THE UNIVERSITY OF QUEENSLAND
AUSTRALIA

CREATE CHANGE

# Evidence for selection using human GWAS data

Naomi R Wray

naomi.wray@uq.edu.au

THE UNIVERSITY OF QUEENSLAND AUSTRALIA | Institute for Molecular Bioscience

THE UNIVERSITY OF QUEENSLAND AUSTRALIA | Queensland Brain Institute
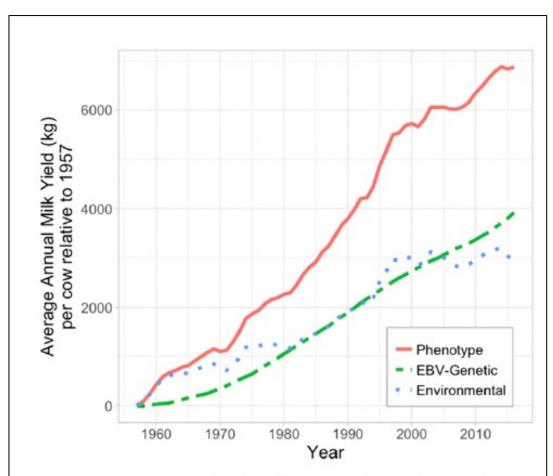
Program in Complex Trait Genomics

@WrayNaomi

PERSPECTIVES

# Polygenic adaptation: a unifying framework to understand positive selection

*Neda Barghi, Joachim Hermisson and Christian Schlötterer*

# Genetic Improvement in Dairy Cattle



**Figure 2** Increase in milk yield in black and white Holstein cattle since 1957. The mean EBV has increased by 3916 or 66 kg per cow per year. The phenotypic and genetic SD of milk yield in 1957 were ~1200 and ~600 kg. Hence, the genetic contribution to milk yield has increased by ~6.5 genetic SD since 1957. Source: Council on Dairy Cattle Breeding (https://queries.uscdcb.com/eval/summary/trend.cfm)

GENETICS | REVIEW

**Complex Trait Prediction from Genome Data: Contrasting EBV in Livestock to PRS in Humans**

Naomi R. Wray,*,†,‡ Kathryn E. Kemper,* Benjamin J. Hayes,‡ Michael E. Goddard,§,** and Peter M. Visscher*,†

The genetic contribution to milk yield has increased 6.5 genetic SD since 1957

70 years; generation interval ~ 5 years

25 million B&W cattle worldwide
Effective population size: ~75

Quant Gen 101

$$A_{child} = \frac{1}{2}A_{dad} + \frac{1}{2}A_{mum} + A_{seg}$$

$$V(A_{seg}) = \frac{1}{2}V(A)$$

Journal of
Animal Breeding and Genetics

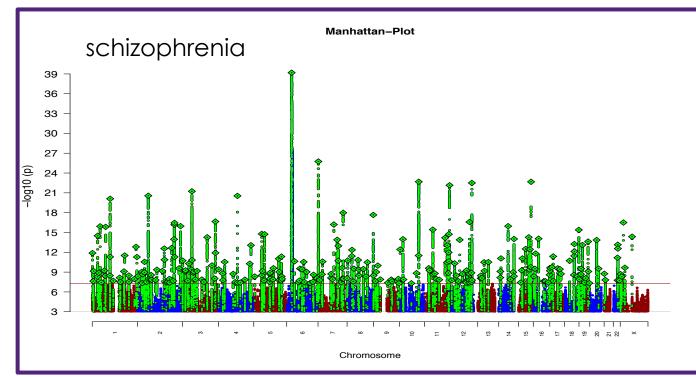J. Anim. Breed. Genet. ISSN 0931-2668

ORIGINAL ARTICLE

**Can more be learned from selection experiments of value in animal breeding programmes? Or is it time for an obituary?**
W.G. Hill

Half the genetic variation in a population is generated by the sampling of genetic material within families

Bill Hill

schizophrenia

Manhattan−Plot



67K ->248

35K->97

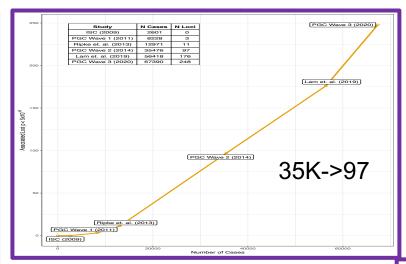| Study | N Cases | N Loci |
|---|---|---|
| ISC (2009) | 2601 | 0 |
| PGC Wave 1 (2011) | 8228 | 3 |
| Ripke et. al. (2013) | 12971 | 11 |
| PGC Wave 2 (2014) | 35476 | 97 |
| Lam et. al. (2019) | 56418 | 176 |
| PGC Wave 3 (2020) | 67390 | 248 |

PGC

248 risk loci identified at genome-wide significance level
We predict thousands are associated

Article

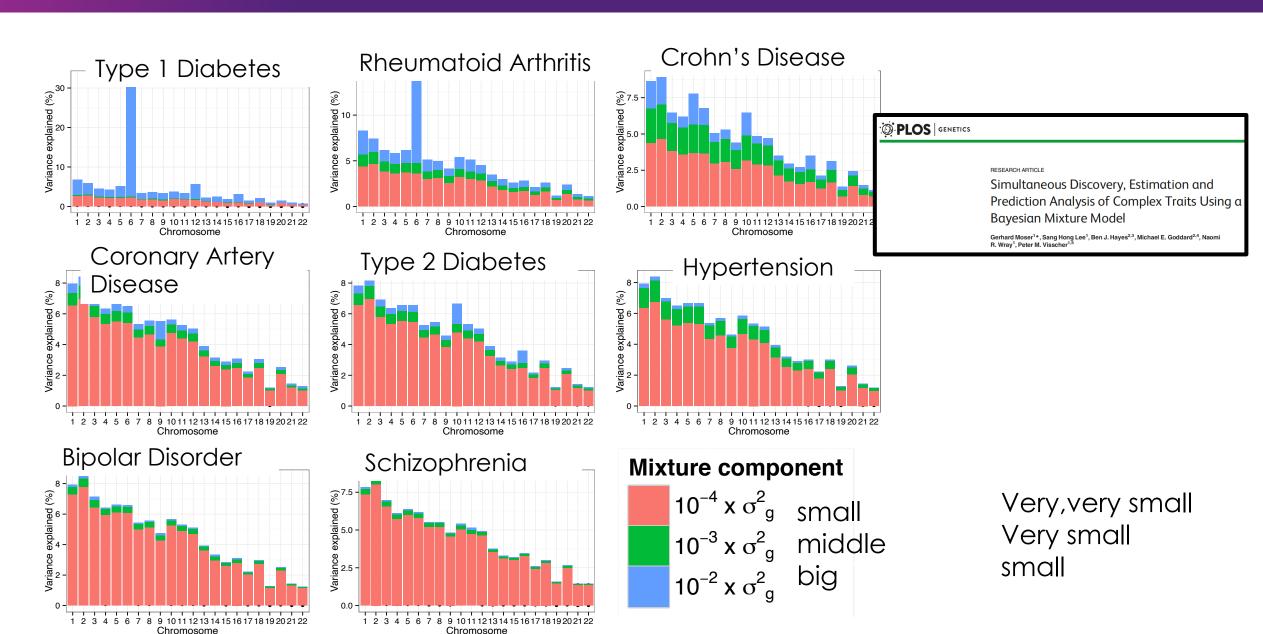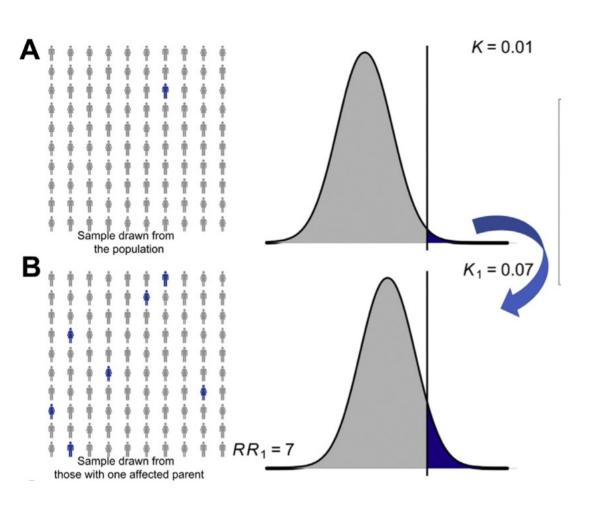**Mapping genomic loci implicates genes and synaptic biology in schizophrenia**

Nature 2022








© Can Stock Photo - csp27761442

5

Type 1 Diabetes

Rheumatoid Arthritis

Crohn's Disease

Coronary Artery Disease

Type 2 Diabetes

Hypertension

Bipolar Disorder

Schizophrenia

**Mixture component**

$10^{-4}$ x $\sigma^2_g$  small

$10^{-3}$ x $\sigma^2_g$  middle

$10^{-2}$ x $\sigma^2_g$  big
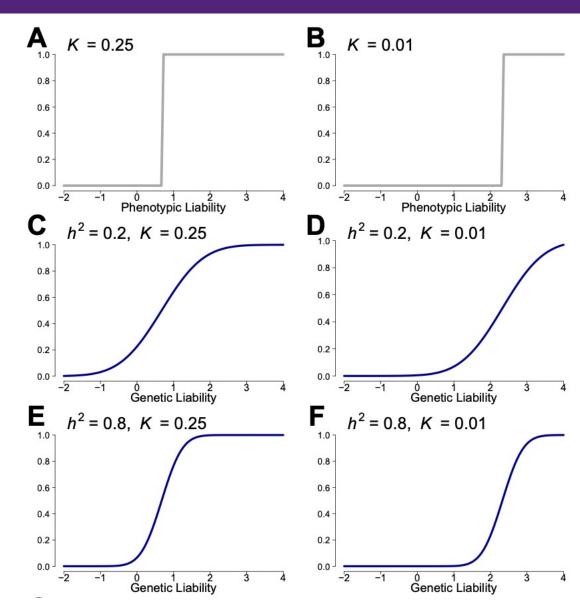
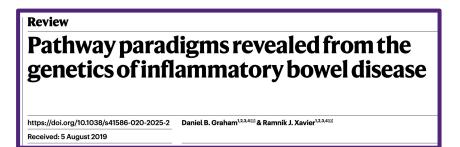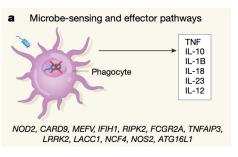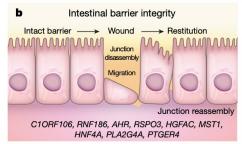Very,very small
Very small
small

THE UNIVERSITY OF QUEENSLAND AUSTRALIA



**Risk in Relatives, Heritability, SNP-Based Heritability, and Genetic Correlations in Psychiatric Disorders: A Review**
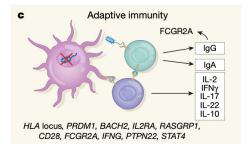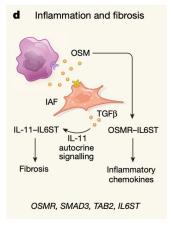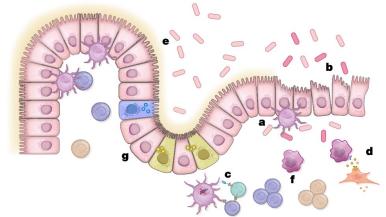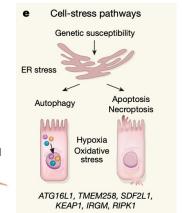
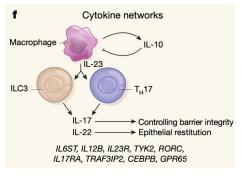Bart M.L. Baselmans, Loïc Yengo, Wouter van Rheenen, and Naomi R. Wray

Robustness

7

# Common diseases are complex

# All traits are polygenic

Lessons from Huntington's Disease



more CAG repeats
earlier age of onset

Number of CAG repeats

Age of onset

Gusella & MacDonald (2009)

9000 HD cases
~10 GWS loci

Much younger age than expected
vs much older age than expected

Article

**CAG Repeat Not Polyglutamine Length Determines Timing of Huntington's Disease Onset**

Genetic Modifiers of Huntington's Disease (GeM-HD) Consortium
Lead contact (James F. Gusella)
Correspondence: gusella@helix.mgh.harvard.edu (James F. Gusella)
https://doi.org/10.1016/j.cell.2019.06.036

Cell

https://www.newswise.com/images/uploads/2011/04/14/retrieve.cfm.jpg

( as expected)



BRCA1 carriers

Kuchenbaecker et al: Evaluation of polygenic risk scores for breast and ovarian cancer risk prediction in BRCA1 and BRCA2 mutation carriers. J Natl Cancer Inst (2017)

# Observational data for selection in humans

**ORIGINAL ARTICLE**

**ONLINE FIRST**

## Fecundity of Patients With Schizophrenia, Autism, Bipolar Disorder, Depression, Anorexia Nervosa, or Substance Abuse vs Their Unaffected Siblings

Robert A. Power, BSc; Simon Kyaga, MD; Rudolf Uher, MD, PhD, MRCPsych; James H. MacCabe, PhD, MRCPsych; Niklas Långström, MD, PhD; Mikael Landen, MD, PhD; Peter McGuffin, FRCP, FRCPsych, PhD; Cathryn M. Lewis, PhD; Paul Lichtenstein, PhD; Anna C. Svensson, PhD



**Figure 1.** Fertility ratios for individuals with schizophrenia, autism, bipolar disorder, depression, anorexia nervosa, and substance abuse. A fertility ratio of 1 (highlighted) represents that of the general population.



**Figure 2.** Fertility ratios for unaffected siblings of individuals with schizophrenia, autism, bipolar disorder, depression, anorexia nervosa, and substance abuse. A fertility ratio of 1 (highlighted) represents that of the general population.

ARTICLE

doi:10.1038/nature21039

Rare and low-frequency coding variants alter human adult height

A full list of authors and affiliations appears in the online version of the paper.

Marouli et al. 2017 (Nature)

Article

Rare coding variants in ten genes confer substantial risk for schizophrenia

Singh et al. 2022 (Nature)

# Signatures of negative selection

## Theoretical prediction



$$effect\ size = \pm\ (selection\ coefficient)^{\tau}$$

Minor allele frequency (MAF)

Eyre-Walker 2010 PNAS
Visscher et al 2012 Mol Psych

## 23 out of 28 traits in UK Biobank (max n=120k) have significant signatures of negative selection.



BayesS (Zeng et al. 2018 Nat Genet)

Minor allele frequency (MAF)

Slide credit: Jian Zeng

nature genetics
ARTICLES
https://doi.org/10.1038/s41588-018-0059-2

**Common schizophrenia alleles are enriched in mutation-intolerant genes and in regions under strong background selection**
2018

Antonio F. Pardiñas[1], Peter Holmans[1], Andrew J. Pocklington[1], Valentina Escott-Price[1],

nature genetics
ANALYSIS
https://doi.org/10.1038/s41588-018-0231-8

**Functional architecture of low-frequency variants highlights strength of negative selection across coding and non-coding annotations**
2019

Steven Gazal[1,2]*, Po-Ru Loh[2,3], Hilary K. Finucane[2,4], Andrea Ganna[2,5,6], Armin Schoech[1,2,7], Shamil Sunyaev[2,3,8] and Alkes L. Price[1,2,7]*

**Research**

# Genome-wide signals of positive selection in human evolution
2014

David Enard,[1] Philipp W. Messer, and Dmitri A. Petrov[1]
*Department of Biology, Stanford University, Stanford, California 94305, USA*

**ARTICLE**

**Extreme Polygenicity of Complex Traits Is Explained by Negative Selection**
2019

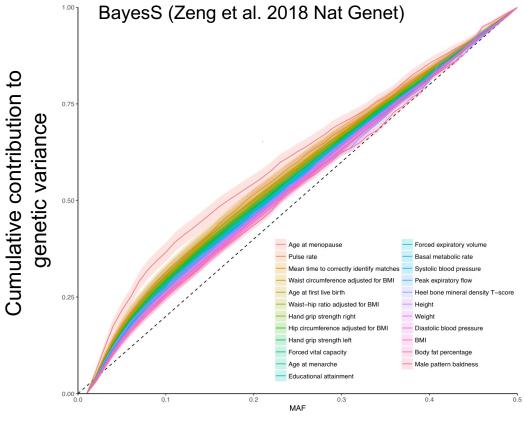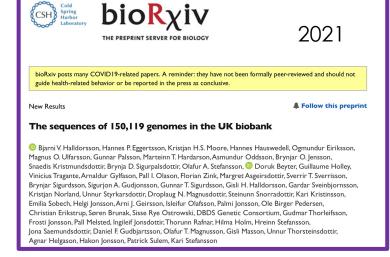Luke J. O'Connor,[1,2,]* Armin P. Schoech,[1] Farhad Hormozdiari,[1] Steven Gazal,[1] Nick Patterson,[3] and Alkes L. Price[1,3,]*

nature genetics
PERSPECTIVE
https://doi.org/10.1038/s41588-019-0383-1

**Measuring intolerance to mutation in human genetics**
2019

Zachary L. Fuller[1]*, Jeremy J. Berg[1], Hakhamanesh Mostafavi[1], Guy Sella[1,2,3] and Molly Przeworski[1,2,3]

PLOS | BIOLOGY
2018

RESEARCH ARTICLE
A population genetic interpretation of GWAS findings for human quantitative traits

Yuval B. Simons[1]*, Kevin Bullaughey[2], Richard R. Hudson[2], Guy Sella[1]*

nature COMMUNICATIONS

ARTICLE
https://doi.org/10.1038/s41467-019-08424-6    OPEN

Quantification of frequency-dependent genetic architectures in 25 UK Biobank traits reveals action of negative selection
2019

Armin P. Schoech[1,2,3], Daniel M. Jordan[4], Po-Ru Loh[3,5], Steven Gazal[1,3], Luke J. O'Connor[1,2,3], Daniel J. Balick[5,6], Pier F. Palamara[7], Hilary K. Finucane[3], Shamil R. Sunyaev[3,5,6] & Alkes L. Price[1,2,3]

CSH Cold Spring Harbor Laboratory    bioRxiv
THE PREPRINT SERVER FOR BIOLOGY
2021

bioRxiv posts many COVID19-related papers. A reminder: they have not been formally peer-reviewed and should not guide health-related behavior or be reported in the press as conclusive.

New Results    Follow this preprint

**The sequences of 150,119 genomes in the UK biobank**

Bjarni V. Halldorsson, Hannes P. Eggertsson, Kristjan H.S. Moore, Hannes Hauswedell, Ogmundur Eiriksson, Magnus O. Ulfarsson, Gunnar Palsson, Marteinn T. Hardarson, Asmundur Oddsson, Brynjar O. Jensson, Snaedis Kristmundsdottir, Brynja D. Sigurpalsdottir, Olafur A. Stefansson, Doruk Beyter, Guillaume Holley, Vinicius Tragante, Arnaldur Gylfason, Pall I. Olason, Florian Zink, Margret Asgeirsdottir, Sverrir T. Sverrisson, Brynjar Sigurdsson, Sigurjon A. Gudjonsson, Gunnar T. Sigurdsson, Gisli H. Halldorsson, Gardar Sveinbjornsson, Kristjan Norland, Unnur Styrkarsdottir, Droplaug N. Magnusdottir, Steinunn Snorradottir, Kari Kristinsson, Emilia Sobech, Helgi Jonsson, Arni J. Geirsson, Isleifur Olafsson, Palmi Jonsson, Ole Birger Pedersen, Christian Erikstrup, Søren Brunak, Sisse Rye Ostrowski, DBDS Genetic Consortium, Gudmar Thorleifsson, Frosti Jonsson, Pall Melsted, Ingileif Jonsdottir, Thorunn Rafnar, Hilma Holm, Hreinn Stefansson, Jona Saemundsdottir, Daniel F. Gudbjartsson, Olafur T. Magnusson, Gisli Masson, Unnur Thorsteinsdottir, Agnar Helgason, Hakon Jonsson, Patrick Sulem, Kari Stefansson

# Bayesian random regression : BayesS

- Bayesian random regression (BayesS)

$$y = \mathbf{1}\mu + \sum_j X_j \beta_j + e$$

where

$$\beta_j \begin{cases} \sim N\left(0, [2p_j q_j]^{\color{red}S} \sigma_\beta^2\right), & \pi \\ = 0, & 1 - \pi \end{cases}$$

**Polygenicity**

Low MAF large effect — Neutral — High MAF large effect

$-$    **0**    $+$

- SNP-based heritability is estimated based on the variance of genetic values
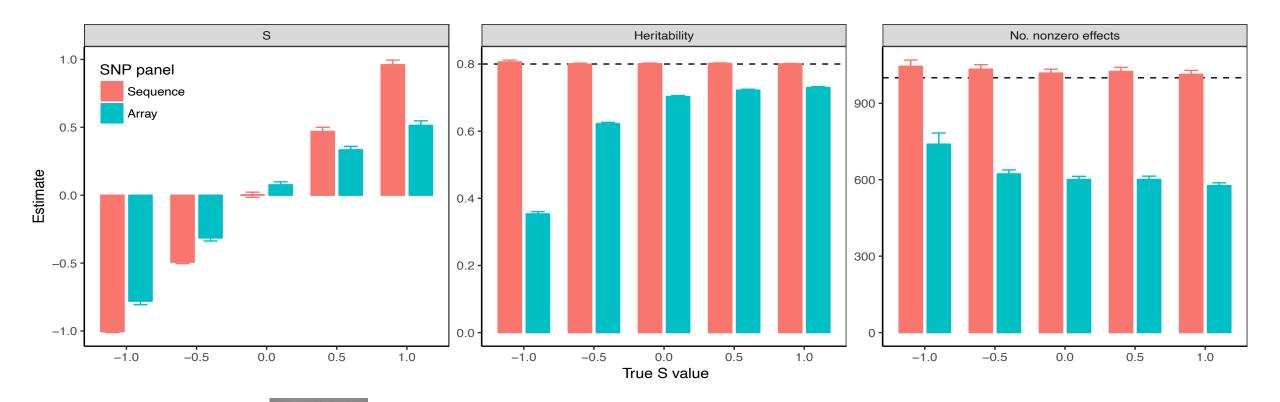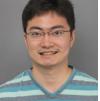- Simultaneously estimate SNP effects and model parameters using MCMC

**Jian Zeng**

**Jian Yang**

# Simulation based on WGS data

- UK10K sequence data, chr 21 & 22, n = 3,642, m ≈ 500k sequence SNPs or 1.5k array SNPs

**Jian Zeng**

THE UNIVERSITY OF QUEENSLAND AUSTRALIA
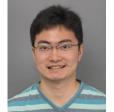
biobank uk

28 traits

- 24 quantitative: anthropometric, cardiovascular, reproductive
- 2 categorical: male pattern baldness (MPB), educational attainment (EA)
- 2 diseases: type 2 diabetes (T2D), major depressive disorder (MDD)

Max N = 126k (unrelated Europeans) for the interim release
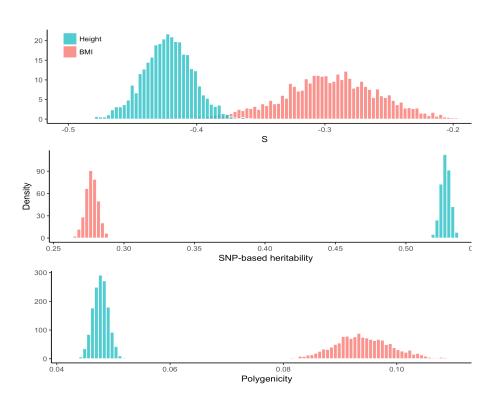
~500k Affymetrix SNPs (MAF > 1%) after QC

Slide credit: Jian Zeng

**Jian Zeng**

ANALYSIS
https://doi.org/10.1038/s41588-018-0101-4

nature genetics

**Signatures of negative selection in the genetic architecture of human complex traits**

Jian Zeng[1], Ronald de Vlaming[2,3], Yang Wu[1], Matthew R. Robinson[1,4], Luke R. Lloyd-Jones[1], Loic Yengo[1], Chloe X. Yap[1], Angli Xue[1], Julia Sidorenko[1,5], Allan F. McRae[1], Joseph E. Powell[1], Grant W. Montgomery[1], Andres Metspalu[5], Tonu Esko[5], Greg Gibson[6], Naomi R. Wray[1,7], Peter M. Visscher[1,7] and Jian Yang[1,7]*
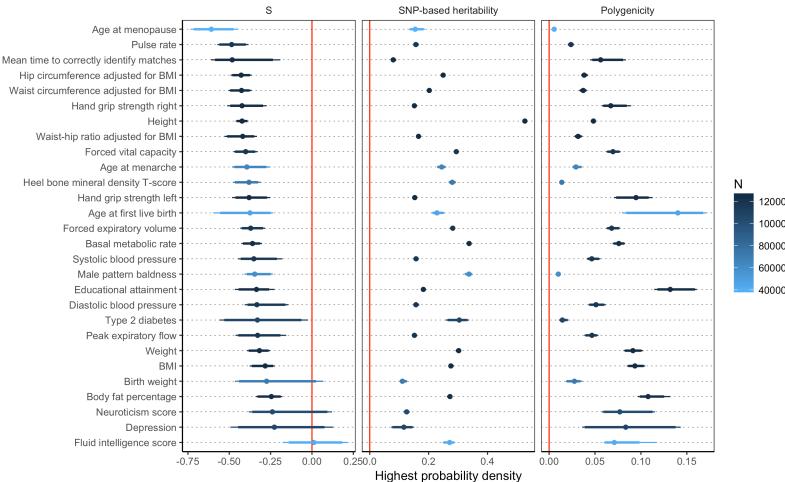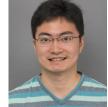
# UK Biobank analysis

## Estimated genetic architecture: height vs. BMI

## Genetic architecture of 28 traits
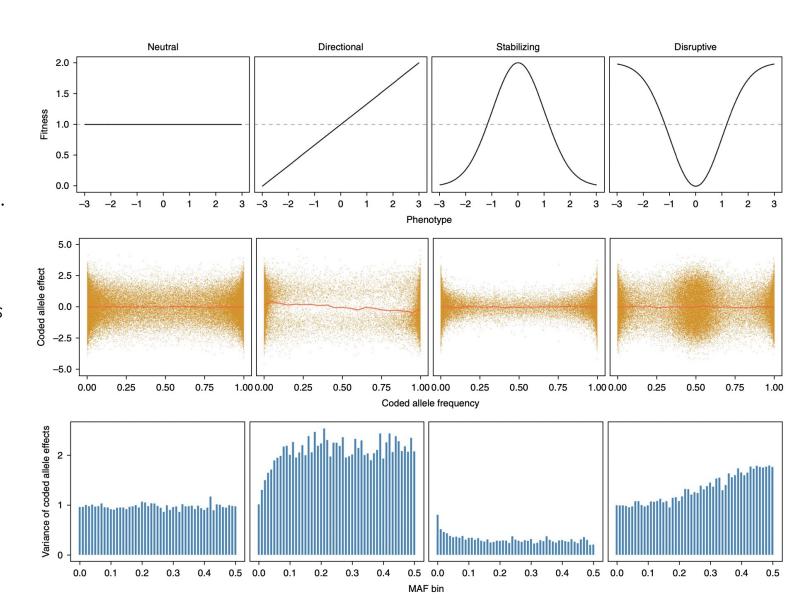


On average 6% of SNPs explain 22% of phenotypic variance

Slide credit: Jian Zeng

Jian Zeng

18

Use SLIM v2.3 to investigate selection models:

- *10-Mb region*
- *mutation rate $1.65 \times 10^{-8}$*
- *new mutations probability*
  - ➢ *0.95 neutral*
  - ➢ *0.05 causal effect sampled from N(0,1).*
- *Phenotype based on the cumulated genotypic values*
  - ➢ *heritability of 0.1 across all causal variants in the current generation.*
- *Evolution of a population of 1,000 individuals over 10,000 generations*
  - ➢ *(equivalent to 10,000 individuals in a population of 100,000 generations).*
- *Burn-in 5,000 generations*
  - ➢ *phenotype did not affect fitness*
  - ➢ *all variants under neutral variation.*
- *Generation 5,001 on*
  - ➢ *standardized phenotype, with mean 0 and variance 1*
  - ➢ *phenotype related to fitness*
- *200 simulation replicates.*
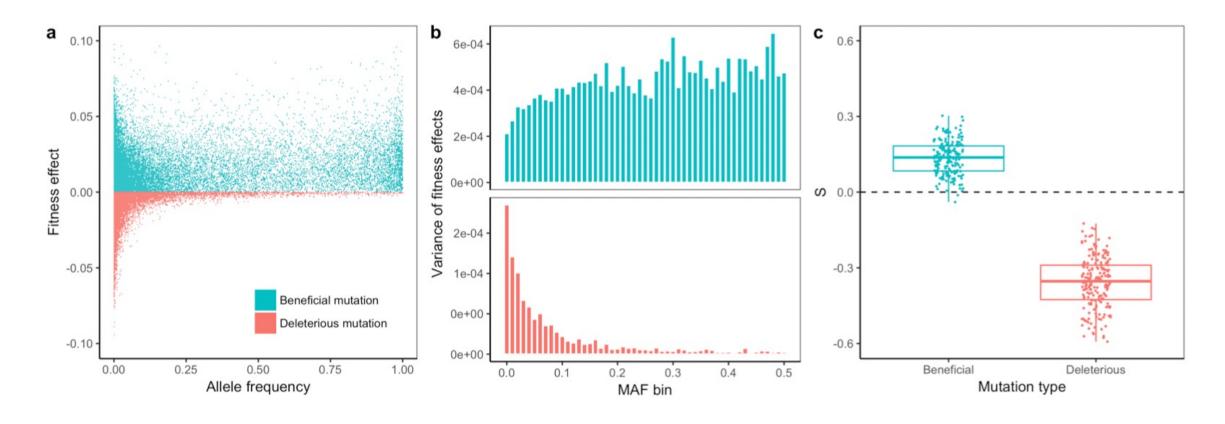- *Results robust to demographic model of bottleneck and expansion*

5% of new mutations beneficial wrt fitness
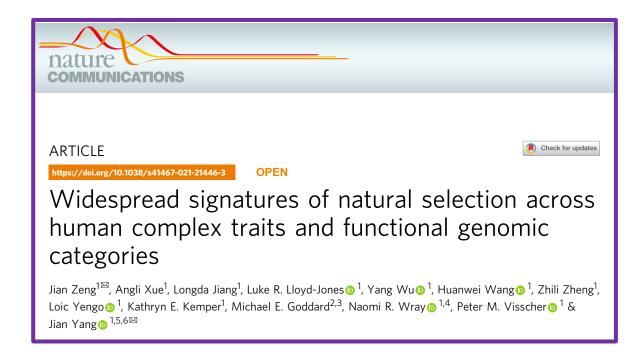OR
5% of new mutation deleterious wrt fitness

THE UNIVERSITY OF QUEENSLAND AUSTRALIA

nature COMMUNICATIONS

ARTICLE

https://doi.org/10.1038/s41467-021-21446-3 OPEN

Check for updates

Widespread signatures of natural selection across human complex traits and functional genomic categories

Jian Zeng[1✉], Angli Xue[1], Longda Jiang[1], Luke R. Lloyd-Jones [1], Yang Wu [1], Huanwei Wang [1], Zhili Zheng[1], Loic Yengo [1], Kathryn E. Kemper[1], Michael E. Goddard[2,3], Naomi R. Wray [1,4], Peter M. Visscher [1] & Jian Yang [1,5,6✉]

## GWAS summary statistics

SNP ID
Chromosome
Base Pair position
Reference Allele
Frequency of reference allele
Effect size of reference allele
Standard error
Sample size

- Most GWAS are meta-analyses from multiple cohorts
- Summary statistics more easily shared than individual level data
- Computational efficiency

**Jian Zeng**

**Jian Yang**

# Summary-data-based model

Consider an individual-data model with a standardised genotype matrix **X**:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$$

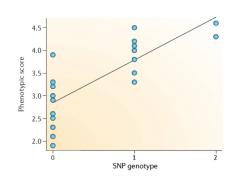Multiply both sides by $\frac{1}{N}\mathbf{X}'$ gives

$$\frac{1}{N}\mathbf{X}'\mathbf{y} = \frac{1}{N}\mathbf{X}'\mathbf{X}\boldsymbol{\beta} + \frac{1}{N}\mathbf{X}'\mathbf{e}$$

$$\mathbf{b} = \mathbf{R}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

$$Var(\boldsymbol{\epsilon}) = \frac{1}{N}\mathbf{R}\sigma_e^2$$

GWAS marginal SNP effects

LD correlation matrix



► Prior for SNP effect:
$$\beta_j \begin{cases} \sim N\left(0, [2p_j q_j]^S \sigma_\beta^2\right), & \pi \\ = 0, & 1 - \pi \end{cases}$$

► $S$ quantifies the relationship between SNP effect size and minor allele frequency (a signature of selection).

► $\pi$ quantifies the proportion of SNPs with nonzero effects (polygenicity).

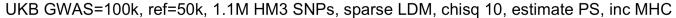► SNP-based heritability:
$$h_{SNP}^2 = \boldsymbol{\beta}'\mathbf{R}\boldsymbol{\beta}/V_P$$

► Implemented in *GCTB* (https://cnsgenomics.com/software/gctb).

Slide credit: Jian Zeng

**Jian Zeng**

# Benchmark SBayesS with BayesS

UKB GWAS=100k, ref=50k, 1.1M HM3 SNPs, sparse LDM, chisq 10, estimate PS, inc MHC
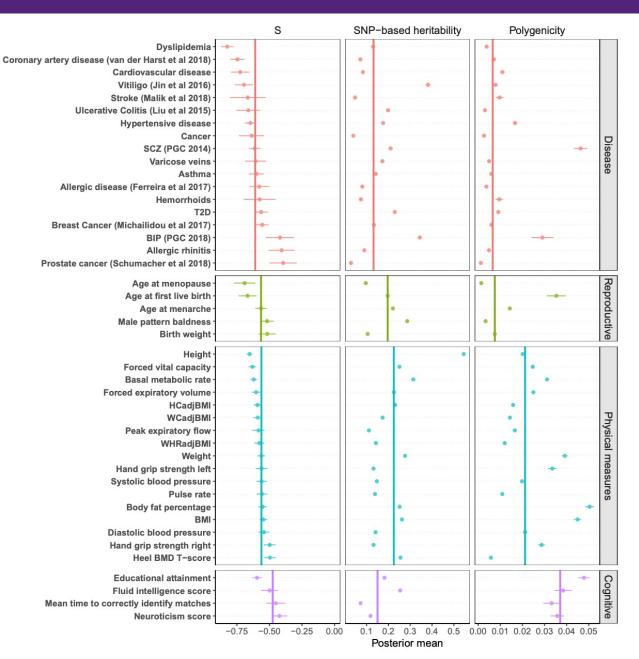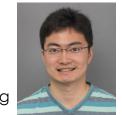
# Genetic architecture of 44 traits

- Full release of UKB + public GWAS summary data (max n = 547k).

- 1.8% of the 1.1 million common HapMap3 SNPs explained 18% of the phenotypic variance.

- The estimate of $S$ was significantly negative ($P<0.001$) in all the traits analysed.

- Median $\hat{S}$ = -0.6 (SD = 0.1).

- Pervasive action of negative selection on the trait-associated variants.

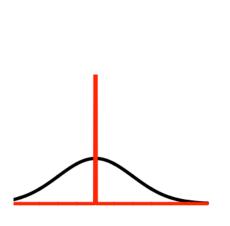- Genetic architecture parameters varied across trait categories.



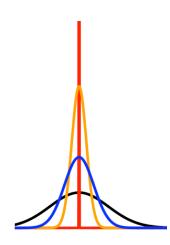Slide credit: Jian Zeng

**Jian Zeng**

Robust to:

- Robust to LD reference sample and LD filtering (but best to make close in ancestry)
- Robust to  Size of LD reference as long as not too small
- Robust to over sampling of cases in case/control studies
- Mostly robust to GWAS sample size (larger sample sizes imply higher polygenicity)
- Mostly robust to modelling of genetic architecture
  - SNP-based heritability is robust
  - S-parameter robust although pattern of differences across traits is changed
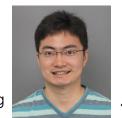  - Polygenicity parameter most sensitive (simulations…)
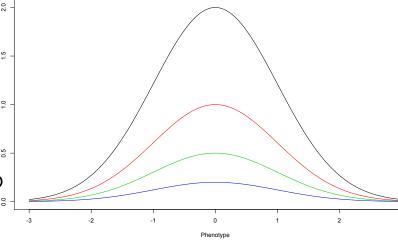
THE UNIVERSITY OF QUEENSLAND
AUSTRALIA

Evolutionary forward simulation (SLiM 3; Haller & Messer 2018)

- 100 MB sequence with a recomb. rate of 1e-5
- Stabilizing selection on phenotypes with different selection strength
  - Selection coefficients sampled from either
    - a normal distribution
    - Mixture of many small and some very large values
- Selection on 10K individuals for 10K generations
- Gravel model for human out-of-Africa evolution (see Ben Haller talk)
- Last generation use two pleiotropic models (Simons or Eyre-Walker) to generate causal effects on focal trait
- Last generation – GWAS on unrelated individuals
  - sample size as UKB
  - SNP density same as real data HapMap 3 SNPs, scaled by genome size
- Estimate S, polygenicity ($\pi$), and SNP-based heritability
- Project real traits into the simulation scenarios

Individual Fitness = $\theta *$ dnorm(Phenotype)/dnorm(0), where $\theta$ is the selection strength



Key parameters varied:
- Mean selection coefficient
- Proportion of mutations that can have a causal effect
- Proportion of phenotypic variance attributed to causal mutations

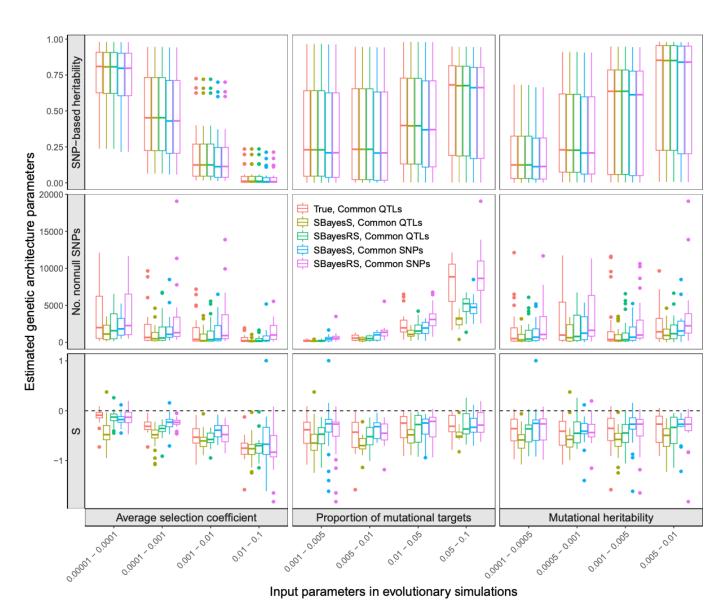Slide credit: Jian Zeng  **Jian Zeng**

28

The results were generally consistent regardless of the use of SNPs

- 36k common SNPs or the actual common causal variants

- Genetic architecture estimation method (SBayesS or SBayesRS)

- simulation model (the Simons et al. or Eyre-Walker model),

- underlying distribution of selection coefficients (mixture or normal distribution)
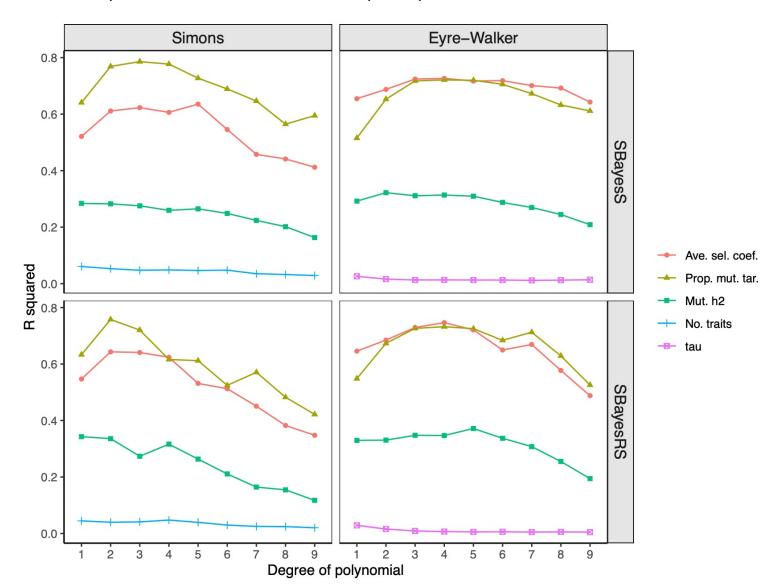
Key differences:

- High s, low $h^2$, low pi, low S

- (large effects are purged)

# Prediction – polynomial regression

Can we predict simulation input parameters from the simulation output estimates?



Input simulation parameters:
- s-bar
- proportion of genome that are mutational targets
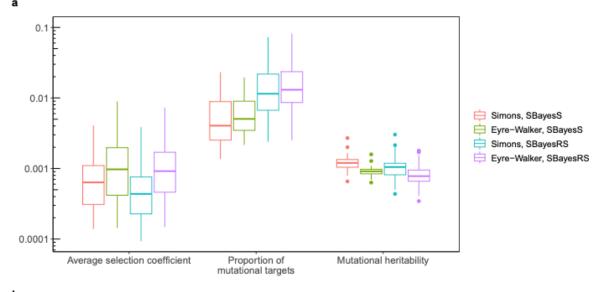- Mutational heritability

Simulation output parameters:
- SNP-based heritability
- Polygenicity parameter
- S coefficient (relationship between allele frequency and effect size

Cross-validation approach
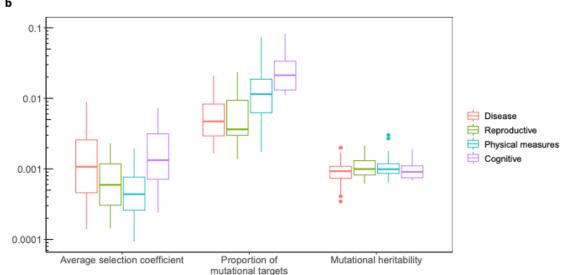
30

Apply polynomial prediction equation to parameter values estimated from real data



Reasonably robust to regression polynomial

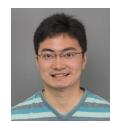Stronger selection on disease and cognitive traits?

The predicted proportion of mutational targets was ~1% on average across traits = ~30 million base pairs of the human genome were mutational targets for a complex trait.

31

# Conclusions

- The action of negative selection is widespread in the genetic architecture of human complex traits
- The strength of negative selection is relatively strong in most traits

- **Interdependence** of underlying evolutionary parameter drive estimated parameters
  - Cannot infer selection strength based solely on S; have to also take SNP-based heritability and polygenicity into account
  - estimated polygenicity π is driven by the mutational target size and selection strength
    - increased average selection coefficient results in decreased estimated π.
    - negative selection removes causal variants of large effects as well as SNPs in LD with them (i.e.,background selection).

- Cognitive trait associated SNPs are under relatively strong selection
- But selection signals detected in the disease associated SNPs are most likely driven by relatively smaller number of mutational targets

- The large estimates of mutational target size per trait implicate widespread pleiotropy across the genome, another study estimated that 90% of GWAS loci affect multiple traits.

**Jian Zeng**          **Jian Yang**

Australian Government
National Health and Medical Research Council
N H M R C

PGC

NIH National Institutes of Health

ISPG INTERNATIONAL SOCIETY OF PSYCHIATRIC GENETICS

FIGHT MND. IT TAKES PEOPLE

AutismCRC

JPND research
EU Joint Programme – Neurodegenerative Disease Research

S.A.L.S.A. Systems Genomics Consortium

MiNDAUS PARTNERSHIP

mnd Research Institute of Australia

BRAIN MEND

Program in Complex Trait Genomics

**The University of Queensland**
Jian Zeng
Jian Yang
Peter Visscher
Luke Lloyd-Jones
Loic Yengo
**University of Melbourne**
Michael Goddard

THE UNIVERSITY OF QUEENSLAND AUSTRALIA | Institute for Molecular Bioscience
THE UNIVERSITY OF QUEENSLAND AUSTRALIA | Queensland Brain Institute