

KSBi-BIML 2024

Bioinformatics & Machine Learning(BIML)
Workshop for Life and Medical Scientists

생명정보학 & 머신러닝 워크숍 (온라인)



Introduction to genome-wide association studies

원홍희 _ 성균관대학교



KSBI
KOREAN SOCIETY FOR
BIOINFORMATICS

| 한국생명정보학회



본 강의 자료는 한국생명정보학회가 주관하는 BIML 2024 워크샵 온라인 수업을 목적으로 제작된 것으로 해당 목적 이외의 다른 용도로 사용할 수 없음을 분명하게 알립니다.

이를 다른 사람과 공유하거나 복제, 배포, 전송할 수 없으며 만약 이러한 사항을 위반할 경우 발생하는 **모든 법적 책임은 전적으로 불법 행위자 본인에게 있음을 경고**합니다.

KSBI-BIML 2024

Bioinformatics & Machine Learning(BIML) Workshop for Life and Medical Scientists

안녕하십니까?

한국생명정보학회가 개최하는 동계 교육 워크숍인 BIML-2024에 여러분을 초대합니다. 생명정보학 분야의 연구자들에게 최신 동향의 데이터 분석기술을 이론과 실습을 겸비해 전달하고자 도입한 전문 교육 프로그램인 BIML 워크숍은 2015년에 시작하여 올해로 벌써 10년 차를 맞이하게 되었습니다. BIML 워크숍은 국내 생명정보학 분야의 최초이자 최고 수준의 교육프로그램으로 크게 인공지능과 생명정보분석 두 개의 분야로 구성되어 있습니다. 올해 인공지능 분야에서는 최근 생명정보 분석에서도 응용이 확대되고 있는 다양한 인공지능 기반 자료모델링 기법들에 대한 현장 강의를 진행될 예정이며, 관련하여 심층학습을 이용한 단백질구조예측, 유전체분석, 신약개발에 대한 이론과 실습 강의를 함께 제공될 예정입니다. 또한 단일세포오믹스, 공간오믹스, 메타오믹스, 그리고 롱리드염기서열 자료 분석에 대한 현장 강의는 많은 연구자의 연구 수월성 확보에 큰 도움을 줄 것으로 기대하고 있습니다.

올해 BIML의 가장 큰 변화는 최근 연구 수요가 급증하고 있는 의료정보자료 분석에 대한 현장 강의를 추가하였다는 것입니다. 특히 의료정보자료 분석을 많이 수행하시는 의과학자 및 의료정보 연구자들께서 본 강좌를 통해 많은 도움을 받으실 수 있기를 기대하고 있습니다. 또한 다양한 생명정보학 분야에 대한 온라인 강좌 프로그램도 점차 증가하고 있는 생명정보 분석기술의 다양화에 발맞추기 위해 작년과 비교해 5강좌 이상을 신규로 추가했습니다. 올해는 무료 강좌 5개를 포함하여 35개 이상의 온라인 강좌가 개설되어 제공되며, 연구 주제에 따른 연관된 강좌 추천 및 강연료 할인 프로그램도 제공되며, 온라인을 통한 Q&A 세션도 마련될 예정입니다. BIML-2024는 국내 주요 연구 중심 대학의 전임 교원이자 각 분야 최고 전문가들의 강의로 구성되었기에 해당 분야의 기초부터 최신 연구 동향까지 포함하는 수준 높은 내용의 강의를 될 것이라 확신합니다.

BIML-2024을 준비하기까지 너무나 많은 수고를 해주신 운영위원회의 정성원, 우현구, 백대현, 김태민, 김준일, 김상우, 장혜식, 박종은 교수님과 KOBIC 이병욱 박사님께 커다란 감사를 드립니다. 마지막으로 부족한 시간에도 불구하고 강의 부탁을 흔쾌히 허락하시고 훌륭한 현장 강의와 온라인 강의를 준비하시는데 노고를 아끼지 않으신 모든 강사분들께 깊은 감사를 드립니다.

2024년 2월

한국생명정보학회장 이 인 석

Introduction to genome-wide association studies

전장유전체연관분석(GWAS, genome-wide association studies)은 인간 질병이나 형질과 연관된 유전 변이를 발굴하고 유전적 조성을 규명하는 대표적인 연구 방법론이다. 그 동안 전세계에서 진행된 대규모 GWAS 연구들은 다양한 형질과 연관된 유전 변이를 발굴하였고 이러한 변이들은 형질의 유전력을 상당 부분 설명하게 되었다. 나아가, 대규모 GWAS 분석 결과(GWAS summary statistics)가 공유됨에 따라, 유전력(heritability), 질병 간 유전적 상관성(genetic correlation), 다인자유전점수(polygenic risk score), 멘델리안 무작위법(Mendelian randomization) 등 여러 post-GWAS 분석이 가능하게 되었고 질병의 유전적 조성을 이해하는데 핵심적인 정보를 제공하고 있다.

본 강의에서는 GWAS를 중심으로 한 유전체 분석의 배경, 이론 및 분석 방법론 등을 소개하고, 복합 질환에서 최근 GWAS 연구 결과를 소개하고자 한다. 이를 통해 GWAS 기반의 연구를 해석하기 위한 기초 지식을 쌓고, 나아가 GWAS 분석 및 GWAS 결과의 응용 연구를 위한 핵심 역량을 갖추는 것을 목표로 한다.

강의는 다음의 내용을 포함한다:

- 유전체 분석을 위한 개념
- GWAS 분석의 이론과 방법론
- Post-GWAS 분석의 이론과 방법론
- 대표적인 연구 결과의 소개

* 참고강의교재:

Tam et al. Benefits and limitations of genome-wide association studies, Nature Reviews Genetics, 20:467-484, 2019.

Balding. A tutorial on statistical methods for population association studies, Nature Reviews Genetics, 7:781-791, 2006.

이종극, 질병 유전체 분석법 3판

* 강의 난이도: 초급

* 강의: 원홍희 교수 (성균관대학교 삼성융합의과학원)

Curriculum Vitae

Speaker Name: Hong-Hee Won, Ph.D.



► Personal Info

Name Hong-Hee Won
Title Associate Professor
Affiliation Sungkyunkwan University

► Contact Information

Address 81, Irwon-Ro, Gangnam-Gu, Seoul, 06351
Email wonhh@skku.edu
Phone Number 010-6326-3452

Research Interest

Population genomics, genome-wide association study, polygenic risk score

Educational Experience

2002 B.S. in Computer Science, Yonsei University, Korea
2004 M.S. in Computer Science, Yonsei University, Korea
2011 Ph.D. in Bioinformatics, KAIST, Korea

Professional Experience

2004-2012 Research Scientist, Samsung Biomedical Research Institute and Samsung Medical Center, Korea
2012-2015 Research Fellow, Massachusetts General Hospital, Harvard Medical School, and Broad Institute of MIT and Harvard, USA
2016-2020 Assistant Professor, Sungkyunkwan University, Samsung Medical Center, Korea
2020- Associate Professor, Sungkyunkwan University, Samsung Medical Center, Korea

Selected Publications (5 maximum)

1. Kim S, et al. Shared genetic architectures of subjective well-being in East Asian and European ancestry populations, *Nature Human Behaviour*, 6(7):1014-1026, 2022.
2. Kim M, et al. Association between adiposity and cardiovascular outcomes: an umbrella review and meta-analysis of observational and Mendelian randomization studies, *European Heart Journal*, 42(34):3388-3403, 2021.
3. Khera AV, et al. Association of rare and common variation in the lipoprotein lipase gene with coronary artery disease, *Journal of the American Medical Association JAMA*, 317(9):937-46, 2017.
4. Do R, et al. Exome sequencing identifies rare LDLR and APOA5 alleles conferring risk for myocardial infarction, *Nature* 518:102-106, 2015.
5. Stitzel NO, et al. Inactivating mutations in NPC1L1 and protection from coronary heart disease, *New England journal of medicine NEJM*, 371(22):2072-2082, 2014.

Introduction to genome-wide association studies

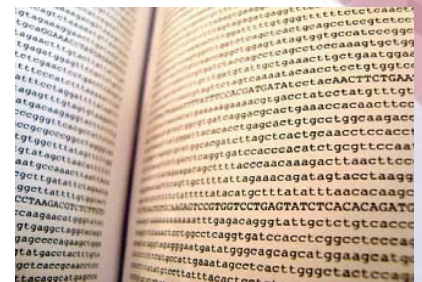
원홍희, Ph.D.

honghee.won@gmail.com

성균관대학교 삼성융합의과학원
삼성서울병원

The Human Genome

- Instruction manual for human cells
- A book with 3.2 billion letters in 23 chapters or chromosomes
- 20,000 genes, exome (1% of the genome)
- 99.9% identical, 4 million letters are different
 - Variation, variant, mutation, polymorphism



Genetic variation affects phenotype

(and risk for disease)

- Genetic variants
 - Pathogenic variants
 - Disease-causing, deleterious, damaging
 - Usually rare (<1%)
 - Often, referred to as “Mutations”
 - Neutral variants
 - Non-disease causing, but may affect disease susceptibility
 - Usually common (>5%)
 - Often, referred to as “Polymorphisms”
 - SNP (single nucleotide polymorphism)

누구나 수백만의 germline 유전 변이를 갖고 있다.

- Single nucleotide variants (SNV)
 - 단일 염기 변이 : 4백만개/사람
- Multi-nucleotide variants
 - Small insertions/deletions (indels) : 50만개/사람
 - Large copy number variants (CNVs)
 - Inversions
 - Translocations
 - Aneuploidy

일생동안
변하지않는다

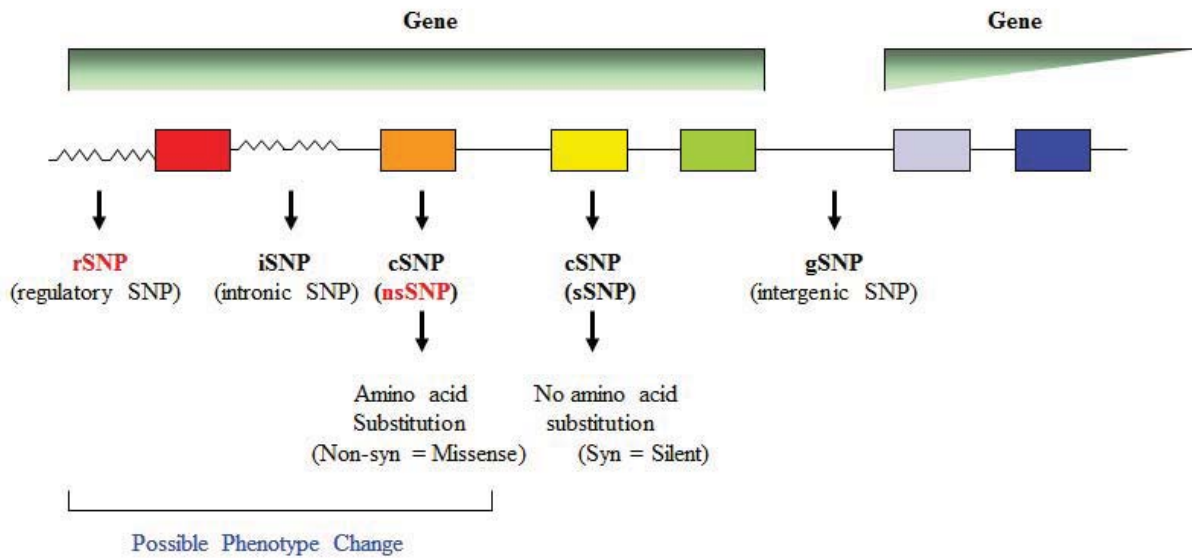
```
AAATAGCACCGTTAGC  
AAATAGCCCCGTTAGC
```

SNV 예

```
AAATAGCACCGTTAGC  
AAATA-----GTTAGC
```

indels 예

Classification of SNP by Location



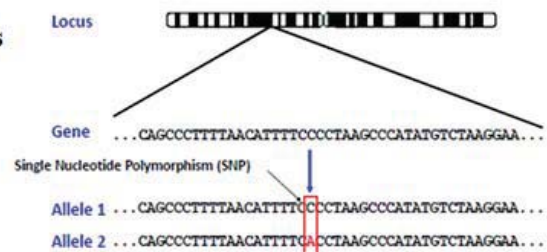
질병유전체분석법3

Definitions: Genes and Genotypes

Gene: A unit of inheritance that is transmitted from parents to offspring, or A region of DNA that codes for a specific product (protein or RNA)

Locus (loci): The place where a particular gene resides on a chromosome
Gene \approx Locus (interchangeable)

Allele: Different forms of a gene, or Variants of DNA at a given locus



Genotype: The diploid (pair) of alleles (except the X chromosome in males)

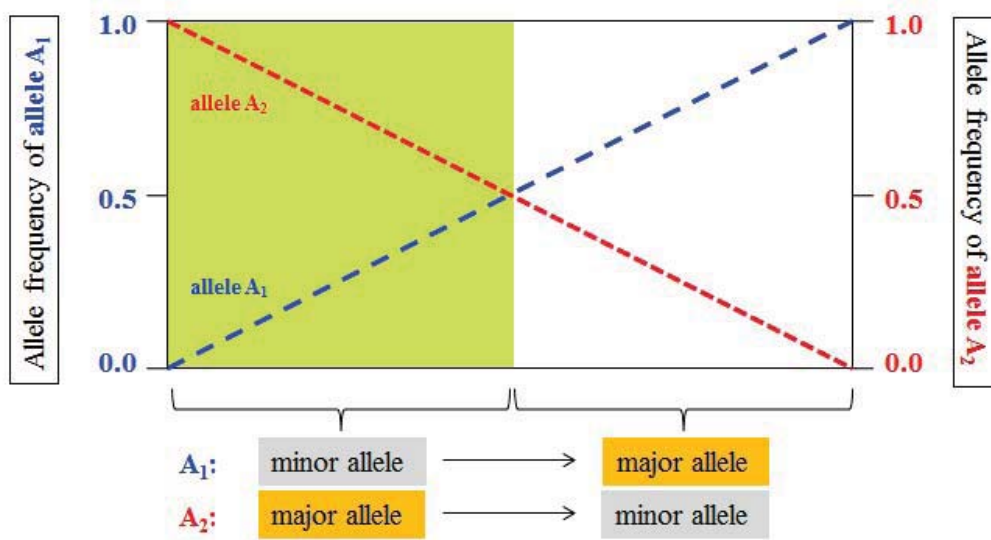
In diploid individual, homozygote (wild type)= 2 identical (wild type) alleles
heterozygote (wild + mutant type)= 2 different alleles
homozygote (mutant type)= 2 identical (mutant type) alleles

In haploid (egg or sperm cells), cells have only one copy of a gene

질병유전체분석법3

Minor Allele Frequency (MAF)

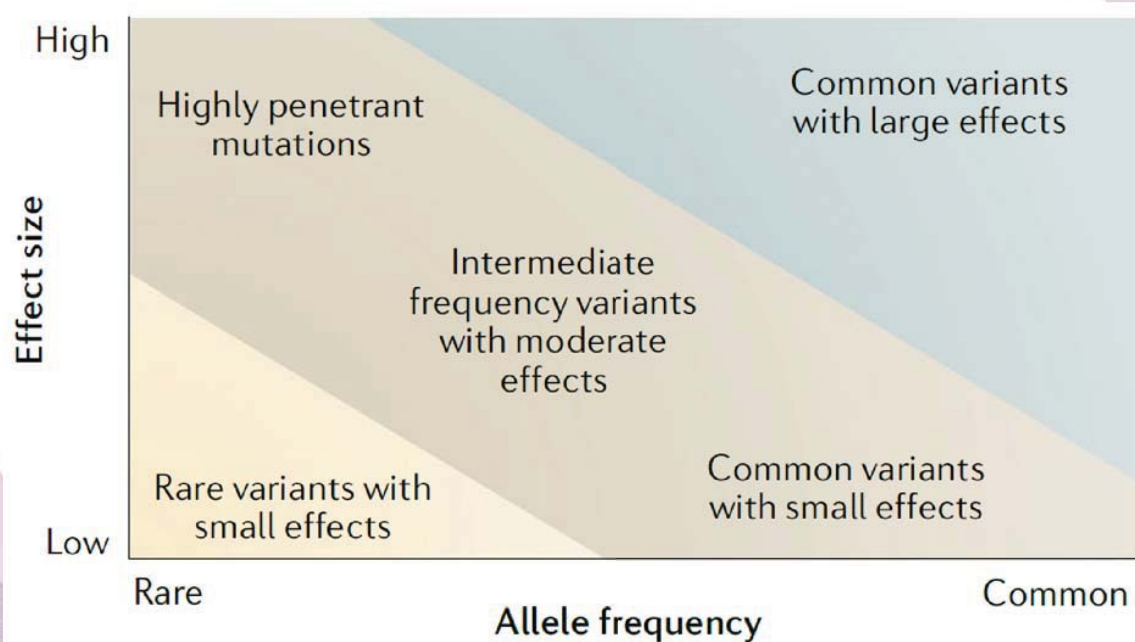
SNP : bi-allele (A_1/A_2) DNA marker



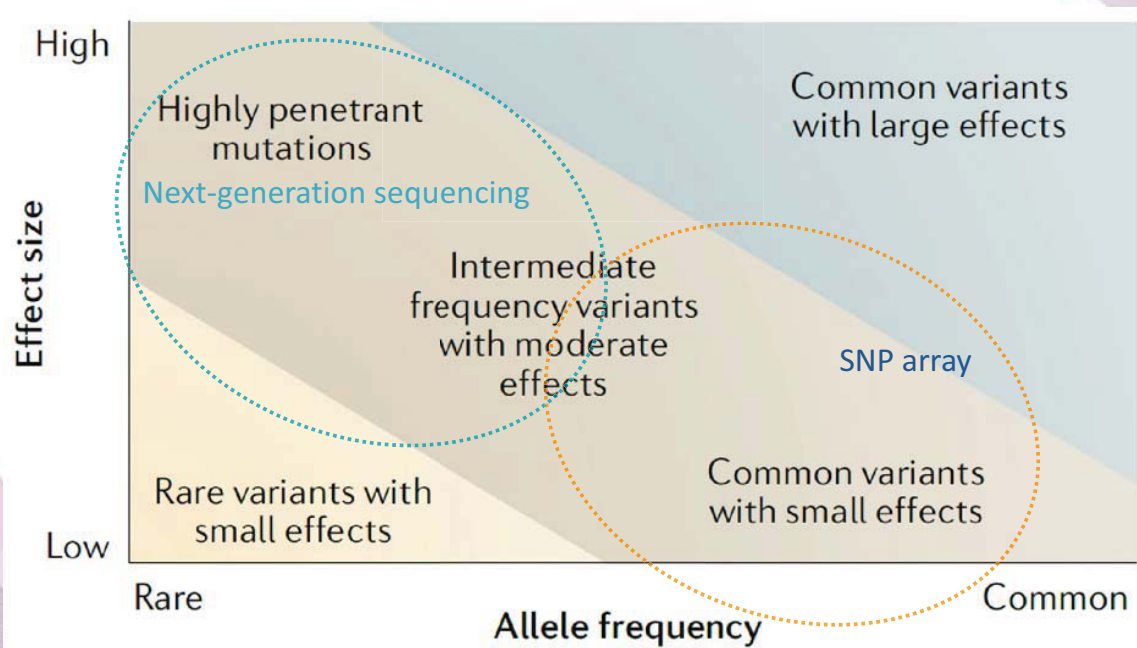
- **Minor Allele Frequency (MAF)** = a standard index for genetic diversity of a SNP marker.
- * Maximum MAF = 0.5

질병유전체분석법3

Variants by frequency and effect size



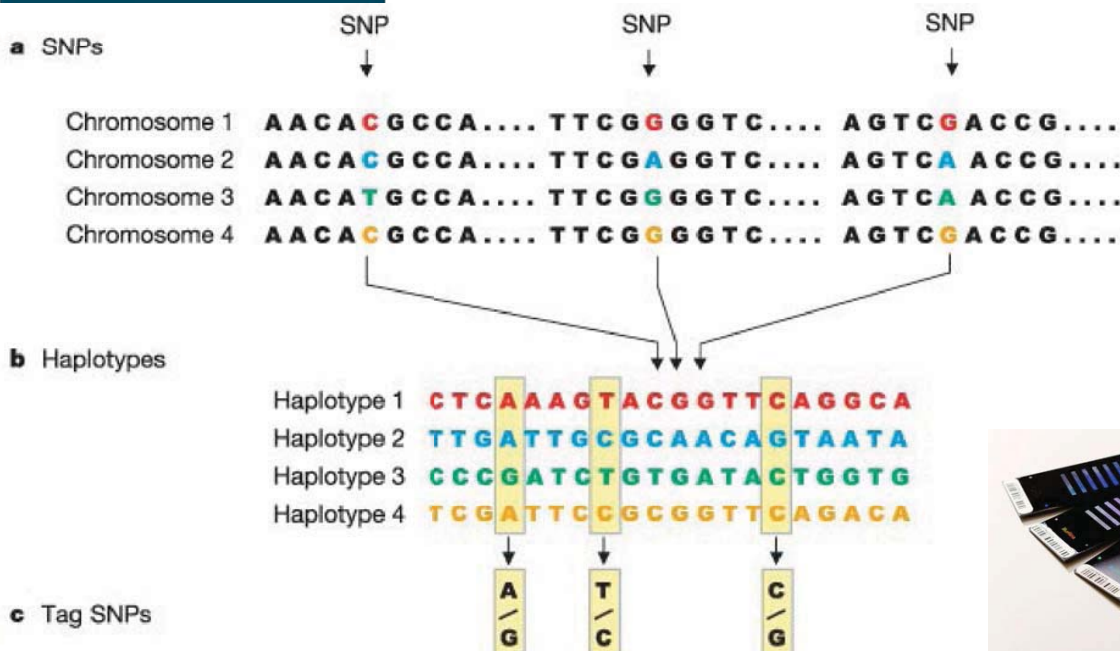
Variants by frequency and effect size



9

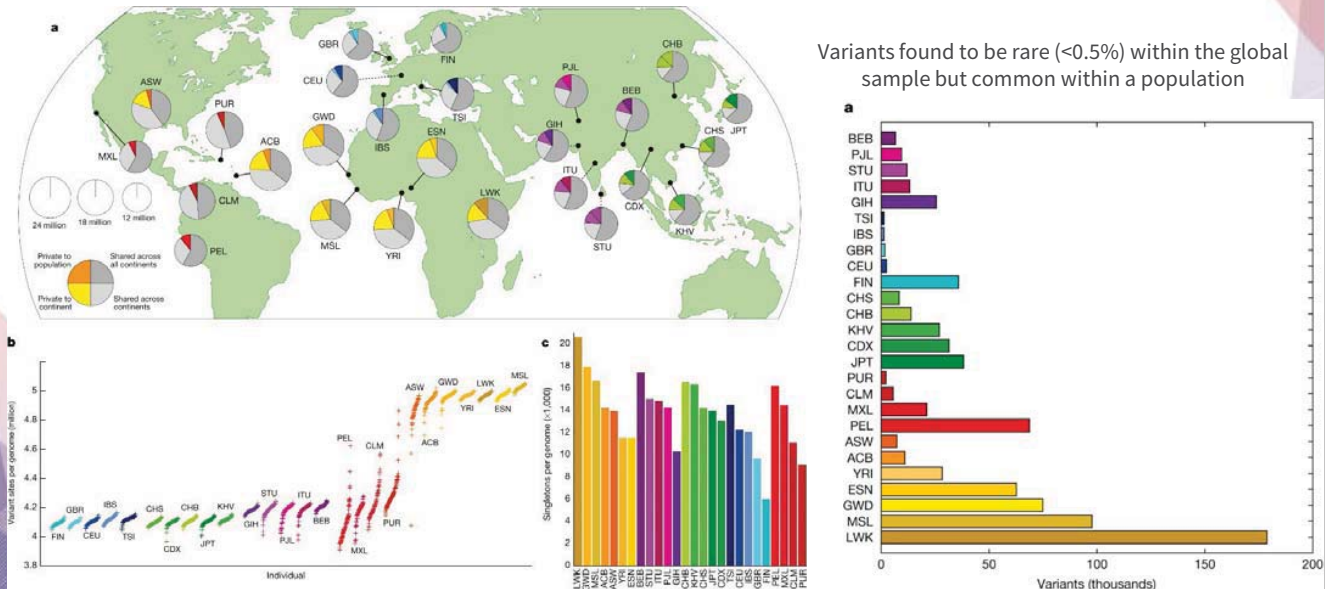


Int. HapMap Project (2002-2009)



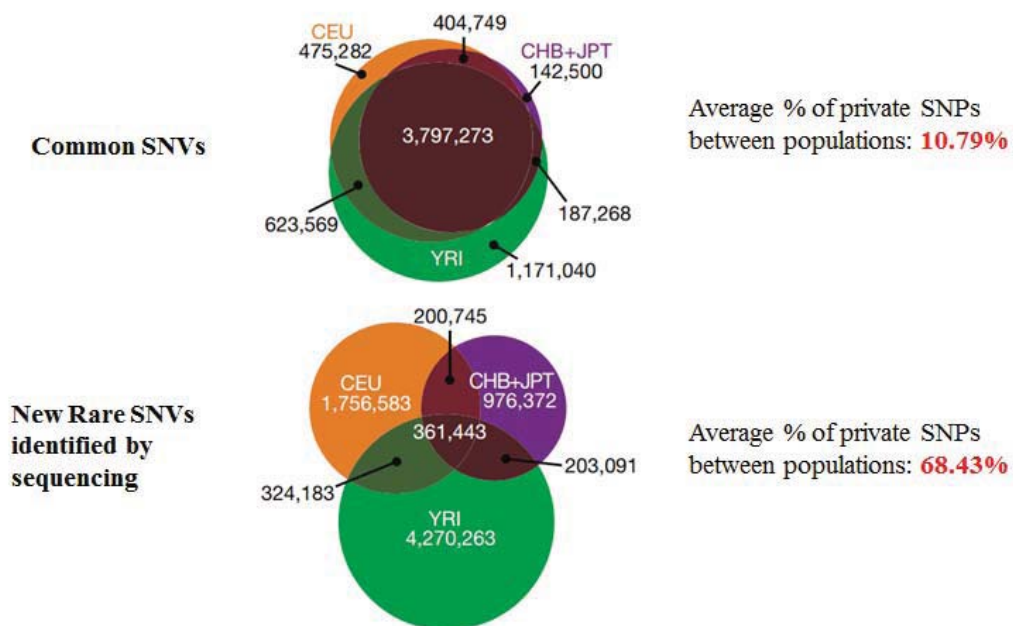
1000 Genomes Project

- Sequenced the genomes of 2,504 individuals from 26 populations in Africa (AFR), East Asia (EAS), Europe (EUR), South Asia (SAS), and the Americas (AMR)



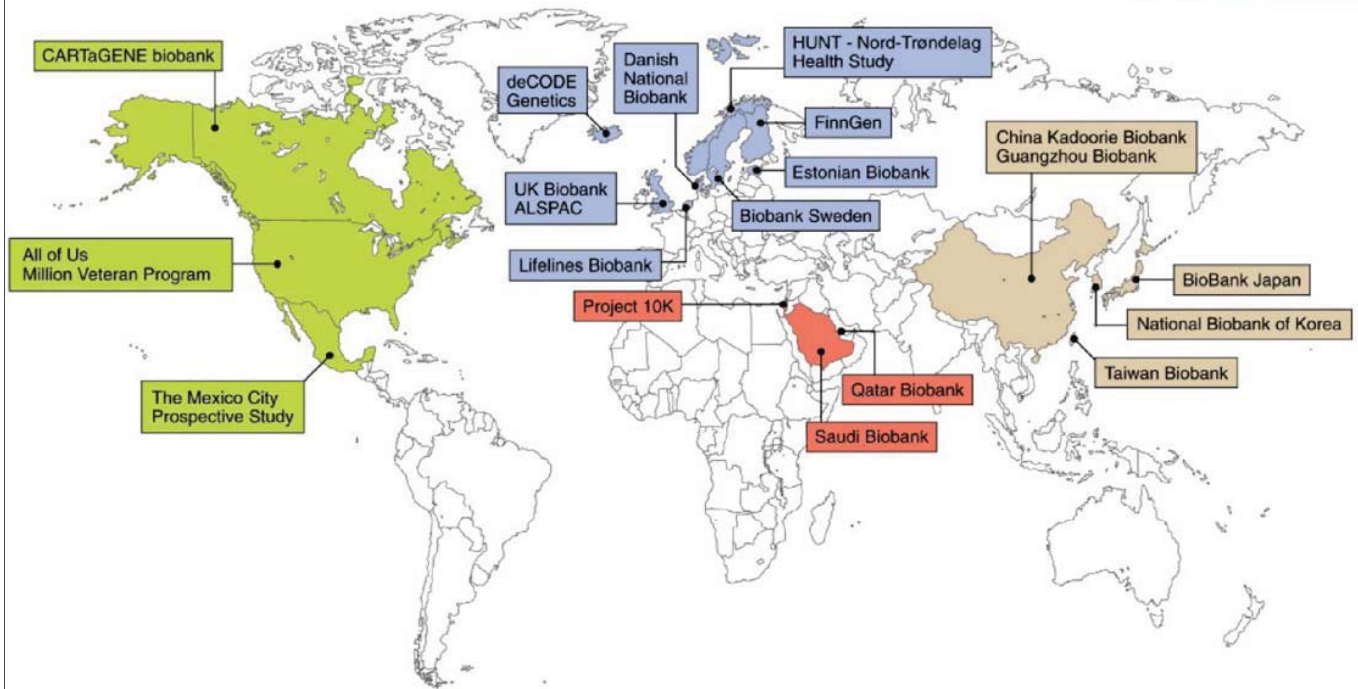
A Auton et al. *Nature* 526, 68-74 (2015) doi:10.1038/nature15393

Rare SNVs are Much More Likely to be Population-Specific



The 1000 genomes project consortium. *Nature* 467:1061-1074 (2010)

Global efforts made for genomic data – biobank



Nature Medicine 26, 29–38 (2020)

Genome-wide association study (GWAS)

Summary of GWAS analysis and tools

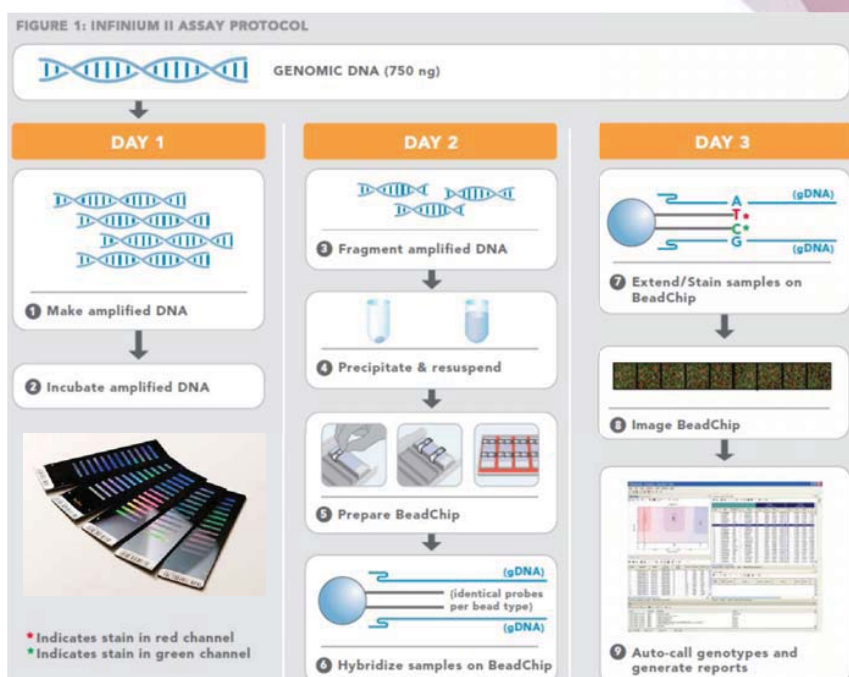
- Quality control
 - Sample QC: PLINK
 - Variant QC: PLINK
 - Related samples: KING
 - PCA of genetic ancestry: EIGENSTRAT(smartpca)
- Imputation
 - Haplotype Reference Consortium: Michigan Imputation Server
 - TOPMed Imputation Server
- Association analysis
 - Logistic/linear regression (unrelated): PLINK
 - Mixed effects regression (including related): SAIGE, BOLT-LMM, REGENIE
- Visualization
 - QQ plot: CM-PLOT
 - Manhattan plot: CM-PLOT

SNP arrays provide fast and accurate genotyping of about a million of genetic variants

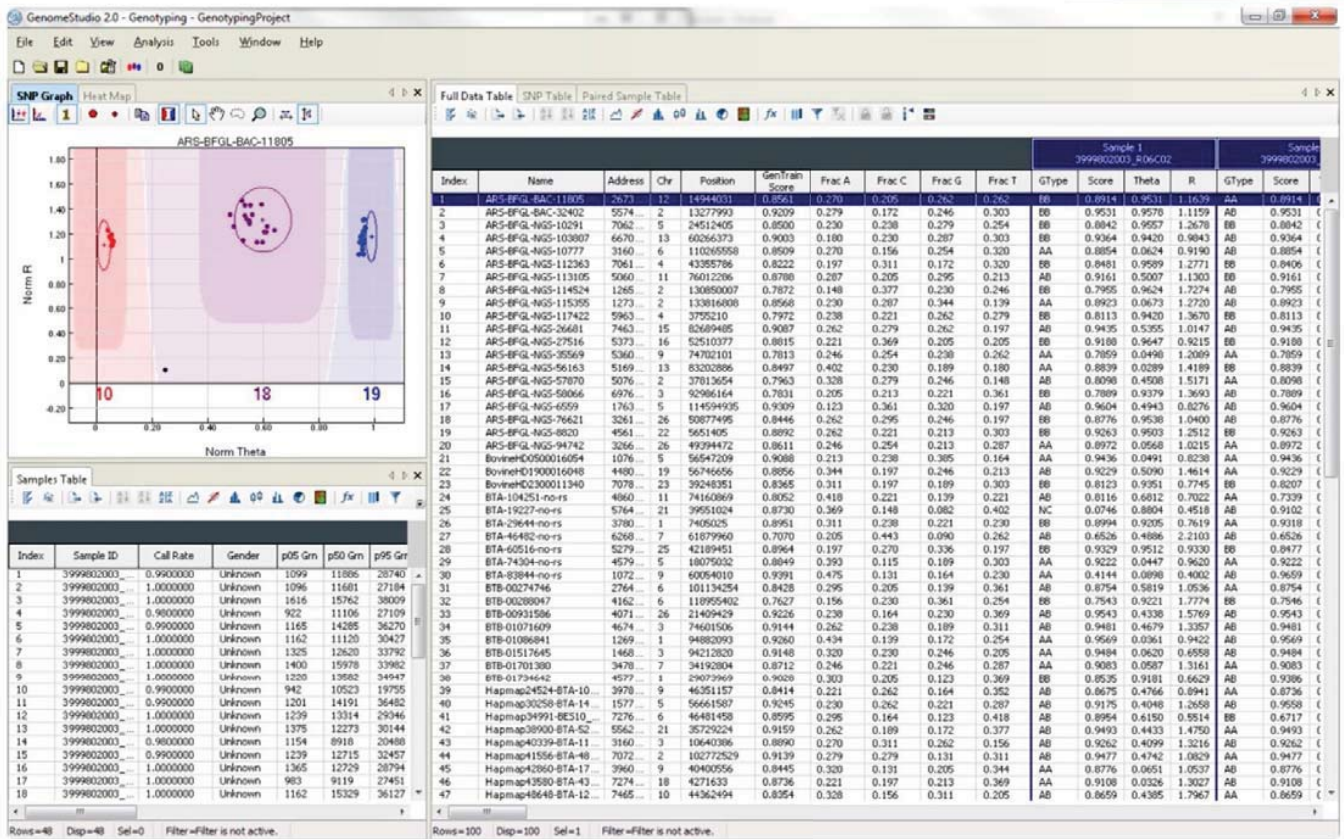


Thermo Fisher (Affymetrix)

Illumina

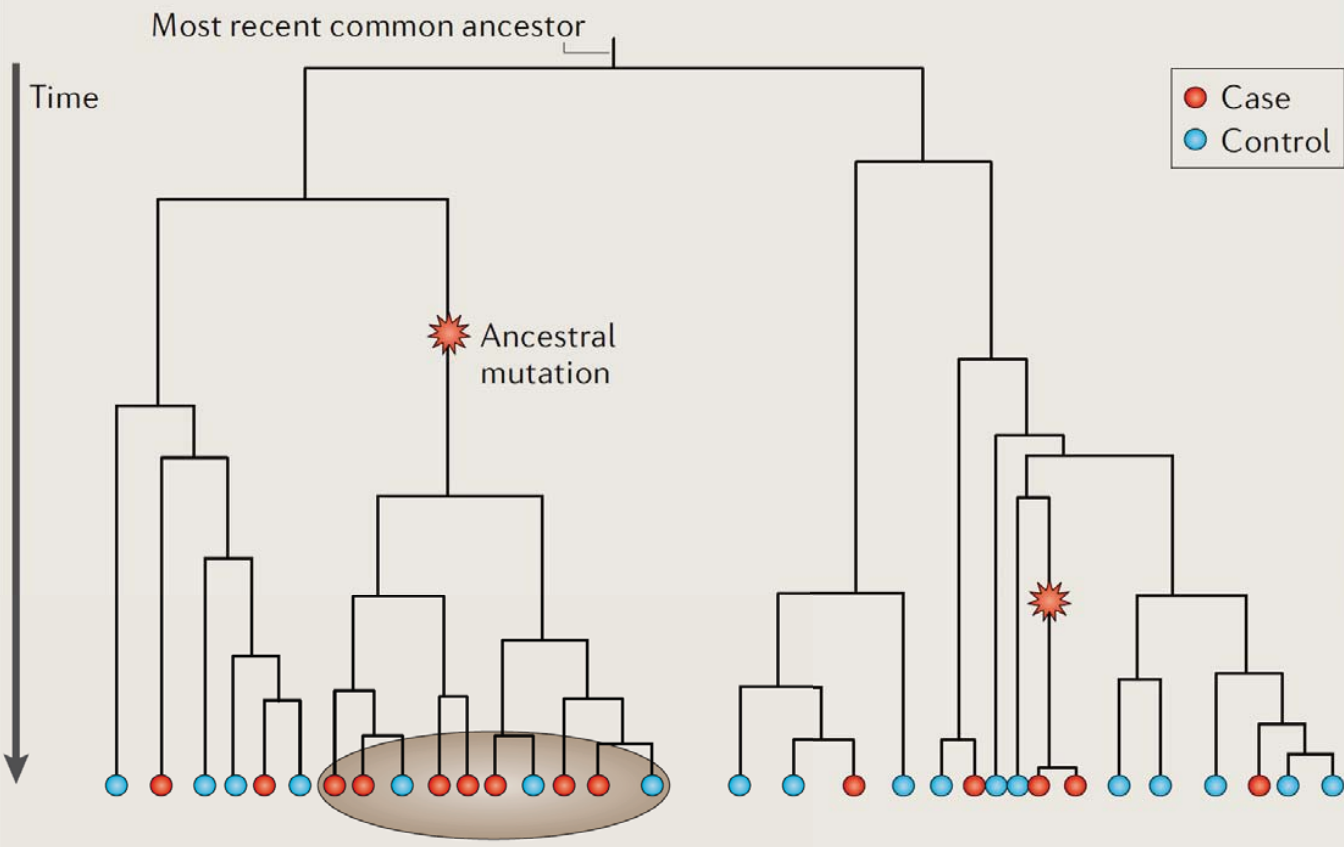


Illumina genotyping software

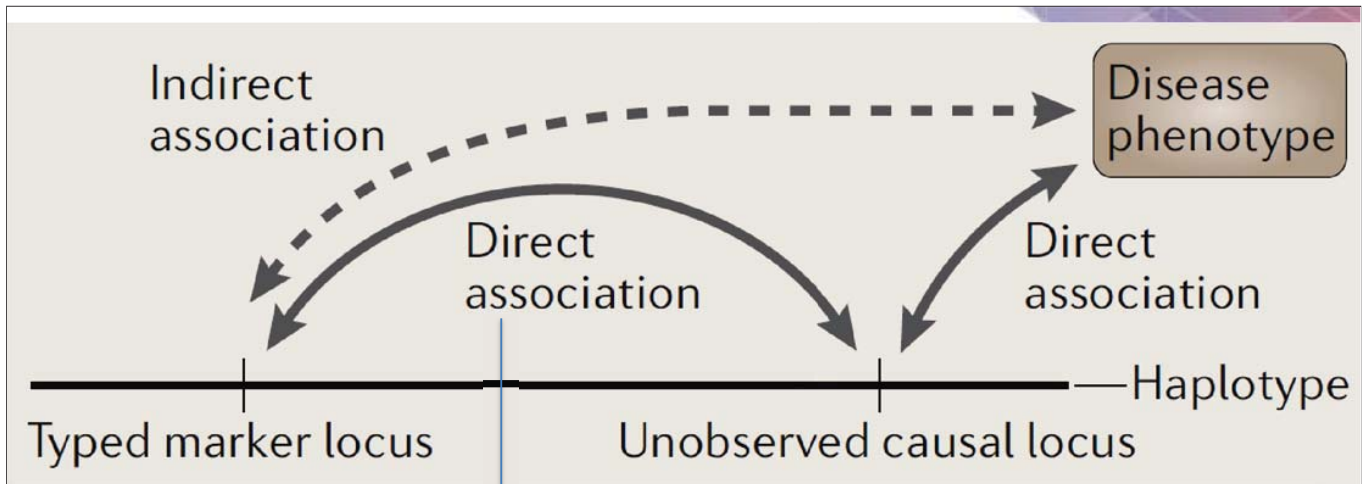


Genotypes are called for each sample (dot) by their signal intensity (Norm R) and Allele Frequency (Norm Theta) relative to canonical cluster positions (dark shading) for a given SNP marker.

Box 1 | Rationale for association studies



Ref: Balding, A tutorial on statistical methods for population association studies, *Nature Reviews Genetics*, 2006.



High linkage disequilibrium (LD)

- the degree to which alleles at two loci are associated
- correlation between two variants

$$r^2 = \frac{D^2}{p_1 * p_2 * q_1 * q_2}$$

- $D = x_{11} - p_1 * q_1$
- $r^2 > 0.8$ considered "high" and $r^2 = 1$ "perfect LD"

Haplotype	Frequency	Allele	Frequency
A_1B_1	x_{11}	A_1	$p_1 = x_{11} + x_{12}$
A_1B_2	x_{12}	A_2	$p_2 = x_{21} + x_{22}$
A_2B_1	x_{21}	B_1	$q_1 = x_{11} + x_{21}$
A_2B_2	x_{22}	B_2	$q_2 = x_{12} + x_{22}$

	A_1	A_2	Total
B_1	$x_{11} = p_1q_1 + D$	$x_{21} = p_2q_1 - D$	q_1
B_2	$x_{12} = p_1q_2 - D$	$x_{22} = p_2q_2 + D$	q_2
Total	p_1	p_2	

If two loci are in linkage equilibrium, $D=0$.
If two loci are in linkage disequilibrium, $D>0$.

Ref: Balding, A tutorial on statistical methods for population association studies, *Nature Reviews Genetics*, 2006.

Quality control of data is very important

- Sample QC
- Variant QC
- Population structure

PROTOCOL

Data quality control in genetic case-control association studies

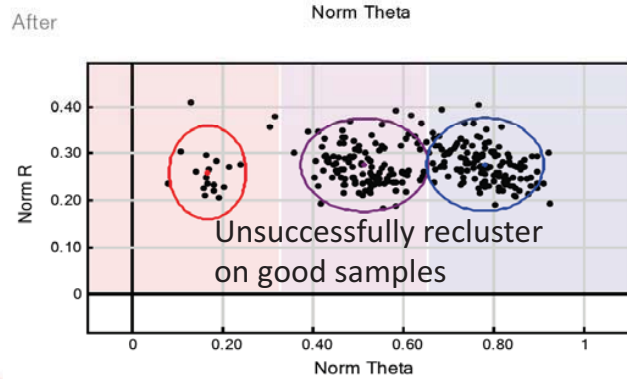
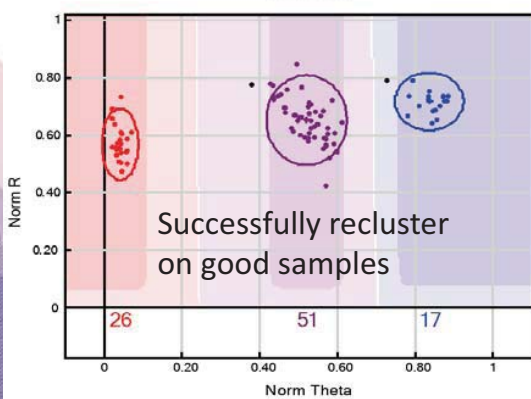
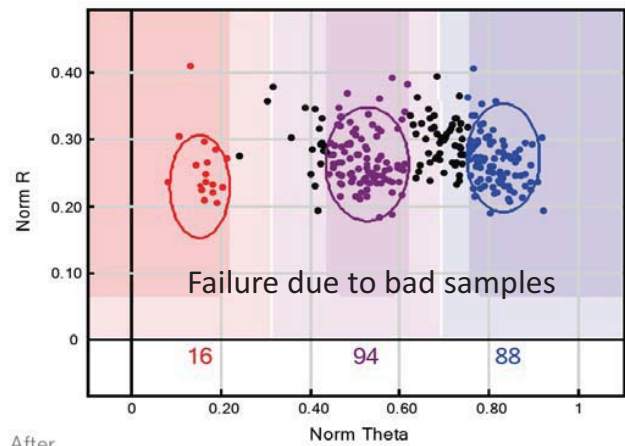
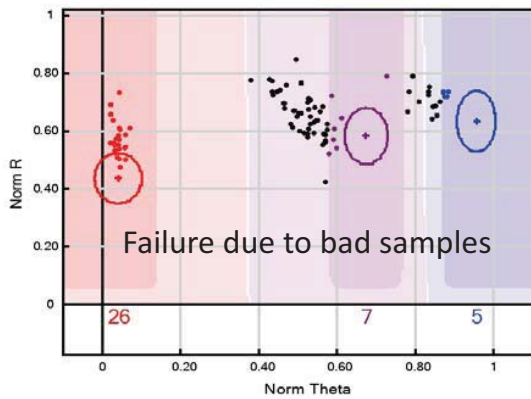
Carl A Anderson^{1,2}, Fredrik H Pettersson¹, Geraldine M Clarke¹, Lon R Cardon³, Andrew P Morris¹ & Krina T Zondervan¹

¹Genetic and Genomic Epidemiology Unit, Wellcome Trust Centre for Human Genetics, University of Oxford, Oxford, UK. ²Statistical Genetics, Wellcome Trust Sanger Institute, Wellcome Trust Genome Campus, Hinxton, Cambridge, UK. ³GlaxoSmithKline, King of Prussia, Pennsylvania, USA. Correspondence should be addressed to C.A.A. (carlanderson@sanger.ac.uk) or K.T.Z. (krinuz@well.ox.ac.uk).

Published online 26 August 2010; doi:10.1038/nprot.2010.116

1564 | VOL.5 NO.9 | 2010 | NATURE PROTOCOLS

Because clustering is not perfect for many reasons...



WTSI QC Pipeline

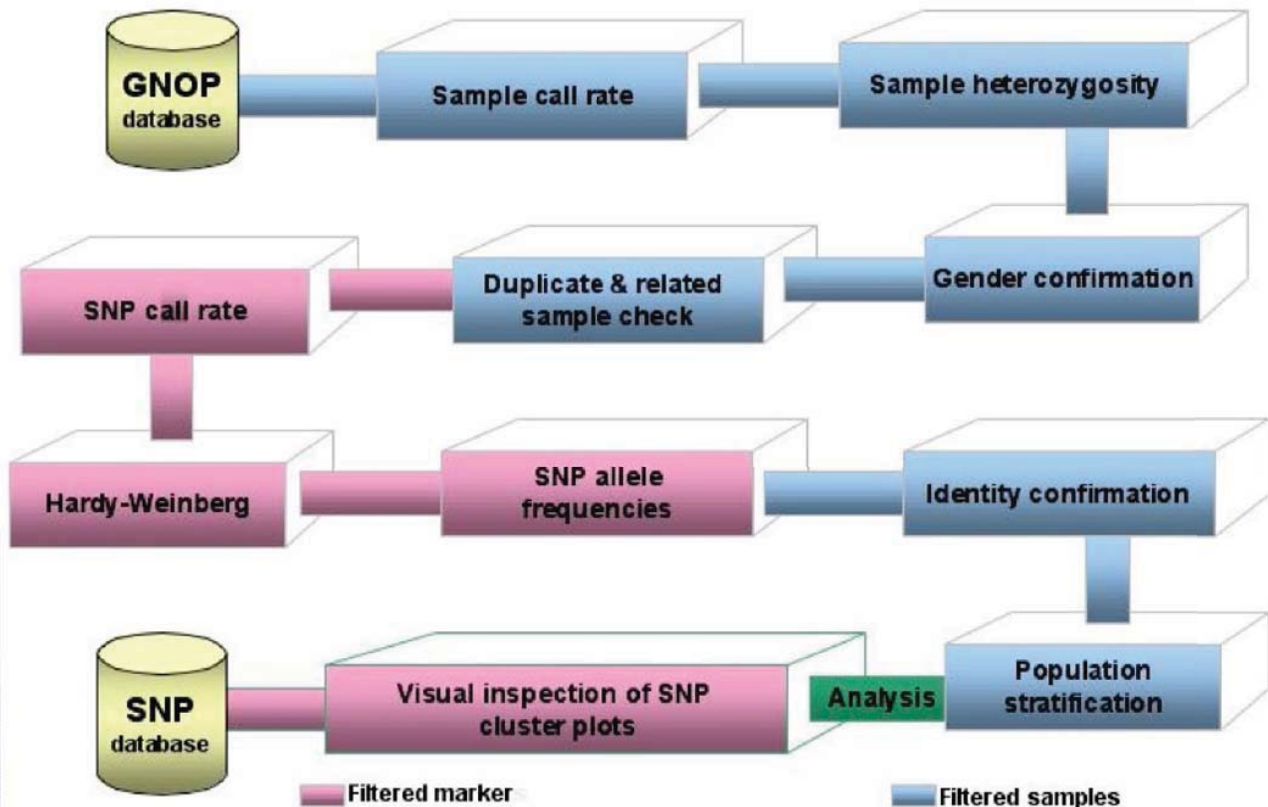
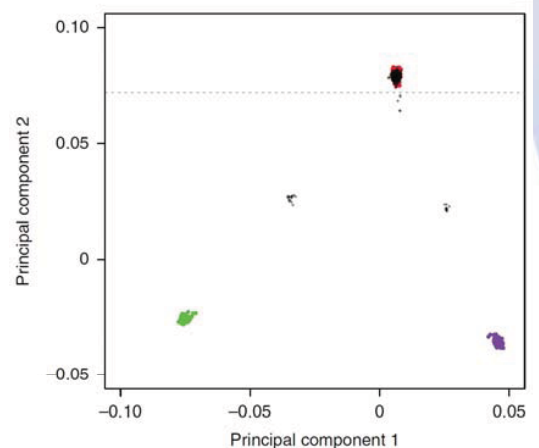
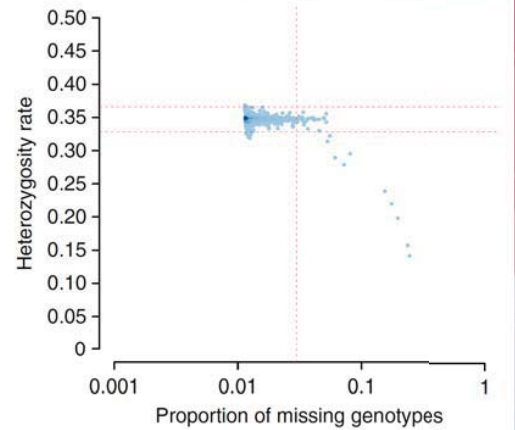


Figure 2: A chronological overview of the pipeline

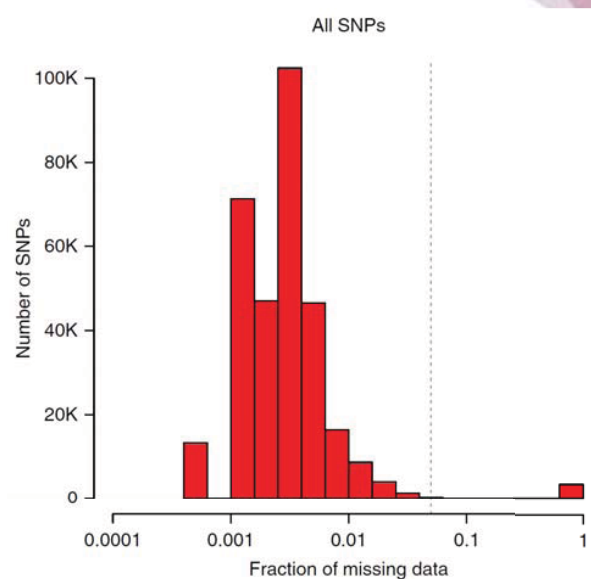
Sample QC

- High batch effects
- Low call rate (<95%)
- Excess of heterozygosity (>mean +5sd or <-5sd)
- Sex mismatch (btw. reported and estimated)
- Population stratification (sub-structure)
- (Hidden) familial relationships



Variant QC

- High batch effects
- Low call rate (<98%)
- Hardy-Weinberg equilibrium (HWE) $p < 1e-06$
- Low minor allele frequency (MAF) <1%



Hardy-Weinberg equilibrium (HWE) test

- Test whether observed genotype counts are deviated from expectations (Hardy-Weinberg equilibrium)
 - Deviations indicate genotyping error, (non-random mating, genetic drift, natural selection, etc.)

Phenotype	White-spotted (AA)	Intermediate (Aa)	Little spotting (aa)	Total
Number	1469	138	5	1612

$$p = \frac{2 \times \text{obs}(AA) + \text{obs}(Aa)}{2 \times (\text{obs}(AA) + \text{obs}(Aa) + \text{obs}(aa))} = \frac{1469 \times 2 + 138}{2 \times (1469 + 138 + 5)} = 0.954$$

$$q = 1 - p$$

$$= 1 - 0.954$$

$$= 0.046$$

Exp(AA) = $p^2 n = 0.954^2 \times 1612 = 1467.4$
Exp(Aa) = $2pq n = 2 \times 0.954 \times 0.046 \times 1612 = 141.2$
Exp(aa) = $q^2 n = 0.046^2 \times 1612 = 3.4$

$$\chi^2 = \sum \frac{(O - E)^2}{E}$$

$$= \frac{(1469 - 1467.4)^2}{1467.4} + \frac{(138 - 141.2)^2}{141.2} + \frac{(5 - 3.4)^2}{3.4}$$

Ref: HW principal from Wikipedia

QC using PLINK

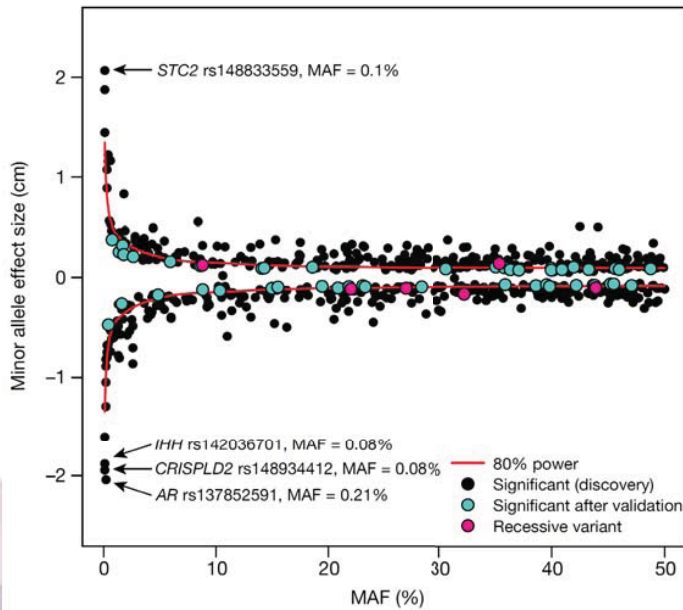
- remove SNPs with MAF < 0.01 : **--maf 0.01**
- remove SNPs with missingness rate ≥ 0.02 (call rate < 0.98) : **--geno 0.02**
- remove SNPs with HWE test P-value < 1e-06 : **--hwe 1e-06**
- remove samples with missingness rate ≥ 0.05 (call rate < 0.95) : **--mind 0.05**

```
plink --bfile gwas --maf 0.01 --mind 0.05 --geno 0.02 --hwe 1e-06 --make-bed --out QC/gwas.1
```

Feature	As summary	As inclusion criteria
Missingness per individual	--missing	--mind N
Missingness per marker	--missing	--geno N
Allele frequency	--freq	--maf N
Hardy-Weinberg equilibrium	--hardy	--hwe N

<https://zzz.bwh.harvard.edu/plink/>

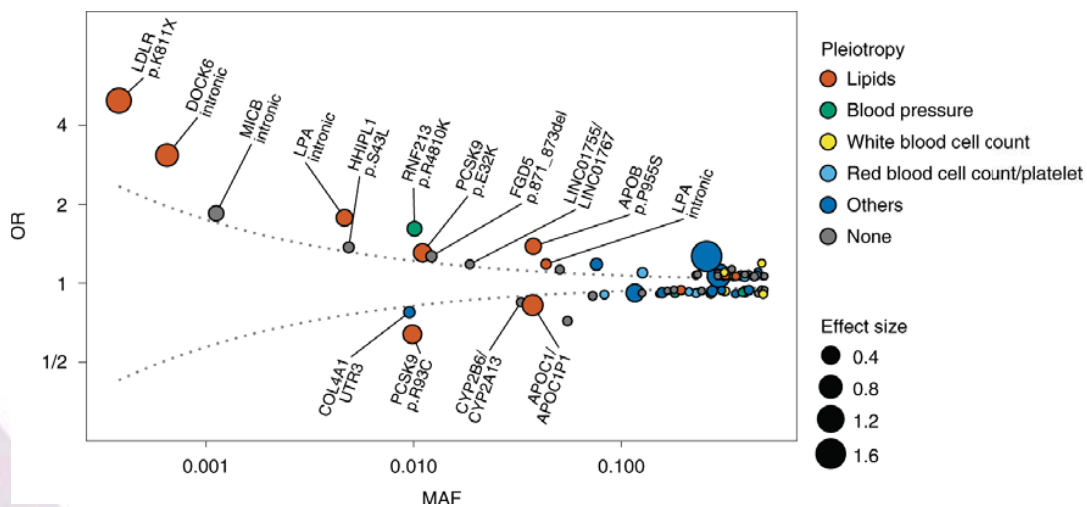
Effects of rare and low-frequency variants on height



- 458,927 individuals
- 697 known loci explained **23.3%** of height heritability
- New loci explained additional **4.1%**
- Rare variants give an increase of **1-2 cm** per allele

Nature 2017 Feb

Effects of rare and low-frequency variants on coronary artery disease

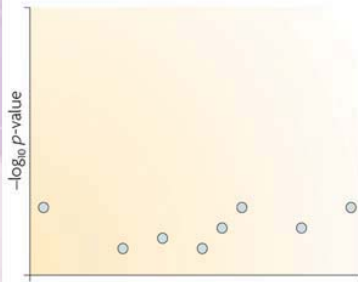


Nature Genetics 52, 1169–1177 (2020)

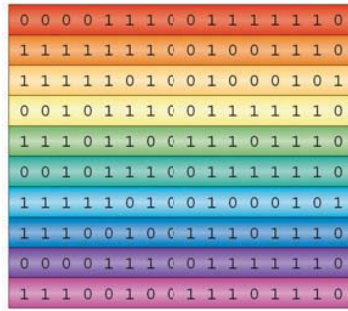
1 of 2 approaches for low-frequency variants:

Statistical imputation : more variants

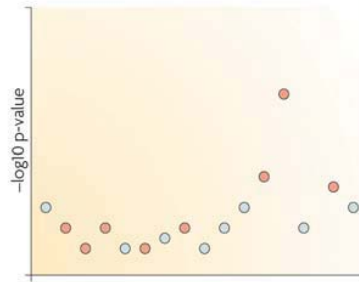
b Testing association at typed SNPs may not lead to a clear signal



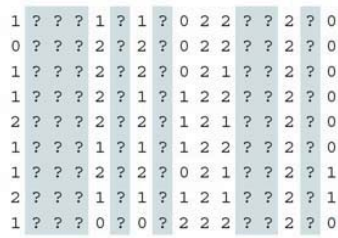
d Reference set of haplotypes, for example, HapMap



f Testing association at imputed SNPs may boost the signal



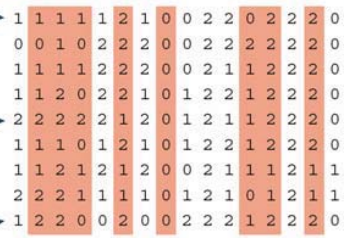
a Genotype data with missing data at untyped SNPs (grey question marks)



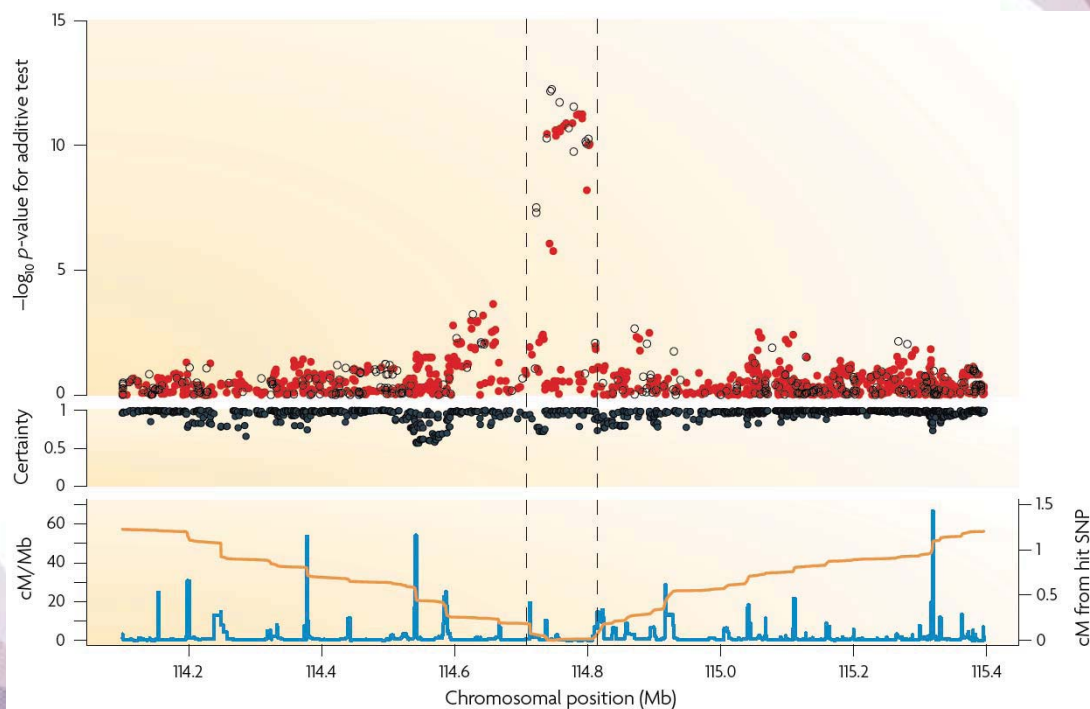
c Each sample is phased and the haplotypes are modelled as a mosaic of those in the haplotype reference panel



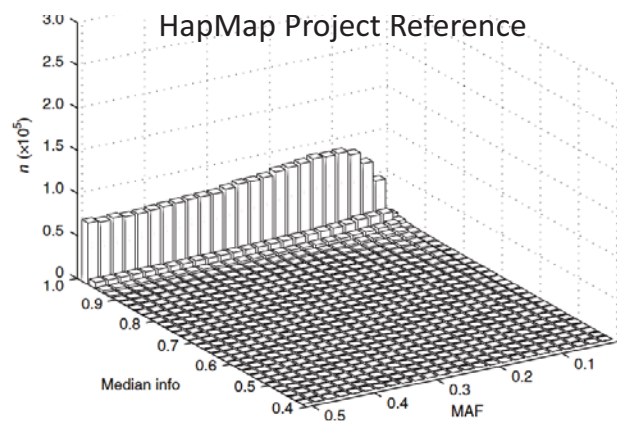
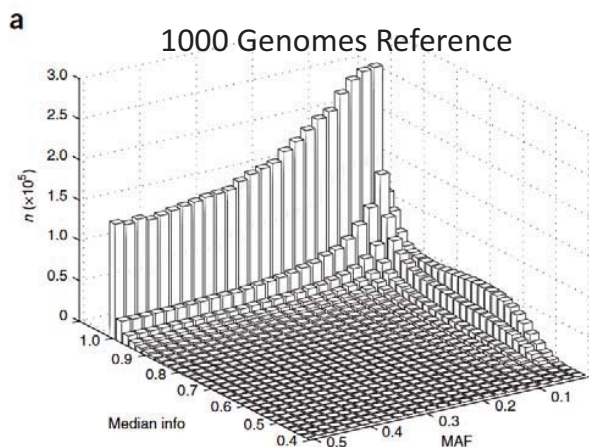
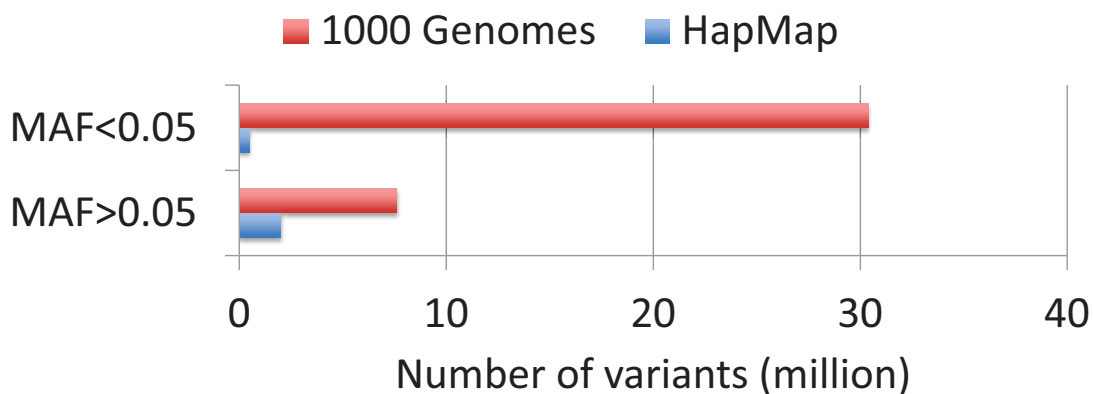
e The reference haplotypes are used to impute alleles into the samples to create imputed genotypes (orange)



Imputed SNPs : good candidates for replication



1000 Genomes based imputation



Imputation reference

The Haplotype Reference Consortium

[OVERVIEW](#) [PARTICIPATING COHORTS](#) [USING THE RESOURCE](#) [CONTACT](#) [SITE LIST](#)

Overview

Goal The Haplotype Reference Consortium (HRC) will create a large reference panel of human haplotypes by combining together sequencing data from multiple cohorts.

Uses The reference panel will be used for [genotype imputation](#) and [phasing](#) in other cohorts, typically [genome-wide association studies \(GWAS\)](#), where genotypes are available from genome-wide SNP microarrays.

Benefits By combining together multiple cohorts, the reference panels produced by the project will be as large as possible in terms of both number of haplotypes, and numbers of variants. This will increase the accuracy of the genotype imputation, especially at low-frequency variants, and the number of imputable variants, thus increasing the power of GWAS.

Ancestry Initially, the reference panel will contain haplotypes from individuals with predominantly European ancestry, although the HRC will include the 1000 Genomes Project data. In the future, we envisage the reference panel increasing in size and consisting of samples from a more diverse set of world-wide populations.

Timelines The first release will consist of 64,976 haplotypes at 39,235,157 SNPs, all with an estimated minor allele count of ≥ 5 . The first release will become accessible early summer 2015.

Run <https://imputationserver.sph.umich.edu/index.html>

Name: optional job name

Reference Panel (Details): HRC r1.1 2016 (GRCh37/hg19)

Input Files (VCF): File Upload

Select Files

Multiple files can be selected by using the **ctrl** / **cmd** or **shift** keys.

Array Build: GRCh37/hg19

Please note that the final SNP coordinates always match the reference build.

rsq Filter: off

Phasing: Eagle v2.4 (phased output)

Population: Other/Mixed

Mode: Quality Control & Imputation

AES 256 encryption

Imputation Server encrypts all zip files by default. Please note that AES encryption does not work with standard unzip programs. Use 7z instead.

I will not attempt to re-identify or contact research participants.

I will report any inadvertent data release, security breach or other data management incident of which I become aware.

Submit Job

TOPMed Imputation server

NIH National Heart, Lung, and Blood Institute | BioData CATALYST | TOPMed Imputation Server | Home About Help Contact | Sign up Login




TOPMed Imputation Server

Free Next-Generation Genotype Imputation Service

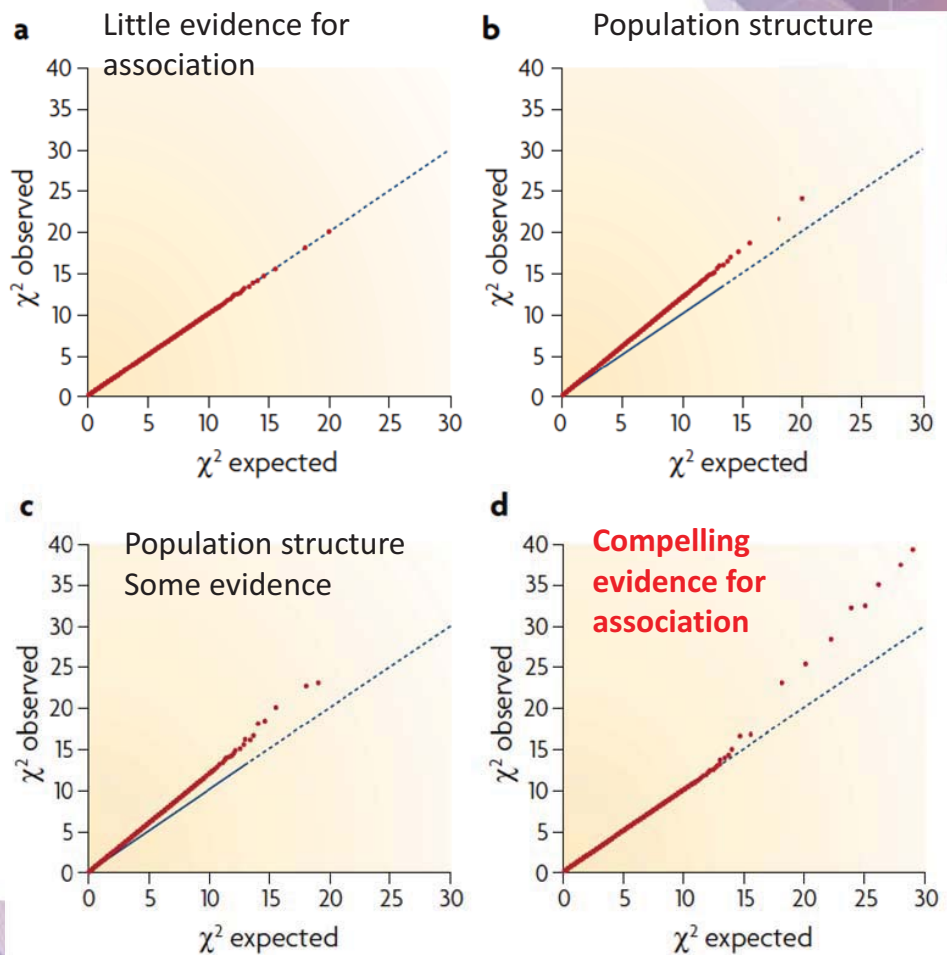
Built from 97,256 deeply sequenced human genomes, this panel contains 308,107,085 genetic variants

41.2M Imputed Genomes	3313 Registered Users	4 Running Jobs
--------------------------	--------------------------	-------------------

The easiest way to impute genotypes

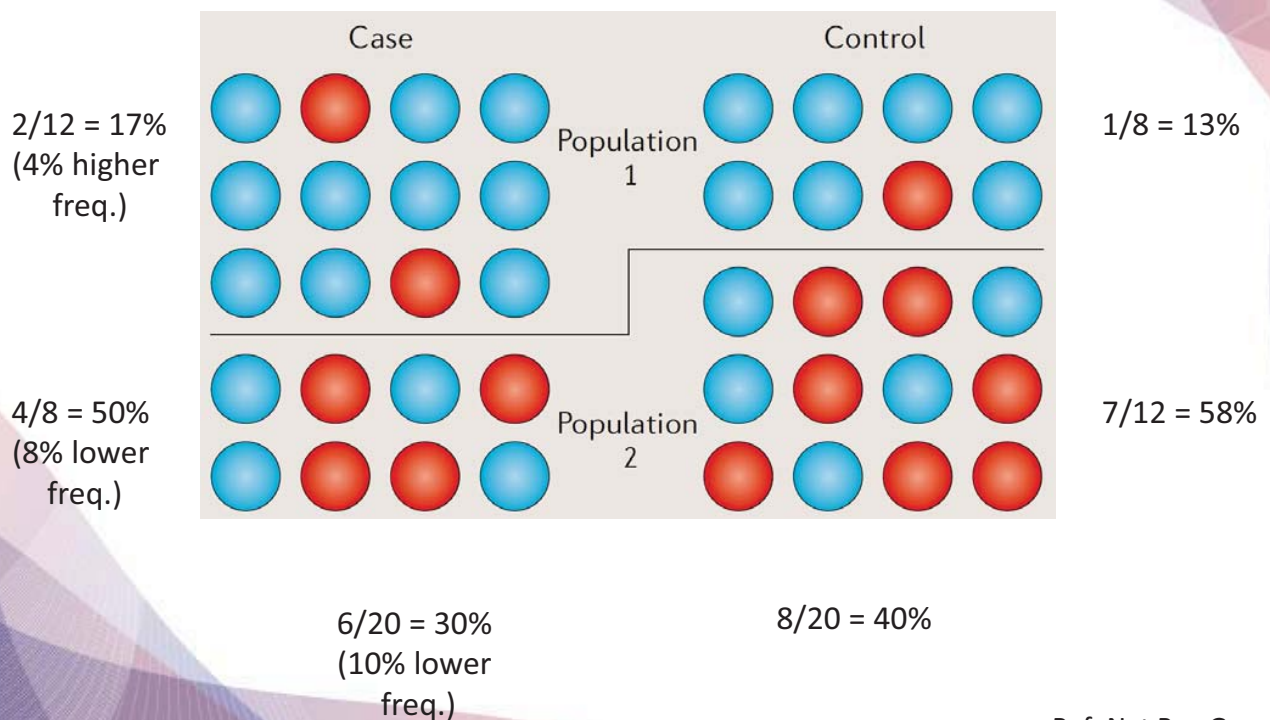
-  Upload your genotypes to our secured service.
-  Choose a reference panel. We will take care of pre-phasing and imputation.
-  Download the results. All results are encrypted with a one-time password. After 7 days, all results are deleted from our server.

Quantile-quantile plot



Ref: Nat Rev Genet

Spurious associations due to population structure



Ref: Nat Rev Genet

Genetics of chopstic use

successful-use-of-selected-hand instruments gene' (SUSHI)

Sample 3: Americans + Chinese

$$\chi^2 = 34.2 \quad P = 4.9 \times 10^{-9}$$

Allele	Use of chopsticks		
	Yes	No	Total
A1	640	340	980
A2	400	100	500
Total	1040	440	1480

Sample 1: Americans

$$\chi^2 = 0 \quad P = 1$$

Allele	Use of chopsticks		
	Yes	No	Total
A1	320	320	640
A2	80	80	160
Total	400	400	800

Sample 2: Chinese

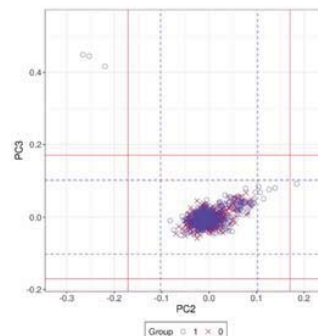
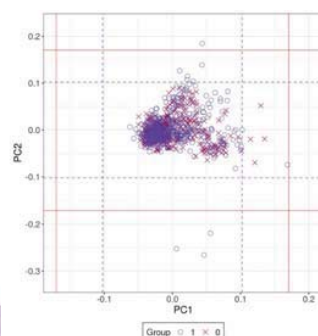
$$\chi^2 = 0 \quad P = 1$$

Allele	Use of chopsticks		
	Yes	No	Total
A1	320	20	340
A2	320	20	340
Total	640	40	680

Ref: Taru Tukiainen

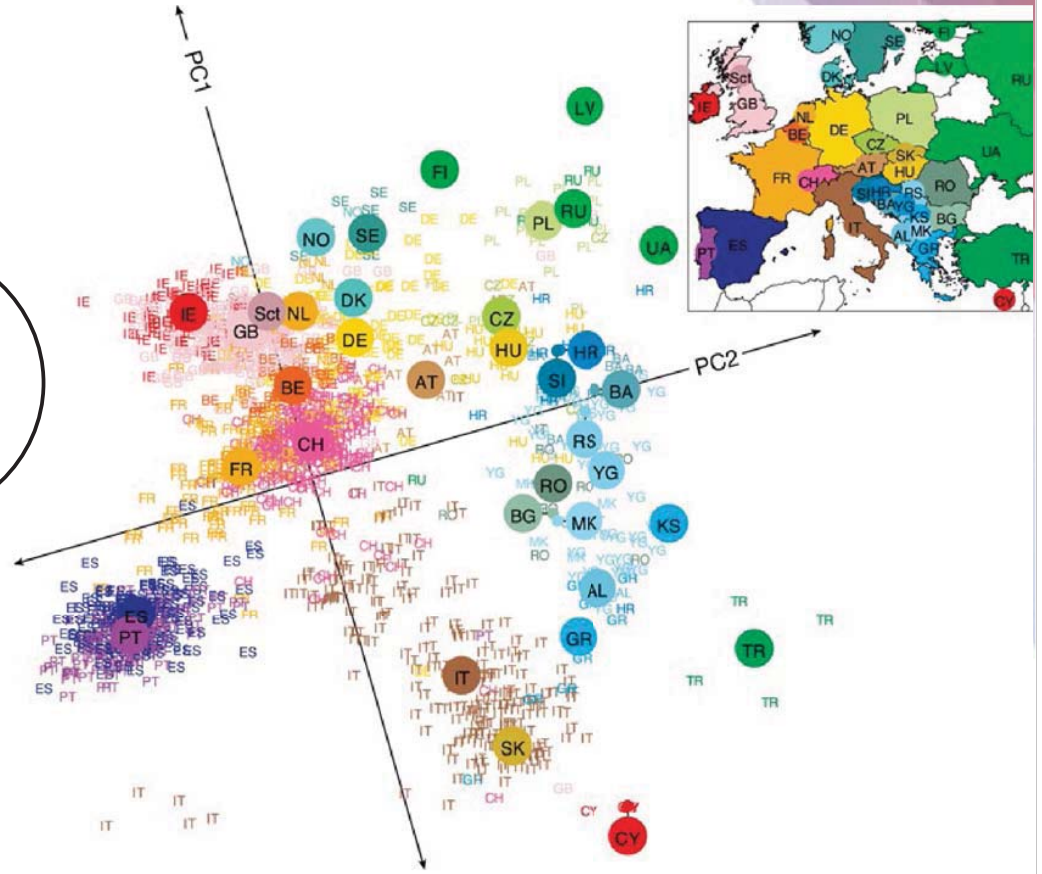
Principal component analysis

- Objective
 - Detect sub-population and any individuals of different ancestry
- Tools
 - smartpca tool of EIGENSOFT software (or using PLINK)
- Solution
 - Check if cases and controls are well overlaid. If not, systematic or technical differences between cases and controls might exist
 - Remove outliers (e.g. $>|5\sigma|$) or include 10 or 20 principal components as covariates in GWAS analyses



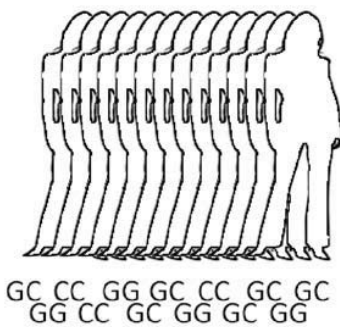
Principle component analysis

PCA plot



Genome-wide association study

Patients



SNP1

Cases

Count of G:
2104 of 4000

Frequency of G:
52.6%

SNP2

Cases

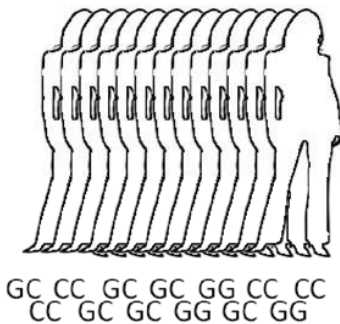
Count of G:
1648 of 4000

Frequency of G:
41.2%

SNP...

Repeat for all SNPs

Controls



Controls

Count of G:
2676 of 6000

Frequency of G:
44.6%

Controls

Count of G:
2532 of 6000

Frequency of G:
42.2%

P-value:
 $5.0 \cdot 10^{-15}$

P-value:
0.33

Basic association test

- H_0 : Frequency of 'A1' is *independent* of case/control status.

	A1	A2
Cases	w	x
Controls	y	z

Odds Ratio (OR): Odds of Allele occurring in cases to the odds of Allele occurring in controls:

$$\frac{w/x}{y/z} = \frac{wz}{xy}$$

$$\chi^2 = (O-E)^2/E$$

[Pearson's chi-Square]

In PLINK, OR > 1 implies A1 is at higher frequency in cases relative to controls.

Note that this is not uniform across all analytical platforms.

Ref: Chris Cotsapas

Regression analysis

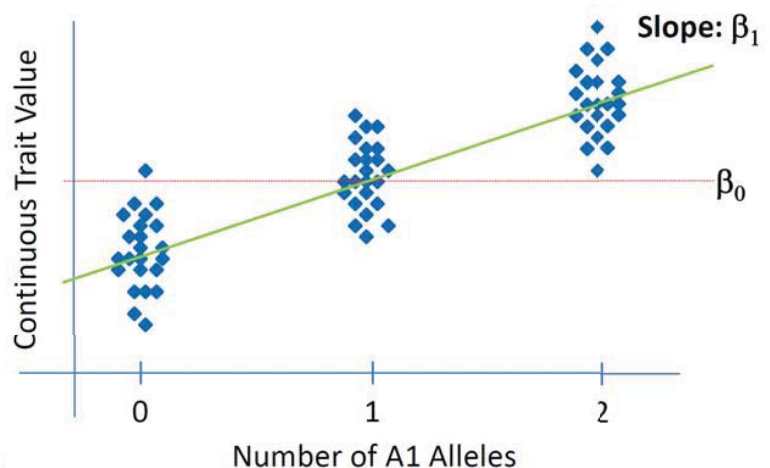
- Analysis of the relationship between a dependent or outcome variable (phenotype) with one or more independent or predictor variables (SNP genotype)

Linear Regression Equation

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$$

Logistic Regression Equation

$$\ln\left(\frac{p_i}{(1-p_i)}\right) = \beta_0 + \beta_1 X_i + \varepsilon_i$$



Ref: Chris Cotsapas

PLINK software (A to Z)

<https://zzz.bwh.harvard.edu/plink/>

- Google PLINK
- Quality control
- Data management
 - `.ped` / `.map`
- Summary stats
- Population stratification
- Association tests
 - Regression, Dominant/Recessive/Trend, Fisher's Exact
- Etc.

```
FAM001 1 0 0 1 2 A A G G A C
FAM001 2 0 0 1 2 A A A G 0 0
...
```

```
Family ID
Individual ID
Paternal ID
Maternal ID
Sex (1=male; 2=female; other=unknown)
Phenotype
```

```
1 rs123456 0 1234555
1 rs234567 0 1237793
1 rs224534 0 -1237697
1 rs233556 0 1337456
...
```

Mixed effects model

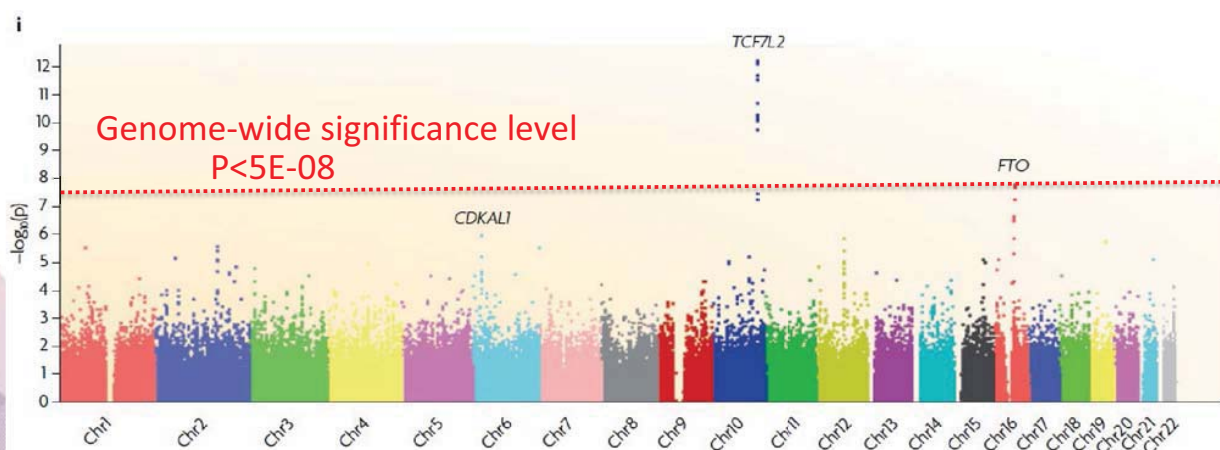
- Mixed effects model
 - $Y = \text{SNP} + \text{sex} + \text{age} + \text{PCs} + \text{Kinship} + e$
- Fixed effects
 - SNP, sex, age, PCs
- Random effects
 - Kinship matrix (due to relatedness)
- Tools
 - Binary (disease): SAIGE, REGENIE
 - Continuous (BMI, blood pressure et al.): BOLT-LMM, REGENIE

Genome-wide significance level

- Multiple-testing (comparisons) problem
 - the problem that arises when many null hypotheses are tested; some significant results are likely even if all the hypotheses are false
- Bonferroni's correction (more stringent method)
 - $0.05 / \#$ of tested variants (usually assuming 1M)
 - $0.05 / 1,000,000 = 5E-08$
- False discovery rate (less stringent method)

<http://pipoli.com>

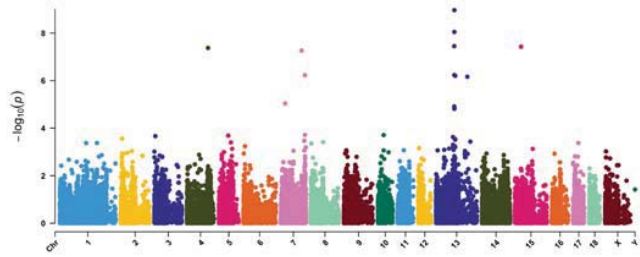
Manhattan plot



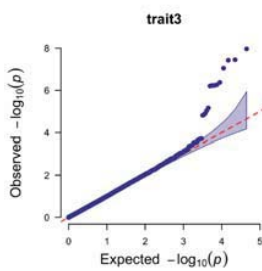
A high-quality drawing tool designed for Manhattan plot of genomic analysis



Circular-Manhattan plot



Rectangular-Manhattan plot



Q-Q plot

```
> install.packages("CMplot")
> library("CMplot")
```

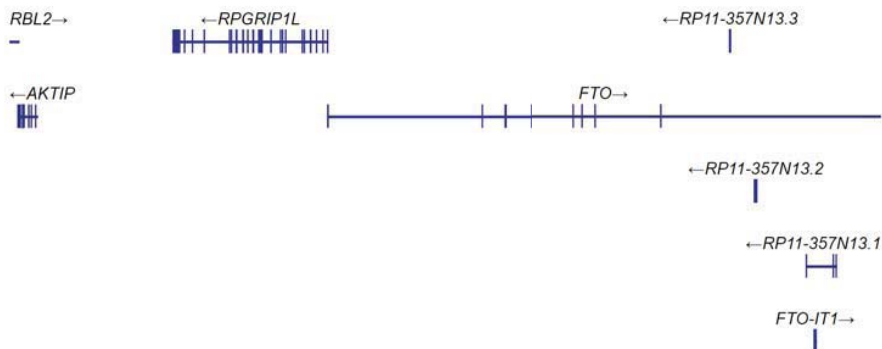
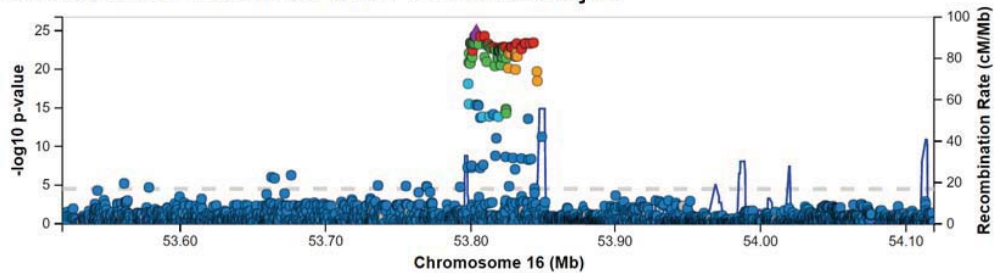
```
## Rectangular Manhattan plot
> cmplot(pig60K, type="p", plot.type="m", LOG10=TRUE,
threshold=NULL, file="jpg", memo="", dpi=300, file.output=TRUE,
verbose=TRUE, width=14, height=6, 축.labels.angle=45)
```

```
## QQ plot
> cmplot(pig60K, plot.type="q", box=FALSE, file="jpg", memo="",
dpi=300, conf.int=TRUE, conf.int.col=NULL, threshold.col="red",
threshold.lty=2, file.output=TRUE, verbose=TRUE, width=5,
height=5)
```

<https://github.com/YinLiLin/CMplot>

Regional plot of association

Scott RA 2017 - DIAGRAM 1000G T2D meta-analysis



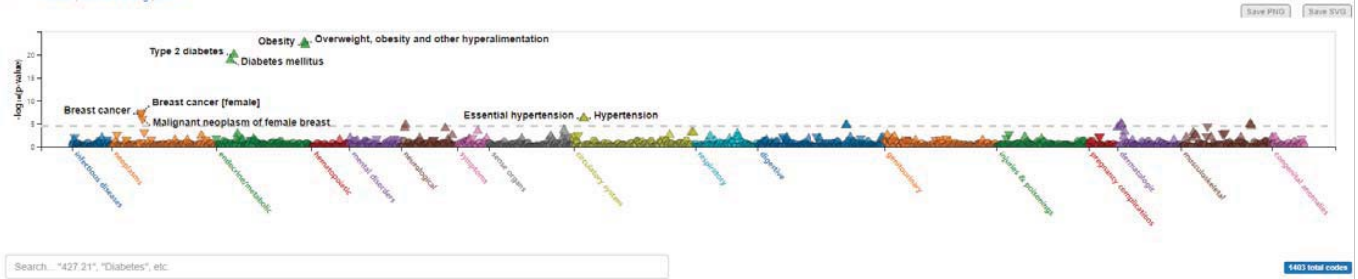
LocusZoom, a visualization tool of GWAS results (<http://locuszoom.org/>)

Phenome-wide analysis (PheWAS)

<https://pheweb.org/UKB-SAIGE/>

16 : 53,821,125 A / G (rs17817712)

Nearest gene: *FTO*
AF: 0.39
View on UCSC, GWAS Catalog, dbSNP

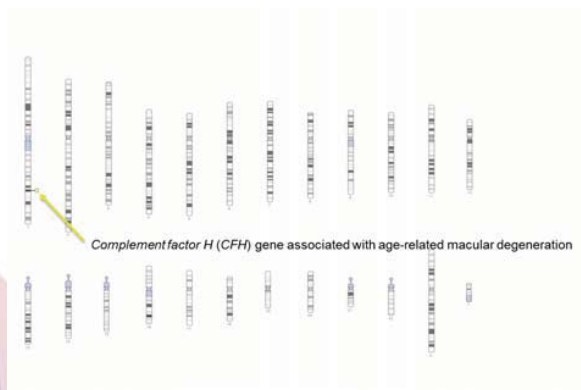


Search: "427.21", "Diabetes", etc. 1463 total codes

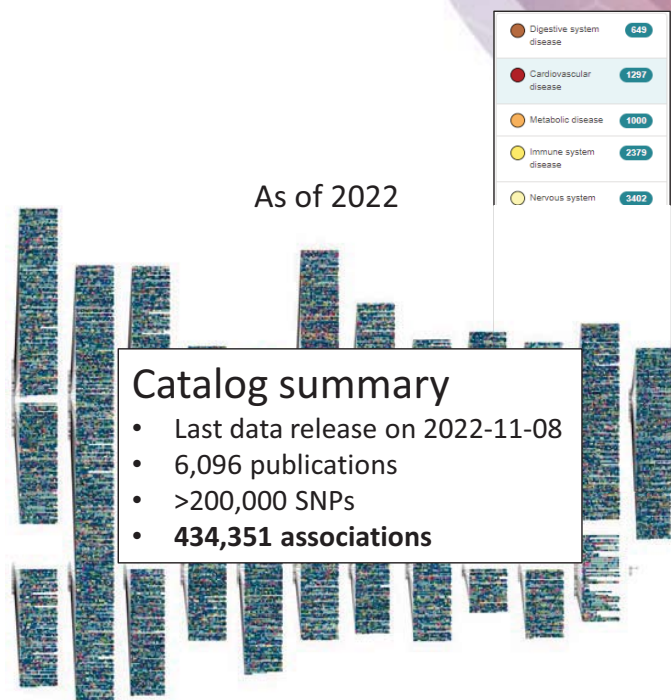
Category	Phenotype	P-value	Effect Size (se)	Number of samples
endocrine/metabolic	Overweight, obesity and other hyperalimentation	1.9e-23	0.14 (0.014)	10968 / 397993
endocrine/metabolic	Obesity	7.3e-23	0.14 (0.015)	10799 / 397993
endocrine/metabolic	Type 2 diabetes	9.6e-21	0.11 (0.011)	18945 / 388756
endocrine/metabolic	Diabetes mellitus	1.4e-19	0.10 (0.011)	20203 / 388756
neoplasms	Breast cancer [female]	4.3e-8	-0.074 (0.014)	12671 / 388549
neoplasms	Breast cancer	7.1e-8	-0.073 (0.013)	12898 / 388549
circulatory system	Hypertension	7.1e-7	0.033 (0.0066)	77977 / 330366
circulatory system	Essential hypertension	7.3e-7	0.033 (0.0066)	77723 / 330366

GWAS loci at $p < 5E-8$

As of 2005



As of 2022

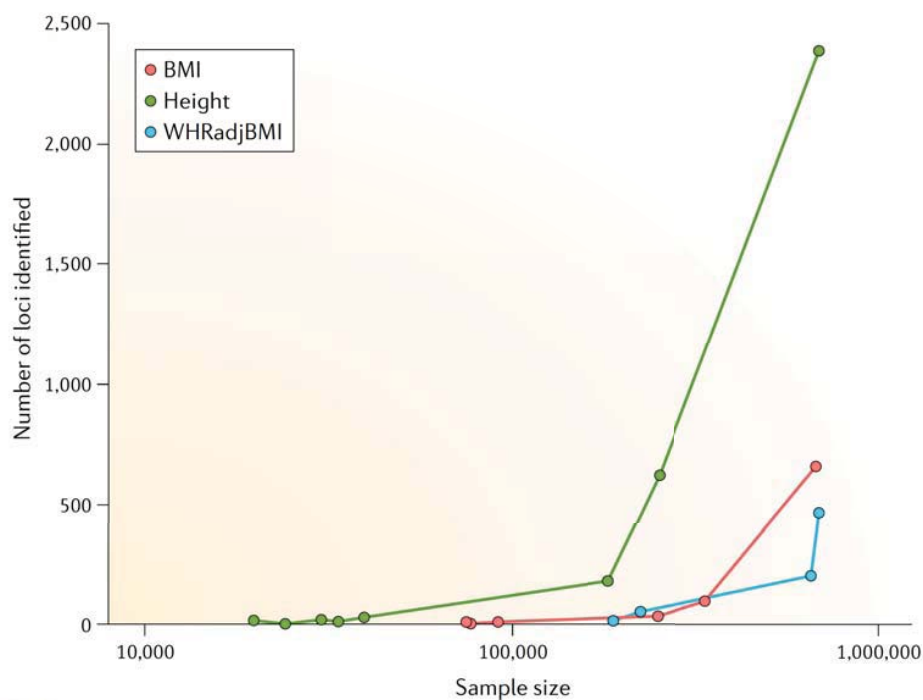


Catalog summary

- Last data release on 2022-11-08
- 6,096 publications
- >200,000 SNPs
- **434,351 associations**

<https://www.ebi.ac.uk/gwas>

Number of loci identified increases as a function of GWAS sample size



Summary of GWAS analysis and tools

- Quality control
 - Sample QC: PLINK
 - Variant QC: PLINK
 - Related samples: KING
 - PCA of genetic ancestry: EIGENSTRAT(smartpca)
- Imputation
 - Haplotype Reference Consortium: Michigan Imputation Server
 - TOPMed Imputation Server
- Association analysis
 - Logistic/linear regression (unrelated): PLINK
 - Mixed effects regression (including related): SAIGE, BOLT-LMM, REGENIE
- Visualization
 - QQ plot: CM-PLOT
 - Manhattan plot: CM-PLOT

Summary

- SNP arrays and statistical imputation provide fast and accurate genotyping of about a million of genetic variants
- Sample-level and variant-level quality control is very important to remove technical errors and false positive findings
- GWAS have identified >200,000 variants associated with various human traits/diseases

Post-GWAS analysis

Summary of post-GWAS analysis and tools

- Understanding genetic architecture
 - SNP-based heritability: LDSC, *GCTA (if genotype available)*
 - Genetic correlation: LDSC (same ancestry), POPCORN or S-LDXR (transancestry)
 - SNP heritability in specific tissues or cells: LDSC-SEG
- Finding causal variants, genes, and pathways
 - Fine-mapping (causal variants): CAVIAR, FINEMAP, PAINTOR, SUSIE
 - eQTL and colocalization analysis (genes): COLOC2
 - Pathway enrichment analysis (pathways or gene sets): MAGMA
- *Identifying individuals at high genetic risk (genotype required)*
 - Polygenic risk score: PRSICE-2, LDPRED, PRS-CS
- Inferring causality between traits
 - Mendelian randomization: MR-BASE, TwoSampleMR (R package)

GWAS summary statistics are publicly available



- Detailed GWAS results of all variants
 - SNP(rsID), effect allele, OR or beta, SE, P value, etc.
- GWAS Catalog
 - <https://www.ebi.ac.uk/gwas>
- GWAS Atlas
 - <https://atlas.ctglab.nl>
- UK Biobank
 - <https://github.com/weizhouUMICH/SAIGE>
- Consortium websites
 - CARDIoGRAMplusC4D
<http://www.cardiogramplusc4d.org/data-downloads>
 - Diabetes DIAGRAM Consortium
<http://diagram-consortium.org/downloads.html>

<https://www.ebi.ac.uk/gwas/>



GWAS Catalog

Home

Search

Diagram

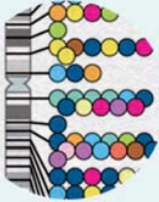
Download

Documentation

About



National Human Genome Research Institute



GWAS Catalog

The NHGRI-EBI Catalog of published genome-wide association studies

Search the catalog



Examples: breast cancer, rs7329174, Yang, 2q37.1, HBS1L, 6:16000000-25000000

Search

Search the Catalog in a number of ways, including by trait, SNP identifier, study and gene.

Diagram

Explore an interactive visualisation of all SNP-trait associations with genome-wide significance ($p \leq 5 \times 10^{-8}$).

Download

Download a full copy of the GWAS Catalog in spreadsheet format and current and older versions of GWAS diagram in SVG format.

Documentation

Including FAQs, our curation process, training materials, related resources and a list of abbreviations.

Summary statistics

A list of all studies for which summary statistics are available in the Catalog.

Ancestry

An introduction to our ancestry curation process.

Summary statistics contain most GWAS results

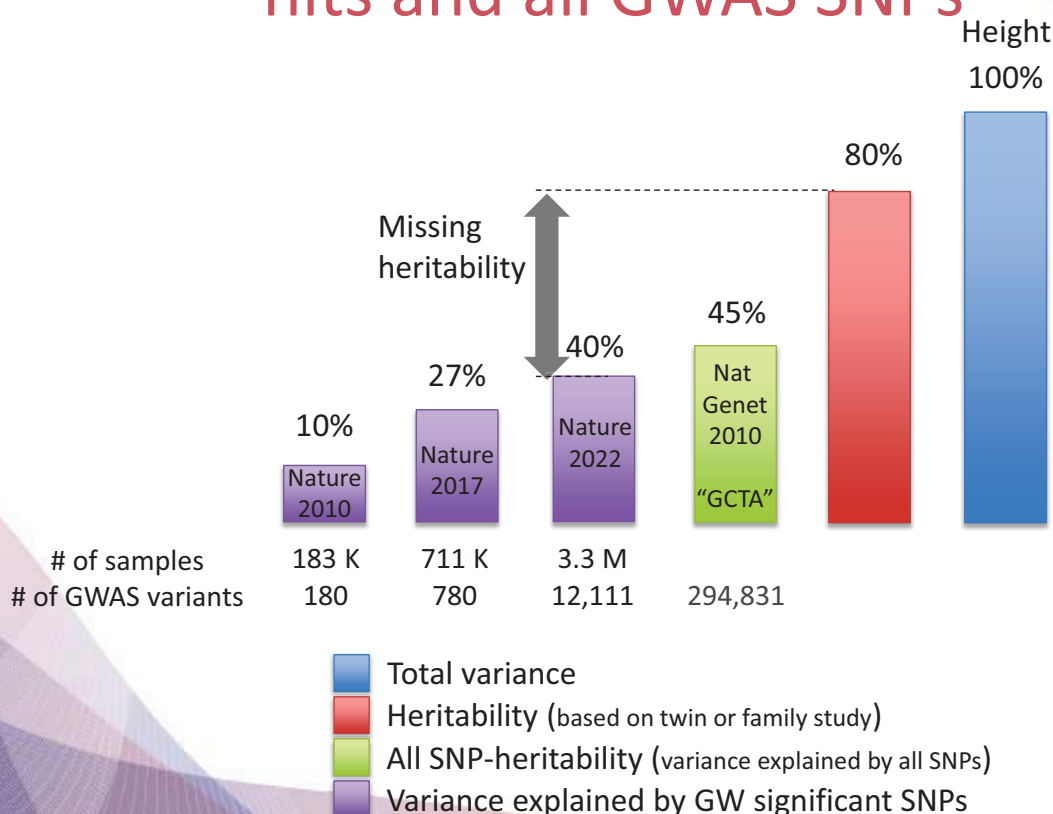
chrom	pos	snpid	ref	alt	ac	af	num_cases	num_controls	beta	sebeta	Tstat	pval
1	16071	rs541172944	G	A	39.843	5.00E-05	650	399970	-2.62	7.55	-0.046	7.28E-01
1	16280	rs866639523	T	C	124	0.000155	650	399970	-2.99	4.06	-0.182	4.61E-01
1	49298	rs10399793	T	C	499790.227	0.623771	650	399970	-0.0468	0.0984	-23.4	6.34E-01
1	54353	rs140052487	C	A	285.302	0.000356	650	399970	-1.22	3.19	-0.12	7.03E-01
1	54564	rs558796213	G	T	121.776	0.000152	650	399970	-0.224	2.89	-0.0269	9.38E-01
1	54591	rs561234294	A	G	79.153	9.90E-05	650	399970	-2.91	6.75	-0.064	6.66E-01
1	54676	rs2462492	C	T	321190.055	0.400866	650	399970	0.039	0.0975	16.9	6.89E-01
1	55326	rs3107975	T	C	6698.62	0.00836	650	399970	-1	0.552	-3.29	6.89E-02

Full information for all the variants (~ several millions)

What can we do with summary statistics?

- Estimate the heritability of traits
- Estimate the genetic correlations among traits
- Test associations between genes and traits
- Infer causality between two traits using MR
- Use for weights of SNPs for disease prediction using polygenic risk score (PRS)
- And more..

Heritability is explained in part by GWAS hits and all GWAS SNPs

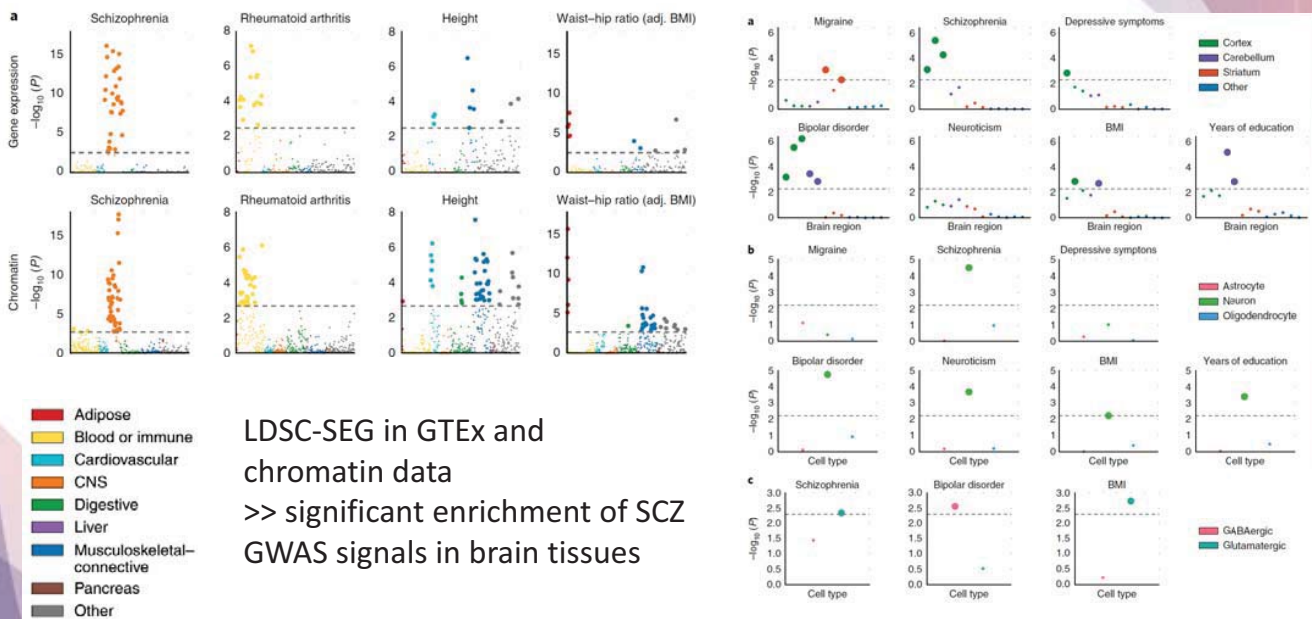


Heritability of GWAS hits and all GWAS SNPs

Trait or Disease	h^2 Pedigree Studies	h^2 GWAS Hits ^a	h^2 All GWAS SNPs ^b
Type 1 diabetes	0.9 ⁹⁸	0.6 ^{99, c}	0.3 ¹²
Type 2 diabetes	0.3–0.6 ¹⁰⁰	0.05–0.10 ³⁴	
Obesity (BMI)	0.4–0.6 ^{101,102}	0.01–0.02 ³⁶	0.2 ¹⁴
Crohn's disease	0.6–0.8 ¹⁰³	0.1 ¹¹	0.4 ¹²
Ulcerative colitis	0.5 ¹⁰³	0.05 ¹²	
Multiple sclerosis	0.3–0.8 ¹⁰⁴	0.1 ⁴⁵	
Ankylosing spondylitis	>0.90 ¹⁰⁵	0.2 ¹⁰⁶	
Rheumatoid arthritis	0.6 ¹⁰⁷		
Schizophrenia	0.7–0.8 ¹⁰⁸	0.01 ⁷⁹	0.3 ¹⁰⁹
Bipolar disorder	0.6–0.7 ¹⁰⁸	0.02 ⁷⁹	0.4 ¹²
Breast cancer	0.3 ¹¹⁰	0.08 ¹¹¹	
Von Willebrand factor	0.66–0.75 ^{112,113}	0.13 ¹¹⁴	0.25 ¹⁴
Height	0.8 ^{115,116}	0.1 ¹³	0.5 ^{13,14}
Bone mineral density	0.6–0.8 ¹¹⁷	0.05 ¹¹⁸	
QT interval	0.37–0.60 ^{119,120}	0.07 ¹²¹	0.2 ¹⁴
HDL cholesterol	0.5 ¹²²	0.1 ⁵⁷	
Platelet count	0.8 ¹²³	0.05–0.1 ⁵⁸	

The American Journal of Human Genetics 90, 7–24, January 13, 2012

Heritability enrichment of specifically expressed genes in tissues and cell types



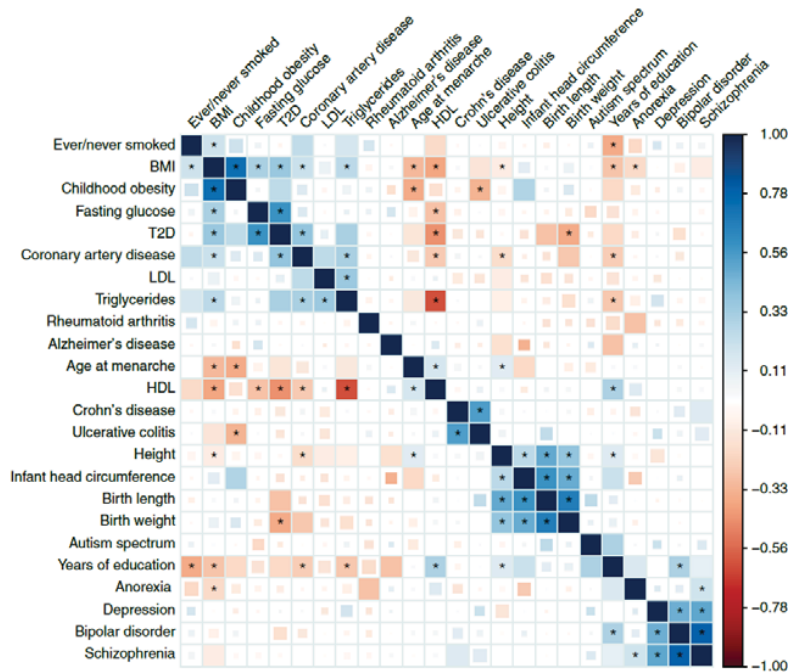
LDSC-SEG in GTEx and chromatin data
 >> significant enrichment of SCZ GWAS signals in brain tissues

Significant enrichment in glutamatergic neurons in cortex

<https://github.com/bulik/ldsc/wiki/Cell-type-specific-analyses>

Nature Genetics 50, 621–629 (2018)

Genetic correlation among diseases

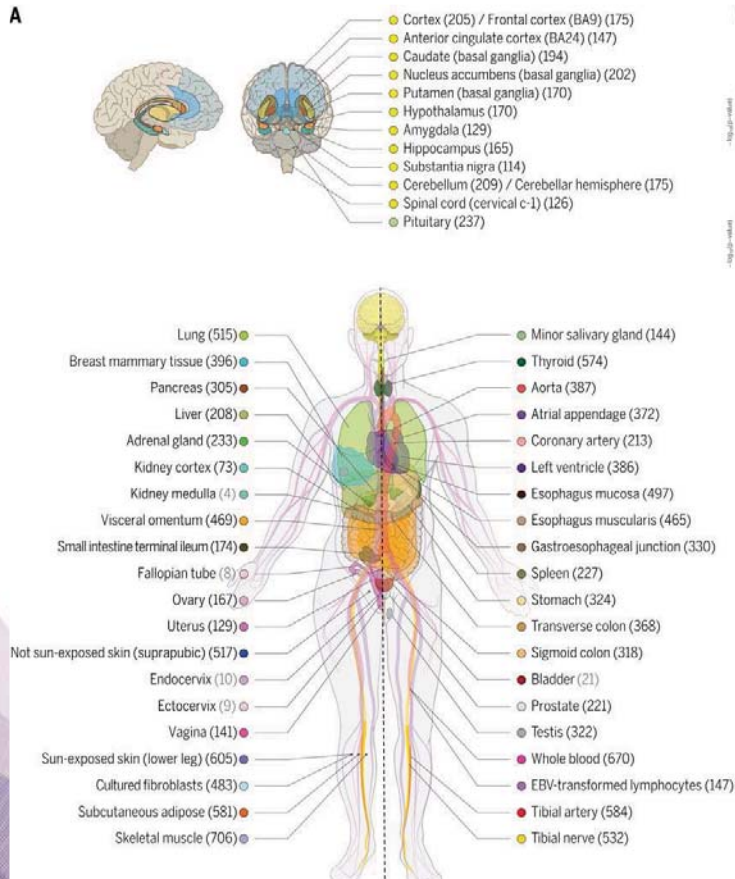


<https://github.com/bulik/ldsc/wiki/Heritability-and-Genetic-Correlation>

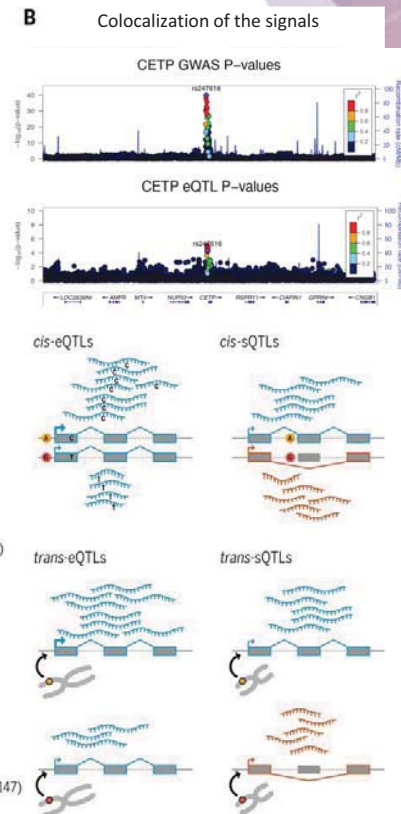
Nat Genet. 2015 Nov; 47(11): 1236–1241.

GTEx (genotype-tissue expression)

A



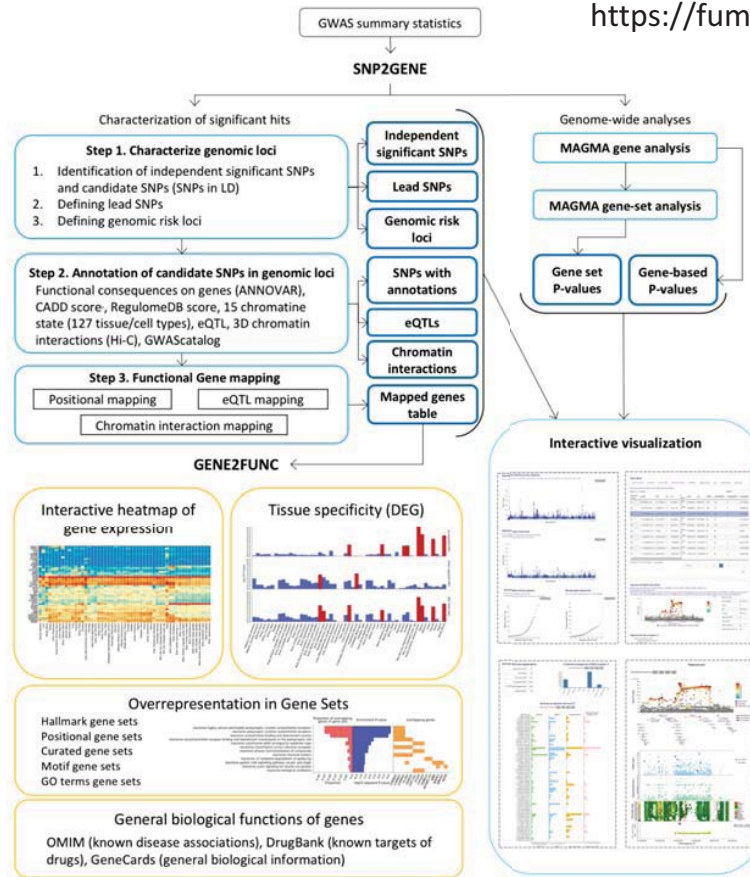
B



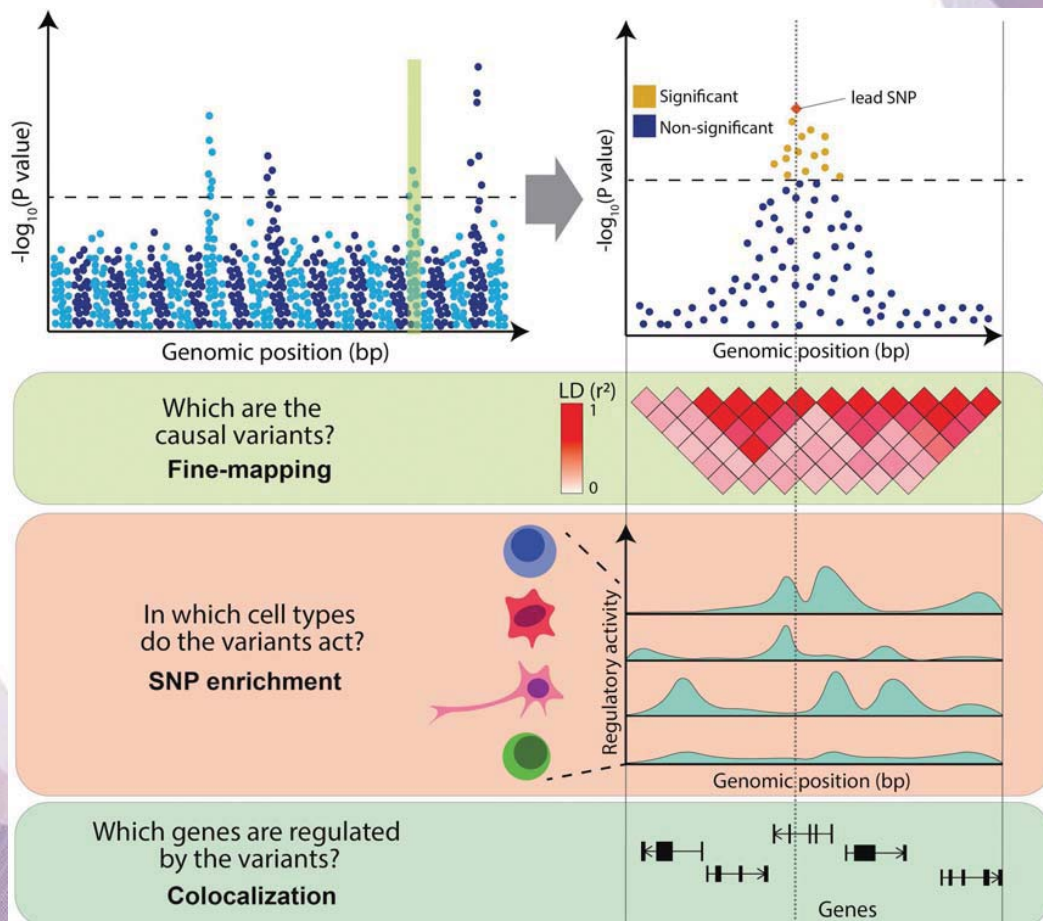
<https://gtexportal.org/home/>

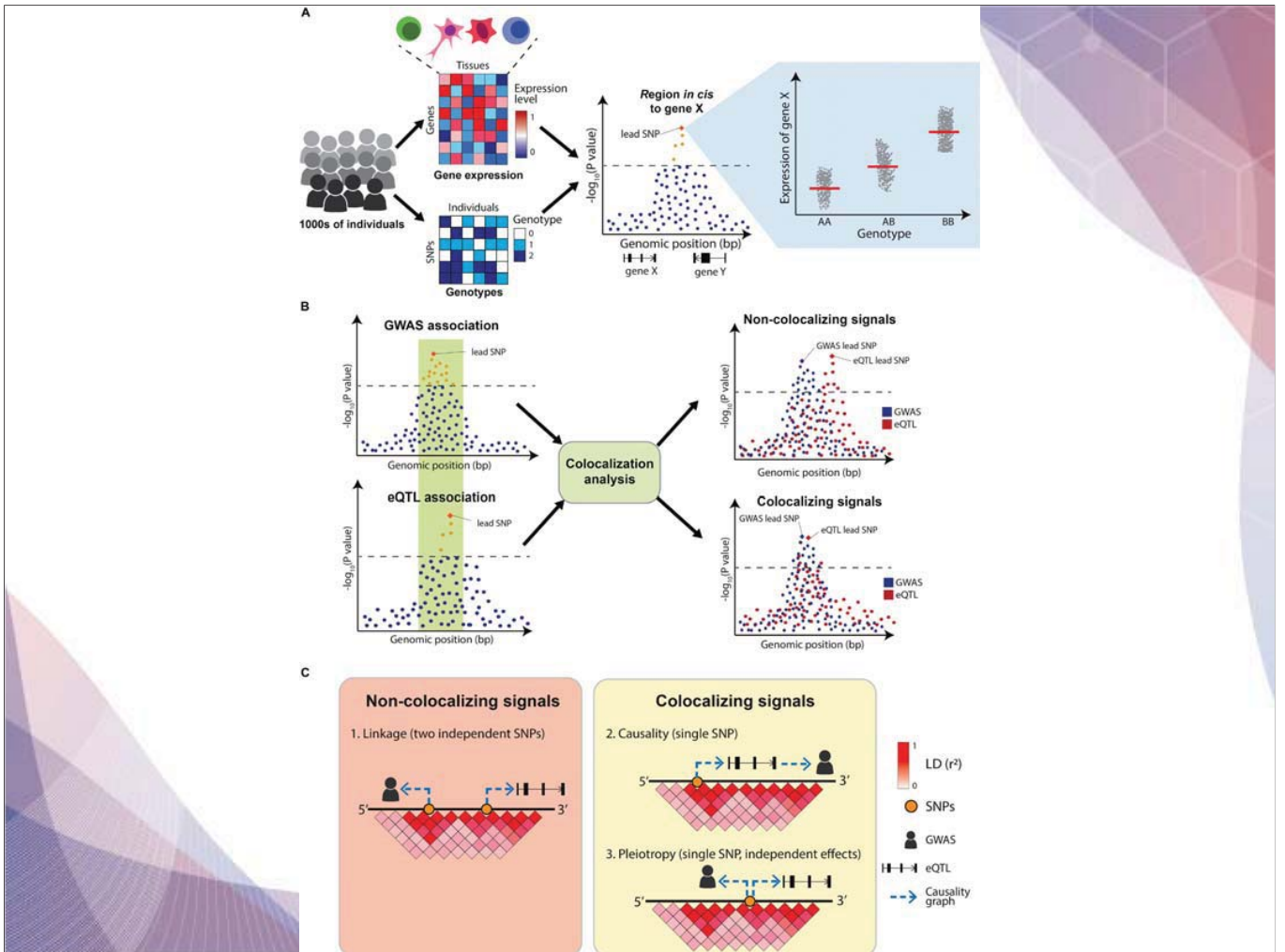
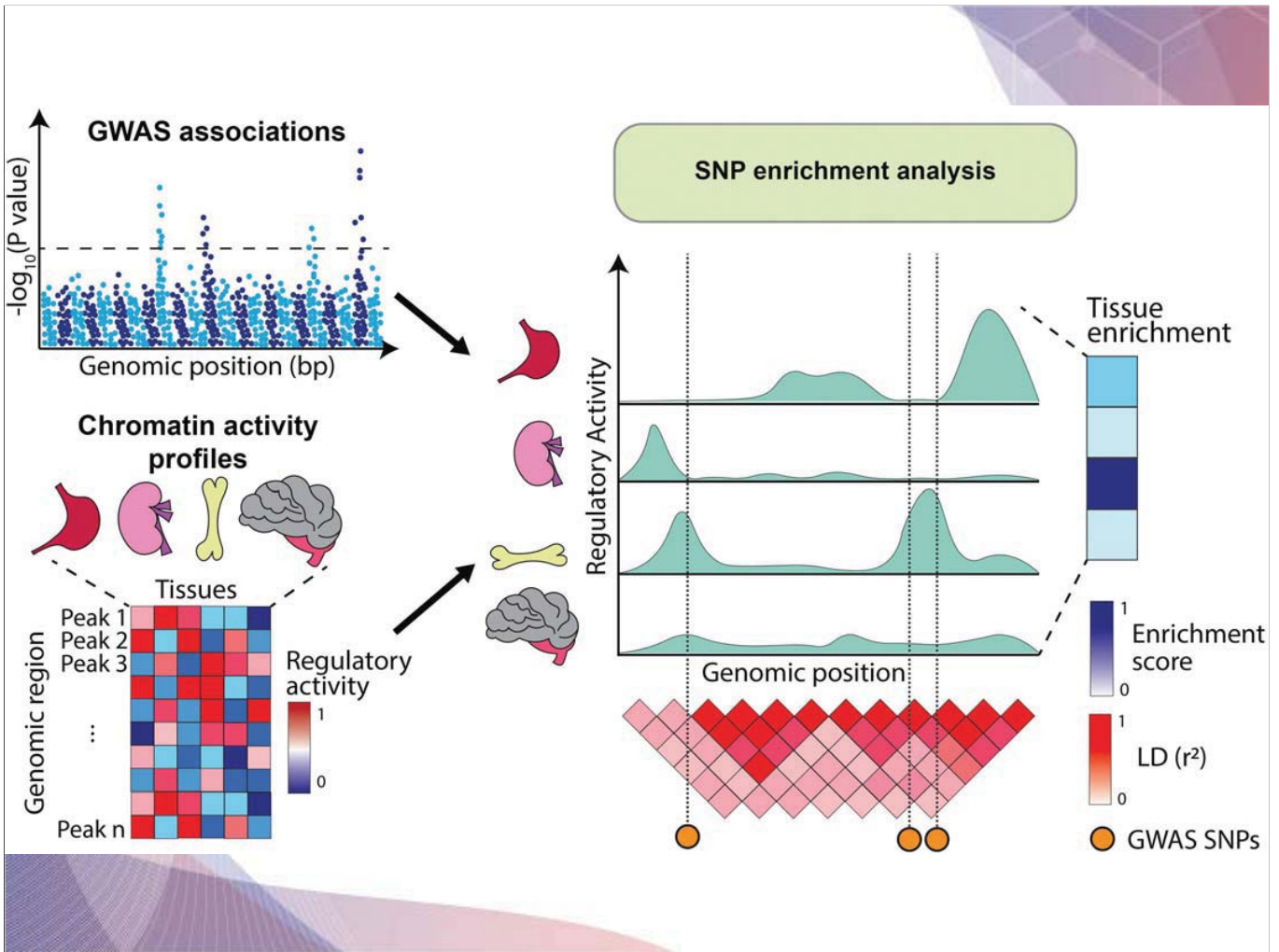
FUMA for various post-GWAS analyses

<https://fuma.ctglab.nl/tutorial>



Review paper: <https://www.frontiersin.org/articles/10.3389/fgene.2020.00424/full>

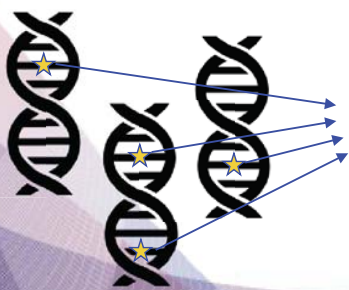




Polygenic risk score: high risk if one has many risk variants



- PRS based on GWAS
 - Sum of risk allele counts across GWAS variants
 - Weighted (effect size) sum of risk allele counts across GWAS variants



Complex disease/phenotypes

“Polygenicity”

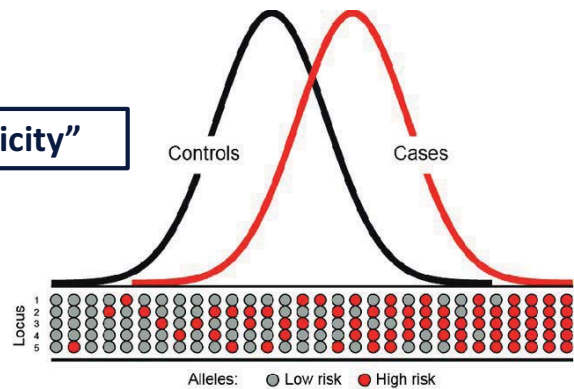


Image Ref: 2014.Genes_ Whiffin and Houlston

PRS (Pruning and Thresholding)

$$PRS(P_T) = \sum_{i=1}^M \mathbb{1}_{\{P_i < P_T\}} \tilde{\beta}_i g_i$$

of SNPs remaining after LD-clumping

↓
normalized marginal effect size estimates

↓
P-value threshold

Summary statistics from GWAS (independent SNP list)

	Effect_Allele	beta	se	P-value
SNP1	A	1.7	0.26	3.11E-11
SNP2	G	1.3	0.25	9.96E-08
SNP3	A	-0.7	0.15	1.53E-06
SNP4	C	-1.2	0.27	4.41E-06
SNP5	G	1.8	0.82	1.41E-02
SNP6	T	0.4	0.41	1.65E-01
...				
...				

P-value < 1e-05

Genotype data

	SNP1	SNP2	SNP3	SNP4	...
indiv1	AT	GG	AC	AA	...
indiv2	AT	TT	AC	CC	...
indiv3	AA	GG	AC	AA	...
indiv4	TT	GT	CC	AC	...
indiv5	AT	TT	AC	AC	...

	SNP1	SNP2	SNP3	SNP4	...
indiv1	1	2	1	0	...
indiv2	1	0	1	2	...
indiv3	2	2	1	0	...
indiv4	0	1	0	1	...
indiv5	1	0	1	1	...

PRS (P+T)

$$PRS(P_T) = \sum_{i=1}^M \mathbb{1}_{\{P_i < P_T\}} \tilde{\beta}_i g_i$$

M → # of SNPs remaining after LD-clumping
 $\tilde{\beta}_i$ → normalized marginal effect size estimates
 P_T → P-value threshold

Summary statistics from GWAS (independent SNP list)

	Effect Allele	beta	se	P-value
SNP1	A	1.7	0.26	3.11E-11
SNP2	G	1.3	0.25	9.96E-08
SNP3	A	-0.7	0.15	1.53E-06
SNP4	C	-1.2	0.27	4.41E-06
SNP5	G	1.8	0.82	1.41E-02
SNP6	T	0.4	0.41	1.65E-01
...				
...				

P-value < 1e-05

PRS for indiv1

$$\begin{aligned}
 &= g_1 * \beta_1 + g_2 * \beta_2 + g_3 * \beta_3 + g_4 * \beta_4 \\
 &= 1 * 1.7 + 2 * 1.3 + 1 * (-0.7) + 0 * (-1.2) \\
 &= 3.6
 \end{aligned}$$

Genotype data

	SNP1	SNP2	SNP3	SNP4	...
indiv1	1	2	1	0	...
indiv2	1	0	1	2	...
indiv3	2	2	1	0	...
indiv4	0	1	0	1	...
indiv5	1	0	1	1	...

	PRS
indiv1	3.6
indiv2	-1.4
indiv3	5.3
indiv4	0.1
indiv5	-0.2

71

nature
genetics

LETTERS

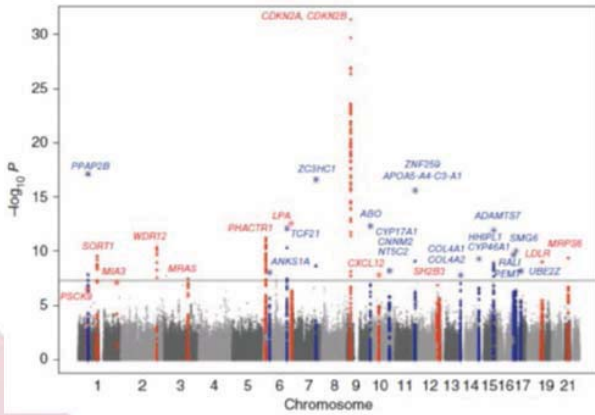
<https://doi.org/10.1038/s41588-018-0183-z>

Genome-wide polygenic scores for common diseases identify individuals with risk equivalent to monogenic mutations

Amit V. Khera^{1,2,3,4,5}, Mark Chaffin^{4,5}, Krishna G. Aragam^{1,2,3,4}, Mary E. Haas⁴, Carolina Roselli⁴, Seung Hoan Choi⁴, Pradeep Natarajan^{2,3,4}, Eric S. Lander⁴, Steven A. Lubitz^{2,3,4}, Patrick T. Ellinor^{2,3,4} and Sekar Kathiresan^{1,2,3,4*}

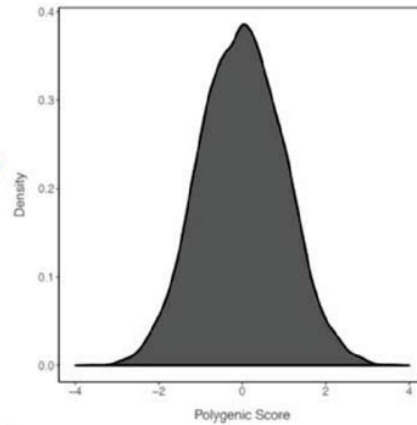
Khera et al. Nat Genet 2018

Polygenic risk score for CAD using 6 million variants is normally distributed



Summary statistics of 6 million variants from a previous large GWAS for coronary artery disease (CAD)

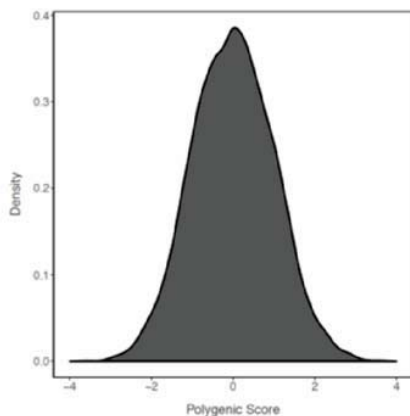
Polygenic score of 6.6 million common variants



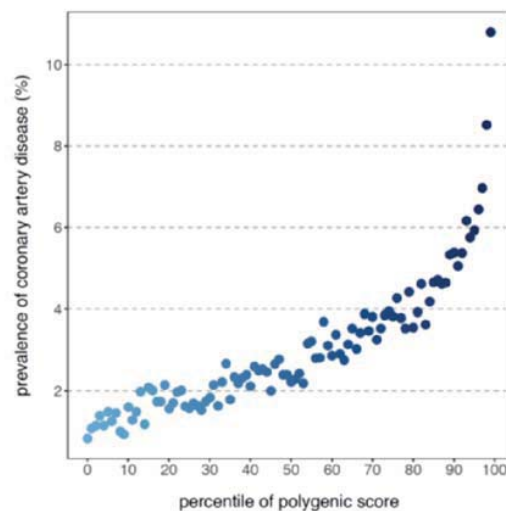
Calculating PRS in independent, testing samples

The empirical risk of CAD rising sharply in the right tail of the distribution

Polygenic score of 6.6 million common variants



Percentile



Calculating PRS in independent, testing samples

Prevalence of CAD according to the percentile of the GPS_{CAD}.

The genetic prediction accuracy was far lower for other populations than for European populations

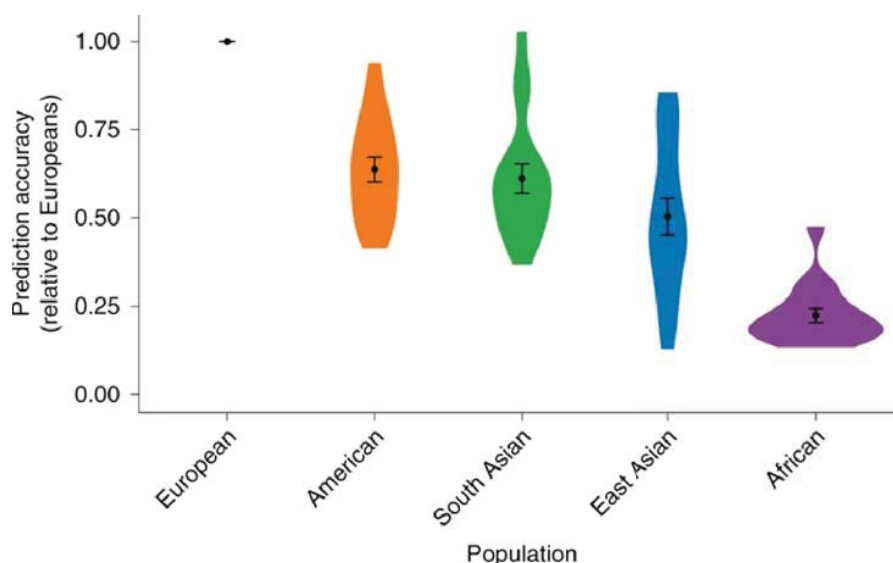
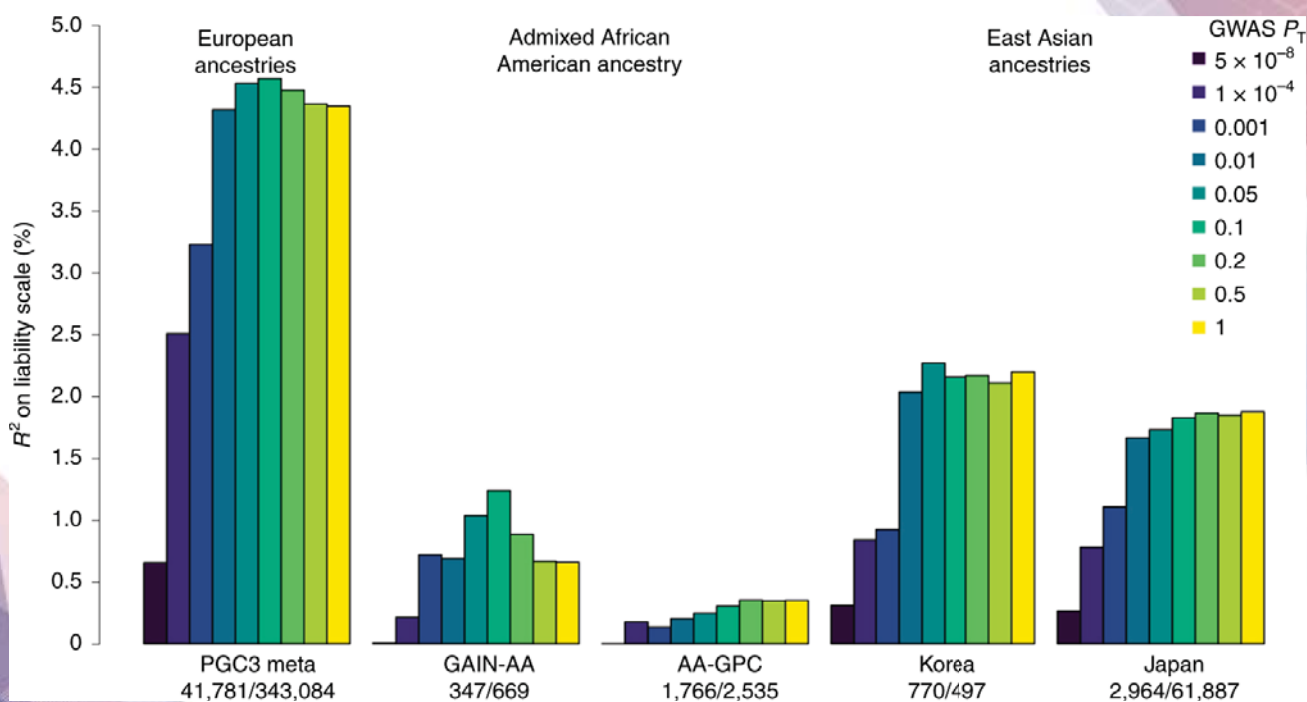


Fig. 3: Prediction accuracy relative to European-ancestry individuals across 17 quantitative traits and 5 continental populations in the UKBB.

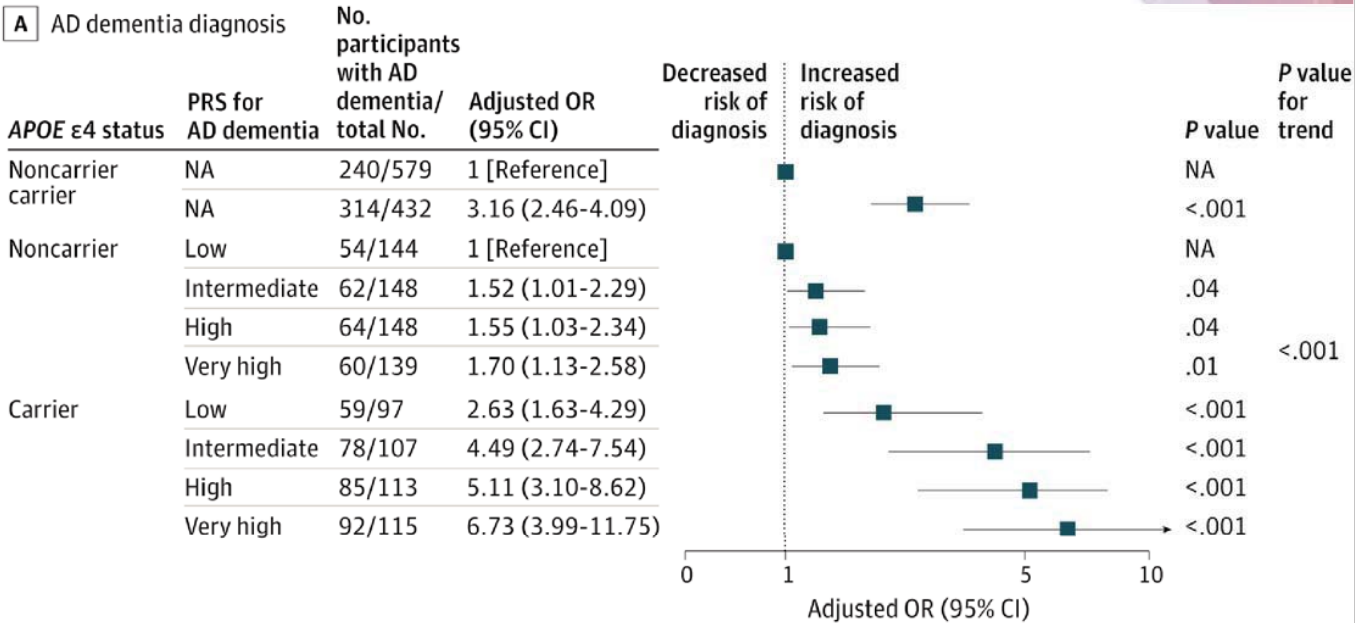
Martin et al. Nat Genet (2019)

PRS for bipolar disorder across ancestries

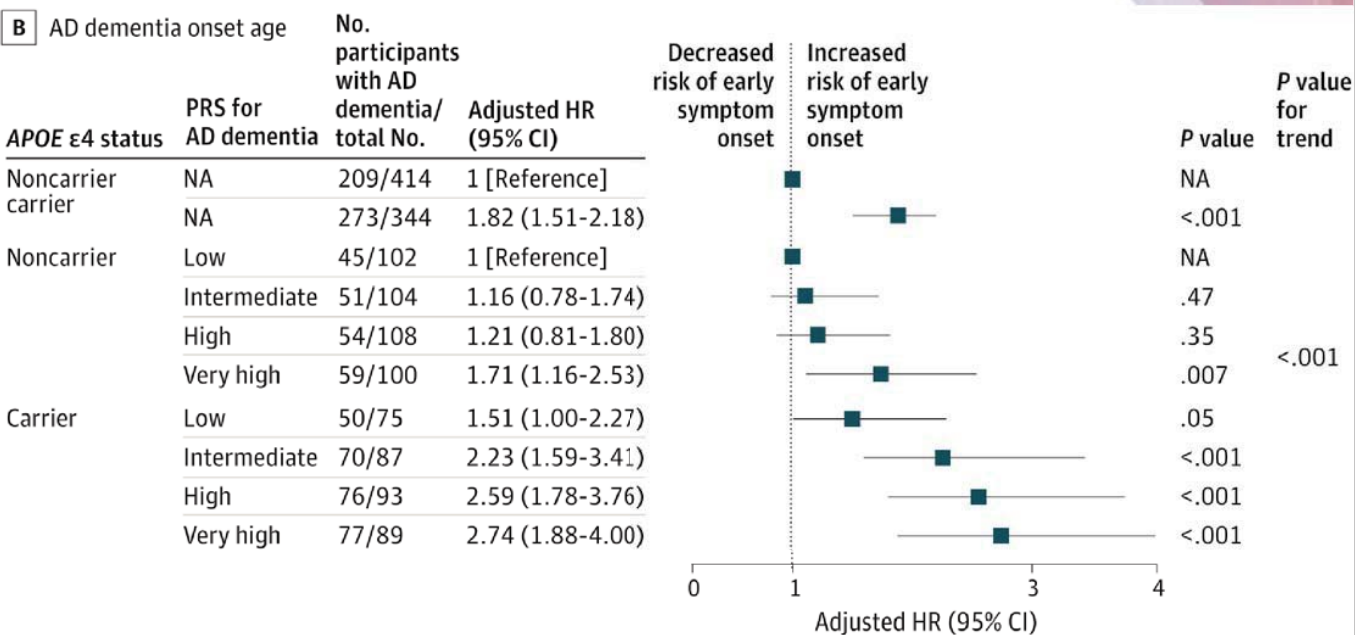


Nature Genetics 53, 817–829 (2021)

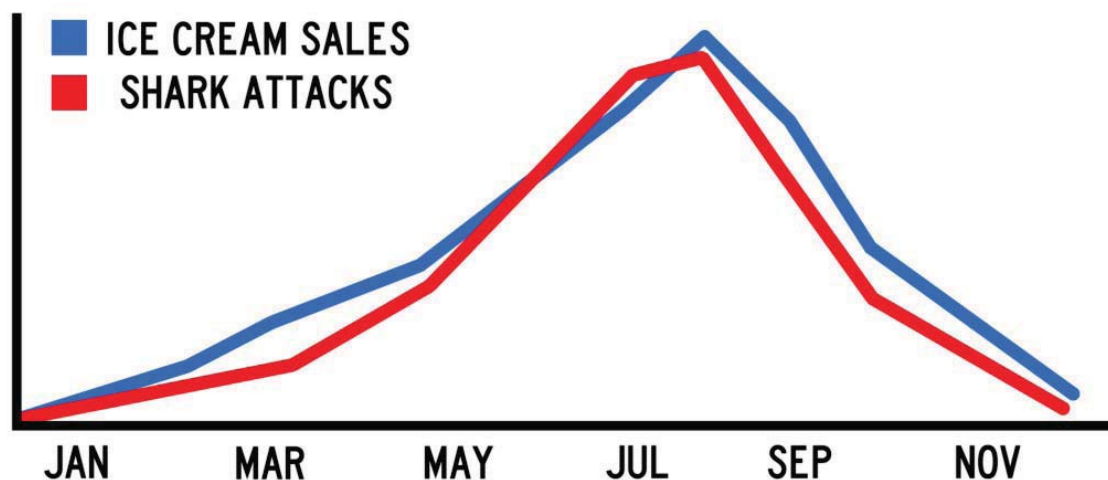
PRS for Alzheimer's disease for Koreans



PRS for Alzheimer's disease for Koreans



CORRELATION IS NOT CAUSATION!

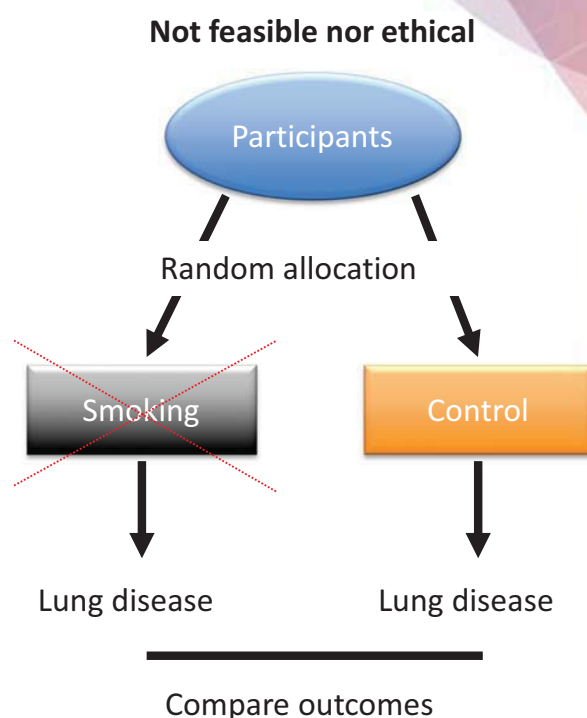


Both ice cream sales and shark attacks increase when the weather is hot and sunny, but they are not caused by each other (they are caused by good weather, with lots of people at the beach, both eating ice cream and having a swim in the sea)

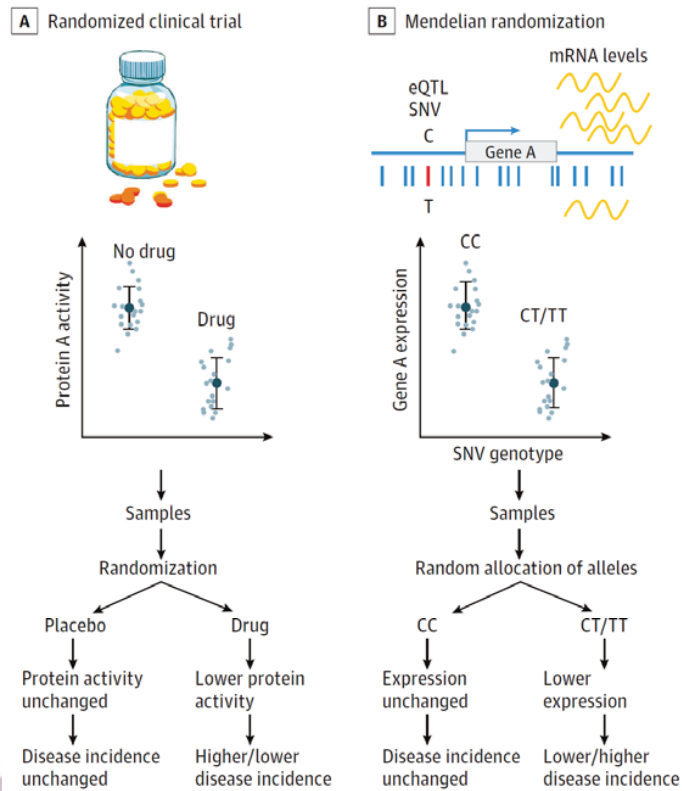
<https://datassist.com/why-journalists-love-causation-and-how-statisticians-can-help/>

RCT is the gold-standard design to infer causality, but

- Exceedingly expensive and time-consuming efforts
- High failure rates (>50% fail owing to lack of efficacy)
- Not always feasible or ethical to conduct

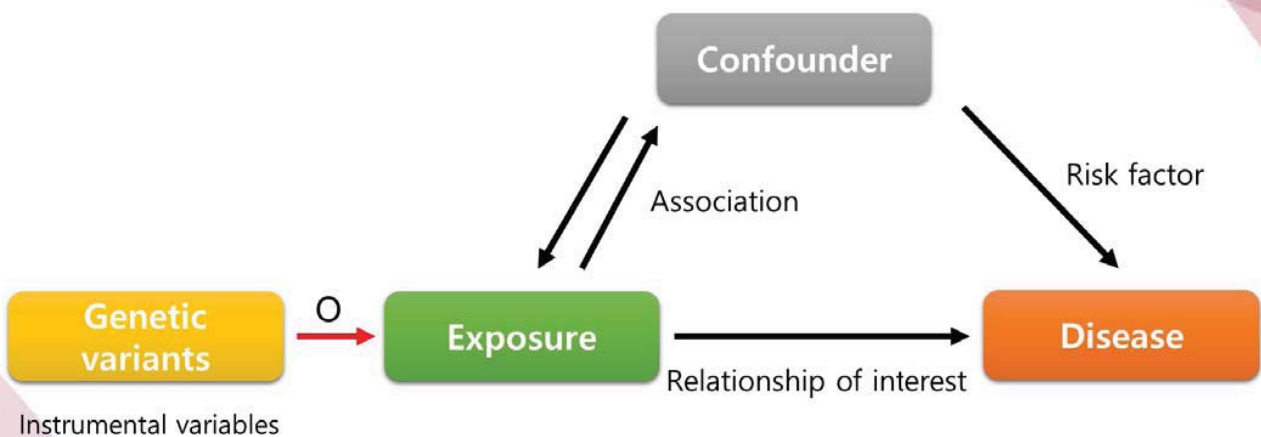


Analogy between RCT and MR



JAMA Psychiatry 2021;78(6):623-631

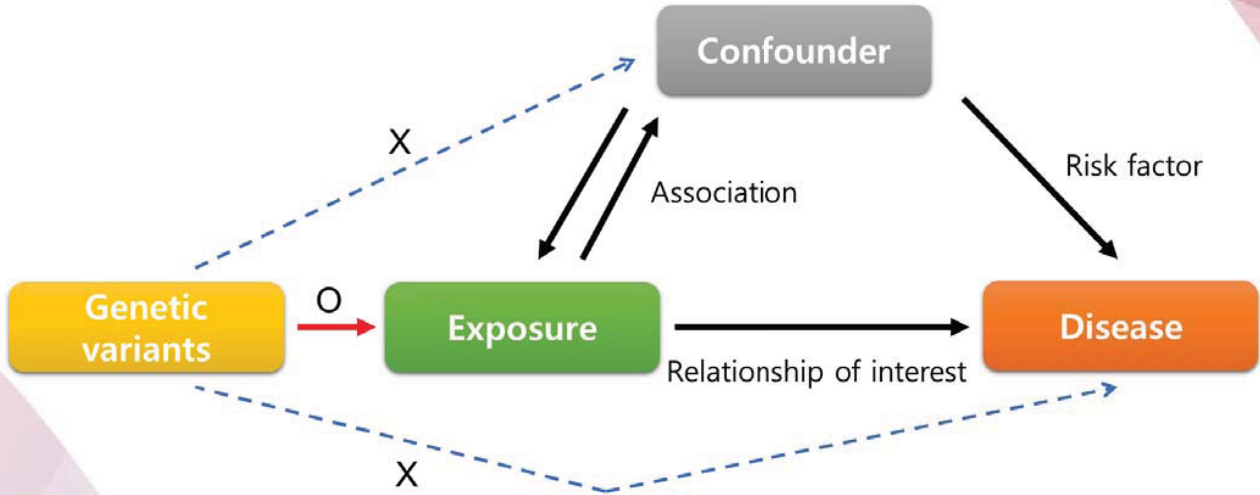
Mendelian randomization (MR)



Assumption I: Genetic variants are robustly associated with the exposure

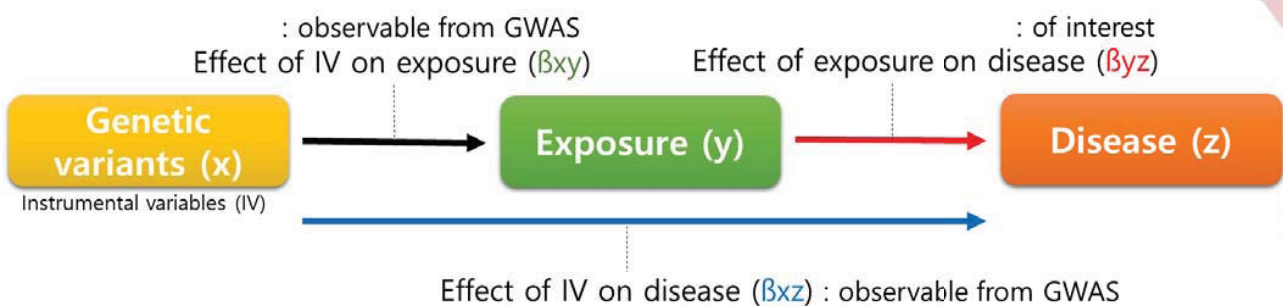
Mendelian randomization (MR)

Assumption II: Genetic variants are NOT associated with the confounder



Assumption III: Genetic variants are associated with the outcome through the exposure (vertical pleiotropy)

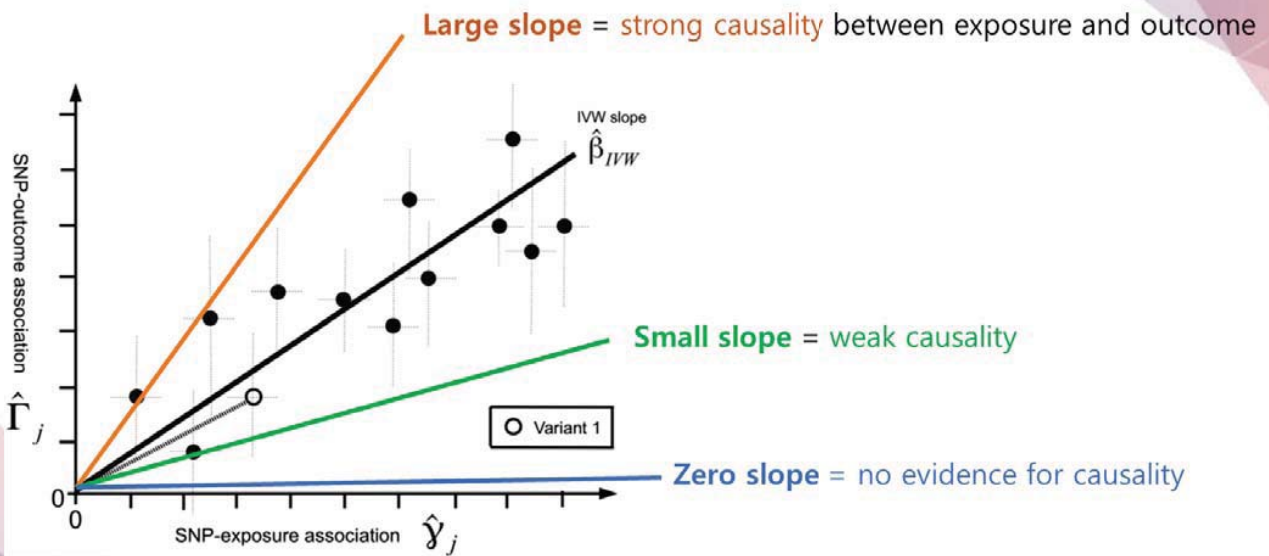
Mendelian randomization (MR)



$$\text{Effect of exposure on disease} = \frac{\text{Effect of IV on disease } (\beta_{xz})}{\text{Effect of IV on exposure } (\beta_{xy})}$$

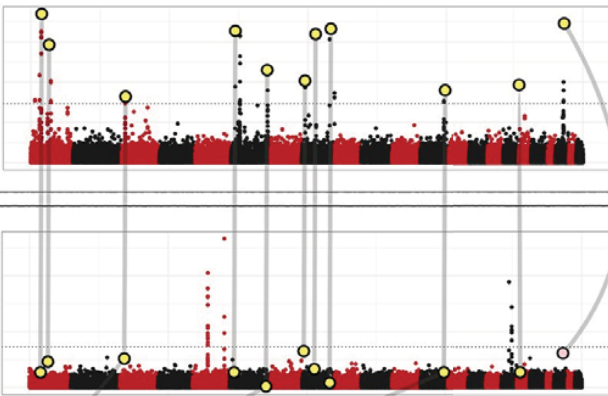
(β_{yz})

Two-sample MR

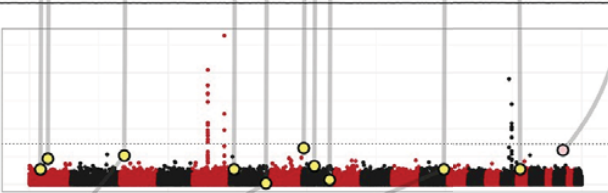


$$\text{Effect of exposure on disease} = \frac{\text{Effect of IV on disease } (\beta_{xz})}{\text{Effect of IV on exposure } (\beta_{xy})} (\beta_{yz})$$

Obtain instruments from exposure GWAS



Extract SNP effects from outcome GWAS



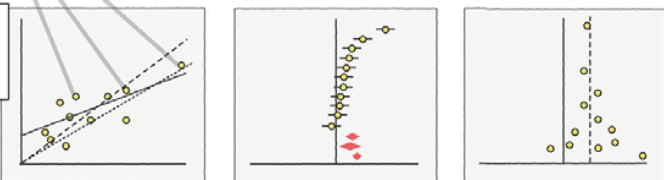
LD Proxies
If an exposure instrument is not available in the outcome GWAS then look for LD proxies in 1000 genomes

Harmonise exposure and outcome effects

SNP	Exposure GWAS				Outcome GWAS			
	Effect	Effect allele	Other allele	Effect allele frequency	Effect	Effect allele	Other allele	Effect allele frequency
rs12345	0.132	A	G	0.28	0.022	A	G	0.26
rs23456	-0.485	G	T	0.41	0.056	T	G	0.61
rs34567	0.203	G	C	0.11	-0.046	G	C	0.88

SNP	Exposure GWAS				Outcome GWAS			
	Effect	Effect allele	Other allele	Effect allele frequency	Effect	Effect allele	Other allele	Effect allele frequency
rs12345	0.132	A	G	0.28	0.022	A	G	0.26
rs23456	-0.485	G	T	0.41	-0.056	G	T	0.39
rs34567	0.203	G	C	0.11	0.046	G	C	0.12

MR estimates and sensitivity analyses



Which modifiable risks are causally associated with AD?

RESEARCH

OPEN ACCESS **Modifiable pathways in Alzheimer's disease: Mendelian randomisation analysis**

Suzanna C Larsson,^{1,2} Matthew Taylor,² Basim Malik,³ Martin Dichgans,^{3,4,5} Stephen Burgess,^{6,7} Hugh S Markus,⁷ for the CoSTRAM Consortium, on behalf of the International Genomics of Alzheimer's Project

ABSTRACT
OBJECTIVE To determine which potentially modifiable risk factors, including socioeconomic, lifestyle/dietary, cardiometabolic, and inflammatory factors, are associated with Alzheimer's disease.
DESIGN Mendelian randomisation study using genetic variants associated with the modifiable risk factors as instrumental variables.
SETTING International Genomics of Alzheimer's Project.
PARTICIPANTS 57 008 cases of Alzheimer's disease and 37 154 controls.
MAIN RESULTS Odds ratio of Alzheimer's per genetically predicted increase in each modifiable risk factor estimated with Mendelian randomisation analysis.
RESULTS This study included analysis of 26 potentially modifiable risk factors. A Bonferroni corrected threshold of $P < 0.002$ was considered to be significant, and $P < 0.05$ was considered suggestive of evidence for a potential association. Genetically predicted educational attainment was significantly associated with Alzheimer's. The odds ratios were 0.89 (95% confidence interval 0.84 to 0.93), $P = 2.4 \times 10^{-6}$ per year of education completed and 0.74 (0.63 to 0.86), $P = 8.0 \times 10^{-5}$ per year of college/university. The contrasted trait intelligence had a suggestive association with Alzheimer's (per genetically predicted 1 SD higher intelligence: 0.73, 0.57 to 0.93, $P = 0.01$). There was suggestive evidence for potential associations between genetically predicted higher quantity of smoking (per 10 cigarettes a day: 0.84, 0.69 to 0.99, $P = 0.04$) and 25-hydroxyvitamin D concentrations (per 30% higher levels: 0.92, 0.85 to 0.99, $P < 0.001$) and lower odds of Alzheimer's and between higher coffee consumption (per one cup a day: 1.26, 1.05 to 1.51, $P = 0.01$) and higher odds of Alzheimer's. Genetically predicted alcohol consumption, serum folate, serum vitamin B₁₂, homocysteine, cardiometabolic factors, and C-reactive protein were not associated with Alzheimer's disease.
CONCLUSION These results provide support that higher educational attainment is associated with a reduced risk of Alzheimer's disease.
Introduction Alzheimer's disease is the leading cause of dementia. The final hallmarks are amyloid plaques and neurofibrillary tangles.⁸ The amyloid cascade hypothesis implies that accumulation of amyloid β impairs neuronal dysfunction and cell death in the brain.⁹ An alternative theory—the vascular hypothesis—implies cerebral hypoperfusion as the primary impetus that drives oxidative stress, deposition of amyloid β , neuroinflammation, blood brain barrier breakdown, cognitive decline, and neurodegeneration.¹⁰
 Apart from increasing age and the apolipoprotein E (APOE) $\epsilon 4$ allele, the causes of Alzheimer's disease are largely unknown, and treatment trials have been disappointing.¹¹ This has led to increasing interest in the potential for reducing Alzheimer's by targeting modifiable risk factors. Conventional observational studies have consistently shown that low educational attainment is associated with an increased risk,¹² and it has been estimated that 10% of cases are potentially attributable to low education.¹³ Incoherence evidence from conventional observational studies includes that obesity, hypertension, and hypercholesterolaemia in middle and older ages, smoking, low vitamin D and folate concentrations, hyperhomocysteinaemia, and high C-reactive protein concentrations are associated with increased risk, whereas physical activity, a healthy diet, moderate alcohol drinking, and coffee consumption are associated with decreased risk (table A in appendix 1).^{14–16} A 2010 state of the science report concluded that there was insufficient evidence to support the association with any modifiable factors with risk.¹⁷ Available evidence in large part

WHAT IS ALREADY KNOWN ON THIS TOPIC
 Conventional observational studies have shown that educational attainment is associated with the risk of Alzheimer's disease.
 Evidence for the associations between lifestyle behaviours and cardiometabolic factors and risk of Alzheimer's disease is inconclusive.
 Available data on modifiable risk factors in relation to Alzheimer's disease are primarily from conventional observational studies, which are vulnerable to confounding and reverse causation bias.

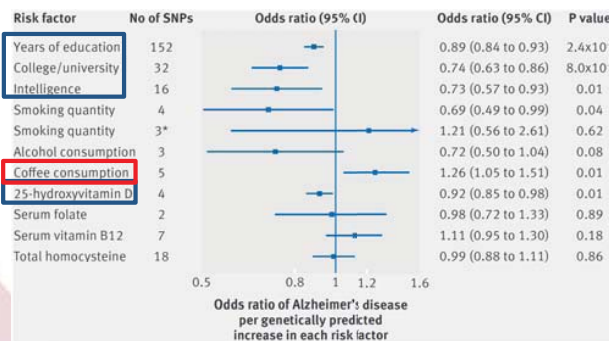
WHAT THIS STUDY ADDS
 A Mendelian randomisation approach shows that a genetic predisposition towards longer education is associated with lower odds of Alzheimer's disease. This study found suggestive evidence of positive associations between higher intelligence, smoking, and concentrations of 25-hydroxyvitamin D and lower odds of Alzheimer's disease and between higher coffee consumption and higher odds of Alzheimer's disease.

thebmj | *BMJ* 2017;355:g5375 | doi:10.1136/bmj.g5375

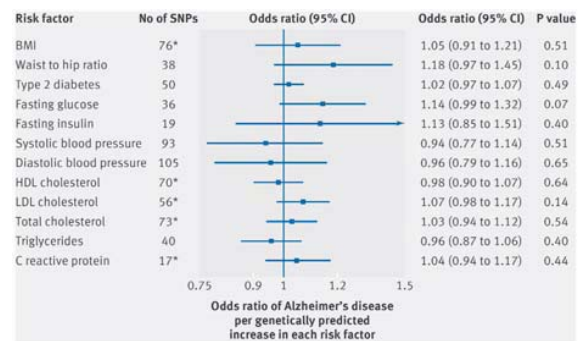
- Modifiable risk factors
 - Selected for the most consistent evidence for an association with Alzheimer's disease in meta-analyses of prospective observational studies
 - 24 socioeconomic, lifestyle/dietary, cardiometabolic, and inflammatory factors were included

Which modifiable risks are causally associated with AD?

Educational attainment, intelligence, and lifestyle and dietary factors



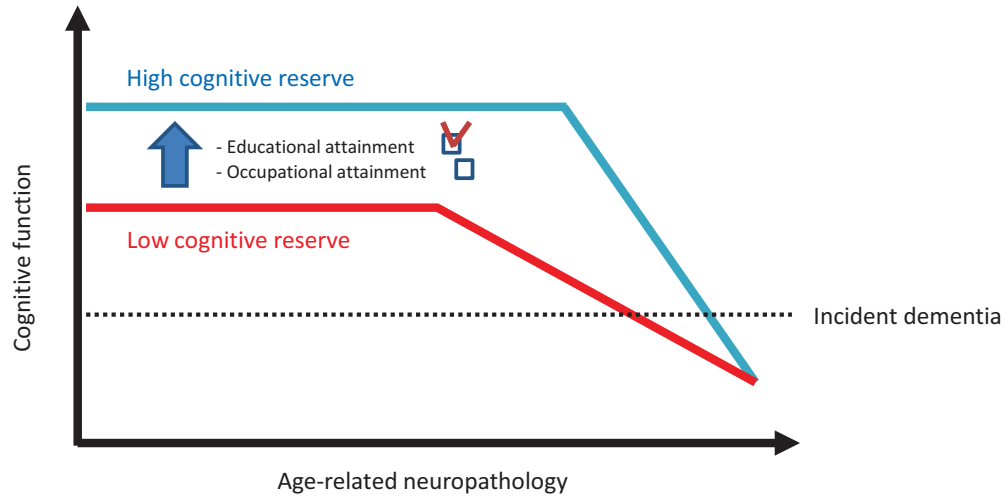
Cardiometabolic and inflammatory factors



CONCLUSION

These results provide support that higher educational attainment is associated with a reduced risk of Alzheimer's disease.

Cognitive reserve



Does occupational attainment also protect against AD?

BRAIN

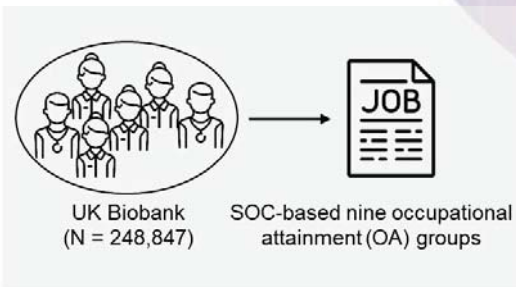
ACCEPTED MANUSCRIPT

Genome-wide association study of occupational attainment as a proxy for cognitive reserve

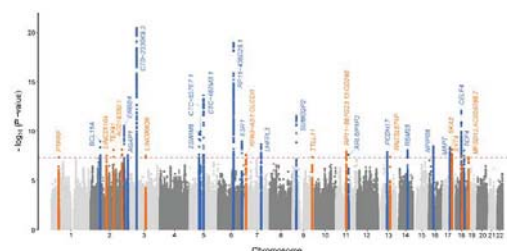
Hyunwoong Ko, Soyeon Kim, Kiwon Kim, Sang-Hyuk Jung, Injeong Shim, Soojin Cha, Hyewon Lee, Beomsu Kim, Joohyun Yoon, Tae Hyon Ha, Seyul Kwak, Jae Myeong Kang, Jun-Young Lee, Jinho Kim, Woong-Yang Park, Kwangsik Nho, Doh Kwan Kim, Woojae Myung ✉, Hong-Hee Won ✉

Brain, awab351, <https://doi.org/10.1093/brain/awab351>

Published: 06 October 2021 [Article history](#)



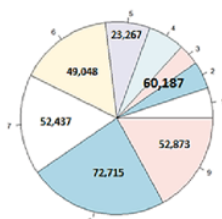
Genome-wide association analysis of OA



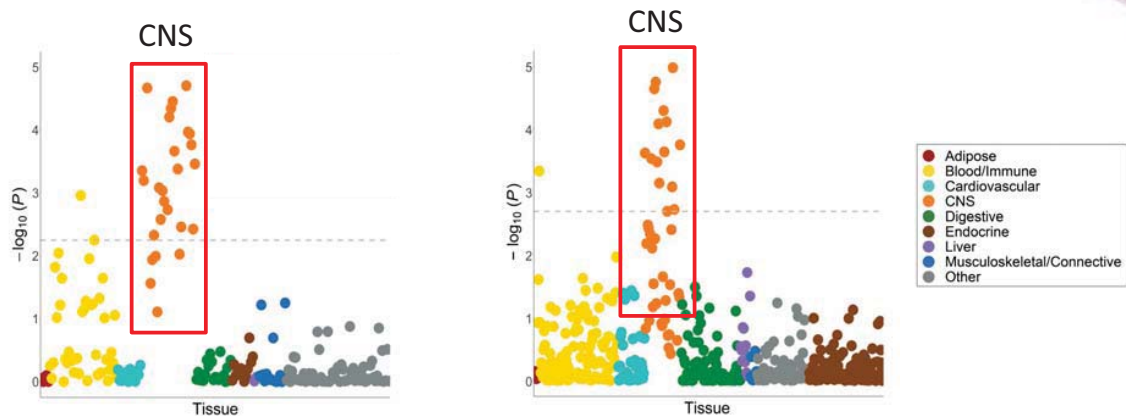
30 significant loci (12 novel variants)
SNP-based heritability: 8.5% (s.e. = 0.4%)

Job Levels

9. Managers and Senior Officials	52,873
8. Professional Occupations	72,715
7. Associate Professional and Technical Occupations	52,437
6. Administrative and Secretarial Occupations	49,048
5. Skilled Trades Occupations	23,267
4. Personal Service Occupations	18,974
3. Sales and Customer Service Occupations	11,077
2. Process, Plant and Machine Operatives	14,179
1. Elementary Occupations	15,957



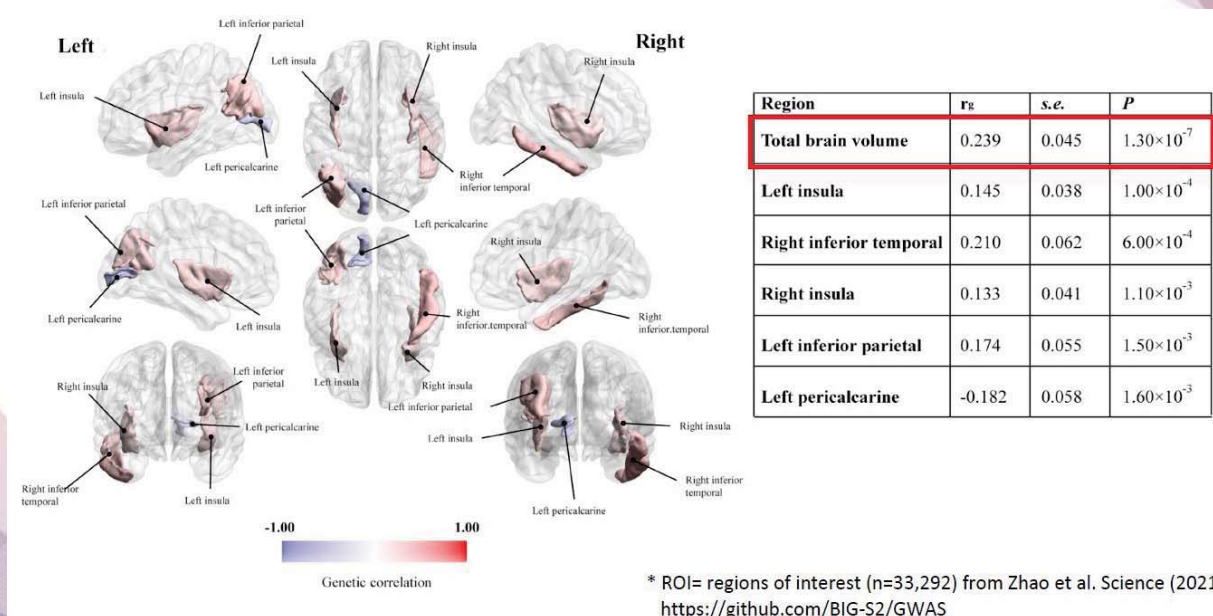
Partitioned heritability was enriched in the central nervous system and brain tissues



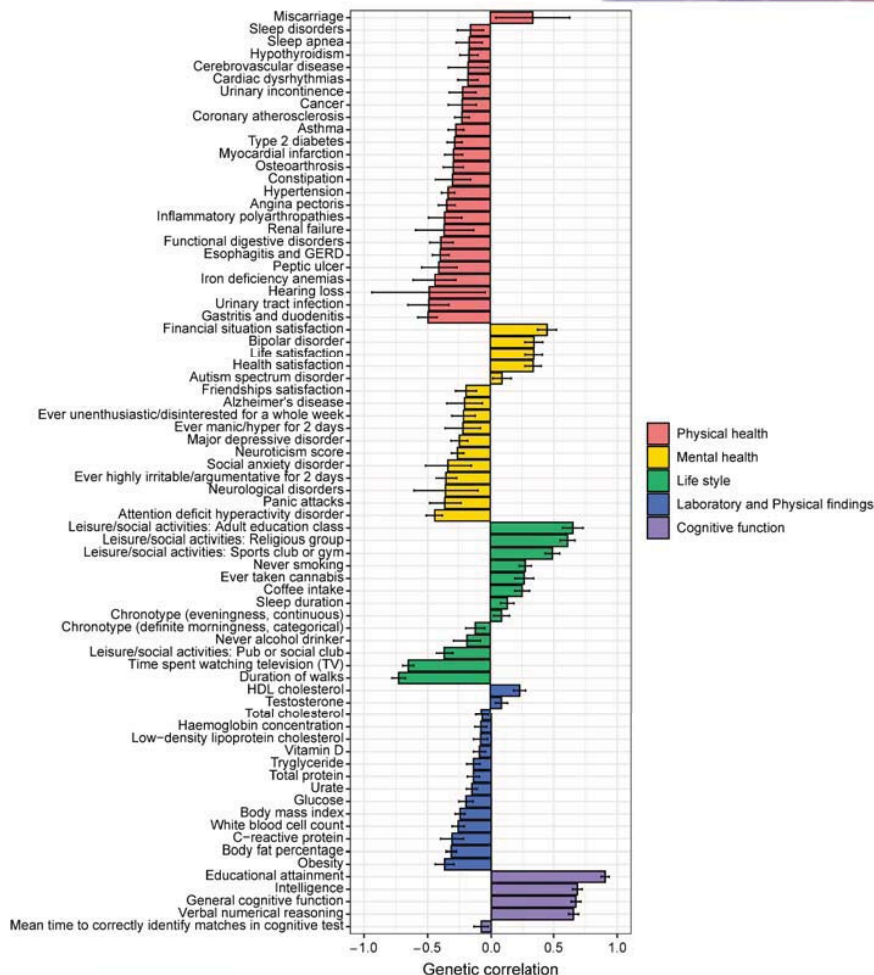
Name	Coefficient	Coefficient std error	Coefficient P value	Coefficient Pval FDR
Neuron	5.12E-09	1.60E-09	0.000664049	0.001992147
Oligodendrocyte	9.08E-10	1.53E-09	0.276438945	0.414658417
Astrocyte	-2.79E-09	1.29E-09	0.984970189	0.984970189

Abbreviation: FDR, false discovery rate

Total brain volume was genetically correlated with occupational attainment



Genetic correlation



MR between occupational attainment and Alzheimer's disease

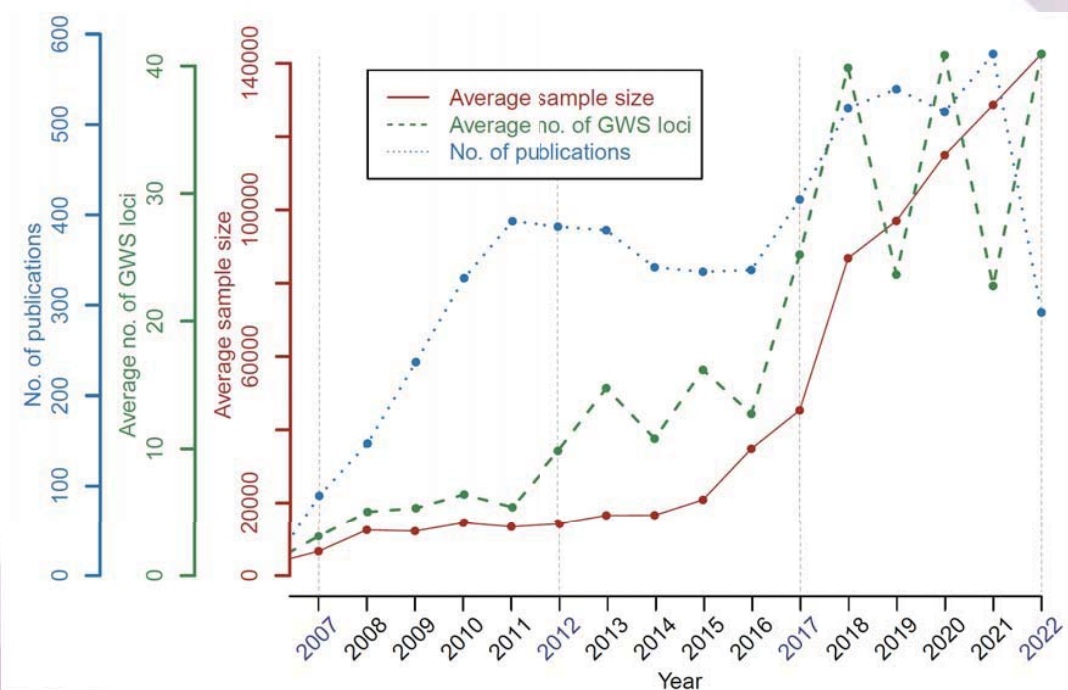
Method	<i>n</i> SNPs	OR (95% CI)	<i>P</i>
Primary MR for occupational attainment on risk of Alzheimers disease			
Inverse variance weighted	18	0.78 (0.65 to 0.92)	4.26×10^{-3}
Weighted median		0.73 (0.57 to 0.92)	9.10×10^{-3}
MR-Egger (P for pleiotropy = 0.90)		0.73 (0.27 to 1.95)	0.54
Sensitivity analysis for occupational attainment on risk of Alzheimer's disease after the exclusion of pleiotropic SNPs			
Inverse variance weighted	11	0.72 (0.57 to 0.91)	5.54×10^{-3}
Weighted median		0.72 (0.53 to 0.97)	0.03
MR-Egger (P for pleiotropy = 0.97)		0.70 (0.12 to 4.00)	0.70
Sensitivity analysis for independent effect of occupational attainment on risk of Alzheimer's disease by multivariate MR controlling for educational attainment			
Exposure: Occupational attainment			
Inverse variance weighted	69	0.72 (0.54 to 0.95)	0.02
Median based		0.68 (0.48 to 0.97)	0.04
MR-Egger (P for pleiotropy ^a = 0.21)		0.63 (0.45 to 0.89)	8.27×10^{-3}
Exposure: Educational attainment			
Inverse variance weighted	69	1.08 (0.62 to 1.91)	0.78
Median based		1.10 (0.53 to 2.30)	0.79
MR-Egger (P for pleiotropy ^a = 0.21)		0.63 (0.23 to 1.74)	0.38

CI = confidence interval; IV = instrumental variable; OR = odds ratio; SNP = single nucleotide polymorphism.

Summary of post-GWAS analysis and tools

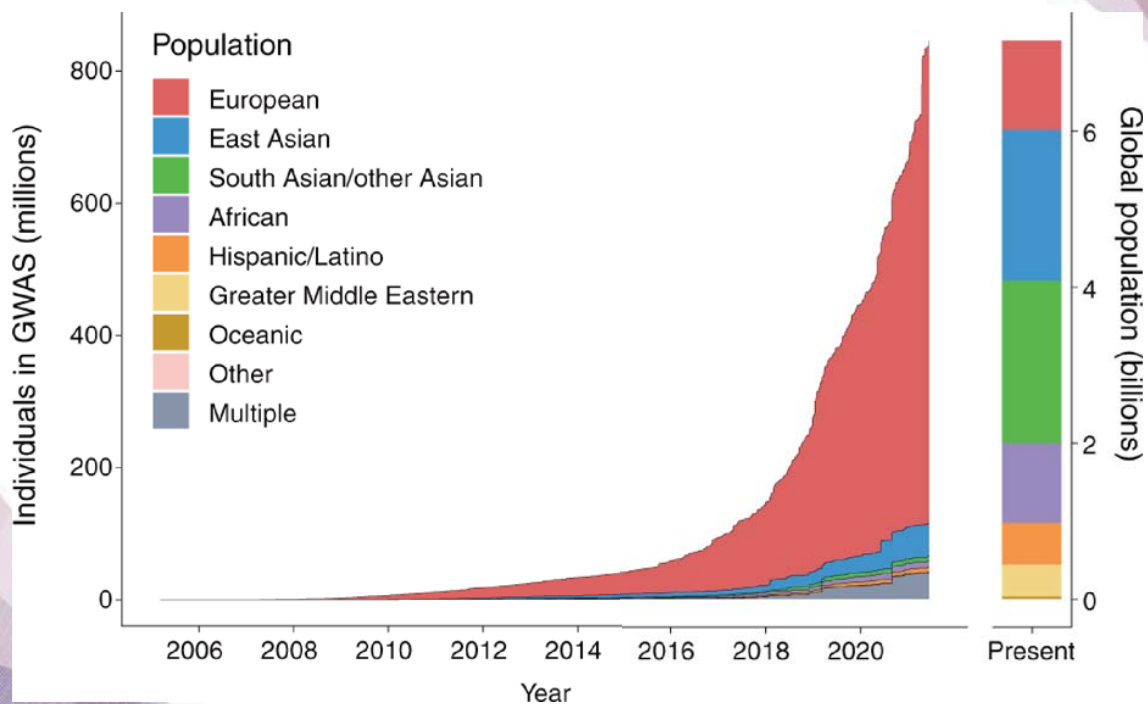
- Understanding genetic architecture
 - SNP-based heritability: LDSC, *GCTA (if genotype available)*
 - Genetic correlation: LDSC (same ancestry), POPCORN or S-LDXR (transancestry)
 - SNP heritability in specific tissues or cells: LDSC-SEG
- Finding causal variants, genes, and pathways
 - Fine-mapping (causal variants): CAVIAR, FINEMAP, PAINTOR, SUSIE
 - eQTL and colocalization analysis (genes): COLOC2
 - Pathway enrichment analysis (pathways or gene sets): MAGMA
- *Identifying individuals at high genetic risk (genotype required)*
 - Polygenic risk score: PRSICE-2, LDPRED, PRS-CS
- Inferring causality between traits
 - Mendelian randomization: MR-BASE, TwoSampleMR (R package)

15 years of GWAS discovery



AJHG 110, 1–16, February 2, 2023

As of June 2021, the vast majority (86%) of genomics studies have been conducted in individuals of European descent



Summary

- Common variants account for a large portion of heritability
- Post-GWAS analyses use GWAS summary statistics that are publicly available
- Post-GWAS analyses reveal the genetic architecture of human traits
- Omics data with GWAS are helpful in identifying target genes
- However, the current imbalance between ancestries may limit the clinical utility of genomics in non-European populations

감사합니다.

honghee.won@gmail.com

성균관대학교 삼성융합의과학원

삼성서울병원