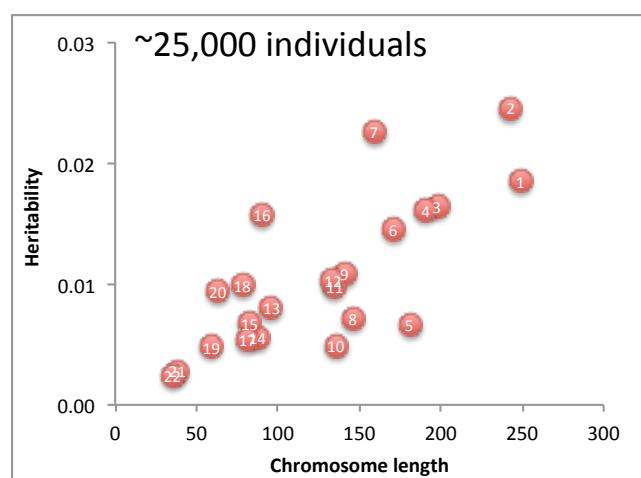
Partitioning the genetic variance into individual chromosomes

Height



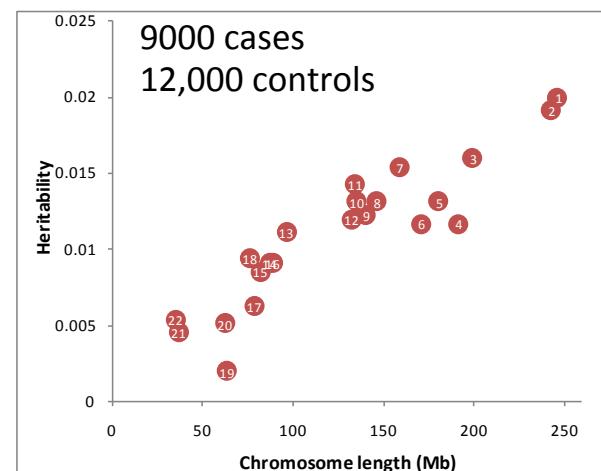
Yang et al. 2011 Nat Genet

BMI



Yang et al. unpublished

Schizophrenia

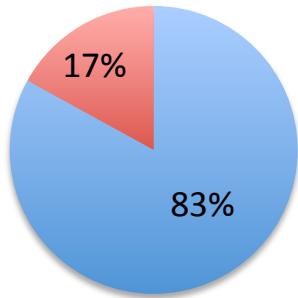


Lee et al. 2012 Nat Genet

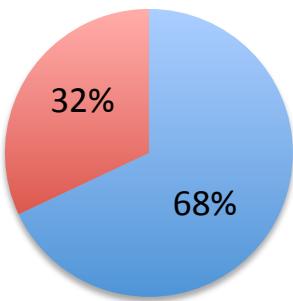
Genetic variants distributed across the whole genome

Partitioning the genetic variance based on functional annotation

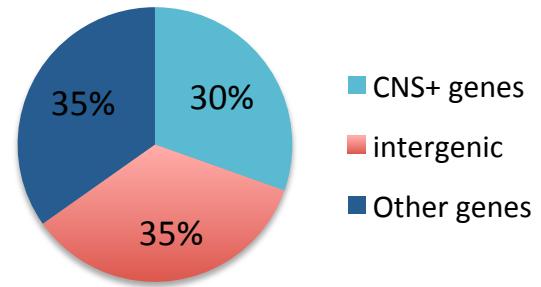
Height



BMI



Schizophrenia



Yang et al. 2011 Nat Genet

Lee et al. 2012 Nat Genet

Genetic signals are enriched in or close to functional genes

More ...

- Bivariate analysis – estimating the genetic correlation between two traits or two diseases using SNP data (Deary et al. 2012 Nature; Lee et al. 2013 Nat Genet)
- Fitting a mixture distribution rather than a single normal distribution to the random effects – e.g. Zhou et al. 2013 PLoS Genet

Linear model (simple regression based association test)

$$y = b_0 + x_1 b_1 + e$$

y = trait value; x_1 = SNP genotype (0, 1 or 2)

$$\hat{b}_1 = \mathbf{X}_1^\top \mathbf{y} / (\mathbf{X}_1^\top \mathbf{X}_1) = \text{cov}(x_1, y) / \text{var}(x_1)$$

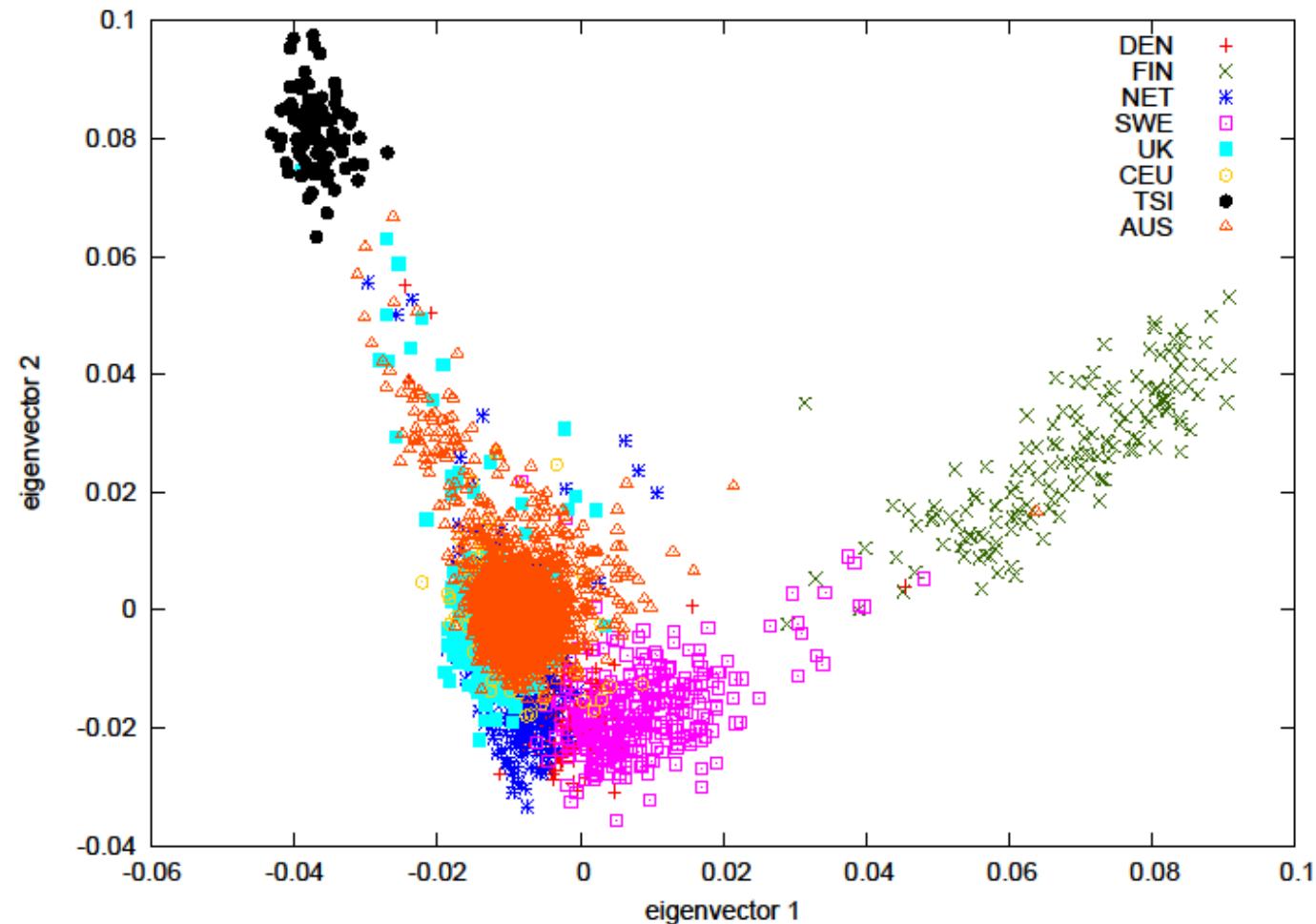
$$\text{SE}(\hat{b}_1) = \sigma_e^2 / [\sqrt{n} \text{ var}(x_1)]$$

Assumption: e is independent and identically distributed

Issues:

- 1) Relatedness: there are relatives in the sample – inflated false positive rate
- 2) Population stratification: individuals of different ancestries – spurious association; e.g. trait = eating with chopsticks, data = a random sample of US population.

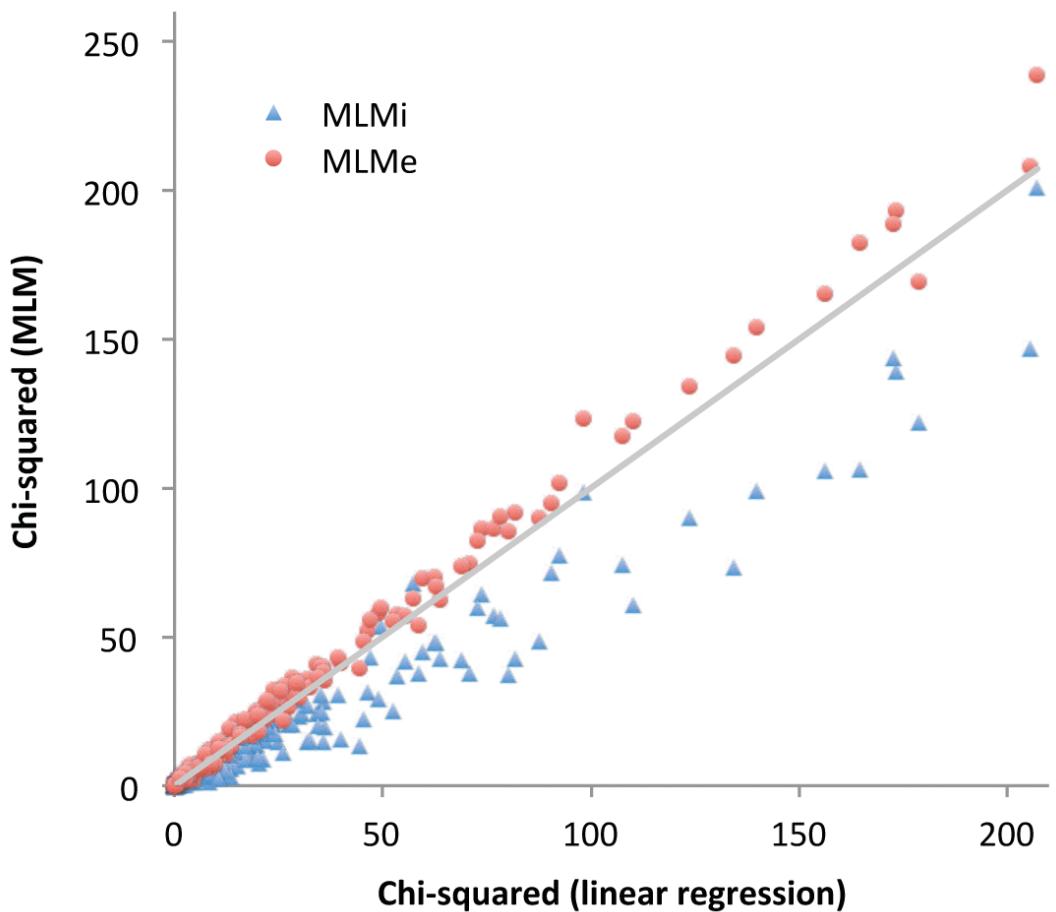
Population stratification estimated from SNP data



Solution: MLM based association analysis

- $\mathbf{y} = \mathbf{X}\mathbf{b} + \mathbf{Z}\mathbf{u} + \mathbf{e}$ or $\mathbf{y} = \mathbf{X}\mathbf{b} + \mathbf{g} + \mathbf{e}$
 $\mathbf{V} = \text{var}(\mathbf{y}) = \mathbf{A}\sigma^2_g + \mathbf{I}\sigma^2_e$
- Testing for fixed effects given sample structure
 $\hat{\mathbf{b}} = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$
 $\text{var}(\hat{\mathbf{b}}) = \sigma^2_e (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1}$
- Issue: a SNP is fitted twice.

Excluding the SNP from calculating the genetic relationship matrix



Software tool

<http://gump.qimr.edu.au/gcta/>

GCTA

a tool for Genome-wide Complex Trait Analysis

[Overview](#)

[Download](#)

[Tutorial](#)

[FAQ](#)

[Options](#)

Overview

GCTA (Genome-wide Complex Trait Analysis) is designed to handle genome-wide association analysis by explaining by genome- or chromosome-wide SNPs for complex traits.

It was developed by [Kang-Mou Lee](#), [Mike Goddard](#) and [Peter Visscher](#) and is maintained by [Peter Visscher](#).

GCTA is currently supported by the Wellcome Trust Medical Research Council. GCTA currently supports the following platforms:

REPORT

GCTA: A Tool for Genome-wide Complex Trait Analysis

Jian Yang,^{1,*} S. Hong Lee,¹ Michael E. Goddard,^{2,3} and Peter M. Visscher¹

Acknowledgements

Complex Traits Genomics Group (UQ)

- Peter Visscher
- Naomi Wray
- Hong Lee

University of Melbourne

- Mike Goddard

QIMR cohort

- Nick Martin
- Grant Montgomery



Australian Government
National Health and
Medical Research Council

N H M R C

GENEVA Consortium

- Teri Manolio
- Bruce Weir

dbGaP



Australian Government
Australian Research Council



**THE UNIVERSITY
OF QUEENSLAND**
AUSTRALIA

The Australian Neurogenetics Conference

**at the Queensland Brain Institute (QBI), The
University of Queensland, on September 11th and
12th, 2014**

<http://web.qbi.uq.edu.au/anc2014/>