

Introduction to Genome-wide Association Studies (GWAS)

Teresa Ferreira and Nilufer Rahmioglu

Wellcome Centre for Human Genetics/Big Data Institute

GWAS Course Overview

- Session 1 (9:30-10:30) Historical Perspective and Concepts
- Session 2 (10:30-11:30) Steps involved in Genome-wide Association Studies

Break: 11:30-12:00

- Session 3 (12:00-13:00): Post-GWAS Analyses: Characterising GWAS Loci

Lunch-break: 13:00-14:00

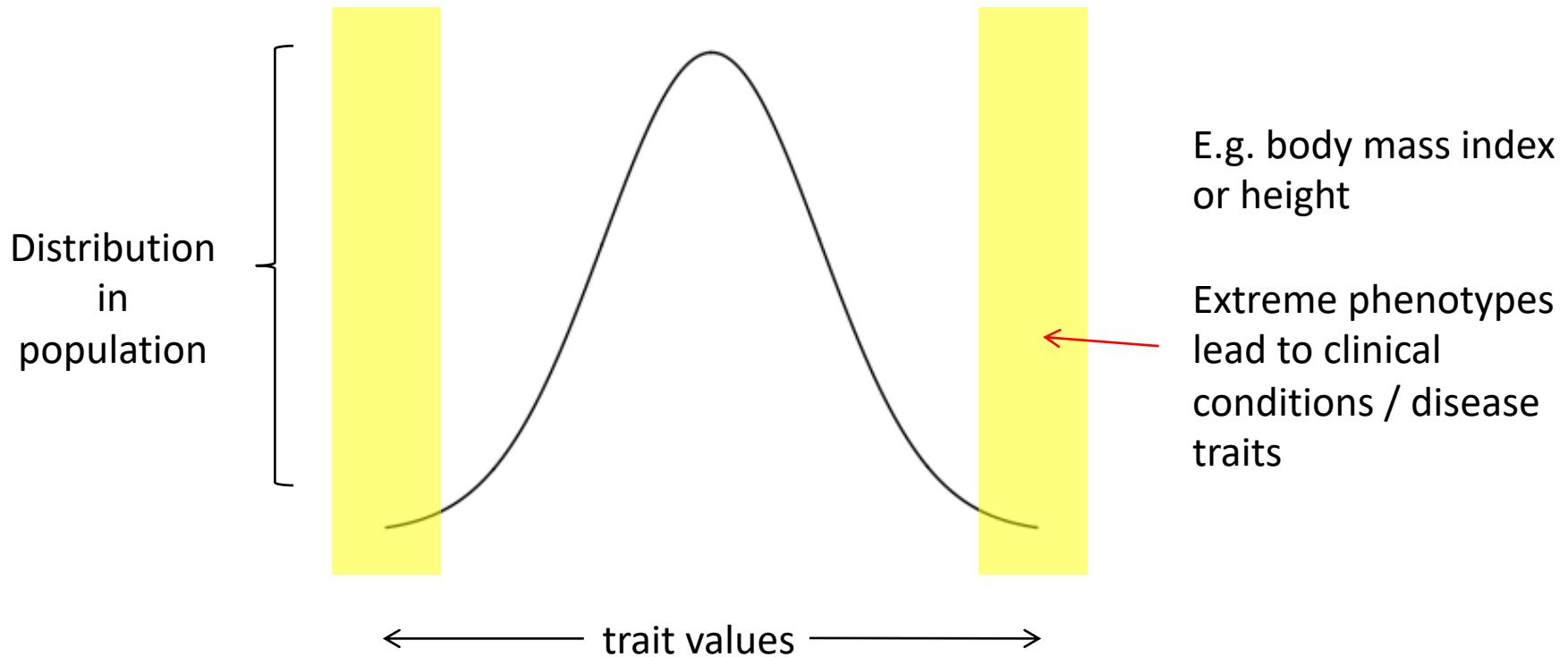
- Session 4 (14:00-16:30): Practical

Session 1: Historical Perspective and Concepts

Learning Objectives

- Understand the definition of a complex trait
- The concepts behind identification of disease variants
- Historical perspective on studying the human genome
- How to study the underlying genetic variants for Monogenic Disease? Linkage studies
- How to study the underlying genetic variants for Complex Disease? Association studies
- GWAS: For identification of common variants for complex traits

A complex trait

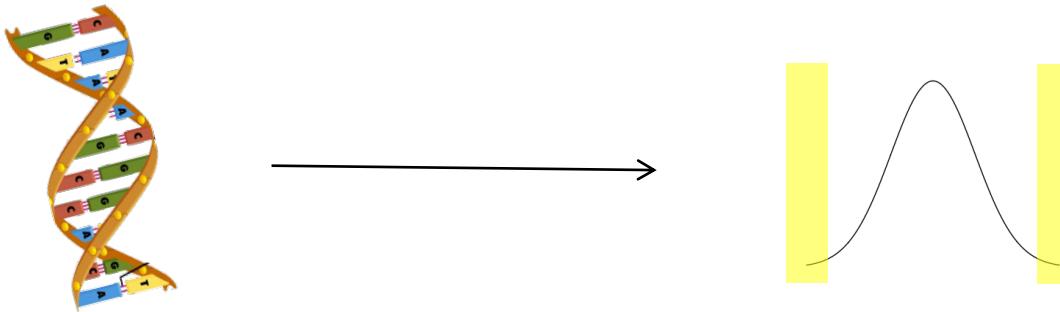


Variation due to age, sex, environmental factors (e.g. diet), and **genetic variation**. May be an effect of **multiple common variants that slightly alter normal physiological processes**.

Why find “disease genes”?

Genetic factors are particularly interesting because (unlike environmental factors) they are:

- inherited at birth
- essentially unchanging
- (often) easily measurable



This makes inferences about *causation* particularly simple.

e.g. compare:

“cyclists tend to be taller” => causation could plausibly work either way

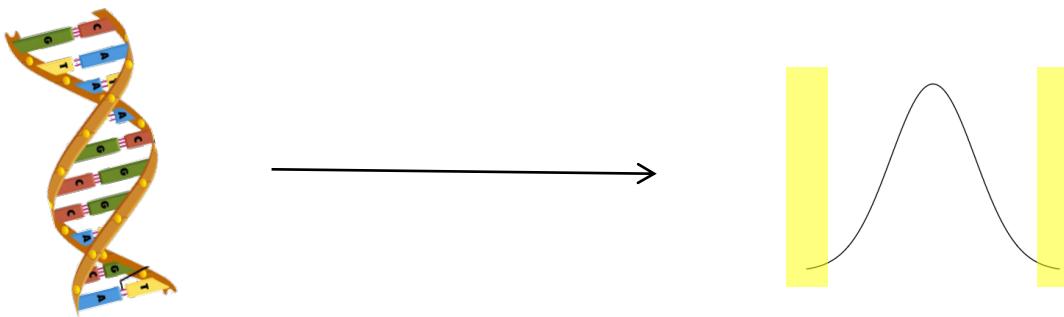
“people with genotype AA tend to be taller” => implies causation (all else being equal)

Inheritance = nature’s randomised control trial

Why find “disease genes”?

Genetic factors are particularly interesting because (unlike environmental factors) they are:

- inherited at birth
- essentially unchanging
- (often) easily measurable

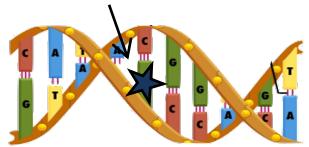


Reasons to look for disease genes:

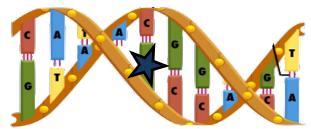
- Identify drug targets
- Predict risk of disease
- Personalised medicine (e.g. stratified by likely treatment response)
- Gene therapy?
- etc...
- ...understand the biology of disease

The circle of genetic causation

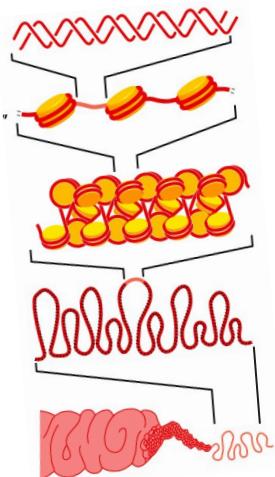
(a causal mutation)



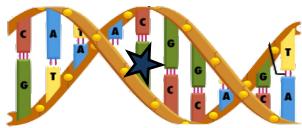
The circle of genetic causation



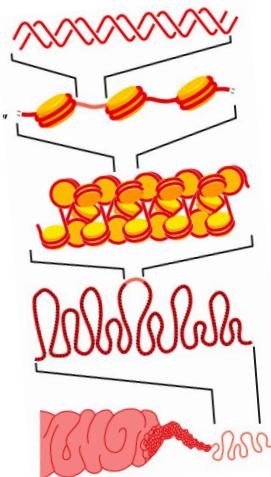
DNA gets physically packaged
up into chromosomes...



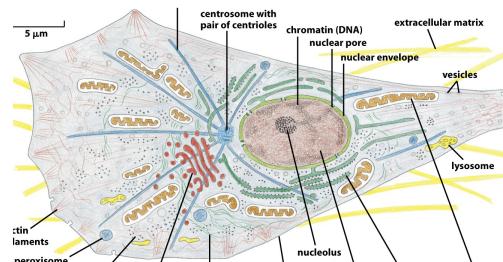
The circle of genetic causation



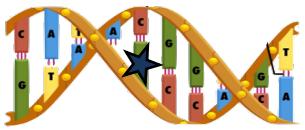
DNA gets physically packaged
up into chromosomes...



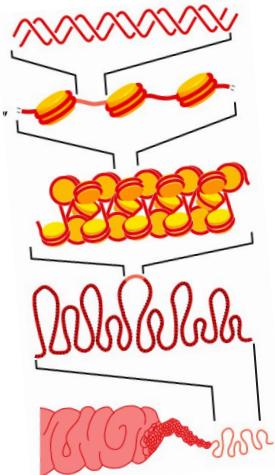
...inside cells, where it is
transcribed to form proteins
and other molecules...



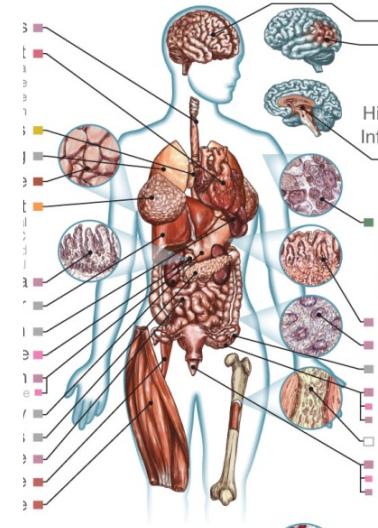
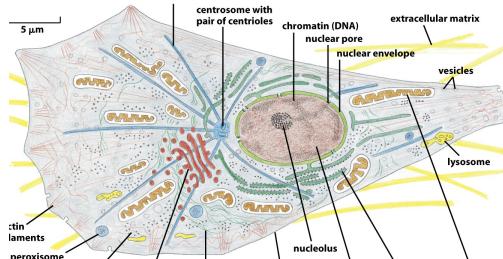
The circle of genetic causation



DNA gets physically packaged up into chromosomes...



...inside cells, where it is **transcribed** to form proteins and other molecules...

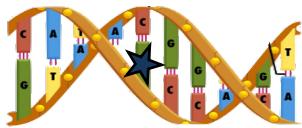


...that combine to make individuals...

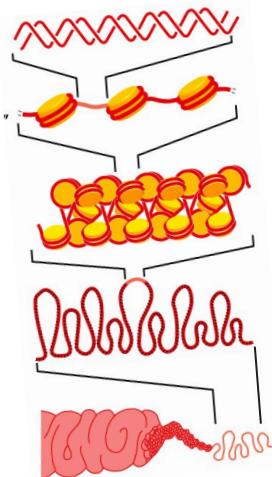


...that affect how the cells behave, forming different organs...

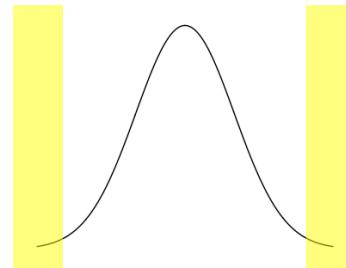
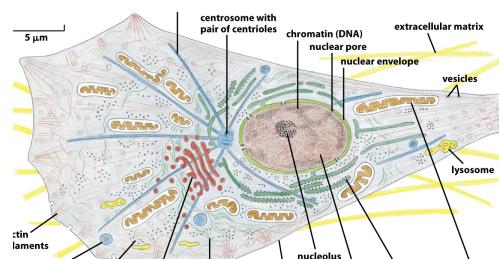
The circle of genetic causation



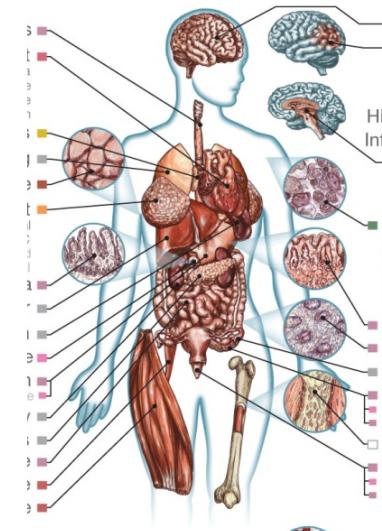
DNA gets physically packaged up into chromosomes...



...inside cells, where it is **transcribed** to form proteins and other molecules...



...whose success is affected by the traits they have...



...that combine to make individuals...



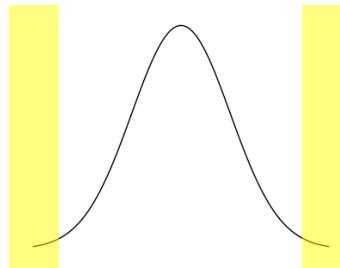
...that affect how the cells behave, forming different organs...

The circle of genetic causation

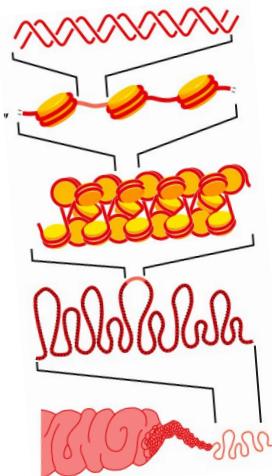
...passing on DNA, with
mutations and
recombination, to new
generations...



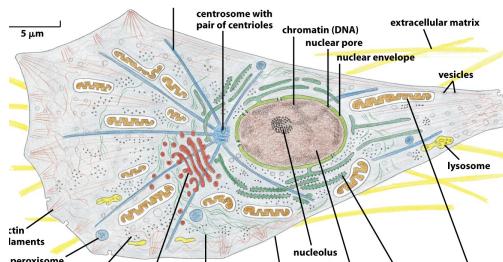
...whose success is affected by
the traits they have...



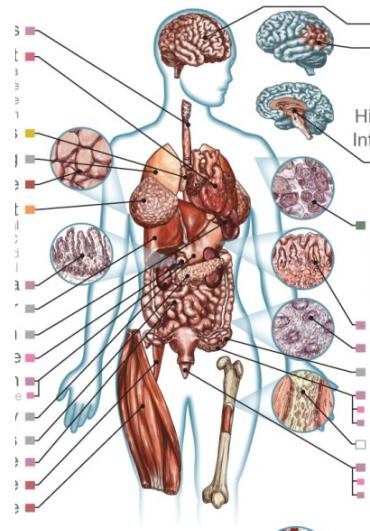
...that gets physically packaged
up into chromosomes...



...inside cells, where it is
transcribed to form proteins
and other molecules...

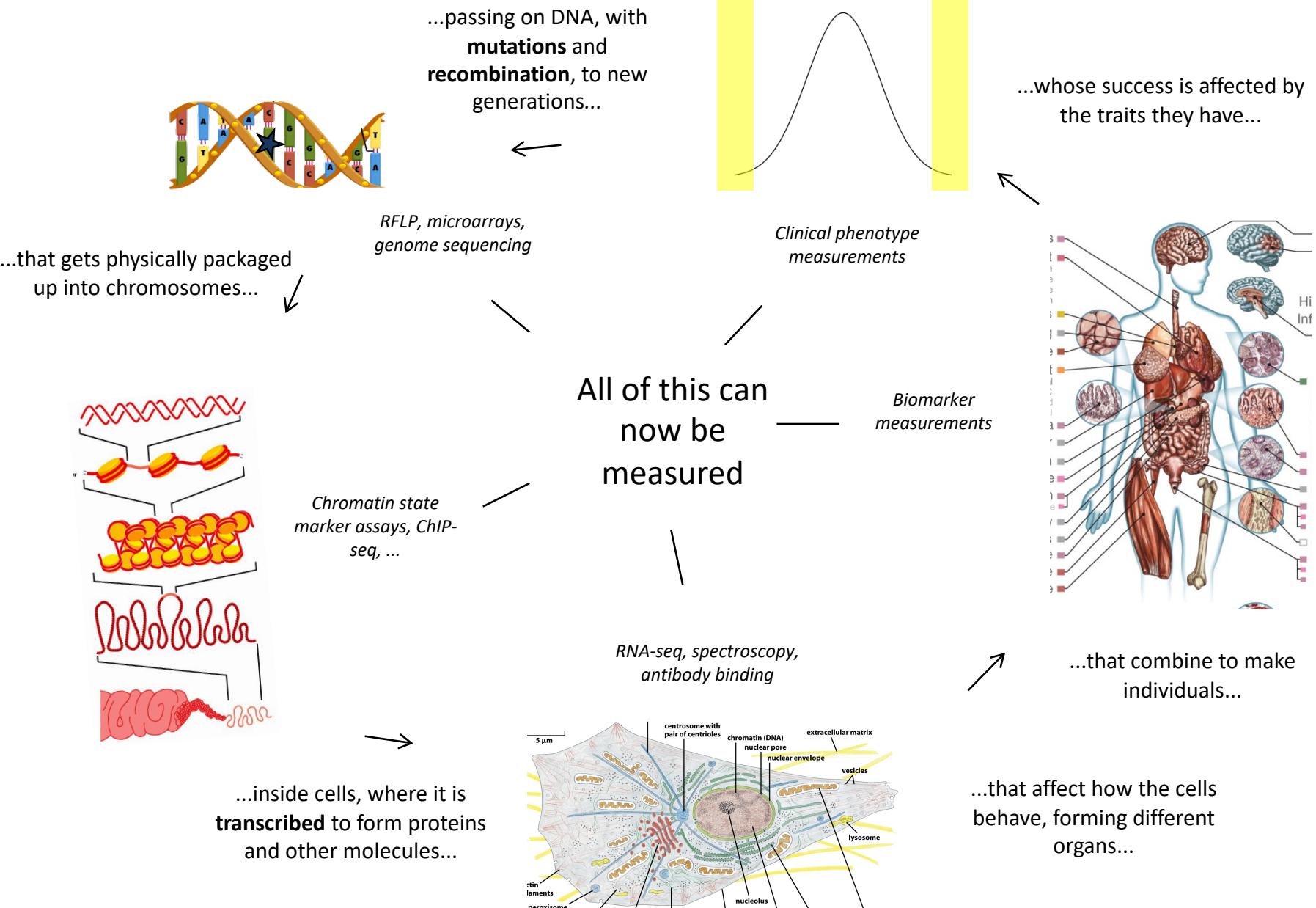


...that combine to make
individuals...



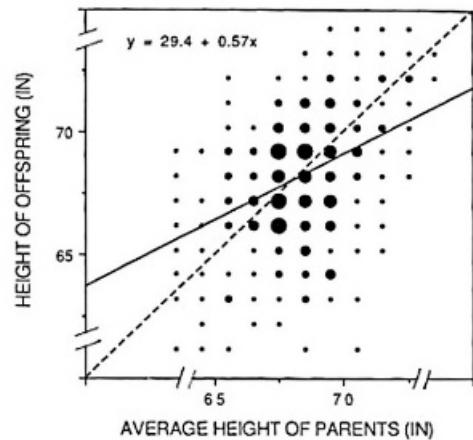
...that affect how the cells
behave, forming different
organs...

The circle of genetic causation



Finding disease genes in practice

- I'm going to assume we've got a trait that we've established is heritable



Demonstration by Francis Galton that human height is heritable (height of parents predicts height of offspring).

We want to find genetic variants influencing it.
How?

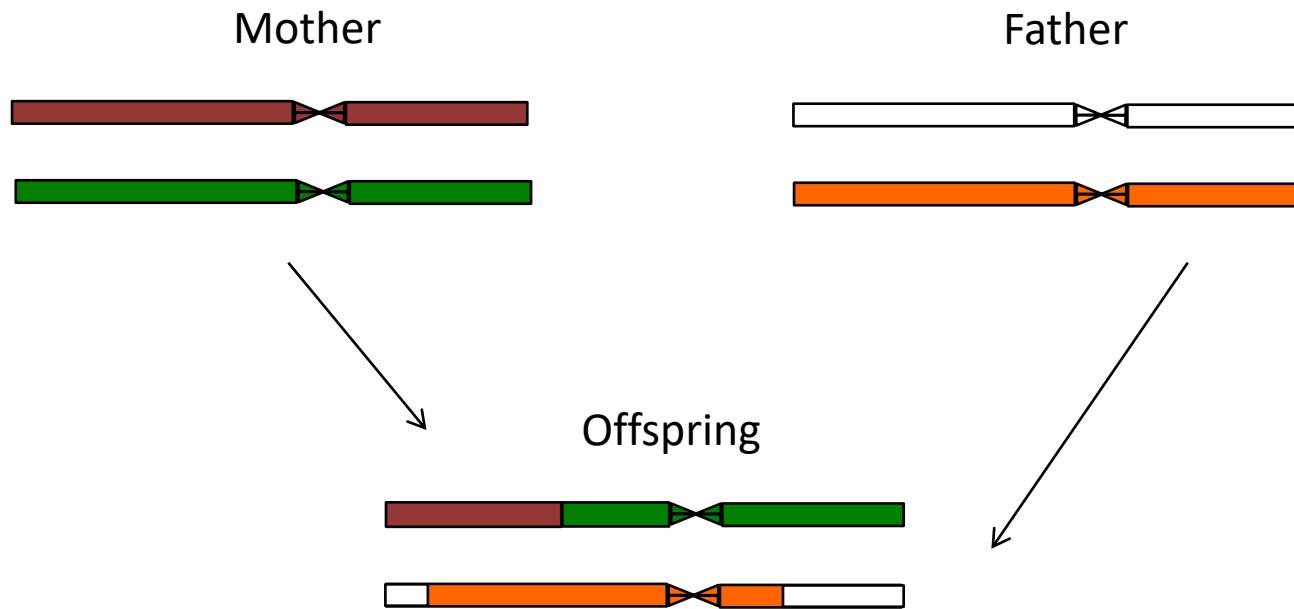
Finding needles in the haystack

The human genome is 3.2 billion base pairs long

We want to find a small number of ‘causal’ genetic mutations in there. How?

Luckily, nature has given us a way to narrow down on specific regions of the genome.

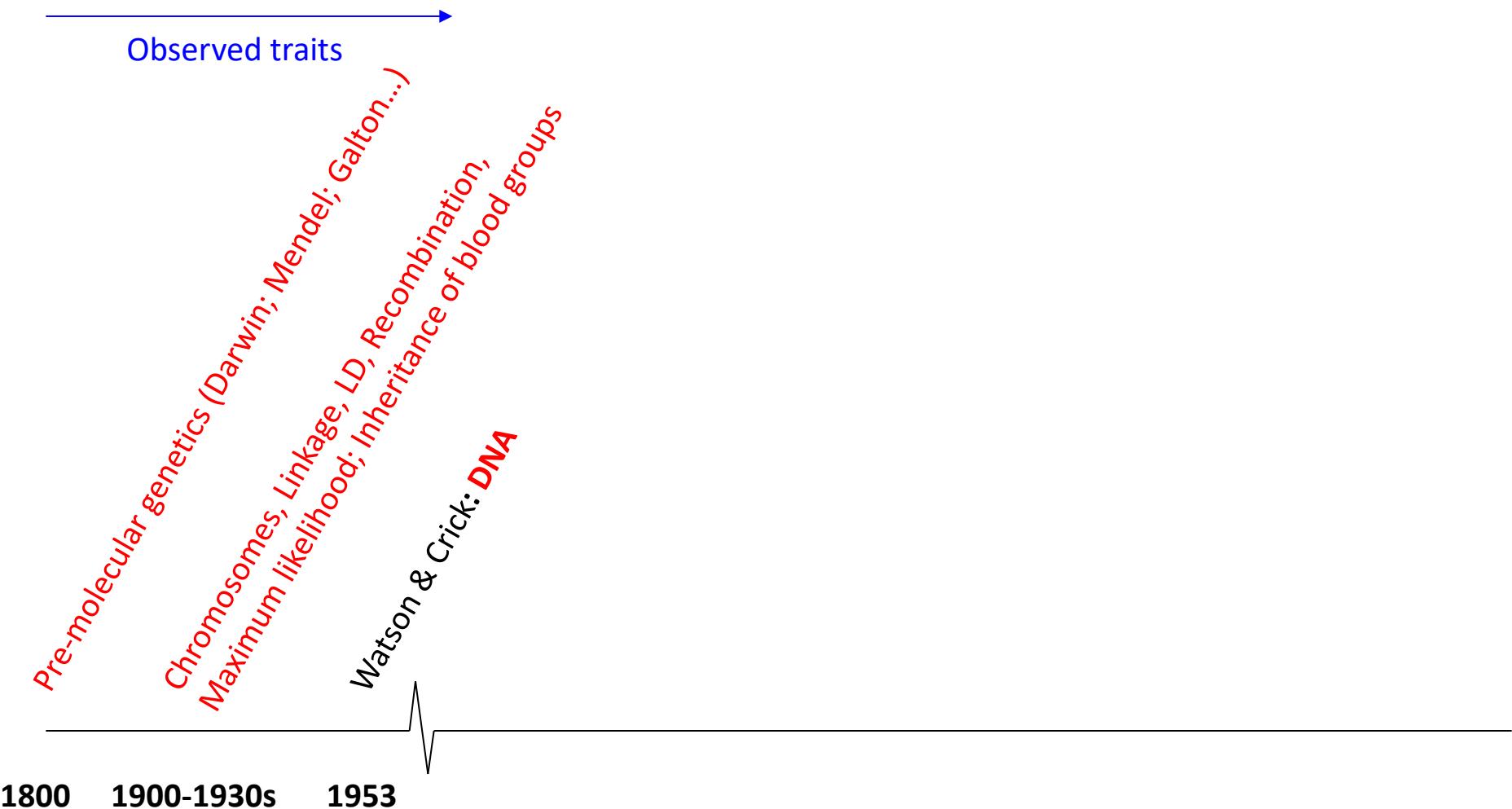
Recombination



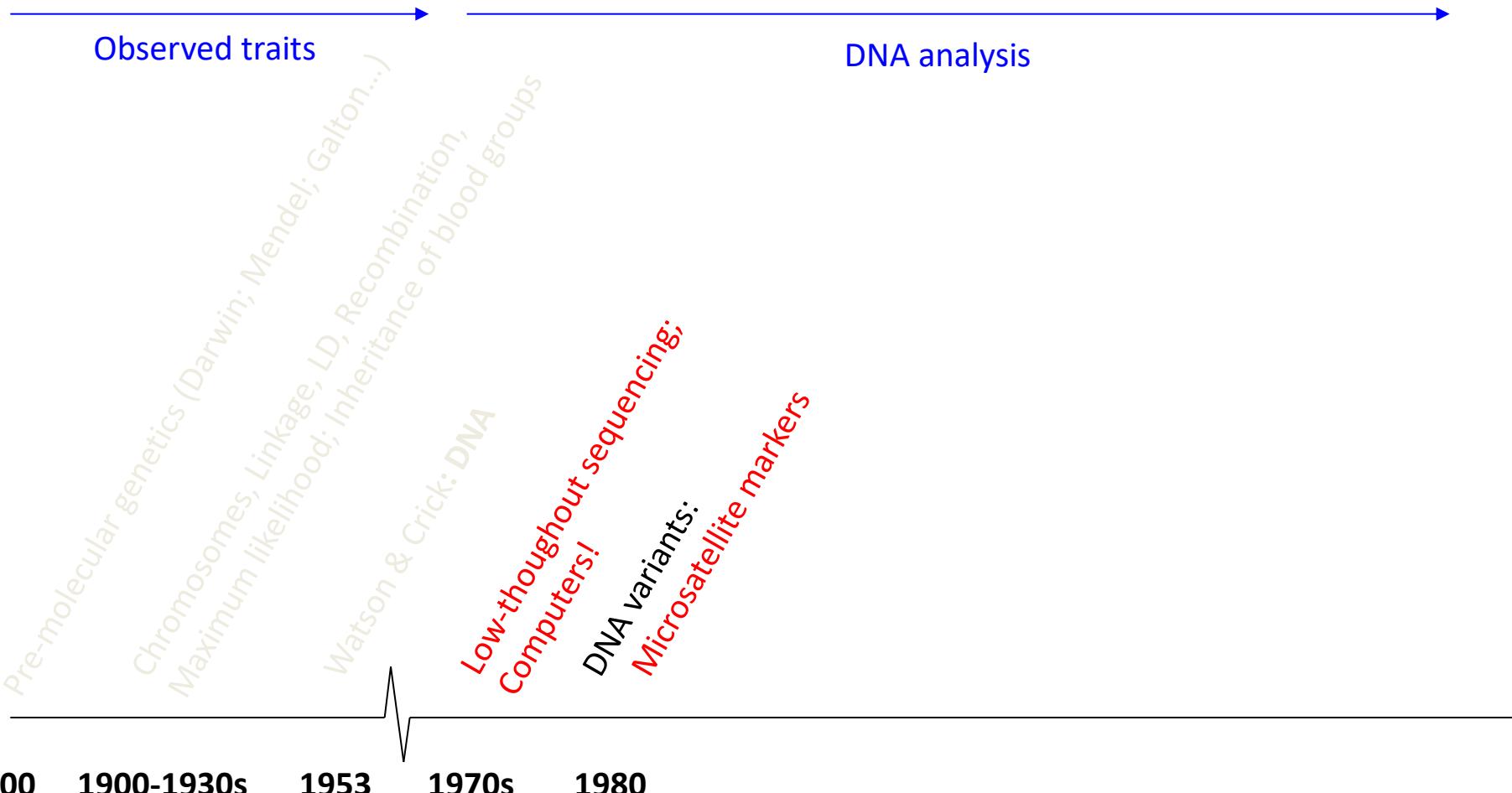
Genetic recombination breaks up the DNA into segments.
You inherit a mosaic of segments of your parents' DNA.

Recombination = nature's magnifying glass

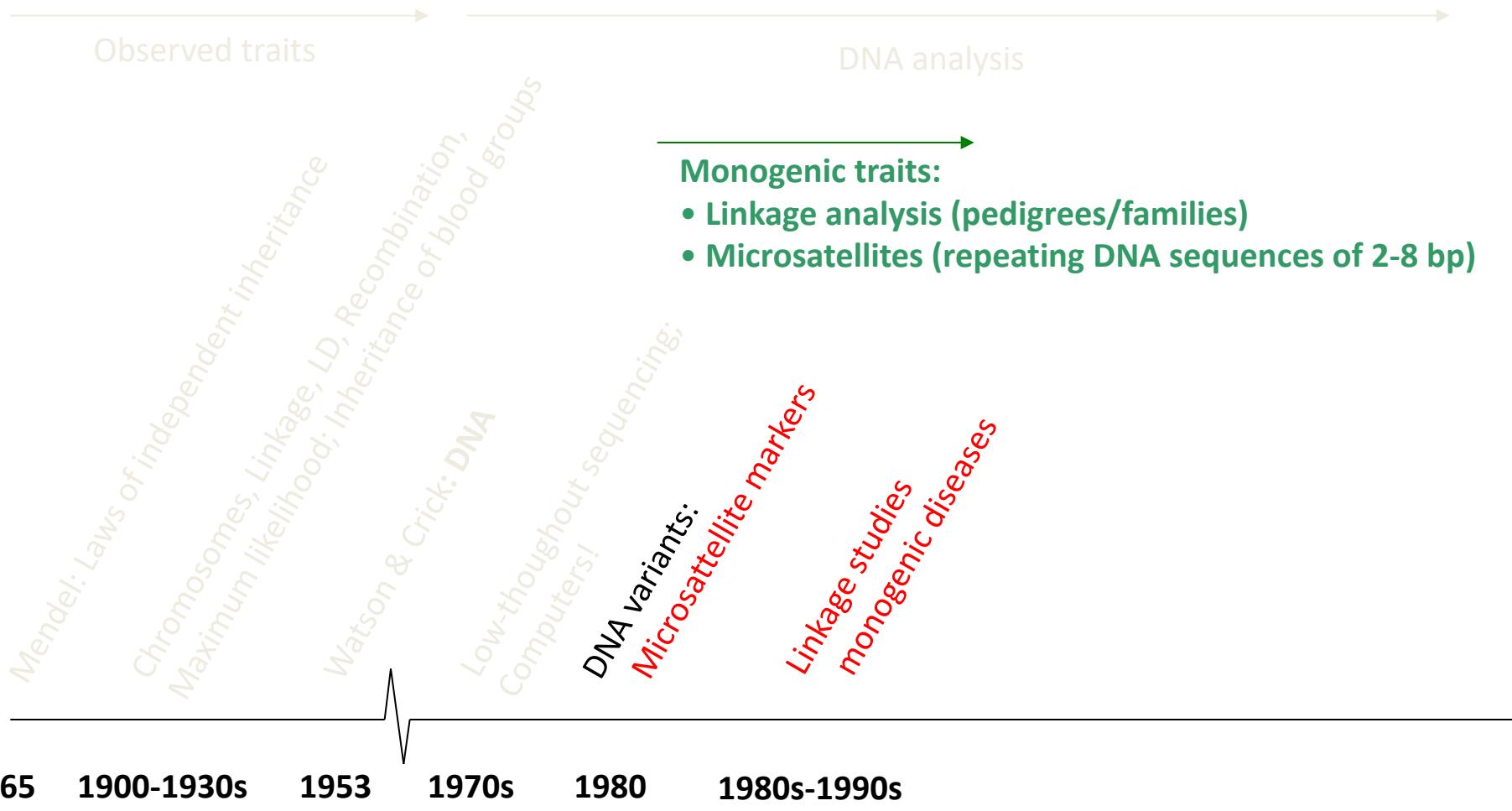
Studying the human genome



Studying the human genome



Studying the human genome



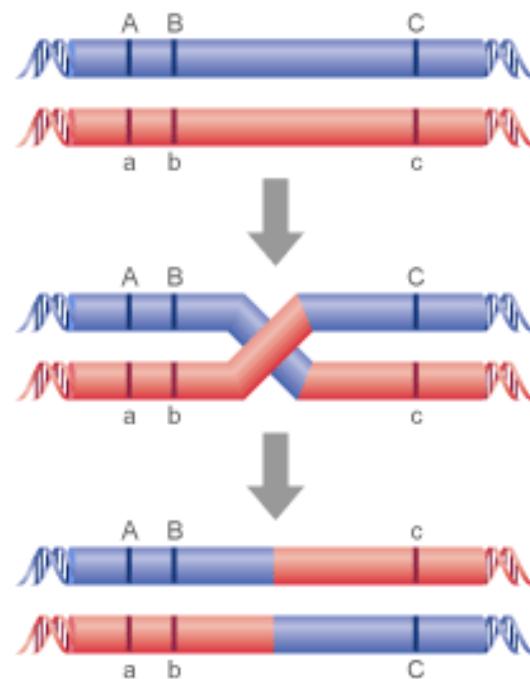
'Mendelian' vs. Complex diseases (traits)

Mendelian (mono-genic)

- Caused by single, **rare** mutations that confer a **very high risk**
- Disease primarily identified in families
- Examples: Huntington's chorea; Cystic fibrosis
- Typical study design: Family-based **Linkage Studies**

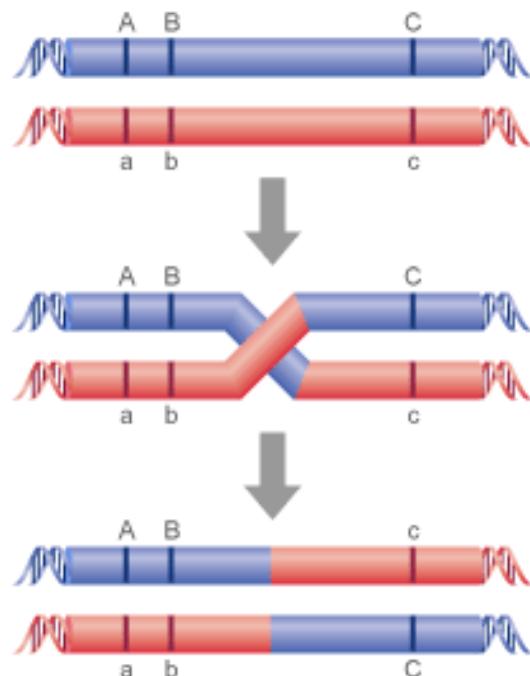
Linkage Study

- Aims to find a (large) section of chromosome that harbours a disease-predisposing gene variant, across **families**
- Utilises chromosomal recombination events during meioses across generations



Linkage Study

- Aims to find a (large) section of chromosome that harbours a disease-predisposing gene variant, across **families**
- Utilises chromosomal recombination events during meioses across generations



Process:

- Genotype ~ 400 microsatellite markers across the genome in families with multiple affecteds;
- Compare whether the allelic status of these variants are shared between affecteds in a family;
- Sum evidence of sharing across families - calculate if sharing is significantly more than expected under Mendelian laws of inheritance → evidence for *linkage*

Linkage Study

- Very successful for monogenic disorders, when carrying a variant means you develop disease (e.g. Cystic fibrosis, Huntingdon's chorea)
Low success rate for complex diseases:
caused by >1 gene; variant=susceptibility; environment
- Most powerful to detect variants responsible for disease in families, but otherwise **rare in the general population**

'Mendelian' vs. Complex diseases (traits)

Mendelian (mono-genic)

- Caused by single, **rare** mutations that confer a **very high risk**
- Disease primarily identified in families
- Examples: Huntingdon's chorea; Cystic fibrosis
- Typical study design: Family-based **Linkage Studies**

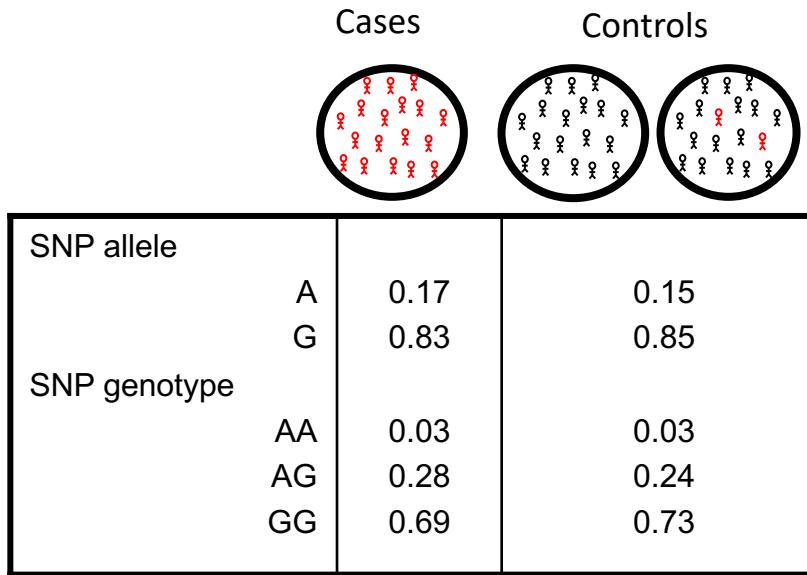
Complex

- Caused by multiple **common** genetic variants + environmental factors, each conferring a **modest increase in susceptibility**
- Disease observed in families as well as sporadic cases
- Examples: Cancers, CVD, Type I/II Diabetes, Body Mass Index.....
- Typical study design: **Association Studies**

NB Sub-types of complex diseases can behave as Mendelian traits,
e.g. Familial breast and ovarian cancer caused by *BRCA1/2* mutations.....

Association Study

Case-control

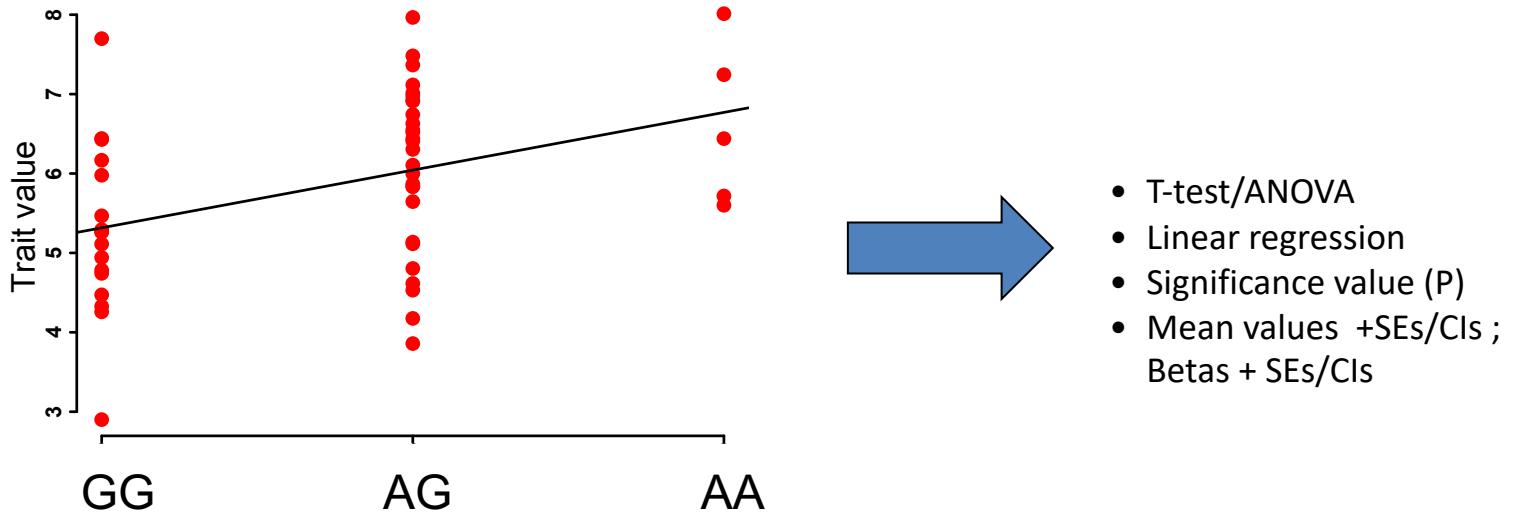


- Chisquare tests
- Logistic regression
- Significance value (P)
- ORs + CIs

- Uses single DNA base pair changes (SNPs)
- Aims to find a **common disease-predisposing** genetic variant (or get close by examining a SNP nearby)
- Powerful method for complex disease gene mapping, provided genetic variant causing disease is not rare

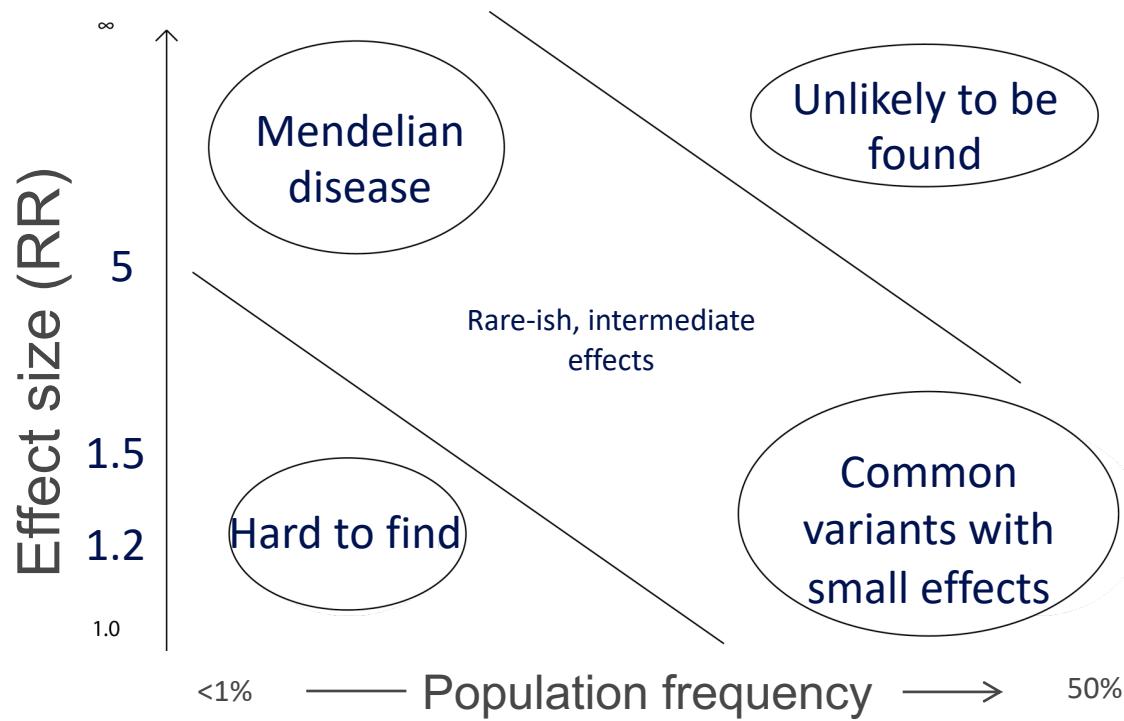
Association Study

Continuous trait



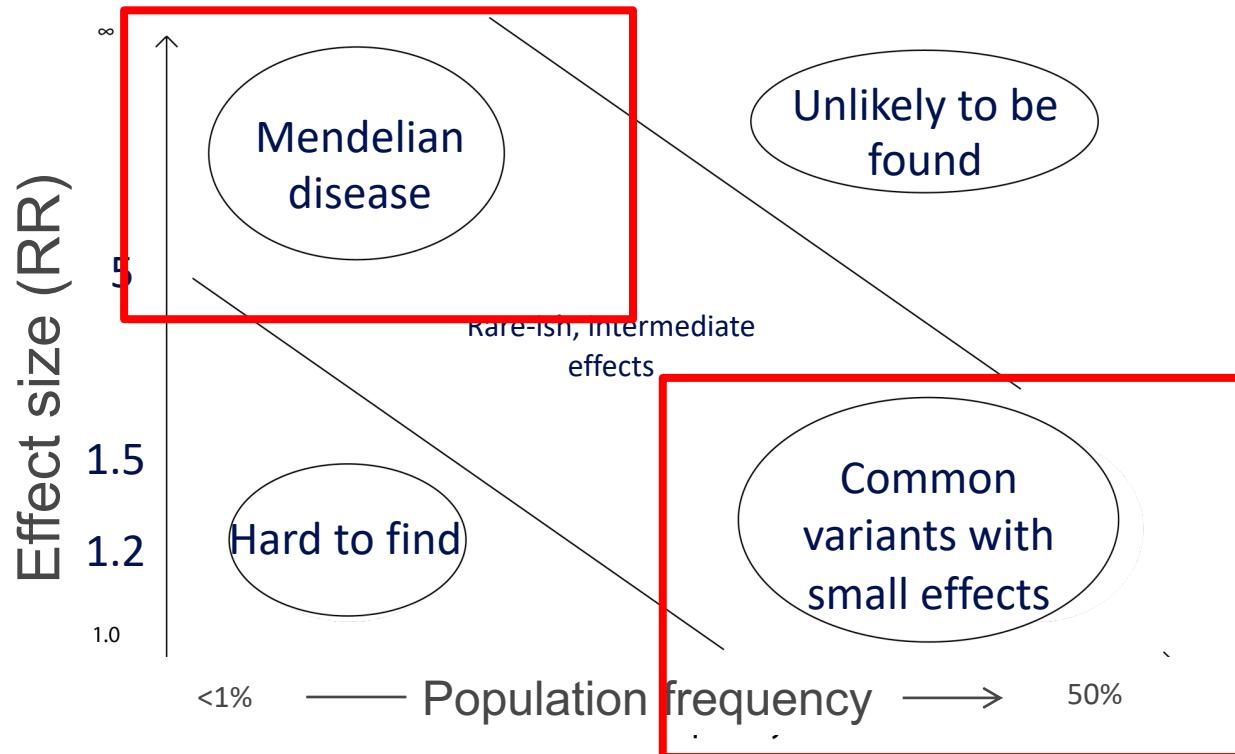
- Uses single DNA base pair changes (SNPs)
- Aims to find a **common trait-affecting** genetic variant (or get close by examining a SNP nearby)
- Powerful method for complex trait mapping, provided genetic variant causing disease is not rare

Complex diseases



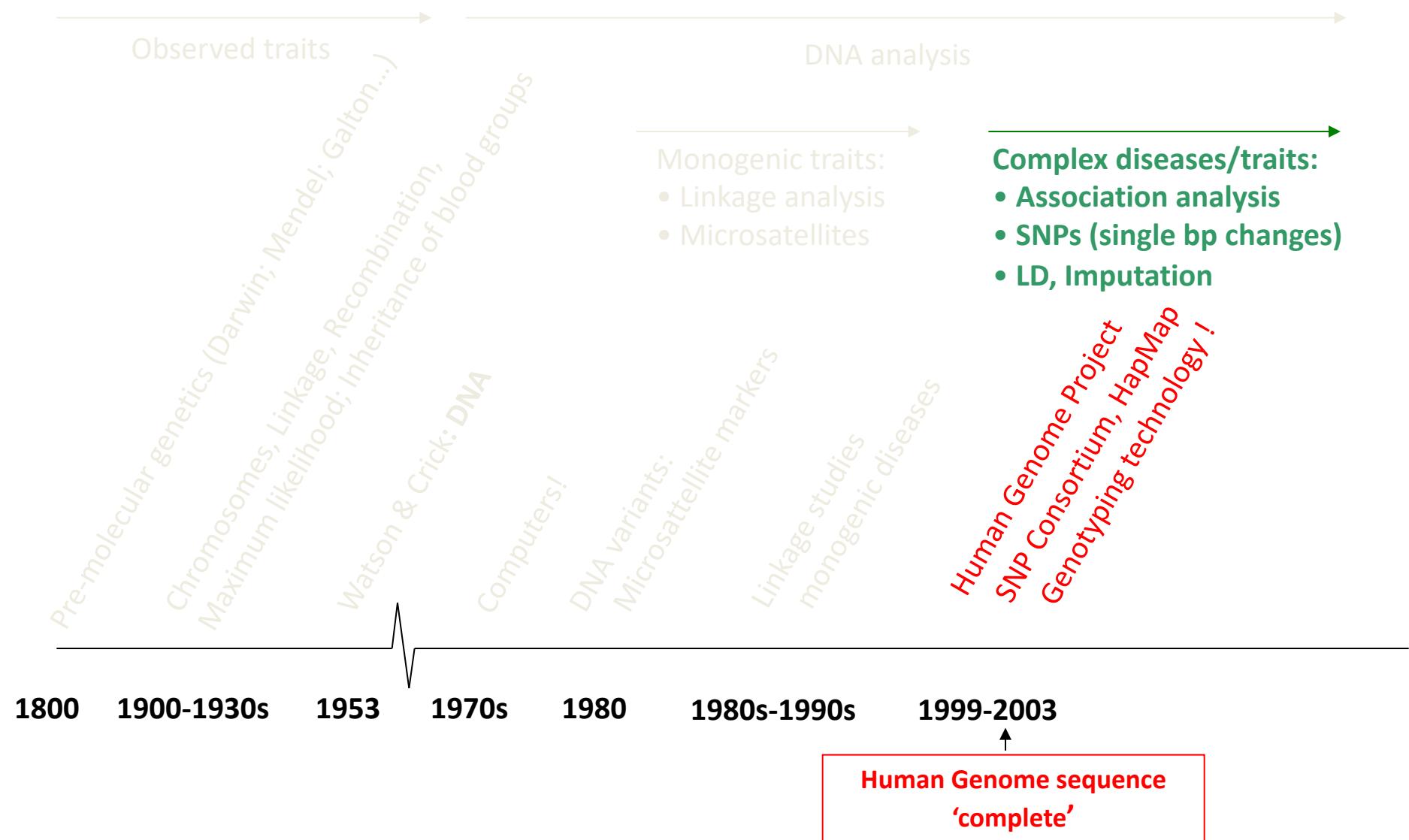
Complex diseases

Where linkage studies are likely to work



Where most complex disease effects are

Studying the human genome



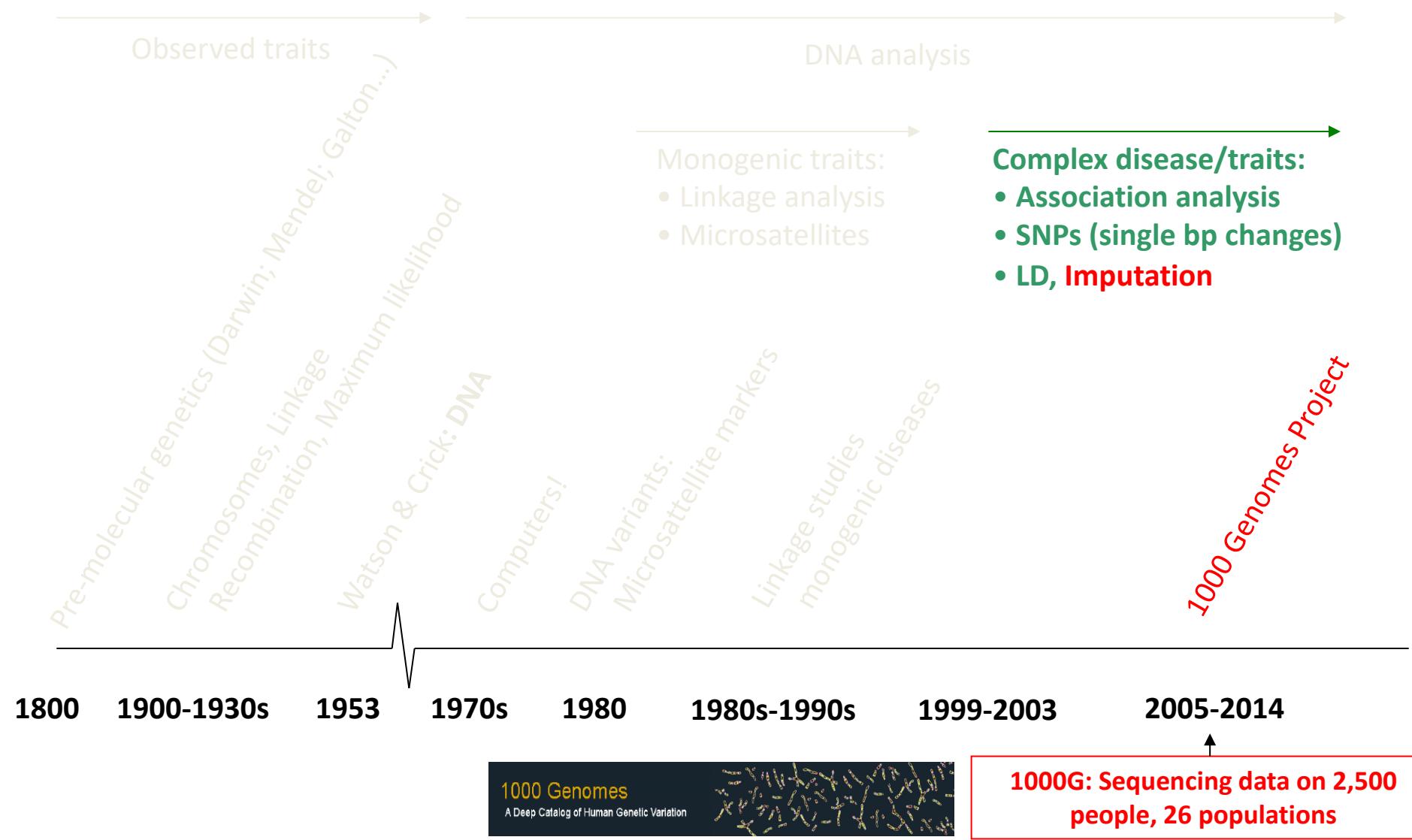
The International HapMap Project: LD

(International HapMap Consortium, Nature 2007; 2009; 2010)

- Common SNPs across the genome (frequency > 5%) are *not independent* when tested in a population, due to common ancestry between individuals, due to **Linkage disequilibrium (LD)**
- SNPs located near each other: the variant status of one predicts ('tags') the other
- The HapMap Project provided a map that shows SNP dependencies across the genome in different ancestral populations (Phase I/II: African-YRI, Asian-CHB/JPT, European— CEU; Phase III: + 9 other)
- We now 'only' need to genotype ~ a few 100K SNPs, to reliably predict the genotype status of all common SNPs in the genome, in a given population
- Means that **when you find an association with a SNP, you have not found the causal variant**, but are most likely 'tagging' the causal variant



Studying the human genome



From ~ 2000: Major developments in genetics knowledge and technology

*Utility: GWAS to detect **common variants** underlying complex diseases*

Human Genome Project

- Human genome: ~ 3 billion base-pairs
- SNPs: 1 per 100-300 bp : ~ 10M common SNPs
(freq >0.05)

Finishing the euchromatic sequence of the human genome

International Human Genome Sequencing Consortium*

NATURE | VOL 431 | 21 OCTOBER 2004 | www.nature.com/nature

Do we need to genotype 10M SNPs for each individual?

- Unnecessary – reference panel for imputation
- Too costly

A map of human genome sequence variation containing 1.42 million single nucleotide polymorphisms

The International SNP Map Working Group*

NATURE | VOL 409 | 15 FEBRUARY 2001 | www.nature.com



High through-put genotyping/sequencing

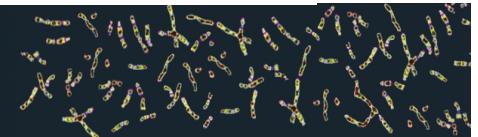


NATURE | VOL 426 | 18/25 DECEMBER 2003

NATURE | VOL 526 | 1 OCTOBER 2015

1000 Genomes

A Deep Catalog of Human Genetic Variation



Technology improvements

Cost and coverage



- Affymetrix 100K
- Affymetrix 500K
- Affymetrix 6.0 (~1M SNPs)
- ...
- Illumina 650Y
- Affymetrix UK Biobank array
- Illumina 1M
- Illumina 2.5M
- Illumina 5M
- ...

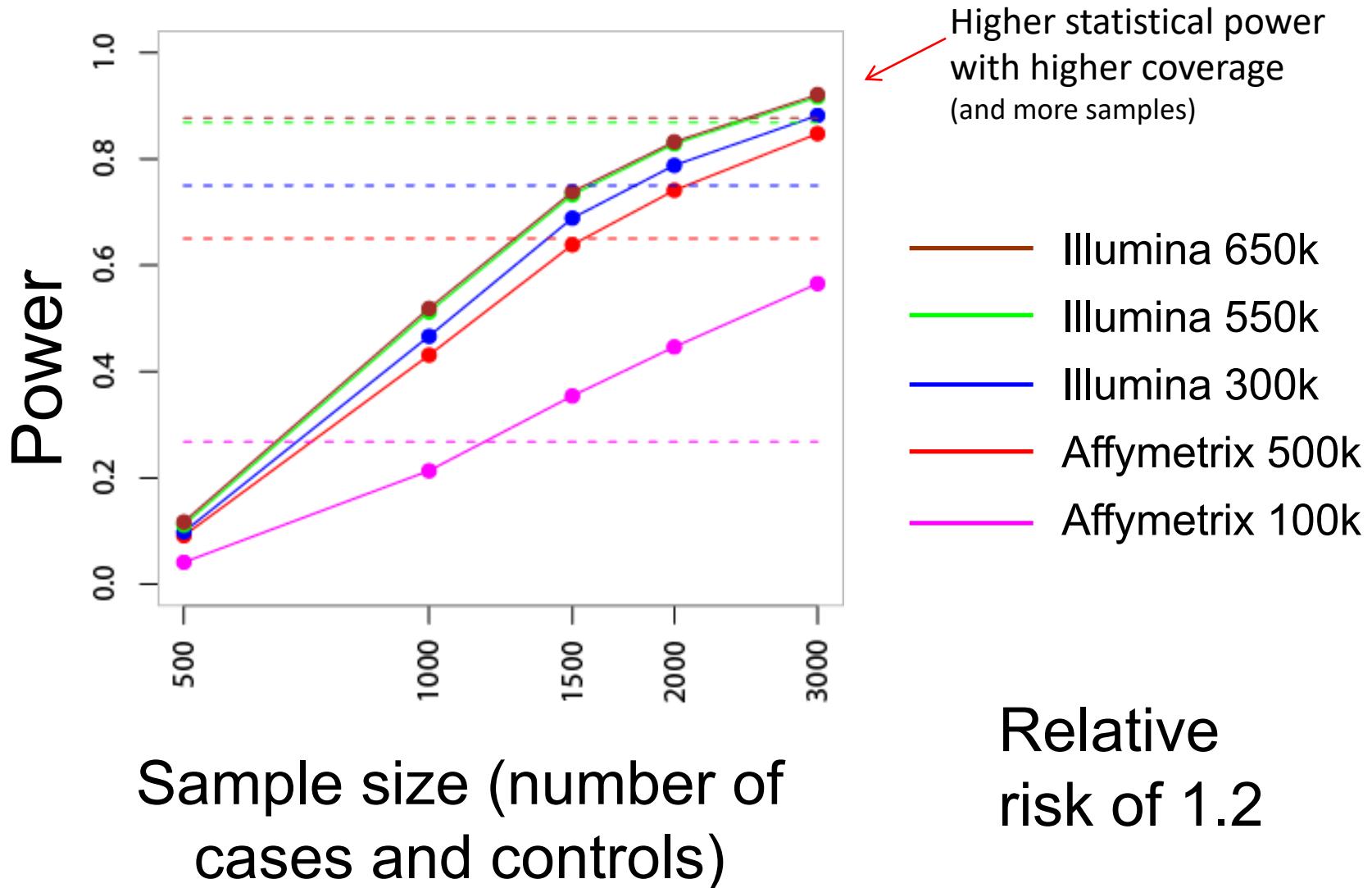


Costs are decreasing with time.

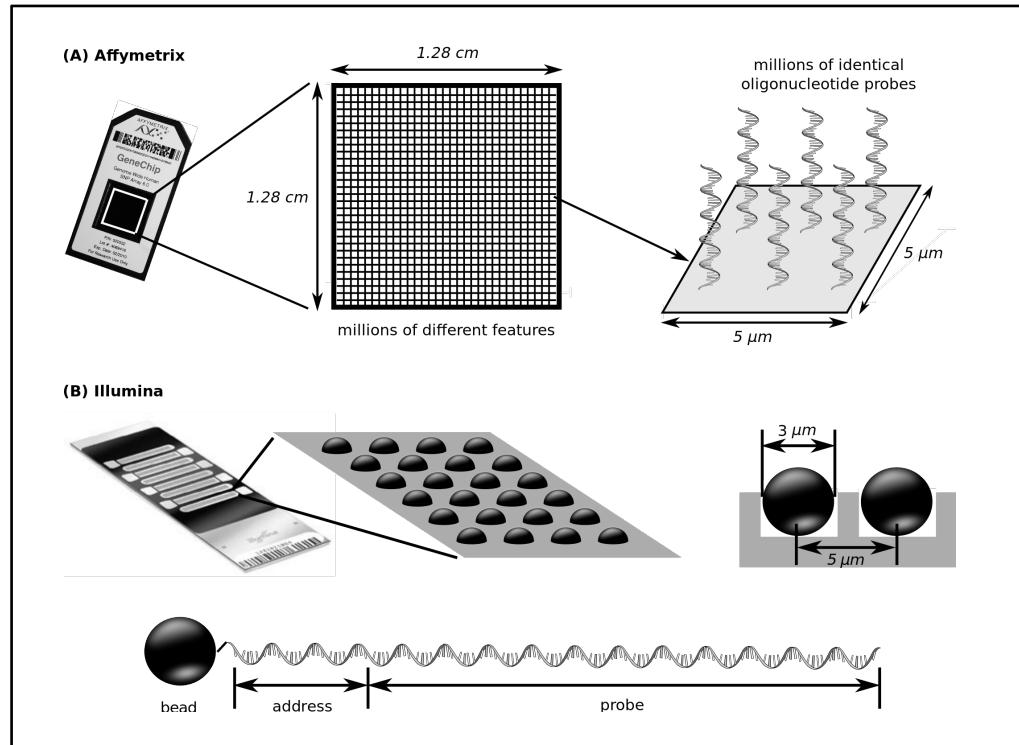
Example

- The UK Biobank has typed ~500,000 individuals on the Affymetrix UK Biobank array (containing ~800k SNPs).
- This array might now cost ~£20 per sample. So this project would cost in the order of £10,000,000 for genotyping.

Power to find weak effects



How a microarray works

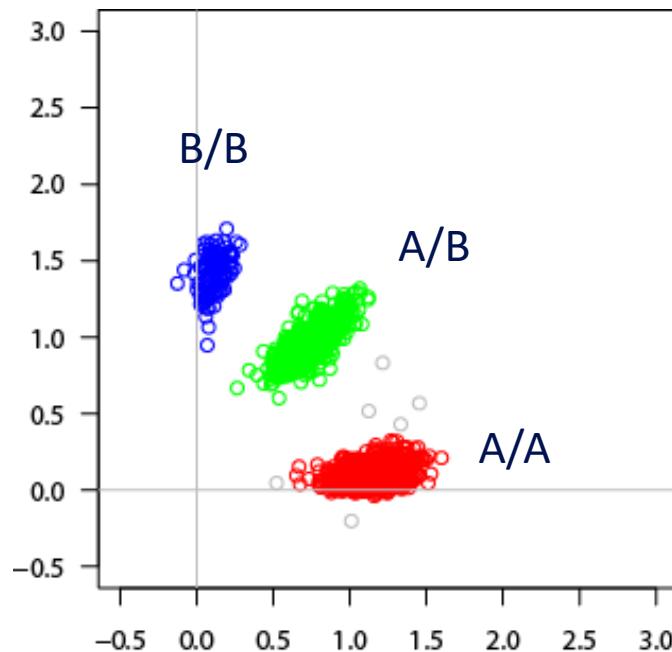


Wash the DNA over and let it hybridise to millions of probes – one for each SNP

Flourescent markers are then attached. A picture is taken of the array.

How a microarray works

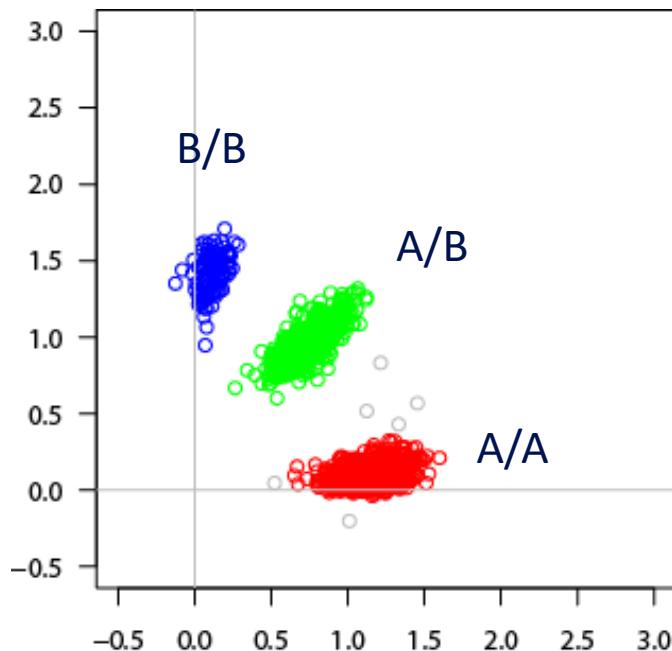
For each SNP, you get back
this:



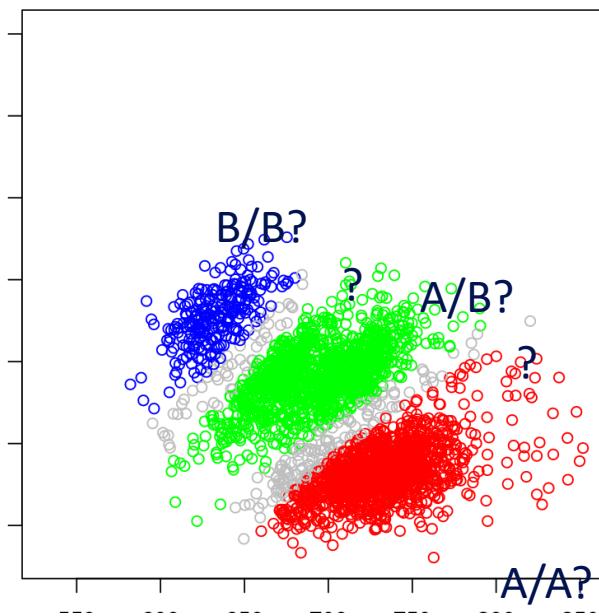
Each dot represents DNA from one individual.
X axis = image intensity for 1st allele probe
Y axis = image intensity for 2nd allele probe

How a microarray works

For each SNP, you get back this:



Or this if you're less lucky:



Each dot represents DNA from one individual.
X axis = image intensity for 1st allele probe
Y axis = image intensity for 2nd allele probe

Genome-wide Association Study (GWAS) Recipe

Genotype 100,000s **common**
SNPs in 1000s of cases+controls



Quality-control analyses:
e.g. genotype calling,
population biases



At each SNP test for allele frequency
difference between cases& controls
(χ^2 , logistic regression)



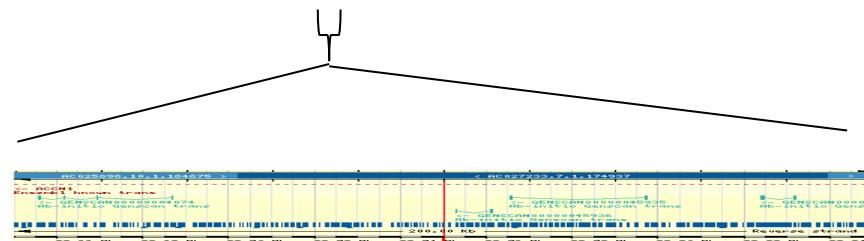
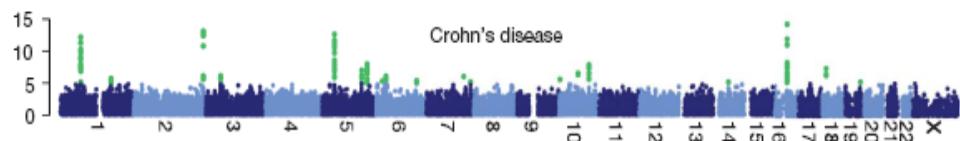
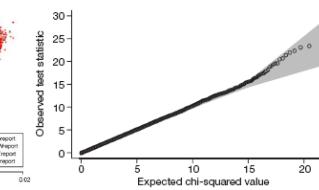
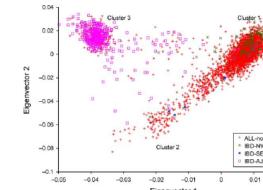
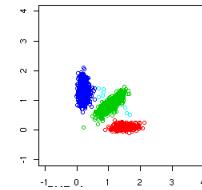
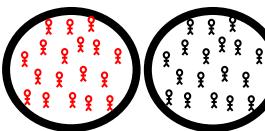
Identify significant associations,
nominal p-value (5×10^{-8})



Assess genomic info: genes, SNP
density, regulatory regions, etc



Genotype selected SNPs in different
case+control samples of same
population: replication/meta-analysis



More on this in Session 2

2005-7: first GWAS emerging....

Complement Factor H Polymorphism in Age-Related Macular Degeneration

Robert J. Klein,¹ Caroline Zeiss,^{2*} Emily Y. Chew,^{3*}
Jen-Yue Tsai,^{4*} Richard S. Sackler,¹ Chad Haynes,¹

Alice K. Henning,⁵ John Paul SanGiovanni,³ Shrikant M. Mane,⁶
Susan T. Mayne,⁷ Michael B. Bracken,⁷ Frederick L. Ferris,³
Jurg Ott,¹ Colin Barnstable,² Josephine Hoh^{7†}

SCIENCE VOL 308 15 APRIL 2005

Age-related Macular Degeneration (AMD):
100K SNPs in 96 cases and 50 controls!!

Rs380390, CC(CG vs. GG:
OR = 4.6 (2.0-11.0), P=4.1 x 10⁻⁸

2005-7: first GWAS emerging....

Complement Factor H Polymorphism in Age-Related Macular Degeneration

Robert J. Klein,¹ Caroline Zeiss,^{2*} Emily Y. Chew,^{3*} Jen-Yue Tsai,^{4*} Richard S. Sackler,¹ Chad Haynes,¹

Alice K. Henning,⁵ John Paul SanGiovanni,³ Shrikant M. Mane,⁶ Susan T. Mayne,⁷ Michael B. Bracken,⁷ Frederick L. Ferris,³ Jurg Ott,¹ Colin Barnstable,² Josephine Hoh^{7†}

SCIENCE VOL 308 15 APRIL 2005

Age-related Macular Degeneration (AMD):
100K SNPs in 96 cases and 50 controls!!

Rs380390, CC(CG vs. GG:
OR = 4.6 (2.0-11.0), P=4.1 x 10⁻⁸

Genome-wide association study of 14,000 cases of seven common diseases and 3,000 shared controls

The Wellcome Trust Case Control Consortium*

NATURE|Vol 447|7 June 2007

A genome-wide association study identifies novel risk loci for type 2 diabetes

Robert Sladek^{1,2,4}, Ghislain Rocheleau^{1*}, Johan Rung^{4*}, Christian Dina^{5*}, Lishuang Shen¹, David Serre¹, Philippe Boutilier⁵, Daniel Vincent⁴, Alexandre Belisle⁴, Samy Hadjadj⁶, Beverley Balkau⁷, Barbara Heude⁷, Guillaume Charpentier⁸, Thomas J. Hudson^{4,9}, Alexandre Montpetit⁴, Alexey V. Pshezhetsky¹⁰, Marc Prentki^{10,11}, Barry I. Posner^{2,12}, David J. Balding¹³, David Meyre⁵, Constantin Polychronakos^{1,3} & Philippe Froguel^{3,14}

NATURE|Vol 445|22 February 2007

A Genome-Wide Association Study Identifies *IL23R* as an Inflammatory Bowel Disease Gene

Richard H. Duerr,^{1,2} Kent D. Taylor,^{3,4} Steven R. Brant,^{5,6} John D. Rioux,^{7,8} Mark S. Silverberg,⁹ Mark J. Daly,^{8,10} A. Hillary Steinhart,⁹ Clara Abraham,¹¹ Miguel Regueiro,¹ Anne Griffiths,¹² Themistocles Dassopoulos,⁵ Alain Bitton,¹³ Huiying Yang,^{3,4} Stephan Targan,^{4,14} Lisa Wu Datta,⁵ Emily O. Kistner,¹⁵ L. Philip Schumm,¹⁵ Annette T. Lee,¹⁶ Peter K. Gregersen,¹⁶ M. Michael Barmada,² Jerome I. Rotter,^{3,4} Dan L. Nicolae,^{11,17} Judy H. Cho^{18*}

SCIENCE VOL 314 1 DECEMBER 2006

~2,000 cases for each of 7 diseases + ~3,000 shared controls
Bipolar disorder, CAD, Crohn's, Hypertension, RA, Type I/II Diabetes

ORs ranging from 1.3 – 2.0 (5-15: MHC)

2007: Wellcome Trust Case-Control Consortium

Genome-Wide Association Across Major Human Diseases

DESIGN

Collaboration amongst 26 UK disease investigators

2000 cases each from 7 diseases

GENOTYPING

Affymetrix 500k SNPs + 175k Perlegen highly selected SNPs

CASES

1. Type 1 Diabetes
2. Type 2 Diabetes
3. Crohn's Disease
4. Coronary Heart Disease
5. Hypertension
6. Bipolar Disorder
7. Rheumatoid Arthritis

Genome-wide association study of 14,000 cases of seven common diseases and 3,000 shared controls

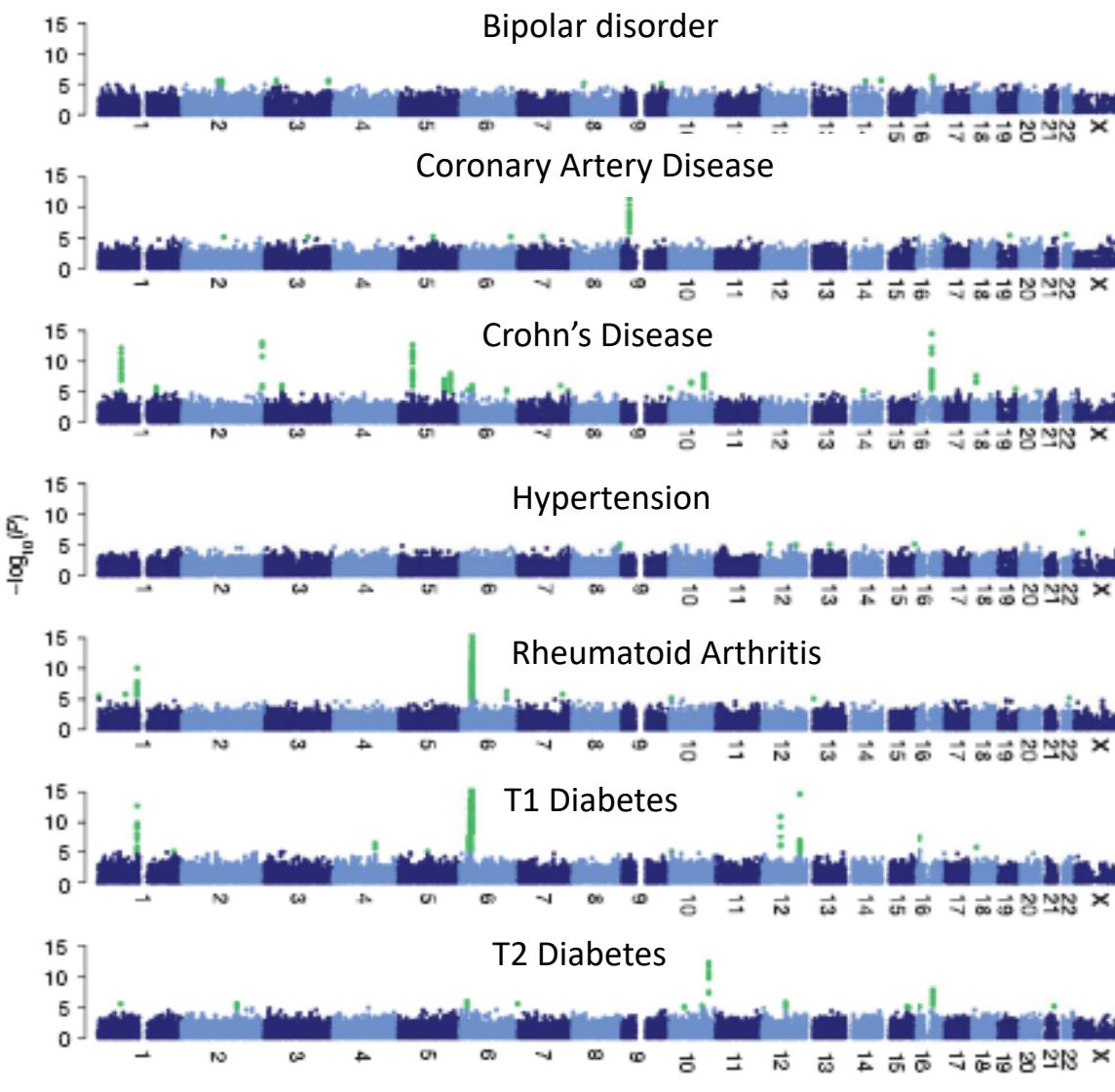
The Wellcome Trust Case Control Consortium*

NATURE | Vol 447 | 7 June 2007

CONTROLS

1. UK Controls A (n=1500: 1958 BC)
2. UK Controls B (n=1500: NBS)

Wellcome Trust Case-Control Consortium



Significant hits

16p12

9p21

1p31 5p13 10q24

2q37 5q33 16q12

3p21 10q21 18p11

-

6 (MHC)

1p13

7q32 (females only!)

6 (MHC) 12q13

1p13 12p13

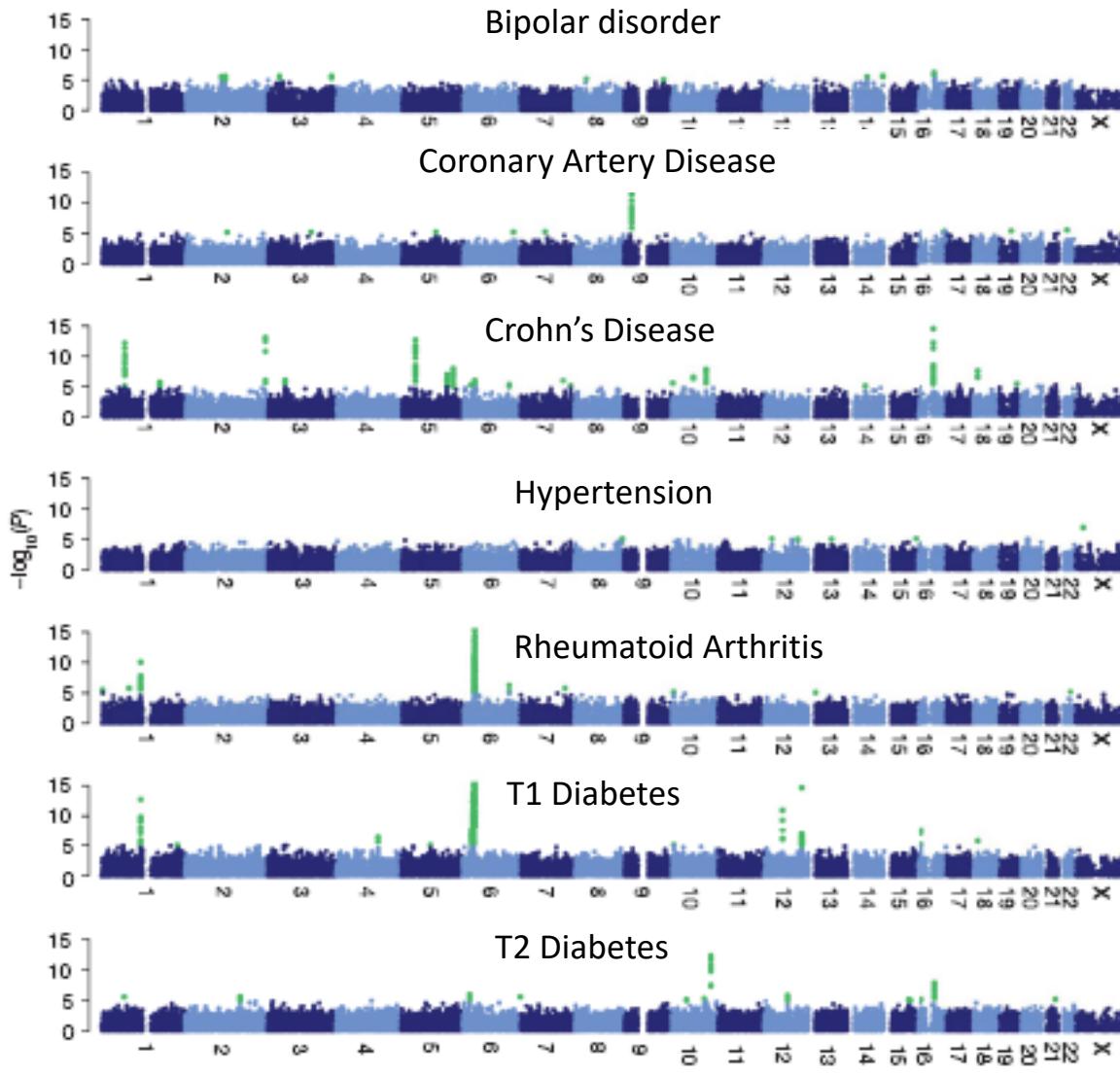
4q27 12q24 16p13

6p22 10q25

16q12

Wellcome Trust Case-Control Consortium

Replicated!



Significant hits

16p12

9p21

1p31 5p13 10q24
2q37 5q33 16q12
3p21 10q21 18p11

6 (MHC)

1p13

7q32 (females only)

6 (MHC) 12q13

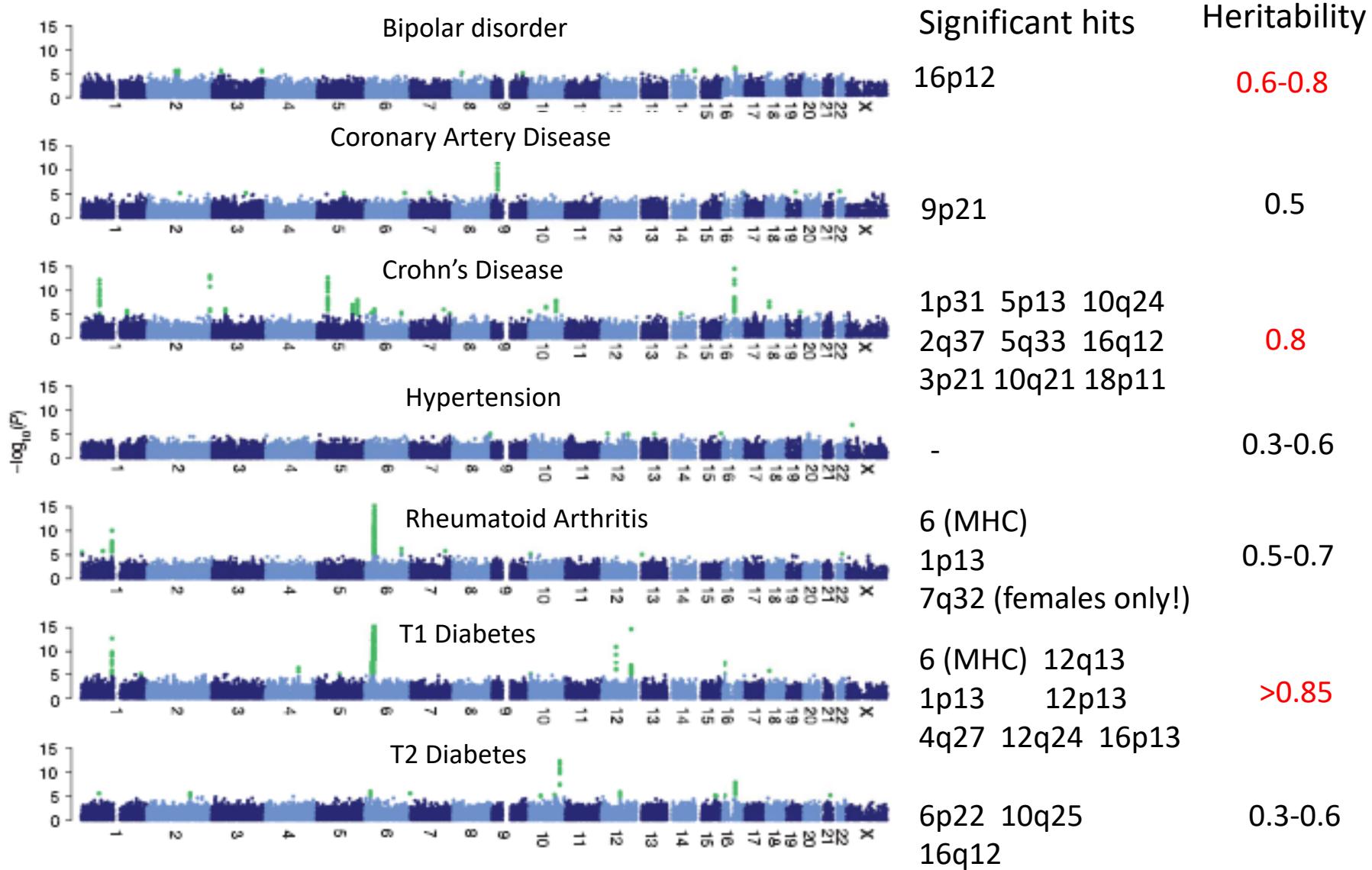
1p13 12p13

4q27 12q24 16p13

6p22 10q25

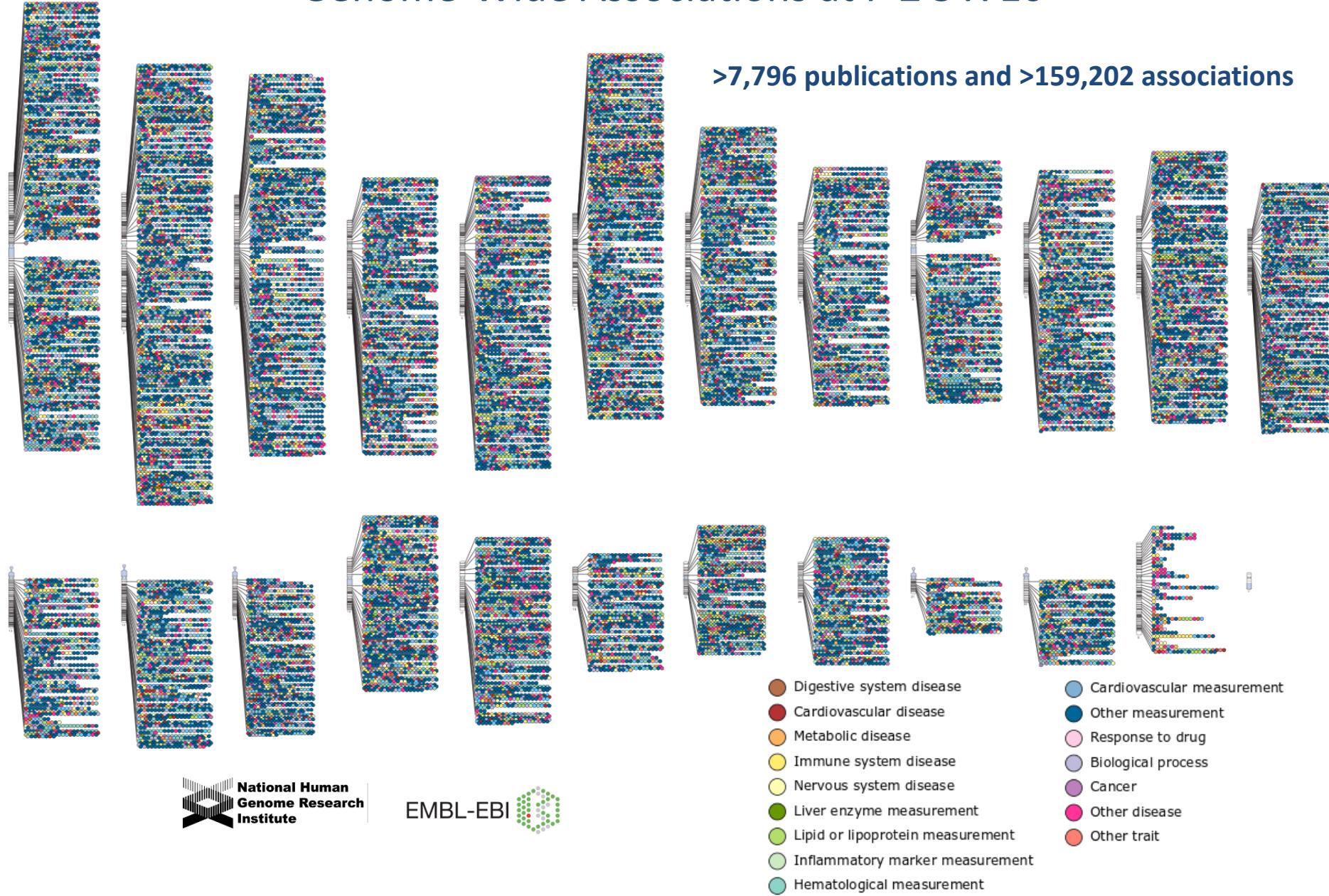
16q12

Wellcome Trust Case-Control Consortium



NHGRI GWA study Catalog: published Genome-Wide Associations at $P \leq 5 \times 10^{-8}$

>7,796 publications and >159,202 associations



Complex diseases: Rare variants?

- Low-frequency (MAF 0.5-5%) and rare (MAF<0.5%) variation may contribute to the heritability of complex traits:
 - IFIH1 and type 1 diabetes;
 - MYH6 and sick sinus syndrome.
- Assaying rare variants: Exome genotyping arrays or re-sequencing
- Rare variants are expected to have larger effects on complex traits than common variants.
- Single variant tests lack power to detect association.
- Statistical methods focus on the accumulation of minor alleles at rare variants (mutational load) within the same “genomic unit”.
 - Burden tests: assumes same direction of effect of all rare variants
 - Dispersion tests: allow for different direction of effect

Summary

It is clinically useful and interesting to look for the genetic variants contributing to human traits

A combination of theory, understanding of population genetics, and technology has made it possible to carry out GWAS analysis.

GWAS are an appropriate tool to identify common variants of relatively small risk (<1.5) for complex traits/disease.

- Requires tens of thousands of samples
- Good coverage of the SNP markers