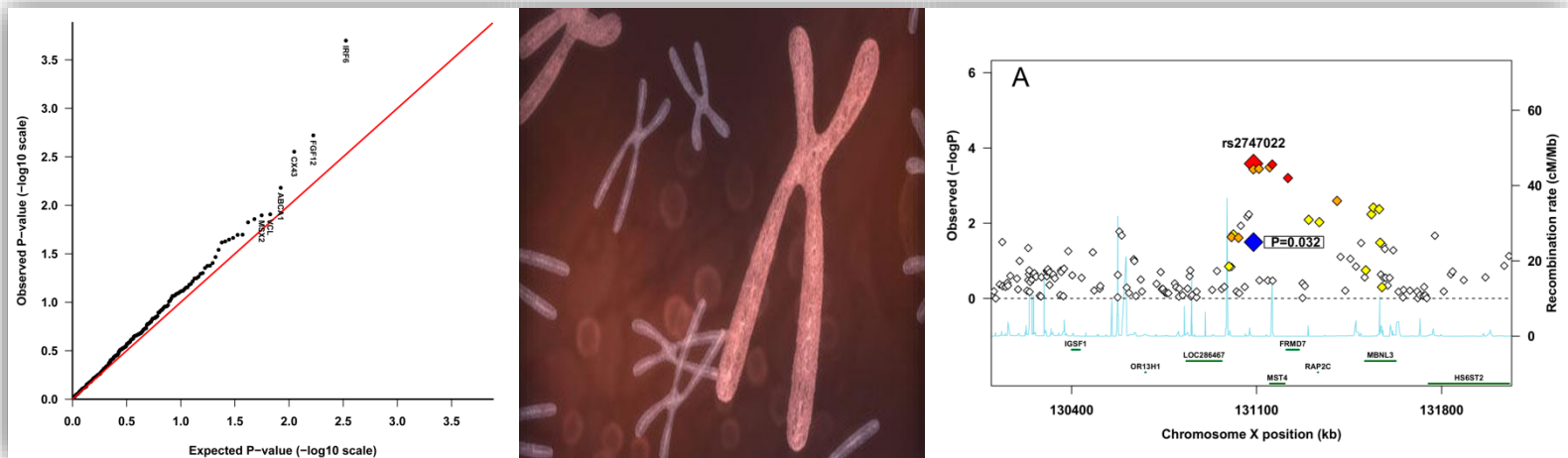


# GENERAL INTRO TO GENETIC EPIDEMIOLOGY

## — LECTURE 1, PARTS 1 & 2 —

Anil Jugessur

Senior scientist, Norwegian Institute of Public Health, Oslo



# LECTURE OUTLINE

## General introduction to genetic epidemiology (lecture I)

- Part I
  - What's a complex trait?
  - Genetic basis of complex traits
- Part II
  - Genetic approaches to studying complex traits
  - Candidate-gene analysis, GWAS, and GWAMA



# LECTURE OUTLINE

## General introduction to genetic epidemiology (lecture I)

### ○ Part I

- What's a complex trait?
- Genetic basis of complex traits

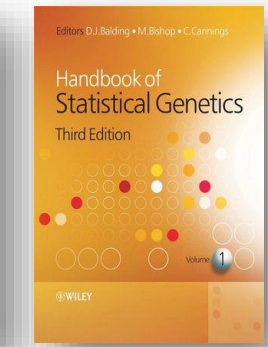
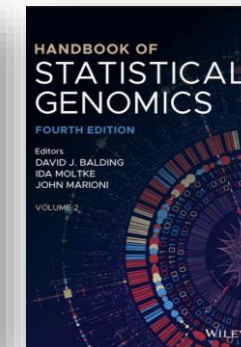
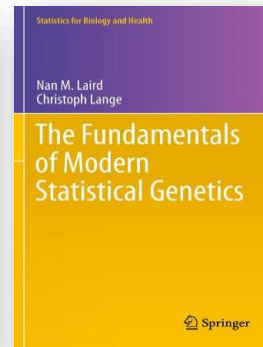
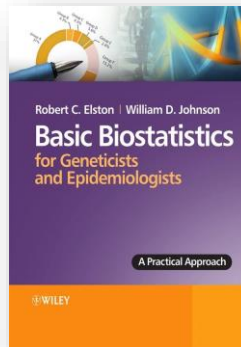
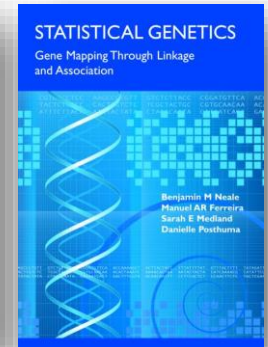
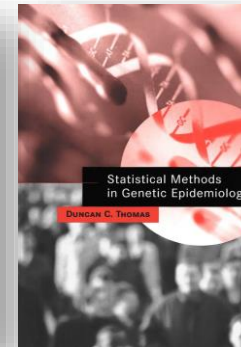
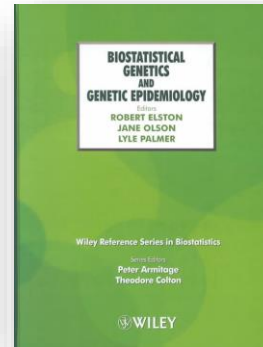
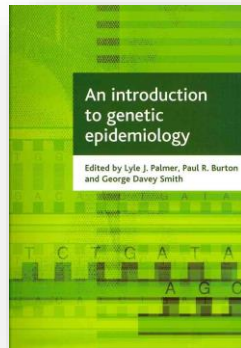
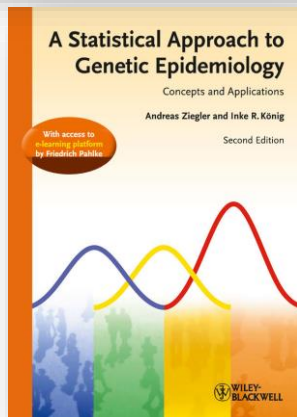
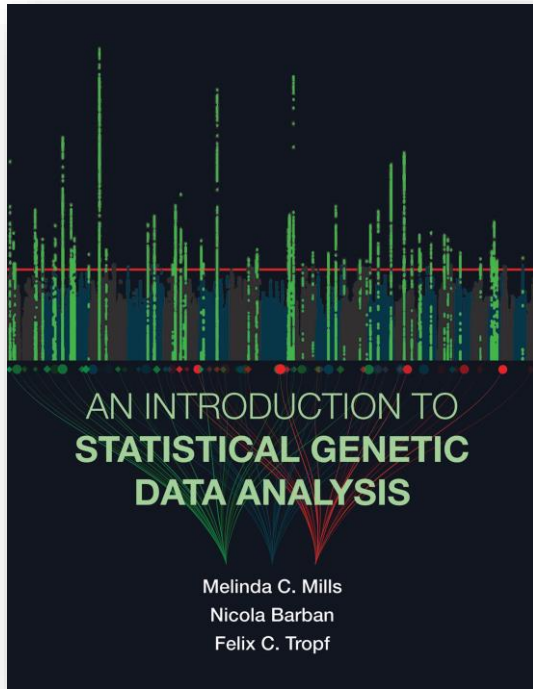
### ○ Part II

- Genetic approaches to studying complex traits
- Candidate-gene analysis, GWAS, and GWAMA



# COURSE BOOK

## Course book

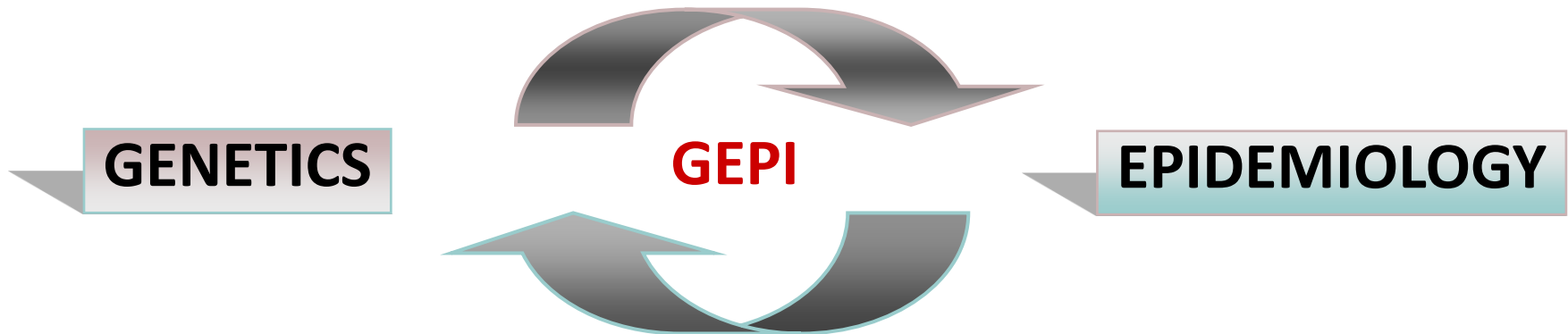


# WHAT IS GENETIC EPIDEMIOLOGY?

---

In broad terms:

«The application of genetic principles and techniques to answering epidemiological questions»



# LOTS OF DEFINITIONS OUT THERE...

**Table 1–1.** Some definitions of genetic epidemiology

*N. E. Morton and C. S. Chung (1978):* “A science that deals with the etiology, distribution, and control of disease in groups of relatives, and with inherited causes of disease in populations.”

*R. Ward (1979):* “The primary objective of the genetic epidemiologist will be to identify the genetic contribution to the etiological pathway.”

*B. H. Cohen (1980):* Genetic epidemiology is defined “as examining the role of genetic factors, along with the environmental contributors to disease, and at the same time, giving equal attention to the differential impact of environmental agents, nonfamilial as well as familial, on different genetic backgrounds.”

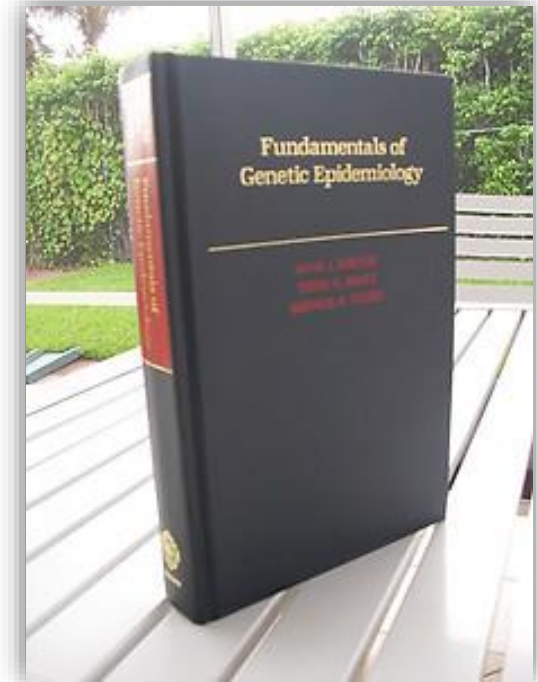
*P. Phillippe (1982):* “Genetic epidemiology studies the interaction between genetic and environmental factors at the origin of disease.”

*M.C. King et al. (1984):* “Genetic epidemiology is the study of how and why diseases cluster in families and ethnic groups.”

*D.C. Rao (1984):* “Genetic epidemiology is an emerging field with diverse interests, one that represents an important interaction between the two parent disciplines: genetics and epidemiology. Genetic epidemiology differs from epidemiology by its explicit consideration of genetic factors and family resemblance; it differs from population genetics by its focus on disease; it also differs from medical genetics by its emphasis on population aspects.”

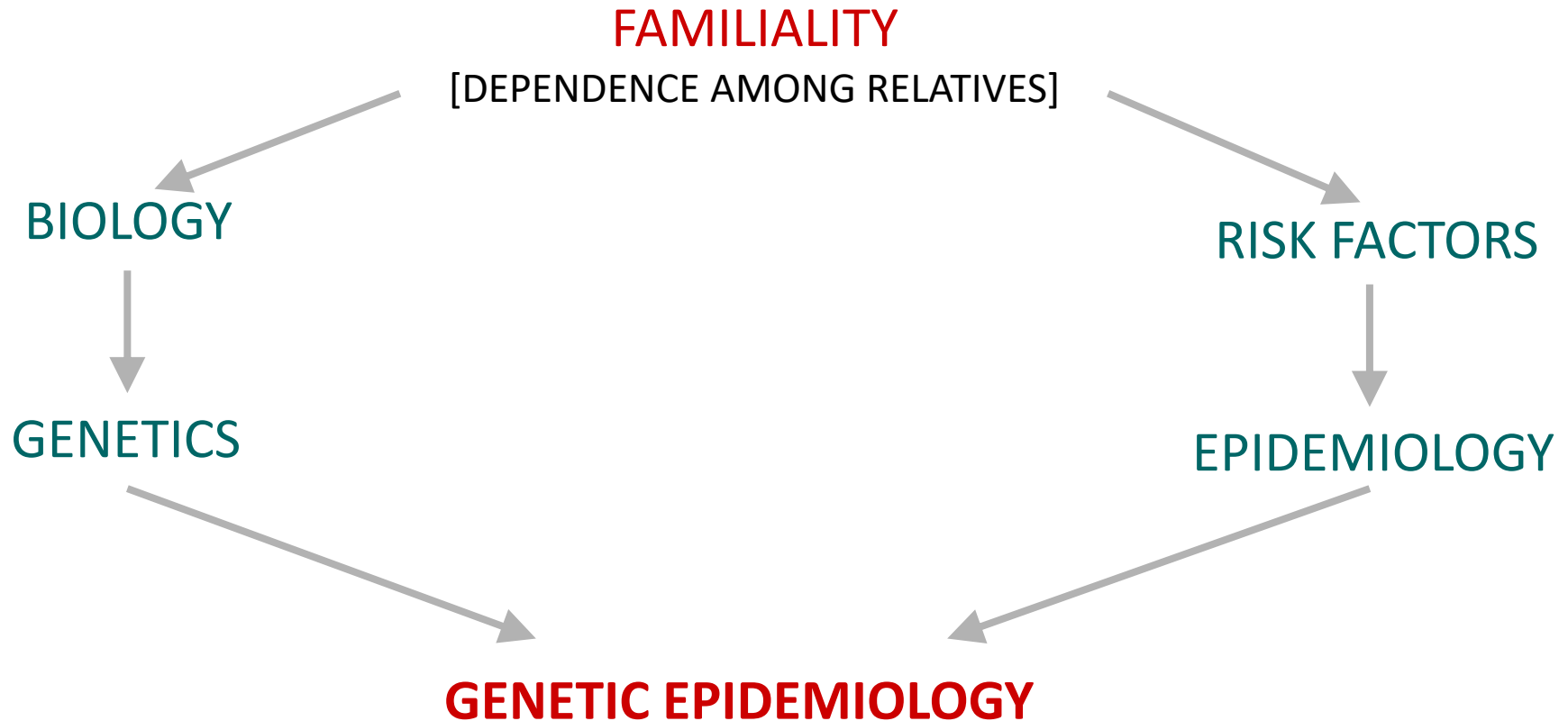
*D.F. Roberts (1985):* argues the distinction of genetic epidemiology from epidemiology in general. Genetic epidemiology “is not merely the application of the central concept of epidemiology, the study of the distribution of disease in space and time, to genetic disease. Instead, in genetic epidemiology, the concept is extended to include the additional variables of the genetic structure of the population, with the object of elucidating the etiology of disease in which there may be a genetic component.”

*E.A. Thompson (1986a):* “Genetic epidemiology is the analysis of the familial distributions of traits, with a view to understanding any possible genetic basis.”



Prof of biostats @ WASH-U.

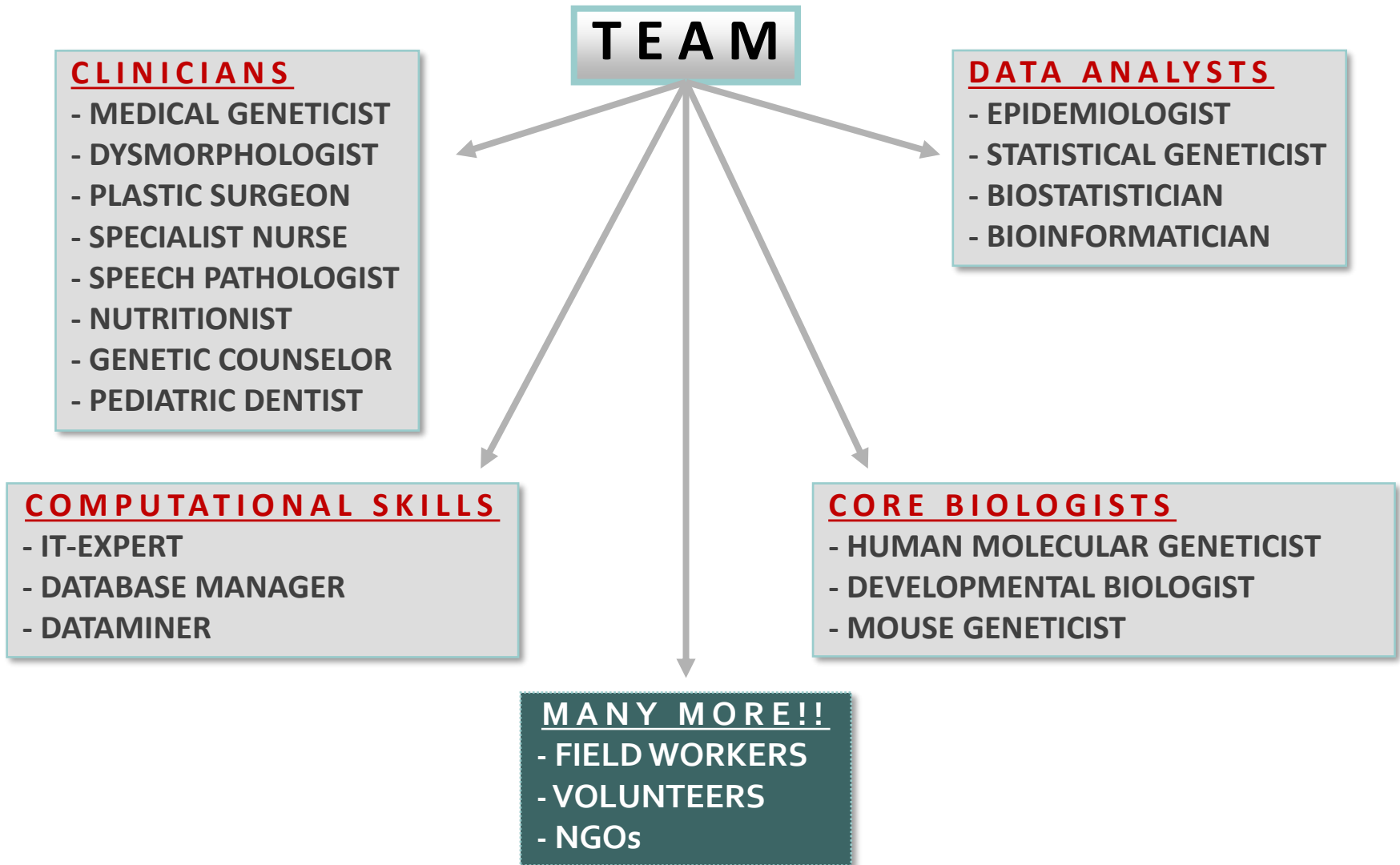
# GROWING CONVERGENCE OF DIFFERENT FIELDS



"Less divergence in terminology and methodology, but an increased conversation, collaboration and convergence across the fields."

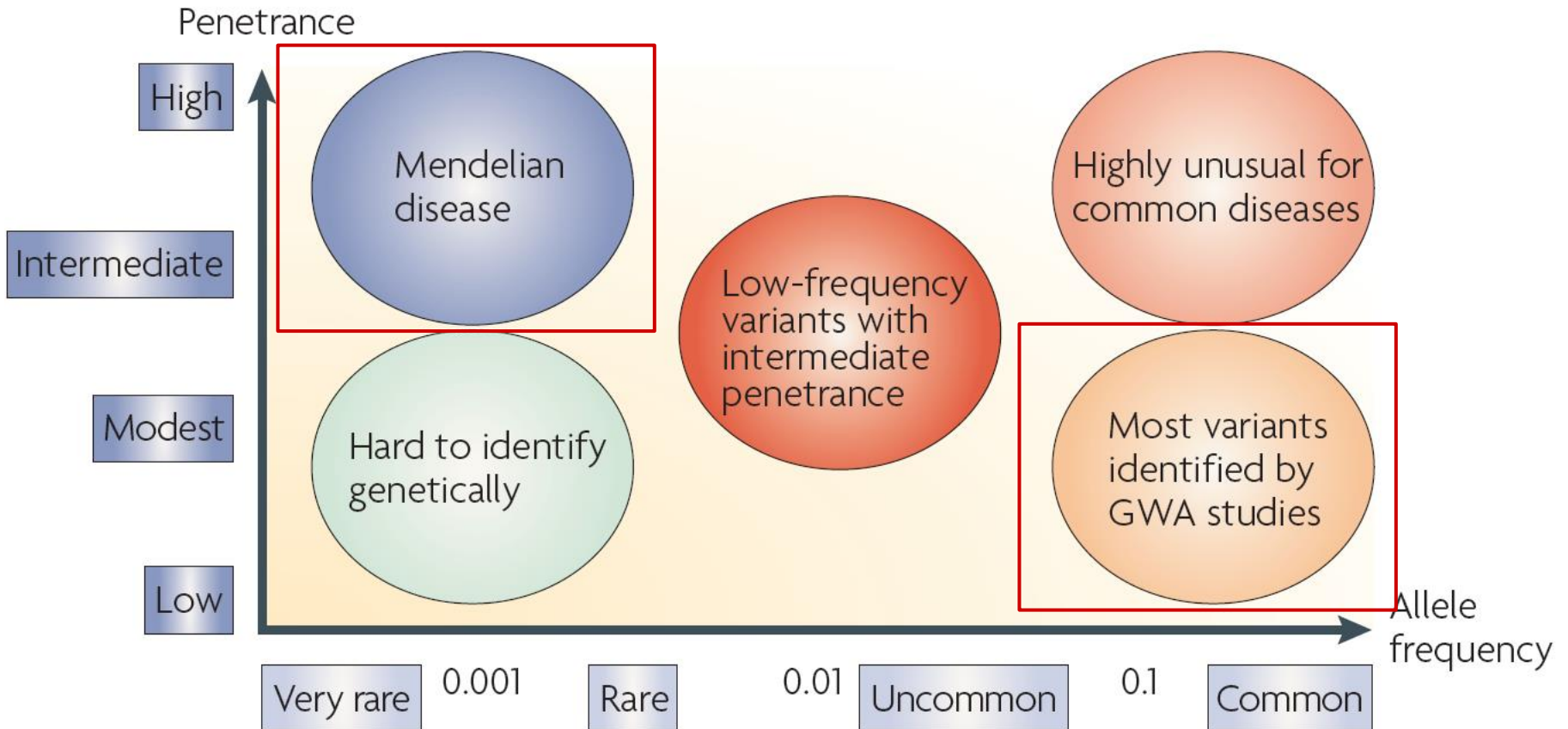
# BUILDING A TEAM

– E.G. FOR A STUDY OF BIRTH DEFECTS –





# Mendelian disorder vs. Complex trait



# WHAT'S A COMPLEX TRAIT?

## Different categories of disease causation

TRAIT «Purely environmental»

Non-genetic

**Multifactorial/Complex**

Polygenic

Oligogenic

Digenic

Monogenic

Chromosomal

***«Purely genetic»***

Etiological spectrum

Physical, chemical, nutritional

Infectious agents / pathogens

Multiple genes and environmental factors

Multiple genes with small additive effects

A few genes with large effects

Mendelian / single-gene disorders

Abnormal chromosome numbers and structural causes

# COMMON FEATURES OF COMPLEX TRAITS

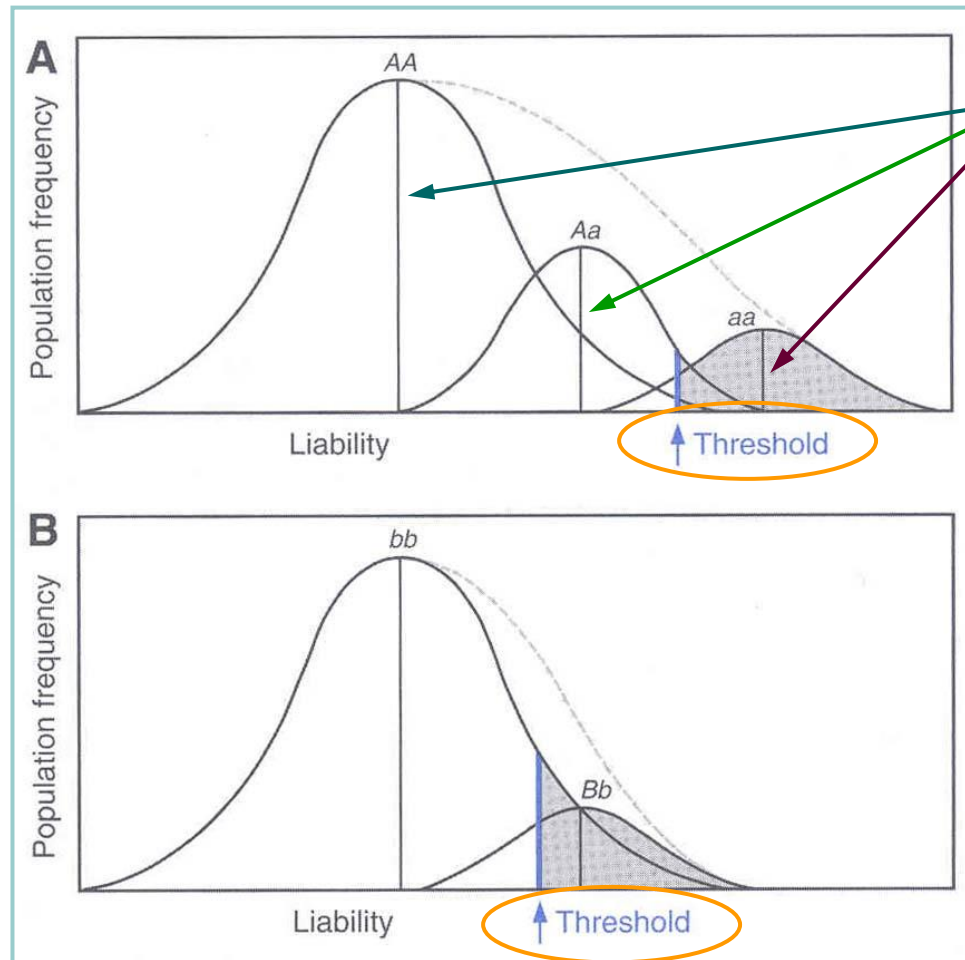
---

- Unlike Mendelian diseases, complex traits are relatively common
- Heterogeneity at several levels:
  - Genetic heterogeneity:
    - «locus» and «allelic» heterogeneity
- Incomplete penetrance  $\Rightarrow$  *not all individuals with the mutant genotype express the phenotype*
- Effect of a gene can be masked by:
  - Phenocopies  $\Rightarrow$  *an environmentally-caused phenotype mirrors a genetically-caused trait*
  - Pleiotropy  $\Rightarrow$  *the mutant genotype affects different traits or organs*
- Complex interactions:
  - «gene-gene» and «gene-environment» interactions
- Stochastic effects  $\Rightarrow$  *random or chance events; biological processes are error-prone!*

# THE CONCEPT OF «LIABILITY»

Liability is an underlying continuous variable comprising both genetic and non-genetic effects.

FIGURE: An idealized distribution of liability in individuals with various genotypes.



Mean liability for each genotype

**Threshold** = value in the liability that determines whether a disease will be expressed or not.

Anyone with liability greater than the threshold manifests the disease.

Recessive allele  
'a' ↑ ses liability

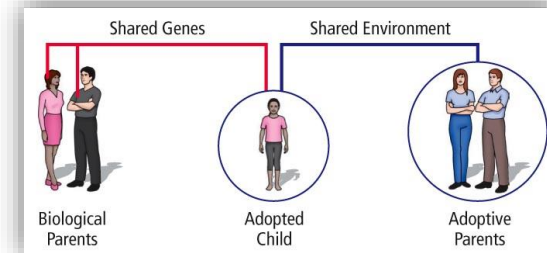
Dominant allele  
'B' ↑ ses liability

# THE CONCEPT OF «HERITABILITY» - CH. 1

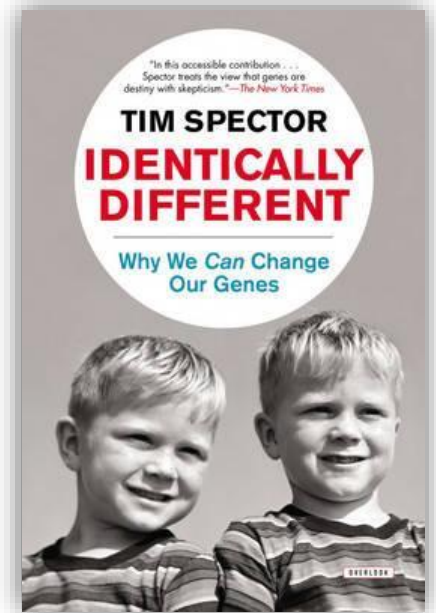
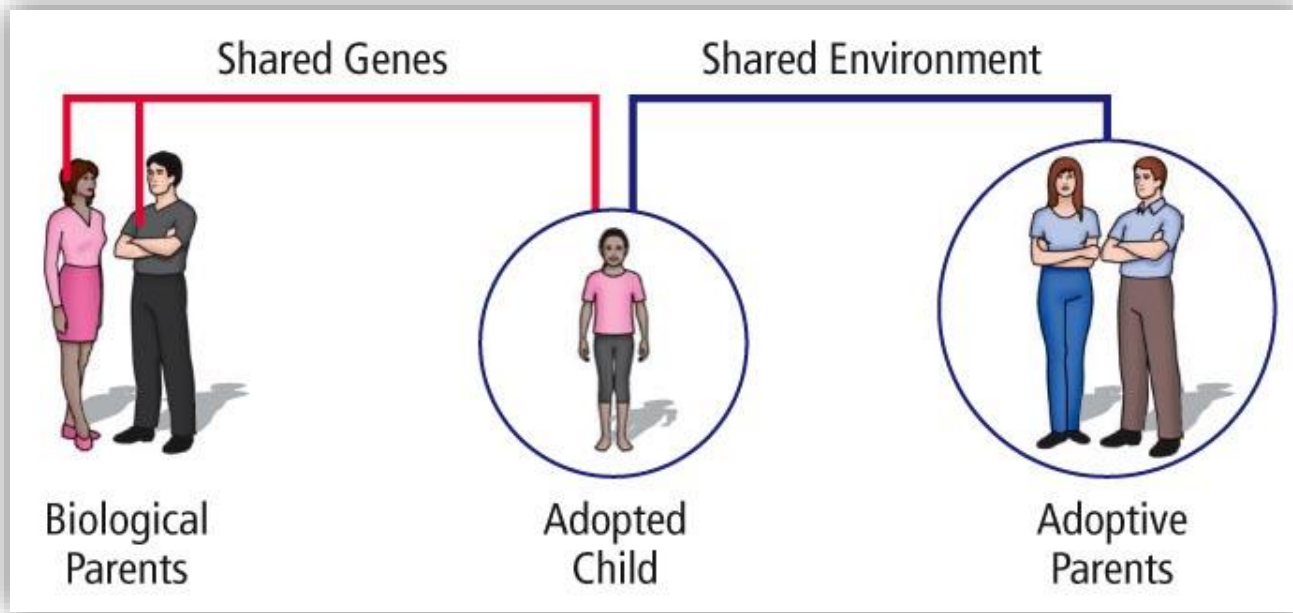
- Heritability ( $H^2$ ) is the proportion of phenotypic variance ( $V_p$ ) attributable to genetic differences.
- Broad-sense vs. narrow-sense heritability
  - **Broad-sense heritability** is the proportion of variance in a phenotype ( $V_p$ ) attributable to the total genetic variance ( $V_g$ ).  $H^2 = V_g/V_p$ , where  $V_p = V_g + V_e$
  - **Narrow-sense heritability** is the proportion of  $V_p$  attributable to additive genetic variance ( $V_a$ ); i.e.,  $H^2 = V_a/V_p$
- Additive vs. non-additive genetic effects
  - **Additive effects**: 2 or more genes contribute to a phenotype, or when alleles in a single gene combine so that their combined effects on the phenotype are equal to the sum of their individual effects.
  - **Non-additive effects** can be dominance ( $V_d$ ) or epistasis ( $V_i$ )
    - **Dominance**: The effect of one allele masks the effect of a second allele at the same locus; e.g., allele  $A$  dominates allele  $a$ .
    - **Epistasis**: An allele at one locus affects the expression of another allele at a different locus.

# IS THERE A GENETIC BASIS TO COMPLEX DISEASES?

- Study whether the disease clusters in families:
  - Familial aggregation studies:
    - Relatives share a greater proportion of their alleles
      - Affected individuals will tend to cluster in families.
    - Recurrence risk measured as relative risk ratio ( $\lambda_r$ )
      - $\lambda_r = [\text{risk to relatives of type } r] \div [\text{Population risk}]$
    - Cannot establish that the disease is hereditary
      - Environmental factors could also cause this clustering!
  - Adoption studies:
    - If a trait has a genetic influence, the risk of disease should be higher in biological relatives than in adopted relatives living in the same household.
  - Twin studies:
    - Compare concordance in MZ vs. DZ twins
      - If MZ twins show close to 100% concordance but DZ twins show significantly less:  $\Rightarrow$  the trait has a strong genetic basis.
      - If MZ twins shows moderate concordance (40-60%) but still significantly higher than DZ twins  $\Rightarrow$  both environmental and genetic components are likely involved in the disease.



# IMPORTANCE OF SHARED ENVIRONMENT!



# ASSESSING EVIDENCE OF FAMILIAL AGGREGATION

Usual to look at two types of correlations between relative pairs:

- «INTER»class correlation
  - Involves two different classes of relatives:
    - E.g. husband-wife, parent-offspring, brother-sister, grandparent-grandchild, etc.
- «INTRA»class correlation
  - Involves only a single class of relatives:
    - E.g. brother-brother, sister-sister, etc.





# AN EXAMPLE

Fingerprint data: count the number of ridges to explore degree of familiarity.

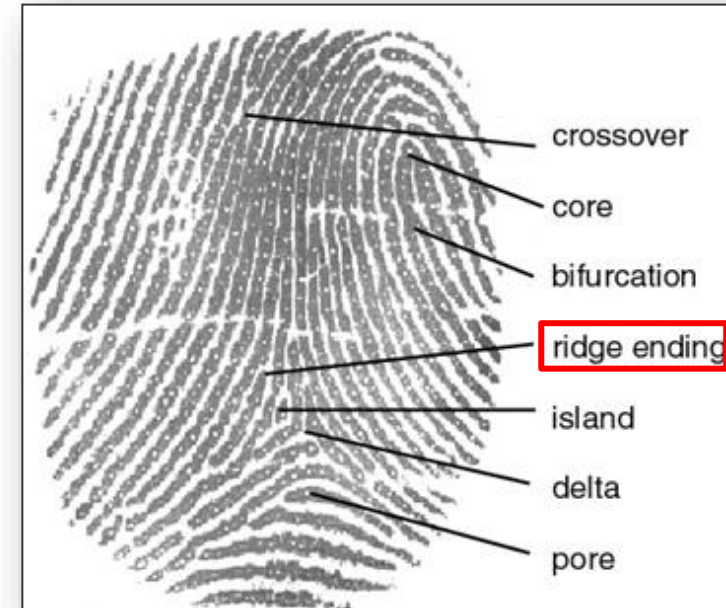
## 2 scenarios:

### ○ Dataset I:

- Parent-offspring correlation:  $0.48 \pm 0.04$
- Sibling correlation:  $0.50 \pm 0.04$
- Spouse correlation:  $0.05 \pm 0.07$

### ○ Dataset II:

- Parent-offspring correlation:  $0.22 \pm 0.01$
- Sibling correlation:  $0.39 \pm 0.01$
- Spouse correlation:  $0.15 \pm 0.02$



# AN EXAMPLE – CONTD...

- **Dataset I:**
  - **Parent-offspring** correlation: **0.48 ± 0.04**
  - **Sibling** correlation: **0.50 ± 0.04**
  - **Spouse** correlation: **0.05 ± 0.07**
- **Dataset II:**
  - **Parent-offspring** correlation **0.22 ± 0.01**
  - **Sibling** correlation: **0.39 ± 0.01**
  - **Spouse** correlation: **0.15 ± 0.02**

- Positive correlation coefficients suggest familial aggregation for this trait
- Strong degree of familiarity in **Dataset I**.
  - Sibling correlation is slightly higher than parent-offspring correlation
    - Consistent with siblings sharing more of their environment than parents & offspring
  - Don't see same degree of correlation in the spouse group
    - Consistent with a less genetic sharing between spouses.
- In **Dataset II**, higher spouse correlation may be due to shared spousal environment (perhaps some assortative mating..?)
- Overall, there seems to be stronger environmental influences in Dataset II.

# LECTURE OUTLINE

## General introduction to genetic epidemiology (lecture I)

### Part I

What's a complex trait?

Genetic basis of complex traits

### ○ Part II

- Genetic approaches to studying complex traits
- Candidate-gene analysis, GWAS, and GWAMA



# GENETIC APPROACHES TO INVESTIGATING A COMPLEX TRAIT

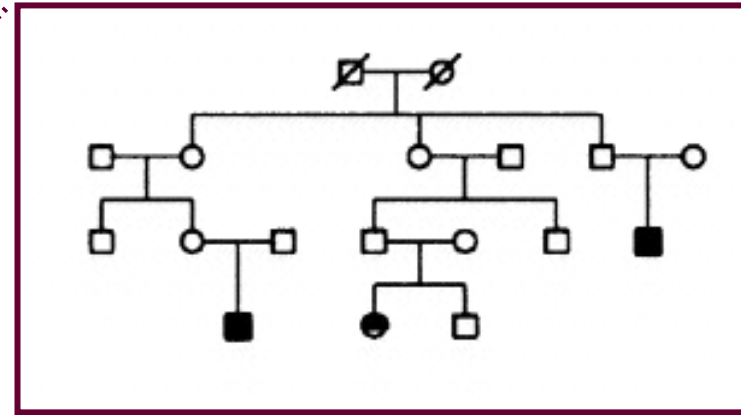
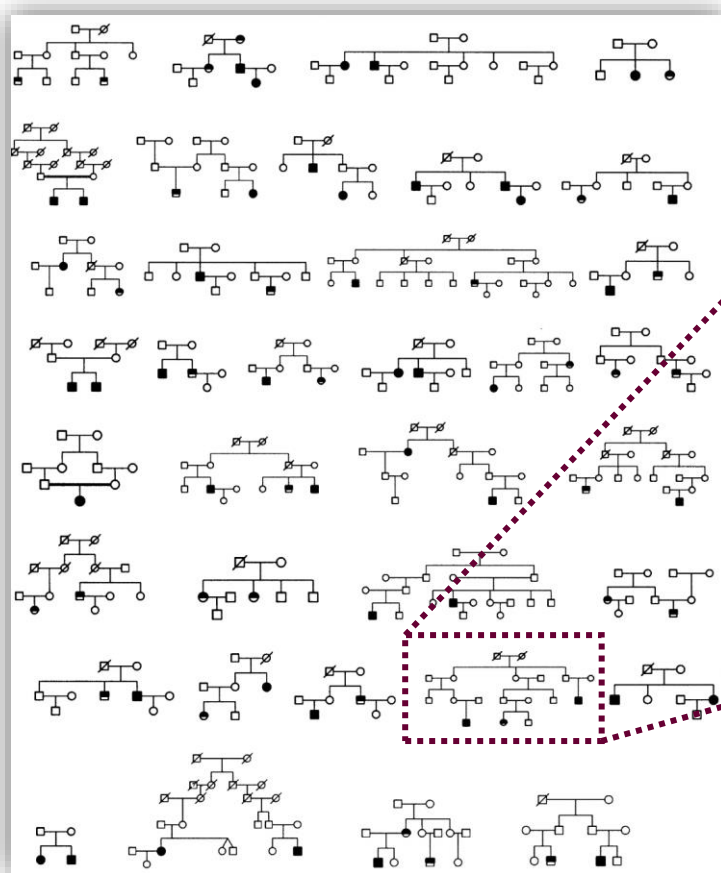
---

Once we have found evidence for a genetic component:

- **Linkage studies** in families with multiple affected members ('multiplex')
  - Test for cosegregation of a marker with the disease to see if the genetic marker and disease gene are physically linked
  - Problematic for complex diseases because of a lack of multiplex families
- **Allele-sharing studies** in affected relative pairs
  - Apply model-free methods on smaller subunits within multiplex families
  - «Identity by descent» (IBD) methods
    - Knowledge of transmission not required (non-parametric, or model-free)
    - Reasonable power to detect genes of fairly modest effects
- **Linkage disequilibrium** approaches
  - Exploit how genetic markers are correlated on chromosomes.

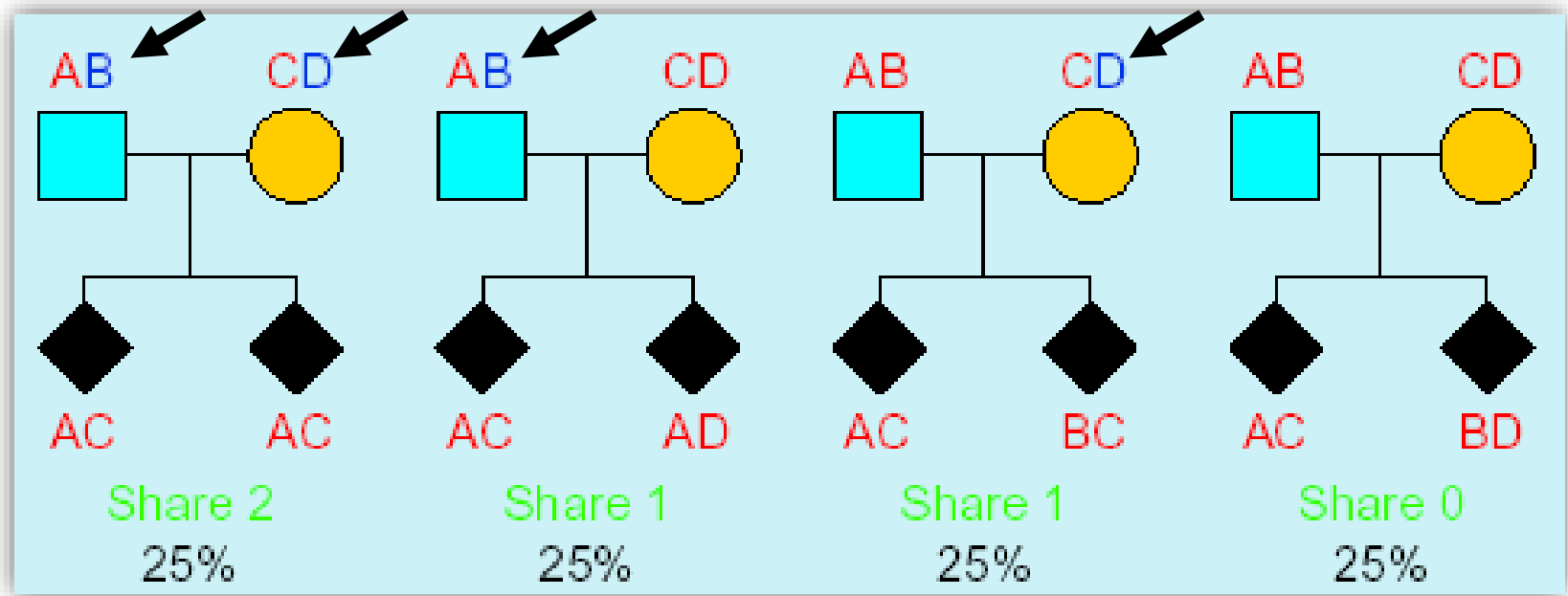
# LINKAGE STUDIES IN MULTIPLEX FAMILIES

Genomewide linkage analyses can be performed using around 400 microsatellite markers distributed with an average spacing of 10 cM for genomewide coverage.



# ALLELE-SHARING STUDIES

**Main idea:** If affected pairs inherit a particular chromosomal fragment more often than would be expected by chance alone – this shows linkage!



2 (25 %)

1 (50 %)

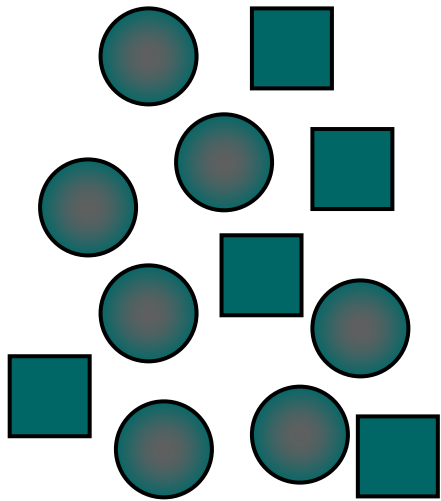
0 (25 %)

No. of parental alleles shared (% of Mendelian proportion)

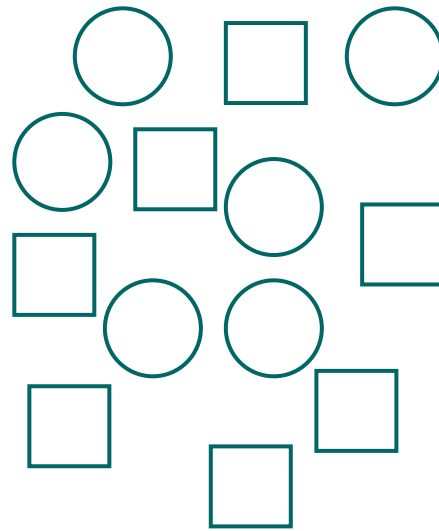
Deviations from these expected proportions  $\Rightarrow$  evidence of linkage

# LINKAGE DISEQUILIBRIUM (LD) APPROACHES

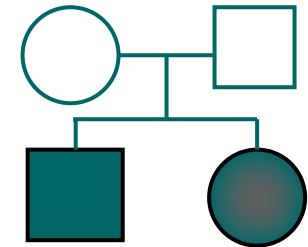
- **Either case-control or family-based**
  - Compare marker allele frequencies between a case and a control population
  - With family data, non-transmitted parental alleles are used as control alleles.
    - Test for deviations from the expected 50% transmission of an allele from parents to offspring.



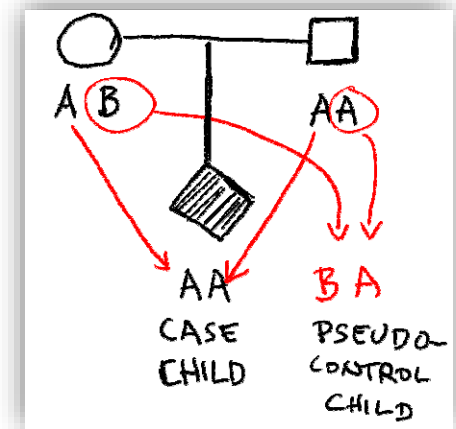
Case (disease)



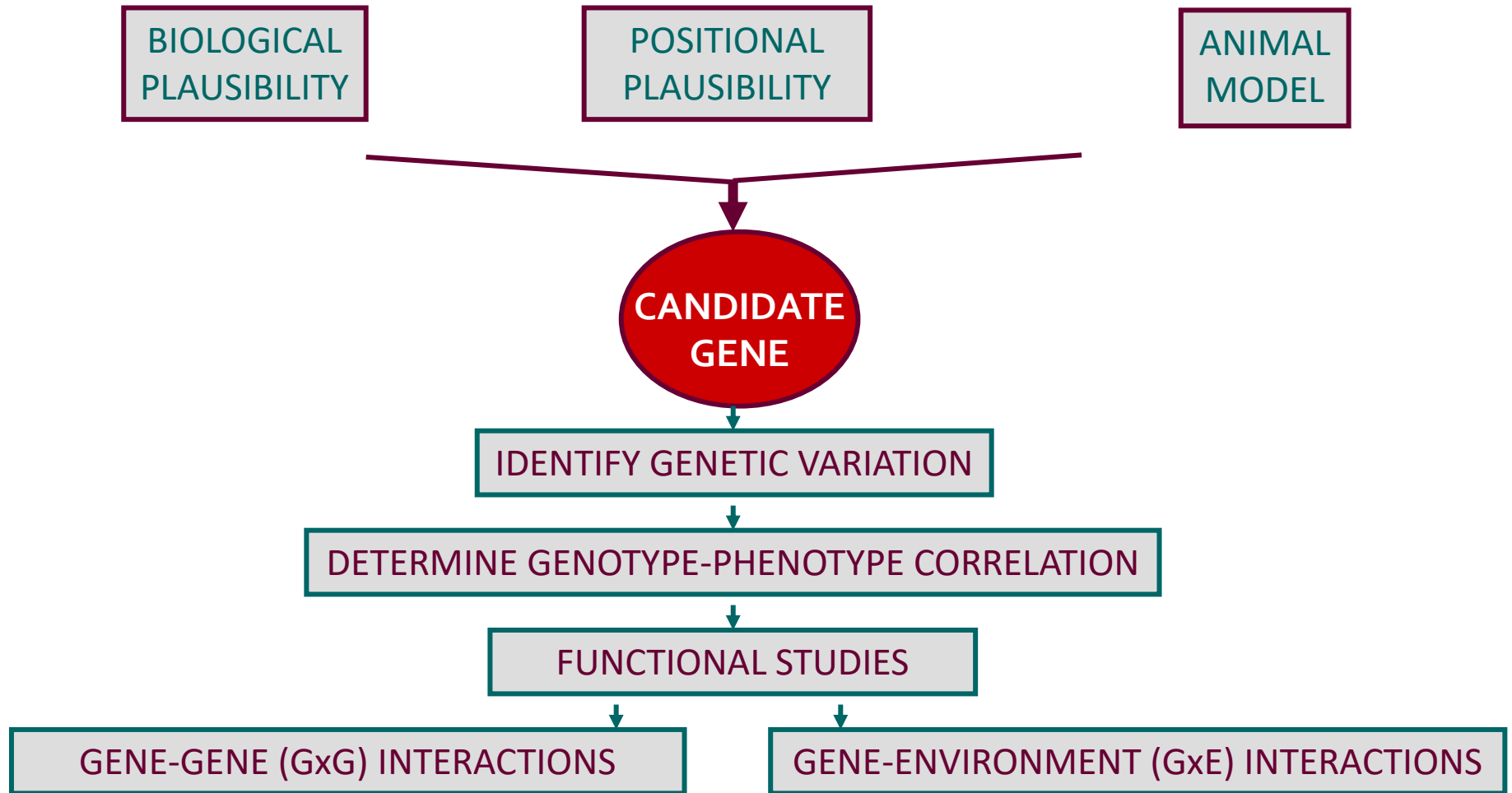
Control (healthy)



VS.



# THE CANDIDATE-GENE APPROACH

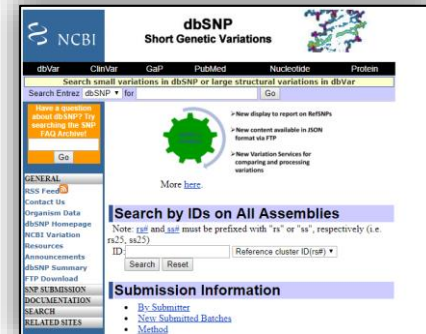
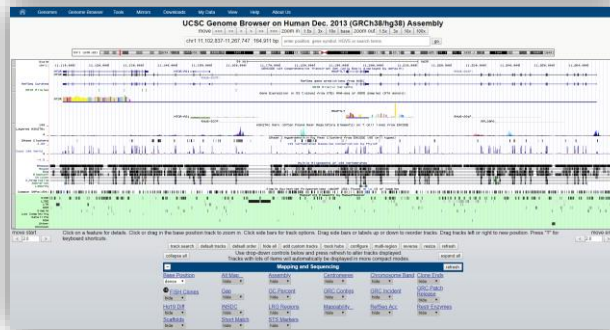
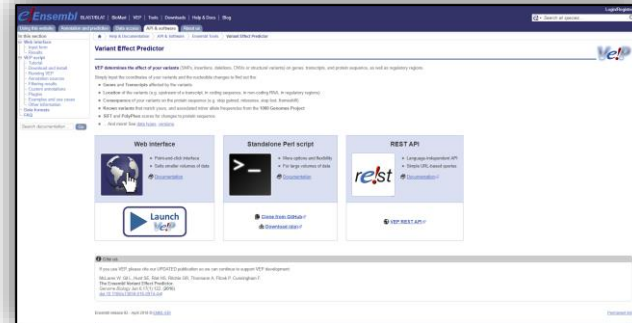
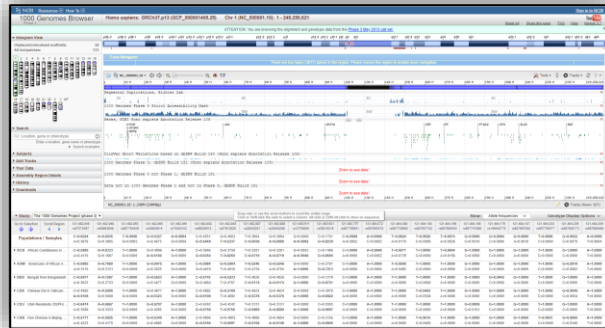




# Selecting SNPs for candidate-gene analysis

## ■ Databases for selection/evaluation of SNPs:

- 1000 Genomes, *e!Ensembl*, UCSC's genome browser, and dbSNP, etc..



## ■ Criteria for prioritizing SNP selection:

- Prior association with the trait being studied
- Minor allele frequency (MAF) of at least 5% to capture common variants
- Preference for coding SNPs and SNPs in regulatory regions – functional!
- SNPs with «haplotype-tagging» properties

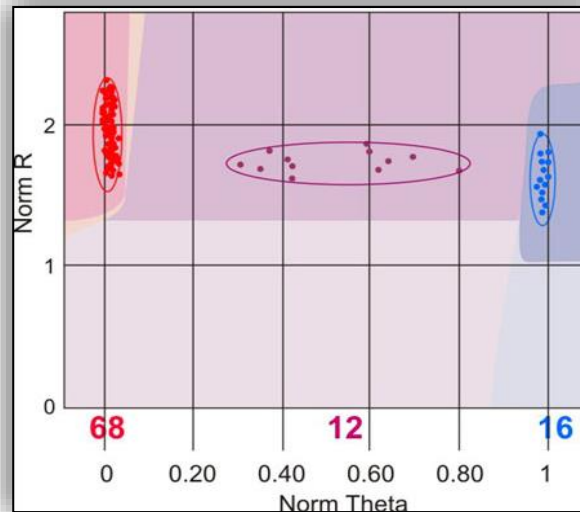
# SNP custom-assay and genotyping

- SNP assays can be designed by ILLUMINA™
  - A customized full panel of X number of SNPs in Y number of candidate genes.
- Outsource the genotyping (and QC) to a core facility: e.g Microarray facility (Oslo), Sanger Institute (UK), DeCode genetics (Iceland), etc..

Illumina iScan system



E.g. of genotype calling



Genomics Core facility Oslo

**Genomics Core Facility**  
Oslo University Hospital and Helse Sør-Øst

Services Courses Resources Publications Contact

Sequencing Microarrays Sample Requirements Analysis

**Proven Solutions. Quality Provider.**

The **Genomics Core Facility** has for almost 20 years provided state-of-the-art laboratory technology and high-throughput genomic services to the South-East Health Region, as well as the Norwegian scientific community. Today, our core facility offers an extensive set of technologies to study genome structure, dynamics and function using Illumina high-throughput sequencing technology and different commercial microarray platforms. In addition to laboratory services, the core facility delivers bioinformatics analysis, providing a comprehensive solution for high-throughput genomic analysis. Our services are equally accessible to all users from our health region, following a first come first served policy.

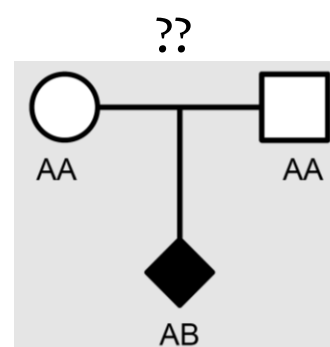
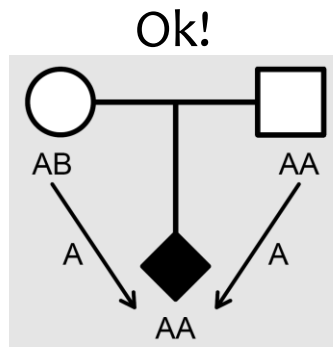
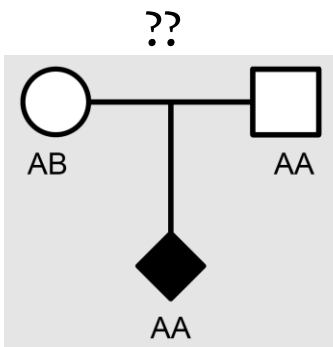
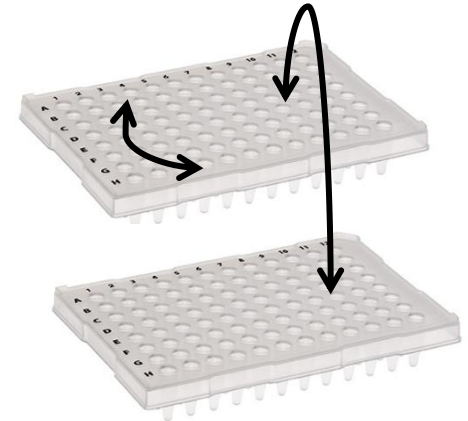
**News**  
Jan 18, 2016  
NEW EMAIL  
Due to a recent hacker attack to the OUS IT system, we currently are unable to receive or send emails to the re-research domain. Until further notice please use this email to contact us: "genomics.cf.oslo@gmail.com"  
[More]

Feb 24, 2016  
New HiSeq 4000 installed  
The HiSeq 4000 builds upon the existing HiSeq 2500 platform using the new HiSeq X patterned flow cell technology, providing unparalleled speed and performance. The dual-flow cell HiSeq 4000 System delivers the highest throughput and lowest price per sample across multiple applications. The new sequencer will provide users with faster turnaround time (run time is 3 days compared to 6-11 days in HiSeq 2500) and higher quality, more data per run and longer reads (150 bp paired-end).  
[More]

Contact Information:  
genomics.cf.oslo@gmail.com

# Data Quality Control (Prelude to Marc's lecture on Tuesday)

- **Assess within/between plate genotype reproducibility**  
⇒ SNP is deemed to have failed if <95% of samples generate a genotype at the locus
- **Exclude all SNPs with MAF <1%.**  
⇒ Low statistical power in association analysis
- **Remove all SNPs that show deviation from HWE.**  
⇒ Systematic genotyping errors, sample mix-ups, latent population substructure, or a biological effect (e.g., natural selection).
- **Screen for Mendelian inconsistencies within families.**  
⇒ Sample switches or misidentified paternity/maternity



# GENOME-WIDE ASSOCIATION STUDIES – CH.4

- Hypothesis-free (agnostic) compared to candidate-gene approach
    - Looks for association across the entire genome using high-resolution SNP arrays (0.5-2.5 mill).
  - What have we learnt?
    - Many association signals are not in genes previously thought to be associated with the disease.
    - Some associations are in areas that weren't even known before.
- ⇒ Provide new insights into biology and disease mechanism 😊

## Signals in «gene deserts»:

Prostate cancer; CL/P  
Crohn's disease

8q24  
5p13.1; 1q31.2; 10p21

## Signals in common (pleiotropy):

Diabetes/CHD/Melanoma  
Prostate/breast/colon cancers; CL/P  
Crohn's disease/Psoriasis  
Crohn's disease/T1DM

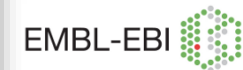
*CDKN2A/2B*  
8q24  
*IL23R*  
*PTPN2*

# Published GWAS through Dec 2012 at $p \leq 5 \times 10^{-8}$ for 17 trait categories



## 17 trait categories

- Digestive system disease
- Cardiovascular disease
- Metabolic disease
- Immune system disease
- Nervous system disease
- Liver enzyme measurement
- Lipid or lipoprotein measurement
- Inflammatory marker measurement
- Hematological measurement
- Body measurement
- Cardiovascular measurement
- Other measurement
- Response to drug
- Biological process
- Cancer
- Other disease
- Other trait

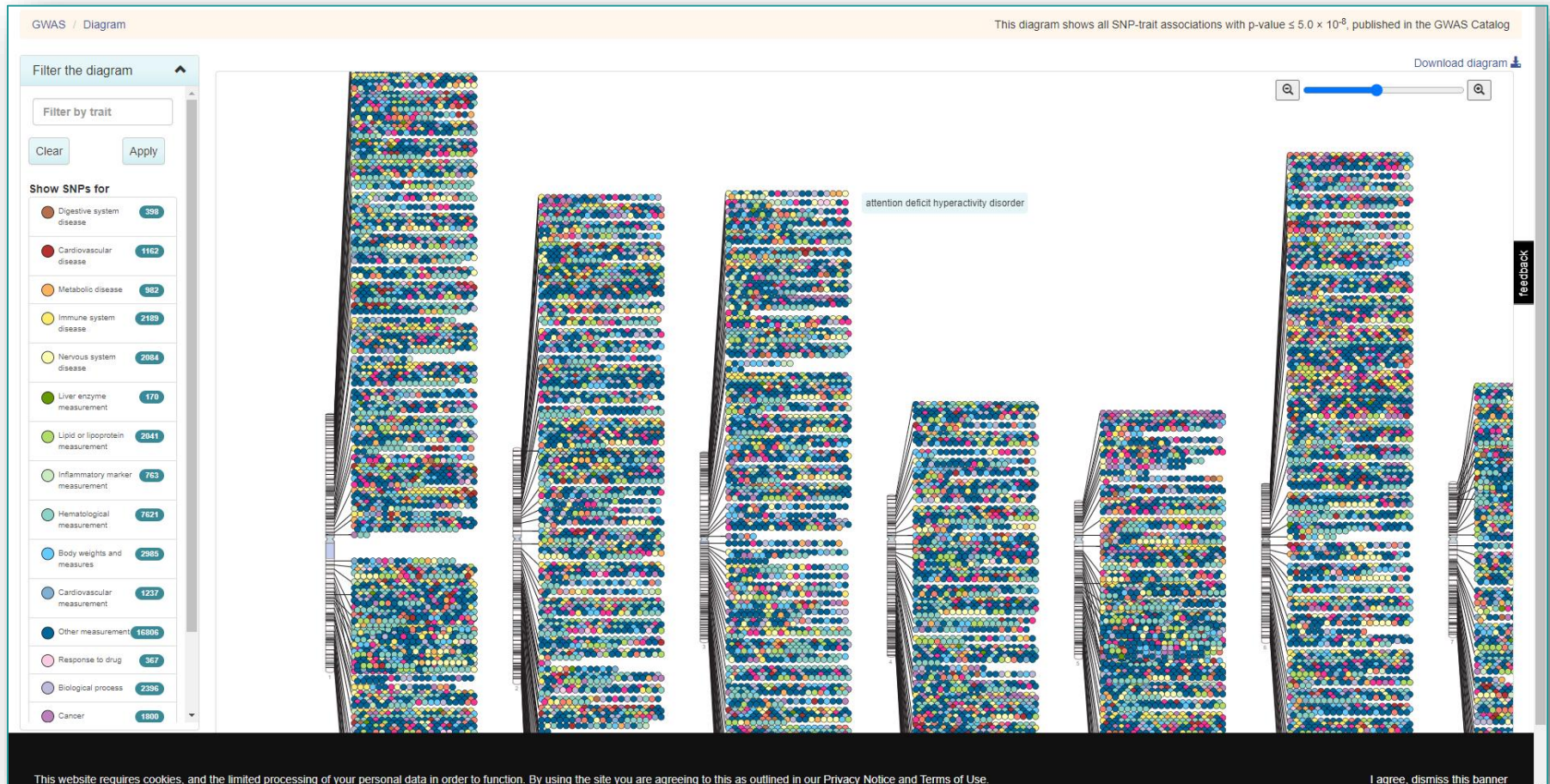


NHGRI GWA Catalog at [www.genome.gov/GWASudies](http://www.genome.gov/GWASudies)

Other useful sites that catalog GWAS (interactive): <https://www.ebi.ac.uk/gwas/>

# Interactive GWAS catalog at EBI

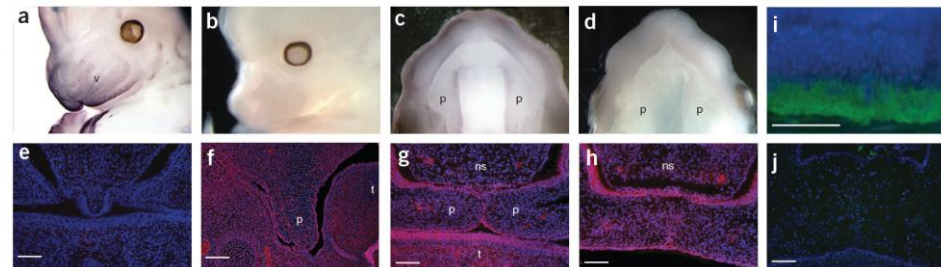
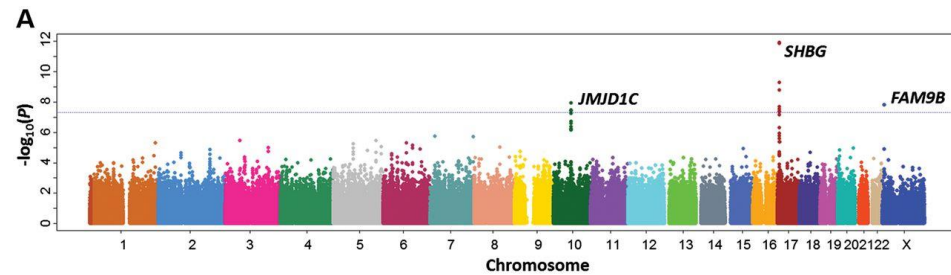
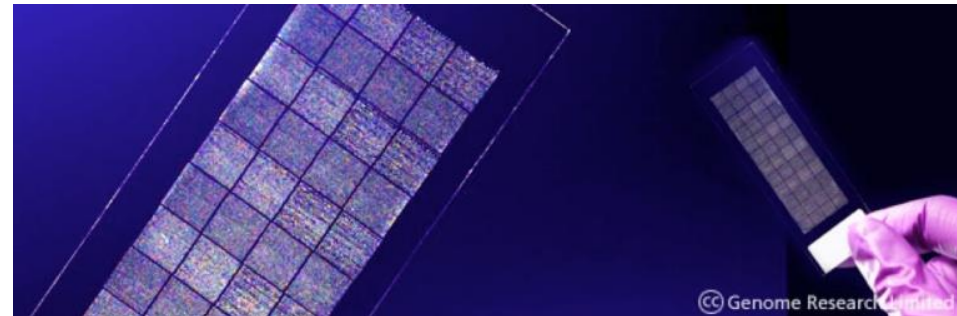
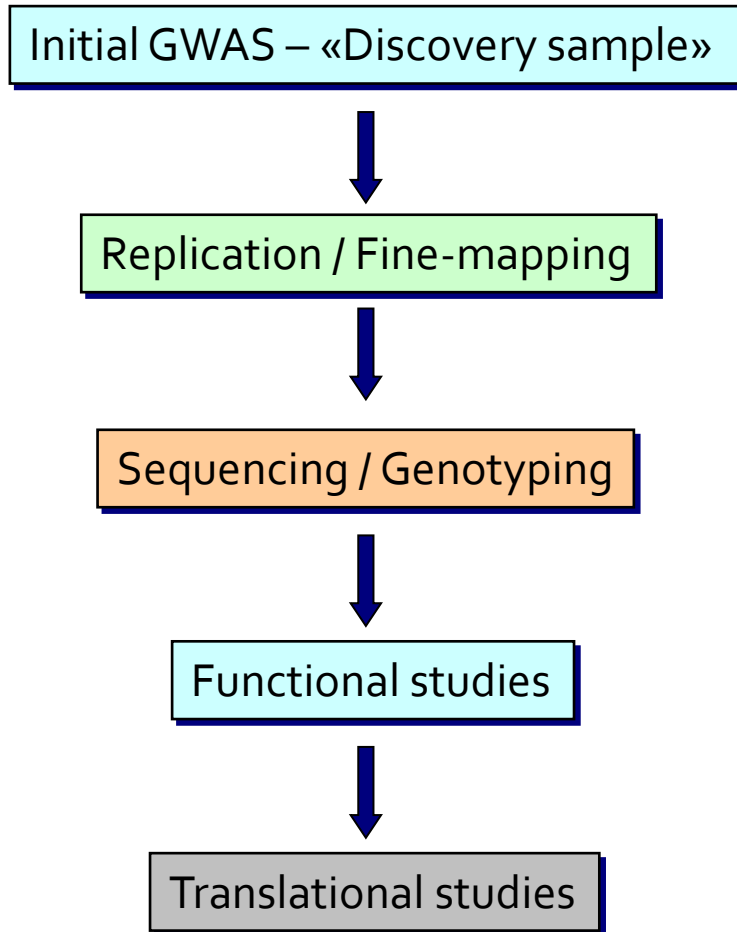
Interactive diagram shows all SNP-trait associations with genome-wide significant p-value  $\leq 5.0 \times 10^{-8}$



NHGRI GWA Catalog at [www.genome.gov/GWASudies](http://www.genome.gov/GWASudies)

Other useful sites that catalog GWAS (interactive): <https://www.ebi.ac.uk/gwas/diagram>

# TYPICAL GWAS WORKFLOW (CH. 4, P 79)



- Most of the GWAS findings so far have not led to any major clinical applications.
- HOPE -- New therapies, improved diagnostics, better prevention, better public health, & precision medicine.

# GWAS – WHAT ARE THE CRITERIA FOR SUCCESS?

- **Costs** and **availability of large samples** are major limitations
  - Useful to meta-analyze summary statistics from multiple cohorts (GWAMA)
- **Strict quality control** throughout the process (Marc Vaudel's Tuesday lecture) + Stringent significance thresholds + Importance of replication
- **Data sharing** between several research groups is an effective way of increasing power to find new genes and loci.
  - But **control for confounders** is even more important when using data from different cohorts participating in a large consortium
- **Disease heterogeneity** is a problem.
  - The more narrowly/precisely the phenotype is defined, the better the odds for identifying a causal variant (but not always!)
- Current methods are not well developed to identify **rare variants** (MAF <1%) that are perhaps associated with higher disease penetrance.



# WHAT CAN WE DO?

## Improving the resolution of current GWAS studies

- Larger sample sizes
- Endo- and sub-phenotypes
- Non-European
- Disease pleiotropy

## Clinical translation

- Prospective studies
- Aggregate risk scores

GWAS

## Exploring the full spectrum of genetic variation

- Rare variants (HapMap3, 1000G and direct sequencing)
- Structural variants (CNVs & indels)
- Epigenetic variation
- Parent-of-origin effects etc

## Understanding function

- Functional genome annotation
- eQTLs
- Model organisms

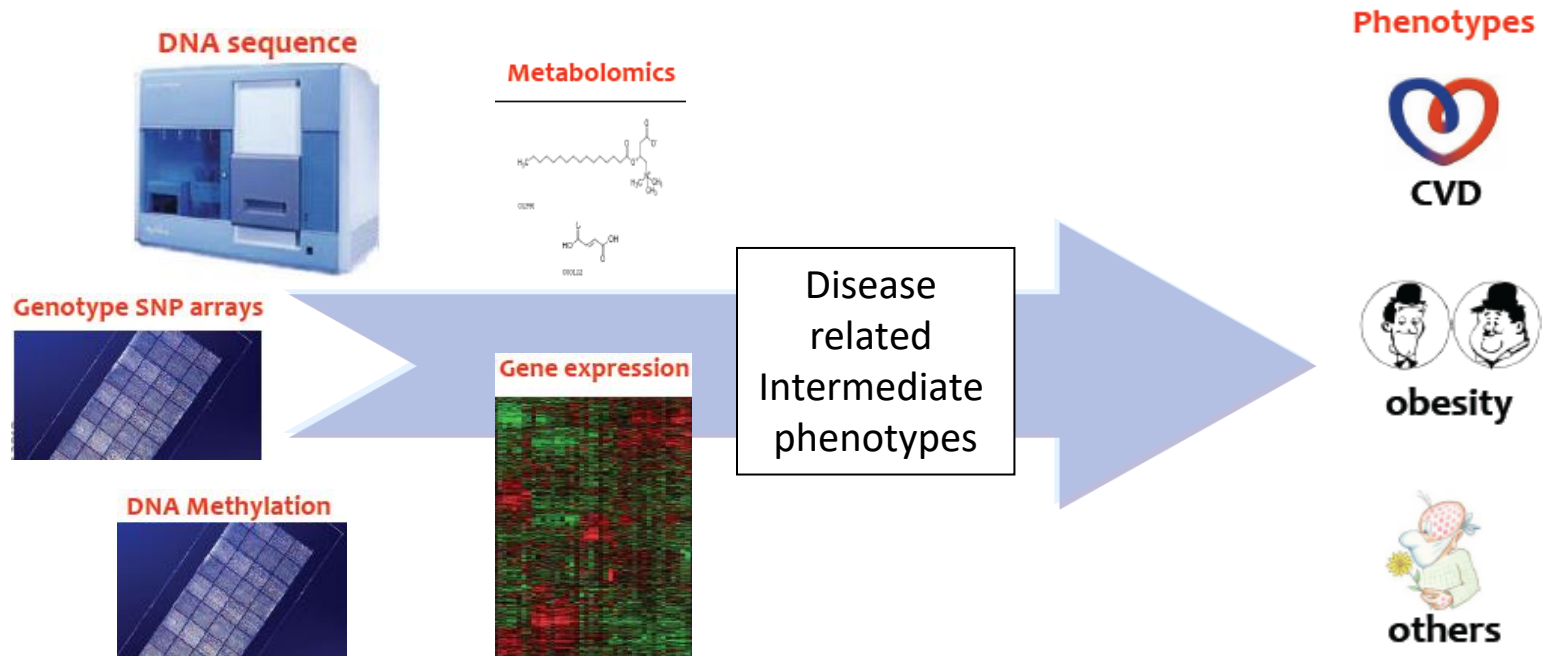
# WHOLE GENOME/EXOME SEQUENCING

## ■ Two main objectives:

- Build a comprehensive catalog of genetic variation containing both common and rare genetic variants
- Test these variants for association with disease.

## ■ Potential applications:

- Sequence based imputations in GWAS data ([Marc Vaudel's Tuesday lecture](#))
- Analyze cohorts with clearly defined phenotypes and map Mendelian diseases



# META-GWAS ANALYSES – A SHORT PRIMER

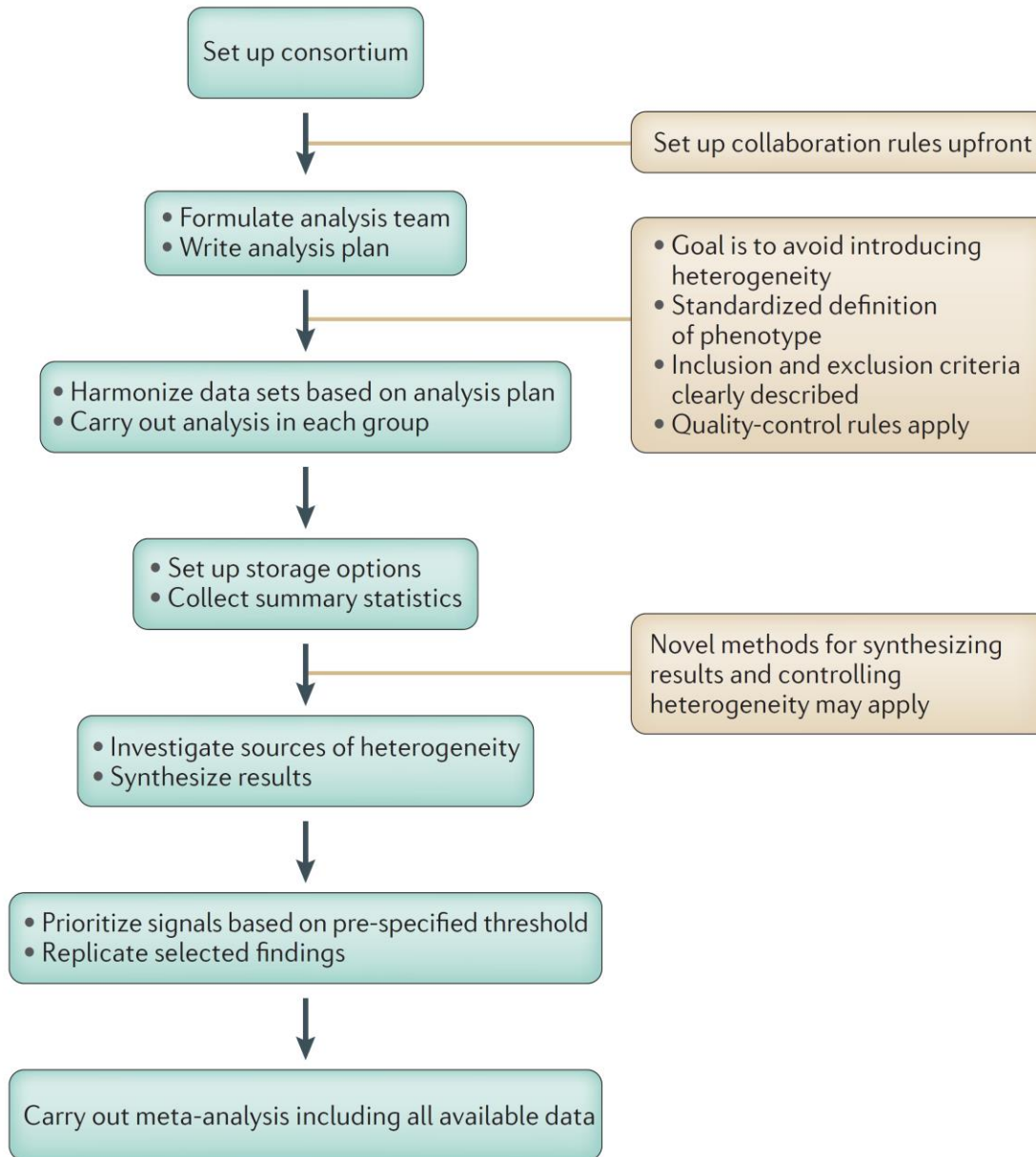
## ○ 1) DISCOVERY PHASE

- ✿ Analyze GWAS results from different cohorts (Consortia)
  - ✿ Increase statistical power through increasing sample size
  - ✿ Beware of heterogeneity (PCA, stratified analyses, inflation, QQ plot)
- ✿ Analysis of data at an aggregated level, i.e. not individual-level data.
  - ✿ Many Ethic Committees have an issue with sharing individual-level data.
  - ✿ Meta-GWAS analysis offers a good compromise
- ✿ Different cohorts perform genome-wide imputation using the same imputation panel to harmonize genotype data across cohorts
  - ✿ Harmonizes genotyping platforms (standardization)
  - ✿ Lots more SNPs to analyze ⇒ More statistical power

## 2) REPLICATION PHASE

- ✿ Invite more cohorts for replication
  - ✿ Confirmation of original findings in discovery phase

# STAGES IN A META-GWAS ANALYSIS



# EXAMPLES OF CONSORTIA

Social Science Genetic Association Consortium

Home About Us Research Events News Contact

Welcome to the Social Science Genetic Association Consortium (SSGAC).

The SSGAC is a cooperative enterprise among medical researchers and social scientists that coordinates genetic association studies for social science outcomes and provides a platform for interdisciplinary collaboration and cross-fertilization of ideas. The SSGAC also tries to promote the collection of harmonized and well-measured phenotypes.

Click here to learn about our upcoming training sessions on Social Science Genomics and Genome-Wide Data Analysis!

Current Initiatives

SSGAC in the News

The SSGAC is currently conducting large-scale

"The Genetics of Staying in School", *The Atlantic*, June 21, 2012

EUROPEAN NETWORK OF GENOMIC AND GENETIC EPIDEMIOLOGY

ENGAGE

SEVENTH FRAMEWORK PROGRAMME

Home  
Objectives  
Partners  
Work Packages  
Events  
Training  
Press & Publications  
Resources  
Contact  
FAQ



Young Investigator Profiles

**ABOUT**

ENGAGE (European Network for Genetic and Genomic Epidemiology) is a research project funded with 12 million euros by the European Commission under the 7th Framework Programme-Health Theme. The project duration is five years, starting from January 1st, 2008.

The ENGAGE Consortium has brought together 24 leading research organizations and two biotechnology and pharmaceutical companies across Europe and in Canada and Australia.

ENGAGE aims to translate the wealth of data emerging from large-scale research in genetic and genomic epidemiology from European (and other) population cohorts into information relevant to future clinical applications. The concept of ENGAGE is to enable European researchers to identify large numbers of novel susceptibility genes that influence metabolic, behavioural and cardiovascular traits, and to study the interactions between genes and life style factors.

The ENGAGE consortium will integrate and analyse one of the largest ever human genetics dataset (more than 80,000 genome-wide association scans and DNAs and serum/plasma samples from over 600,000 individuals).

One goal is to demonstrate that the findings from ENGAGE can be used as diagnostic indicators for common diseases that will help us to understand better risk factors, disease progression and why people differ in responses to treatment.


**NEWS**

ENGAGE Flagship Paper: 'The Role of Adiposity in Cardiometabolic Traits: A Mendelian Randomization Analysis' (Fall T et al, Pedersen NL, McCarthy MI, Ingelsson E, Prokopenko I for ENGAGE, 25 June 2013).

ENGAGE Paper: 'Data sharing in large research consortia: experiences and recommendations from ENGAGE' (Budín-Ljøsne I et al, June 2013)

ENGAGE ESHG Satellite Meeting 'Beyond GWAS: Biological and Clinical Insights from Research in European Biobanks', June 10th, Paris

ENGAGE Paper: 'GWAS of 126,559



COHORTS FOR HEART AND AGING RESEARCH IN GENOMIC EPIDEMIOLOGY

**CHARGE Consortium**


*The Cohorts for Heart and Aging Research in Genomic Epidemiology (CHARGE) Consortium was formed by phenotyped longitudinal cohort studies.*

Its founding member cohorts include:

- [Age, Gene, Environment, Susceptibility Study -- Reykjavik](#)
- [Atherosclerosis Risk in Communities Study](#)
- [Cardiovascular Health Study](#)
- [Framingham Heart Study](#)
- [Rotterdam Study](#)

Additional core cohorts include:

- [Coronary Artery Risk Development in Young Adults](#)
- [Family Heart Study](#)
- [Health, Aging, and Body Composition Study](#)
- [Jackson Heart Study](#)
- [Multi-Ethnic Study of Atherosclerosis](#)



**EAGLE Consortium**


The EARly Genetics and Lifecourse Epidemiology (EAGLE) Consortium is a consortium of pregnancy and birth cohorts that aims to collaborate to investigate the genetic basis of phenotypes in antenatal and early life and childhood.

EAGLE covers a broad range of pathways and phenotypes, and will integrate closely with the DOHaD (developmental origins of health and disease) community.

All participating cohorts (1958 British Birth Cohort; ALSPAC; CHOP; COPSAC; DBC; Exeter Family Study; Generation R; HBCS; LISA+; MoBa; NTR; NFBC 66; Project Viva; Raine) have GWAS data available by July 1st 2009.

EAGLE working groups and leaders are listed below:

- **Antenatal Growth** (Vincent Jaddoe and Craig Pennell)



**EAGLE**  
EARly Genetics & Lifecourse Epidemiology Consortium

# Welcome to the Social Science Genetic Association Consortium (SSGAC).

The SSGAC is a cooperative enterprise among medical researchers and social scientists that coordinates genetic association studies for social science outcomes and provides a platform for interdisciplinary collaboration and cross-fertilization of ideas. The SSGAC also tries to promote the collection of harmonized and well-measured phenotypes.



Recent Events: [Russell Sage Foundation Summer Institute in Social-Science Genomics, 2021](#)  
[Click here to learn about the 2018 Polygenic Prediction and its Application in Social Science Conference](#)

## Current Initiatives



The SSGAC is currently conducting:

- Within-family GWAS of multiple phenotypes
- Ongoing updates of the Polygenic Index Repository
- Developing methods for multi-ancestry genetic analysis

Please contact us if you are interested in joining these initiatives.

MORE

## Data



To locate and download summary data from past studies of the SSGAC, click the link below.

MORE

## SSGAC in the News



["Many Genes Play a Role in Educational Attainment, Enormous Genetic Study Finds"](#)  
*The New York Times*, July 23, 2018

["Why Study the Genetics of Staying in School?"](#)  
*The Atlantic*, July 23, 2018.

MORE

Check us out on Twitter: [@thessgac](#)

### Team Co-Founders

- [Daniel Benjamin](#)  
University of California, Los Angeles (UCLA)
- [David Cesarini](#)  
New York University
- [Philippe Koellinger](#)  
Vrije Universiteit Amsterdam

### Steering Committee

- [Daniel Benjamin](#)  
University of California, Los Angeles (UCLA)
- [David Cesarini](#)  
New York University
- [Philippe Koellinger](#)  
Vrije Universiteit Amsterdam
- [Patrick Turley](#)  
University of Southern California (USC)

# STANDARD OPERATION PROTOCOL (SOP)

---

## 1) STANDARD OPERATING PROTOCOL (SOP) in «Discovery Phase»

- ✿ Background of the proposed Meta-GWAS analysis (GWAMA)
  - ✿ Goals of the initiative
- ✿ Trait definition and instructions for phenotype harmonization
  - ✿ A detailed definition of the trait (not all cohorts have same measures)
  - ✿ Eligibility and sample inclusion/exclusion criteria
- ✿ Genotypes and imputation
  - ✿ Imputation with chosen panel (HapMap Phase II CEU Panel, 1000 Genomes, HRC)
  - ✿ Filters to be applied before imputation (SNP call >95%, HWE  $p > 10e-6$ , MAF >5%)
- ✿ Analysis details
  - ✿ Specification of models to be used in the analysis
  - ✿ Linear regression/Logistic regression, Include PCA for correcting for stratification
- ✿ Results file formats
  - ✿ Format to report GWAS results from individual cohorts

# REPORTING OF RESULTS

Variable name (case sensitive!!)	Description
SNPID	SNP ID as rs number
Chr	Chromosome number (1-22).
position	physical position for the reference sequence (indicate build 35/36 in readme file)
coded_all	Coded allele, also called modelled allele (in example of A/G SNP in which AA=0, AG=1 and GG=2, the coded allele is G)
noncoded_all	The other allele
strand_genome	+ or -, representing either the positive/forward strand or the negative/reverse strand of the human genome reference sequence; to clarify which strand the coded_all and noncoded_all are on
Beta	Beta estimate from genotype-phenotype association, at least 5 decimal places – ‘NA’ if not available
SE	Standard error of beta estimate, to at least 5 decimal places – ‘NA’ if not available
Pval	<i>p</i> -value of test statistic, here just as a double check – ‘NA’ if not available
AF_coded_all	Allele frequency for the coded allele – ‘NA’ if not available
HWE_pval	Exact test Hardy-Weinberg equilibrium <i>p</i> -value -- only directly typed SNPs, NA for imputed
callrate	Genotyping call rate after exclusions
n_total	Total sample with phenotype and genotype for SNP
imputed	1/0 coding; 1=imputed SNP, 0=if directly typed
used_for_imp	1/0 coding; 1=used for imputation, 0=not used for imputation
oevar_imp*	Observed divided by expected variance for imputed allele dosage -- NA otherwise
avpostprob**	Average posterior probability for imputed SNP allele dosage (applies to best-guess genotype imputation)

\* oevar\_imp is called  $r^2$  in Mach, proper\_info in Impute and  $R^2$  in Beagle.

\*\* avpostprob is called Quality in Mach, certainty in Impute and Beagle does not give this statistic.



## PDB 289 (Study of prematurity – PI Bo Jacobsson)

	A	B	C	D	E	F	G	H	I	J	K	L	
1	<b>SAMPLE INFORMATION</b>												
2													
3	Country					Norway							
4	Sampling scheme					Population-based, nested case-control	e.g., family-based, clinically-selected (please specify selected phenotype), etc.						
5	SNP chip					Illumina 660W quad							
6	Pre-imputation QC												
7		Marker filters:											
8		MAF >			0.5	5% recommended							
9		Call rate >			95%	95% recommended							
10		HWE exact test at p >			0.001	10E-06 recommended							
11		Removed subjects with:											
12		Overall call rates <			98%								
13	Imputation & association procedure												
14		Imputation software					PLINK 1.07	please specify version number					
15		Reference sample					HAPMAP II CEU	HapMap phase II CEU recommended					
16		NCBI build					NCBI 36.2	e.g. NCBI 36.2					
17		Association software					PLINK 1.07	please specify version number					
18	Study contacts (name, email):												
19		Data analyst:	Ronny Myhre, Astanand Jugessur and Håkon Gjessing										
20		Primary contact:	Bo Jacobsson (PI; bo.jacobsson@obgyn.gu.se) and Astanand Jugessur (astanand.jugessur@fhi.no)										
21		Other contact(s):	Per Magnus										
22	Additional notes					x	e.g., non-standard covariates included in analyses (please specify)						
23													
24	<b>SAMPLE DEMOGRAPHICS</b>												
25													
26	N					Females	1338	Controls	678	sPTD cases	660		
27	Age at reporting												
28		Mean					28.7	28.9	28.4				
29		St. Dev.					3.5	3.6	3.6				
30		Range					14 (20-34)	14 (20-34)	14 (20-34)				
31	Birth Year												
32		Mean					1974.7	1974.4	1975				
33		St. Dev.					3.8	3.4	4.0				
34		Range					20 (1966-1986)	19 (1966-1985)	20 (1966-1986)				
35	Race (N per category)												
36		American Indian or Alaska Native					na	Note that only individuals of European heritage should be included in the analysis					
37		Asian					na						
38		Native Hawaiian or Other Pacific Islander					na						
39		Black					na						
40		White					na						
41	Ethnicity (N per category)												
42		Hispanic or Latino					na						
43		Not Hispanic or Latino					na						

# EXAMPLE OF A GWAMA

## 1) Trait proposed for a GWAMA: «Aggressive behavior»

- ✿ SOP describes the goal of the proposed GWAMA
  - ✿ **Goal:** large-scale meta-GWAS on Aggressive behavior
  - ✿ **Merit:** Findings will help identify to what extent the effect of the SNP(s) changes with age, instrument, or the rater of the behavior.
- ✿ Trait definition and instructions for phenotype harmonization
  - ✿ Phenotype data at different ages (3 to 18 yrs) and as rated by different raters (parental, self and/or teacher ratings) to be included in a single analysis
  - ✿ **Instruments:** A variety of psychometric instruments (e.g. CBCL, SDQ, ASR, YSR)
  - ✿ Sample size threshold for inclusion: at least 1000 subjects.
  - ✿ Limit analyses to subjects of European ancestry.
- ✿ Genotypes and imputation
  - ✿ Imputation with chosen panel (1000 Genomes)
  - ✿ **Software for imputation:** IMPUTE, MACH, MINIMAC or BEAGLE.
  - ✿ Filters to be applied before imputation (SNP call >95%, HWE  $p > 10e-6$ , MAF >5%)
- ✿ Analysis
  - ✿ For cohorts providing a single phenotype measure: Run the GWA using linear Reg.
  - ✿ **Covariates:** sex, Z-score of age at time of assessment, Age<sup>2</sup> (Z-transformed, then squared), the first 5 PCs, Study-specific covariates (study site, batch effects etc.)

# EXAMPLE OF A META-GWAS – CONTD...

## 1) Instructions for genotype handling (pre-imputation QC):

- ✿ Exclude SNPs with:
  - ✿ MAF <1%
  - ✿ SNP call rate <95%
  - ✿ Failure of HWE exact test at  $p < 1e-6$
  - ✿ Poor clustering on visual inspection of intensity plots.
  - ✿ Wrong sex, aberrant genotype (XXY), known 1st or 2nd degree relatives in sample

## 2) Imputation:

- ✿ Use 1000 genomes Phase I release and coordinates as used in GRCh37
- ✿ Imputation software: IMPUTE or MACH
- ✿ Use servers for imputation: Michigan imputation server or Sanger Institute in UK
- ✿ Provide per-SNP quality indicators (proper\_info in IMPUTE,  $r^2$ .hat in MACH)

## 3) Analysis:

- ✿ Perform association test using MACH2QTL or SNPTEST

# EXAMPLE OF A META-GWAS – CONTD...

## Uploading data

- ✿ To a secure server using secure transfer protocol (sftp)
  - ✿ Download and Install an sftp software; e.g. Filezilla or WinScp
  - ✿ Upload a «README.txt» file with a brief description of data uploaded, the date, the human genome reference sequence used for strand reference, and scale of Beta estimates.
  - ✿ Prepare a file named «STUDY.PHEN.DATE.txt»
    - ✿ Study=Cohort, PHEN=phenotype, Date=DDMMYYYY (date file was prepared)

## Meta-analysis:

- ✿ Usually done by the lead analysts from the cohort(s) initiating this GWAMA
- ✿ Software: METAL or GWAMA

# A FEW EXAMPLES...

## Papers on EA based including the MoBa dataset (PDB 289)

2013 (N=101,069 individuals; 3 GW sig SNPs; 2% variance)

**ScienceExpress** Reports

### GWAS of 126,559 Individuals Identifies Genetic Variants Associated with Educational Attainment

All authors with their affiliations appear at the end of this paper.

A genome-wide association study of educational attainment was conducted in a discovery sample of 101,069 individuals and a replication sample of 25,490. Three independent SNPs are genome-wide significant ( $r^2=0.20913$ ,  $r^2=11654700$ ,  $r^2=4681268$ ), and all three replicate. Estimated effect sizes are small ( $R^2 = 0.02\%$ ), approximately 1 month of schooling per allele. A linear polygenic score from all measured SNPs accounts for  $\approx 2\%$  of the variance in both educational attainment and cognitive function. Genes in the region of the loci have previously been associated with health, cognitive, and central nervous system phenotypes, and bioinformatics analyses suggest the involvement of the anterior caudate nucleus. These findings provide promising candidate SNPs for follow-up work, and our effect size estimates can anchor power analyses in social-science genetics.

used at an age at which very likely to have education (over 95% of at least 30, (5)). On a have 13.3 years of 23.1% have a college of pooling of GWAS re-conducted analyses w/ to the HapMap 2 CEU set. To guard ag stratification, the first components of the gen included as controls in level analyses. All GWAS results were q cross-checked, and in single genomic co- ple-size weighting sc independent analysis ca. At the cohort level evidence of general values. As in previous comes, tests, (1,3), the

2018 (N=1.1 mill individuals; 1271 GW sig SNPs; 11-13% variance)

naturegenetics

Explore content About the journal Publish with us Subscribe

nature > naturegenetics > articles > article

Article | Published: 23 July 2018

### Gene discovery and polygenic prediction from a genome-wide association study of educational attainment in 1.1 million individuals

James J. Lee<sup>1</sup>, Robbee Wedow<sup>2</sup>, David Cesarini<sup>3</sup>, + show authors

Nature Genetics 50, 1112–1121 (2018) | Cite this article

64k Accesses | 727 Citations | 2042 Altmetric | Metrics

**Abstract**

Here we conducted a large-scale genetic association analysis of educational attainment in a sample of approximately 1.1 million individuals and identify 1,271 independent genome-wide significant SNPs. For the SNPs taken together, we found evidence of heterogeneous effects across environments. The SNPs implicate genes involved in brain-development processes and neuron-to-neuron communication. In a separate analysis of the X chromosome, we identify 10 independent genome-wide significant SNPs and estimate a SNP heritability of around 0.3% in both men and women, consistent with partial dosage compensation. A joint (multi-phenotype) analysis of educational attainment and three related cognitive phenotypes generates polygenic scores that explain 11–13% of the variance in educational attainment and 7–10% of the variance in cognitive performance. This prediction accuracy substantially increases the utility of polygenic scores as tools in research.

2015 (N=293,723 individuals; 74 GW sig SNPs)

LETTER

doi:10.1038/nature17671

### Genome-wide association study identifies 74 loci associated with educational attainment

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

Educational attainment is strongly influenced by social and other environmental factors, but genetic factors are estimated to account for at least 20% of the variation across individuals<sup>1</sup>. Here we report the results of a genome-wide association study (GWAS) for educational attainment that extends our earlier discovery sample<sup>1,2</sup> of 101,069 individuals to 293,723 individuals, and a replication study in an independent sample of 111,349 individuals from the UK Biobank. We identify 74 genome-wide significant loci associated with the number of years of schooling completed. Single-nucleotide polymorphisms associated with educational attainment are disproportionately found in genomic regions regulating gene expression in the fetal brain. Candidate genes are preferentially

Our meta-analysis identified 74 approximately independent genome-wide significant loci. For each locus, we define the 'lead SNP' in the SNP in the genomic region that has the smallest *P* value (Supplementary Information section 1.6.1). Figure 1 shows a Manhattan plot with the lead SNPs highlighted. This includes the three SNPs that reached genome-wide significance in the discovery stage of our previous GWAS meta-analysis of educational attainment<sup>1</sup>. The quantile-quantile (Q-Q) plot of the meta-analysis (Extended Data Fig. 1) exhibits inflation ( $\lambda_{GC} = 1.28$ ), as expected under polygenicity<sup>3</sup>. Extended Data Fig. 2 shows the estimated effect sizes of the lead SNPs. The estimates range from 0.014 to 0.048 standard deviations per allele (2.7 to 9.0 weeks of schooling), with incremental *R*<sup>2</sup> in the

2022 (N= ~3 mill individuals; 3952 GW sig SNPs; 12-16% of variance)

naturegenetics

ARTICLES

https://doi.org/10.1038/s41588-022-01016-z

Check for updates

### OPEN Polygenic prediction of educational attainment within and between families from genome-wide association analyses in 3 million individuals

Aysu Okbay<sup>1,2,3,9,10</sup>, Yeda Wu<sup>2</sup>, Nancy Wang<sup>3</sup>, Hariharan Jayashankar<sup>1</sup>, Michael Bennett<sup>2,3</sup>, Seyed Moeen Nehzati<sup>1</sup>, Julia Sidorenko<sup>2,3</sup>, Hyeokmoon Kweon<sup>1</sup>, Grant Goldman<sup>1</sup>, Tamara Gjorgjieva<sup>2,3</sup>, Yunxuan Jiang<sup>1</sup>, Barry Hicks<sup>1</sup>, Chao Tian<sup>1</sup>, David A. Hinds<sup>2,3</sup>, Rafael Ahlsgoer<sup>1</sup>, Patrik K. E. Magnusson<sup>2,3</sup>, Sven Oskarsson<sup>2,3</sup>, Caroline Hayward<sup>2,3</sup>, Archie Campbell<sup>2,3,10</sup>, David J. Porteous<sup>2,3,10,11</sup>, Jeremy Freese<sup>12</sup>, Pamela Herd<sup>13</sup>, Z3andMe Research Team<sup>14</sup>, Social Science Genetic Association Consortium<sup>15</sup>, Chelsea Watson<sup>16</sup>, Jonathan Jala<sup>17</sup>, Dalton Conley<sup>18</sup>, Philipp D. Koellinger<sup>15</sup>, Magnus Johannesson<sup>19</sup>, David Laibson<sup>20</sup>, Michelle N. Meyer<sup>19</sup>, James J. Lee<sup>19</sup>, Augustine Kong<sup>20</sup>, Loic Yengo<sup>2,19</sup>, David Cesarini<sup>21,22,23,24</sup>, Patrick Turley<sup>23,24,25</sup>, Peter M. Visscher<sup>2,2,9,26</sup>, Jonathan P. Beauchamp<sup>25,26</sup>, Daniel J. Benjamin<sup>2,3,4,26,26,26</sup> and Alexander I. Young<sup>4,26,26,26,26</sup>

We conduct a genome-wide association study (GWAS) of educational attainment (EA) in a sample of ~3 million individuals and identify 3,952 approximately uncorrelated genome-wide significant single-nucleotide polymorphisms (SNPs). A genome-wide polygenic predictor, or polygenic index (PGI), explains 12–16% of EA variance and contributes to risk prediction for ten diseases. Direct effects (i.e., controlling for parental PGIs) explain roughly half the PGI's magnitude of association with EA and other phenotypes. The correlation between mate-pair PGIs is far too large to be consistent with phenotypic assortment alone, implying additional assortment on PGI-associated factors. In an additional GWAS of dominance deviations from the additive model, we identify no genome-wide significant SNPs, and a separate X-chromosome additive GWAS identifies 57.

2015

Molecular Psychiatry (2015) 20, 735–743  
© 2015 Macmillan Publishers Limited. All rights reserved 1359-4184/15

ORIGINAL ARTICLE

### The association between lower educational attainment and depression owing to shared genetic effects? Results in ~25 000 subjects

WJ Peyrot<sup>1</sup>, SH Lee<sup>2</sup>, Y Milaneschi<sup>1</sup>, A Abdellaoui<sup>3</sup>, EM Byrne<sup>4</sup>, T Esko<sup>5,6</sup>, EJC de Geus<sup>7</sup>, G Heman<sup>8,9,10</sup>, JJ Hottenga<sup>7</sup>, S Kloiber<sup>7</sup>, DJ Levinson<sup>11</sup>, S Lucae<sup>12</sup>, Major Depressive Disorder Working Group of the Psychiatric GWAS Consortium (Corporate Collaborator), NG Martin<sup>9</sup>, SE Medland<sup>9</sup>, A Metspalu<sup>13,14</sup>, M Milani<sup>15</sup>, MM Nothen<sup>16</sup>, JB Potkin<sup>17</sup>, M Rietschel<sup>12</sup>, CA Rietveld<sup>13,14</sup>, S Ripke<sup>15</sup>, J Shi<sup>16</sup>, Social Science Genetic Association Consortium (Corporate Collaborator), G Willemsen<sup>18</sup>, Z Zhu<sup>19</sup>, DI Boomsma<sup>7</sup>, NR Wray<sup>20</sup> and BWJ Penninx<sup>1</sup>

2016

CrossMark

### Genetic variants linked to education predict longevity

Ricardo E. Marioni<sup>1,2,3,4,5,6,7,8,9,10,11,12</sup>, Stuart J. Ritchie<sup>13,14,15</sup>, Peter K. Joshi<sup>16</sup>, Saskia P. Hagenaars<sup>17,18</sup>, Aysu Okbay<sup>19,20</sup>, Krista Fischer<sup>21</sup>, Mark J. Adams<sup>22</sup>, W. David Hill<sup>23</sup>, Gail Davies<sup>24</sup>, Social Science Genetic Association Consortium<sup>25</sup>, Reka Nagy<sup>26</sup>, Carmen Amador<sup>27</sup>, Kristi Lilja<sup>28</sup>, Andres Metspalu<sup>29</sup>, David C. Liewald<sup>30</sup>, Archie Campbell<sup>31</sup>, James F. Wilson<sup>32</sup>, Caroline Hayward<sup>33</sup>, Tatu Esko<sup>34,35</sup>, David J. Porteous<sup>36</sup>, Catharine R. Gale<sup>37,38</sup>, and Ian J. Deary<sup>39,40</sup>

<sup>1</sup>Centre for Cognitive Ageing and Cognitive Epidemiology, University of Edinburgh, Edinburgh EH8 9JZ, United Kingdom; <sup>2</sup>Medical Genetics Section, Centre for Genetic and Experimental Medicine, Institute of Genetics and Molecular Medicine, University of Edinburgh, Edinburgh EH4 2XU, United Kingdom; <sup>3</sup>Institute for Molecular Bioscience, University of Queensland, Brisbane, QLD 4072, Australia; <sup>4</sup>Department of Psychology, University of Edinburgh, Edinburgh EH8 9JZ, United Kingdom; <sup>5</sup>Usher Institute for Population Health Sciences and Informatics, University of Edinburgh, Edinburgh EH8 40X, United Kingdom; <sup>6</sup>Division of Psychiatry, University of Edinburgh, Edinburgh EH8 9YL, United Kingdom; <sup>7</sup>Department of Applied Economics, Erasmus School of Economics, Erasmus University, 3002 PA Rotterdam, The Netherlands; <sup>8</sup>Department of Epidemiology, Erasmus Medical Center, 3015 CE Rotterdam, The Netherlands; <sup>9</sup>Erasmus University Rotterdam Institute for Behavior and Biology, Rotterdam 3062 PA, The Netherlands; <sup>10</sup>Estonian Genome Centre, University of Tartu, Tartu, 51010, Estonia; <sup>11</sup>Medical Research Council Human Genetics Unit, Institute of Genetics and Molecular Medicine, University of Edinburgh, Edinburgh EH4 2XU, United Kingdom; <sup>12</sup>Institute of Mathematical Statistics, University of Tartu, Tartu, 50009, Estonia; <sup>13</sup>Broad Institute, Cambridge, MA 02142; <sup>14</sup>Department of Endocrinology, Children's Hospital Boston, Boston, MA 02115; and <sup>15</sup>Medical Research Council Lifecourse Epidemiology Unit, University of Southampton, Southampton SO17 1BJ, United Kingdom

<http://uis.unesco.org/en/topic/international-standard-classification-education-isced>

**unesco**  
Institute for Statistics

EXPLORE THEMES > EDUCATION & LITERACY >

# International Standard Classification of Education (ISCED)

MY PINS FRANÇAIS

The world's education systems vary widely in terms of structure and curricular content. Consequently, it can be difficult to compare national education systems with those of other countries or to benchmark progress towards national and international goals.

The International Standard Classification of Education (ISCED 2011) provides a comprehensive framework for organising education programmes and qualification by applying uniform and internationally agreed definitions to facilitate comparisons of education systems across countries. ISCED is a widely-use...

[READ MORE](#)

## In Focus

**QUESTIONNAIRE**

**ISCED 2011**

International Standard Classification of Education (ISCED) 2011

[DOWNLOAD \(1.14 MB\)](#) 28/03/2018 PDF

## Latest News

- **UIS data release features new SDG 4 indicators and disaggregated dimensions**  
With just a decade remaining to achieve the Sustainable Development Goals for Education (SDG 4), relevant, timely, and quality data is key to monitoring progress and informing policy responses. This...  
22/09/2021
- **WHAT'S NEXT? Lessons on Education Recovery**  
Since the outbreak of the COVID-19 pandemic in 2020, most governments worldwide have implemented policies to contain the disease's spread. While incurring high economic costs, restrictive procedures...  
08/07/2021

**unesco**  
INSTITUTE  
for  
STATISTICS  
United Nations  
Educational, Scientific and  
Cultural Organisation

# International Standard Classification of Education ISCED 2011



OPEN

# Polygenic prediction of educational attainment within and between families from genome-wide association analyses in 3 million individuals

Aysu Okbay<sup>1,197,198</sup>  , Yeda Wu<sup>2</sup>, Nancy Wang<sup>3</sup>, Hariharan Jayashankar<sup>3</sup>, Michael Bennett <sup>3</sup>, Seyed Moeen Nehzati<sup>4</sup>, Julia Sidorenko <sup>2</sup>, Hyeokmoon Kweon<sup>1</sup>, Grant Goldman<sup>3</sup>, Tamara Gjorgjieva <sup>3</sup>, Yunxuan Jiang<sup>5</sup>, Barry Hicks<sup>5</sup>, Chao Tian<sup>5</sup>, David A. Hinds <sup>5</sup>, Rafael Ahlskog<sup>6</sup>, Patrik K. E. Magnusson <sup>7</sup>, Sven Oskarsson <sup>6</sup>, Caroline Hayward <sup>8</sup>, Archie Campbell <sup>9,10</sup>, David J. Porteous <sup>9,10,11</sup>, Jeremy Freese<sup>12</sup>, Pamela Herd<sup>13</sup>, 23andMe Research Team\*, Social Science Genetic Association Consortium\*, Chelsea Watson<sup>4</sup>, Jonathan Jala<sup>4</sup>, Dalton Conley<sup>14</sup>, Philipp D. Koellinger<sup>1,15</sup>, Magnus Johannesson<sup>16</sup>, David Laibson<sup>17</sup>, Michelle N. Meyer<sup>18</sup>, James J. Lee<sup>19</sup>, Augustine Kong<sup>20</sup>, Loic Yengo<sup>2,198</sup>, David Cesarini<sup>3,21,22,198</sup>, Patrick Turley<sup>23,24,198</sup>, Peter M. Visscher<sup>2,198</sup> , Jonathan P. Beauchamp<sup>25,198</sup>, Daniel J. Benjamin <sup>3,4,26,198</sup>  and Alexander I. Young<sup>4,26,197,198</sup> 

**We conduct a genome-wide association study (GWAS) of educational attainment (EA) in a sample of ~3 million individuals and identify 3,952 approximately uncorrelated genome-wide-significant single-nucleotide polymorphisms (SNPs). A genome-wide polygenic predictor, or polygenic index (PGI), explains 12-16% of EA variance and contributes to risk prediction for ten diseases. Direct effects (i.e., controlling for parental PGIs) explain roughly half the PGI's magnitude of association with EA and other phenotypes. The correlation between mate-pair PGIs is far too large to be consistent with phenotypic assortment alone, implying additional assortment on PGI-associated factors. In an additional GWAS of dominance deviations from the additive model, we identify no genome-wide-significant SNPs, and a separate X-chromosome additive GWAS identifies 57.**

## Main findings of 2022 paper

- $N \approx 3$  mill individuals
- 3952 GW significant SNPs identified
- GW polygenic predictor (PGI) explains 12-16% of EA variance
- The PGI contributes to risk prediction for 10 diseases



# LECTURE OUTLINE

## General introduction to genetic epidemiology (lecture I)

- Part I
  - What's a complex trait?
  - Genetic basis of complex traits
- Part II
  - Genetic approaches to studying complex traits
  - Candidate-gene analysis, GWAS, and GWAMA

