

# Lecture 19. Disease Genetics

Michael Schatz

April 10 2019

JHU 600.749: Applied Comparative Genomics



# Preliminary Project Report

---

Assignment Date: April 8, 2019

Due Date: Monday, April 15, 2019 @ 11:59pm

Each team should submit a PDF of your preliminary project proposal (2 to 3 pages) to GradeScope by 11:59pm on Monday April 15.

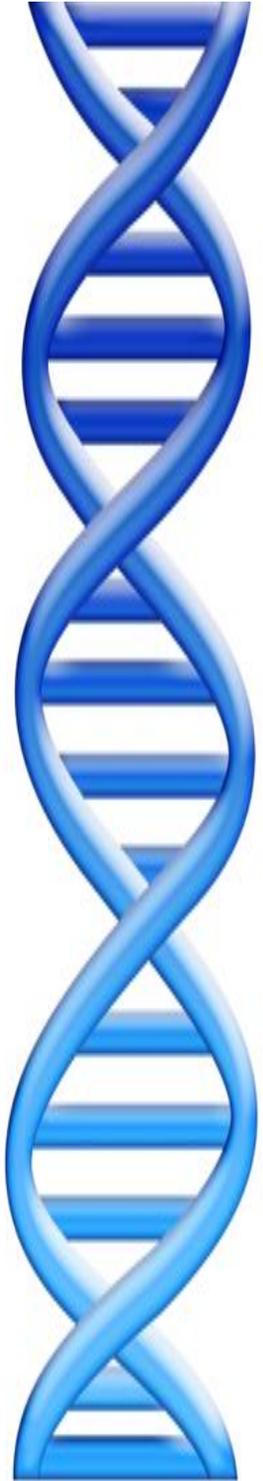
The preliminary report should have at least:

- Title of your project
- List of team members and email addresses
- 1 paragraph abstract summarizing the project
- 1+ paragraph of Introduction
- 1+ paragraph of Methods that you are using
- 1+ paragraph of Results, describing the data evaluated and any any preliminary results
- 1+ paragraph of Discussion (what you have seen or expect to see)
- 1+ figure showing a preliminary result
- 5+ References to relevant papers and data

The preliminary report should use the Bioinformatics style template. Word and LaTeX templates are available at [https://academic.oup.com/bioinformatics/pages/submission\\_online](https://academic.oup.com/bioinformatics/pages/submission_online)

Later, you will present your project in class starting the week of April 24. You will also submit your final written report (5-7 pages) of your project by May 15

Please use Piazza if you have any general questions!

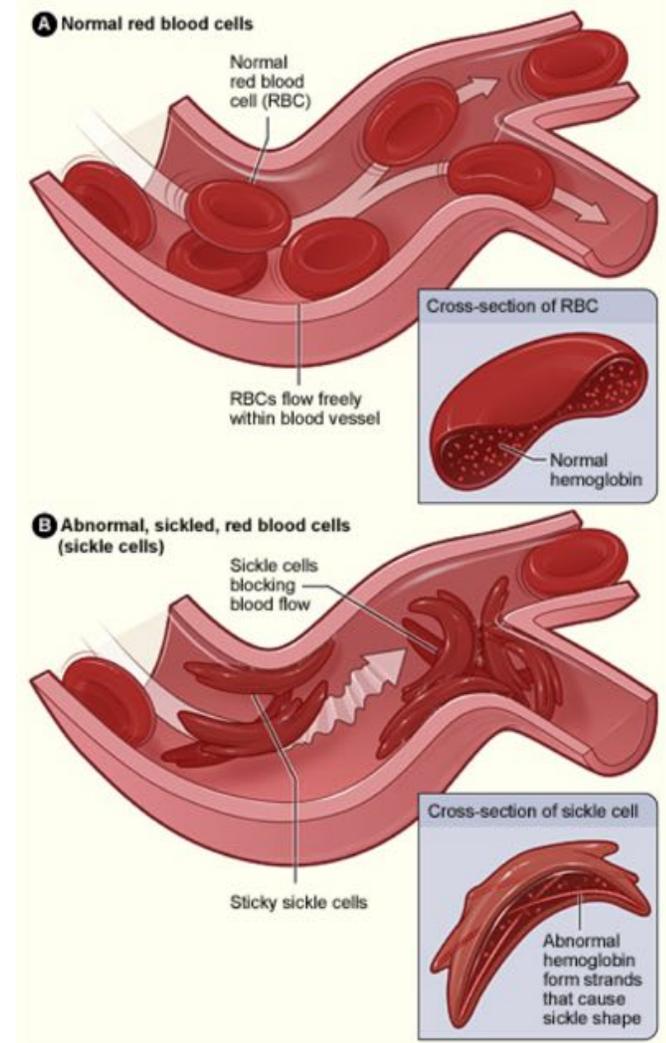


**Part I:**

**Pre-genome Era**

# Sickle Cell Anaemia

- Sickle-cell anaemia (SCA) is an abnormality in the oxygen-carrying protein haemoglobin (hemoglobin S) found in red blood cells. First modern clinical description in 1910s
- **The genetic basis of sickle cell disease is an A-to-T transversion in the sixth codon of the HBB gene.**
- The mutation was actually found in the protein sequence first in the 1950s! Occurs when a person inherits two abnormal copies of the haemoglobin gene, one from each parent. Interestingly, heterozygous patients also incur a resistance to malaria infection, contributing to its prevalence in Africa where malaria infections remain a major disease



**OMIM: SICKLE CELL ANEMIA**

<https://www.omim.org/entry/603903>

# Huntington's Disease

---

## **A polymorphic DNA marker genetically linked to Huntington's disease**

**James F. Gusella<sup>\*</sup>, Nancy S. Wexler<sup>†||</sup>, P. Michael Conneally<sup>†</sup>, Susan L. Naylor<sup>§</sup>,  
Mary Anne Anderson<sup>\*</sup>, Rudolph E. Tanzi<sup>\*</sup>, Paul C. Watkins<sup>\*\*</sup>, Kathleen Ottina<sup>\*</sup>,  
Margaret R. Wallace<sup>‡</sup>, Alan Y. Sakaguchi<sup>§</sup>, Anne B. Young<sup>|</sup>, Ira Shoulson<sup>|</sup>,  
Ernesto Bonilla<sup>|</sup> & Joseph B. Martin<sup>\*</sup>**

<sup>\*</sup> Neurology Department and Genetics Unit, Massachusetts General Hospital and Harvard Medical School, Boston, Massachusetts 02114, USA

<sup>†</sup> Hereditary Disease Foundation, 9701 Wilshire Blvd, Beverly Hills, California 90212, USA

<sup>‡</sup> Department of Medical Genetics, Indiana University Medical Center, Indianapolis, Indiana 46223, USA

<sup>§</sup> Department of Human Genetics, Roswell Park Memorial Institute, Buffalo, New York 14263, USA

<sup>|</sup> Venezuela Collaborative Huntington's Disease Project<sup>\*</sup>

---

*Family studies show that the Huntington's disease gene is linked to a polymorphic DNA marker that maps to human chromosome 4. The chromosomal localization of the Huntington's disease gene is the first step in using recombinant DNA technology to identify the primary genetic defect in this disorder.*

---

# Huntington's Disease

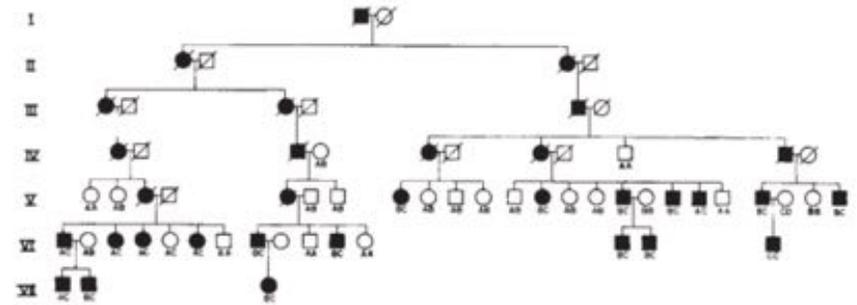
## A polymorphic D to H

James F. Gusella\*, Nan  
Mary Anne Anderson\*,  
Margaret R. Wallace  
Er

\* Neurology Department and Genetics Unit, M  
† Hereditary Disease I  
‡ Department of Medical Ge  
§ Department of Human C  
Ive

Family studies show that the Huntington  
chromosome 4. The chromosomal loca  
DNA technology to identify the primar

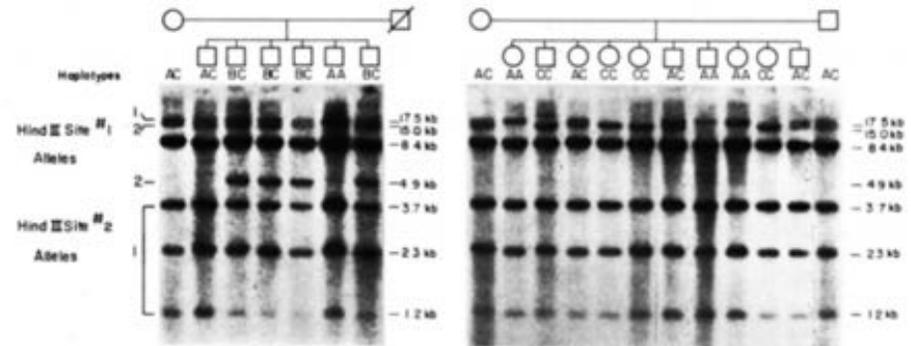
**Fig. 2** Pedigree of the Venezuelan Huntington's disease family. This pedigree represents a small part of a much larger pedigree that will be described in detail elsewhere. Permanent EBV-transformed lymphoblastoid cell lines were established from blood samples of these individuals (unpublished data). DNA prepared from the lymphoblastoid lines will be used to determine the phenotype of each individual at the G8 locus as described in Fig. 3. The data were analysed for linkage to the Huntington's disease gene using the program LIPED<sup>17</sup> with a correction for the late age of onset<sup>2</sup>. Because of the high frequency of the Huntington's disease gene in this population some of the spouses of affected individuals have also descended from identified Huntington's disease gene carriers. In none of these cases, however, was the unaffected individual at significantly greater risk for Huntington's disease than a member of the general population. Although a number of younger at-risk individuals were also analysed as part of this study, for the sake of these family members the data are not shown due to their predictive nature. The data are available upon request if confidentiality can be assured.



**Fig. 3** Hybridization of the G8 Probe to *Hind*III-digested human genomic DNA.

**Methods:** DNA was prepared as described<sup>23</sup> from lymphoblastoid cell lines derived from members of two nuclear families. 5 µg of each DNA was digested to completion with 20 units of *Hind*III in a volume of 30 µl using the buffer recommended by the supplier. The DNAs were fractionated on a 1% horizontal agarose gel in TBE buffer (89 mM Tris, pH 8, 89 mM Na borate, 2 mM Na EDTA) for 18 h. *Hind*III-digested λC1857 DNA was loaded in a separate lane as a size marker.

The gels were stained with ethidium bromide (0.5 µg ml<sup>-1</sup>) for 30 min and the DNA was visualized with UV light. The gels were incubated for 45 min in 1 M NaOH with gentle shaking and for two successive 20 min periods in 1 M Tris, pH 7.6, 1.5 M NaCl. DNA from the gel was transferred in 20×SSC (3 M NaCl, 0.3 M Na citrate) by capillary action to a positively charged nylon membrane. After overnight transfer, agarose clinging to the filters was removed by washing in 3×SSC and the filters were air dried and baked for 2 h under vacuum at 80 °C. Baked filters were prehybridized in 500 ml 6×SSC, 1×Denhardt's solution (0.02% bovine serum albumin, 0.02% polyvinylpyrrolidone, 0.02% Ficoll), 0.3% SDS and 100 µg ml<sup>-1</sup> denatured salmon sperm DNA at 65 °C for 18 h. Prehybridized filters were washed extensively at room temperature in 3×SSC until no evidence of SDS remained. Excess liquid was removed from the filters by blotting on Whatman 3MM paper and damp filters were placed individually in heat-sealable plastic bags. 5 ml of hybridization solution (6×SSC, 1×Denhardt's solution, 0.1% SDS, 100 µg ml<sup>-1</sup> denatured salmon sperm DNA) containing approximately 5×10<sup>8</sup> c.p.m. of nick-translated G8 DNA (specific activity ~2×10<sup>8</sup> c.p.m. µg<sup>-1</sup>)<sup>24</sup> was added to each bag which was then sealed and placed at 65 °C for 24–48 h. Filters were removed from the bags and washed at 65 °C for 30 min each in 3×SSC, 2×SSC, 1×SSC and 0.3×SSC. The filters were dried and exposed to X-ray film (Kodak XR-5) at -70 °C with a Dupont Cronex intensifying screen for 1 to 4 days. The haplotypes observed in each individual were determined from the alleles seen for each *Hind*III RFLP (site 1 and 2) as explained in Fig. 4.



# Huntington's Disease

Cell, Vol. 72, 971-983, March 26, 1993, Copyright © 1993 by Cell Press

## A Novel Gene Containing a Trinucleotide Repeat That Is Expanded and Unstable on Huntington's Disease Chromosomes

The Huntington's Disease Collaborative Research Group\*

### Summary

The Huntington's disease (HD) gene has been mapped in 4p16.3 but has eluded identification. We have used haplotype analysis of linkage disequilibrium to spotlight a small segment of 4p16.3 as the likely location of the defect. A new gene, IT15, isolated using cloned trapped exons from the target area contains a polymorphic trinucleotide repeat that is expanded and unstable on HD chromosomes. A (CAG)<sub>n</sub> repeat longer than the normal range was observed on HD chromosomes from all 75 disease families examined, comprising a variety of ethnic backgrounds and 4p16.3 haplotypes. The (CAG)<sub>n</sub> repeat appears to be located within the coding sequence of a predicted ~348 kd protein that is widely expressed but unrelated to any known gene. Thus, the HD mutation involves an unstable DNA segment, similar to those described in fragile X syndrome, spino-bulbar muscular atrophy, and myotonic dystrophy, acting in the context of a novel 4p16.3 gene to produce a dominant phenotype.

### Introduction

Huntington's disease (HD) is a progressive neurodegenerative disorder characterized by motor disturbance, cognitive loss, and psychiatric manifestations (Martin and Gusella, 1986). It is inherited in an autosomal dominant fashion and affects ~1 in 10,000 individuals in most populations of European origin (Harper et al., 1991). The hallmark of HD is a distinctive choreic movement disorder that typically has a subtle, insidious onset in the fourth to fifth decade of life and gradually worsens over a course of 10 to 20 years until death. Occasionally, HD is expressed in juveniles, typically manifesting with more severe symptoms including rigidity and a more rapid course. Juvenile onset of HD is associated with a preponderance of paternal transmission of the disease allele. The neuropathology of HD also displays a distinctive pattern, with selective loss of neurons that is most severe in the caudate and putamen. The biochemical basis for neuronal death in HD has not yet been explained, and there is consequently no treatment effective in delaying or preventing the onset and progression of this devastating disorder.

The genetic defect causing HD was assigned to chromosome 4 in 1983 in one of the first successful linkage analyses using polymorphic DNA markers in humans (Gusella

# Huntington's Disease

Cell, Vol. 72, 971-983, March 26, 1993, Copyright © 1993 by Cell Press

## A Novel Gene Containing a Trinucleotide Repeat That Is Expanded and Unstable on Huntington's Disease Chromosomes

The Huntington's Disease Collaborative Research Group\*

### Summary

The Huntington's disease (HD) gene has been mapped in 4p16.3 but has eluded identification. We have used haplotype analysis of linkage disequilibrium to spotlight a small segment of 4p16.3 as the likely location of the defect. A new gene, IT15, isolated using cloned trapped exons from the target area contains a polymorphic trinucleotide repeat that is expanded and unstable on HD chromosomes. A (CAG)<sub>n</sub> repeat longer than the normal range was observed on HD chromosomes from all 75 disease families examined, comprising a variety of ethnic backgrounds and 4p16.3 haplotypes. The (CAG)<sub>n</sub> repeat appears to be located within the coding sequence of a predicted ~348 kd protein that is widely expressed but unrelated to any known gene. Thus, the HD mutation involves an unstable DNA segment, similar to those described in fragile X syndrome, spino-bulbar muscular atrophy, and myotonic dystrophy, acting in the context of a novel 4p16.3 gene to produce a dominant phenotype.

### Introduction

Huntington's disease (HD) is a progressive disorder characterized by motor, cognitive, and psychiatric manifestations (Huntington, 1986). It is inherited in an autosomal dominant fashion and affects ~1 in 10,000 individuals of European origin (Harper et al., 1986). A distinctive choreic movement disorder is a hallmark of HD that typically has a subtle, insidious onset in the fifth decade of life and gradually worsens over a period of 10 to 20 years until death. Occurrence in juveniles, typically manifested by severe symptoms including rigidity and chorea, is associated with a pattern of paternal transmission of the disease. The pathology of HD also displays a distinctive selective loss of neurons that is most prominent in the caudate and putamen. The biochemical basis of the disease in HD has not yet been explained, and frequently no treatment effective in delaying the onset and progression of this disease is available.

The genetic defect causing HD was first identified in 1983 in one of the first successful gene clones using polymorphic DNA markers

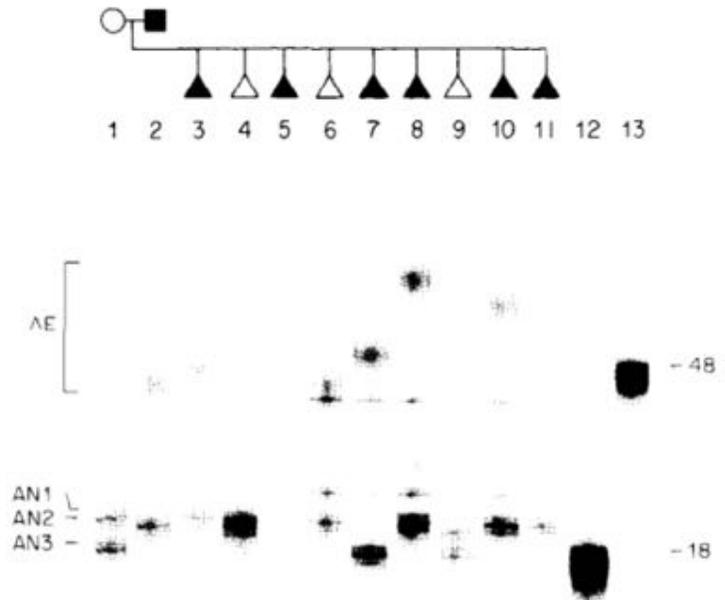


Figure 6. PCR Analysis of the (CAG)<sub>n</sub> Repeat in a Venezuelan HD Sibship with Some Offspring Displaying Juvenile Onset  
Results of PCR analysis of a sibship in the Venezuelan HD pedigree are shown. Affected individuals are represented by closed symbols. Progeny are shown as triangles, and the birth order of some individuals has been changed for confidentiality. AN1, AN2, and AN3 mark the positions of the allelic products from normal chromosomes. AE marks the range of PCR products from the HD chromosome. The intensity of background constant bands, which represent a useful reference for comparison of the above PCR products, varies with slight differences in PCR conditions. The PCR products from cosmids L191F1 and GUS72-2130 are loaded in lanes 12 and 13 and have 18 and 48 CAG repeats, respectively.

# Human disease genes

Gerardo Jimenez-Sanchez\*, Barton Childs\* & David Valle\*†

\* Department of Pediatrics, McKusick-Nathans Institute of Genetic Medicine, and † Howard Hughes Medical Institute, Johns Hopkins University School of Medicine, Baltimore, Maryland 21205, USA

**The complete human genome sequence will facilitate the identification of all genes that contribute to disease. We propose that the functional classification of disease genes and their products will reveal general principles of human disease. We have determined functional categories for nearly 1,000 documented disease genes, and found striking correlations between the function of the gene product and features of disease, such as age of onset and mode of inheritance. As knowledge of disease genes grows, including those contributing to complex traits, more sophisticated analyses will be possible; their results will yield a deeper understanding of disease and an enhanced integration of medicine with biology.**

**T**o test the proposal that classifying disease genes and their products according to function will provide general insight into disease processes<sup>1,2</sup>, we have compiled and classified a list of disease genes. To assemble the list, we began with 269 genes identified in a survey of the 7th edition of *Metabolic and Molecular Bases of Inherited Disease*<sup>2</sup>. We then searched the 'morbid map' and allelic variants listed in the *Online Mendelian Inheritance in Man*<sup>3</sup> (OMIM), an online resource documenting human diseases and their associated genes

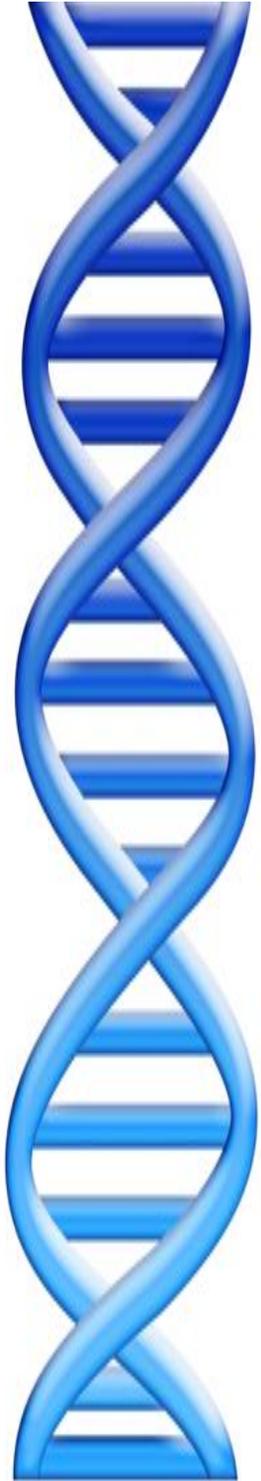
([www.ncbi.nlm.nih.gov](http://www.ncbi.nlm.nih.gov)), and increased the total disease gene set to 923. This sample included genes that cause monogenic disease (97% of the sample) and genes that increase susceptibility for complex traits. We excluded genes associated only with somatic genetic disease (such as non-inherited forms of cancer) or the mitochondrial genome.

## Functional classification

We categorized each disease gene according to the function of its

## Human disease genes

Jimenez-Sanchez, G., Childs, B. & Valle, D. (2001) *Nature* 409, 853–855



**Part 2:**

**Post-genome  
Inherited Diseases**

“Genome-wide linkage analysis has also been carried out for many common diseases and quantitative traits, for which the aforementioned characteristics of Mendelian diseases might not apply. In some cases, genomic regions that show significant linkage to the disease have been identified, leading to the discovery of variants that contribute to susceptibility to diseases such as inflammatory bowel disease (IBD), schizophrenia and type 1 diabetes.

***However, for most common diseases, linkage analysis has achieved only limited success, and the genes discovered usually explain only a small fraction of the overall heritability of the disease.”***

***Genome-wide association studies for common diseases and complex traits***

Hirschhorn and Daly (2005) Nature Review Genetics

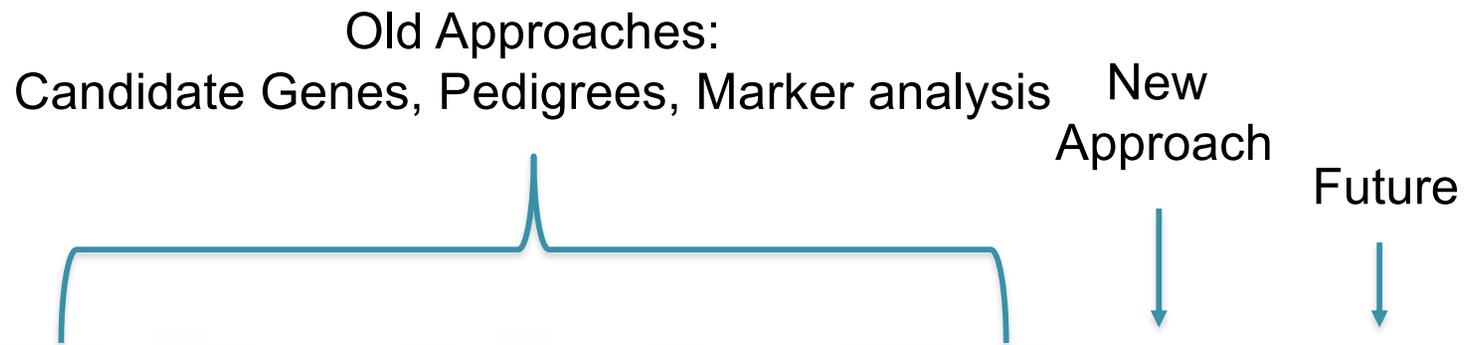


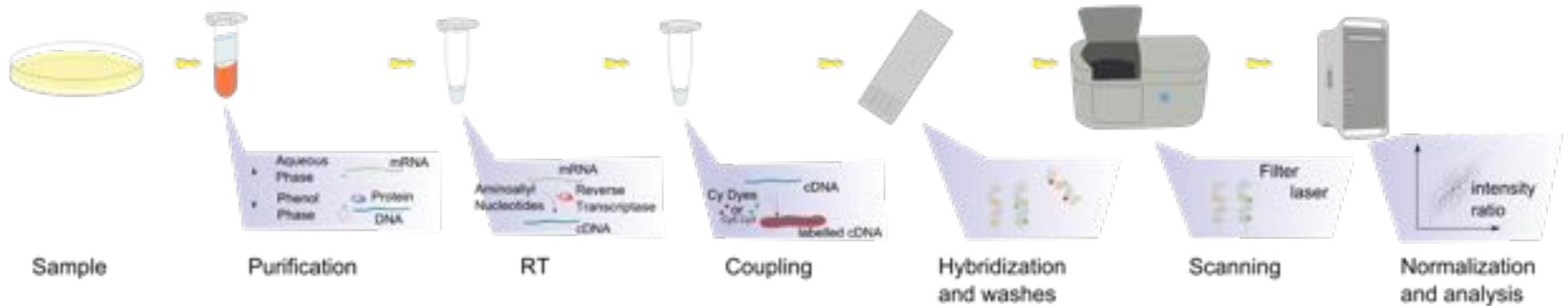
Table 1 | **Approaches to identifying variants underlying complex traits and common diseases**

Potential advantages	Association*	Resequencing*	Linkage <sup>†</sup>	Admixture <sup>†</sup>	Missense SNPs <sup>‡</sup>	Association <sup>†</sup>	Resequencing <sup>†</sup>
No prior information regarding gene function required	-	-	+	+	+	+	+
Localization to small genomic region	+	+	-	-	+	+	+
Inexpensive	+	-	+	+	+/-	-	Prohibitive
Families not required	+	+	-	+	+	+	+
No assumptions necessary regarding type of variant involved	+	-	+	+	-	+	+
Not susceptible to effects of stratification <sup>§</sup>	-/+	-/+	+	+	-/+	-/+	-/+
No requirement for variation of allele frequency among populations	+	+	+	-	+	+	+
Sufficient power to detect common alleles (MAFs>5%) of modest effect	+	-	-/+	+	+	+	+
Ability to detect rare alleles (MAFs<1%)	-	+	+	-	-	-	+
Reasonable track record for common diseases	+	-/+	+/-	N/A	N/A	N/A	N/A
Tools for analysis available	+	+	+	+	+	+/-	-

\*Candidate-gene studies. <sup>†</sup>Genome-wide studies. <sup>‡</sup>Association and resequencing studies are immune to stratification if they use family-based designs. Symbols indicate whether the potential advantage in the left column applies completely (+), partially (+/-), weakly (-/+) or not at all (-). MAF, minor allele frequency; N/A, not yet attempted.

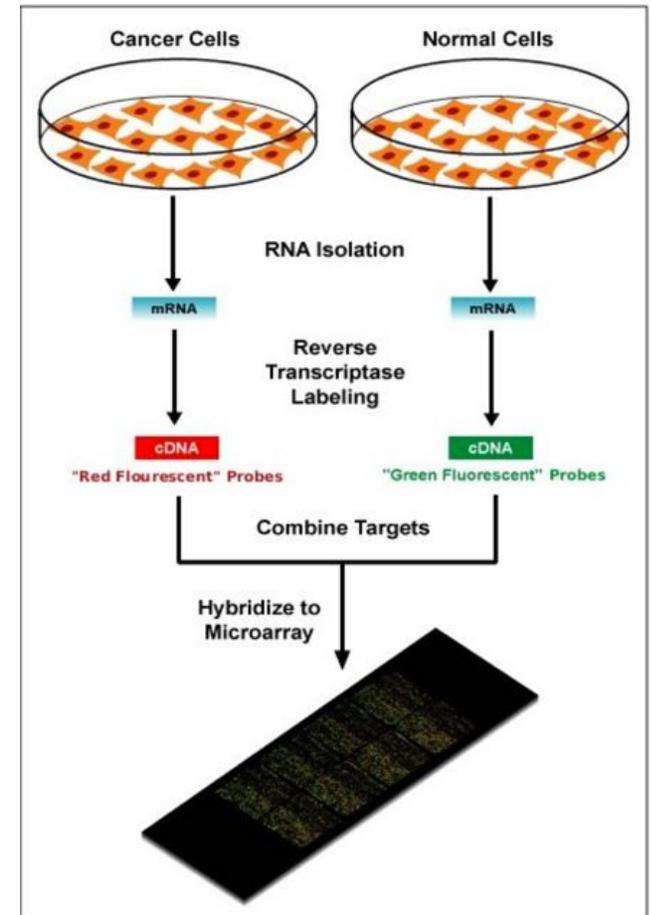
## **Genome-wide association studies for common diseases and complex traits**

Hirschhorn and Daly (2005) Nature Review Genetics



A DNA microarray is a collection of microscopic DNA “spots” attached to a solid surface.

- DNA microarrays can measure the expression levels of large numbers of genes simultaneously or to genotype multiple regions of a genome.
- Each DNA spot contains picomoles (10<sup>-12</sup> moles) of a specific DNA sequence, known as probes (or reporters or oligos).
- Very cost effective (~\$10) for millions of probes at once



AFFYMETRIX®

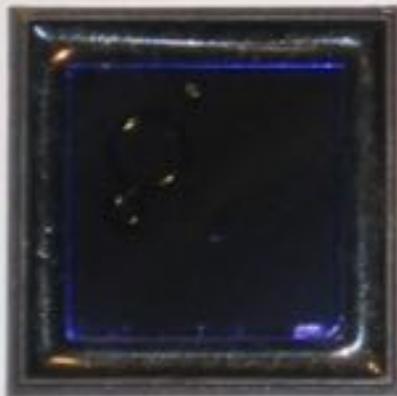


@52001900552056050706401000885092

GeneChip®

J  
2

Human Genome  
U133 Plus 2.0



P/N: 520019

Lot #: 4010008

Exp. Date: 05/07/06

For Research Use Only

AFFYMETRIX®

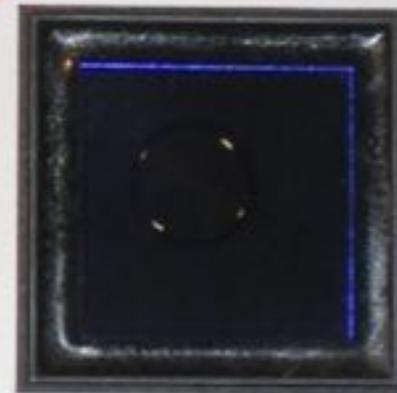


@52002900588841112406401511225336

GeneChip®

1  
5

Mouse Genome  
430 2.0 Array



P/N: 520029

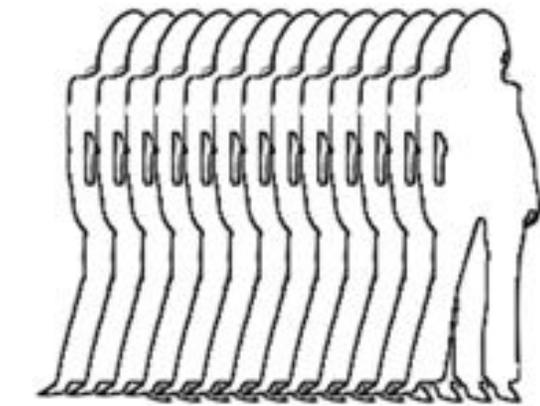
Lot #: 4015112

Exp. Date: 11/24/06

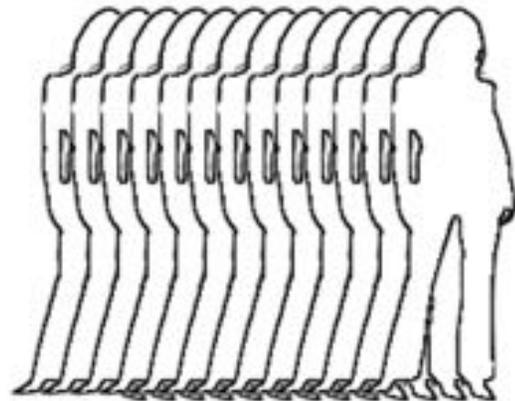
For Research Use Only



# Genome Wide Association (GWAS)



GC CC GG GC CC GC GC  
GG CC GC GG GC GG



GC CC GC GC GG CC CC  
CC GC GC GG GC GG

SNP1

**Cases**

Count of G:  
2104 of 4000

Frequency of G:  
52.6%

**Controls**

Count of G:  
2676 of 6000

Frequency of G:  
44.6%

SNP2

**Cases**

Count of G:  
1648 of 4000

Frequency of G:  
41.2%

**Controls**

Count of G:  
2532 of 6000

Frequency of G:  
42.2%

SNP...

*Repeat for all  
SNPs*

Are these significant  
differences in frequencies?

# Pearson's Chi-squared test

The value of the test-statistic is

$$\chi^2 = \sum_{i=1}^n \frac{(O_i - E_i)^2}{E_i} = N \sum_{i=1}^n \frac{(O_i/N - p_i)^2}{p_i}$$

where

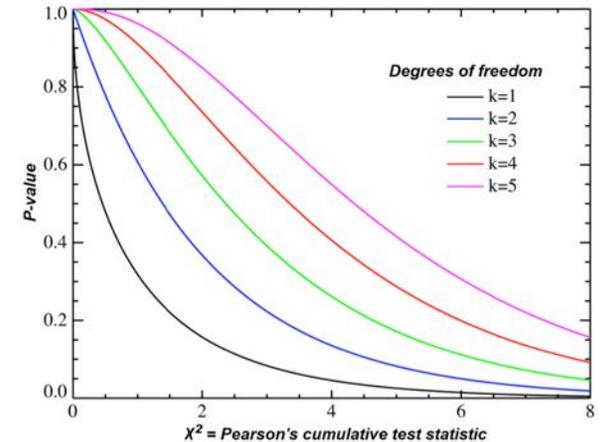
$\chi^2$  = Pearson's cumulative test statistic, which asymptotically approaches a  $\chi^2$  distribution.

$O_i$  = the number of observations of type  $i$ .

$N$  = total number of observations

$E_i = Np_i$  = the expected (theoretical) frequency of type  $i$ , asserted by the null hypothesis that the fraction of type  $i$  in the population is  $p_i$

$n$  = the number of cells in the table.



$$P(\chi_P^2(\{p_i\}) > T) \sim C \int_{\sum_{i=1}^{m-1} y_i^2 > T} \left\{ \prod_{i=1}^{m-1} dy_i \right\} \prod_{i=1}^{m-1} \exp \left[ -\frac{1}{2} \left( \sum_{i=1}^{m-1} y_i^2 \right) \right]$$

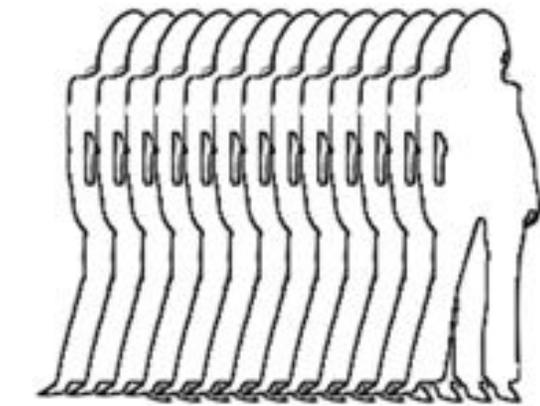
	has G	Not G	Marginal Row Totals
<b>Cases</b>	2104 (1912) [19.28]	1896 (2088) [17.66]	4000
<b>Controls</b>	2676 (2868) [12.85]	3324 (3132) [11.77]	6000
<b>Marginal Column Totals</b>	4780	5220	10000 (Grand Total)

Cases/hasG expected:  $4000 * (4780/10000) = 1912$  expected

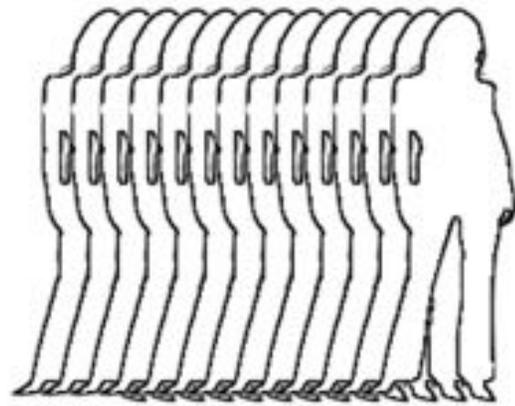
Cases/hasG squared deviation:  $(2104 - 1912)^2 / 1912 = 19.28$  deviation

The chi-square statistic is  $19.28+17.66+12.85+11.77 = 61.56$ . The p-value is  $5e-15$

# Genome Wide Association (GWAS)



GC CC GG GC CC GC GC  
GG CC GC GG GC GG



GC CC GC GC GG CC CC  
CC GC GC GG GC GG

SNP1

**Cases**

Count of G:  
2104 of 4000

Frequency of G:  
52.6%

**Controls**

Count of G:  
2676 of 6000

Frequency of G:  
44.6%

**P-value:**

$5.0 \cdot 10^{-15}$

SNP2

**Cases**

Count of G:  
1648 of 4000

Frequency of G:  
41.2%

**Controls**

Count of G:  
2532 of 6000

Frequency of G:  
42.2%

**P-value:**

0.33

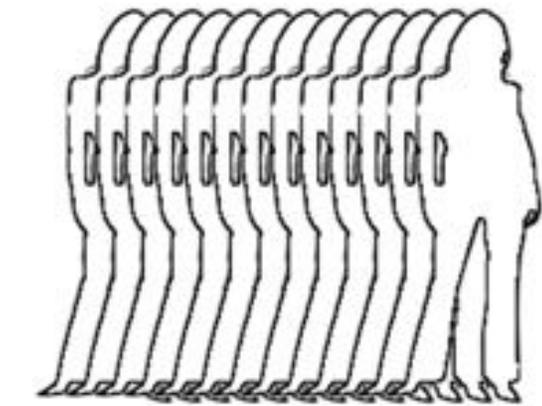
SNP...

*Repeat for all  
SNPs*

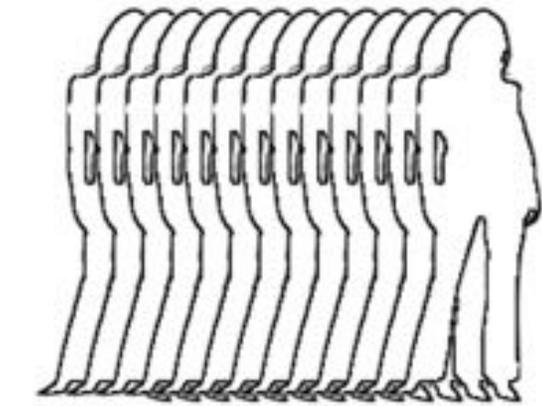
Chi-squared or  
similar test



# Genome Wide Association (GWAS)



GC CC GG GC CC GC GC  
GG CC GC GG GC GG



GC CC GC GC GG CC CC  
CC GC GC GG GC GG

SNP1

**Cases**

Count of G:  
2104 of 4000

Frequency of G:  
52.6%

**Controls**

Count of G:  
2676 of 6000

Frequency of G:  
44.6%

**P-value:**

$5.0 \cdot 10^{-15}$

SNP2

**Cases**

Count of G:  
1648 of 4000

Frequency of G:  
41.2%

**Controls**

Count of G:  
2532 of 6000

Frequency of G:  
42.2%

**P-value:**

0.33

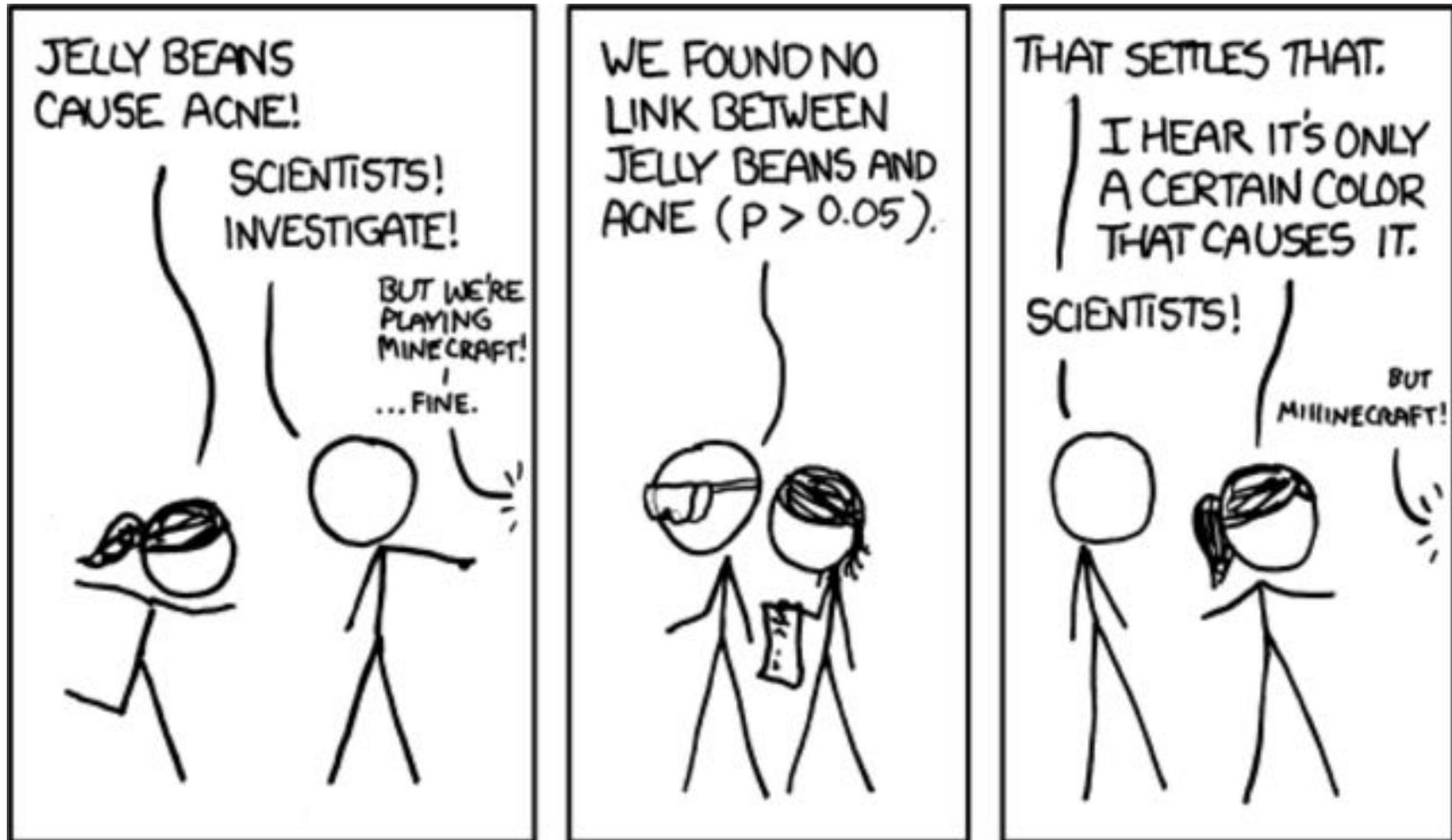
SNP...

*Repeat for all  
SNPs*

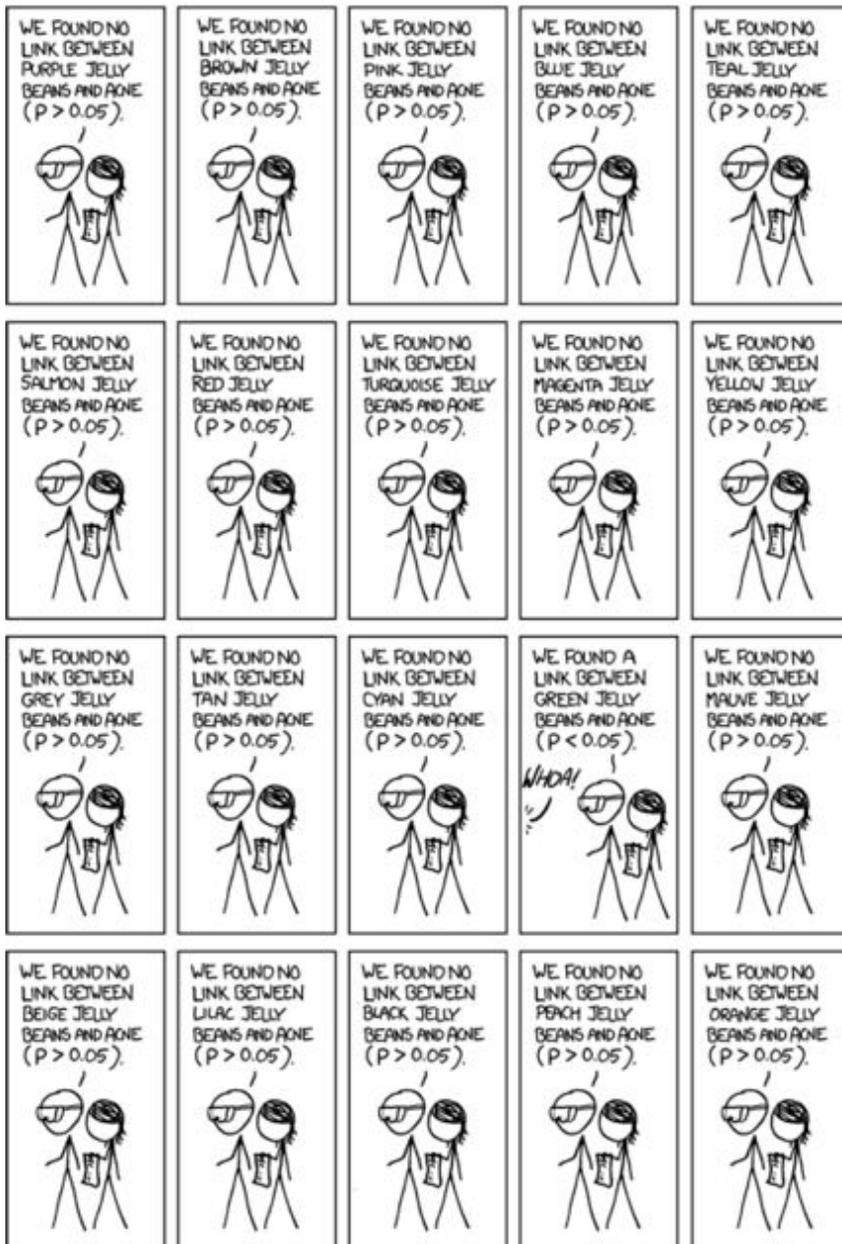
With a (much) larger population, this might be a significant difference in rate:  
 $25320/60000 \Rightarrow$   
 $p = 5e-7$

Chi-squared or similar test

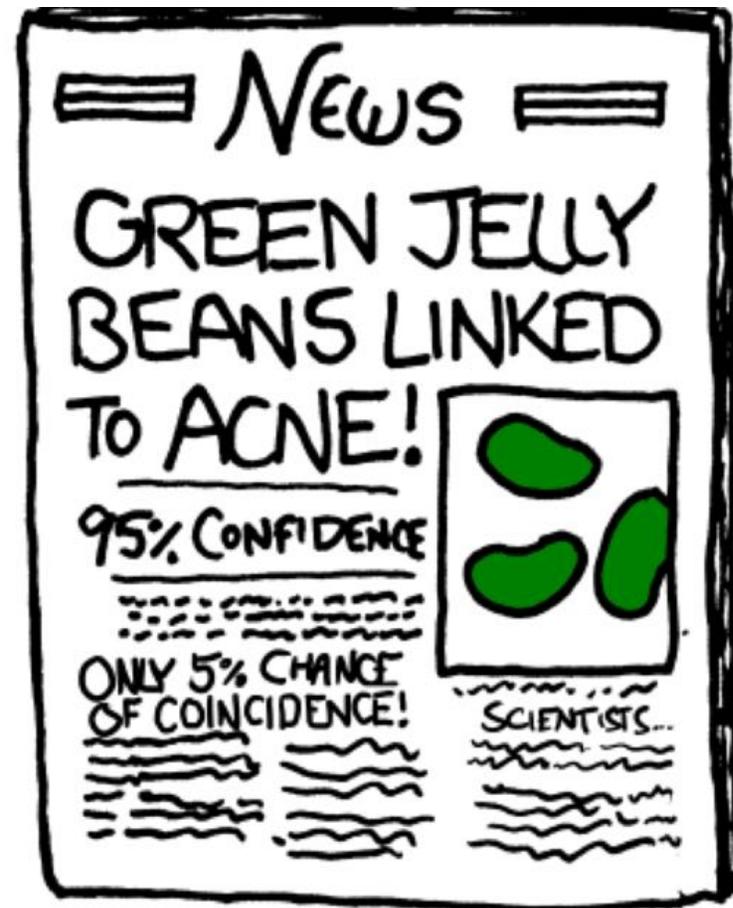
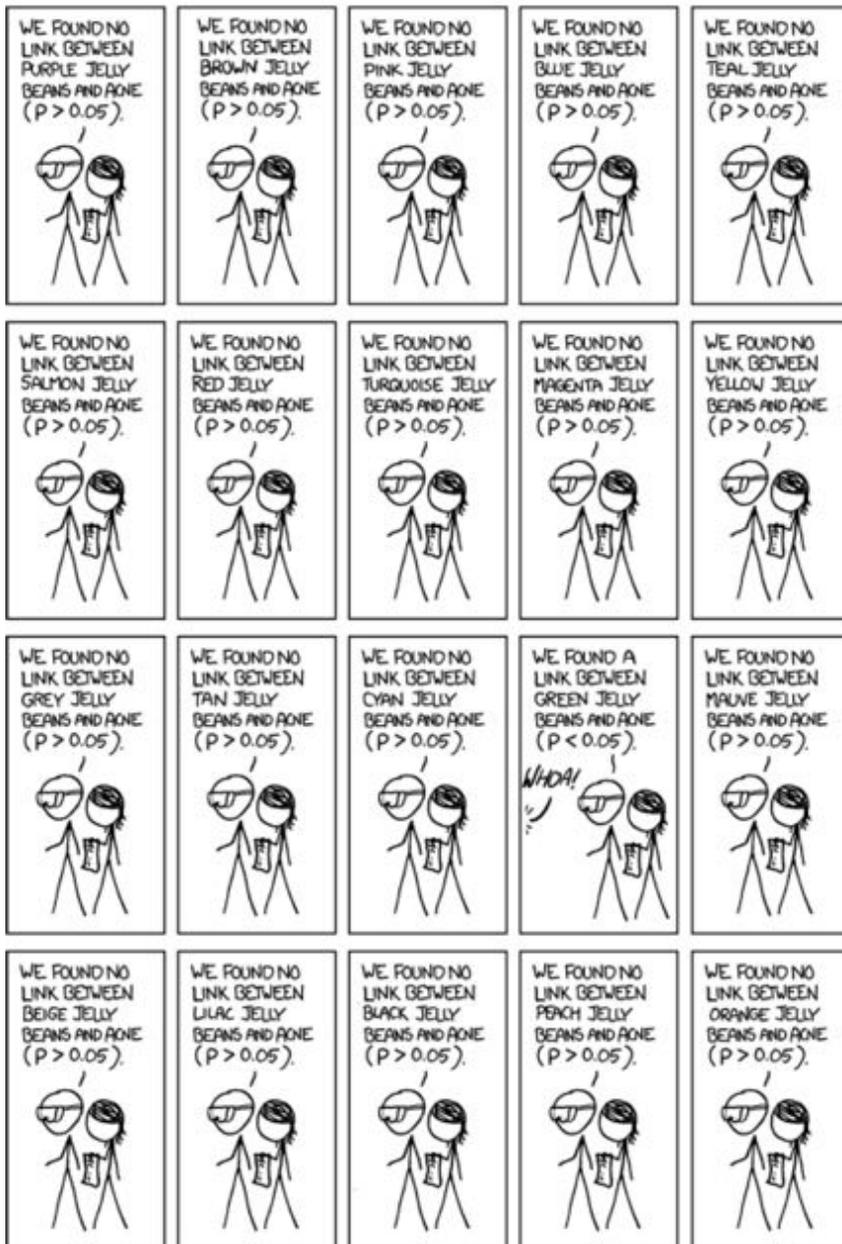
# The curse of multiple testing



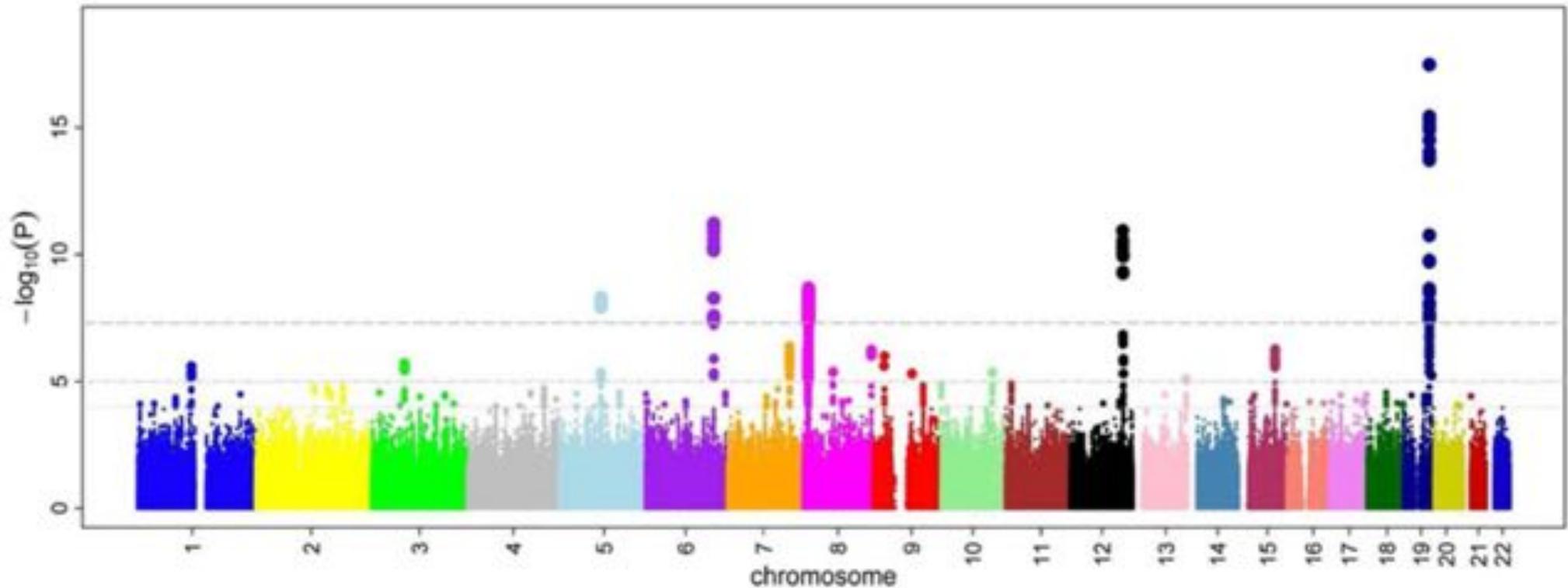
# The curse of multiple testing



# The curse of multiple testing



# Manhattan Plot



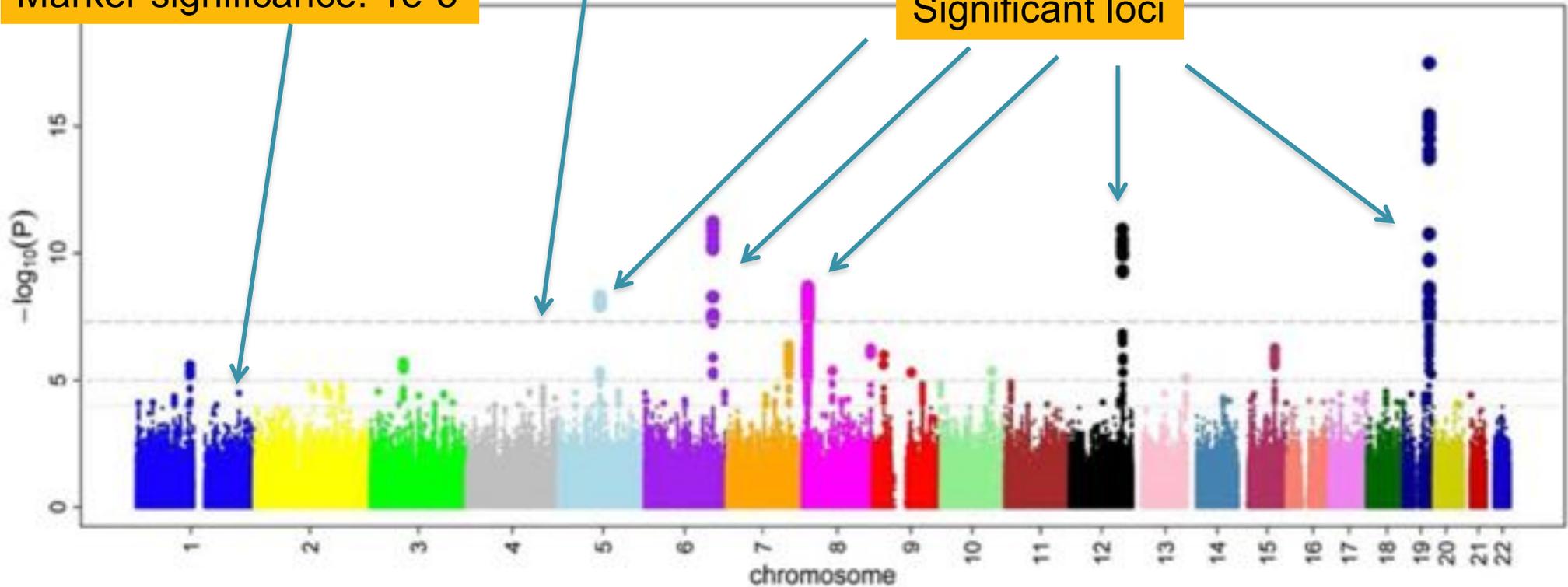
***Four Novel Loci (19q13, 6q24, 12q24, and 5q14) Influence the Microcirculation In Vivo***  
Ikram et al (2010) PLOS Genetics. doi: 10.1371/journal.pgen.1001184

# Manhattan Plot

Genome-wide significance:  $5e-8$

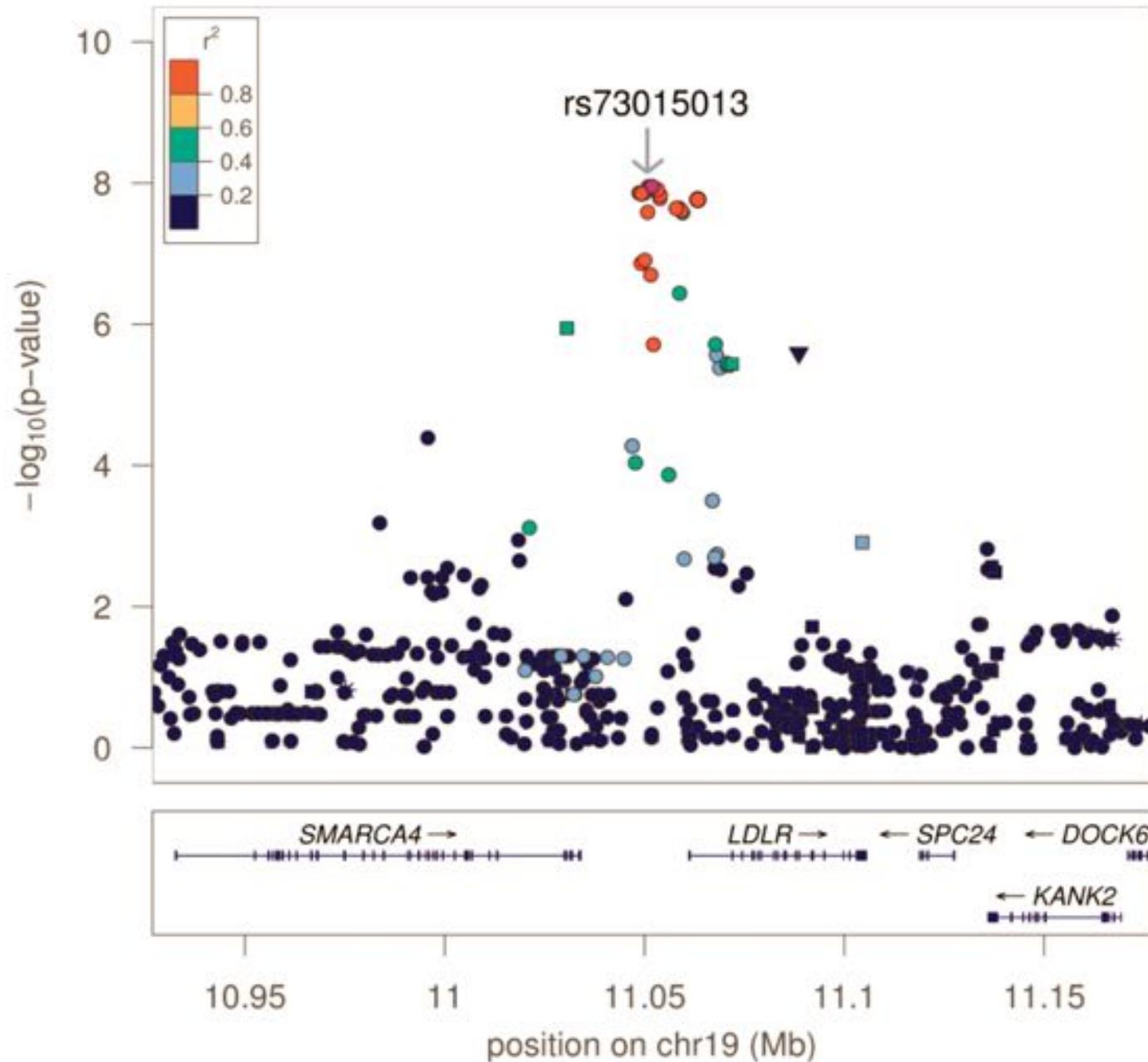
Marker significance:  $1e-5$

Significant loci

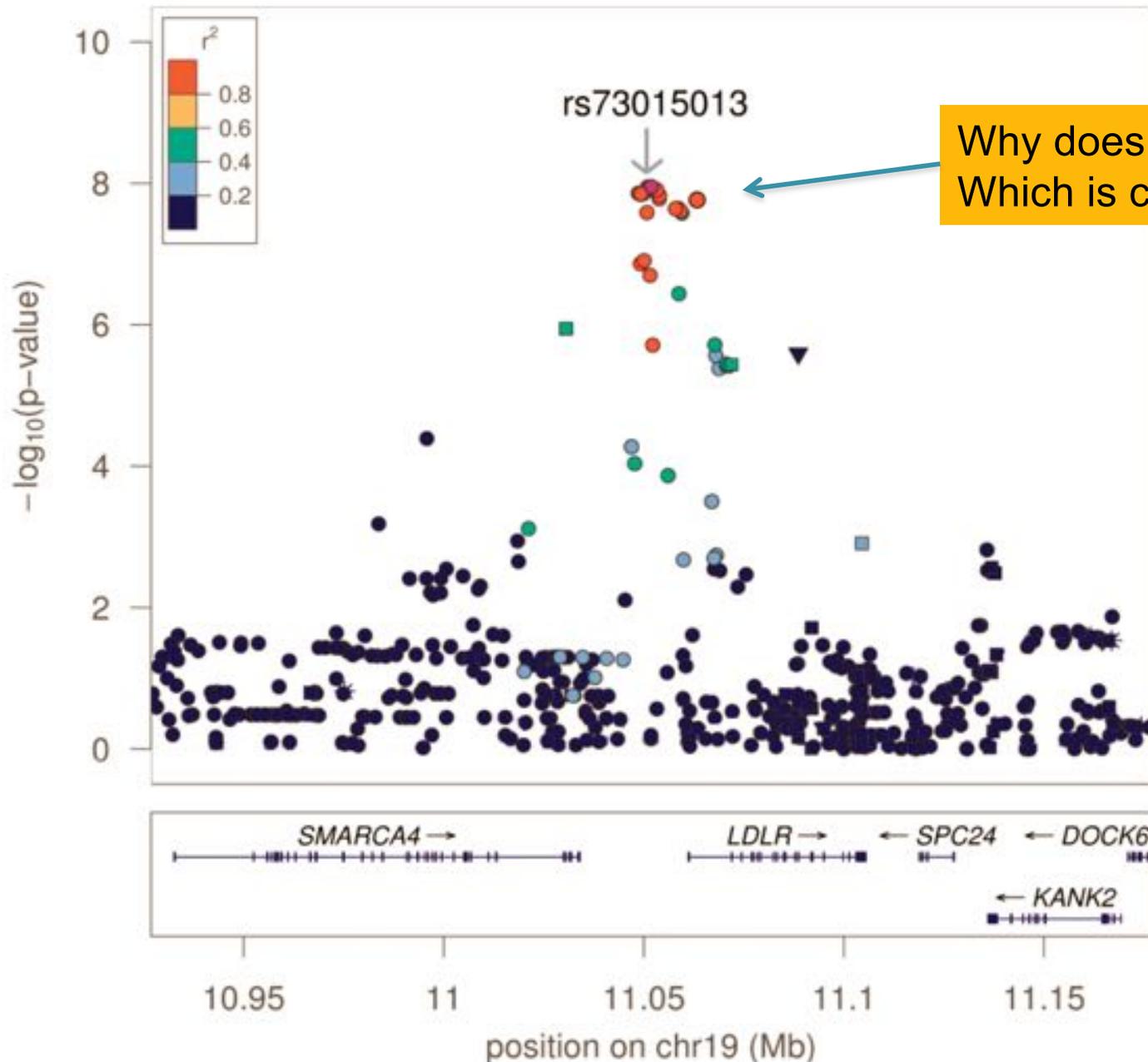


***Four Novel Loci (19q13, 6q24, 12q24, and 5q14) Influence the Microcirculation In Vivo***  
Ikram et al (2010) PLOS Genetics. doi: 10.1371/journal.pgen.1001184

# Regional Association Plot



# Regional Association Plot



# First published GWAS

## Complement Factor H Polymorphism in Age-Related Macular Degeneration

Robert J. Klein,<sup>1</sup> Caroline Zeiss,<sup>2\*</sup> Emily Y. Chew,<sup>3\*</sup>  
Jen-Yue Tsai,<sup>4\*</sup> Richard S. Sackler,<sup>1</sup> Chad Haynes,<sup>1</sup>  
Alice K. Henning,<sup>5</sup> John Paul SanGiovanni,<sup>3</sup> Shrikant M. Mane,<sup>6</sup>  
Susan T. Mayne,<sup>7</sup> Michael B. Bracken,<sup>7</sup> Frederick L. Ferris,<sup>3</sup>  
Jurg Ott,<sup>1</sup> Colin Barnstable,<sup>2</sup> Josephine Hoh<sup>7†</sup>

Age-related macular degeneration (AMD) is a major cause of blindness in the elderly. We report a genome-wide screen of 96 cases and 50 controls for polymorphisms associated with AMD. Among 116,204 single-nucleotide polymorphisms genotyped, an intronic and common variant in the complement factor H gene (*CFH*) is strongly associated with AMD (nominal *P* value  $<10^{-7}$ ). In individuals homozygous for the risk allele, the likelihood of AMD is increased by a factor of 7.4 (95% confidence interval 2.9 to 19). Resequencing revealed a polymorphism in linkage disequilibrium with the risk allele representing a tyrosine-histidine change at amino acid 402. This polymorphism is in a region of *CFH* that binds heparin and C-reactive protein. The *CFH* gene is located on chromosome 1 in a region repeatedly linked to AMD in family-based studies.

Age-related macular degeneration (AMD) is the leading cause of blindness in the developed world. Its incidence is increasing as the elderly population expands (1). AMD is characterized by progressive destruction of the retina's central region (macula), causing central field visual loss (2). A key feature of AMD is the formation of extracellular deposits called drusen concentrated in and around the macula behind the retina between the retinal pigment epithelium (RPE) and the choroid. To date, no therapy for this disease has proven to be broadly effective. Several risk factors have been linked to AMD, including age, smoking, and family history (3). Candidate-gene studies

have not found any genetic differences that can account for a large proportion of the overall prevalence (2). Family-based whole-genome linkage scans have identified chromosomal regions that show evidence of linkage to AMD (4–8), but the linkage areas have not been resolved to any causative mutations.

Like many other chronic diseases, AMD is caused by a combination of genetic and environmental risk factors. Linkage studies are not as powerful as association studies for the identification of genes contributing to the risk for common, complex diseases (9). However, linkage studies have the advantage of searching the whole genome in an unbiased manner

without presupposing the involvement of particular genes. Searching the whole genome in an association study requires typing 100,000 or more single-nucleotide polymorphisms (SNPs) (10). Because of these technical demands, only one whole-genome association study, on susceptibility to myocardial infarction, has been published to date (11).

**Study design.** We report a whole-genome case-control association study for genes involved in AMD. To maximize the chance of success, we chose clearly defined phenotypes for cases and controls. Case individuals exhibited at least some large drusen in a quantitative photographic assessment combined with evidence of sight-threatening AMD (geographic atrophy or neovascular AMD). Control individuals had either no or only a few small drusen. We analyzed our data using a statistically conservative approach to correct for the large number of SNPs tested, thereby guaranteeing that the probability of a false positive is no greater than our reported *P* values.

We used a subset of individuals who participated in the Age-Related Eye Disease Study (AREDS) (12). From the AREDS

<sup>1</sup>Laboratory of Statistical Genetics, Rockefeller University, 1230 York Avenue, New York, NY 10021, USA. <sup>2</sup>Department of Ophthalmology and Visual Science, Yale University School of Medicine, 330 Cedar Street, New Haven, CT 06520, USA. <sup>3</sup>National Eye Institute, Building 10, CRC, 10 Center Drive, Bethesda, MD 20892–1204, USA. <sup>4</sup>Biological Imaging Core, National Eye Institute, 9000 Rockville Pike, Bethesda, MD 20892, USA. <sup>5</sup>The EMMES Corporation, 401 North Washington Street, Suite 700, Rockville MD 20850, USA. <sup>6</sup>W. M. Keck Facility, Yale University, 300 George Street, Suite 201, New Haven, CT 06511, USA. <sup>7</sup>Department of Epidemiology and Public Health, Yale University School of Medicine, 60 College Street, New Haven, CT 06520, USA.

\*These authors contributed equally to this work.  
†To whom correspondence should be addressed.  
E-mail: josephine.hoh@yale.edu

# First published GWAS

## Complement Factor H Polymorphism in Age-Related Macular Degeneration

Robert J. Klein,<sup>1</sup> Caroline Zeiss,<sup>2\*</sup> Emily Y. Chew,<sup>3\*</sup> Jen-Yue Tsai,<sup>4\*</sup> Richard S. Sackler,<sup>1</sup> Chad Haynes,<sup>1</sup> Alice K. Henning,<sup>5</sup> John Paul SanGiovanni,<sup>3</sup> Shrikant M. Mane,<sup>6</sup> Susan T. Mayne,<sup>7</sup> Michael B. Bracken,<sup>7</sup> Frederick L. Ferris,<sup>3</sup> Jurg Ott,<sup>1</sup> Colin Barnstable,<sup>2</sup> Josephine Hoh<sup>7†</sup>

Age-related macular degeneration (AMD) is a major cause of blindness in the elderly. We report a genome-wide screen of 96 cases and 50 controls for polymorphisms associated with AMD. Among 116,204 single-nucleotide polymorphisms genotyped, an intronic and common variant in the complement factor H gene (*CFH*) is strongly associated with AMD (nominal  $P$  value  $<10^{-7}$ ). In individuals homozygous for the risk allele, the likelihood of AMD is increased by a factor of 7.4 (95% confidence interval 2.9 to 19). Resequencing revealed a polymorphism in linkage disequilibrium with the risk allele representing a tyrosine-histidine change at amino acid 402. This polymorphism is in a region of *CFH* that binds heparin and C-reactive protein. The *CFH* gene is located on chromosome 1 in a region repeatedly linked to AMD in family-based studies.

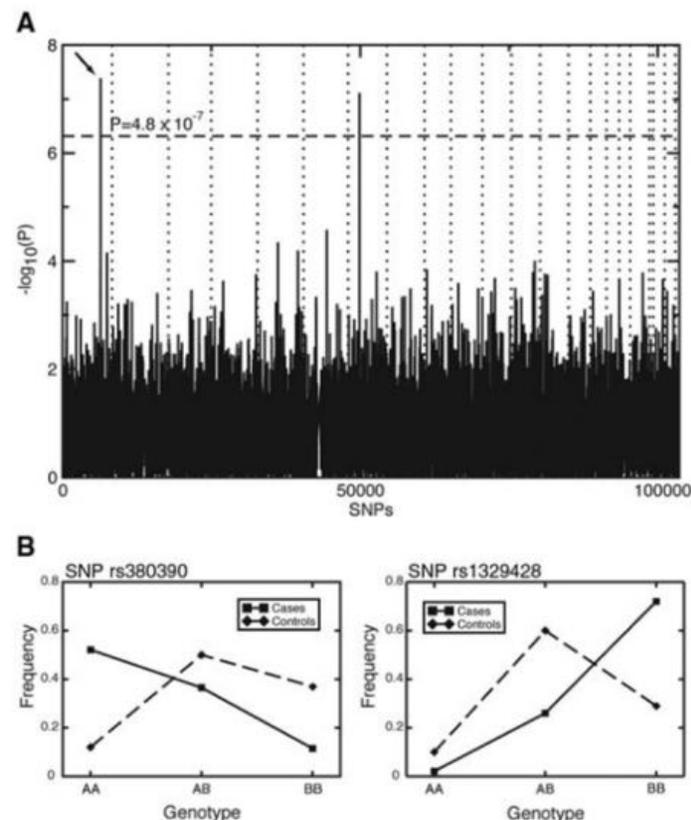
Age-related macular degeneration (AMD) is the leading cause of blindness in the developed world. Its incidence is increasing as the elderly population expands (1). AMD is characterized by progressive destruction of the retina's central region (macula), causing central field visual loss (2). A key feature of AMD is the formation of extracellular deposits called drusen concentrated in and around the macula behind the retina between the retinal pigment epithelium (RPE) and the choroid. To date, no therapy for this disease has proven to be broadly effective. Several risk factors have been linked to AMD, including age, smoking, and family history (3). Candidate-gene studies

have not found any genetic differences that account for a large proportion of the overall prevalence (2). Family-based whole-genome linkage scans have identified chromosomal regions that show evidence of linkage to AMD (4–8), but the linkage areas have not been resolved to any causative mutations.

Like many other chronic diseases, AMD caused by a combination of genetic and environmental risk factors. Linkage studies are not as powerful as association studies for the identification of genes contributing to the risk for common, complex diseases (9). However, linkage studies have the advantage of searching the whole genome in an unbiased manner

without presupposing the involvement of particular genes. Searching the whole genome in an association study requires typing 100,000 or more single-nucleotide polymorphisms (SNPs) (10). Because of these technical demands, only one whole-genome association study, on susceptibility to myocardial infarction, has been published to date (11).

**Study design.** We report a whole-genome case-control association study for genes involved in AMD. To maximize the chance of



**Fig. 1.** (A)  $P$  values of genome-wide association scan for genes that affect the risk of developing AMD.  $-\log_{10}(P)$  is plotted for each SNP in chromosomal order. The spacing between SNPs on the plot is uniform and does not reflect distances between SNPs on the chromosomes. The dotted horizontal line shows the cutoff for  $P = 0.05$  after Bonferroni correction. The vertical dotted lines show chromosomal boundaries. The arrow indicates the peak for SNP rs380390, the most significant association, which was studied further. (B) Variation in genotype frequencies between cases and controls.

# GWAS Catalog

As of 2019-03-22, the GWAS Catalog contains 3886 publications and 130,397 associations.



<http://www.ebi.ac.uk/gwas/diagram>

# ClinVar

ACTGATGGTATGGGGCCAAGAGATATATCT  
CAGGTACGGCTGTCATCACTTAGACCTCAC  
CAGGGCTGGGCATAAAAGTCAGGGCAGAGC  
CCATGGTGCATCTGACTCCTCAGGAGAAGT  
GCAGGTTGGTATCAAGGTTACAAGACAGGT  
GGCACCTGACTCTCTGCTTATTGGTCTAT

## ClinVar

ClinVar aggregates information about genomic variation and its relationship to human health.

### Using ClinVar

- [About ClinVar](#)
- [Data Dictionary](#)
- [Downloads/FTP site](#)
- [FAQ](#)
- [Contact Us](#)
- [RSS feed/What's new?](#)
- [Factsheet](#)

### Tools

- [ACMG Recommendations for Reporting of Incidental Findings](#)
- [ClinVar Submission Portal](#)
- [Submissions](#)
- [Variation Viewer](#)
- [Clinical Remapping - Between assemblies and RefSeqGenes](#)
- [RefSeqGene/LRG](#)

### Related Sites

- [ClinGen](#)
- [GeneReviews®](#)
- [GTR®](#)
- [MedGen](#)
- [OMIM®](#)
- [Variation](#)

### Submitter highlights

We gratefully acknowledge those who have submitted data and provided advice during the development of ClinVar. Subscribe to our [RSS feed](#) and follow us on [Twitter](#) to receive announcements of the release of new datasets. More [information about our submitters](#) is available, as well as a list of submitters with [the number of records each has submitted](#).

### Disclaimer

- ClinVar is a freely accessible, public archive of reports of the relationships among human variations and phenotypes, with supporting evidence
- Currently has 295k mutations
- Most (179k) variants have uncertain affect, only 23 have “4 stars” of significance

# OMIM

Secure https://www.omim.org

About Statistics Downloads Contact Us MIMmatch Donate Help

50 YEARS  
OMIM  
Human Genetics Knowledge  
for the World

**OMIM®**  
Online Mendelian Inheritance in Man®  
An Online Catalog of Human Genes and Genetic Disorders  
Updated April 7, 2017

Search OMIM for clinical features, phenotypes, genes, and more...

Advanced Search : OMIM, Clinical Synopses, Gene Map | Search History  
Need help? : Example Searches, OMIM Search Help, OMIM Tutorial  
Mirror site : [mirror.omim.org](http://mirror.omim.org)

OMIM is supported by a grant from NHGRI, licensing fees, and generous contributions from people like you.

Make a donation!

Institute of Genetic Medicine  
JOHNS HOPKINS MEDICINE

Follow us on Twitter

IRDIRC

NOTE: OMIM is intended for use primarily by physicians and other professionals concerned with genetic disorders, by genetics researchers, and by advanced students in science and medicine. While the OMIM database is open to the public, users seeking information about a personal medical or genetic condition are urged to consult with a qualified physician for diagnosis and for answers to personal questions.  
OMIM® and Online Mendelian Inheritance in Man® are registered trademarks of the Johns Hopkins University.  
Copyright © 1966-2017 Johns Hopkins University.

- For many different diseases and phenotypes, lists what are all of the known genetic associations
- Has records for nearly all genes, ~5k different conditions with known molecular basis, ~1k with unknown basis, ~1k with questionable basis
- Started at JHU 50 years ago 😊

# Biological insights from 108 schizophrenia-associated genetic loci

Schizophrenia Working Group of the Psychiatric Genomics Consortium\*

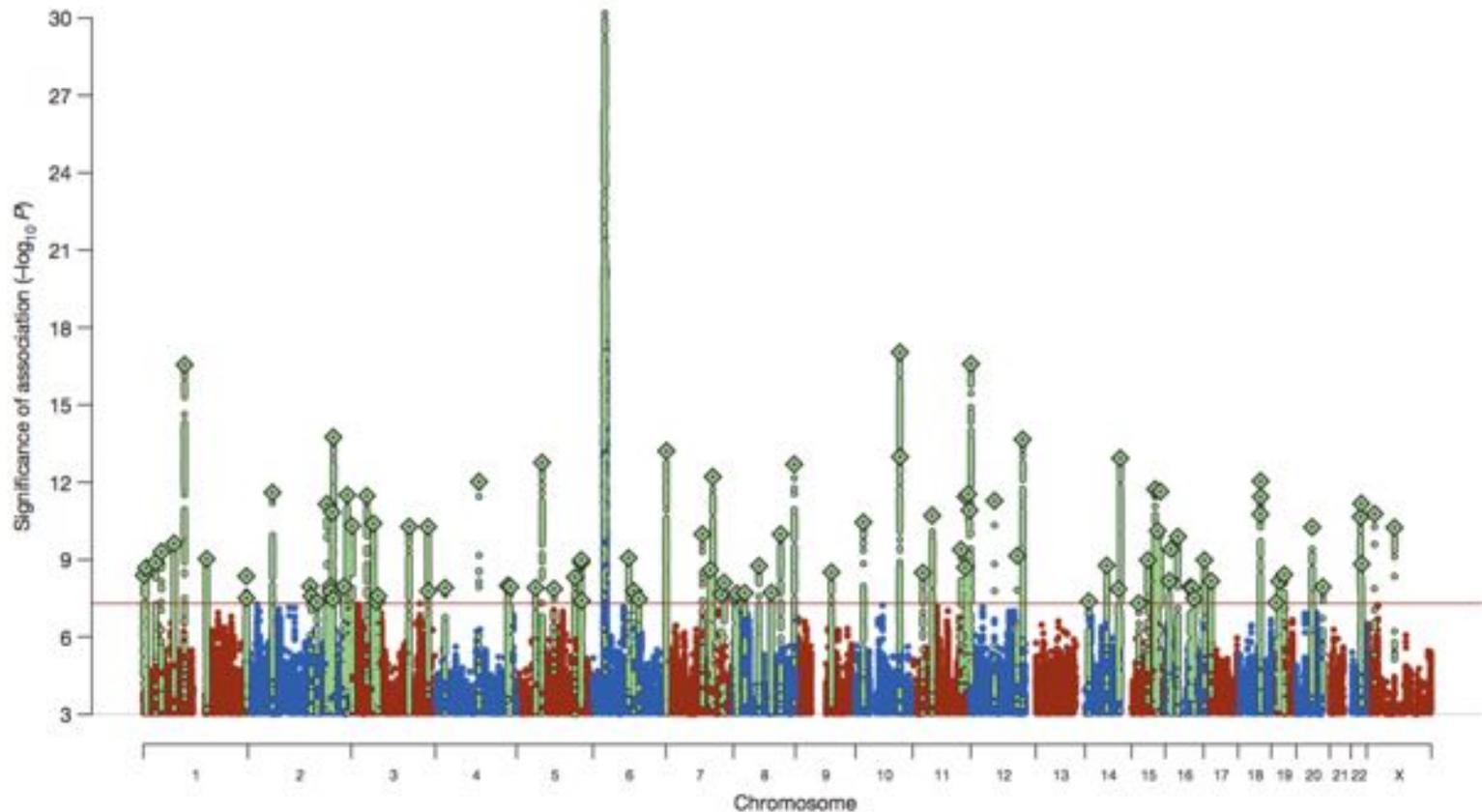
Schizophrenia is a highly heritable disorder. Genetic risk is conferred by a large number of alleles, including common alleles of small effect that might be detected by genome-wide association studies. Here we report a multi-stage schizophrenia genome-wide association study of up to 36,989 cases and 113,075 controls. We identify 128 independent associations spanning 108 conservatively defined loci that meet genome-wide significance, 83 of which have not been previously reported. Associations were enriched among genes expressed in brain, providing biological plausibility for the findings. Many findings have the potential to provide entirely new insights into aetiology, but associations at *DRD2* and several genes involved in glutamatergic neurotransmission highlight molecules of known and potential therapeutic relevance to schizophrenia, and are consistent with leading pathophysiological hypotheses. Independent of genes expressed in brain, associations were enriched among genes expressed in tissues that have important roles in immunity, providing support for the speculated link between the immune system and schizophrenia.

# Biological insights from 108

## schizophrenia

Schizophrenia W

Schizophrenia alleles of small effect sizes span the genome. The findings are consistent with previous reports of polygenic inheritance and several genes of relevance to schizophrenia in brain, associated with support for the



**Figure 1 | Manhattan plot showing schizophrenia associations.** Manhattan plot of the discovery genome-wide association meta-analysis of 49 case control samples (34,241 cases and 45,604 controls) and 3 family based association studies (1,235 parent affected-offspring trios). The x axis is chromosomal

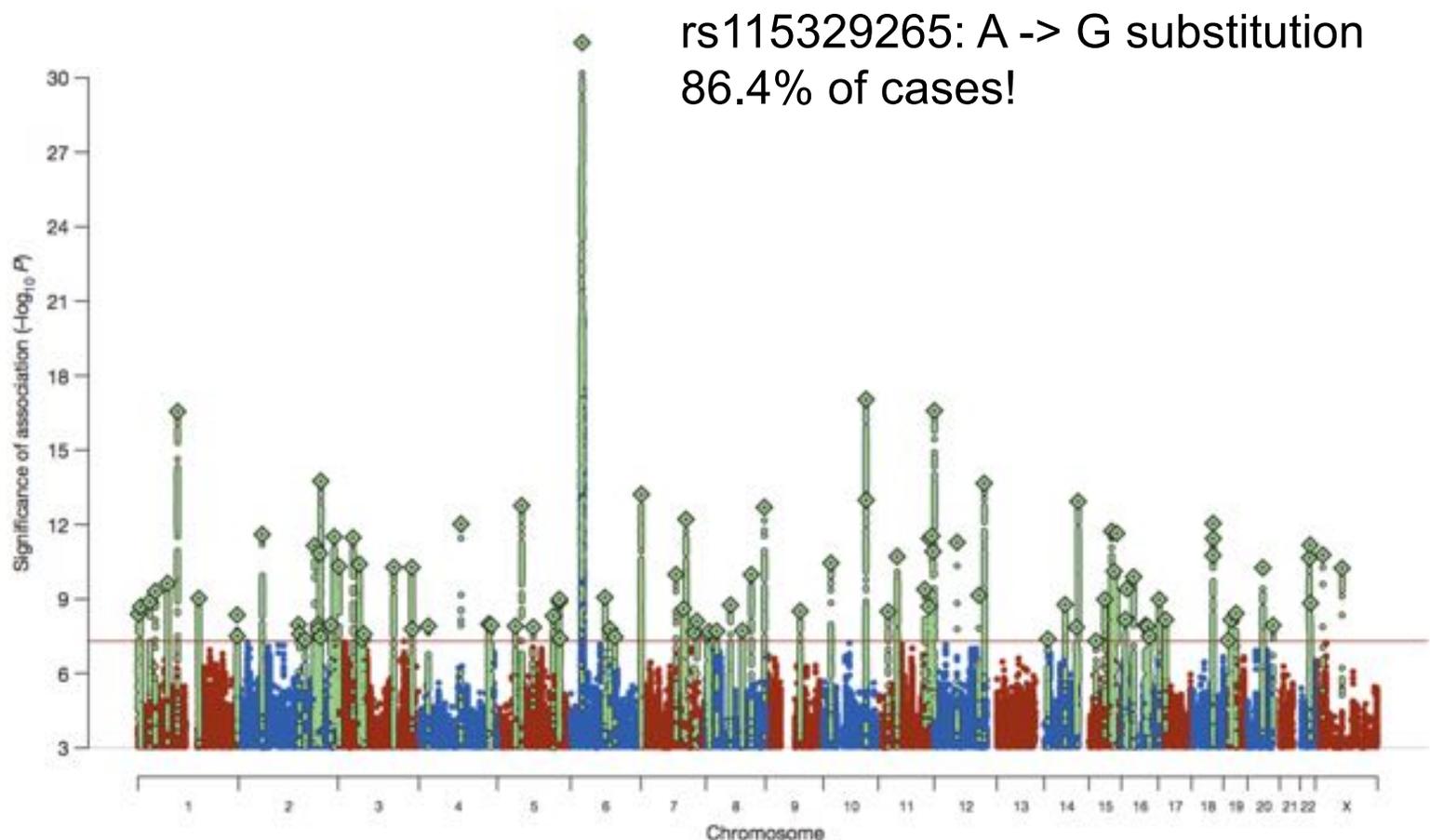
position and the y axis is the significance ( $-\log_{10} P$ ; 2-tailed) of association derived by logistic regression. The red line shows the genome-wide significance level ( $5 \times 10^{-8}$ ). SNPs in green are in linkage disequilibrium with the index SNPs (diamonds) which represent independent genome-wide significant associations.

# Biological insights from 108

## schizo

Schizophrenia W

Schizophrenia alleles of small schizophrenia genetic associations span previously reported findings. Many and several genes are relevant to schizophrenia in brain, associated support for the



**Figure 1 | Manhattan plot showing schizophrenia associations.** Manhattan plot of the discovery genome-wide association meta-analysis of 49 case control samples (34,241 cases and 45,604 controls) and 3 family based association studies (1,235 parent affected-offspring trios). The x axis is chromosomal

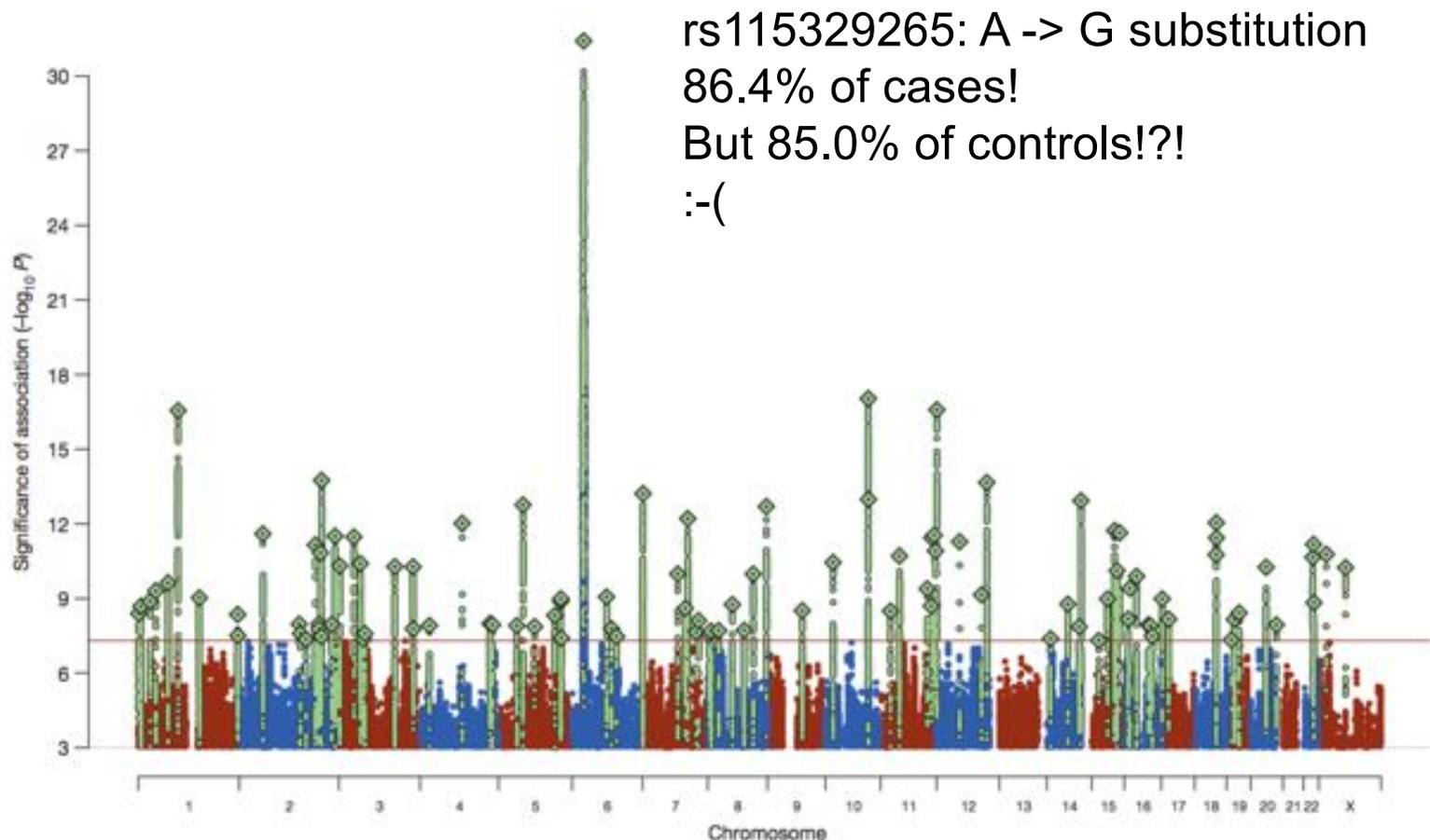
position and the y axis is the significance ( $-\log_{10} P$ ; 2-tailed) of association derived by logistic regression. The red line shows the genome-wide significance level ( $5 \times 10^{-8}$ ). SNPs in green are in linkage disequilibrium with the index SNPs (diamonds) which represent independent genome-wide significant associations.

# Biological insights from 108

## schizo

Schizophrenia W

Schizophrenia alleles of small schizophrenia genetic associations span previously reported findings. Many and several genes are relevant to schizophrenia in brain, associated support for the



**Figure 1 | Manhattan plot showing schizophrenia associations.** Manhattan plot of the discovery genome-wide association meta-analysis of 49 case control samples (34,241 cases and 45,604 controls) and 3 family based association studies (1,235 parent affected-offspring trios). The x axis is chromosomal

position and the y axis is the significance ( $-\log_{10} P$ ; 2-tailed) of association derived by logistic regression. The red line shows the genome-wide significance level ( $5 \times 10^{-8}$ ). SNPs in green are in linkage disequilibrium with the index SNPs (diamonds) which represent independent genome-wide significant associations.

Bi  
SC

Schizo

Schi  
allel  
phre  
ciati  
prev  
the  
and  
relev  
in br  
supp

Compared to the brains of healthy individuals, those of people with schizophrenia have higher expression of a gene called *C4*, according to a paper published in Nature today (January 27). The gene encodes an immune protein that moonlights in the brain as an eradicator of unwanted neural connections (synapses). The findings, which suggest increased synaptic pruning is a feature of the disease, are a direct extension of genome-wide association studies (GWASs) that pointed to the major histocompatibility (MHC) locus as a key region associated with schizophrenia risk.

“The MHC [locus] is the first and the strongest genetic association for schizophrenia, but many people have said this finding is not useful,” said psychiatric geneticist Patrick Sullivan of the University of North Carolina School of Medicine who was not involved in the study.

-Ruth Williams, The Scientist

plot of the discovery genome-wide association meta-analysis of 49 case control samples (34,241 cases and 45,604 controls) and 3 family based association studies (1,235 parent affected-offspring trios). The x axis is chromosomal

derived by logistic regression. The red line shows the genome-wide significance level ( $5 \times 10^{-8}$ ). SNPs in green are in linkage disequilibrium with the index SNPs (diamonds) which represent independent genome-wide significant associations.

# GWAS In Crisis

**Table 1.** Replication and non-replication in associations found by GWA studies of complex diseases published until the end of 2006

Phenotype	Genome-wide association study characteristics				Identified gene/SNPs	Replication status (January 2007)
	platform (SNPs/analyzed)	design	stratification control	n		
Age-related macular degeneration	Affymetrix 100k (116204/103611)	UCC; then sequencing of region	Genomic control, F-ratio	146	<i>CFH</i> /Intronic rs380390; then sequencing showing exonic rs106170 (Y420H) 2kb upstream of 41-kb haplotype block	Meta-analysis of 11 studies (n = 8,991): OR 2.49 and 6.15 (heterozygotes and homozygotes respectively), <b>no large between-study inconsistency in effect sizes; also replicated in large Dutch cohort</b> (n = 5,681); several studies on Asian populations claim no association
Obesity	Affymetrix 100k (116204/86604)	Family-based, 2-stage, followed by mapping 100 neighboring SNPs	Family-based design	694, then up to 923	<i>INSIG2</i> /rs7566605 10kb upstream of the transcription start site	Replication in the same publication in 3 of 4 independent populations of n = 9,881 subjects with modest between-study heterogeneity; 7 more independent populations with over 21,000 subjects total <b>failed to replicate the association</b> ; no effect and no heterogeneity across the independent replication teams
Parkinson disease	Perlegen (248535/198345)	Family-based, second stage with matched case-controls	Family-based design; matching at second stage; also genomic control	443 sib-pairs, then 664	Thirteen genes/ 13 different SNPs identified from analysis of both stages; none with genome-wide significance	Several small replication studies and a large collaborative consortium (n = 12,208) <b>failed to replicate any of the 13 proposed SNPs</b> ; null results were consistent across the teams participating in the consortium
Myocardial infarction	Random gene-based (92788/67671)	UCC	None (just Japanese nationality)	752 (only 94 cases)	<i>LTA</i> /Haplotype of 5 SNPs (2 in <i>LTA</i> and 3 in adjacent genes); the two <i>LTA</i> SNPs had association in larger sample and then Thr26Asn had also functional assay support	Replication in the same publication in additional 1,133 cases and two control groups (n = 1,006 and 872); association not replicated in subsequent ISIS-4 case-control study and meta-analysis (n = 18,325) shows <b>no association (non-significant OR 1.07)</b> without significant between-study heterogeneity vs. 1.77 in originally proposed association for recessive model)
Age-related macular degeneration	Affymetrix 100k (116204/97824)	UCC; then sequencing of region	Genomic control, F-ratio	226	<i>HTRA1</i> /Intragenic rs10490924; then sequencing showing promoter rs11200638 6kb downstream	Independent study (n = 890) published in the same issue starting from dense mapping of locus showing consistent effects with OR 1.90 and 7.51 for heterozygotes and homozygotes, respectively

## Non-Replication and Inconsistency in the Genome-Wide Association Setting

Ioannidis (2007) Hum Hered 2007;64:203–213 <https://doi.org/10.1159/000103512>

# Missing Heritability

NEWS FEATURE PERSONAL GENOMES

NATURE | Vol 456 | 6 November 2008



## The case of the missing heritability

When scientists opened up the human genome, they expected to find the genetic components of common traits and diseases. But they were nowhere to be seen. **Brendan Maher** shines a light on six places where the missing loot could be stashed away.

If you want to predict how tall your children might one day be, a good bet would be to look in the mirror, and at your mate. Studies going back almost a century have estimated that height is 80–90% heritable. So if 29 centimetres separate the tallest 5% of a population from the shortest, then genetics would account for as many as 27 of them.

This year, three groups of researchers scoured the genomes of huge populations (the largest study<sup>1</sup> looked at more than 30,000 people) for genetic variants associated with the height differences. More than 40 turned up.

But there was a problem: the variants had tiny effects. Altogether, they accounted for little more than 5% of height's heritability — just 6 centimetres by the calculations above.



Even though these genome-wide association studies (GWAS) turned up dozens of variants, they did “very little of the prediction that you would do just by asking people how tall their parents are”, says Joel Hirschhorn at the Broad Institute in Cambridge, Massachusetts, who led one of the studies.

Height isn't the only trait in which genes have gone missing, nor is it the most important. Studies looking at similarities between identical and fraternal twins estimate heritability at more than 90% for autism<sup>2</sup> and more than 80% for schizophrenia<sup>3</sup>. And genetics makes a major contribution to disorders such as obesity, diabetes and heart disease. GWAS, one of the most celebrated techniques of the past five years, promised to deliver many of the genes involved (see “Where's the reward?”, page 20). And to some extent they have, identifying more than 400 genetic variants that

contribute to a variety of traits and common diseases. But even when dozens of genes have been linked to a trait, both the individual and cumulative effects are disappointingly small and nowhere near enough to explain earlier estimates of heritability. “It is the big topic in the genetics of common disease right now”, says Francis Collins, former head of the National Human Genome Research Institute (NHGRI) in Bethesda, Maryland. The unexpected results left researchers at a point “where we all had to scratch our heads and say, ‘Huh?’”, he says.

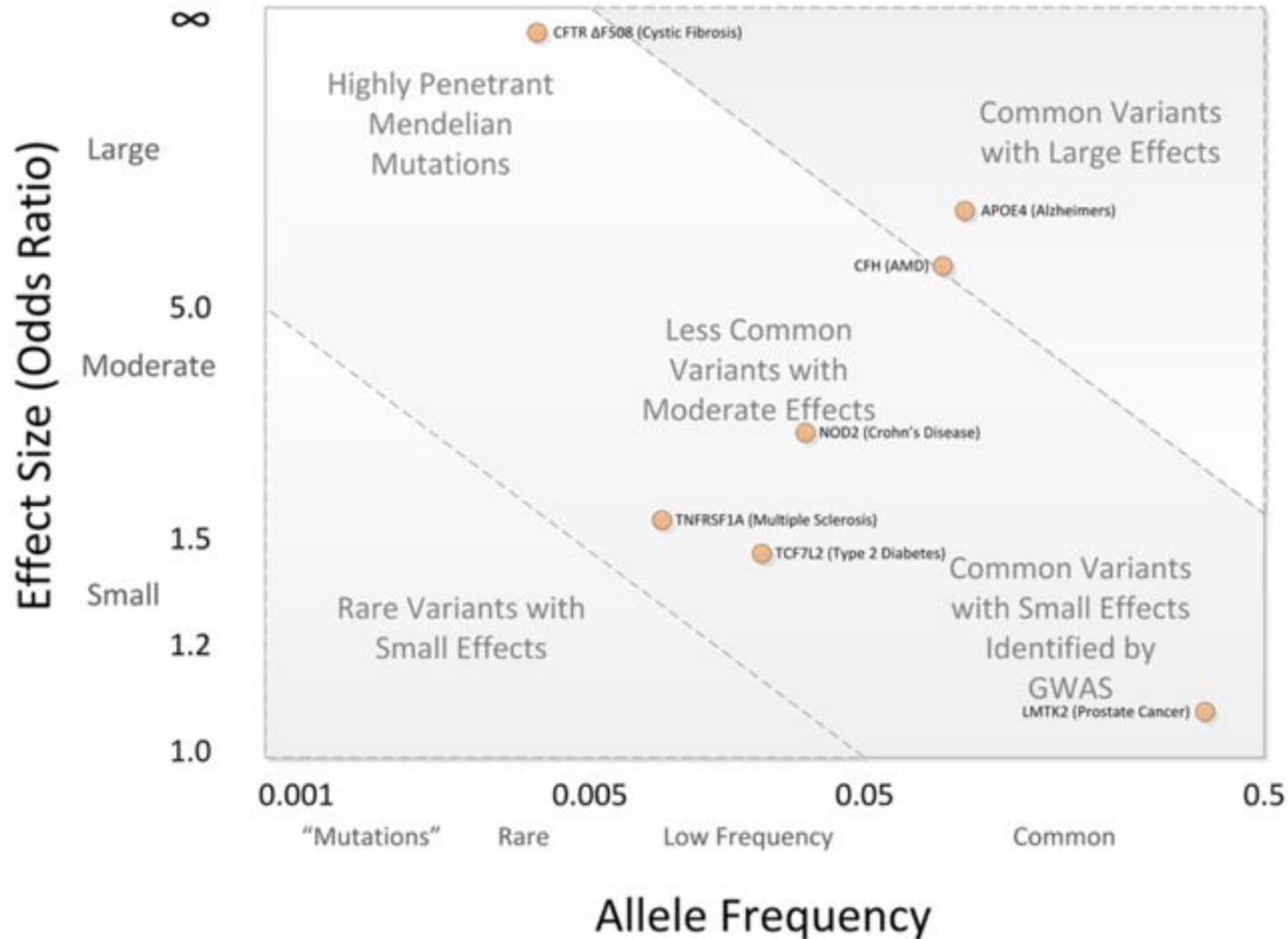
Although flummoxed by this missing heritability, geneticists remain optimistic that they can find more of it. “These are very early days, and there are things that are doable in the next year or two that may well explain another sizeable chunk of heritability”, says Hirschhorn. So where might it be hiding?

ILLUSTRATION BY D. HARRIS

“Three groups of researchers scoured the genomes of huge populations (>30,000 people) for genetic variants associated with the height differences. More than 40 turned up. **But there was a problem: the variants had tiny effects.** Altogether, they accounted for little more than 5% of height's heritability”

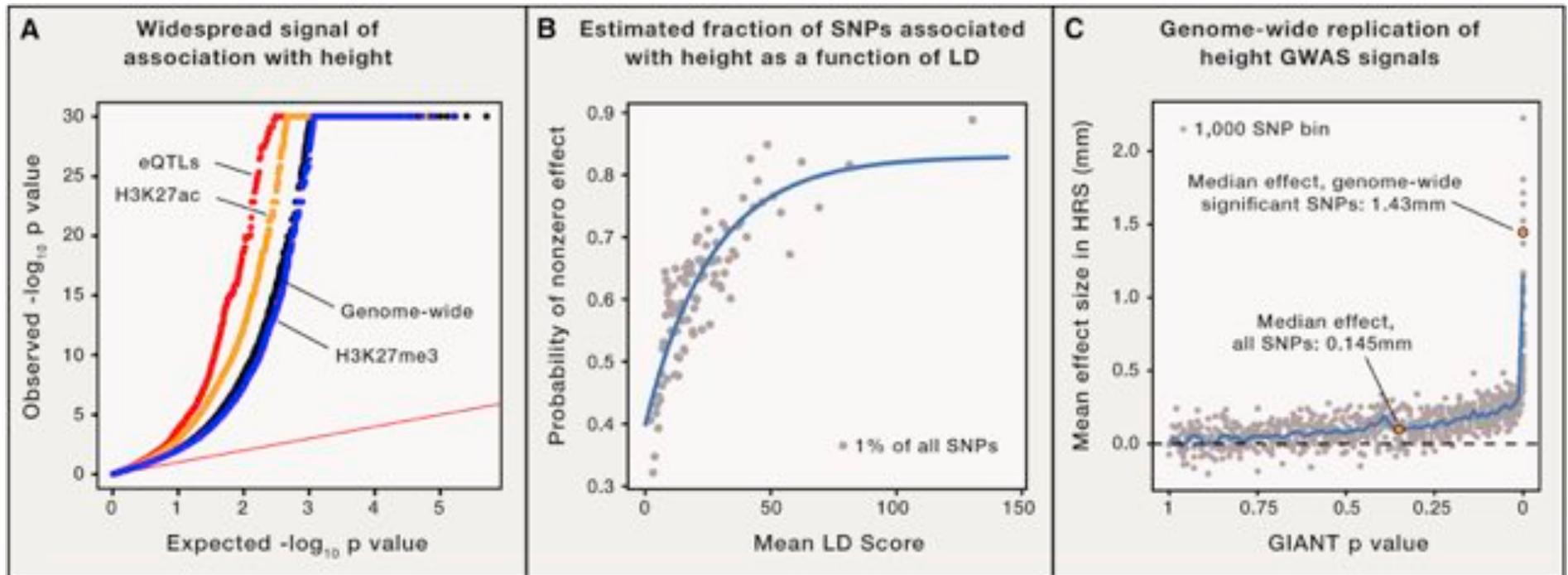
- **Rare, moderately penetrant or common, weakly penetrant variants?**
- **CNVs and SVs?**
- **Epistasis (multiple genes working together)?**
- **Epigenetic effects, especially in utero?**

# Penetrance & Allele Frequency



***Penetrance: The proportion of individuals with a specific genotype who manifest the genotype at the phenotypic level.***

# Omnigenics

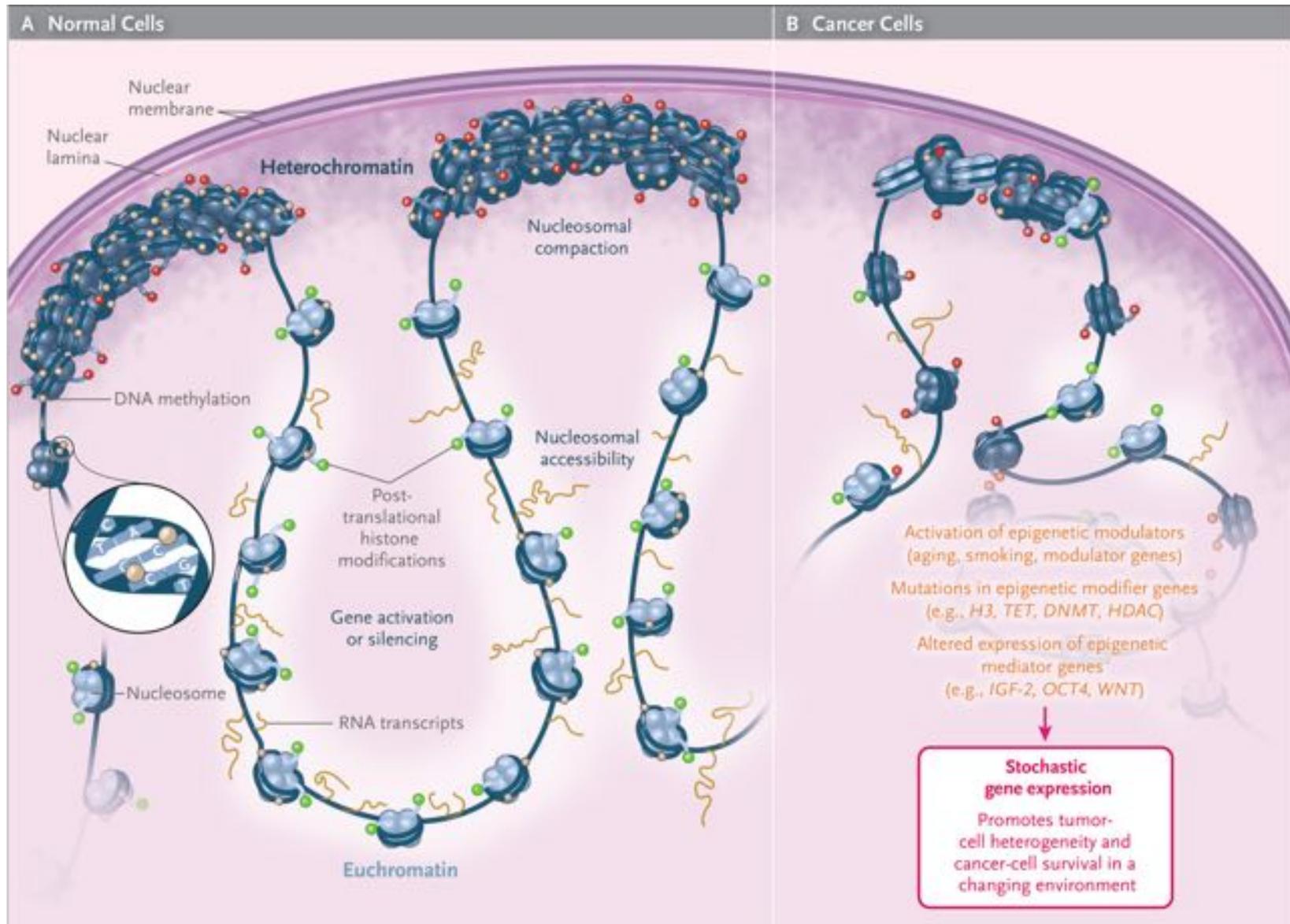


A central goal of genetics is to understand the links between genetic variation and disease. Intuitively, one might expect disease-causing variants to cluster into key pathways that drive disease etiology. But for complex traits, association signals tend to be spread across most of the genome—including near many genes without an obvious connection to disease. **We propose that gene regulatory networks are sufficiently interconnected such that all genes expressed in disease-relevant cells are liable to affect the functions of core disease-related genes and that most heritability can be explained by effects on genes outside core pathways. We refer to this hypothesis as an “omnigenic” model.**

## *An Expanded View of Complex Traits: From Polygenic to Omnigenic*

Boyle, Li, Pritchard (2017) Cell. <https://doi.org/10.1016/j.cell.2017.05.038>

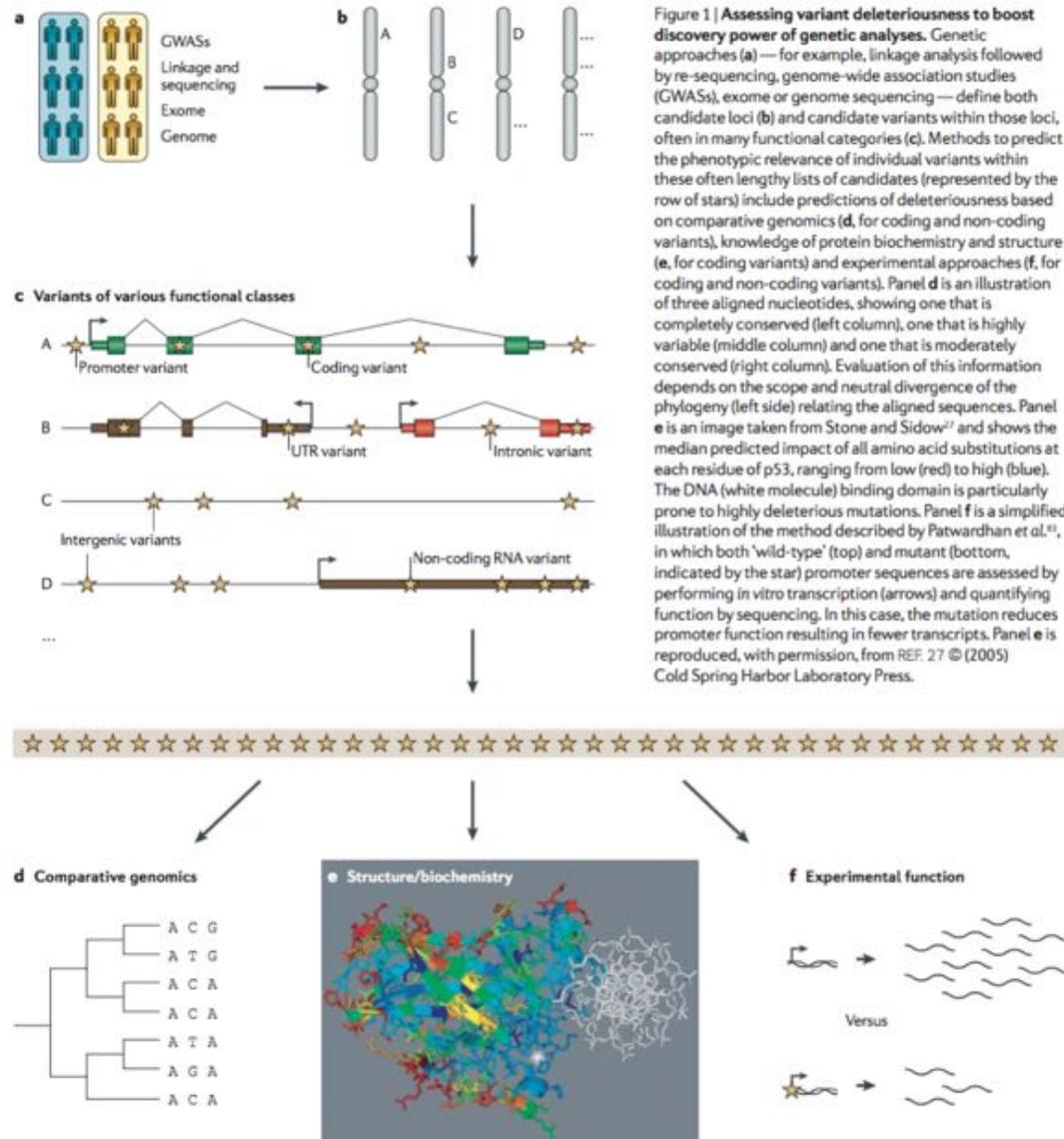
# Epigenetic Factors



## The Key Role of Epigenetics in Human Disease Prevention and Mitigation

Feinberg (2018) NEJM. doi: 10.1056/NEJMra1402513

# Needles in stacks of needles



**Needles in stacks of needles: finding disease-causal variants in a wealth of genomic data**  
Cooper & Shendure (2011) Nature Reviews Genetics.

## Methods

---

# Predicting Deleterious Amino Acid Substitutions

Pauline C. Ng<sup>1,2</sup> and Steven Henikoff<sup>1,3,4</sup>

<sup>1</sup>Fred Hutchinson Cancer Research Center, Seattle, Washington 98109, USA; <sup>2</sup>Department of Bioengineering, University of Washington, Seattle, Washington 98105, USA; <sup>3</sup>Howard Hughes Medical Institute, Seattle, Washington 98109, USA

Many missense substitutions are identified in single nucleotide polymorphism (SNP) data and large-scale random mutagenesis projects. Each amino acid substitution potentially affects protein function. We have constructed a tool that uses sequence homology to predict whether a substitution affects protein function. *SIFT*, which sorts intolerant from tolerant substitutions, classifies substitutions as tolerated or deleterious. A higher proportion of substitutions predicted to be deleterious by *SIFT* gives an affected phenotype than substitutions predicted to be deleterious by substitution scoring matrices in three test cases. Using *SIFT* before mutagenesis studies could reduce the number of functional assays required and yield a higher proportion of affected phenotypes. *SIFT* may be used to identify plausible disease candidates among the SNPs that cause missense substitutions.

***SIFT Key Idea:*** Substituting one amino acid for another with another with very similar biochemical properties is probably less significant than a more dissimilar substitution. Learn those similarities by comparing orthologs across species

# A probabilistic disease-gene finder for personal genomes

Mark Yandell,<sup>1,3,4</sup> Chad Huff,<sup>1,3</sup> Hao Hu,<sup>1,3</sup> Marc Singleton,<sup>1</sup> Barry Moore,<sup>1</sup> Jinchuan Xing,<sup>1</sup> Lynn B. Jorde,<sup>1</sup> and Martin G. Reese<sup>2</sup>

<sup>1</sup>Department of Human Genetics, Eccles Institute of Human Genetics, University of Utah and School of Medicine, Salt Lake City, Utah 84112, USA; <sup>2</sup>Omicia, Inc., Emeryville, California 94608, USA

VAAST (the Variant Annotation, Analysis & Search Tool) is a probabilistic search tool for identifying damaged genes and their disease-causing variants in personal genome sequences. VAAST builds on existing amino acid substitution (AAS) and aggregative approaches to variant prioritization, combining elements of both into a single unified likelihood framework that allows users to identify damaged genes and deleterious variants with greater accuracy, and in an easy-to-use fashion. VAAST can score both coding and noncoding variants, evaluating the cumulative impact of both types of variants simultaneously. VAAST can identify rare variants causing rare genetic diseases, and it can also use both rare and common variants to identify genes responsible for common diseases. VAAST thus has a much greater scope of use than any existing methodology. Here we demonstrate its ability to identify damaged genes using small cohorts ( $n = 3$ ) of unrelated individuals, wherein no two share the same deleterious variants, and for common, multigenic diseases using as few as 150 cases.

[Supplemental material is available for this article.]

**VAAST Key Idea:** Evaluate amino acid substitutions in evolution AND allele frequencies in 1000 genomes project

---

# A general framework for estimating the relative pathogenicity of human genetic variants

Martin Kircher<sup>1,5</sup>, Daniela M Witten<sup>2,5</sup>, Preti Jain<sup>3,4</sup>, Brian J O’Roak<sup>1,4</sup>, Gregory M Cooper<sup>3</sup> & Jay Shendure<sup>1</sup>

Current methods for annotating and interpreting human genetic variation tend to exploit a single information type (for example, conservation) and/or are restricted in scope (for example, to missense changes). Here we describe Combined Annotation–Dependent Depletion (CADD), a method for objectively integrating many diverse annotations into a single measure (C score) for each variant. We implement CADD as a support vector machine trained to differentiate 14.7 million high-frequency human-derived alleles from 14.7 million simulated variants. We precompute C scores for all 8.6 billion possible human single-nucleotide variants and enable scoring of short insertions-deletions. C scores correlate with allelic diversity, annotations of functionality, pathogenicity, disease severity, experimentally measured regulatory effects and complex trait associations, and they highly rank known pathogenic variants within individual genomes. The ability of CADD to prioritize functional, deleterious and pathogenic variants across many functional categories, effect sizes and genetic architectures is unmatched by any current single-annotation method.

comparable, making it difficult to evaluate the relative importance of distinct variant categories or annotations. Third, annotation methods trained on known pathogenic mutations are subject to major ascertainment biases and may not be generalizable. Fourth, it is a major practical challenge to obtain, let alone to objectively evaluate or combine, the existing panoply of partially correlated and partially overlapping annotations; this challenge will only increase in size as large-scale projects such as the Encyclopedia of DNA Elements (ENCODE)<sup>11</sup> continually increase the amount of relevant data available. The net result of these limitations is that many potentially relevant annotations are ignored, while the annotations that are used are applied and combined in *ad hoc* and subjective ways that undermine their usefulness.

Here we describe a general framework, Combined Annotation–Dependent Depletion (CADD), for integrating diverse genome annotations and scoring any possible human single-nucleotide variant (SNV) or small insertion-deletion (indel) event. The basis of CADD is to contrast the annotations of fixed or nearly fixed derived alleles in humans with those of simulated variants. Deleterious variants—that is, variants that reduce organismal fitness—are depleted by natural selection in fixed but not simulated variation. CADD therefore

**CADD Key Idea:** Evaluate amino acid substitutions AND allele frequencies in 1000 genomes project AND ENCODE regions AND ... (63 annotations total :)

---

# A method for calculating probabilities of fitness consequences for point mutations across the human genome

Brad Gulko<sup>1</sup>, Melissa J Hubisz<sup>2</sup>, Ilan Gronau<sup>2,3</sup> & Adam Siepel<sup>1-3</sup>

We describe a new computational method for estimating the probability that a point mutation at each position in a genome will influence fitness. These ‘fitness consequence’ (fitCons) scores serve as evolution-based measures of potential genomic function. Our approach is to cluster genomic positions into groups exhibiting distinct ‘fingerprints’ on the basis of high-throughput functional genomic data, then to estimate a probability of fitness consequences for each group from associated patterns of genetic polymorphism and divergence. We have generated fitCons scores for three human cell types on the basis of public data from ENCODE. In comparison with conventional conservation scores, fitCons scores show considerably improved prediction power for *cis* regulatory elements. In addition, fitCons scores indicate that 4.2–7.5% of nucleotides in the human genome have influenced fitness since the human-chimpanzee divergence, and they suggest that recent evolutionary turnover has had limited impact on the functional content of the genome.

roles<sup>16-19</sup> by getting at fitness directly through observations of evolutionary change. In essence, the ‘experiment’ considered by these methods is the one conducted directly on genomes by nature over millennia, and the outcomes of interest are the presence or absence of fixed mutations.

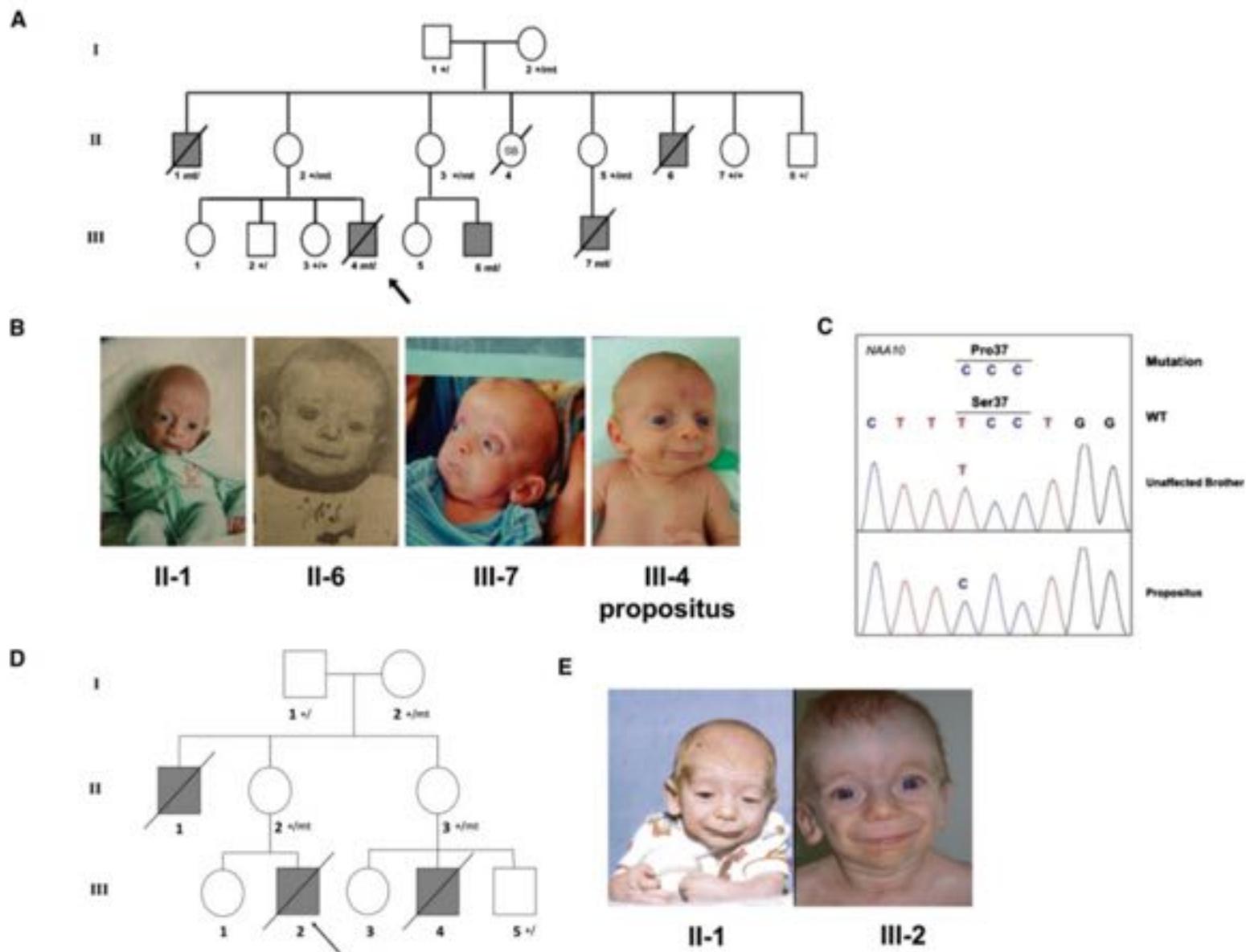
These conservation-based methods, however, depend critically on the assumption that genomic elements are present at orthologous locations and maintain similar functional roles over relatively long evolutionary time periods. Evolutionary turnover may cause inconsistencies between sequence orthology and functional homology that substantially limit this type of analysis. Consequently, investigators have developed two major alternative strategies for the identification and characterization of functional elements. The first strategy is to augment information about interspecies conservation with information about genetic polymorphism<sup>20-28</sup>. The shorter evolutionary time scales associated with intraspecies variation make this approach more robust to evolutionary turnover and less sensitive to errors in alignment and orthology detection. Polymorphic sites tend to be sparse along the genome, however, so this approach requires some type

**fitCons Key Idea:** Evaluate amino acid substitutions AND allele frequencies in 1000 genomes project AND aggregate by ENCODE regions

# Using VAAST to Identify an X-Linked Disorder Resulting in Lethality in Male Infants Due to N-Terminal Acetyltransferase Deficiency

Alan F. Rope,<sup>1</sup> Kai Wang,<sup>2,19</sup> Rune Evjenth,<sup>3</sup> Jinchuan Xing,<sup>4</sup> Jennifer J. Johnston,<sup>5</sup> Jeffrey J. Swensen,<sup>6,7</sup> W. Evan Johnson,<sup>8</sup> Barry Moore,<sup>4</sup> Chad D. Huff,<sup>4</sup> Lynne M. Bird,<sup>9</sup> John C. Carey,<sup>1</sup> John M. Opitz,<sup>1,4,6,10,11</sup> Cathy A. Stevens,<sup>12</sup> Tao Jiang,<sup>13,14</sup> Christa Schank,<sup>8</sup> Heidi Deborah Fain,<sup>15</sup> Reid Robison,<sup>15</sup> Brian Dalley,<sup>16</sup> Steven Chin,<sup>6</sup> Sarah T. South,<sup>1,7</sup> Theodore J. Pysker,<sup>6</sup> Lynn B. Jorde,<sup>4</sup> Hakon Hakonarson,<sup>2</sup> Johan R. Lillehaug,<sup>3</sup> Leslie G. Biesecker,<sup>5</sup> Mark Yandell,<sup>4</sup> Thomas Arnesen,<sup>3,17</sup> and Gholson J. Lyon<sup>15,18,20,\*</sup>

We have identified two families with a previously undescribed lethal X-linked disorder of infancy; the disorder comprises a distinct combination of an aged appearance, craniofacial anomalies, hypotonia, global developmental delays, cryptorchidism, and cardiac arrhythmias. Using X chromosome exon sequencing and a recently developed probabilistic algorithm aimed at discovering disease-causing variants, we identified in one family a c.109T>C (p.Ser37Pro) variant in *NAA10*, a gene encoding the catalytic subunit of the major human N-terminal acetyltransferase (NAT). A parallel effort on a second unrelated family converged on the same variant. The absence of this variant in controls, the amino acid conservation of this region of the protein, the predicted disruptive change, and the co-occurrence in two unrelated families with the same rare disorder suggest that this is the pathogenic mutation. We confirmed this by demonstrating a significantly impaired biochemical activity of the mutant hNaa10p, and from this we conclude that a reduction in acetylation by hNaa10p causes this disease. Here we provide evidence of a human genetic disorder resulting from direct impairment of N-terminal acetylation, one of the most common protein modifications in humans.



**Figure 2. Pedigree Drawing and Pictures of Families 1 and 2**

(A) Pedigree drawing for family 1. The most recent deceased individual, III-4, is the most well-studied subject in the family and is indicated by an arrow. Genotypes are marked for those in which DNA was available and tested. The following abbreviations are used: SB, stillborn; +, normal variant; mt, rare mutant variant.

(B) Pictures of four affected and deceased boys in this family, showing the aged appearance.

(C) Sanger sequencing results of *NAA10* in individual III-4 from family 1.

(D) Pedigree for family 2. Individual III-2 is the most well-studied subject in the family and is indicated by an arrow.

(E) Picture of individuals II-1 and III-2 in family 2 at ~1 year of age.

# Genetic Basis of Autism Spectrum Disorders



## ***Complex disorders of brain development***

- Characterized by difficulties in social interaction, verbal and nonverbal communication and repetitive behaviors.
- Have their roots in very early brain development, and the most obvious signs of autism and symptoms of autism tend to emerge between 2 and 3 years of age.

## ***U.S. CDC identify around 1 in 68 American children as on the autism spectrum***

- Ten-fold increase in prevalence in 40 years, only partly explained by improved diagnosis and awareness.
- Studies also show that autism is four to five times more common among boys than girls.
- Specific causes remain elusive

## **What is Autism?**

<http://www.autismspeaks.org/what-autism>

# Autism is NOT caused by vaccines

EARLY REPORT

## Early report

### Ileal-lymphoid-nodular hyperplasia, non-specific colitis, and pervasive developmental disorder in children

A J Wakefield, S H Murch, A Anthony, J Linnell, D M Casson, M Malik, M Berelowitz, A P Dillon, M A Thomson, P Harvey, A Valentine, S E Davies, J A Walker-Smith

#### Summary

**Background** We investigated a consecutive series of children with chronic enterocolitis and regressive developmental disorder.

**Methods** 12 children (mean age 6 years [range 3–10], 11 boys) were referred to a paediatric gastroenterology unit with a history of normal development followed by loss of acquired skills, including language, together with diarrhoea and abdominal pain. Children underwent gastroenterological, neurological, and developmental assessment and review of developmental records, ileocolonoscopy and biopsy sampling, magnetic-resonance imaging (MRI), electroencephalography (EEG), and lumbar puncture were done under sedation. Barium follow-through radiography was done where possible. Biochemical, haematological, and immunologic profiles were examined.

**Findings** Onset of behavioural symptoms was associated by the parents, with measles, mumps, and rubella vaccination in eight of the 12 children, with measles infection in one child, and otitis media in another. All 12 children had intestinal abnormalities ranging from lymphoid nodular hyperplasia to granulomatous inflammation. Histology showed patchy chronic inflammation in 11 children and reactive ileal lymphoid hyperplasia in seven, but no granulomas. Behavioural disorders included autism (nine), disintegrative psychosis (one), and possible postviral or vaccinal encephalitis (two). There were no focal neurological abnormalities and MRI and EEG tests were normal. Abnormal laboratory results were significantly raised urinary methylmalonic acid compared with age-matched controls (0.03), low haemoglobin in four children, and low serum IgA in three children.

**Interpretation** The identical associated gastrointestinal disease and developmental regression in a group of previously normal children, which was generally associated in time to possible environmental triggers.

Lancet 1998; **351**: 637–41  
See Commentary page

**Inflammatory Bowel Disease Study Group, University Departments of Medicine and Histopathology** (A J Wakefield *msc*, A Anthony *msc*, J Linnell *md*, A P Dillon *msc*, S E Davies *msc*) and **the University Departments of Paediatric Gastroenterology** (S H Murch *msc*, D M Casson *msc*, M Malik *msc*, M A Thomson *msc*, J A Walker-Smith *msc*), **Child and Adolescent Psychiatry** (M Berelowitz *msc*), **Neurology** (P Harvey *msc*), and **Radiology** (A Valentine *msc*), **Royal Free Hospital and School of Medicine, London NW3 2QG, UK**

Correspondence to: Dr A J Wakefield

#### Introduction

We saw several children who, after a period of apparent normality, lost acquired skills, including communication. They all had gastrointestinal symptoms, including abdominal pain, diarrhoea, and vomiting and, in some cases, food intolerance. We describe the clinical findings, and gastrointestinal features, of these children.

#### Patients and methods

12 children, consecutively referred to the department of paediatric gastroenterology with a history of a pervasive developmental disorder with loss of acquired skills and intestinal symptoms (abdominal pain, bloating and food intolerance), were investigated. All children were admitted to the ward for a week, accompanied by their parents.

#### Clinical investigations

Each child took history, including details of immunisations and exposure to infectious diseases, and assessed the children. In 11 cases the history was obtained by the senior clinician (JW-S). Neurological and psychiatric assessments were done by consultant staff (PH, MB) with HMS-4 criteria.<sup>1</sup> Developmental records included a review of prospective developmental records from parents, health visitors, and general practitioners. Four children did not undergo psychiatric assessment in hospital; all had been assessed professionally elsewhere, so these assessments were used as the basis for their behavioural diagnosis.

After bowel preparation, ileocolonoscopy was performed by SHM or MAT under sedation with midazolam and pethidine. Paired frozen and formalin-fixed mucosal biopsy samples were taken from the terminal ileum; ascending, transverse, descending, and sigmoid colons, and from the rectum. The procedure was recorded by video or still images, and were compared with images of the previous seven consecutive paediatric colonoscopies (four normal colonoscopies and three on children with ulcerative colitis), in which the physician reported normal appearances in the terminal ileum. Barium follow-through radiography was possible in some cases.

Also under sedation, cerebral magnetic-resonance imaging (MRI), electroencephalography (EEG) including visual, brain stem auditory, and sensory evoked potentials (where compliance made these possible), and lumbar puncture were done.

#### Laboratory investigations

Thyroid function, serum long-chain fatty acids, and cerebrospinal-fluid lactate were measured to exclude known causes of childhood neurodegenerative disease. Urinary methylmalonic acid was measured in random urine samples from eight of the 12 children and 14 age-matched and sex-matched normal controls, by a modification of a technique described previously.<sup>2</sup> Chromatograms were scanned digitally on computer, to analyse the methylmalonic-acid zones from cases and controls. Urinary methylmalonic-acid concentrations in patients and controls were compared by a two-sample *t* test. Urinary creatinine was estimated by routine spectrophotometric assay.

Children were screened for antiendomyxial antibodies and boys were screened for fragile-X if this had not been done

THE JOURNAL OF PEDIATRICS • www.jpeds.com

ORIGINAL  
ARTICLES

### Increasing Exposure to Antibody-Stimulating Proteins and Polysaccharides in Vaccines Is Not Associated with Risk of Autism

Frank DeStefano, MD, MPH<sup>1</sup>, Cristofer S. Price, ScM<sup>2</sup>, and Eric S. Weintraub, MPH<sup>1</sup>

**Objective** To evaluate the association between autism and the level of immunologic stimulation received from vaccines administered during the first 2 years of life.

**Study design** We analyzed data from a case-control study conducted in 3 managed care organizations (MCOs) of 256 children with autism spectrum disorder (ASD) and 752 control children matched on birth year, sex, and MCO. In addition to the broader category of ASD, we also evaluated autistic disorder and ASD with regression. ASD diagnoses were validated through standardized in-person evaluations. Exposure to total antibody-stimulating proteins and polysaccharides from vaccines was determined by summing the antigen content of each vaccine received, as obtained from immunization registries and medical records. Potential confounding factors were ascertained from parent interviews and medical charts. Conditional logistic regression was used to assess associations between ASD outcomes and exposure to antigens in selected time periods.

**Results** The aOR (95% CI) of ASD associated with each 25-unit increase in total antigen exposure was 0.999 (0.994-1.003) for cumulative exposure to age 3 months, 0.999 (0.997-1.001) for cumulative exposure to age 7 months, and 0.999 (0.998-1.001) for cumulative exposure to age 2 years. Similarly, no increased risk was found for autistic disorder or ASD with regression.

**Conclusion** In this study of MCO members, increasing exposure to antibody-stimulating proteins and polysaccharides in vaccines during the first 2 years of life was not related to the risk of developing an ASD. (*J Pediatr* 2013;163:561-7).

The initial concerns that vaccines may cause autism were related to the measles, mumps, and rubella vaccine<sup>1</sup> and thimerosal-containing vaccines.<sup>2</sup> In 2004, a comprehensive review by the Institute of Medicine concluded that the evidence favors rejection of possible causal associations between each of these vaccine types and autism.<sup>3</sup> Nonetheless, concerns about a possible link between vaccines and autism persist,<sup>4</sup> with the latest concern centering on the number of vaccines administered to infants and young children.<sup>5</sup> A recent survey found that parents' top vaccine-related concerns included administration of too many vaccines during the first 2 years of life, administration of too many vaccines in a single doctor visit, and a possible link between vaccines and learning disabilities, such as autism.<sup>6</sup> All of the foregoing concerns were reported by 30%-36% of all survey respondents, and were reported by 55%-90% of parents who indicated that their children would receive some, but not all, of the vaccines on the recommended schedule. Another recent survey found that more than 10% of parents of young children refuse or delay vaccinations, with most believing that delaying vaccine doses is safer than providing them in accordance with the Centers for Disease Control and Prevention's recommended vaccination schedule.<sup>7</sup>

Using the number of antibody-stimulating proteins and polysaccharides contained in vaccines as a measure, we evaluated the association between the level of immunologic stimulation received from vaccines during the first 2 years of life and the risk of developing an autism spectrum disorder (ASD), including specific ASD subtypes.

#### Methods

We performed a secondary analysis of publicly available data from a case-control study designed to examine potential associations between exposure to thimerosal-containing injections and ASD.<sup>8</sup> The study was conducted in 3 managed care organizations (MCOs). Data sources for the original study included MCO computerized data files, abstraction of biological mothers' and children's medical charts, and standardized telephone interviews with biological mothers. Case children underwent standardized in-person assessment to verify case status.

AD Autistic disorder  
ADI-R Autism Diagnostic Interview-Revised  
ADOS Autism Diagnostic Observation Schedule  
ASD Autism spectrum disorder  
MCO Managed care organization  
SCQ Social Communication Questionnaire

From the <sup>1</sup>Immunization Safety Office, Centers for Disease Control and Prevention, Atlanta, GA and <sup>2</sup>Abt Associates Inc, Bethesda, MD

Funded by a contract from the Centers for Disease Control and Prevention to America's Health Insurance Plans (AHIP), and by subcontracts from AHIP to Abt Associates, Inc. The findings and conclusions in this study are those of the authors and do not necessarily represent the official position of the Centers for Disease Control and Prevention. The authors declare no conflicts of interest.

0022-3476/\$ - see front matter. Copyright © 2013 Mosby Inc. All rights reserved. <http://dx.doi.org/10.1016/j.jpeds.2013.08.021>

# Autism is NOT caused by vaccines

EARLY REPORT

Early report

THE JOURNAL OF PEDIATRICS • www.jpeds.com

ORIGINAL  
ARTICLES

Increasing Exposure to Antibody-Stimulating Proteins and Polysaccharides

The GMC hearings, which began in July 2007, centered on Wakefield's 1998 report. Many studies have found no connections [5,6], but sensational publicity caused immunization rates in the UK to drop more than 10 percent and have left lingering doubts among parents worldwide.

The GMC began investigating after learning from Deer that Wakefield had failed to declare he had been paid £55,000 to advise lawyers representing parents who believed that the vaccine had harmed their children. The GMC found that Wakefield had:

- Improperly obtained blood for research purposes from normal children attending his son's birthday party, paid them £5 for their discomfort, and [later joked during a lecture about having done this](#).
- Subjected autistic children to colonoscopy, lumbar punctures, and other tests without approval from a research review board.
- Failed to disclose that he had filed a patent for a vaccine to compete with the MMR
- Starting a child on an experimental product called Transfer Factor, which he planned to market.

S H Murch *MD*, D M Casson *MD*, M Malik *MD*,  
M A Thomson *MD*, J A Walker-Smith *MD*, **Child and Adolescent  
Psychiatry** (M Berelowitz *MD*), **Neurology** (P Harvey *MD*), and  
**Radiology** (A Valentine *MD*), **Royal Free Hospital and School of  
Medicine, London NW3 2QG, UK**

Correspondence to: Dr A J Wakefield

and controls. Urinary methylmalonic-acid concentrations in patients and controls were compared by a two-sample *t* test. Urinary creatinine was estimated by routine spectrophotometric assay.

Children were screened for antidiemysal antibodies and boys were screened for fragile-X if this had not been done

ADI-R Autism Diagnostic Interview-Revised  
ADOS Autism Diagnostic Observation Schedule  
ASD Autism spectrum disorder  
MCO Managed care organization  
SCQ Social Communication Questionnaire

Associate, Inc. The findings and conclusions in this study are those of the authors and do not necessarily represent the official position of the Centers for Disease Control and Prevention. The authors declare no conflicts of interest.

0022-3476/\$ - see front matter. Copyright © 2013 Mosby Inc. All rights reserved. <http://dx.doi.org/10.1016/j.jpeds.2013.02.001>

# Autism is NOT caused by antidepressants

Research

JAMA | Original Investigation

## Association Between Serotonergic Antidepressant Use During Pregnancy and Autism Spectrum Disorder in Children

Hilary K. Brown, PhD; Joel G. Ray, MD, MSc, FRCPC; Andrew S. Wilton, MSc; Yona Lunsky, PhD, CPsych; Tara Gomes, MSc; Simone N. Vigod, MD, MSc, FRCPC

**IMPORTANCE** Previous observations of a higher risk of child autism spectrum disorder with serotonergic antidepressant exposure during pregnancy may have been confounded.

**OBJECTIVE** To evaluate the association between serotonergic antidepressant exposure during pregnancy and child autism spectrum disorder.

**DESIGN, SETTING, AND PARTICIPANTS** Retrospective cohort study. Health administrative data sets were used to study children born to mothers who were receiving public prescription drug coverage during pregnancy in Ontario, Canada, from 2002-2010, reflecting 4.2% of births. Children were followed up until March 31, 2014.

**EXPOSURES** Serotonergic antidepressant exposure was defined as 2 or more consecutive maternal prescriptions for a selective serotonin reuptake inhibitor or serotonin-norepinephrine reuptake inhibitor between conception and delivery.

**MAIN OUTCOMES AND MEASURES** Child autism spectrum disorder identified after the age of 2 years. Exposure group differences were addressed by inverse probability of treatment weighting based on derived high-dimensional propensity scores (computerized algorithm used to select a large number of potential confounders) and by comparing exposed children with unexposed siblings.

**RESULTS** There were 35 906 singleton births at a mean gestational age of 38.7 weeks (50.4% were male, mean maternal age was 26.7 years, and mean duration of follow-up was 4.95 years). In the 2837 pregnancies (7.9%) exposed to antidepressants, 2.0% (95% CI, 1.6%-2.6%) of children were diagnosed with autism spectrum disorder. The incidence of autism spectrum disorder was 4.51 per 1000 person-years among children exposed to antidepressants vs 2.03 per 1000 person-years among unexposed children (between-group difference, 2.48 [95% CI, 2.33-2.62] per 1000 person-years; hazard ratio [HR], 2.16 [95% CI, 1.64-2.86]; adjusted HR, 1.59 [95% CI, 1.17-2.17]). After inverse probability of treatment weighting based on the high-dimensional propensity score, the association was not significant (HR, 1.61 [95% CI, 0.997-2.59]). The association was also not significant when exposed children were compared with unexposed siblings (incidence of autism spectrum disorder was 3.40 per 1000 person-years vs 2.05 per 1000 person-years, respectively; adjusted HR, 1.60 [95% CI, 0.69-3.74]).

**CONCLUSIONS AND RELEVANCE** In children born to mothers receiving public drug coverage in Ontario, Canada, in utero serotonergic antidepressant exposure compared with no exposure was not associated with autism spectrum disorder in the child. Although a causal relationship cannot be ruled out, the previously observed association may be explained by other factors.

[Editorial page 1533](#)  
[Related article page 1553](#)  
[Supplemental content](#)

**Author Affiliations:** Women's College Research Institute, Women's College Hospital, Toronto, Ontario, Canada (Brown, Vigod); Department of Psychiatry, University of Toronto, Toronto, Ontario, Canada (Brown, Lunsky, Vigod); Institute for Clinical Evaluative Sciences, Toronto, Ontario, Canada (Brown, Ray, Wilton, Lunsky, Gomes, Vigod); Department of Obstetrics and Gynecology, St Michael's Hospital, Toronto, Ontario, Canada (Ray); Department of Medicine, University of Toronto, Toronto, Ontario, Canada (Ray); Centre for Addiction and Mental Health, Toronto, Ontario, Canada (Lunsky); Leslie Dan Faculty of Pharmacy, University of Toronto, Toronto, Ontario, Canada (Gomes).

**Corresponding Author:** Simone N. Vigod, MD, MSc, FRCPC, Women's College Hospital, 76 Grenville St, Toronto, ON M5S 1B2, Canada (simone.vigod@wchospital.ca).

JAMA. 2017;317(15):1544-1552. doi:10.1001/jama.2017.3415

1544

Copyright 2017 American Medical Association. All rights reserved.

jama.com

Research

JAMA | Original Investigation

## Associations of Maternal Antidepressant Use During the First Trimester of Pregnancy With Preterm Birth, Small for Gestational Age, Autism Spectrum Disorder, and Attention-Deficit/Hyperactivity Disorder in Offspring

Ayesha C. Sujan, MA, Martin E. Rickert, PhD, A. Sara Öberg, MD, PhD; Patrick D. Quinn, PhD; Sonia Hernández-Díaz, MD, PhD; Catarina Almqvist, MD, PhD; Paul Lichtenstein, PhD; Henrik Larsson, PhD; Brian M. D'Onofrio, PhD

**IMPORTANCE** Prenatal antidepressant exposure has been associated with adverse outcomes. Previous studies, however, may not have adequately accounted for confounding.

**OBJECTIVE** To evaluate alternative hypotheses for associations between first-trimester antidepressant exposure and birth and neurodevelopmental problems.

**DESIGN, SETTING, AND PARTICIPANTS** This retrospective cohort study included Swedish offspring born between 1996 and 2012 and followed up through 2013 or censored by death or emigration. Analyses controlling for pregnancy, maternal and paternal covariates, as well as sibling comparisons, timing of exposure comparisons, and paternal comparisons, were used to examine the associations.

**EXPOSURES** Maternal self-reported first-trimester antidepressant use and first-trimester antidepressant dispensations.

**MAIN OUTCOMES AND MEASURES** Preterm birth (<37 gestational weeks), small for gestational age (birth weight <2 SDs below the mean for gestational age), and first inpatient or outpatient clinical diagnosis of autism spectrum disorder and attention-deficit/hyperactivity disorder in offspring.

**RESULTS** Among 1 580 629 offspring (mean gestational age, 279 days; 48.6% female, 1.4% [n = 22 544] with maternal first-trimester self-reported antidepressant use) born to 943 776 mothers (mean age at childbirth, 30 years), 6.98% of exposed vs 4.78% of unexposed offspring were preterm, 2.54% of exposed vs 2.19% of unexposed were small for gestational age, 5.28% of exposed vs 2.14% of unexposed were diagnosed with autism spectrum disorder by age 15 years, and 12.63% of exposed vs 5.46% of unexposed were diagnosed with attention-deficit/hyperactivity disorder by age 15 years. At the population level, first-trimester exposure was associated with all outcomes compared with unexposed offspring (preterm birth odds ratio [OR], 1.47 [95% CI, 1.40-1.55]; small for gestational age OR, 1.35 [95% CI, 1.06-1.25]; autism spectrum disorder hazard ratio [HR], 2.02 [95% CI, 1.80-2.26]; attention-deficit/hyperactivity disorder HR, 2.21 [95% CI, 2.04-2.39]). However, in models that compared siblings while adjusting for pregnancy, maternal, and paternal traits, first-trimester antidepressant exposure was associated with preterm birth (OR, 1.34 [95% CI, 1.18-1.52]) but not with small for gestational age (OR, 1.01 [95% CI, 0.81-1.25]), autism spectrum disorder (HR, 0.83 [95% CI, 0.62-1.13]), or attention-deficit/hyperactivity disorder (HR, 0.99 [95% CI, 0.79-1.25]). Results from analyses assessing associations with maternal dispensations before pregnancy and with paternal first-trimester dispensations were consistent with findings from the sibling comparisons.

**CONCLUSIONS AND RELEVANCE** Among offspring born in Sweden, after accounting for confounding factors, first-trimester exposure to antidepressants, compared with no exposure, was associated with a small increased risk of preterm birth but no increased risk of small for gestational age, autism spectrum disorder, or attention-deficit/hyperactivity disorder.

JAMA. 2017;317(15):1553-1562. doi:10.1001/jama.2017.3413

[Editorial page 1533](#)  
[Related article page 1544](#)  
[Supplemental content](#)

**Author Affiliations:** Department of Psychological and Brain Sciences, Indiana University, Bloomington (Sujan, Rickert, Quinn, D'Onofrio); Department of Medical Epidemiology and Biostatistics, Karolinska Institutet, Stockholm, Sweden (Öberg, Almqvist, Lichtenstein, Larsson, D'Onofrio); Department of Epidemiology, Harvard T.H. Chan School of Public Health, Boston, Massachusetts (Öberg, Hernández-Díaz); School of Medical Sciences, Örebro University, Örebro, Sweden (Larsson).

**Corresponding Author:** Brian M. D'Onofrio, PhD, Department of Psychological and Brain Sciences, Indiana University, 1101 E 10th St, Bloomington, IN 47405 (bmdonof@indiana.edu).

Copyright 2017 American Medical Association. All rights reserved.

1553

## LETTERS

# A genome-wide linkage and association scan reveals novel loci for autism

Lauren A. Weiss<sup>1,2\*</sup>†, Dan E. Arking<sup>3\*</sup> & The Gene Discovery Project of Johns Hopkins & the Autism Consortium‡

Although autism is a highly heritable neurodevelopmental disorder, attempts to identify specific susceptibility genes have thus far met with limited success<sup>1</sup>. Genome-wide association studies using half a million or more markers, particularly those with very large sample sizes achieved through meta-analysis, have shown great success in mapping genes for other complex genetic traits. Consequently, we initiated a linkage and association mapping study using half a million genome-wide single nucleotide polymorphisms (SNPs) in a common set of 1,031 multiplex autism families (1,553 affected offspring). We identified regions of suggestive and significant linkage on chromosomes 6q27 and 20p13, respectively. **Initial analysis did not yield genome-wide significant associations;** however, genotyping of top hits in additional families revealed an SNP on chromosome 5p15 (between *SEMA5A* and *TAS2R1*) that was significantly associated with autism ( $P = 2 \times 10^{-7}$ ). We also demonstrated that expression of *SEMA5A* is reduced in brains from autistic patients, further implicating *SEMA5A* as an autism susceptibility gene. The linkage regions reported here provide targets for rare variation screening whereas the discovery of a single novel association demonstrates the action of common variants.

For a high-resolution genetic study of autism, we selected families

Before merging, we carefully filtered each data set separately to ensure the highest possible genotype quality for analysis, because technical genotyping artefacts can create false positive findings. We therefore examined the distribution of  $\chi^2$  values for the highest quality data, and used a series of quality control (QC) filters designed to identify a robust set of SNPs, including data completeness for each SNP, Mendelian errors per SNP and per family, and a careful evaluation of inflation of association statistics as a function of allele frequency and missing data (see Methods). As 324 individuals were genotyped at both centres, we performed a concordance check to validate our approach. After excluding one sample mix-up, we obtained an overall genotype concordance between the two centres of 99.7% for samples typed on 500K at Johns Hopkins University and 5.0 at the Broad Institute and 99.9% for samples run on 5.0 arrays at both sites. The combined data set, consisting of 1,031 nuclear families (856 with two parents) and a total of 1,553 affected offspring, was used for genetic analyses (Supplementary Table 1). These data were publicly released in October 2007 and are directly available from AGRE and NIMH.

For linkage analyses, the common AGRE/NIMH data set was further merged with Illumina 550K genotype data generated at the Children's Hospital of Philadelphia (CHOP) and available from AGRE, adding

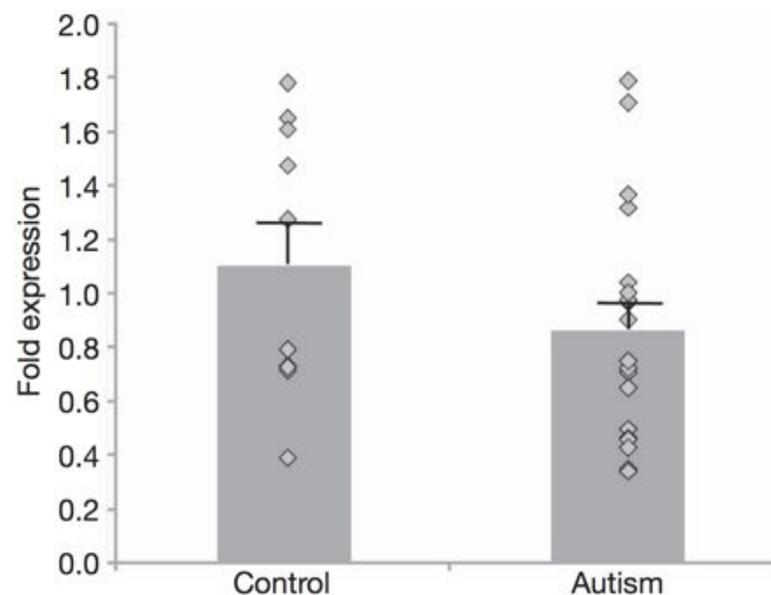
## LETTERS

# A genome-wide linkage and association scan reveals novel loci for autism

Lauren A. Weiss<sup>1,2\*</sup>†, Dan E. Arking<sup>3\*</sup> & The Gene Discovery Project of Johns Hopkins & the Autism Consortium‡

Although autism is a highly heritable neurodevelopmental disorder, attempts to identify specific susceptibility genes have so far met with limited success<sup>1</sup>. Genome-wide association studies using half a million or more markers, particularly those with large sample sizes achieved through meta-analysis, have had great success in mapping genes for other complex genetic disorders. Consequently, we initiated a linkage and association mapping study using half a million genome-wide single nucleotide polymorphisms (SNPs) in a common set of 1,031 multiplex autism families (affected offspring). We identified regions of suggestive association on chromosomes 6q27 and 20p13, respectively. However, linkage analysis did not yield genome-wide significant associations. However, genotyping of top hits in additional families revealed a SNP on chromosome 5p15 (between *SEMA5A* and *TAS1R3*) was significantly associated with autism ( $P = 2 \times 10^{-7}$ ). This SNP demonstrated that expression of *SEMA5A* is reduced in brain tissue from autistic patients, further implicating *SEMA5A* as an autism susceptibility gene. The linkage regions reported here provide targets for rare variation screening whereas the discovery of a significant association demonstrates the action of common variants.

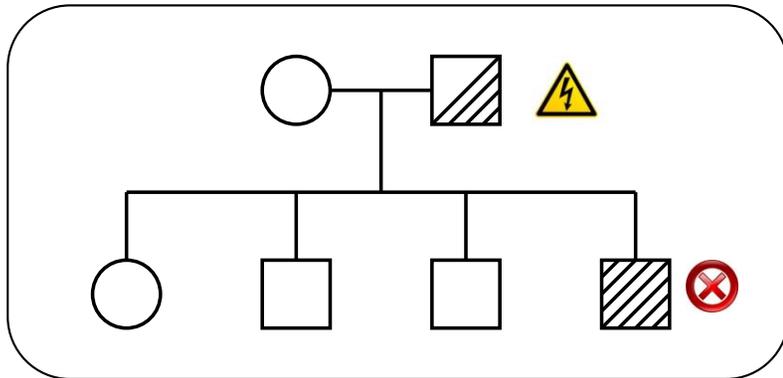
For a high-resolution genetic study of autism, we selected



**Figure 2 | *SEMA5A* expression in autism brains.** *SEMA5A* gene expression is shown relative to *MAP2*. Diamonds indicate individual expression levels for each sample; error bars indicate standard error (s.e.).

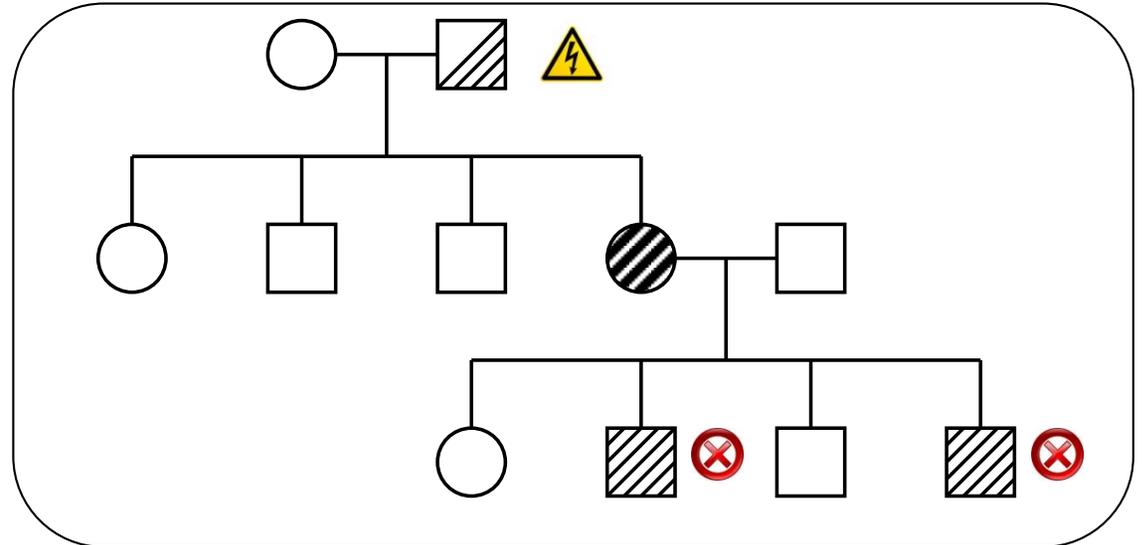
# Unified Model of Autism

## Sporadic Autism



De novo mutations of high penetrance contributes to autism, especially in low risk families with no history of autism.

## Familial Autism



### Legend



Sporadic mutation



Fails to procreate

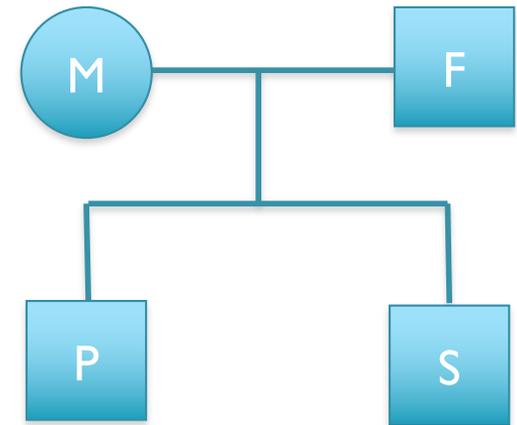
**A unified genetic theory for sporadic and inherited autism**

Zhao et al. (2007) *PNAS*. 104(31)12831-12836.

# De novo mutation discovery and validation

**Concept:** Identify mutations not present in parents.

**Challenge:** Sequencing errors in the child or low coverage in parents lead to false positive de novos



**Reference:** ...TCAAATCCTTTTAATAAAGAAGAGCTGACA...

**Father:** ...TCAAATCCTTTTAATAAAGAAGAGCTGACA...

**Mother:** ...TCAAATCCTTTTAATAAAGAAGAGCTGACA...

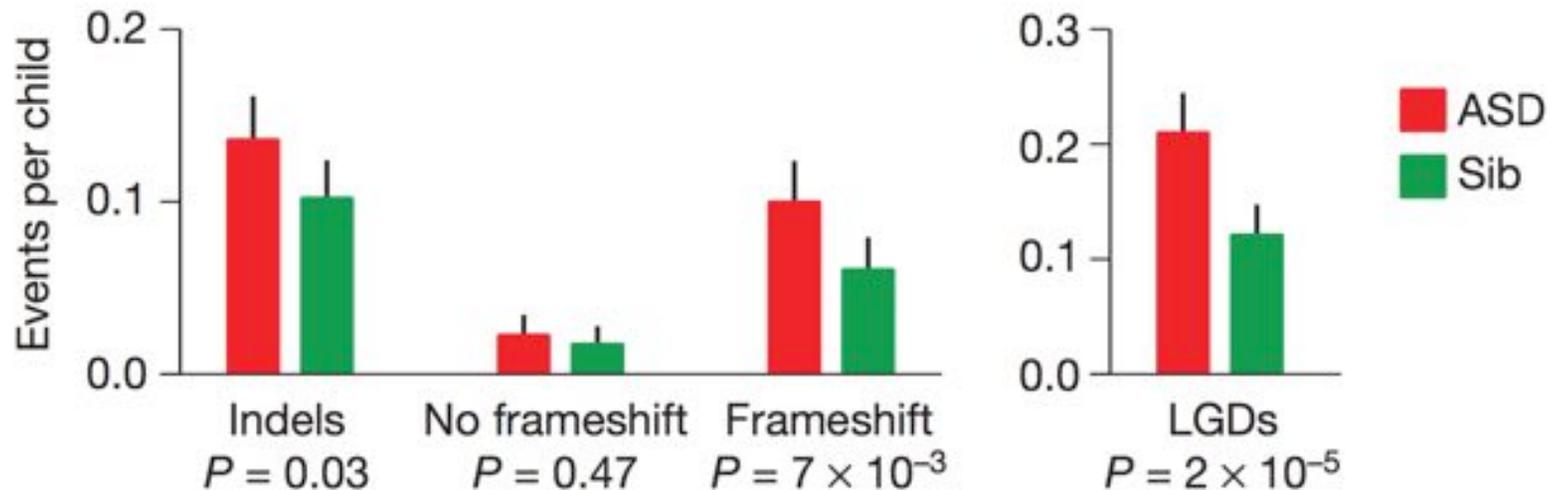
**Sibling:** ...TCAAATCCTTTTAATAAAGAAGAGCTGACA...

**Proband(1):** ...TCAAATCCTTTTAATAAAGAAGAGCTGACA...

**Proband(2):** ...TCAAATCCTTTTAAT\*\*\*AAGAGCTGACA...

4bp heterozygous deletion at chr15:9352406 | CHD2

# De novo Genetics of Autism

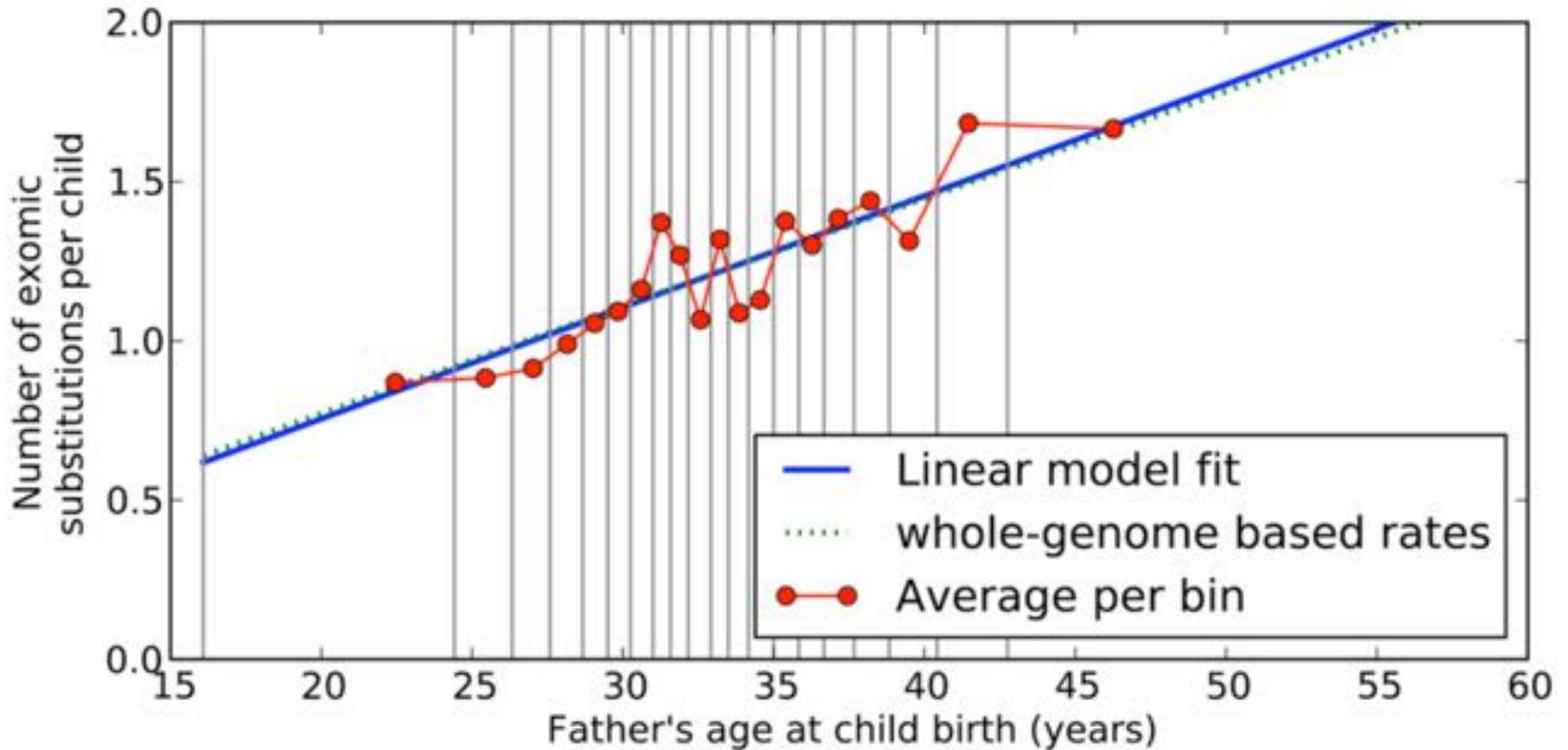


- In 2,500 family quads we see significant enrichment in de novo **likely gene disruptions (LGDs)** in the autistic children
  - Overall rate of de novo mutations basically 1:1
  - 2:1 enrichment in frameshift indels, nonsense mutations
  - Contributed dozens of new autism candidate genes, highly enriched for neuron development or chromatin modifiers

**The contribution of de novo coding mutations to autism spectrum disorder**  
Iossifov et al (2014) *Nature*. doi:10.1038/nature13908

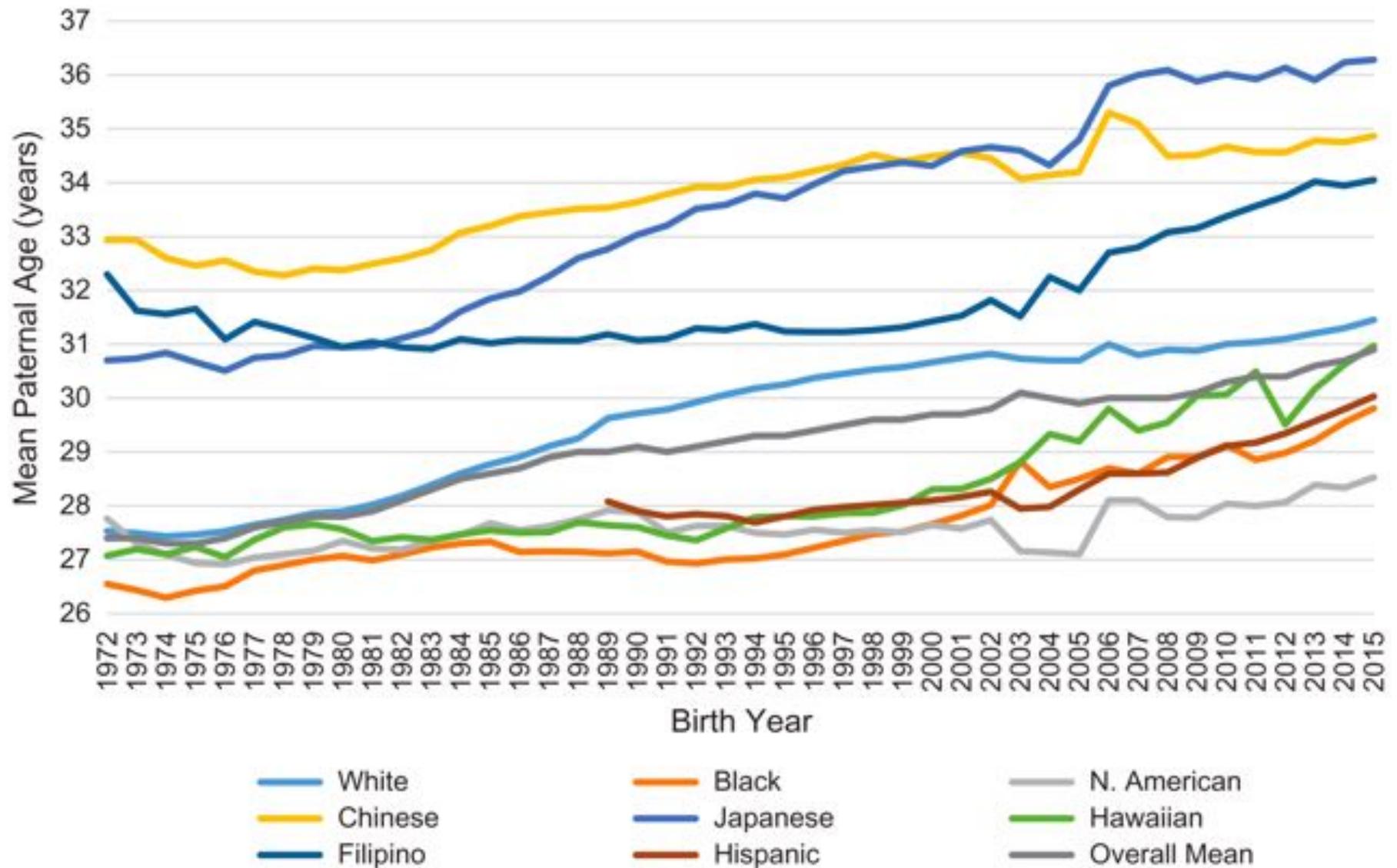


# De novo Mutations in Men



**The contribution of de novo coding mutations to autism spectrum disorder**  
Iossifov *et al* (2014) *Nature*. doi:10.1038/nature13908

# Age of Fatherhood



**The age of fathers in the USA is rising: an analysis of 168 867 480 births from 1972 to 2015**

Khandwala et al (2017) *Human Reproduction*. <https://doi.org/10.1093/humrep/dex267>