

# BIOINFORMATICS

**Kristel Van Steen, PhD<sup>2</sup>**

**Montefiore Institute - Systems and Modeling**

**GIGA - Bioinformatics**

**ULg**

**[kristel.vansteen@ulg.ac.be](mailto:kristel.vansteen@ulg.ac.be)**

## CHAPTER 3: GENOME-WIDE ASSOCIATION STUDIES

### 1 Setting the pace

#### 1.a “The Human Genome Project”

#### 1.b The concept of a genetic marker

#### 1.c The rise of Genome-Wide Association studies (GWAs)

### 2 Components of a GWAs

### 3 Study Design

#### 3.a Marker level

#### 3.b Subject level

#### 3.c Gender level (not considered in this course)

## 4 Pre-analysis Steps

### 3.a Quality-Control: The Travemünde criteria

Hardy-Weinberg equilibrium

### 3.b Linkage disequilibrium

### 3.c Confounding: population stratification

## 5 Testing for Associations

## 6 Replication and Validation

## 7 GWA Interpretation and Follow-Up

# 1 Setting the pace

## 1.a “The Human Genome Project”

genome.gov  
National Human Genome Research Institute  
National Institutes of Health

Research Funding | Research at NHGRI | Health | **Education** | Issues in Genetics | Newsroom | Careers & Training | About | For You

Home > Education > All About The Human Genome Project (HGP)

**Education**

- All About The Human Genome Project (HGP)
- Education Archive
- Fact Sheets
- Genetic Education Resources for Teachers
- NHGRI Webinar Series
- National DNA Day
- Online Genetics Education Resources
- Smithsonian NHGRI Genome Exhibition
- Talking Glossary
- Understanding the Human Genome Project

**All About The Human Genome Project (HGP)**

The Human Genome Project (HGP) was one of the great feats of exploration in history - an inward voyage of discovery rather than an outward exploration of the planet or the cosmos; an international research effort to sequence and map all of the genes - together known as the genome - of members of our species, *Homo sapiens*. Completed in April 2003, the HGP gave us the ability, for the first time, to read nature's complete genetic blueprint for building a human being.

In this section, you will find access to a wealth of information on the history of the HGP, its progress, cast of characters and future.

- [Educational Resources](#)
- [General Information](#)
- [Research](#)
- [Model Organisms](#)

**Educational Resources**

- [An Interactive Timeline of the Human Genome](#) [unlockinglifescode.org]  
An interactive, hyper-linked timeline of genetics that takes the reader from Mendel (1865) to the completion of the mapping of the human genome (2003).
- [The Human Genome: A Decade of Discovery. Creating a Healthy Future](#)  
A workshop for science reporters about the 10th anniversary of the completion of the draft sequence of the human genome and to look at the future of genomic research.
- [Understanding the Human Genome Project](#)  
NHGRI's Online Education Kit
- [An Overview of the Human Genome Project](#)  
A brief overview of the HGP.
- [50 Years of DNA: From Double Helix to Health](#)  
Information about the celebration of the completion of the HGP and the 50th anniversary of the discovery of the

**See Also:**

- [White House Announcement](#)  
June 26, 2000
- [Extramural Research Program](#)
- [Other Federal Agencies Involved in Genomics](#)

**On Other Sites:**

- [Human Genome Resources](#)  
Access to the full human sequence

## Historical overview



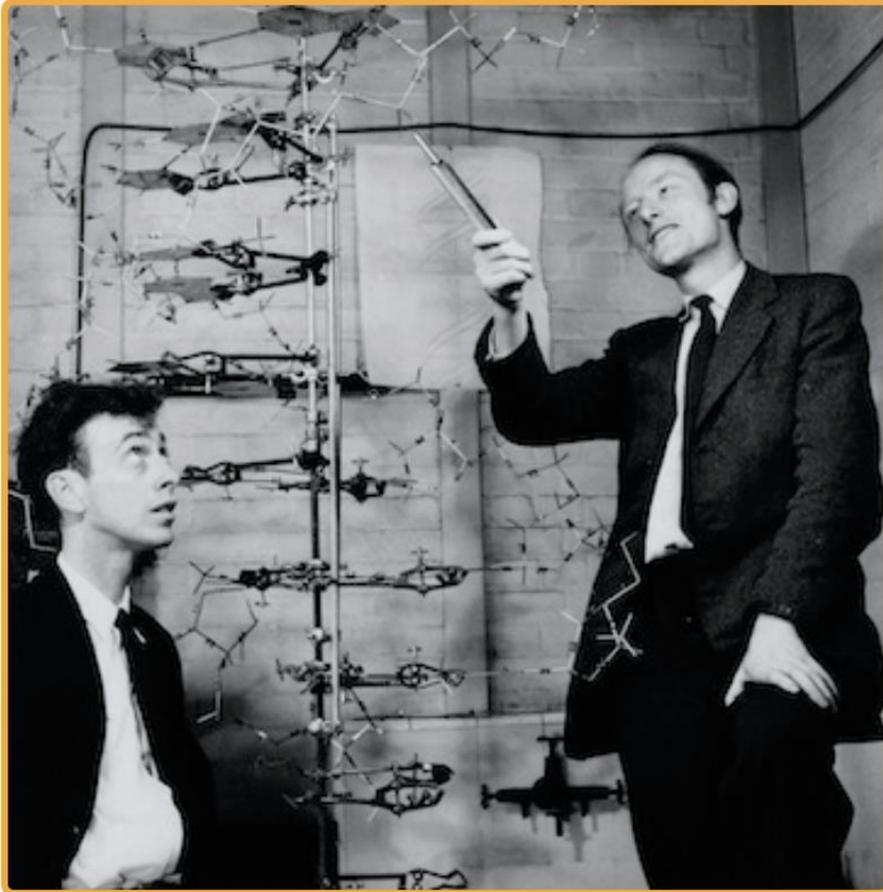
### Gregor Mendel, the father of modern genetics, presents his research on experiments in plant hybridization ✕

Gregor Mendel, a 19th century Augustinian monk, is called the father of modern genetics. He used a monastery garden for crossing pea plant varieties having different heights, colors, pod shapes, seed shapes, and flower positions. Mendel's experiments, between 1856 and 1863, revealed how traits are passed down from parents. For example, when he crossed yellow peas with green peas, all the offspring peas were yellow. But when these offspring reproduced, the next generation was  $\frac{3}{4}$  yellow and  $\frac{1}{4}$  green. Mendel's work, which was presented in 1865, showed that what we now call "genes" determine traits in predictable ways.

---

1865

## Historical overview



### James Watson and Francis Crick discover the double helix structure of DNA

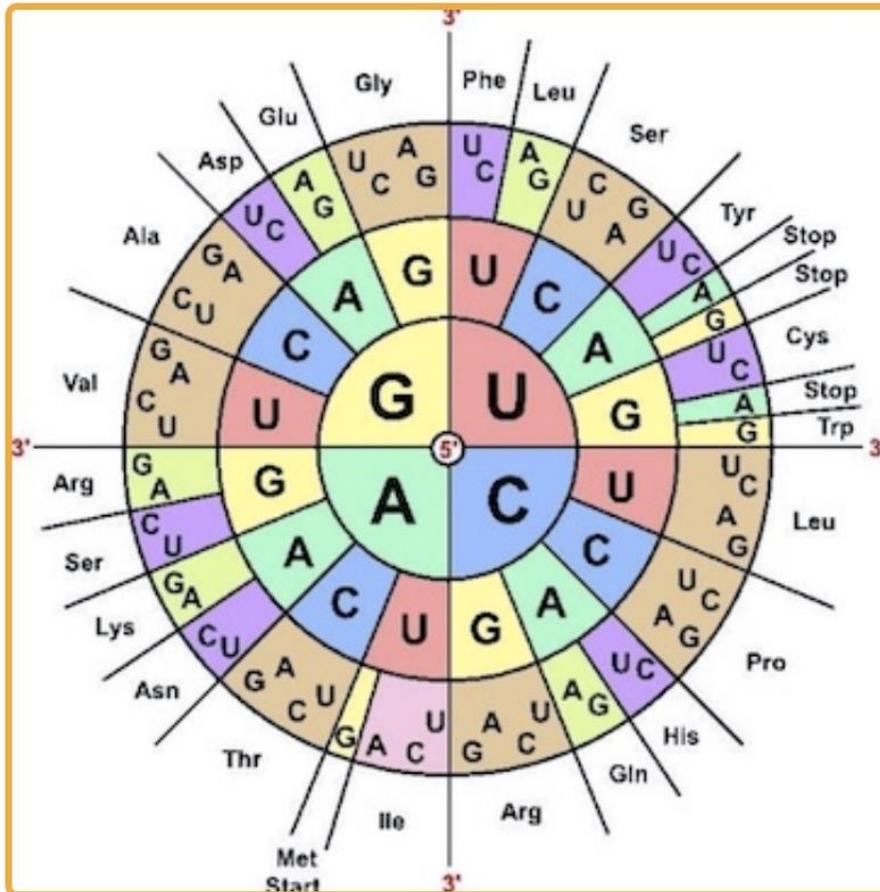


When Francis Crick and James Watson modeled the structure of DNA, they used paper cutouts of the bases (A, C, G, T) and metal scraps from a machine shop. Their model represented DNA as a double helix, with sugars and phosphates forming the outer strands of the helix and the bases pointing into the center. Hydrogen bonds connect the bases, pairing A–T and C–G; and the two strands of the helix are parallel but oriented in opposite directions. Their 1953 paper notes that the model “immediately suggests a possible copying mechanism for the genetic material.”

---

1953

## Historical overview

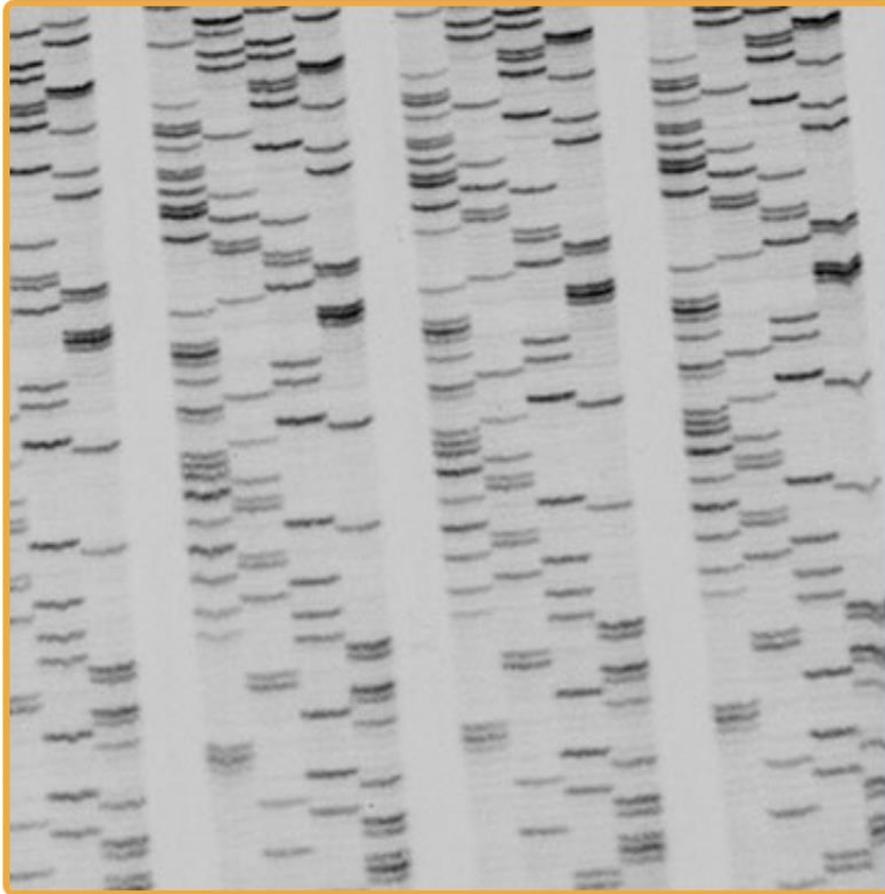


## Marshall Nirenberg cracks the genetic code for protein synthesis

In the early 1960s, Marshall Nirenberg and National Institutes of Health colleagues focused on how DNA directs protein synthesis and the role of RNA in these processes. Their 1961 experiment, using a synthetic messenger RNA (mRNA) strand that contained only uracils (U), yielded a protein that contained only phenylalanines. Identifying UUU (three uracil bases in a row) as the RNA code for phenylalanine was their first breakthrough. Within a few years, Nirenberg's team had cracked the 60 mRNA codons for all 20 amino acids. In 1968, Nirenberg shared the Nobel Prize in Physiology or Medicine for his contributions to breaking the genetic code and understanding protein synthesis.

1961

## Historical overview



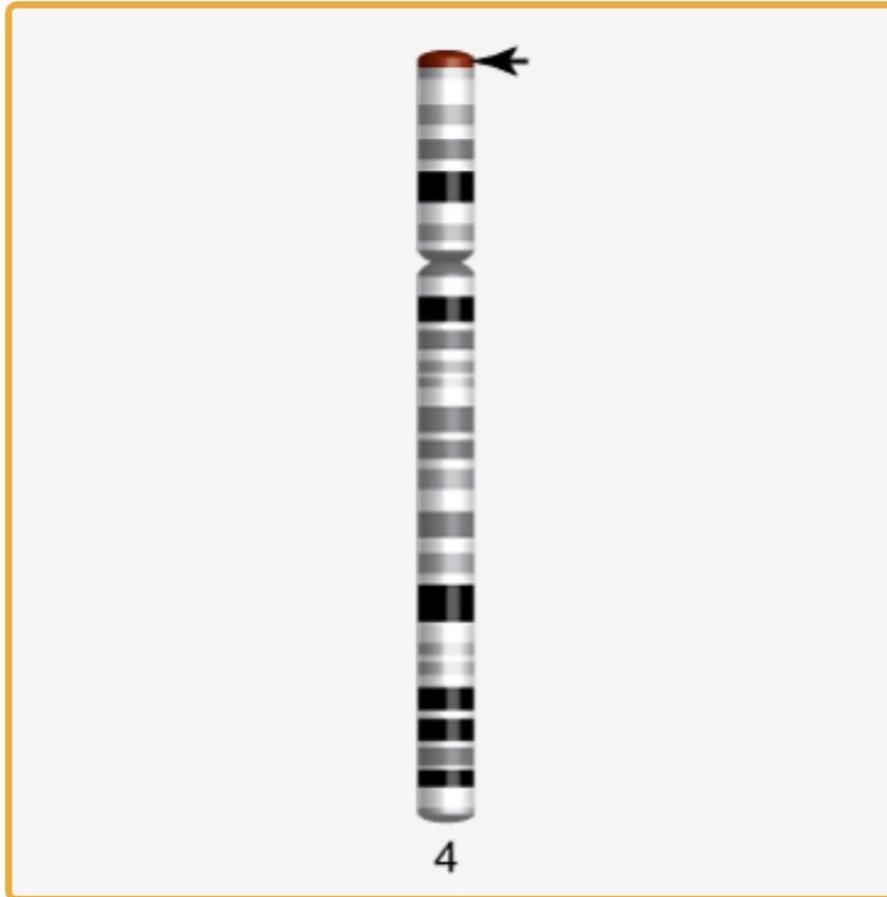
### Frederick Sanger develops rapid DNA sequencing technique

In 1977, Frederick Sanger developed the classical “rapid DNA sequencing” technique, now known as the Sanger method, to determine the order of bases in a strand of DNA. Special enzymes are used to synthesize short pieces of DNA, which end when a selected “terminating” base is added to the stretch of DNA being synthesized. Typically, each of these terminating bases is tagged with a radioactive marker, so it can be identified. Then the DNA fragments, of varying lengths, are separated by how rapidly they move through a gel matrix when an electric field is applied – a technique called electrophoresis. Frederick Sanger shared the 1980 Nobel Prize in Chemistry for his contributions to DNA-sequencing methods.

---

1977

## Historical overview



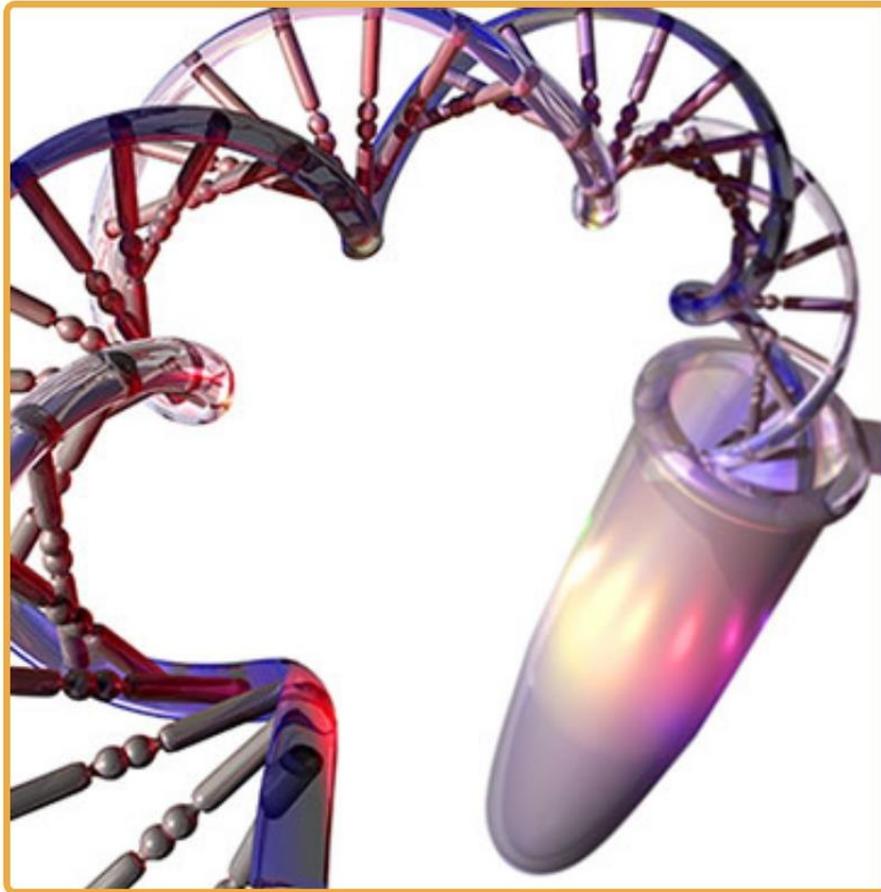
### First genetic disease mapped, Huntington's Disease

Huntington's disease (HD) causes the death of specific neurons in the brain, leading to jerky movements, physical rigidity, and dementia. Symptoms usually appear in midlife and worsen progressively. The location of the HD gene, whose mutation causes Huntington's disease, was mapped to chromosome 4 in 1983, making HD the first disease gene to be mapped using DNA polymorphisms – variants in the DNA sequence. The mutation consists of increasing repetitions of "CAG" in the DNA that codes for the protein huntingtin. The number of CAG repeats may increase when passed from parent to child, leading to earlier HD onset in each generation. The gene was finally isolated in 1993.

---

1983

## Historical overview



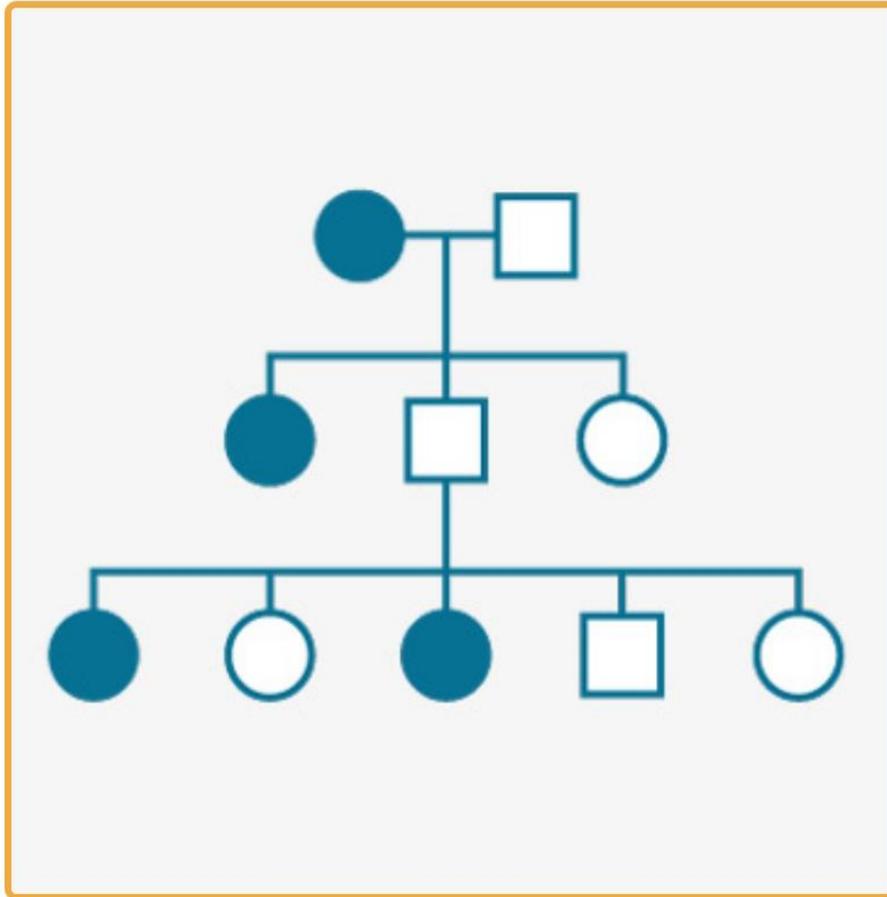
### Invention of polymerase chain reaction (PCR) technology for amplifying DNA

Conceived in 1983 by Kary Mullis, the Polymerase Chain Reaction (PCR) is a relatively simple and inexpensive technology used to amplify or make billions of copies of a segment of DNA. One of the most important scientific advances in molecular biology, PCR amplification is used every day to diagnose diseases, identify bacteria and viruses, and match criminals to crime scenes. PCR revolutionized the study of DNA to such an extent that Dr. Mullis was awarded the Nobel Prize in Chemistry in 1993.

---

1983

## Historical overview



### First evidence provided for the existence of the BRCA1 gene

BRCA1 (BReast CAncer gene 1) is a “tumor suppressor gene,” which normally produces a protein that prevents cells from growing and dividing out of control. However, certain variations of BRCA1 can disrupt its normal function, leading to increased hereditary risk for cancer. The first evidence for existence of the BRCA1 gene was provided in 1990 by the King laboratory at University of California Berkeley. After a heated international race, the gene was finally isolated in 1994. Today, researchers have identified more than 1,000 mutations of the BRCA1 gene, many of them associated with increased risk of cancer, particularly breast and ovarian cancers in women.

1990

## Historical overview



### The Human Genome Project ✕ begins

Beginning in 1984, the U.S. Department of Energy (DOE), National Institutes of Health (NIH), and international groups held meetings about studying the human genome. In 1988, the National Research Council recommended starting a program to map the human genome. Finally, in 1990, NIH and DOE published a plan for the first five years of an expected 15-year project. The project would develop technology for analyzing DNA; map and sequence human and other genomes – including fruit flies and mice; and study related ethical, legal, and social issues.

---

1990

## Historical overview

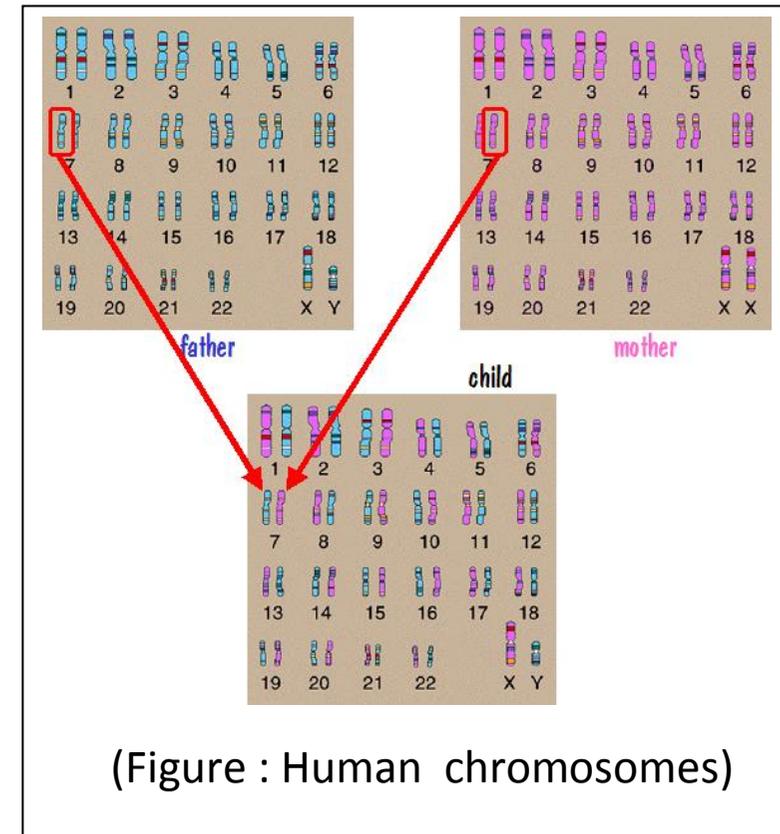


- In June 2000 came the announcement that the majority of the human genome had in fact been sequenced, which was followed by the publication of **90 percent of the sequence of the genome's three billion base-pairs** in the journal *Nature*, in February 2001
- Surprises accompanying the sequence publication included:
  - the relatively small **number of human genes**, perhaps as few as **30,000-35,000**;
  - the complex architecture of human proteins compared to their homologs - similar genes with the same functions - in, for example, roundworms and fruit flies;
  - the lessons to be taught by repeat sequences of DNA.

## Terminology: Genes

- The **gene** is the basic physical unit of inheritance.
- Genes are passed from parents to offspring and contain the information needed to specify traits.
- They are arranged, one after another, on structures called chromosomes.
- A chromosome contains a single, long DNA molecule, only a portion

of which corresponds to a single gene.



## Terminology: Gene Annotation

- An annotation (irrespective of the context) is a note added by way of explanation or commentary.
- **Genome annotation** is the process of identifying the locations of genes and all of the coding regions in a genome and determining what those genes do.
- Once a genome is sequenced, it needs to be annotated to make sense of it  
  
→ links to giving an “interpretation”

## Terminology: Alleles

- Allele: one of several alternative forms of DNA sequence at specific chromosomal location
- Polymorphism: often used to indicate the existence of at least 2 alleles at a single “locus”
- Homozygosity (homozygous): both alleles identical at locus
- Heterozygosity (heterozygous): different alleles at locus
- Genetic marker: polymorphic DNA sequence at single locus

[Mutations ~polymorphisms (see later)]

# Historical overview


National Human Genome Research Institute  
National Institutes of Health

Research Funding
Research at NHGRI
Health
Education
Issues in Genetics
Newsroom
Careers & Training
About
For You

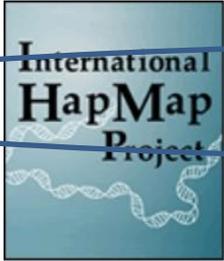



Home > [Education](#) > [Understanding the Human Genome Project](#) > [Dynamic Timeline](#) > [2004-The Future](#) > **2005b: HapMap Project Completed**

Online Education Kit: 2004-The Future

- 2004a: Rat and Chicken Genomes Sequenced
- 2004b: FDA Approves First Microarray
- 2004c: Refined Analysis of Complete Human Genome Sequence
- 2004d: Surgeon General Stresses Importance of Family History
- 2005a: Chimpanzee Genomes Sequenced
- 2005b: HapMap Project Completed**
- 2005c: Trypanosomatid Genomes Sequenced
- 2005d: Dog Genomes Sequenced
- 2006a: The Cancer Genome Atlas (TCGA) Project Started
- 2006b: Second Non-human Primate Genome is Sequenced
- 2006c: Initiatives to Establish the Genetic and Environmental Causes of Common Diseases Launched
- The Future

## 2005: HapMap Project Completed



The International HapMap Consortium published a catalog of human genetic variation that is expected to help speed the identification of genes associated with common diseases such as asthma, cancer, diabetes, and heart disease. While the Human Genome Project focused on the DNA sequence from a single individual, the HapMap project focused on variation in the genome and on human populations. The \$138 million project was a three-year collaboration between more than 200 researchers from Canada, China, Japan, Nigeria and the United States. The new paper described the completion of a Phase I HapMap that contains more than 1 million markers of genetic variation. At the time of the publication, the consortium was nearing completion of a Phase II HapMap that would contain more than 3 million genetic markers.

 Share
  Print

**See Also:**

[2005 Release: International Consortium Completes Map](#)

[International HapMap Project](#)

**On Other Sites:**

[International HapMap Project](#) Web page for the International HapMap Consortium

### More Information

**References:**

The International HapMap Consortium. Integrating ethics and science in the International HapMap Project. *Nature Genetics*, 5: 467-475. 2004. [[Full Text](#)] 

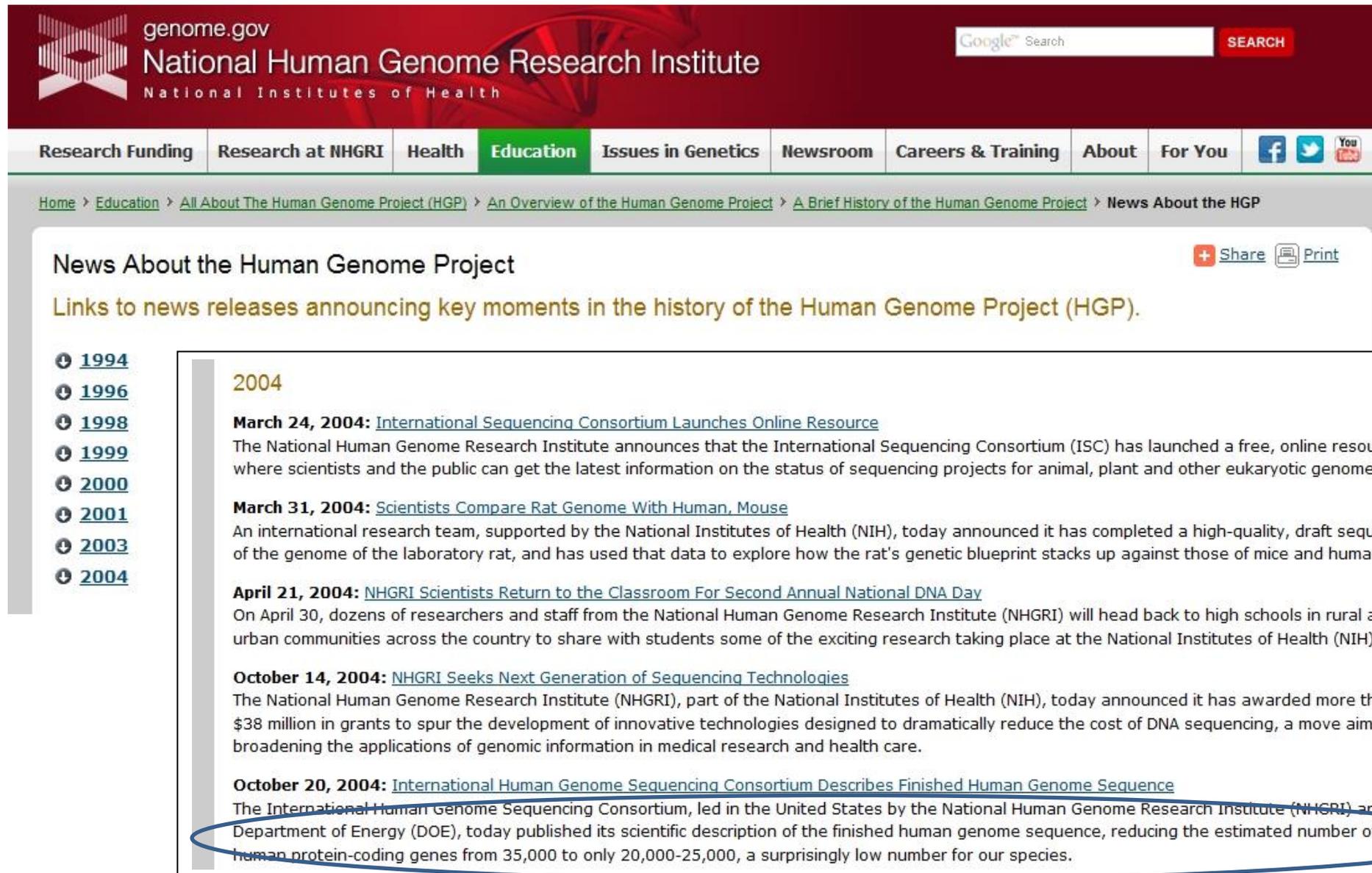
International HapMap Consortium. A haplotype map of the human genome. *Nature*, 437: 1229-1320. 2005. [[Full Text](#)] 

Klein, R.J., Zeiss, C., Chew, E.Y., Tsai, J.Y., et al. Complement factor H polymorphism in age-related macular degeneration. *Science*, 308: 385-389. 2005. [[PubMed](#)]

To view the PDFs on this page, you will need Adobe Reader.



# News about the HGP



The screenshot shows the homepage of the National Human Genome Research Institute (NHGRI) website. The header includes the NHGRI logo, the text "genome.gov National Human Genome Research Institute National Institutes of Health", and a Google search bar. A navigation menu is visible with categories like "Research Funding", "Research at NHGRI", "Health", "Education" (highlighted), "Issues in Genetics", "Newsroom", "Careers & Training", "About", and "For You". Social media icons for Facebook, Twitter, and YouTube are also present.

The main content area is titled "News About the Human Genome Project" and includes a "Share" button and a "Print" button. Below the title is a sub-header: "Links to news releases announcing key moments in the history of the Human Genome Project (HGP)."

A vertical list of years is provided on the left side of the page, with "2004" selected. The main content area displays news releases for the year 2004:

- 2004**
- March 24, 2004:** [International Sequencing Consortium Launches Online Resource](#)  
The National Human Genome Research Institute announces that the International Sequencing Consortium (ISC) has launched a free, online resource where scientists and the public can get the latest information on the status of sequencing projects for animal, plant and other eukaryotic genomes.
- March 31, 2004:** [Scientists Compare Rat Genome With Human, Mouse](#)  
An international research team, supported by the National Institutes of Health (NIH), today announced it has completed a high-quality, draft sequence of the genome of the laboratory rat, and has used that data to explore how the rat's genetic blueprint stacks up against those of mice and humans.
- April 21, 2004:** [NHGRI Scientists Return to the Classroom For Second Annual National DNA Day](#)  
On April 30, dozens of researchers and staff from the National Human Genome Research Institute (NHGRI) will head back to high schools in rural and urban communities across the country to share with students some of the exciting research taking place at the National Institutes of Health (NIH).
- October 14, 2004:** [NHGRI Seeks Next Generation of Sequencing Technologies](#)  
The National Human Genome Research Institute (NHGRI), part of the National Institutes of Health (NIH), today announced it has awarded more than \$38 million in grants to spur the development of innovative technologies designed to dramatically reduce the cost of DNA sequencing, a move aimed at broadening the applications of genomic information in medical research and health care.
- October 20, 2004:** [International Human Genome Sequencing Consortium Describes Finished Human Genome Sequence](#)  
The International Human Genome Sequencing Consortium, led in the United States by the National Human Genome Research Institute (NHGRI) and the Department of Energy (DOE), today published its scientific description of the finished human genome sequence, reducing the estimated number of human protein-coding genes from 35,000 to only 20,000-25,000, a surprisingly low number for our species.

## Historical overview

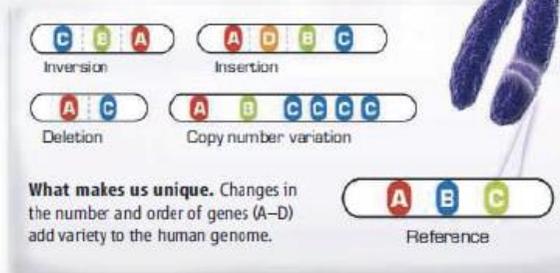
**BREAKTHROUGH OF THE YEAR**

# Human Genetic Variation

Equipped with faster, cheaper technologies for sequencing DNA and assessing variation in genomes on scales ranging from one to millions of bases, researchers are finding out how truly different we are from one another

THE UNVEILING OF THE HUMAN GENOME ALMOST 7 YEARS AGO cast the first faint light on our complete genetic makeup. Since then, each new genome sequenced and each new individual studied has illuminated our genomic landscape in ever more detail. In 2007, researchers came to appreciate the extent to which our genomes differ from person to person and the implications of this variation for deciphering the genetics of complex diseases and personal traits.

Less than a year ago, the big news was triangulating variation between us and our primate cousins to get a better handle on genetic changes along the evolutionary tree that led to humans. Now, we have moved from asking what in our DNA makes us human to striving to know what in my DNA makes me me.



Pennisi 2007 Science 318:1842-3

## Historical overview: associating genetic variation to disease outcomes



### BREAKTHROUGH OF THE YEAR: The Runners-Up

*Science* 314, 1850a (2006);  
DOI: 10.1126/science.314.5807.1850a

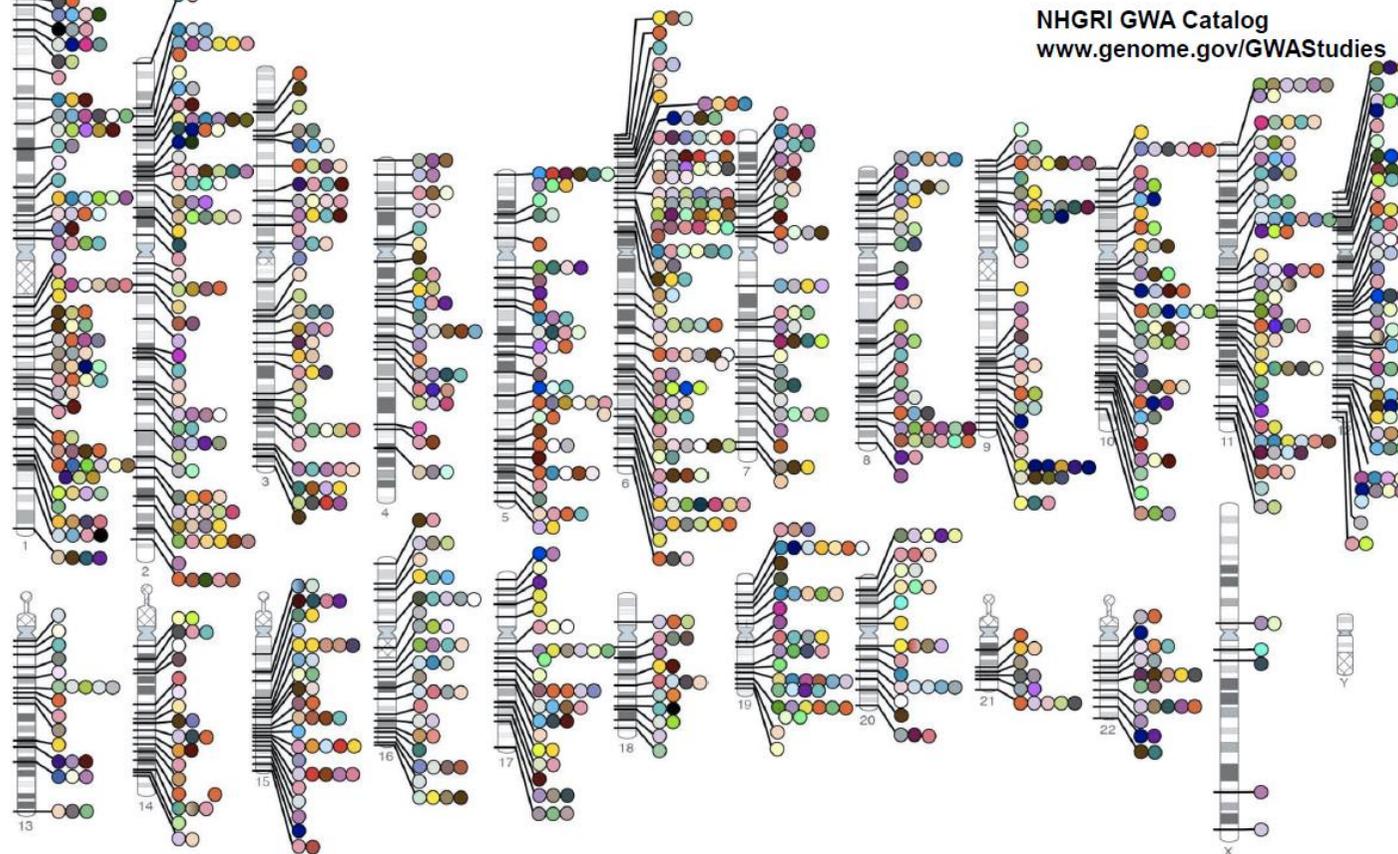
#### Areas to Watch in 2007

**Whole-genome association studies.** The trickle of studies comparing the genomes of healthy people to those of the sick is fast becoming a flood. Already, scientists have applied this strategy to macular degeneration, memory, and inflammatory bowel disease, and new projects on schizophrenia, psoriasis, diabetes, and more are heating up. But will the wave of data and new gene possibilities offer real insight into how diseases germinate? And will the genetic associations hold up better than those found the old-fashioned way?

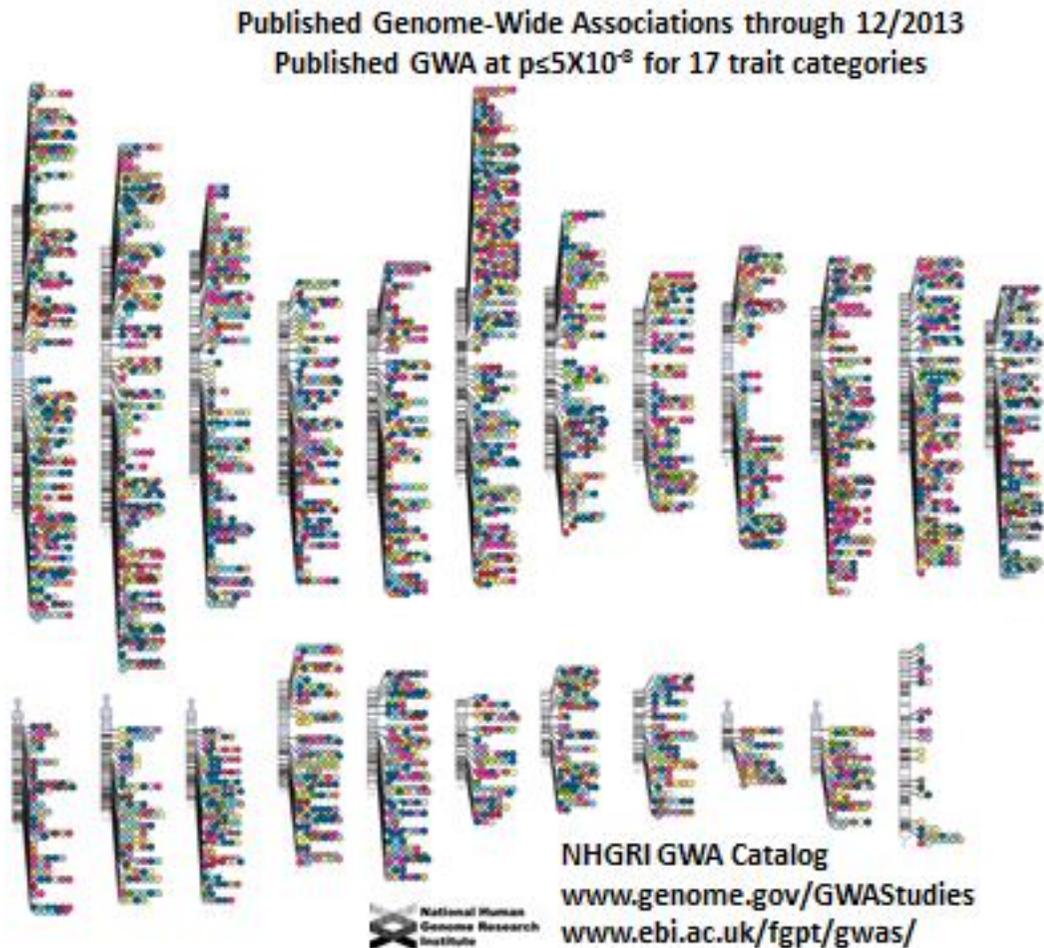


## Historical overview: 210 traits – multiple loci (sites, locations)

Published Genome-Wide Associations through 12/2010,  
1212 published GWA at  $p < 5 \times 10^{-8}$  for 210 traits



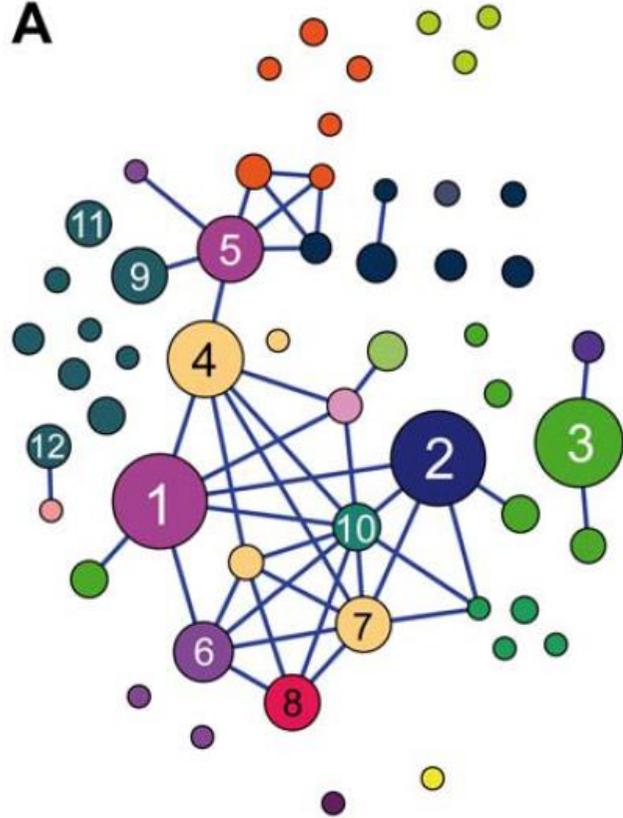
# Historical overview: trait categories



- 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

# Historical overview: inter-relationships (networks)

**A**

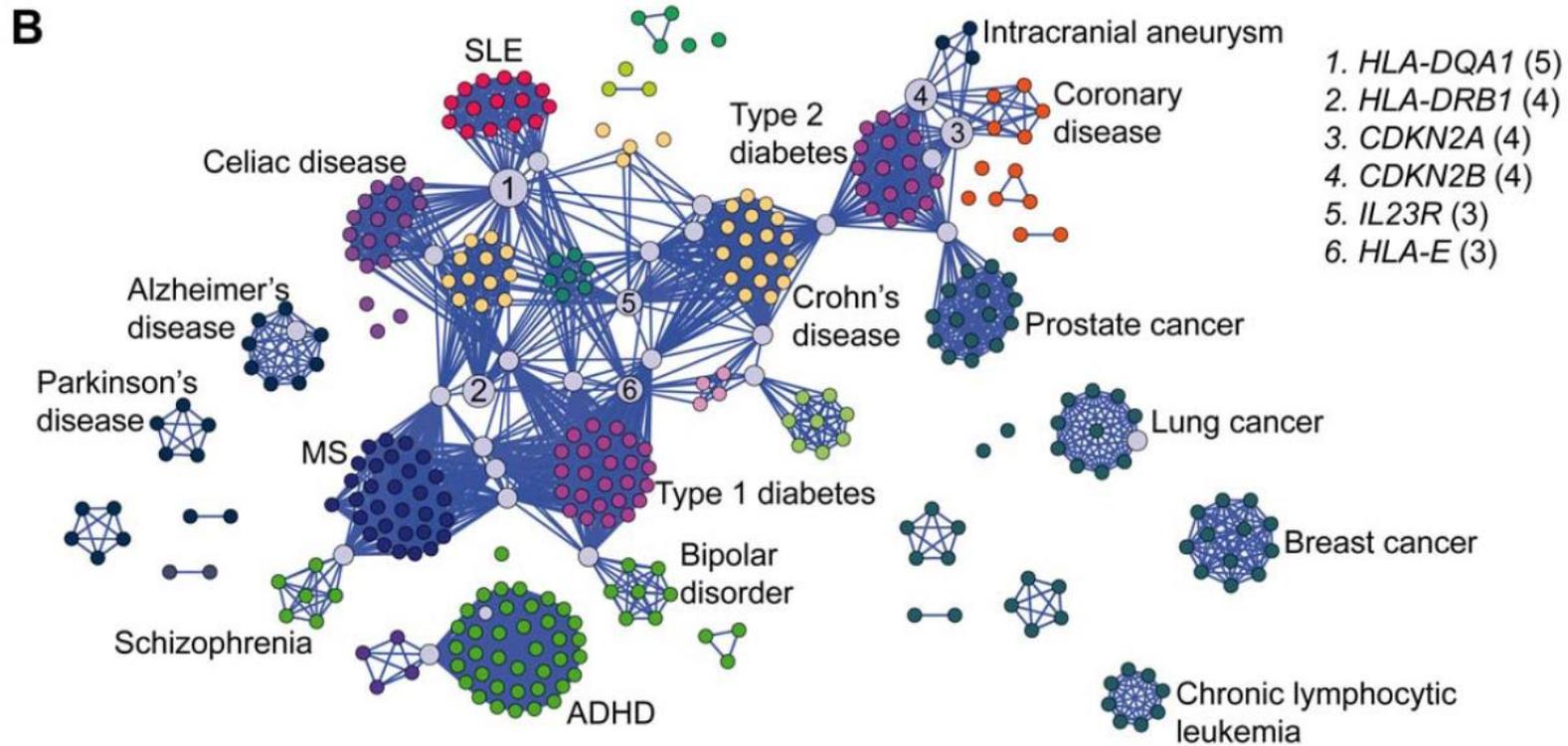


1. Type 1 diabetes (36)
2. Multiple sclerosis (36)
3. ADHD and conduct disorder (33)
4. Crohn's disease (27)
5. Type 2 diabetes (22)
6. Celiac disease (19)
7. Ulcerative colitis(17)
8. Systemic lupus erythematosus (17)
9. Prostate cancer (17)
10. Rheumatoid arthritis (13)
11. Breast cancer (12)
12. Lung cancer (11)

- Cardiovascular diseases (Cv)
- Digestive system diseases
- Endocrine system diseases
- Eye diseases
- Immune system diseases (Is)
- Mental disorders
- Multiple diseases
- Musculoskeletal diseases (Ms)
- Ms, Sc, Is
- Neoplasms
- Nervous system diseases (Ns)
- Ns, Cv
- Ns, Is
- Ns, Ms
- Nutritional and metabolic diseases (Nm)
- Nm, Es, Is
- Skin and connective tissue diseases (Sc)
- Sc, Is
- Urogenital diseases

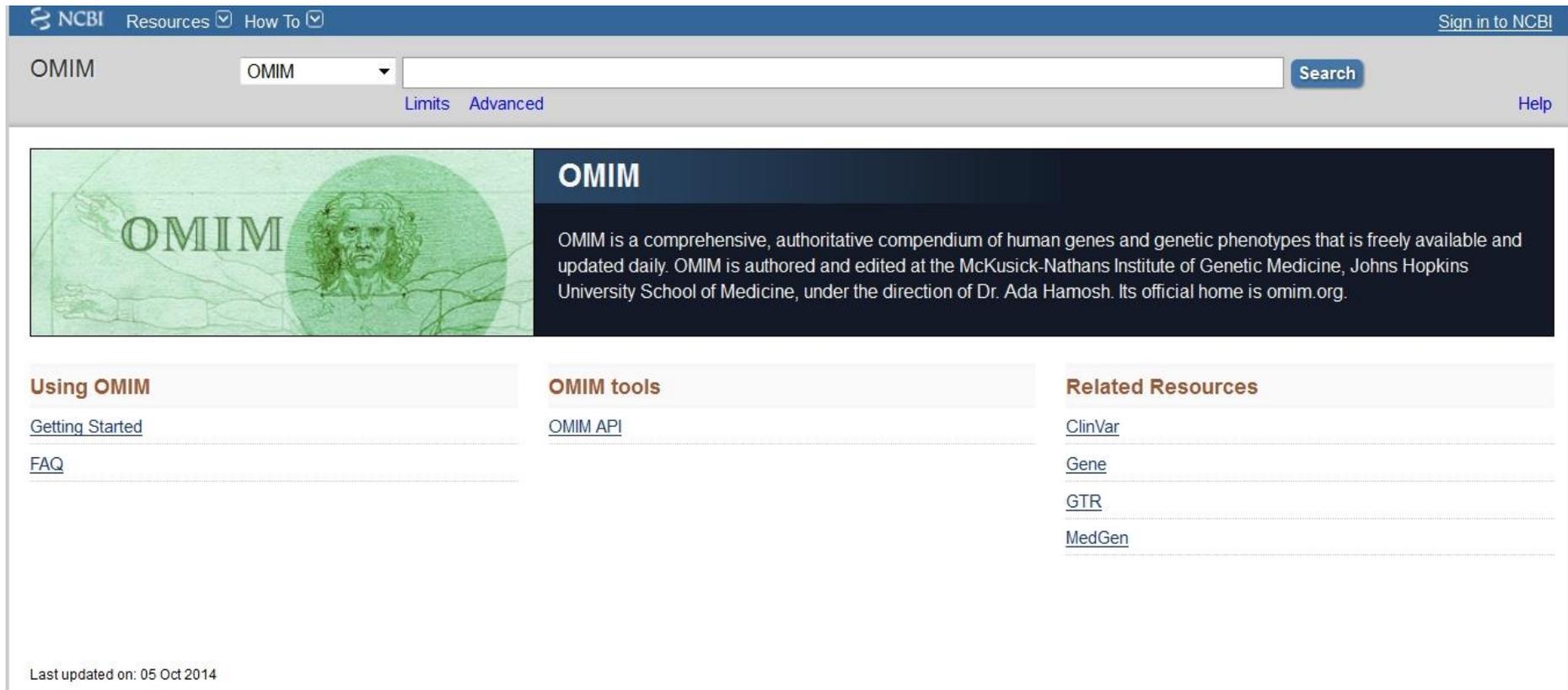
(Barrenas et al 2009: complex disease network – nodes are diseases)

## Historical overview: inter-relationships (networks)



(Barrenas et al 2009: complex disease GENE network – nodes are genes)

## Historical overview: monitoring the progress



The screenshot shows the OMIM website interface. At the top, there is a navigation bar with the NCBI logo, "Resources" and "How To" dropdown menus, and a "Sign in to NCBI" link. Below this is a search bar with "OMIM" selected in a dropdown menu, a search input field, and a "Search" button. There are also links for "Limits" and "Advanced" search options, and a "Help" link.

The main content area features a large banner with the OMIM logo and a portrait of a man. To the right of the banner, the text reads: "OMIM is a comprehensive, authoritative compendium of human genes and genetic phenotypes that is freely available and updated daily. OMIM is authored and edited at the McKusick-Nathans Institute of Genetic Medicine, Johns Hopkins University School of Medicine, under the direction of Dr. Ada Hamosh. Its official home is [omim.org](http://omim.org)."

Below the banner, there are three columns of links:

- Using OMIM**
  - [Getting Started](#)
  - [FAQ](#)
- OMIM tools**
  - [OMIM API](#)
- Related Resources**
  - [ClinVar](#)
  - [Gene](#)
  - [GTR](#)
  - [MedGen](#)

At the bottom left, it says "Last updated on: 05 Oct 2014".

## **OMIM:** molecular dissection of human disease

- Online Mendelian Inheritance in Man (OMIM<sup>®</sup>) is a continuously updated **catalog of human genes and genetic disorders and traits** (i.e. coded phenotypes, where phenotype is any characteristic of the organism), with particular focus on the molecular relationship between genetic variation and phenotypic expression.
- It can be considered to be a phenotypic companion to the Human Genome Project. OMIM is a continuation of Dr. Victor A. McKusick's Mendelian Inheritance in Man, which was published through 12 editions, the last in 1998.
- OMIM is currently biocurated at the McKusick-Nathans Institute of Genetic Medicine, The Johns Hopkins University School of Medicine.
- Frequently asked questions: <http://www.omim.org/help/faq>

# Accessing OMIM

The screenshot shows the NCBI website interface. At the top, there is a navigation bar with 'NCBI Resources' and 'How To' dropdown menus, and a 'Sign in to NCBI' link. Below this is a search bar with the text 'OMIM' entered and a 'Search' button. On the left side, there is a vertical menu with various categories, and 'OMIM' is highlighted in blue. A dropdown menu is open from the 'OMIM' category, listing various resources: dbGaP, dbVar, Epigenomics, EST, Gene, Genome, GEO DataSets, GEO Profiles, GSS, HomoloGene, MedGen, MeSH, NCBI Web Site, NLM Catalog, Nucleotide, OMIM (highlighted), PMC, PopSet, Probe, and Protein. Below the main menu, there is a 'NCBI Facebook page' banner with a 'GO' button and a progress bar showing 4 out of 8 items. On the right side, there are sections for 'Popular Resources' (PubMed, Bookshelf, PubMed Central, PubMed Health, BLAST, Nucleotide, Genome, SNP, Gene, Protein, PubChem) and 'NCBI Announcements' (dbVar now accepts VCF submissions of structural variation data, dbVar, NCBI's database of genomic, New NCBI Insights blog post: NCBI's medical genetics resources, The latest blog post on NCBI Insights, NCBI webinar on E-Utilities October 15th, On October 15th, NCBI will have a webinar entitled "An Introduction to NCBI's").

# Historical overview: exome sequencing, full genome sequencing

NCBI Resources How To Sign in to NCBI

1000 Genomes Browser Homo sapiens: GRCh37.p13 (GCF\_000001405.25) Chr 1 (NC\_000001.10): 1 - 249.3M

Reset All Share this page FAQ Help Version 3.3

**Ideogram View**

**Exon Navigator**

There are too many genes in the region (3570). Please narrow the region to enable exon navigation.

NC\_000001.10: 1..249M (249Mbp)

Segmental Duplications on GRCh37

1000 Genomes Phase 1 Strict Accessibility Mask

Genes, NCBI Homo sapiens Annotation Release 105

ClinVar Short Variations based on dbSNP 141 (Homo sapiens Annotation Release 105)

dbSNP 141 (Homo sapiens Annotation Release 105) HapMap Recombination Rate

dbSNP 141 (Homo sapiens Annotation Release 105) all data

Data not in 1000 Genomes Phase 1, dbSNP 141 (Homo sapiens Annotation Release 105)

Zoom to see data!

Download data for this region

Hide populations with unchecked samples

**Genotypes**

Go to Selection	Scroll Region	114,057,996	114,058,146	114,058,215	114,058,280	114,058,349	114,058,466	114,058,537	114,058,563	114,058,640	114,058,739	114,058,746	114,058,853	114,059,189	114,059,335	114,059,373	114,059,380	114,059,727	114,059,954	114,059,998	114,060,051	114,060,051
rs6685909	rs114841421	rs183004883	rs4838988	rs4839323	rs187330676	rs786891846	rs113381680	rs191384074	rs114985670	rs141798006	rs145542737	rs147746490	rs141083549	rs150301358	rs183913552	rs79656666	rs141847716	rs115616102	rs149798012	rs116		
G=0.2039	C=0.9927	A=0.9995	C=0.7332	T=0.6809	G=0.9995	C=0.9986	G=0.9991	G=0.9986	C=0.9968	G=0.9940	G=0.9991	G=0.9995	C=0.9995	G=0.9995	G=0.9995	G=0.9927	TA=0.9853	G=0.9982	GA=0.9862	T=0		
A=0.7961	T=0.0073	T=0.0005	T=0.2668	C=0.3191	T=0.0005	A=0.0014	T=0.0009	A=0.0014	T=0.0032	A=0.0060	A=0.0009	A=0.0005	A=0.0005	A=0.0005	A=0.0005	A=0.0005	A=0.0073	T=0.0147	A=0.0018	G=0.0138	G=0	
G=0.0656	C=0.9754	A=1.0000	C=0.9590	T=0.8279	G=1.0000	C=1.0000	G=1.0000	G=1.0000	C=0.9836	G=1.0000	G=1.0000	G=1.0000	C=1.0000	G=1.0000	G=1.0000	G=1.0000	G=0.9508	TA=0.9590	G=0.9918	GA=0.9426	T=0	
A=0.9344	T=0.0246	T=0.0000	T=0.0410	C=0.1721	T=0.0000	A=0.0000	T=0.0000	A=0.0000	T=0.0164	A=0.0000	A=0.0492	T=0.0410	A=0.0082	G=0.0574	G=0							

## 1.b The concept of a genetic marker

### The evolution of molecular markers (Schlötterer 2004)

OPINION

#### The evolution of molecular markers — just a matter of fashion?

*Christian Schlötterer*

In less than half a century, molecular markers have totally changed our view of nature, and in the process they have evolved themselves. However, all of the molecular methods developed over the years to detect variation do so in one of only three conceptually different classes of marker: protein variants (allozymes), DNA sequence polymorphism and DNA repeat variation. The latest techniques promise to provide cheap, high-throughput methods for genotyping existing markers, but might other traditional approaches offer better value for some applications?

Being able to distinguish between genotypes that are relevant to a trait of interest is a key goal in genetics. Often, this distinction is not based directly on the trait of interest, but on informative marker systems. A genetic marker provides information about allelic variation at a given locus. The first genetic map of *Drosophila melanogaster* was built by Sturtevant using phenotypic markers<sup>1</sup>. How-

continuous improvement in the way in which we assay genetic variation; that is, the latest marker systems are the most informative ones. Nevertheless, in reviewing the history of molecular markers and their pros and cons, I argue that there are only a few conceptually different classes of marker and that recently developed high-throughput methods might not be unconditionally superior to more traditional approaches.

#### **Allozymes**

The first true molecular markers to be established were allozymes (a term that originates from a contraction of the phrase ‘allelic variants of enzymes’). The principle of allozyme markers is that protein variants in enzymes can be distinguished by native gel electrophoresis according to differences in size and charge caused by amino-acid substitutions. To visualize the allozyme bands, the electrophoretic gels are treated with enzyme-specific stains that contain substrate for the enzyme, cofactors and an oxidized salt (for example, nitro-blue tetra-

sample sizes are typically studied in allozyme surveys. Nevertheless, the number of informative marker loci is too small to use allozymes for mapping and ASSOCIATION STUDIES<sup>8</sup>. Furthermore, surveys of natural variation based on allozymes were often challenged by non-neutral evolution of some of the markers used (see, for example, REFS 9–11).

#### **The arrival of DNA-based markers**

One of the criticisms levelled at allozyme markers is that they are an indirect and insensitive method of detecting variation in DNA. A more direct molecular marker would survey DNA variation itself, rather than rely on variations in the electrophoretic mobility of proteins that the DNA encodes. Another important advantage that DNA-based markers have over allozymes is that they allow the number of mutations between different alleles to be quantified. Given these unambiguous advantages, the arrival of DNA manipulation techniques promoted a shift from enzyme-based to DNA-based markers.

“...the arrival of DNA manipulation techniques promoted a shift from enzyme-based to DNA-based markers.”

Marker	Advantages	Disadvantages
SNPs	<ul style="list-style-type: none"> <li>• Low mutation rate</li> <li>• High abundance</li> <li>• Easy to type</li> <li>• New analytical approaches are being developed at present</li> <li>• Cross-study comparisons are easy; data repositories already exist</li> </ul>	<ul style="list-style-type: none"> <li>• Substantial rate heterogeneity among sites</li> <li>• Expensive to isolate</li> <li>• Ascertainment bias</li> <li>• Low information content of a single SNP</li> </ul>
Microsatellites	<ul style="list-style-type: none"> <li>• Highly informative (large number of alleles, high heterozygosity)</li> <li>• Low ascertainment bias</li> <li>• Easy to isolate</li> </ul>	<ul style="list-style-type: none"> <li>• High mutation rate</li> <li>• Complex mutation behaviour</li> <li>• Not abundant enough</li> <li>• Difficult to automate</li> <li>• Cross-study comparisons require special preparation</li> </ul>
Allozymes	<ul style="list-style-type: none"> <li>• Cheap</li> <li>• Universal protocols</li> </ul>	<ul style="list-style-type: none"> <li>• Requirement for fresh or frozen material</li> <li>• Some loci show protein instability</li> <li>• Limited number of available markers</li> <li>• Potentially direct target of selection</li> </ul>
RAPDs and derivatives	<ul style="list-style-type: none"> <li>• Cheap</li> <li>• Produces a large number of bands, which can then be further characterized individually (for example, converted into single locus markers)</li> </ul>	<ul style="list-style-type: none"> <li>• Low reproducibility</li> <li>• Mainly dominant</li> <li>• Difficult to analyse</li> <li>• Difficult to automate</li> <li>• Cross-study comparisons are difficult</li> </ul>
DNA sequencing	<ul style="list-style-type: none"> <li>• Highest level of resolution possible</li> <li>• Not biased</li> <li>• Cross-study comparisons are easy; data repositories</li> </ul>	<ul style="list-style-type: none"> <li>• Still significantly more expensive than the other techniques</li> </ul>

**Be critical**

(date of publication = 2004)

Hence, it is important to keep the historical time lines and achievements in mind

## The evolution of molecular markers

- Nowadays, **genetic markers represent sequences of DNA** which have been traced to specific locations on the chromosomes and associated with particular traits.
- They demonstrate **polymorphism**, which means that the genetic markers in different organisms of the same species are different.
- A classic example of a genetic marker:  
the area of the DNA which codes for blood type in humans – all humans have and need blood, but the blood of individual humans can be very different as a result of polymorphism in the area of the genome which codes for blood.

## Types of genetic markers – example 1: microsatellites

- Synonymous: short tandem repeat, STR
- Number of repeats varies between individuals
  - Mononucleotide, dinucleotide, trinucleotide, tetranucleotide, non-integer STRs
- Determine allele length (e.g., 133, 136, 139, 142, ...)
- Occurrence in non-coding regions
- High mutation frequency  $\approx 10^{-2} - 10^{-4}$  events per locus per generation
- Not easy to score automatically
- Frequent but not dense enough for some applications

(Ziegler and Van Steen, Brazil 2010)

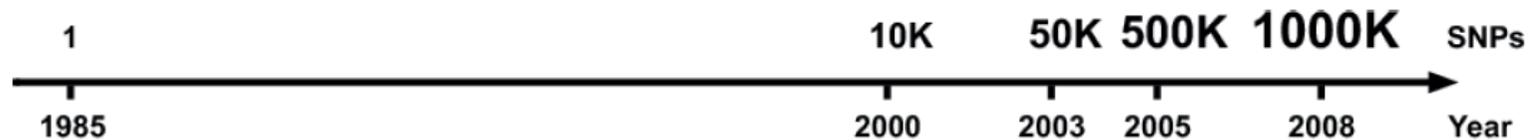
## Types of genetic markers – example 2: single nucleotide polymorphisms

- Variations in single base, i.e., one base substituted by another base
- In theory: four different nucleotides possible at base
- In practice: generally only two different nucleotides observed
- Definition strict and loose:
  - Strict: minor allele frequency  $\geq 1\%$
  - Loose:  $\geq 2$  nucleotides observed in two individuals at position
- Nomenclature:
  - ss-number (submitted SNP number)
  - rs-number: searchable in dbSNP, mapped to external resources, unique
  - rs-numbers do not provide information about possible function of SNP
  - Alternative: nomenclature of Human Genome Variation Society

(Ziegler and Van Steen, Brazil 2010)

## Why are SNPs preferred over STRs?

- SNPs very frequent → dense marker map
- Some SNPs functionally relevant → candidate variations for disease
- SNPs more stable, i.e., lower mutation rate
- Genotyping in highly automated fashion



(Ziegler and Van Steen, Brazil 2010)

## 1.c The rise of Genome-Wide Association studies (GWAs)

### Definition

- A **genome-wide association study** is an approach that involves rapidly scanning markers across the complete sets of DNA, or genomes, of many people to find genetic variations associated with a particular trait.
- **Recall:** a **trait** can be defined as a coded phenotype, a particular characteristic such as hair color, BMI, disease, gene expression intensity level, ...

## Genome-wide association studies for improved public health

- Once new genetic associations are identified, researchers can use the information to develop better strategies to detect, treat and prevent the disease.

The screenshot shows the NIH Office of Extramural Research website. The header includes the U.S. Department of Health & Human Services logo and the text "www.hhs.gov". Below this is the Office of Extramural Research logo and a search bar. A navigation menu includes links for Home, About Grants, Funding, Forms & Deadlines, Grants Policy, News & Events, About OER, and NIH Home. The main content area is titled "Genome-Wide Association Studies (GWAS)" and contains the following text:

The NIH is interested in advancing genome-wide association studies (GWAS) to identify common genetic factors that influence health and disease. For the purposes of this policy, a genome-wide association study is defined as any study of genetic variation across the entire human genome that is designed to identify genetic associations with observable traits (such as blood pressure or weight), or the presence or absence of a disease or condition. Whole genome information, when combined with clinical and other phenotype data, offers the potential for increased understanding of basic biological processes affecting human health, improvement in the prediction of disease and patient care, and ultimately the realization of the promise of personalized medicine. In addition, rapid advances in understanding the patterns of human genetic variation and maturing high-throughput, cost-effective methods for genotyping are providing powerful research tools for identifying genetic variants that contribute to health and disease. The purpose of this Website is to support the implementation of the GWAS Policy.

The NIH will continue to release additional guidance information on this site. Please e-mail [GWAS@mail.nih.gov](mailto:GWAS@mail.nih.gov) with any questions.

**Recent News**

- [NIH Background Fact Sheet on GWAS Policy Update](#) - (08/28/2008) (PDF - 40 KB)
- [NIH Modifications to Genome-Wide Association Studies \(GWAS\) Data Access](#) - (08/28/2008) (PDF - 43 KB)

**Data Access Information**

- [Senior Oversight Committee \(SOC\) Charge and Roster](#) - (07/10/2008) (PDF - 103 KB)
- [Data Access Committees \(DACs\) Charge and Roster](#) - (07/10/2008) (PDF - 50 KB)

The left sidebar contains a list of links under various categories:

- Funding Opportunities**
  - Funding Opportunities (RFAs, PAs) & Notices
  - Unsolicited Applications (Parent Announcements)
  - Research Training & Career Development
  - Small Business (SBIR/STTR)
  - Contract Opportunities
- NIH-Wide Initiatives**
  - Stem Cell Information
  - New and Early Stage Investigators
  - Genome-Wide Association Studies (GWAS)
  - NIH Roadmap for Medical Research
- Global OER Resources**
  - Glossary & Acronyms

# View the GWAs catalogue (<http://www.genome.gov/gwastudies/>)

2317 studies (6/10/2014)

(Entries 1-50 of 2317)

Page 1 of 47 [Next >](#) [Last >>](#)

Date Added to Catalog (since 11/25/08)	First Author/Date/Journal/Study	Disease/Trait	Initial Sample Description	Replication Sample Description	Region	Reported Gene(s)	Mapped Gene(s)	Strongest SNP-Risk Allele	Context	Risk Allele Frequency in Controls	P-value	OR or beta-coefficient and [95% CI]	Platform [SNPs passing QC]	CNV
04/16/14	Chung CM March 03, 2014 <i>Diabetes Metab Res Rev</i> <a href="#">Common quantitative trait locus downstream of RETN gene identified by genome-wide association study is associated with risk of type 2 diabetes mellitus in Han Chinese: a Mendelian randomization effect.</a>	Resistin levels	382 Han Chinese ancestry individuals	559 Han Chinese ancestry individuals	19p13.2	RETN	RETN - C19orf59	rs1423096-G		0.78	$1 \times 10^{-7}$	.322 [0.25-0.40] ug/mL increase	Illumina [NR]	N
10/03/14	Zhang B January 21, 2014 <i>Int J Cancer</i> <a href="#">Genome-wide association study identifies a new SMAD7 risk variant associated with colorectal cancer risk in East Asians.</a>	Colorectal cancer	1,773 East Asian ancestry cases, 2,642 East Asian ancestry controls	6,902 East Asian ancestry cases, 7,862 East Asian ancestry controls	18q21.1	SMAD7	SMAD7	rs7229639-A	intron	0.145	$3 \times 10^{-11}$	1.22 [1.15-1.29]	Affymetrix & Illumina [1,695,815] (imputed)	N
10/06/14	Xie T January 17, 2014 <i>Neurobiol Aging</i> <a href="#">A genome-wide association study combining pathway analysis for typical sporadic</a>	Amyotrophic lateral sclerosis (sporadic)	250 Han Chinese ancestry cases, 250 Han Chinese ancestry controls	NA	<a href="#">View full set of 175 SNPs</a>								Illumina [859,311] (pooled)	N
					NA	RAB9P1	NA	kgp22272527-?		NR	$8 \times 10^{-11}$	NR		
					NA	MYO18B	NA	kgp8087771-?		0.2	$2 \times 10^{-10}$	3.0327 [2.212039-4.157817]		
					12q24.33	GPR133	GPR133	rs11061269-?	intron	0.08	$8 \times 10^{-10}$	3.7761 [2.49-5.74]		
					21q22.3	TMPRSS2	TMPRSS2 -	rs9977018-?		0.05	$2 \times 10^{-9}$	NR		

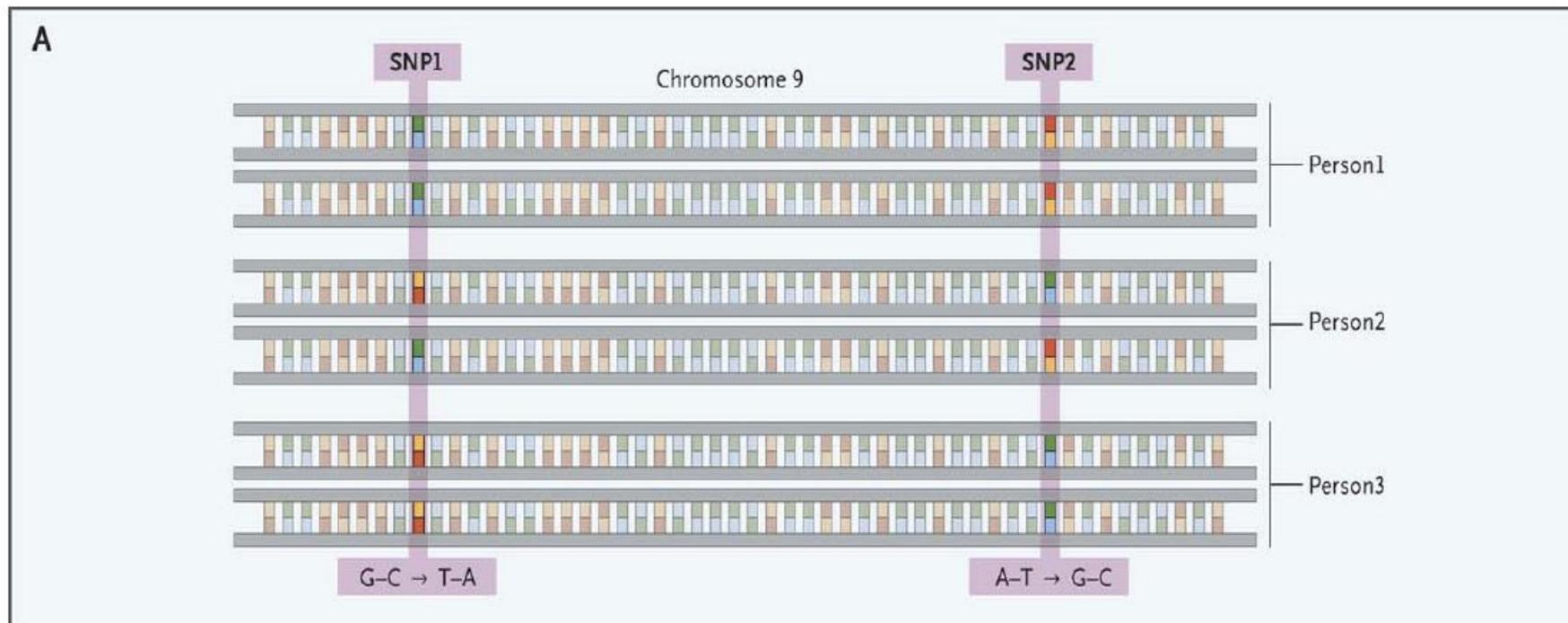
## What do we need to carry out a genome-wide association study?

- The tools include
  - computerized databases that contain the reference human genome sequence,
  - a map of human genetic variation and
  - a set of new technologies that can quickly and accurately analyze (whole-genome) samples for genetic variations that contribute to the onset of a disease.

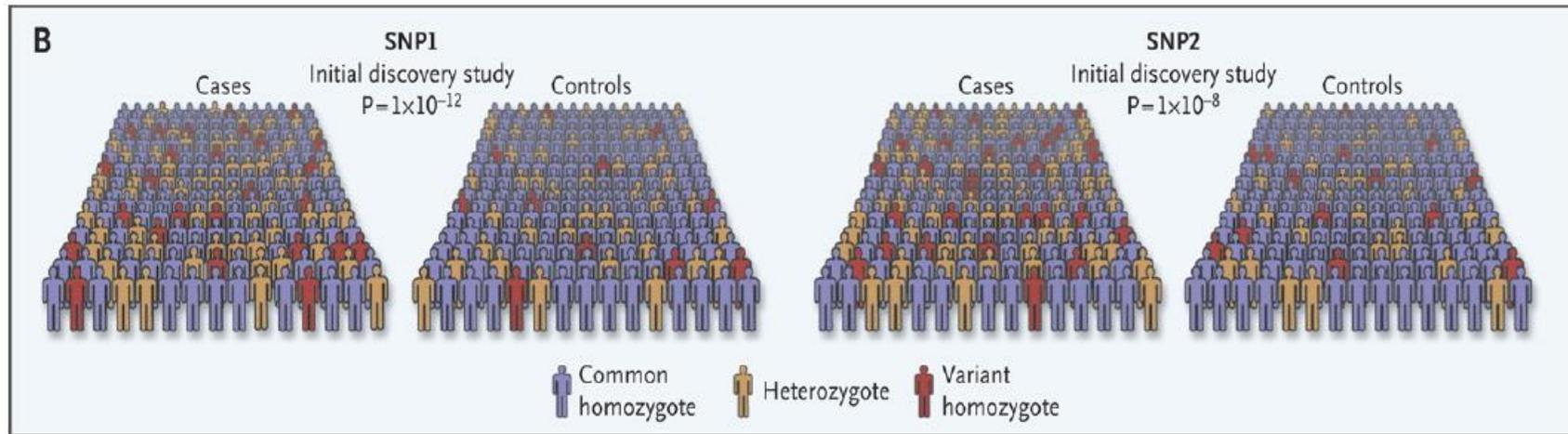
(<http://www.genome.gov/pfv.cfm?pageID=20019523>)

## Genome-wide association studies in practice

The genome-wide association study is typically (but not solely!!!) based on a case-control design in which single-nucleotide polymorphisms (SNPs) across the human genome are genotyped ... (Panel A: small fragment)



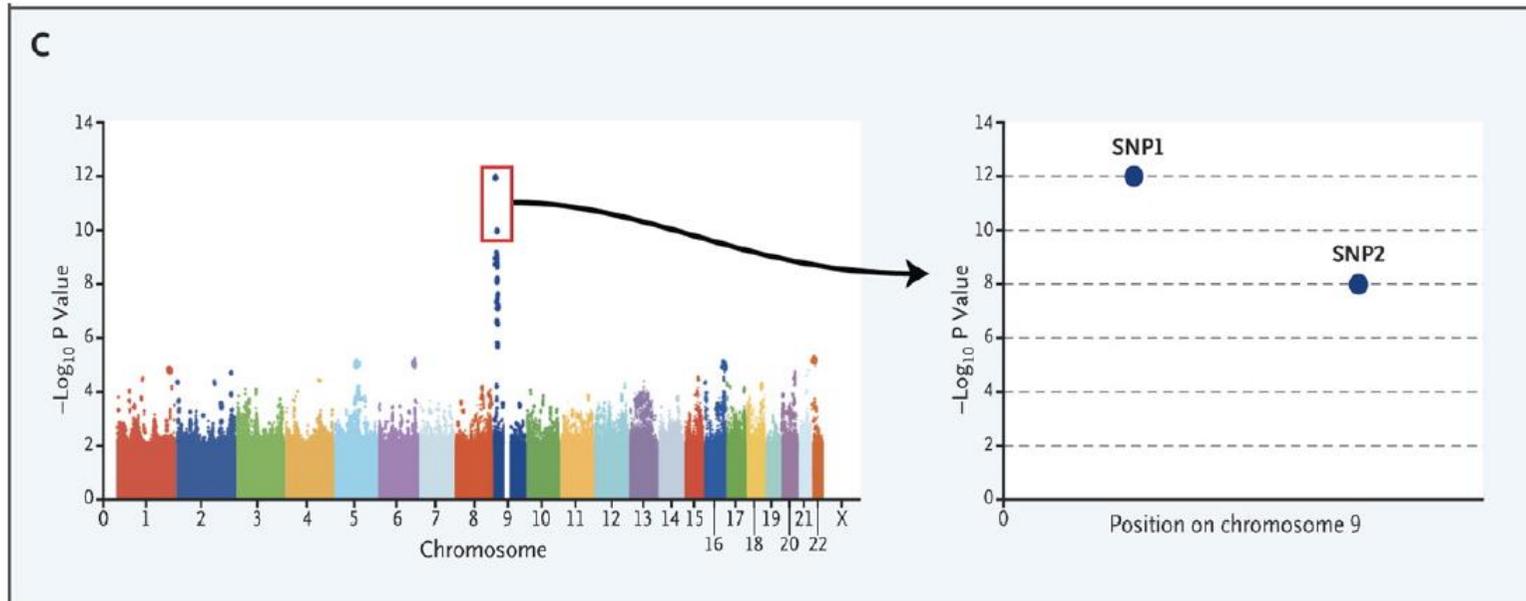
## Genome-wide association studies in practice



- Panel B, the strength of association between each SNP and disease is calculated on the basis of the prevalence of each SNP in cases and controls. In this example, SNPs 1 and 2 on chromosome 9 are associated with disease, with P values of  $10^{-12}$  and  $10^{-8}$ , respectively

(Manolio 2010)

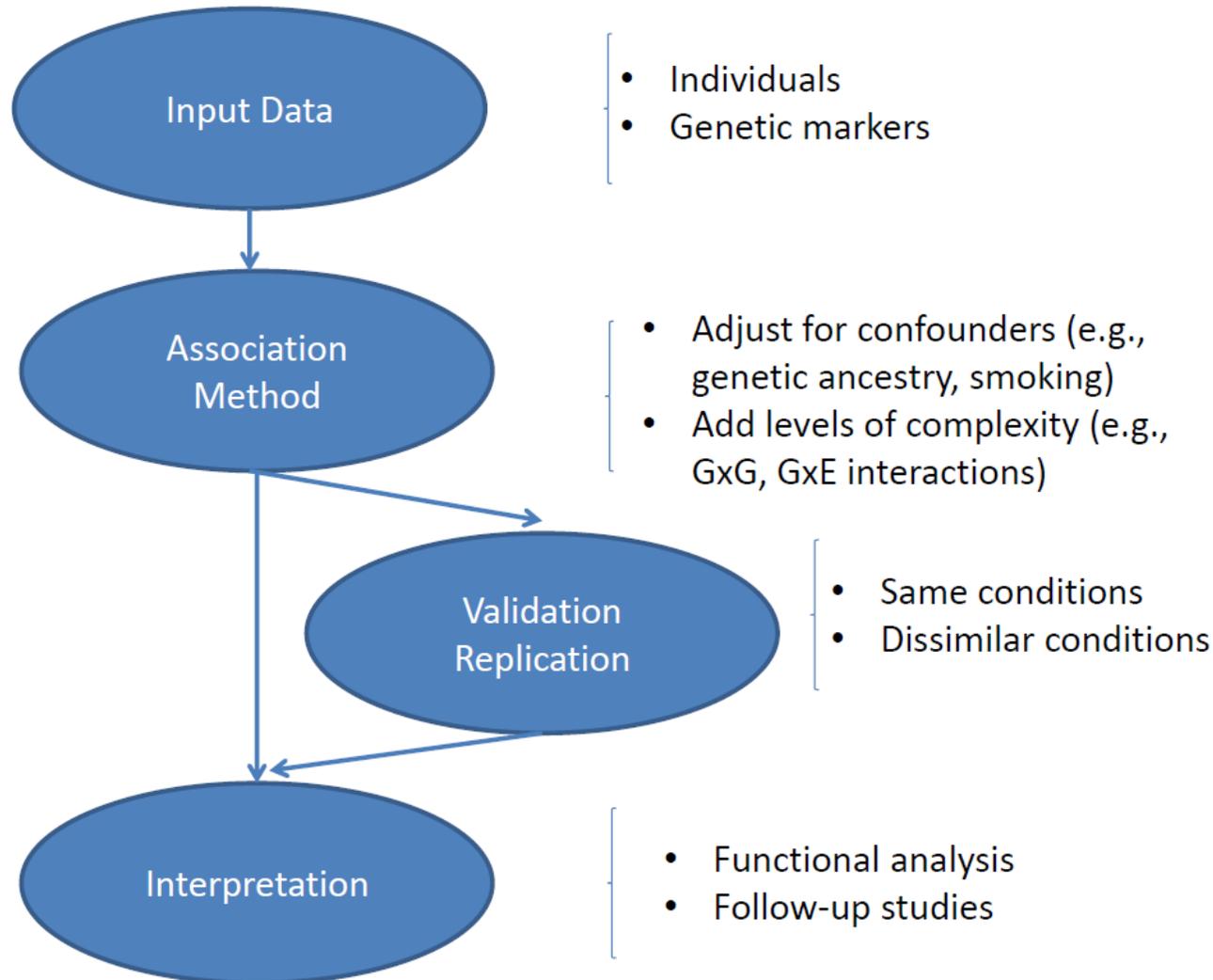
## Genome-wide association studies in practice



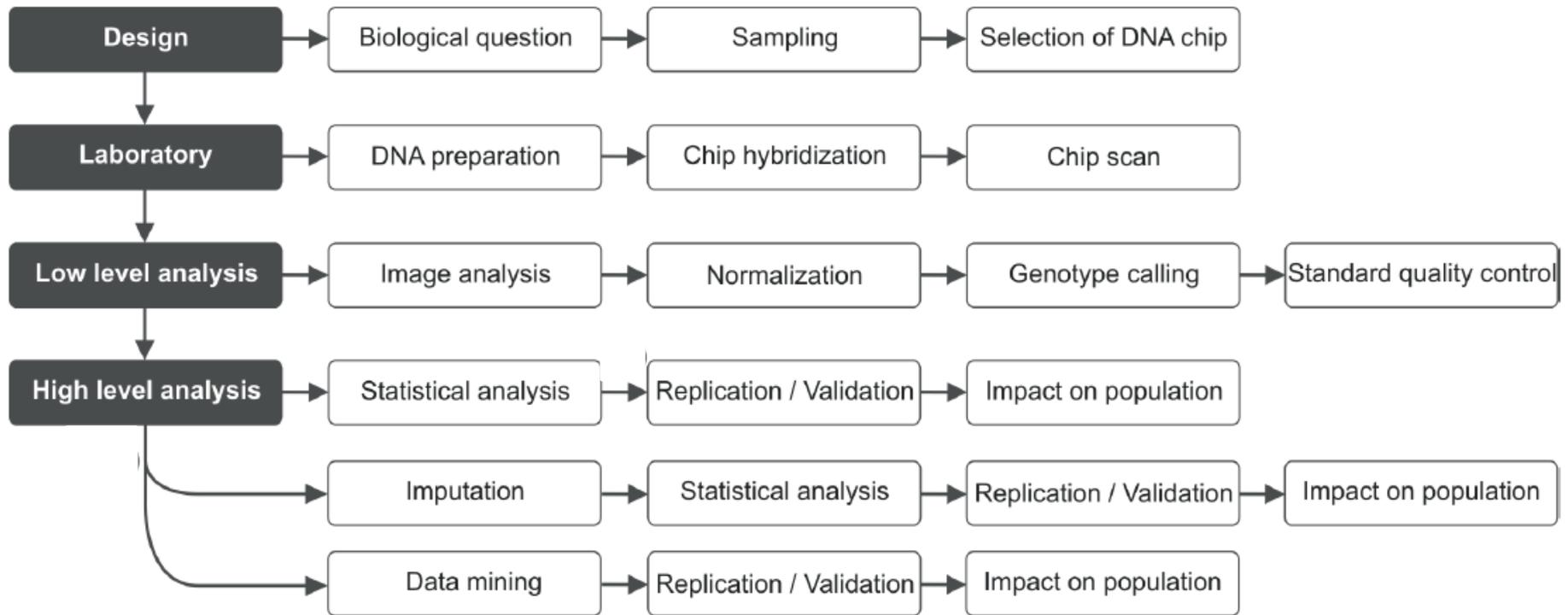
- The plot in Panel C shows the P values for all genotyped SNPs that have survived a quality-control screen (each chromosome, a different color).
- The results implicate a locus on chromosome 9, marked by SNPs 1 and 2, which are adjacent to each other (graph at right), and other neighboring SNPs.

(Manolio 2010)

## 2 Components of a GWAs



## Detailed flow of a genome-wide association study



(Ziegler 2009)

# GWAs belong to Bioinformatics research (1)

**BIOINFORMATICS APPLICATIONS NOTE** Vol. 23 no. 10 2007, pages 1294–1296  
doi:10.1093/bioinformatics/btm108

*Genetics and population analysis*

## GenABEL: an R library for genome-wide association analysis

Yurii S. Aulchenko<sup>1,\*</sup>, Stephan Ripke<sup>2</sup>, Aaron Isaacs<sup>1</sup> and Cornelia M. van Duijn<sup>1</sup>

<sup>1</sup>Department of Epidemiology and Biostatistics, Erasmus MC Rotterdam, Postbus 2040, 3000 CA Rotterdam, The Netherlands and <sup>2</sup>Statistical Genetics Group, Max-Planck-Institute of Psychiatry, Kraepelinstr. 10, D-80804 Munich, Germany

Received on December 3, 2006; revised on February 14, 2007; accepted on March 13, 2007

Advance Access publication March 23, 2007

Associate Editor: Martin Bishop

### ABSTRACT

Here we describe an R library for genome-wide association (GWA) analysis. It implements effective storage and handling of GWA data, fast procedures for genetic data quality control, testing of association of single nucleotide polymorphisms with binary or quantitative traits, visualization of results and also provides easy interfaces to standard statistical and graphical procedures implemented in base R and special R libraries for genetic analysis. We evaluated GenABEL using one simulated and two real data sets. We conclude that GenABEL enables the analysis of GWA data on desktop computers.

**Availability:** <http://cran.r-project.org>

**Contact:** [i.aoultchenko@erasmusmc.nl](mailto:i.aoultchenko@erasmusmc.nl)

With these objectives in mind, we developed the GenABEL software, implemented as an R library. R is a free, open source language and environment for statistical analysis (<http://www.r-project.org/>). Building upon existing statistical analysis facilities allowed for rapid development of the package.

## 2 IMPLEMENTATION

### 2.1 Objective (1)

GWA data storage using standard R data types is ineffective. A SNP genotype for a single person may take four values (AA, AB, BB and missing). Two bits, therefore, are required to store these data. However, the standard R data types occupy 32 bits, leading to an overhead of 1500%, compared to the theoretical optimum. Use of the raw R data format, occupying

# GWAs belong to Bioinformatics research (2)

BIOINFORMATICS

Vol. 26 ISMB 2010, pages i208–i216  
doi:10.1093/bioinformatics/btq191

## Multi-population GWA mapping via multi-task regularized regression

Kriti Puniyani, Seyoung Kim and Eric P. Xing\*

School of Computer Science, Carnegie Mellon University, Pittsburgh, PA, USA

### ABSTRACT

**Motivation:** Population heterogeneity through admixing of different founder populations can produce spurious associations in genome-wide association studies that are linked to the population structure rather than the phenotype. Since samples from the same population generally co-evolve, different populations may or may not share the same genetic underpinnings for the seemingly common phenotype. Our goal is to develop a unified framework for detecting causal genetic markers through a joint association analysis of multiple populations.

**Results:** Based on a multi-task regression principle, we present a multi-population group lasso algorithm using  $L_1/L_2$ -regularized regression for joint association analysis of multiple populations that are stratified either via population survey or computational estimation. Our algorithm combines information from genetic markers across populations, to identify causal markers. It also implicitly accounts for correlations between the genetic markers, thus enabling better control over false positive rates. Joint analysis across populations enables the detection of weak associations common to all populations with greater power than in a separate analysis of each population. At the same time, the regression-based framework allows causal alleles that are unique to a subset of the populations to be correctly identified. We demonstrate the effectiveness of our method on HapMap-simulated and lactase persistence datasets, where we significantly outperform state of the art methods, with greater power for detecting weak associations and reduced spurious associations.

**Availability:** Software will be available at <http://www.sailing.cs.cmu.edu/>

the geographical distribution of the individuals. For example, it has been shown that such heterogeneity is present in the HapMap data (The International HapMap Consortium, 2005) across European, Asian and African populations; and heterogeneity at a finer scale within European ancestry has been found in many genomic regions in the UK samples of Wellcome trust case control consortium (WTCCC) dataset (Wellcome Trust Case Control Consortium, 2007). Although the standard assumption in existing approaches for association mapping is that the effects of causal mutations are likely to be common across multiple populations, the individuals in the same population or geographical region tend to co-evolve, and are likely to possess a population-specific causal allele for the same phenotype. For example, Tishkoff *et al.* (2006) reported that the lactase-persistence phenotype is caused by different mutations in Africans and Europeans. In addition, the same genetic variation has been observed to be correlated with gene-expression levels with different association strengths across different HapMap populations. Our goal is to be able to leverage information across multiple populations, to find causal markers in a multi-population association study.

### 1.1 Highlights of this article

We propose a novel multi-task-regression-based technique that performs a joint GWA mapping on individuals from multiple populations, rather than separate analysis of each population, to detect associated genome variations. The joint inference is achieved by using a multi-population group lasso (MPGL), with an  $L_1/L_2$

# GWAs belong to Bioinformatics research (3)

**BIOINFORMATICS APPLICATIONS NOTE** Vol. 24 no. 1 2008, pages 140–142  
doi:10.1093/bioinformatics/btm549

*Genetics and population analysis*

## GWAsimulator: a rapid whole-genome simulation program

Chun Li<sup>1,\*</sup> and Mingyao Li<sup>2</sup>

<sup>1</sup>Department of Biostatistics, Vanderbilt University School of Medicine, Nashville, TN 37232 and <sup>2</sup>Department of Biostatistics and Epidemiology, University of Pennsylvania School of Medicine, Philadelphia, PA 19104, USA

Received on July 20, 2007; revised on October 10, 2007; accepted on October 29, 2007

Advance Access publication November 15, 2007

Associate Editor: Martin Bishop

### ABSTRACT

**Summary:** GWAsimulator implements a rapid moving-window algorithm to simulate genotype data for case-control or population samples from genomic SNP chips. For case-control data, the program generates cases and controls according to a user-specified multi-locus disease model, and can simulate specific regions if desired. The program uses phased genotype data as input and has the flexibility of simulating genotypes for different populations and different genomic SNP chips. When the HapMap phased data are used, the simulated data have similar local LD patterns as the HapMap data. As genome-wide association (GWA) studies become increasingly popular and new GWA data analysis methods are being developed, we anticipate that GWAsimulator will be an important tool for evaluating performance of new GWA analysis methods.

**Availability:** The C++ source code, executables for Linux, Windows and MacOS, manual, example data sets and analysis program are available at <http://biostat.mc.vanderbilt.edu/GWAsimulator>

**Contact:** [chun.li@vanderbilt.edu](mailto:chun.li@vanderbilt.edu)

**Supplementary information:** Supplementary data are available at *Bioinformatics* online.

### 2 METHODS

The program can generate unrelated case-control (sampled retrospectively conditional on affection status) or population (sampled randomly) data of genome-wide SNP genotypes with patterns of LD similar to the input data.

#### 2.1 Phased input data and control file

The program requires phased data as input. If the HapMap data are used, the number of phased autosomes and X chromosomes are 120 and 90 for both CEU and YRI, 90 and 68 for CHB, and 90 and 67 for JPT. Additional parameters needed by the program should be provided in a control file, including disease model (see Section 2.2), window size (see Section 2.3), whether to output the simulated data (see Section 2.4), and the number of subjects to be simulated.

#### 2.2 Determination of disease model

For simulations of case-control data, a disease model is needed. The program allows the user to specify disease model parameters, including disease prevalence, the number of disease loci, and for each disease locus, its location, risk allele and genotypic relative risk. If the user wants to simulate specific regions, the start and end positions need

# GWAs belong to Bioinformatics research (4)

**BIOINFORMATICS APPLICATIONS NOTE** Vol. 25 no. 5 2009, pages 662–663  
doi:10.1093/bioinformatics/btp017

*Genome analysis*

## AssociationViewer: a scalable and integrated software tool for visualization of large-scale variation data in genomic context

Olivier Martin<sup>1,†</sup>, Armand Valsesia<sup>1,2,†</sup>, Amalio Telenti<sup>3</sup>, Ioannis Xenarios<sup>1</sup>  
and Brian J. Stevenson<sup>1,2,\*</sup>

<sup>1</sup>Swiss Institute of Bioinformatics, <sup>2</sup>Ludwig Institute for Cancer Research, 1015 Lausanne and <sup>3</sup>Institute of Microbiology, University Hospital, University of Lausanne, 1011 Lausanne, Switzerland

Received on September 16, 2008; revised on December 16, 2008; accepted on January 5, 2009

Advance Access publication January 25, 2009

Associate Editor: John Quackenbush

### ABSTRACT

**Summary:** We present a tool designed for visualization of large-scale genetic and genomic data exemplified by results from genome-wide association studies. This software provides an integrated framework to facilitate the interpretation of SNP association studies in genomic context. Gene annotations can be retrieved from Ensembl, linkage disequilibrium data downloaded from HapMap and custom data imported in BED or WIG format. AssociationViewer integrates functionalities that enable the aggregation or intersection of data tracks. It implements an efficient cache system and allows the display of several, very large-scale genomic datasets.

**Availability:** The Java code for AssociationViewer is distributed under the GNU General Public Licence and has been tested on Microsoft Windows XP, MacOSX and GNU/Linux operating systems. It is available from the SourceForge repository. This also includes Java webstart, documentation and example datafiles.

**Contact:** brian.stevenson@licr.org

**Supplementary information:** Supplementary data are available at <http://sourceforge.net/projects/associationview/> online.

represented in BED or WIG format and implements aggregation (union) or intersection of data tracks.

## 2 PROGRAM OVERVIEW

### 2.1 Cache and memory management

With increasing data volumes, efficient resource management is essential. One approach is to store the data in a cache with fast indexing mechanisms to retrieve the data, and to keep in memory only the information that is visualized. We implemented such a system in AssociationViewer. For comparison, loading a single dataset with 500 K SNPs in WGAViewer needs about 224 MB of RAM, whereas loading 10 different datasets (a total of 10 M data points) and displaying all genes on chromosome 1 needs only 50 MB in AssociationViewer.

### 2.2 Data import and export

A typical GWA dataset consists of a list of SNPs with *P*-values derived from an association analysis. In AssociationViewer, such

## 3 Study Design

### Components of a study design for GWA studies

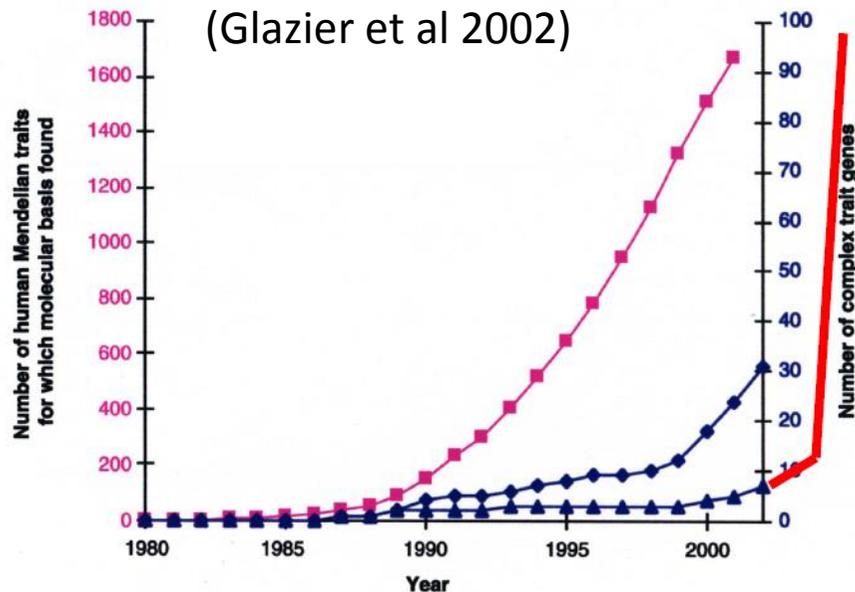
- The design of a genetic association study may refer to
  - study scale:
    - Genetic (e.g., hypothesis-drive, panel of candidate genes)
    - Genomic (e.g., hypothesis-free, genome-wide)
  - marker design:
    - Which markers are most informative in GWAs? Common variants-SNPs and/or Rare Variants (MAF<1%)
    - Which platform is the most promising? Least error-prone? Marker-distribution over the genome?
  - subject design



## Hypothesis: Common Disease – Common Variant (CDCV)

- This hypothesis predicts that the genetic risk for **common diseases will often be due to disease-predisposing alleles with relatively high frequencies**; there will be one or a few predominating disease alleles at each of the major underlying disease loci (Lander, 1996; Chakravarti, 1999; Weiss & Clark, 2002; Becker, 2004).
- The hypothesis speculates that the gene variation underlying susceptibility to common heritable diseases existed within the **founding population of contemporary humans** → explains the success of GWAs?
- Whether the CDCV hypothesis is true for most diseases is yet unknown but there are a few prototypical examples: the APOE e4 allele in Alzheimer disease
- An alternative hypothesis = Common Disease Rare Variant (CDRV) hypothesis

## Number of traits for which molecular basis exists



PINK : Human Mendelian traits

BLUE middle line : All complex traits

BLUE bottom line + red extension:  
Human complex traits

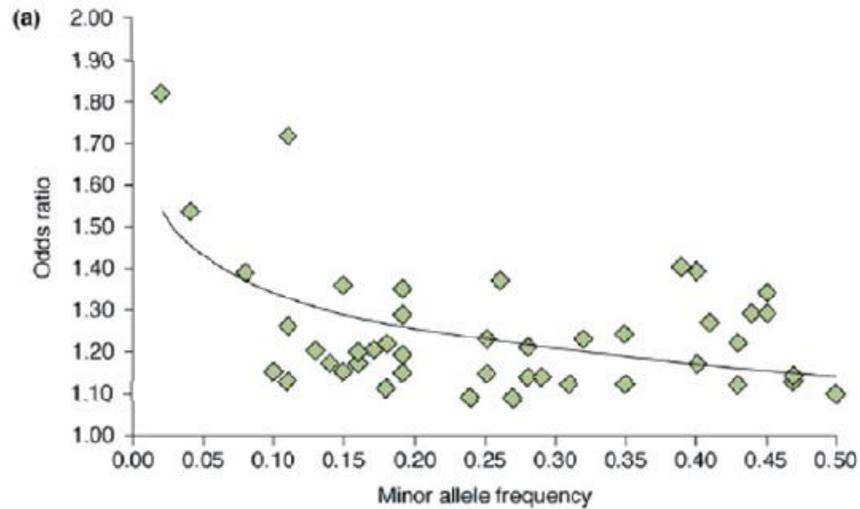
### Complex disease (definition):

The term complex trait/disease refers to any phenotype that

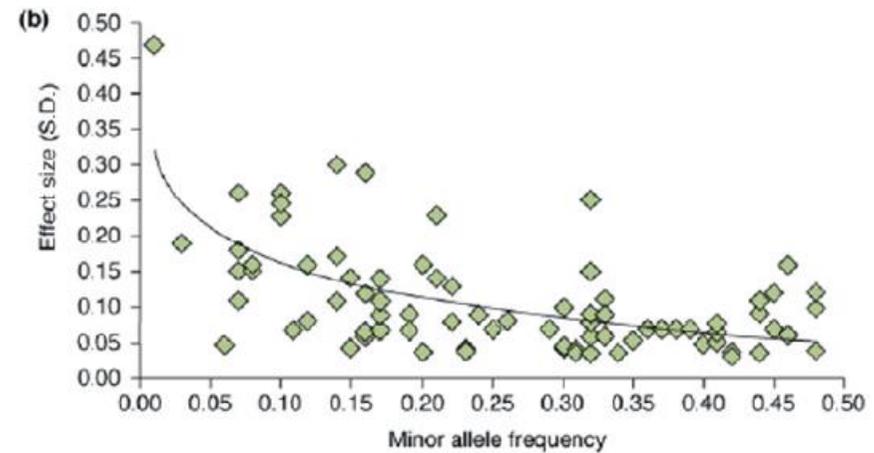
**does NOT exhibit classic Mendelian inheritance attributable to a single gene;**

although they may exhibit familial tendencies (familial clustering, concordance among relatives).

## Dichotomous Traits



## Quantitative Traits



Arking & Chakravarti 2009 Trends Genet

### Food for thought:

- The higher the MAF, the lower the effect size
- Rare variants analysis is in its infancy ....

### 3.b Subject Level

<b>Aim</b>	<b>Selection scheme</b>
<b>Increased effect size</b>	Extreme sampling: Severely affected cases vs. extremely normal controls
<b>Genes causing early onset</b>	Affected, early onset vs. normal, elderly
<b>Genes with large / moderate effect size</b>	Cases with positive family history vs. controls with negative family history
<b>Specific GxE interaction</b>	Affected vs. normal subjects with heavy environmental exposure
<b>Longevity genes</b>	Elderly survivors serve as cases vs. young serve as controls
<b>Control for covariates with strong effect</b>	Affected with favorable covariates vs. normal with unfavorable covariate

Morton & Collins 1998 Proc Natl Acad Sci USA 95:11389

## Popular design 1: cases and controls

### Avoiding bias – checking assumptions:

1. Cases and controls drawn from same population
2. Cases representative for all cases in the population
3. All data collected similarly in cases and controls

### Advantages:

1. Simple
2. Cheap
3. Large number of cases and controls available
4. Optimal for studying rare diseases

### Disadvantages:

1. Population stratification
2. Prone to batch effects and other biases
3. Cases usually mild cases, etc
4. Overestimation of risk for common diseases

## Popular design 2: family-based

### Avoiding bias – checking assumptions:

1. Families representative for population of interest
2. Same genetic background in both parents

### Advantages:

1. Controls immune to population stratification (no association without linkage, no “spurious” (false positive) association)
2. Checks for Mendelian inheritance possible (fewer genotyping errors)
3. Parental phenotyping not required (late onset diseases)

4. Simple logistics for diseases in children
5. Allows investigating imprinting (“bad allele” from father or mother?)

### Disadvantages

1. Cost inefficient
2. Lower power when compared with case-control studies
3. Sensitive to genotyping errors

## 4 Pre-analysis steps

### 3.a Quality control – The Travemünde Criteria

#### Standard file format for GWA studies

Standard data format: tped = transposed ped format file

FamID	PID	FID	MID	SEX	AFF	SNP1 <sub>1</sub>	SNP1 <sub>2</sub>	SNP2 <sub>1</sub>	SNP2 <sub>2</sub>
1	1	0	0	1	1	A	A	G	T
2	1	0	0	1	1	A	C	T	G
3	1	0	0	1	1	C	C	G	G
4	1	0	0	1	2	A	C	T	T
5	1	0	0	1	2	C	C	G	T
6	1	0	0	1	2	C	C	T	T

ped file

Chr	SNP name	Genetic distance	Chromosomal position
1	SNP1	0	123456
1	SNP2	0	123654

map file

## Standard file format for GWA studies (continued)

Chr	SNP	Gen. dist.	Pos	PID 1	PID 2	PID 3	PID 4	PID 5	PID 6						
1	SNP1	0	123456	A	A	A	C	C	C	A	C	C	C	C	C
1	SNP2	0	123654	G	T	G	T	G	G	T	T	G	T	T	T

tfam file: First 6 columns of standard ped file

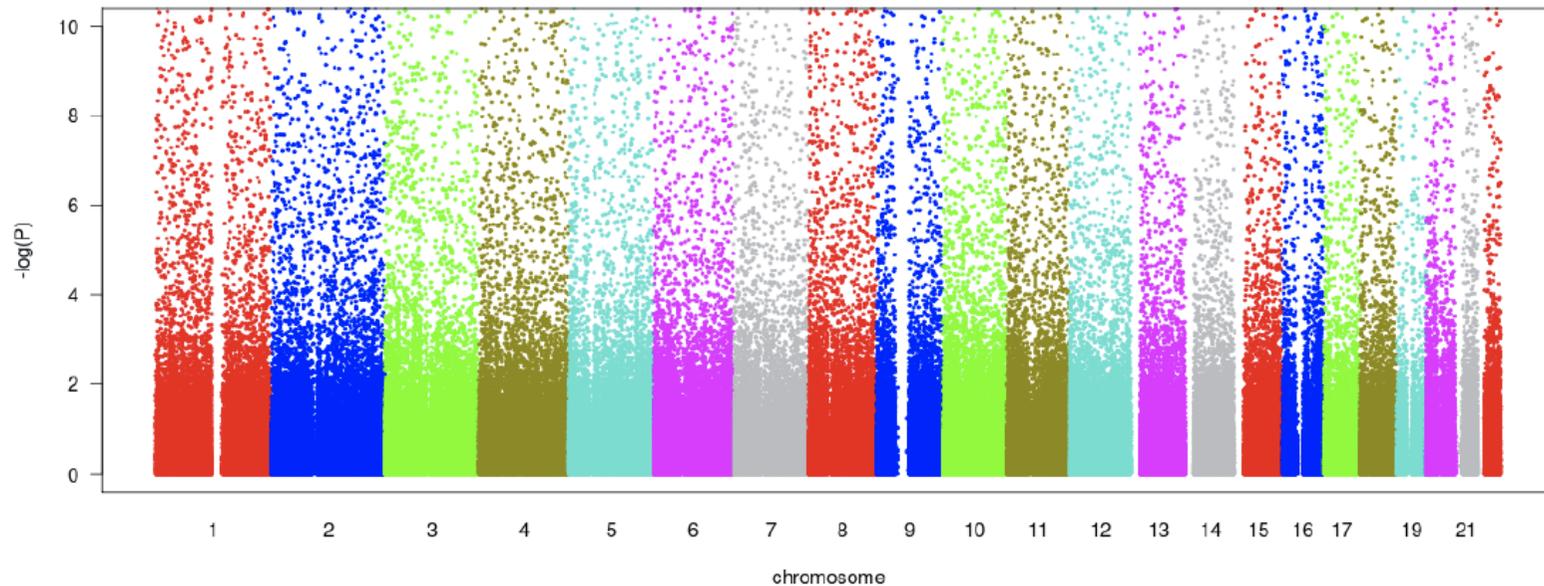
tped file

FamID	PID	FID	MID	SEX	AFF
1	1	0	0	1	1
2	1	0	0	1	1
3	1	0	0	1	1
4	1	0	0	1	2
5	1	0	0	1	2
6	1	0	0	1	2

tfam file

## Why is quality control (QC) important?

**BEFORE QC** → true signals are lost in false positive signals

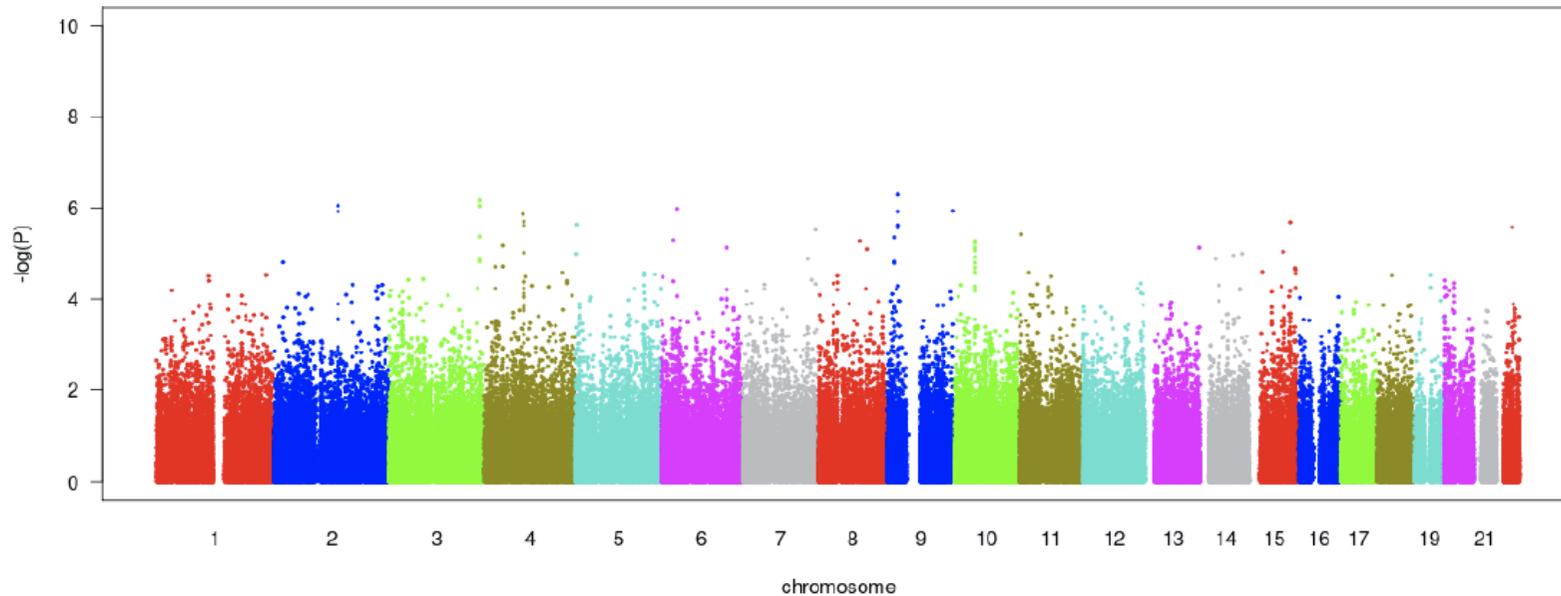


Ger MI FS I, Affymetrix 500k array set, SNPs on chip: 493,840

(Ziegler and Van Steen 2010)

## Why is quality control important?

**AFTER QC** → skyline of Manhattan (→ name of plot: Manhattan plot):



Ger MI FS I, Affymetrix 500k array set, SNPs on chip: 493,840

SNPs passing standard quality control: 270,701

(Ziegler and Van Steen 2010)

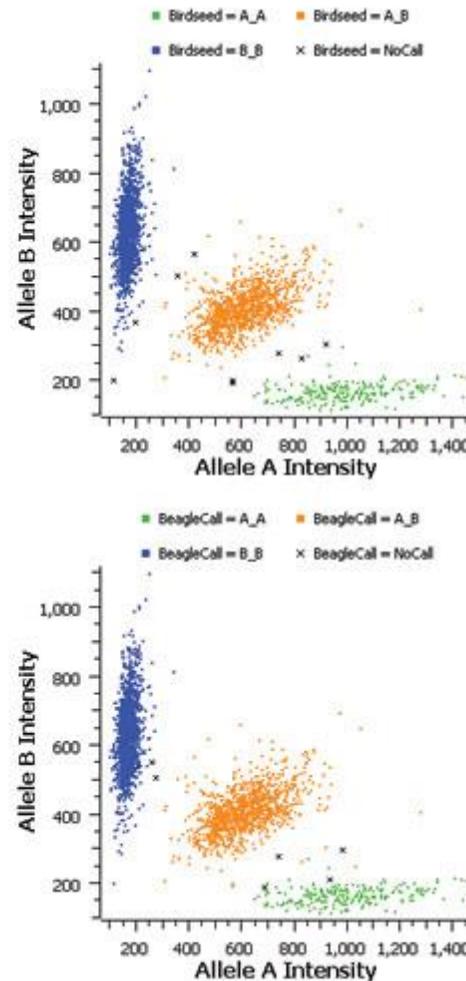
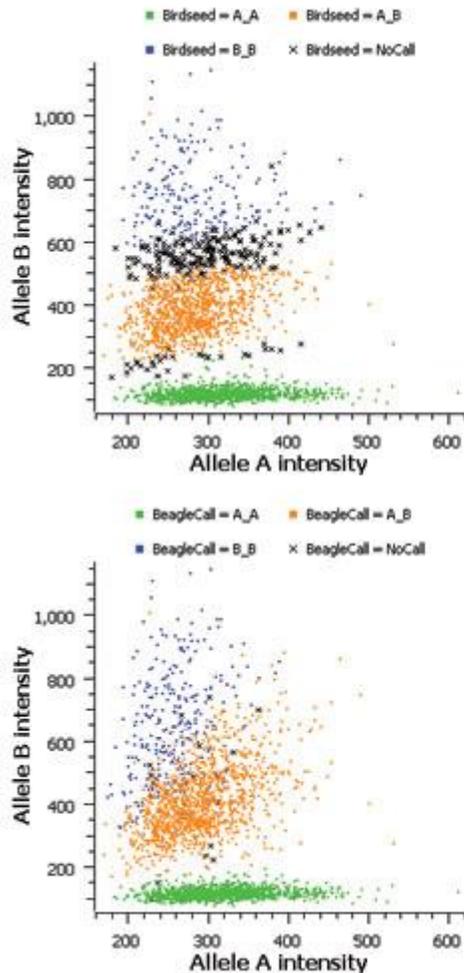
## What is the standard quality control?

- Quality control can be performed on different levels:
  - Subject or sample level
  - Marker level (here: SNP level)
  - X-chromosomal SNP level
- Consensus on how to best QC data has led to the so-called “Travemünde criteria” (obtained in the town Travemünde)

## Example of QC at the marker level

- **Minor allele frequency (MAF):**
  - Genotype calling algorithms perform poorly for SNPs with low MAF
  - Power is low for detecting associations to genetic markers with low MAF (with standard large-sample statistics)
- **Missing frequency (MiF)**
  - 1 minus call rate
  - MiF needs to be investigated separately in cases and controls because differential missingness may bias association results
- **Hardy-Weinberg equilibrium (HWE)**
  - SNPs excluded if substantially more or fewer subjects heterozygous at a SNP than expected (excess heterozygosity or heterozygote deficiency)

## Reading genotype calling cluster plots



**Allele signal intensity cluster plots** for two different SNPs from the same study population.

Upper panels: Birdseed genotypes

Lower panels: BEAGLECALL genotypes.

The plots on the left show a SNP with poor resolution of A\_B and B\_B genotype clusters and the increased clarity of genotype calls that comes from using BEAGLECALL

(Golden Helix Blog)

## What is Hardy-Weinberg Equilibrium (HWE)?

Consider diallelic SNP with alleles  $A_1$  and  $A_2$

- Genotype frequencies

$$P(A_1A_1) = p_{11}, P(A_1A_2) = p_{12}, P(A_2A_2) = p_{22}$$

- Allele frequencies  $P(A_1) = p = p_{11} + \frac{1}{2}p_{12}$ ,  $P(A_2) = q = p_{22} + \frac{1}{2}p_{12}$

If

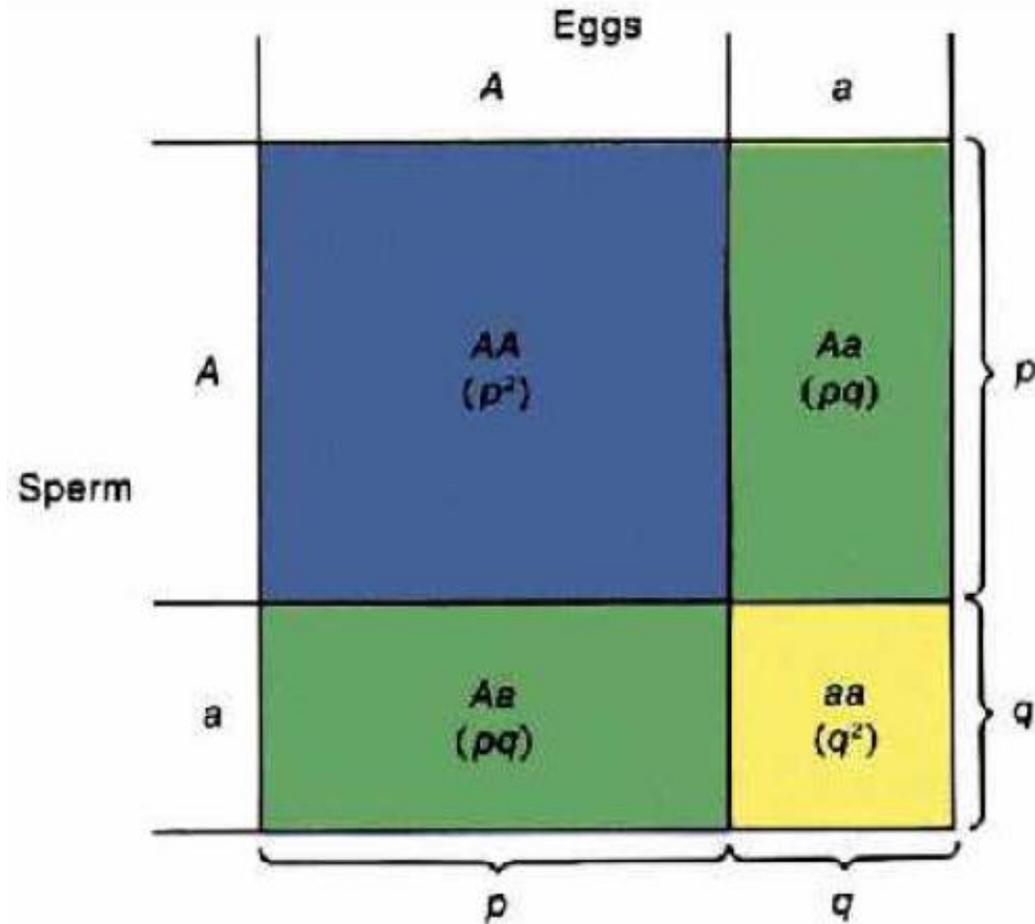
- $P(A_1A_1) = p_{11} = p^2$
- $P(A_1A_2) = p_{12} = 2pq$
- $P(A_2A_2) = p_{22} = q^2$

the population is said to be in HWE at the SNP

(Ziegler and Van Steen 2010)

## What is Hardy-Weinberg Equilibrium (HWE)?

Consider diallelic SNP with alleles A and a



## **Distorting factors to HWE causing evolution to occur**

### **1. Non-random mating**

**2. Mutation** - by definition mutations change allele frequencies causing evolution

**3. Migration** - if new alleles are brought in by immigrants or old alleles are taken out by emigrants then the frequencies of alleles will change causing evolution

**4. Genetic drift** - random events due to small population size (bottleneck caused by storm and leading to reduced variation, migration events leading to founder effects)

**5. Natural selection** – some genotypes give higher reproductive success  
**(Darwin)**

## The Travemünde criteria

Level	Filter criterion	Standard value for filter
Sample level	Call fraction	$\geq 97\%$
	Cryptic relatedness	Study specific
	Ethnic origin	Study specific; visual inspection of principal components
	Heterozygosity	Mean $\pm$ 3 std.dev. over all samples
	Heterozygosity by gender	Mean $\pm$ 3 std.dev. within gender group
SNP level	MAF	$\geq 1\%$
	MiF	$\leq 2\%$ in any study group, e.g., in both cases and controls
	MiF by gender	$\leq 2\%$ in any gender
	HWE	$p < 10^{-4}$

(Ziegler 2009)

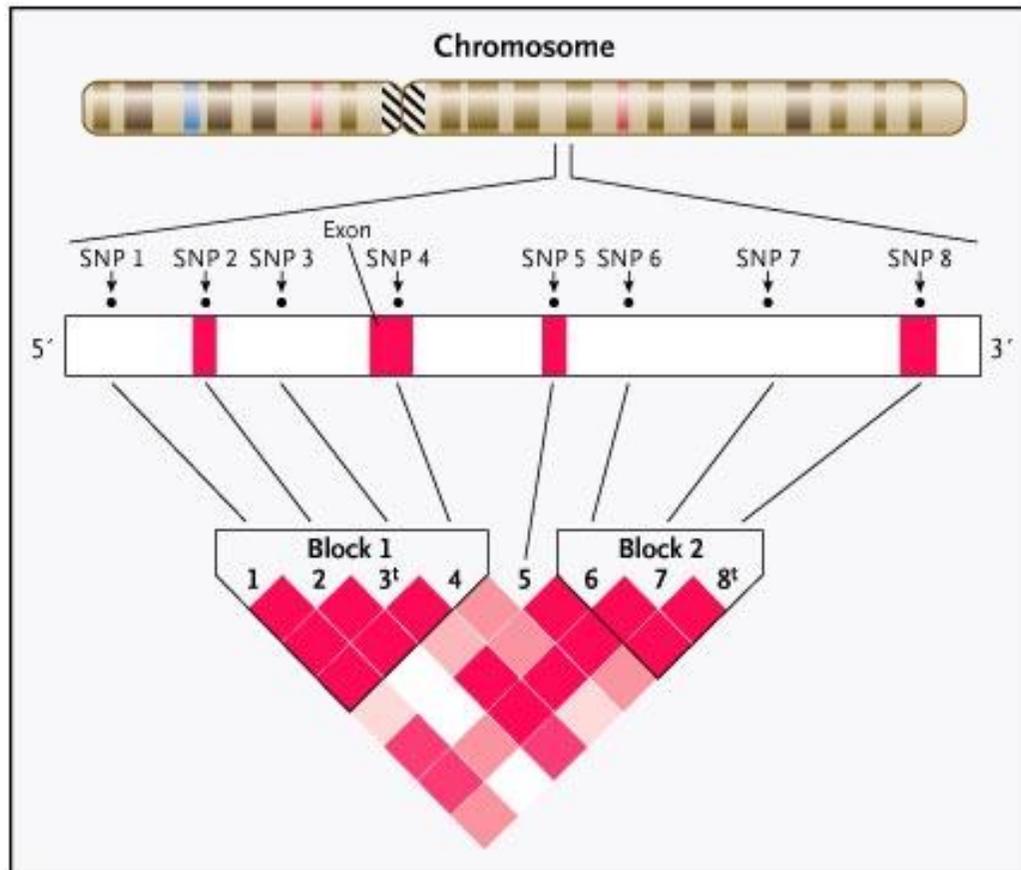
## The Travemünde criteria

Level	Filter criterion	Standard value for filter
<b>SNP level</b>	Difference between control groups	$p > 10^{-4}$ in trend test
	Gender differences among controls	$p > 10^{-4}$ in trend test
<b>X-Chr SNPs</b>	Missingness by gender	No standards available
	Proportion of male heterozygote calls	No standards available
	Absolute difference in call fractions for males and females	No standards available
	Gender-specific heterozygosity	No standard value available

(Ziegler 2009)

## 4.b Linkage disequilibrium

Mapping the “relationships” between SNPs (Christensen and Murray 2007)



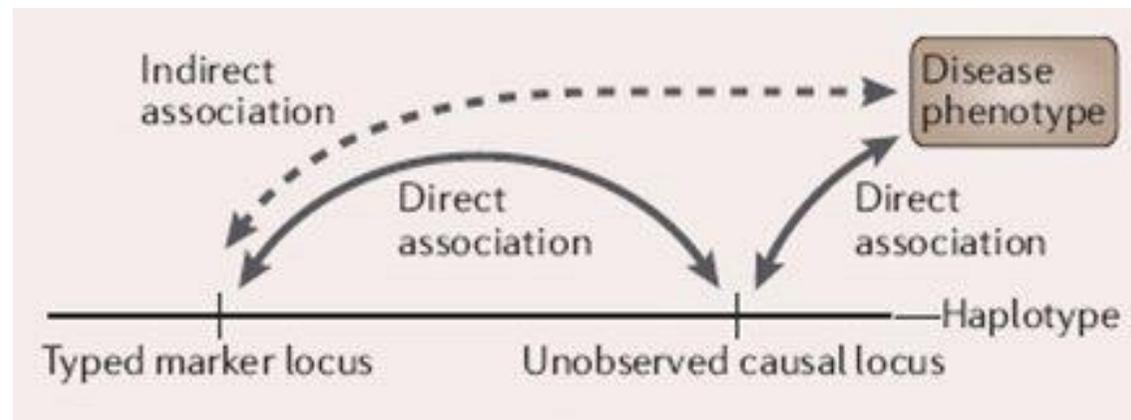
(HaploView software)

## Linkage disequilibrium (LD) between genetic markers

- **Linkage Disequilibrium (LD)** is a measure of co-segregation of alleles in a population – allelic association

Two alleles at different loci that occur together on the same chromosome (or gamete) more often than would be predicted by random chance.

- It is a very important concept for GWAs, since it gives the rationale for performing genetic association studies



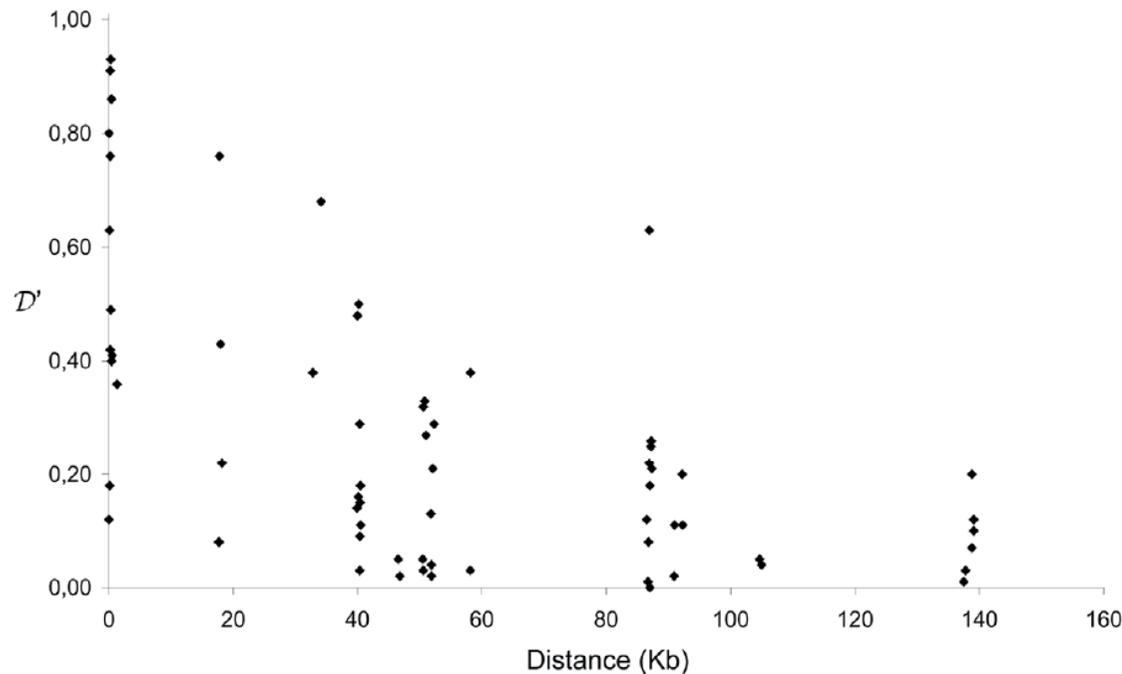
## Distances among SNPs

- The measure  $D$  is defined as the difference between the observed and expected (under the null hypothesis of independence) proportion of “haplotypes” bearing specific alleles at two loci:  $p_{AB} - p_A p_B$

	$A$	$a$
$B$	$p_{AB}$	$p_{aB}$
$b$	$p_{Ab}$	$p_{ab}$

- A **haplotype** is a linear arrangement of alleles on the same chromosome that have been inherited as a unit.
- Because of its interpretation, the **measure  $r^2$  (coefficient of determination)** is most often used for GWAs

## How far does linkage disequilibrium extend?



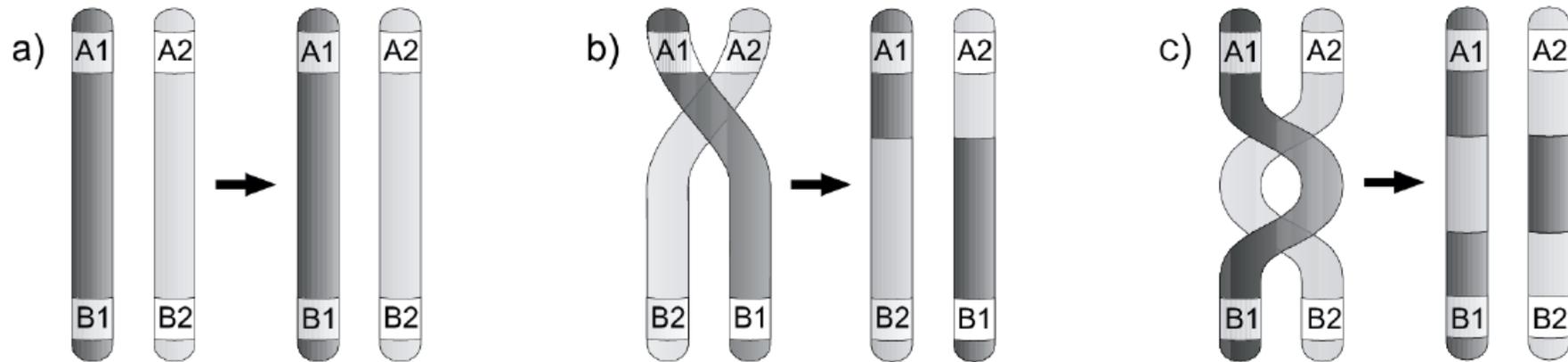
$D'$  (Lewontin's  $D$  prime) is the absolute ratio of  $D$  compared with its maximum value.

$D' = 1$  : complete LD

(Hecker et al 2003)

- LD is usually a function of distance between the two loci. This is mainly because recombination acts to break down LD in successive generations (Hill, 1966).

## Recombination



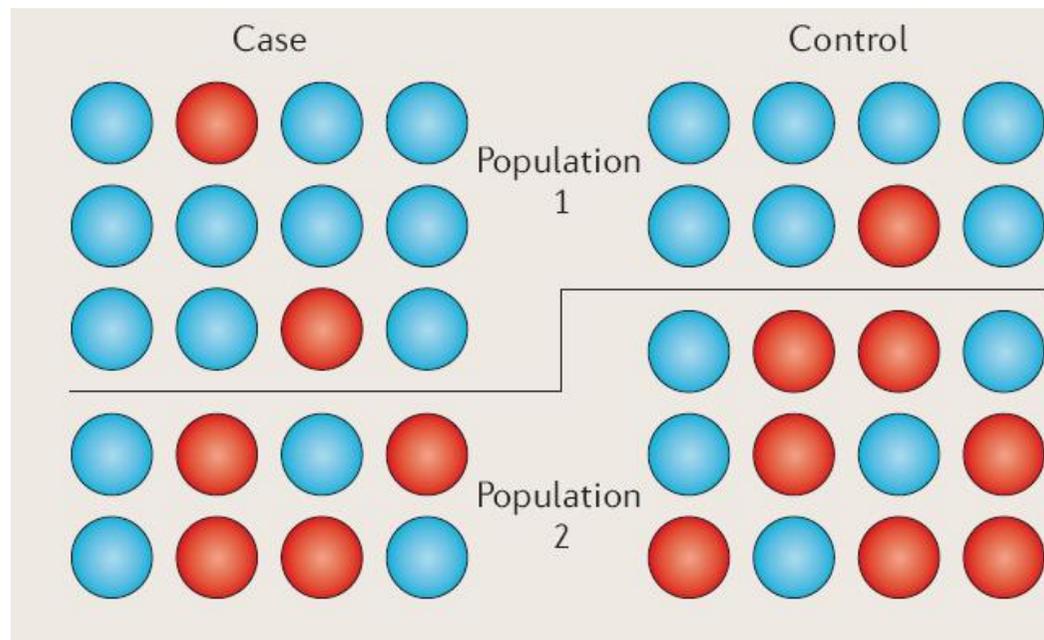
- Relevant measure: recombination fraction (probability of odd number of crossovers) between two chromosomal positions
- Strong correlation between recombination fraction and distance in base pairs

(Ziegler and Van Steen, Brazil 2010)

## 4.c Confounding: population stratification

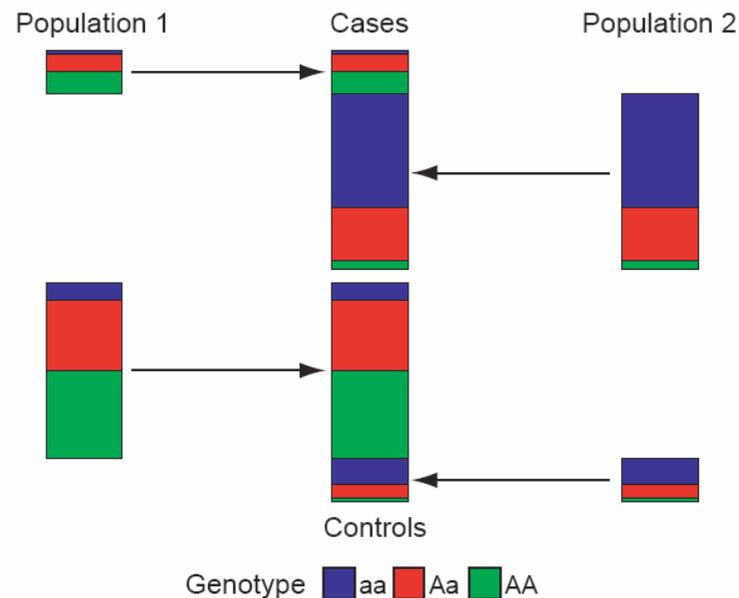
### What is spurious association?

- **Spurious association** refers to false positive association results due to not having accounted for population substructure as a confounding factor in the analysis



## What is spurious association?

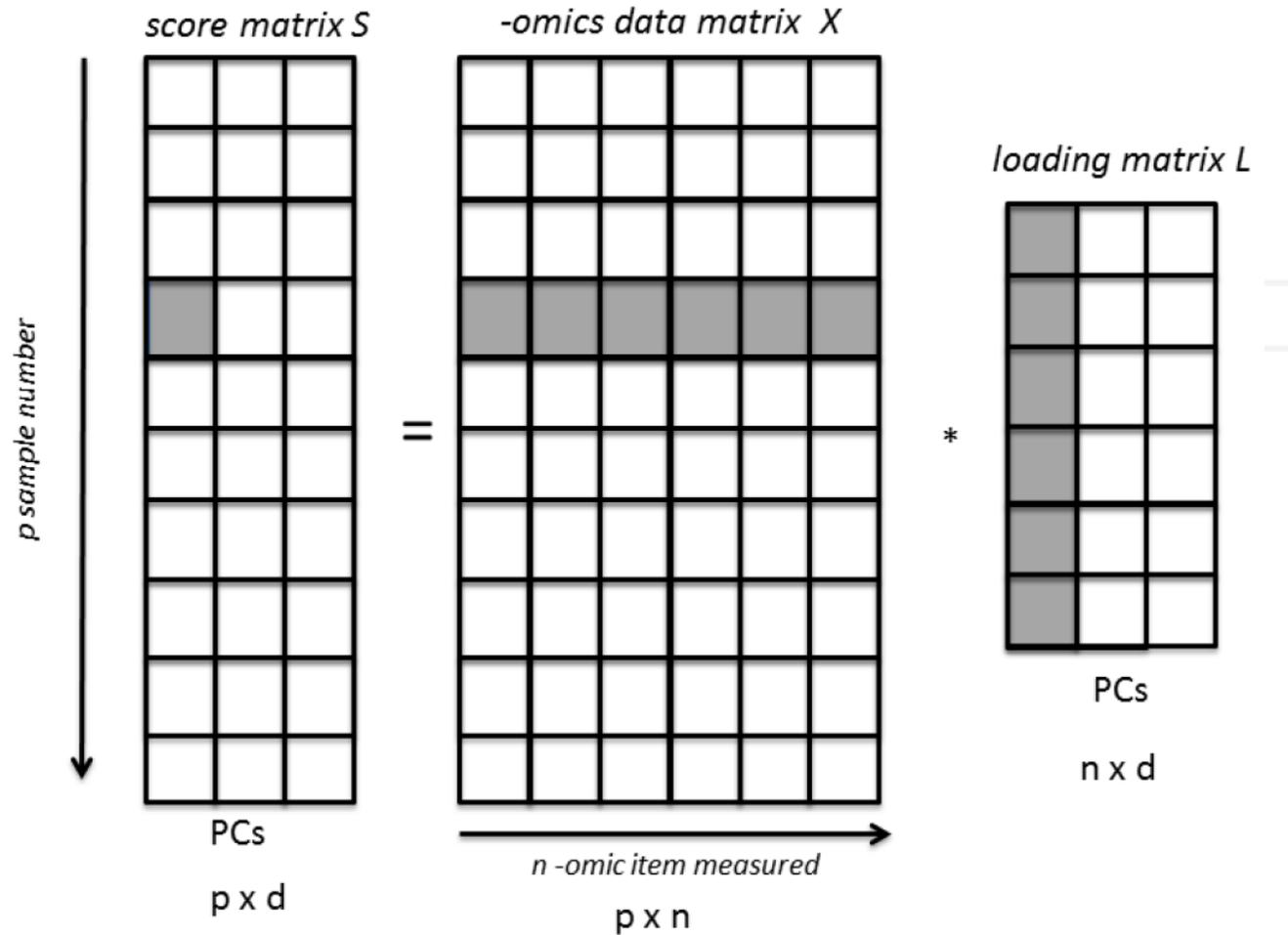
- Typically, there are two characteristics present:
  - A difference in proportion of individual from two (or more) subpopulation in case and controls
  - Subpopulations have different allele frequencies at the locus.



## What are typical methods to deal with population stratification?

- Methods to deal with spurious associations generated by population structure generally require a number (at least  $>100$ ) of widely spaced null SNPs that have been genotyped in cases and controls in addition to the candidate SNPs.
- These methods large group into:
  - **Principal components**
  - Structured association methods: “First look for structure (population clusters) and **second** perform an association **analysis** conditional on the cluster allocation”
  - **Genomic control methods**: “**First analyze** and second downplay association test results for over optimism” → see later

# Principal components (and omics: <http://cdn.intechopen.com/pdfs-wm/30002.pdf>)



$$s_{4,1} = x_{4,1} \cdot l_{1,1} + x_{4,2} \cdot l_{2,1} + x_{4,3} \cdot l_{3,1} + x_{4,4} \cdot l_{4,1} + x_{4,5} \cdot l_{5,1} + x_{4,6} \cdot l_{6,1}$$

## Principal components

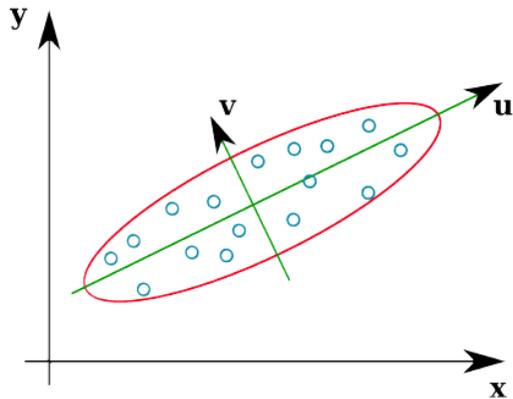


Figure 1: PCA for Data Representation

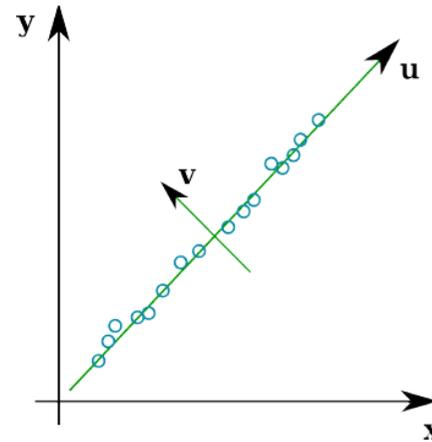


Figure 2: PCA for Dimension Reduction

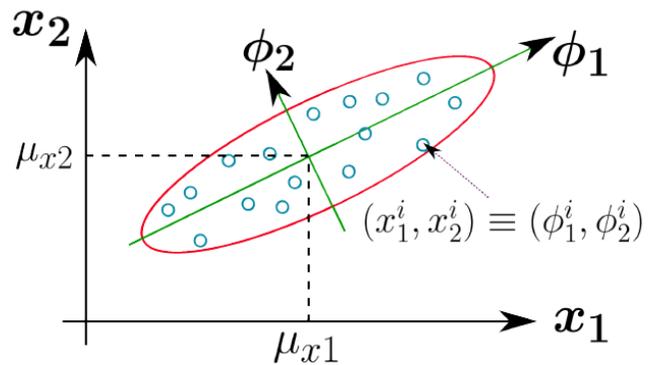
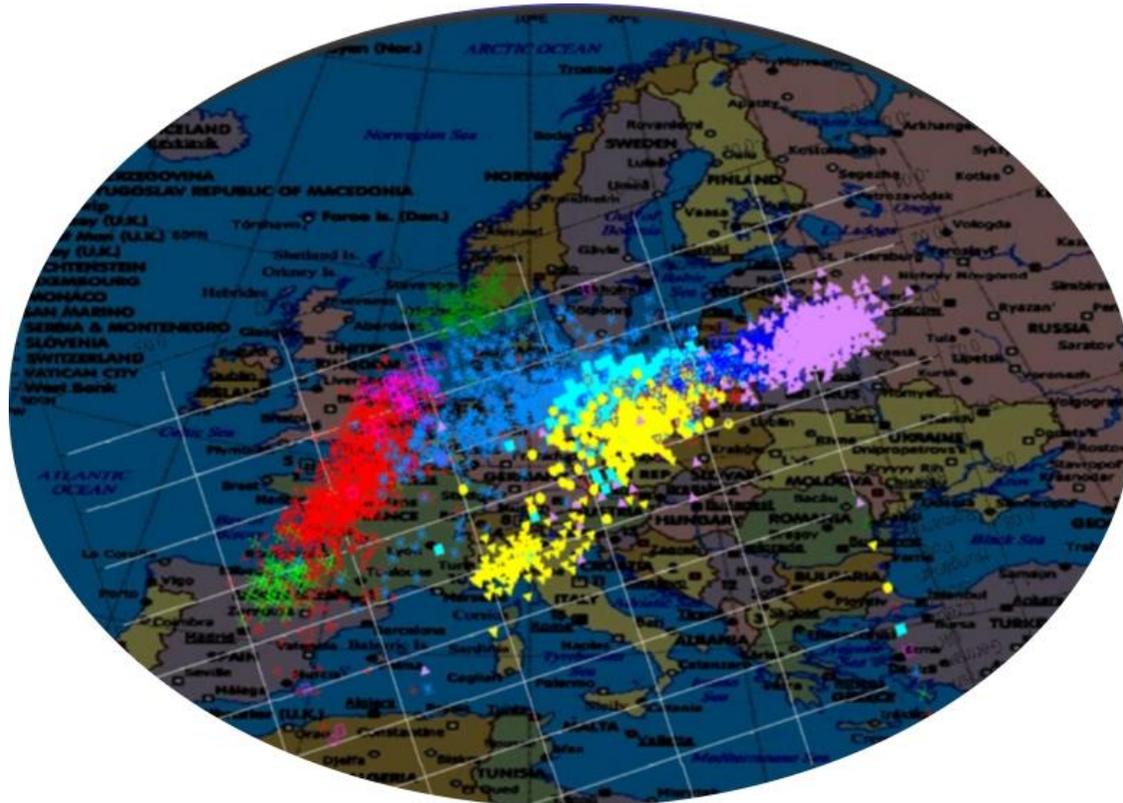


Figure 3: The PCA Transformation

- Find eigenvectors of the covariance matrix for standardized  $(x_1, x_2, \dots)$  [ $\rightarrow$ SNPs]
- These will give you the direction vectors indicated in Fig3 by  $\phi_1$  and  $\phi_2$
- These determine the axes of maximal variation

## Principal components

- In European data, the first 2 principal components “nicely” reflect the N-S and E-W axes !



Y-axis: PC2 (6% of variance); X-axis: PC1 (26% of variance)

## 5 Testing for Associations

### One SNP at a time

	Coding scheme for statistical modeling/testing					
Indiv. genotype	X1	X1	X2	X1	X1	X1
	Additive coding	Genotype coding (general mode of inheritance)		Dominant coding (for a)	Recessive coding (for a)	Advantage Heterozygous
AA	0	0	0	0	0	0
Aa	1	1	0	1	0	1
aa	2	0	1	1	1	0

## Testing for association between case/control status and a SNP

- Fill in the table below and perform a chi-squared test for independence between rows and columns → **genotype test**

	<b>AA</b>	<b>Aa</b>	<b>aa</b>
<b>Cases</b>			
<b>Controls</b>			

Sum of entries =  
cases+controls

- Fill in the table below and perform a chi-squared test for independence between rows and columns → **allelic test (ONLY to be used under HWE)**

	<b>A</b>	<b>a</b>
<b>Cases</b>		
<b>Controls</b>		

Sum of entries is  
 $2 \times (\text{cases} + \text{controls})$

## Testing for association between case/control status and a SNP

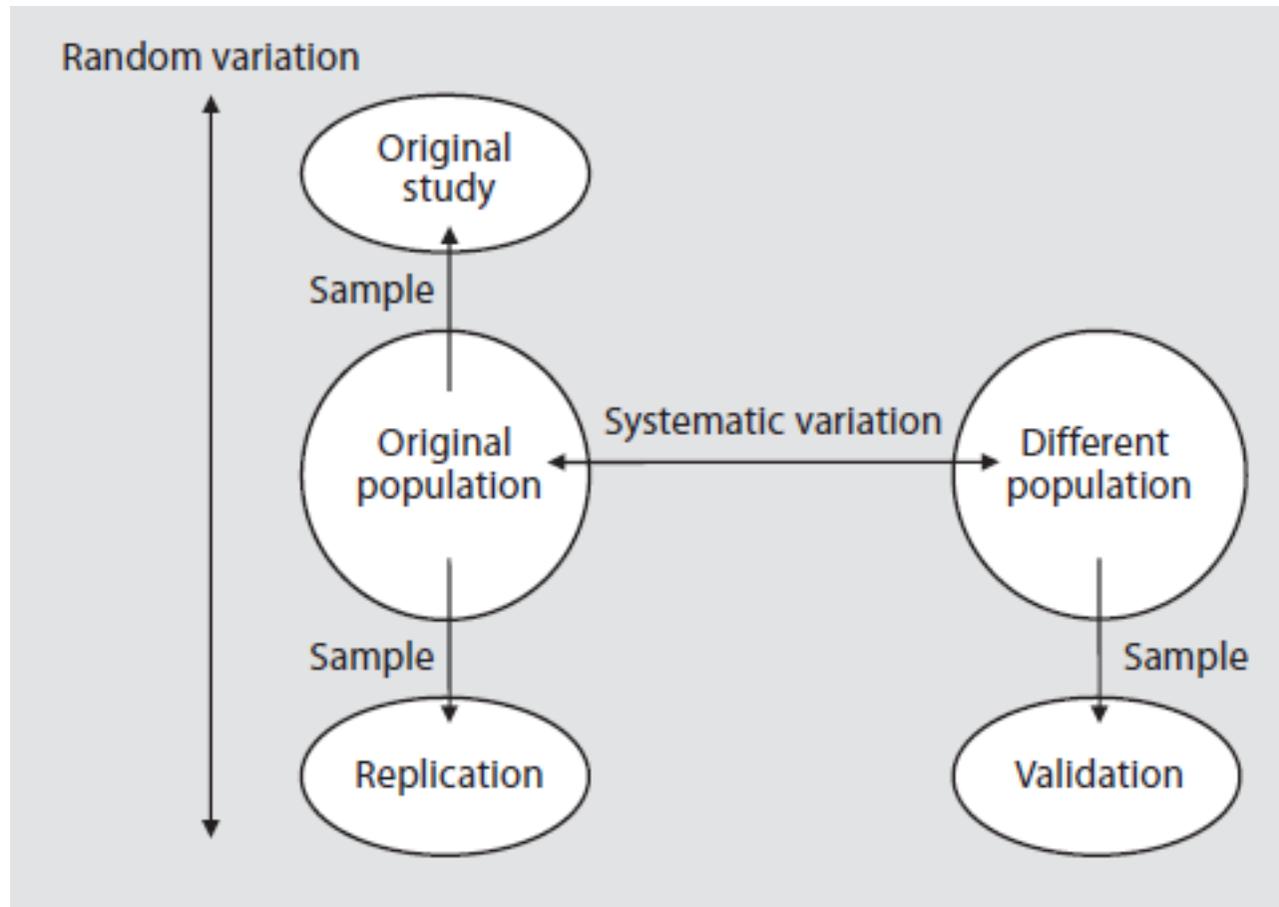
- The genotype test involves a 2df test (note that two variables  $X_1$  and  $X_2$  were needed for genotype coding).
- It has been shown that usually, the additive coding gives adequate power, even when the true underlying mode of inheritance is NOT additive (note that the additive coding can be achieved by only using 1 variable ( $X_1$ )).
- For large sample sizes, a “test for trend” (risk for disease, or average trait increases/decreases with increasing number of “a” copies) theoretically follows a chi-squared distribution with 1df.
- How can we adjust such a test for population substructure?

## Genomic control

- In Genomic Control (GC), a 1-df association test statistic is computed at each of the null SNPs, and a parameter  $\lambda$  is calculated as the empirical median divided by its expectation under the chi-squared 1-df distribution.
- Then the association test is applied at the candidate SNPs, and if  $\lambda > 1$  the test statistics are divided by  $\lambda$ .
  - Under  $H_0$  of no association p-values uniformly distributed
  - In case of population stratification: inflation of test statistics
  - $$\hat{\lambda} = \frac{\text{median}(\chi_1^2, \chi_2^2, \dots, \chi_L^2)}{\text{median}(\mathcal{L}(\chi_1^2))} = \frac{\text{median}(\chi_1^2, \chi_2^2, \dots, \chi_L^2)}{0.456}$$
  - $$\chi_{GC}^2 = \chi^2 / \hat{\lambda}$$

## 6 Replication and Validation

### The difference



(Igl et al. 2009)

## Guidelines for replication studies

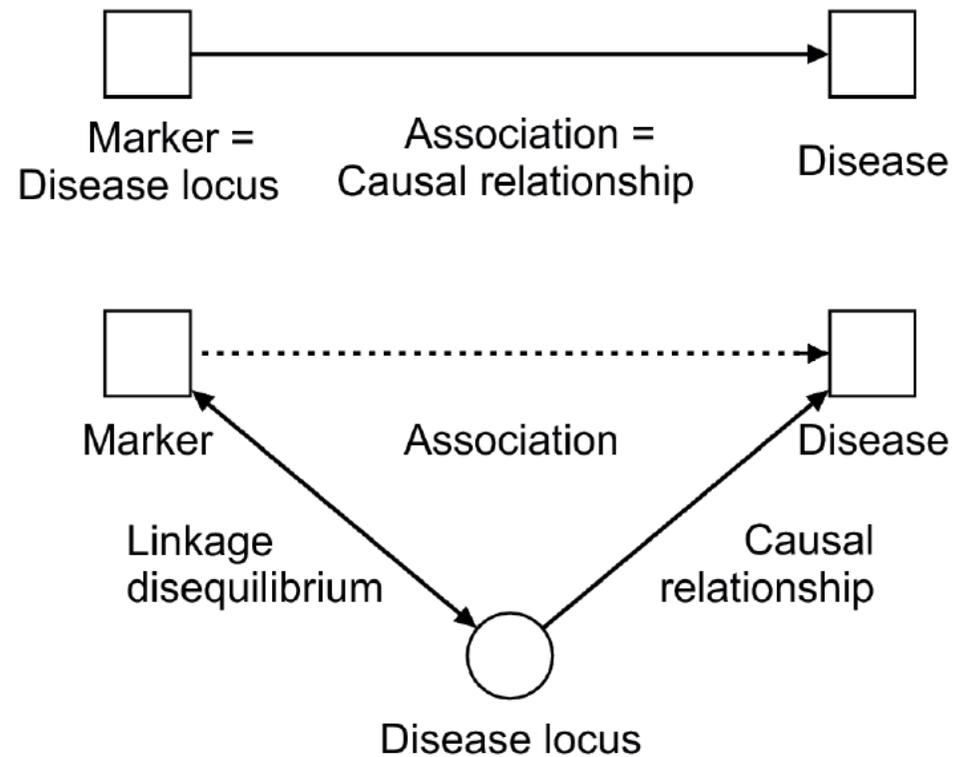
- Replication studies should be of sufficient size to demonstrate the effect
- Replication studies should be conducted in independent datasets
- Replication should involve the same phenotype
- Replication should be conducted in a similar population
- The same SNP should be tested
- The replicated signal should be in the same direction
- Joint analysis should lead to a lower  $p$ -value than the original report
- Well-designed negative studies are valuable

SNPs are most likely to replicate when they:

- Show modest to strong statistical significance
- Have common minor allele frequency
- Exhibit modest to strong **genetic effect size** (~strength of association)

## 7 GWA Interpretation and Follow-Up

### Finding the causal locus

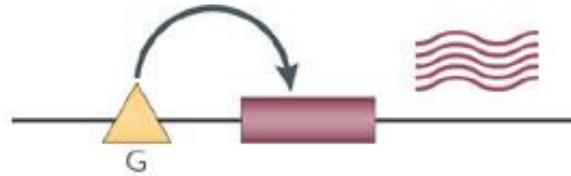
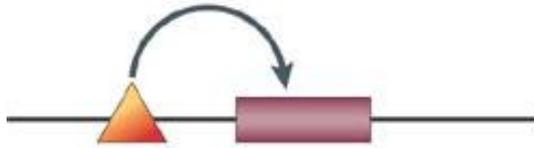


(Ziegler and Van Steen, Brazil 2010)

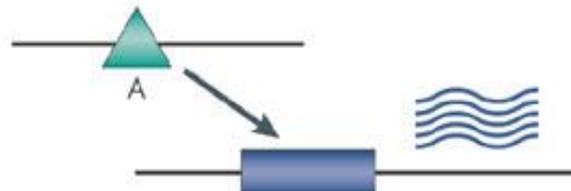
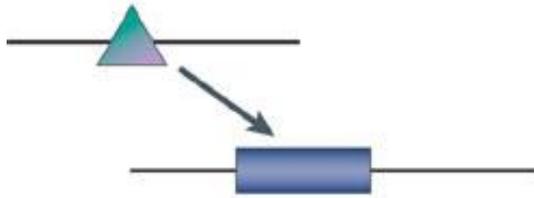


## Finding the causal locus

### a *Cis* (local)



### b *Trans* (distal)



- Cis-acting variants are found close to the target genes and trans-acting variants are located far from the target genes, often on another chromosome.
- Different allelic forms of the cis- and trans-acting variants have different influence on gene expression.

(Cheung and Spielman 2009)

## Finding function

- Evaluating the functional properties of gene sets is commonly used both to verify that the genes implicated in a biological experiment are **functionally relevant** and to **discover unexpected shared functions** between those genes.
- Many **functional annotation databases** have been developed in order to classify genes according their various roles in the cell. E.g., the comprehensive Gene Ontology (GO) is one of the most widely used by many functional enrichment tools

(Glass and Girvan 2014)

# Functional annotation in R



Home » BioViews

## All Packages

### Bioconductor version 2.14 (Release)

Autocomplete biocViews search:

- ▶ Software (824)
  - ▼ AnnotationData (867)
    - ▶ ChipManufacturer (370)
    - ▶ ChipName (195)
    - ▶ CustomArray (2)
    - ▶ CustomCDF (16)
    - ▶ CustomDBSchema (10)
    - ▶ **FunctionalAnnotation (13)**
    - ▶ Organism (529)
    - ▶ PackageType (638)
    - ▶ SequenceAnnotation (2)
  - ▶ ExperimentData (202)

### Packages found under FunctionalAnnotation:

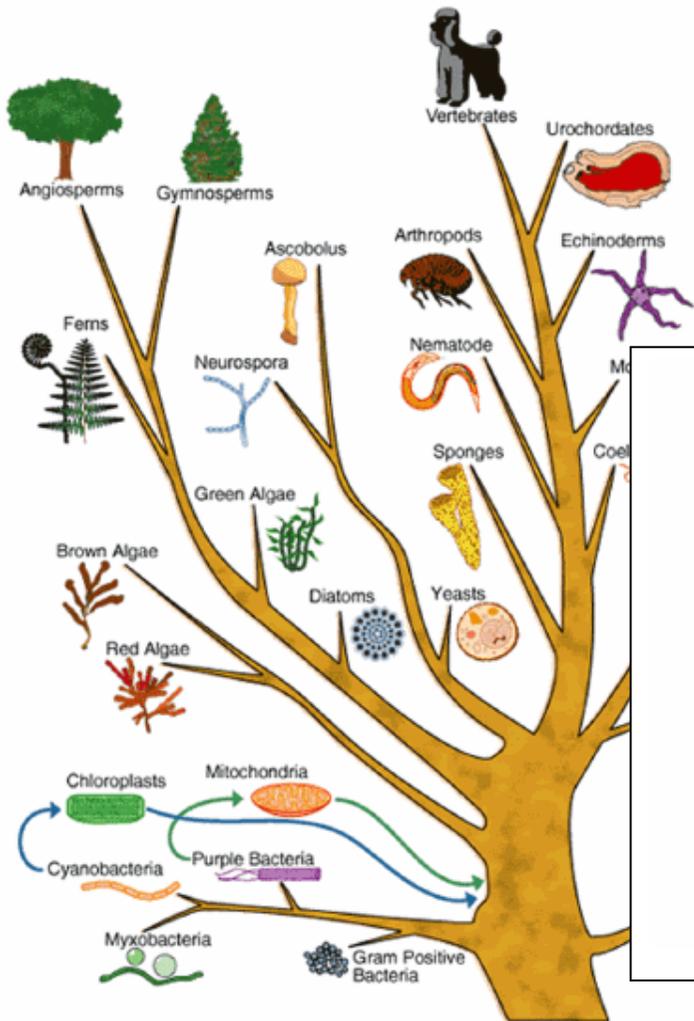
Show  entries Search table:

Package	Maintainer	Title
<a href="#">DO.db</a>	Jiang Li	A set of annotation maps describing the entire Disease Ontology
<a href="#">GO.db</a>	Bioconductor Package Maintainer	A set of annotation maps describing the entire Gene Ontology
<a href="#">humanCHRLOC</a>	Biocore Data Team	A data package containing annotation data for humanCHRLOC
<a href="#">KEGG.db</a>	Bioconductor Package Maintainer	A set of annotation maps for KEGG
<a href="#">MeSH.AOR.db</a>	Koki Tsuyuzaki	A set of annotation maps describing the entire MeSH
<a href="#">MeSH.db</a>	Koki Tsuyuzaki	A set of annotation maps describing the entire MeSH
<a href="#">MeSH.PCR.db</a>	Koki Tsuyuzaki	A set of annotation maps describing the entire MeSH
<a href="#">mirbase.db</a>	James F. Reid	miRBase: the microRNA database
<a href="#">mouseCHRLOC</a>	Biocore Data Team	A data package containing annotation data for mouseCHRLOC
<a href="#">ratCHRLOC</a>	Biocore Data Team	A data package containing annotation data for ratCHRLOC

## Finding function

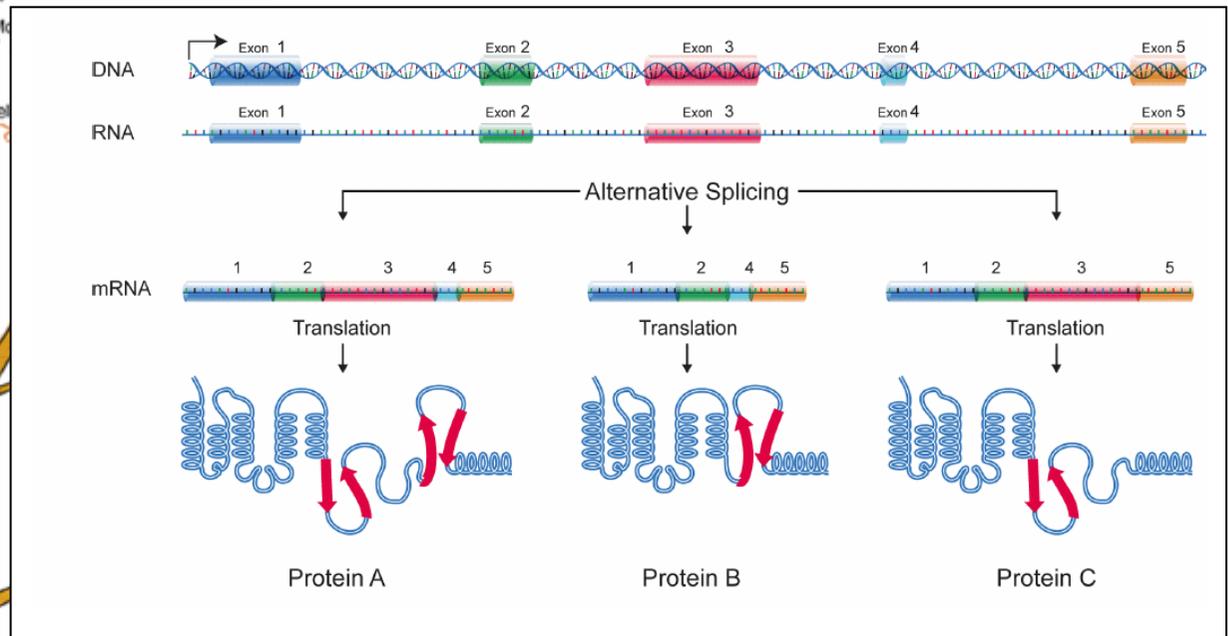
- Once a genomic region of interest has been sequenced, bioinformatic analysis can be used to determine if the sequence is similar to that of a known gene.
  - This is where sequences from model organisms are helpful.
  - For example, let's say we have an unknown human DNA sequence that is associated with the disease cystic fibrosis.
    - A bioinformatic analysis finds a similar sequence from mouse that is associated with a gene that codes for a membrane protein that regulates salt balance.
    - It is a good bet that the human sequence also is part of a gene that codes for a membrane protein that regulates salt balance.
- Links to DNA-seq lectures (sequence comparisons)

# Finding function



Q: "How can something as complicated as a human have only 25 percent more genes than the tiny roundworm *C. elegans*?"

Part of the answer seems to involve **alternative splicing**:



*“The more we find, the more we see, the more we come to learn.  
The more that we explore, the more we shall return.”*

Sir Tim Rice, *Aida*, 2000

## Some References:

- Ziegler A and Van Steen K 2010: IBS short course on “Genome-Wide Association Studies”
- Balding D 2006. A tutorial on statistical methods for population association studies. *Nature Reviews Genetics*, 7, 781-791.
- Kruglyak L 2008. The road to genomewide association studies. *Nature Reviews Genetics* 9: 314-
- Wang et al 2005. Genome-wide association studies: theoretical and practical concerns. *Nature Reviews Genetics* 6: 109-
- Peltonen L and McKusick VA 2001. Dissecting human disease in the postgenomic era. *Science* 291, 1224-1229
- Li 2007. Three lectures on case-control genetic association analysis. *Briefings in bioinformatics* 9: 1-13.
- Rebbeck et al 2004. Assessing the function of genetic variants in candidate gene association studies 5: 589-
- Robinson 2010. Common Disease, Multiple Rare (and Distant) Variants. *PLoS Biology* 8(1): e1000293