

# Genome-wide association studies (GWAS)

Matti Pirinen

FIMM, University of Helsinki

Manuel A. Rivas

University of Oxford and  
The Broad Institute

15 January 2015

NIASC event, Helsinki

# Part I: Common variants

- 1. Motivation
  - What is a “genetic association” ?
  - Why is this important ?
- 2. How to read the genome (for SNPs) ?
- 3. An Example GWAS of Multiple Sclerosis
- 4. Interpreting a GWAS
- 5. Current state of GWAS
  - Criticism

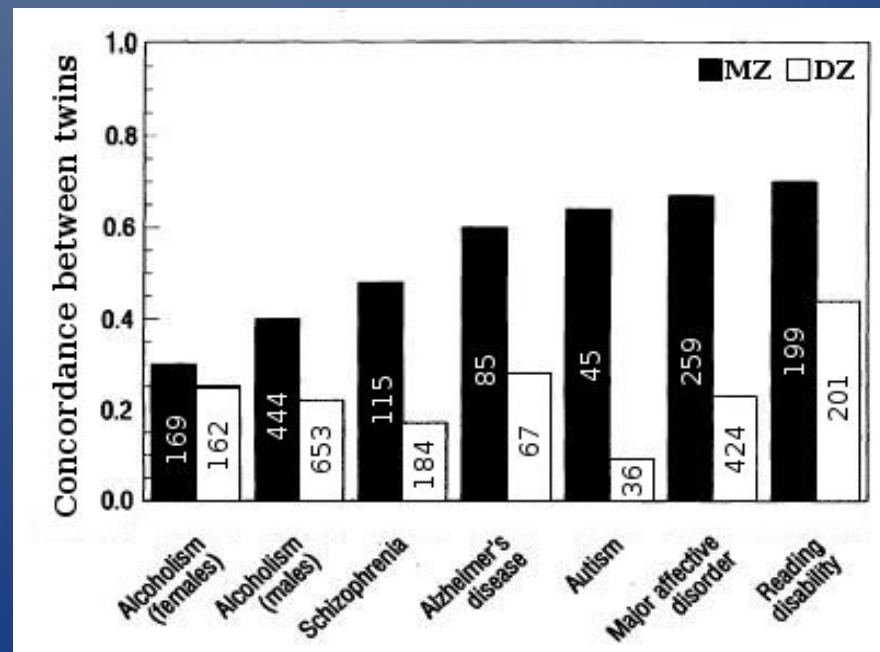
# Part I: Common variants

- 1. Motivation
  - Genetics contribute to traits
  - What is a “genetic association” ?
  - Why is this important ?

# Does genetics affect trait?

## Is the trait *heritable*?

- Do more close relatives have more similar phenotypes (on average)?
  - Twin studies compare monozygotic twins with dizygotic twins (but environment may confound)



# Human Genome

... G C G T T T A C G ... DNA sequence

A Human genome is  $3 \times 10^9$  letters from alphabet {A, C, G, T}

# Single Nucleotide Polymorphism (SNP)

- On average, 1:300 positions has (common) variation in population, called “SNP”

# Single Nucleotide Polymorphism (SNP)

- On average, 1:300 positions has (common) variation in population, called “SNP”

Genomes in population		Individuals in population	
... G C G T T ...	96%	0: GG	~ 92.1%
... G C T T T ...	4%	1: GT	~ 7.7 %
		2: TT	~ 0.2 %

↑  
SNP, alleles: G / T, minor allele frequency (MAF) = 4%

# PCSK9 gene

- Location: Chromosome 1, 55.50 – 55.53 Mb



**Codes protein**



692 Amino acids

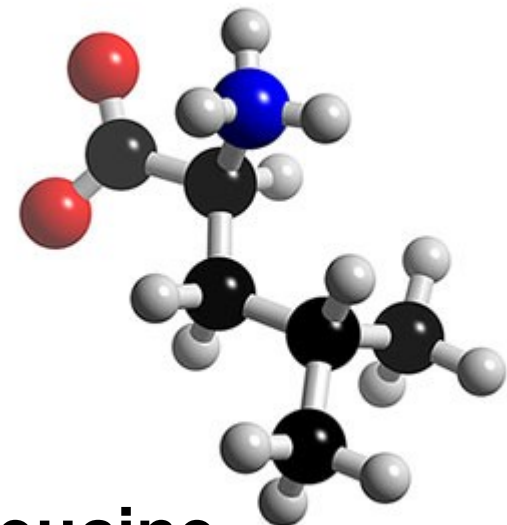
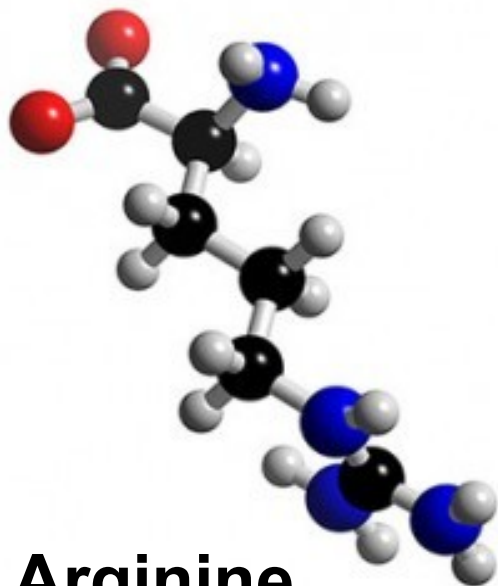


# A SNP in PCSK9 gene

- Alleles: G / T , MAF=4% (Finland)
- Location: Chromosome 1, position 55,505,647

# A SNP in PCSK9 gene

- Alleles: G / T , MAF=4% (Finland)
- Location: Chromosome 1, position 55,505,647
- Function: “missense”, “nonsynonymous”, changes 46th AA from Arginine to Leucine

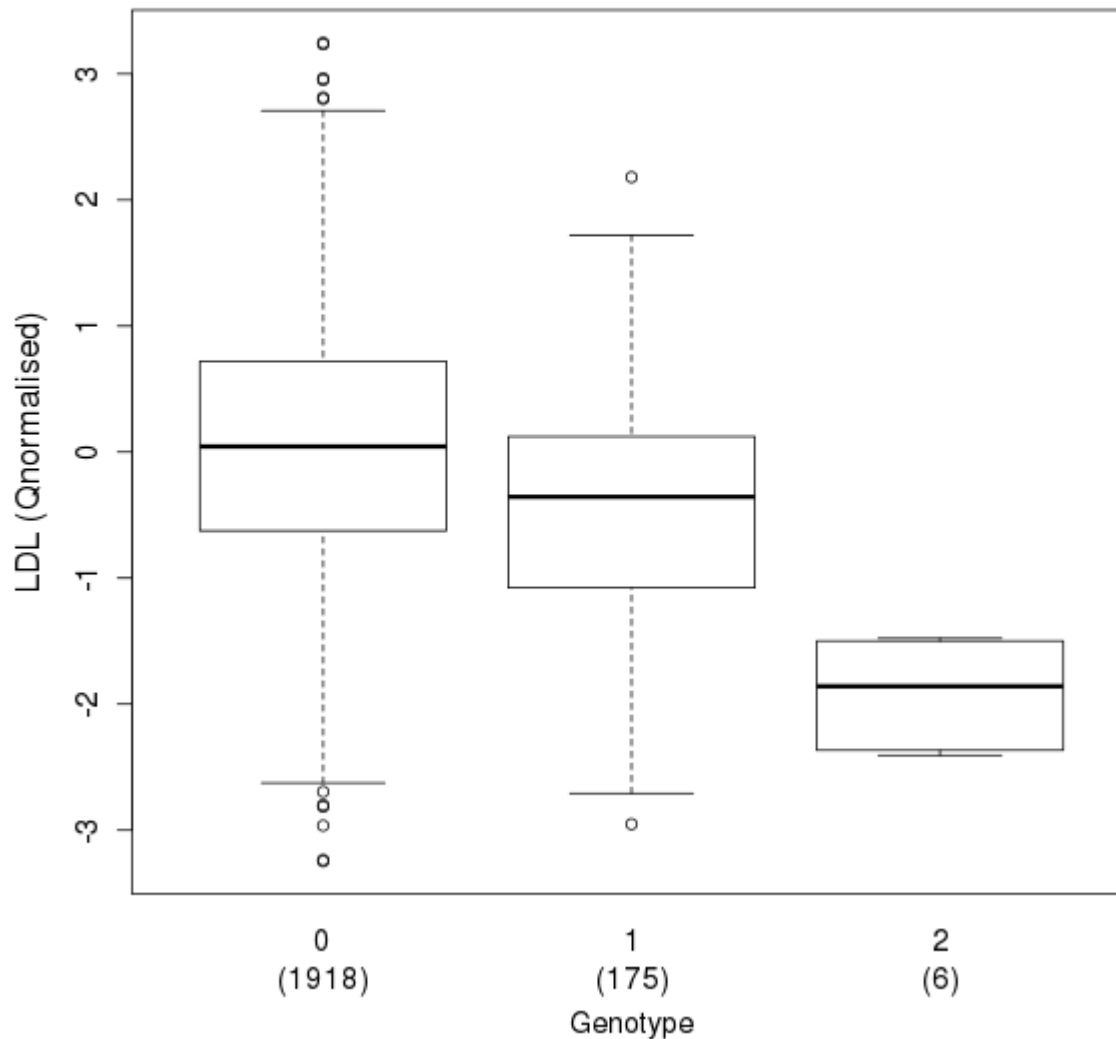


# A SNP in PCSK9 gene

- Alleles: G / T , MAF=4% (Finland)
- Location: Chromosome 1, position 55,505,647
- Function: “missense”, “nonsynonymous”, changes 46th AA from Arginine to Leucine
- Let's see if there is association with LDL cholesterol
  - LDL-C is a risk factor for heart disease

# What is a “genetic association”?

PSCK9 NONSYN;  $b=-0.57$ ,  $p=1e-14$ ;  $N=2099$



GG

GT

TT

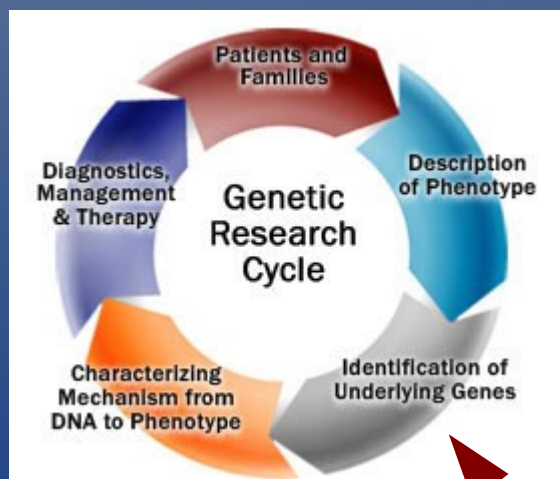
Finn-Metabo-Seq project:  
2099 Finns exome  
seq'ed (Aug 2014)

Boxplot shows  
(1) medians (thick lines),  
(2) interquartile range  
(boxes)  
(3) 1.5 x interquartile  
range (dotted segments)  
(4) outliers (points)

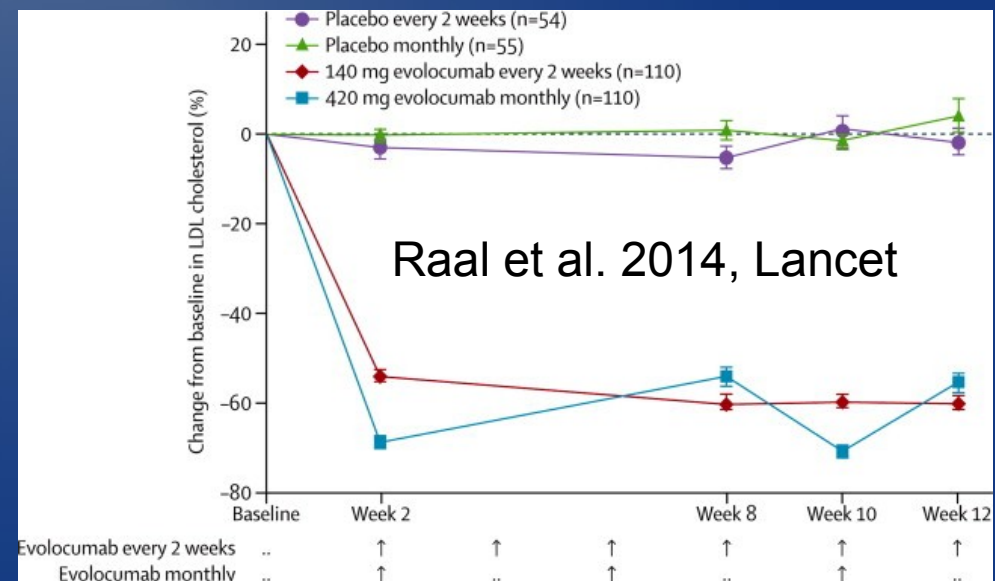
Carriers of allele T have  
lower LDL-C

# Why are genetic associations important?

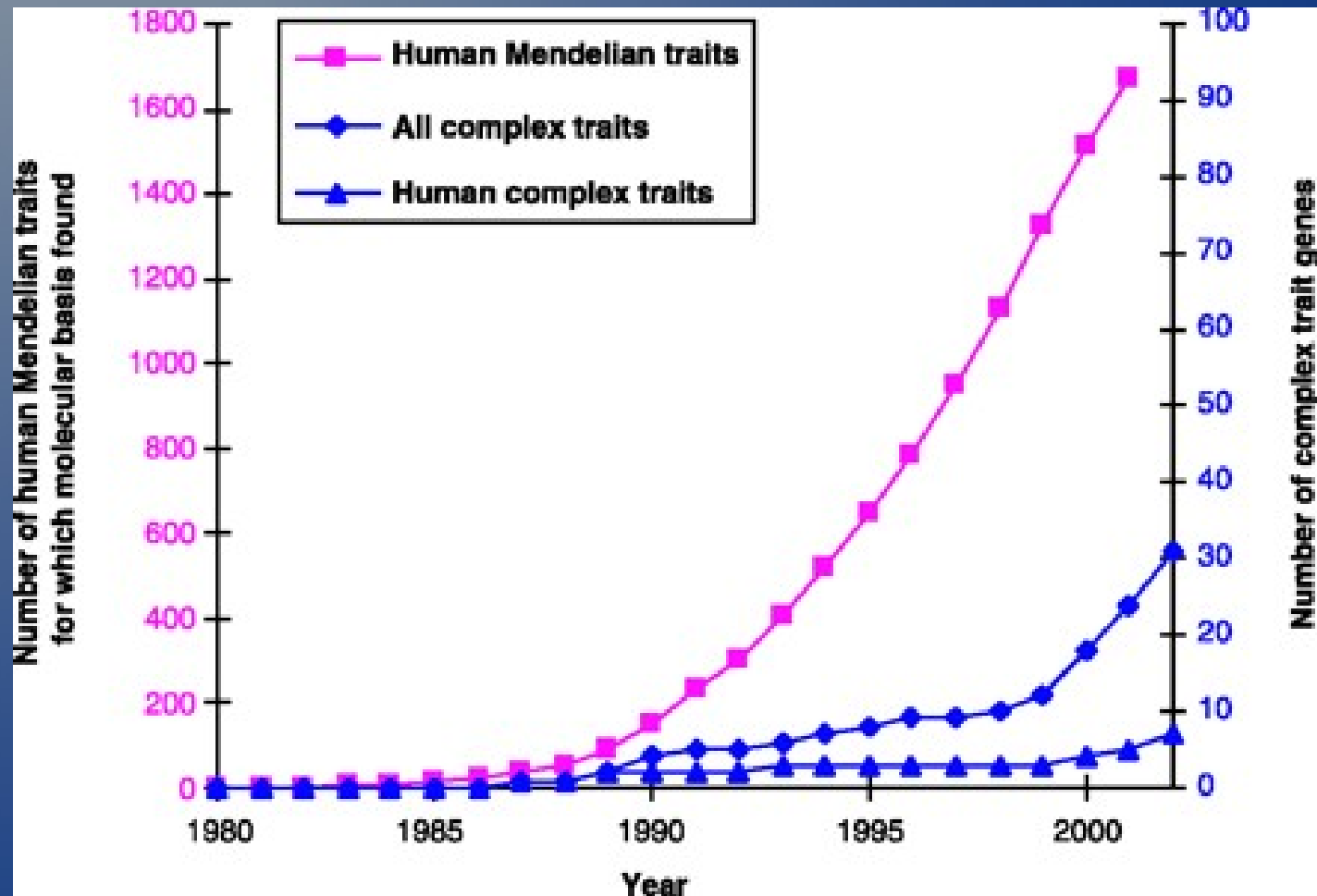
- Hints of biology behind the diseases and traits
  - Later: examples of Multiple sclerosis and schizophrenia
- Hints of targets for therapeutics
  - Inhibition of PCSK9 reduces LDL-C ?



[www.chgr.org](http://www.chgr.org)

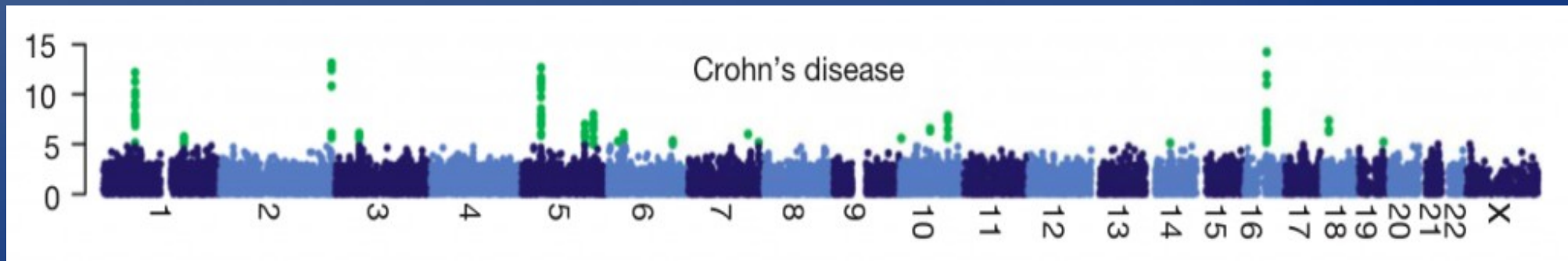


# Large effect variants were found for many Mendelian traits but not for complex traits during 1985-2000



# Genome-wide association studies (~2006 onwards )

- Idea: To look for associations in a detailed map of common variation across the genome
  - “Common disease – common variant”
- Became possible through
  - Technology (SNP arrays, later sequencing)
  - Collaboration (genetics + medicine + lab tech + bioinformatics + statistics)



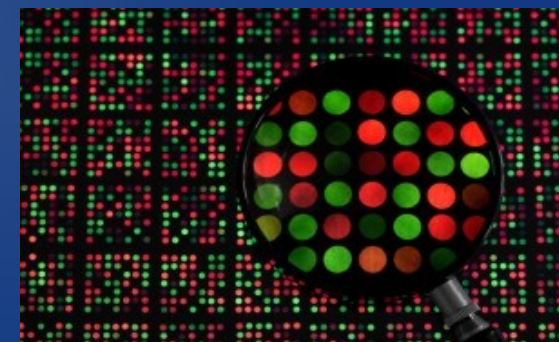
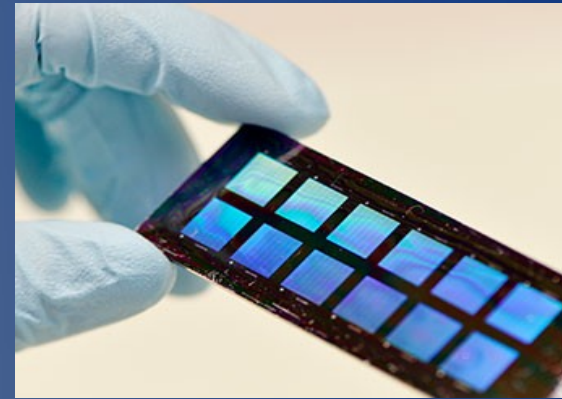
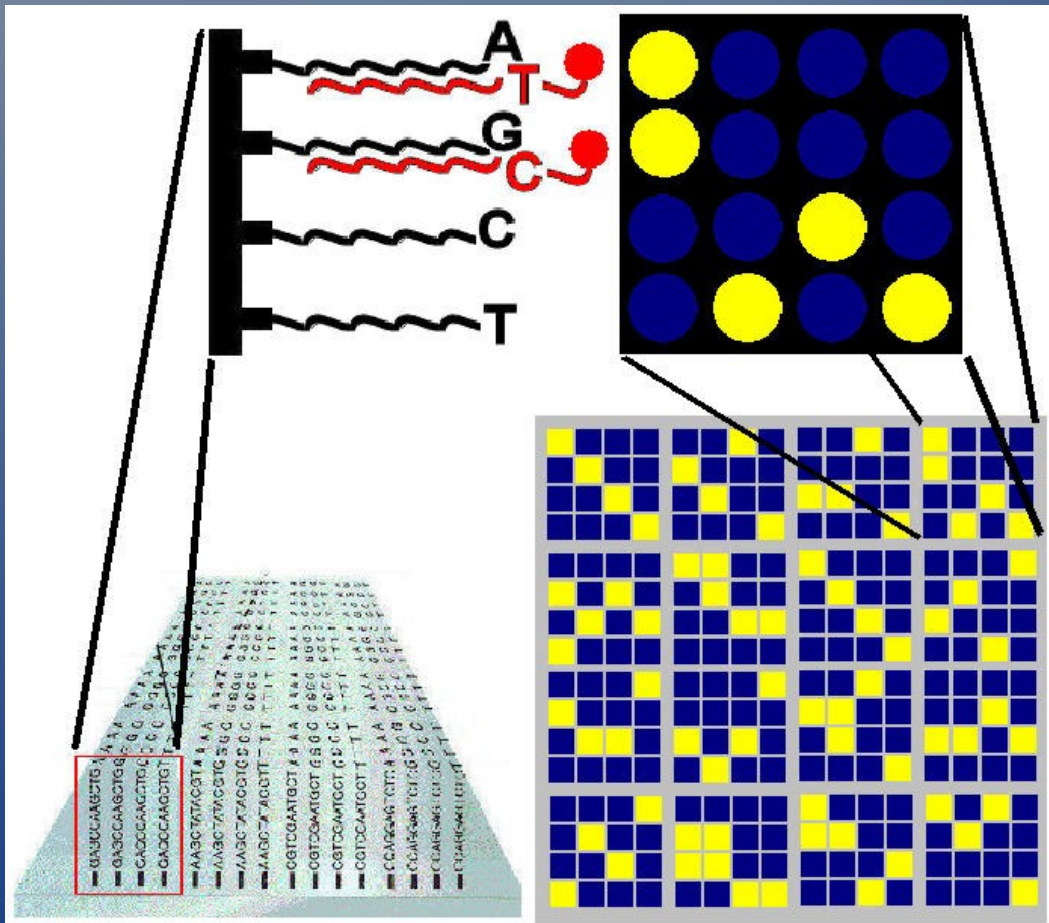
# Part I: Common variants

- 1. Motivation
  - What is a “genetic association” ?
  - Why is this important ?
- 2. How to read the genome (for SNPs) ?



# Human SNP arrays

- Contain probes for several million SNPs
- Price ~50-100 euros/sample

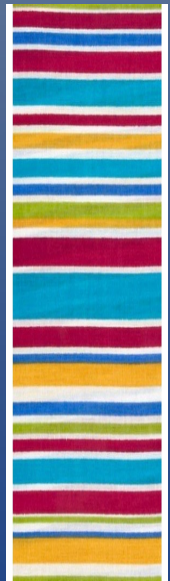


Steven M. Carr  
[www.mun.ca/biology/scarr/DNA\\_Chips.html](http://www.mun.ca/biology/scarr/DNA_Chips.html)

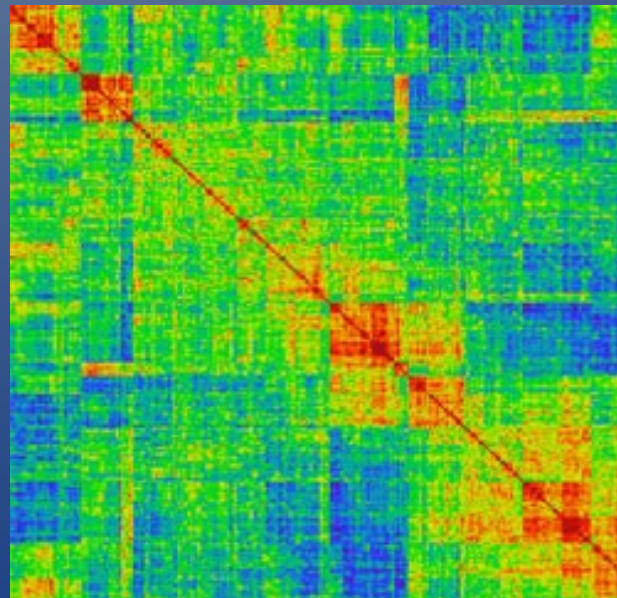
# Does genetics affect trait?

## Is the trait *heritable*?

- With detailed genetic data, the question can be asked with “unrelated” individuals (program GCTA, Yang et al. 2010, Nat Gen)
  - Environment is not an (obvious) confounder



~  
?



Example:  
For height in 14,500  
Finnish samples,  
genetic component  
of common SNPs  
explains 52% of  
variation.

Y

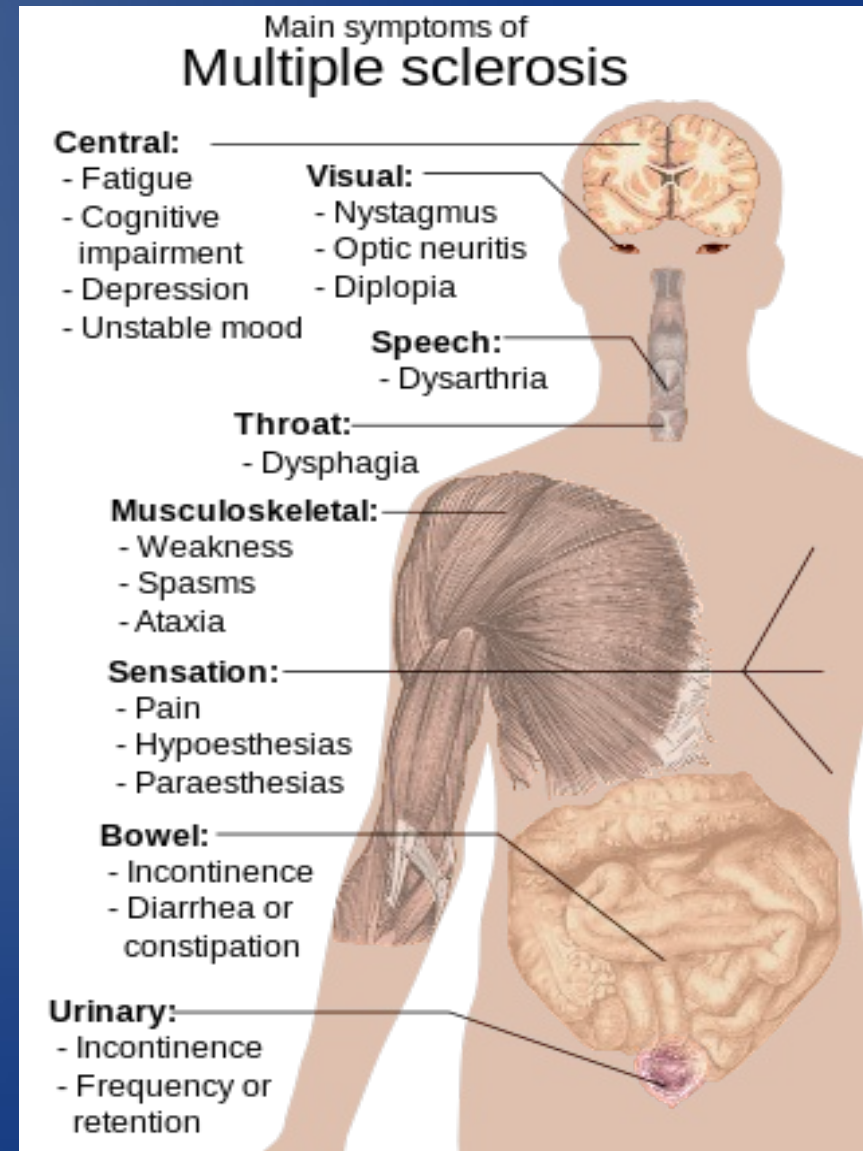
G

# Part I: Common variants

- 1. Motivation
  - What is a “genetic association” ?
  - Why is this important ?
- 2. How to read the genome (for SNPs) ?
- 3. An Example GWAS of Multiple Sclerosis

# An example GWAS: Multiple Sclerosis by WTCCC2

- MS: nerve cells are damaged through inflammation reaction -> wide range of severe symptoms
- A GWAS on 27,000 samples from 15 countries
  - 10,000 cases
  - 17,000 controls
- Nature, August 2011



# Case-control association study

Individuals

$\sim 10^4$   
Cases

$\sim 10^4$   
Controls

Genotypes

$\sim 10^6$  Single nucleotide polymorphisms (SNPs)

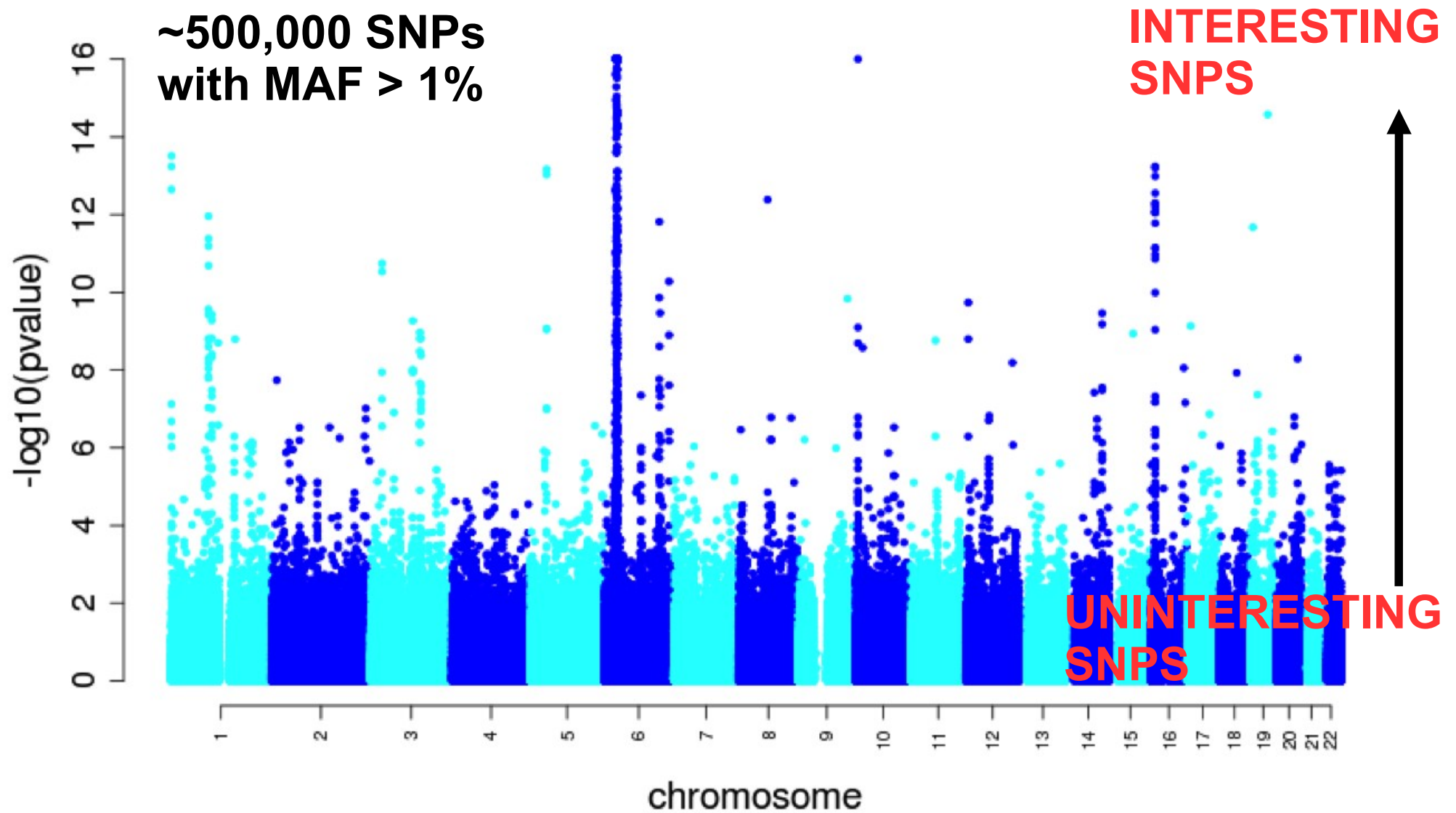


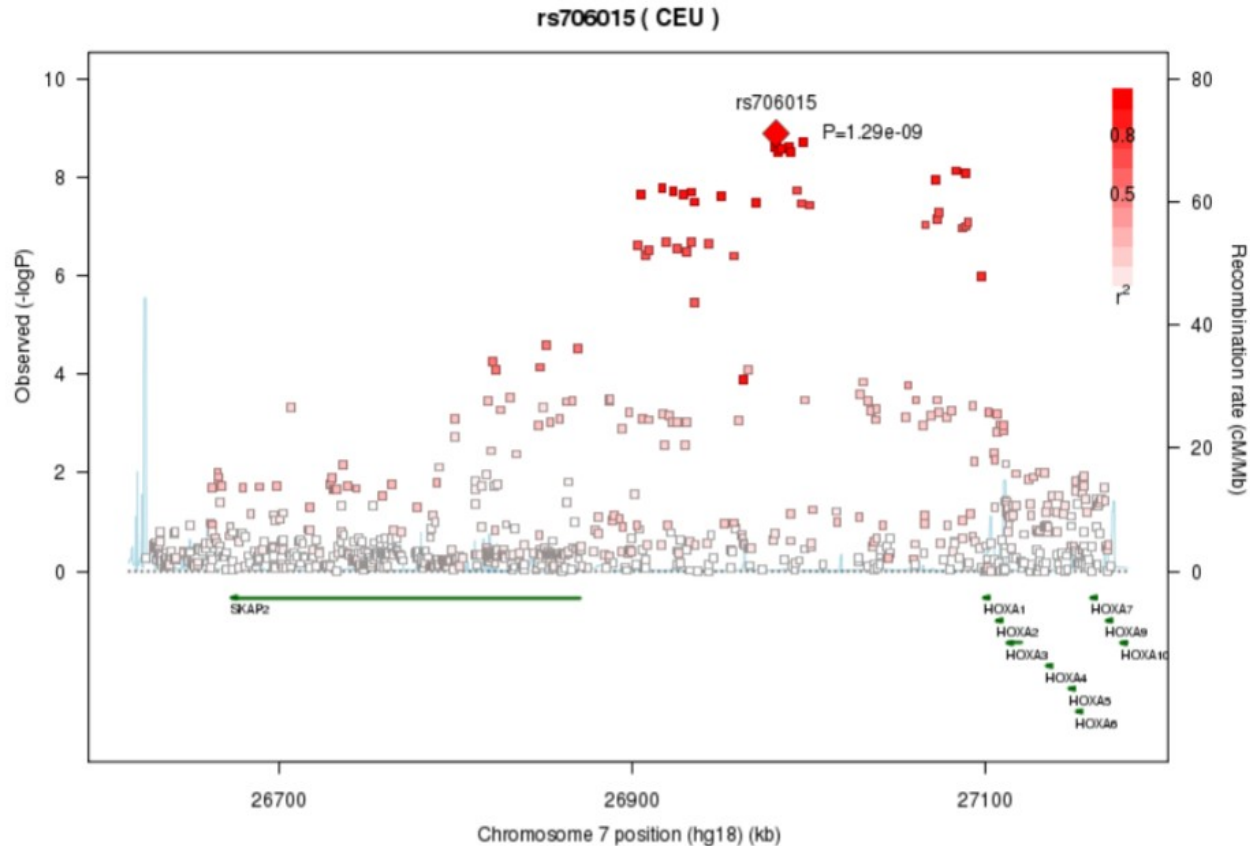
Question

Are the genotype distributions different between cases and controls?



# A genome-wide scan

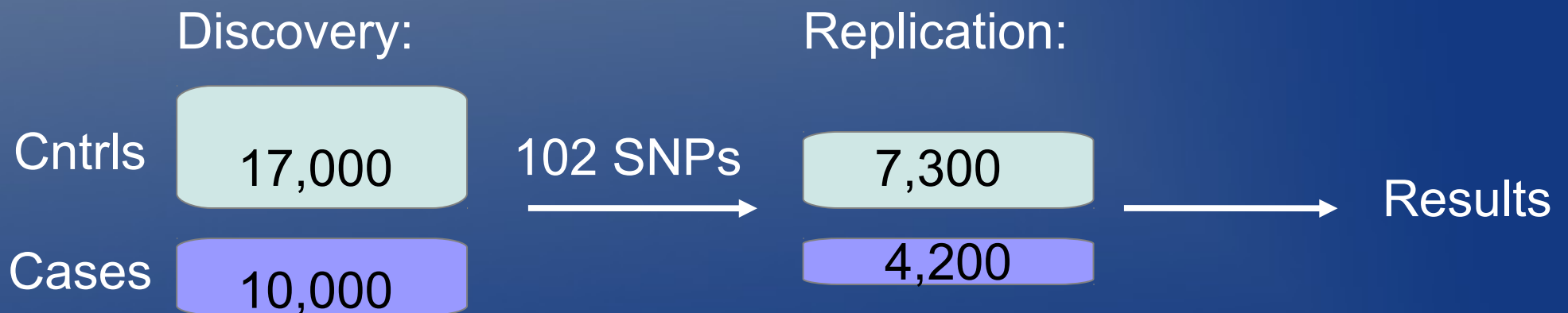




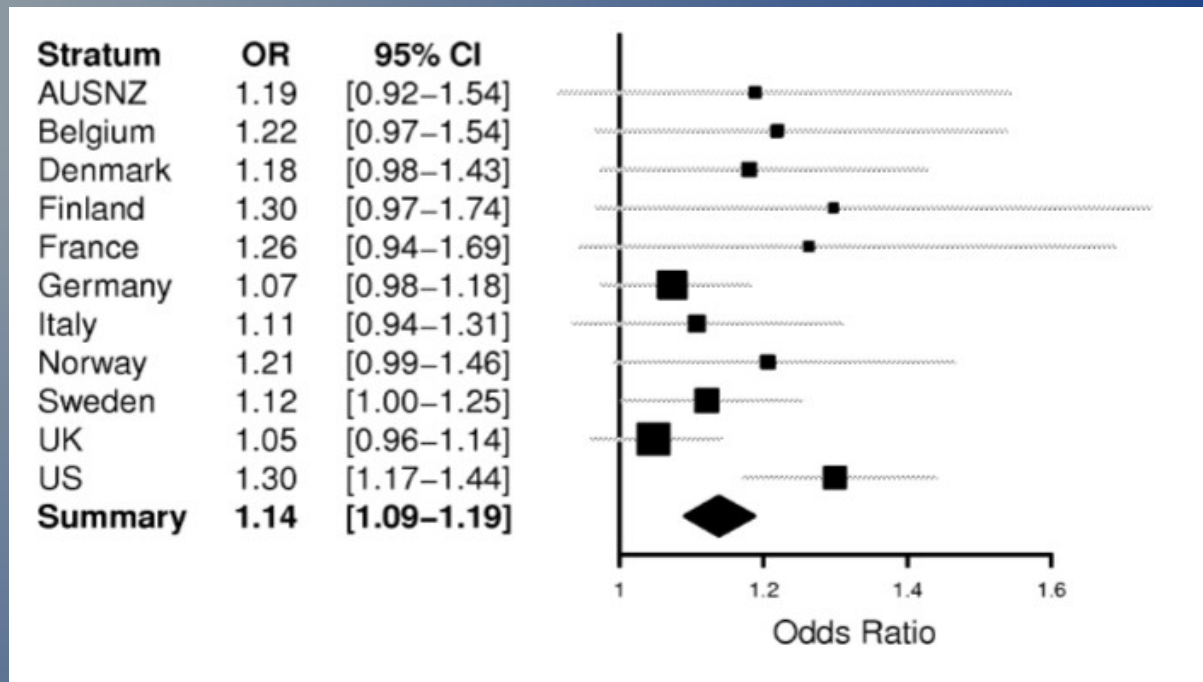
Each "hit SNP" resides in a region with many correlated SNPs and possibly many (or none) genes

# Replication

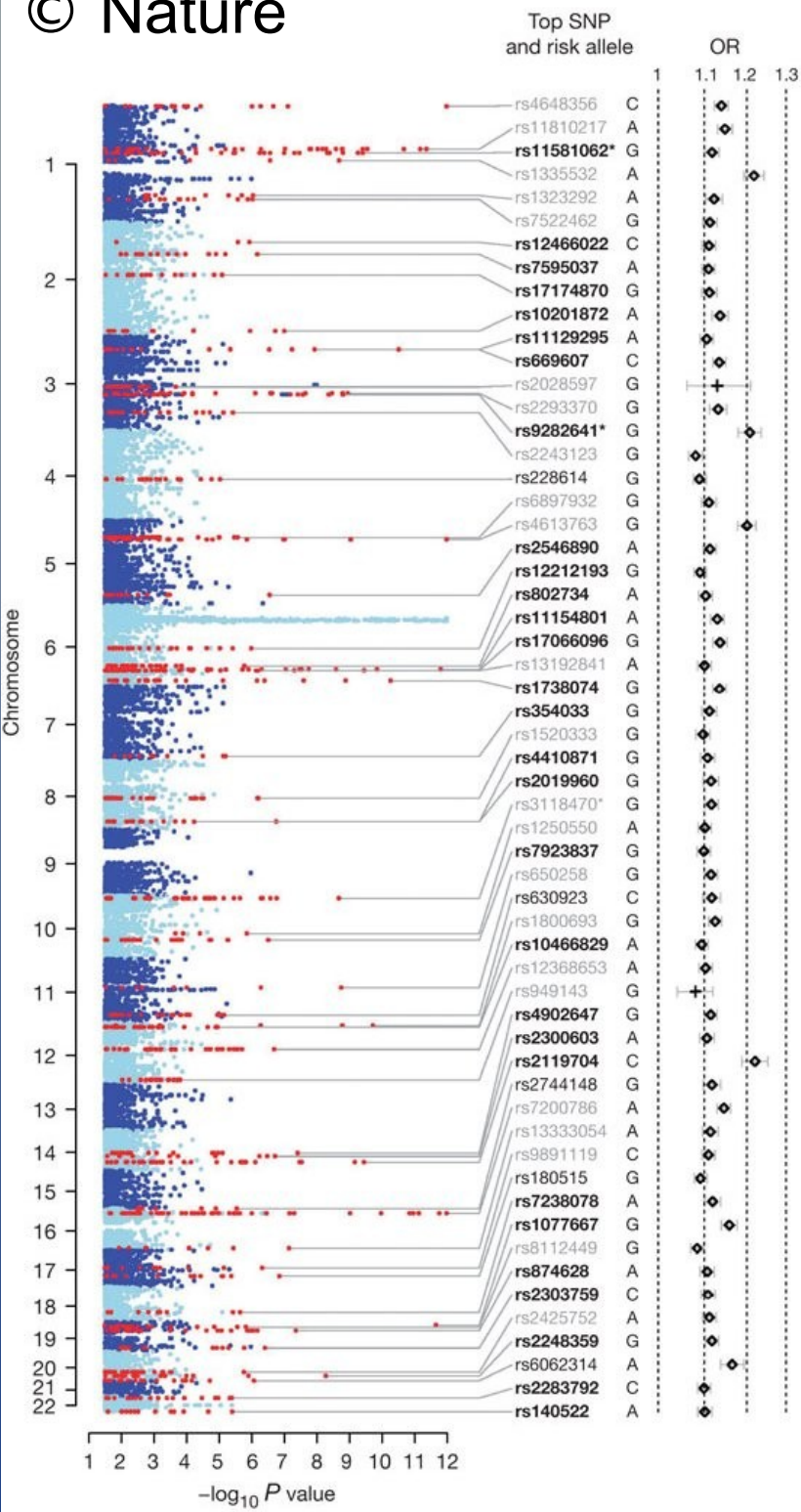
- To verify the most promising associations in an independent data set
- IN MS study, took 102 SNPs to replication





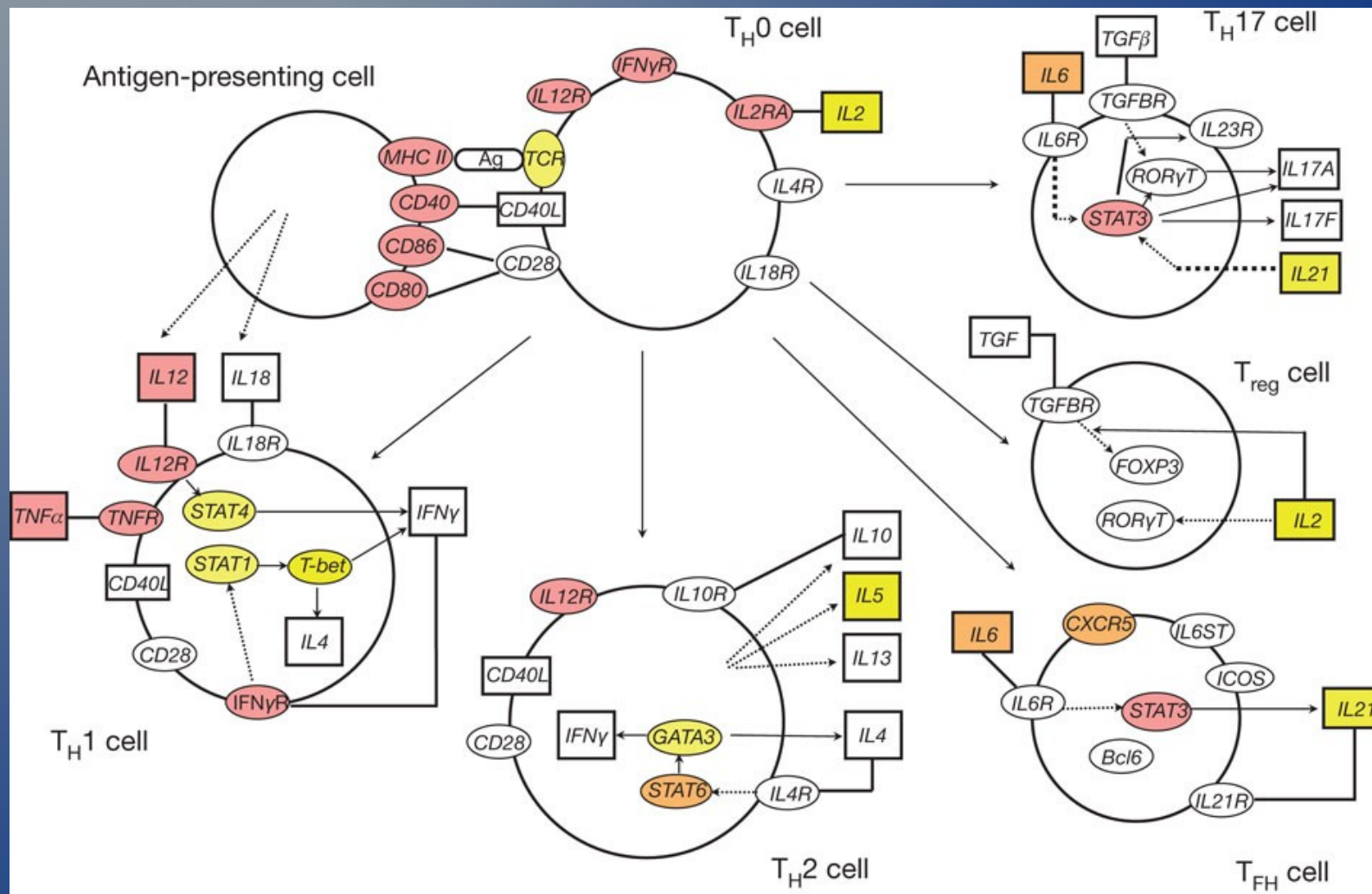


- Consistent effects across
  - populations
  - genotyping platformsbring confidence



- 98 /102 SNPs have consistent effects in replication data
- Over 50 convincing associations with MS
- Immunological genes are over represented among the hits; in particular, T-helper cell differentiation pathway

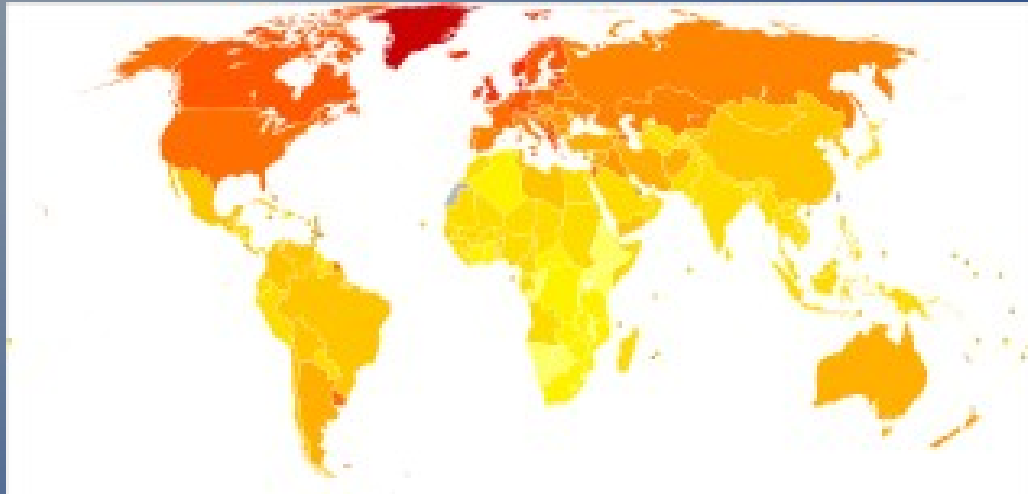
# T-helper cell pathway implicated by GWAS hits



Red, strong signal ;  
orange, medium signal  
yellow, some signal

# MS results

- Two implicated genes involved in vitamin D metabolism



Related to latitudinal variation in MS prevalence?

- Two other implicated genes are already targets of MS therapies
  - VCAM1 for 'natalizumab' and IL2RA for 'daclizumab'



Stephen Sawcer<sup>1\*</sup>, Garrett Hellenthal<sup>2\*</sup>, Matti Pirinen<sup>2\*</sup>, Chris C. A. Spencer<sup>2\*</sup>, Nikolaos A. Patsopoulos<sup>3,4,5</sup>, Loukas Moutsianas<sup>6</sup>, Alexander Dilthey<sup>6</sup>, Zhan Su<sup>2</sup>, Colin Freeman<sup>2</sup>, Sarah E. Hunt<sup>7</sup>, Sarah Edkins<sup>7</sup>, Emma Gray<sup>7</sup>, David R. Booth<sup>8</sup>, Simon C. Potter<sup>7</sup>, An Goris<sup>9</sup>, Gavin Band<sup>2</sup>, Annette Bang Oturai<sup>10</sup>, Amy Strange<sup>2</sup>, Janna Saarela<sup>11</sup>, Céline Belleguez<sup>7</sup>, Bertrand Fontaine<sup>12</sup>, Matthew Gillman<sup>7</sup>, Bernhard Hemmer<sup>13</sup>, Rhian Gwilliam<sup>7</sup>, Frauke Zipp<sup>14,15</sup>, Alagurevathi Jayakumar<sup>7</sup>, Roland Martin<sup>16</sup>, Stephen Leslie<sup>17</sup>, Stanley Hawkins<sup>18</sup>, Eleni Giannoulidou<sup>2</sup>, Sandra D'Alfonso<sup>19</sup>, Hannah Blackburn<sup>7</sup>, Filippo Martinelli Boneschi<sup>20</sup>, Jennifer Liddle<sup>7</sup>, Hanne F. Harbo<sup>21,22</sup>, Marc L. Perez<sup>7</sup>, Anne Spurkland<sup>23</sup>, Matthew J. Waller<sup>7</sup>, Marcin P. Mycko<sup>24</sup>, Michelle Ricketts<sup>7</sup>, Manuel Comabella<sup>25</sup>, Naomi Hammond<sup>7</sup>, Ingrid Kockum<sup>26</sup>, Owen T. McCann<sup>7</sup>, Maria Ban<sup>1</sup>, Pamela Whittaker<sup>7</sup>, Anu Kemppinen<sup>1</sup>, Paul Weston<sup>7</sup>, Clive Hawkins<sup>27</sup>, Sara Widaa<sup>7</sup>, John Zajicek<sup>28</sup>, Serge Dronov<sup>7</sup>, Neil Robertson<sup>29</sup>, Suzannah J. Bumpstead<sup>7</sup>, Lisa F. Barcellos<sup>30,31</sup>, Rathi Ravindrarajah<sup>7</sup>, Roby Abraham<sup>27</sup>, Lars Alfredsson<sup>32</sup>, Kristin Ardlie<sup>4</sup>, Cristin Aubin<sup>7</sup>, Amie Baker<sup>7</sup>, Katharine Baker<sup>29</sup>, Sergio E. Baranzini<sup>33</sup>, Laura Bergamaschi<sup>19</sup>, Roberto Bergamaschi<sup>34</sup>, Allan Bernstein<sup>31</sup>, Achim Berthele<sup>13</sup>, Mike Boggild<sup>35</sup>, Jonathan P. Bradfield<sup>36</sup>, David Brassat<sup>37</sup>, Simon A. Bradley<sup>38</sup>, Dorothea Buck<sup>13</sup>, Helmut Butzkueven<sup>39,40,41,42</sup>, Ruggero Capra<sup>43</sup>, William M. Carroll<sup>44</sup>, Paola Cavalla<sup>45</sup>, Elisabeth G. Celius<sup>21</sup>, Sabine Cepok<sup>13</sup>, Rosetta Chiavacci<sup>36</sup>, Françoise Clerget-Darpoux<sup>46</sup>, Kathleen Clysters<sup>9</sup>, Giancarlo Comi<sup>20</sup>, Mark Cossburn<sup>29</sup>, Isabelle Cournu-Rebeix<sup>12</sup>, Mathew B. Cox<sup>47</sup>, Wendy Cozen<sup>48</sup>, Bruce A. C. Cree<sup>33</sup>, Anne H. Cross<sup>49</sup>, Daniele Cusi<sup>50</sup>, Mark J. Daly<sup>4,51,52</sup>, Emma Davis<sup>53</sup>, Paul I. W. de Bakker<sup>3,4,54,55</sup>, Marc Debouverie<sup>56</sup>, Marie Beatrice D'hooghe<sup>57</sup>, Katherine Dixon<sup>53</sup>, Rita Dobosi<sup>9</sup>, Bénédicte Dubois<sup>9</sup>, David Ellinghaus<sup>58</sup>, Irina Elovaara<sup>59,60</sup>, Federica Esposito<sup>20</sup>, Claire Fontenille<sup>12</sup>, Simon Foote<sup>61</sup>, Andre Franke<sup>58</sup>, Daniela Galimberti<sup>62</sup>, Angelo Ghezzi<sup>63</sup>, Joseph Glessner<sup>36</sup>, Refugia Gomez<sup>33</sup>, Olivier Gout<sup>64</sup>, Colin Graham<sup>65</sup>, Struan F. A. Grant<sup>36,66,67</sup>, Franca Rosa Guerini<sup>68</sup>, Hakon Hakonarson<sup>36,66,67</sup>, Per Hall<sup>69</sup>, Anders Hamsten<sup>70</sup>, Hans-Peter Hartung<sup>71</sup>, Rob N. Heard<sup>8</sup>, Simon Heath<sup>72</sup>, Jeremy Hobart<sup>28</sup>, Muna Hoshi<sup>13</sup>, Carmen Infante-Duarte<sup>73</sup>, Gillian Ingram<sup>29</sup>, Wendy Ingram<sup>28</sup>, Talat Islam<sup>48</sup>, Maja Jagodic<sup>26</sup>, Michael Kabesch<sup>74</sup>, Allan G. Kermodé<sup>44</sup>, Trevor J. Kilpatrick<sup>1,39,40,75</sup>, Cecilia Kim<sup>36</sup>, Norman Klopp<sup>76</sup>, Keijo Koivisto<sup>77</sup>, Malin Larsson<sup>70</sup>, Mark Lathrop<sup>72</sup>, Jeannette S. Lechner-Scott<sup>4,7,78</sup>, Maurizio A. Leone<sup>79</sup>, Virpi Leppä<sup>11,80</sup>, Ulrika Liljedahl<sup>81</sup>, Izaura Lima Bomfim<sup>26</sup>, Robin R. Lincoln<sup>33</sup>, Jenny Link<sup>26</sup>, Jianjun Liu<sup>82</sup>, Aslaug R. Lorentzen<sup>22,83</sup>, Sara Lupoli<sup>50,84</sup>, Fabio Macchiardi<sup>50,85</sup>, Thomas Mack<sup>48</sup>, Mark Marriott<sup>39,40</sup>, Vittorio Martinelli<sup>20</sup>, Deborah Mason<sup>86</sup>, Jacob L. McCauley<sup>87</sup>, Frank Mentch<sup>36</sup>, Inger-Lise Mero<sup>21,83</sup>, Tania Mihalova<sup>27</sup>, Xavier Montalban<sup>25</sup>, John Mothershead<sup>88,89</sup>, Kjell-Morten Myhr<sup>90,91</sup>, Paola Naldi<sup>79</sup>, William Ollier<sup>53</sup>, Alison Page<sup>92</sup>, Aarno Palotie<sup>7,11,93,94</sup>, Jean Pelletier<sup>95</sup>, Laura Piccio<sup>49</sup>, Trevor Pickersgill<sup>29</sup>, Fredrik Piehl<sup>26</sup>, Susan Pobywajlo<sup>3</sup>, Hong L. Quach<sup>30</sup>, Patricia P. Ramsay<sup>30</sup>, Mauri Reunanen<sup>96</sup>, Richard Reynolds<sup>97</sup>, John D. Rioux<sup>98</sup>, Mariaemma Rodegher<sup>20</sup>, Sabine Roesner<sup>16</sup>, Justin P. Rubio<sup>39</sup>, Ina-Maria Rückert<sup>76</sup>, Marco Salvetti<sup>99</sup>, Erika Salvi<sup>50,100</sup>, Adam Santaniello<sup>33</sup>, Catherine A. Schaefer<sup>31</sup>, Stefan Schreiber<sup>58,101</sup>, Christian Schulze<sup>102</sup>, Rodney J. Scott<sup>47</sup>, Finn Selberg<sup>10</sup>, Krzysztof W. Selmaj<sup>24</sup>, David Sexton<sup>103</sup>, Ling Shen<sup>31</sup>, Brigid Simms-Acuna<sup>31</sup>, Sheila Skidmore<sup>1</sup>, Patrick M. A. Sleiman<sup>36,66</sup>, Cathrine Smestad<sup>21</sup>, Per Soelberg Sørensen<sup>10</sup>, Helle Bach Søndergaard<sup>10</sup>, Jim Stankovich<sup>61</sup>, Richard C. Strange<sup>27</sup>, Anna-Maija Sulonen<sup>11,80</sup>, Emilie Sundqvist<sup>26</sup>, Ann-Christine Syvänen<sup>81</sup>, Francesca Taddeo<sup>100</sup>, Bruce Taylor<sup>61</sup>, Jenefer M. Blackwell<sup>104,105</sup>, Pentti Tienari<sup>106</sup>, Elvira Bramer<sup>107</sup>, Ayman Tourbah<sup>108</sup>, Matthew A. Brown<sup>109</sup>, Ewa Tronczynska<sup>24</sup>, Juan P. Casas<sup>110</sup>, Niall Tubridy<sup>4,111</sup>, Aiden Corvin<sup>112</sup>, Jane Vickery<sup>28</sup>, Janusz Jankowski<sup>113</sup>, Pablo Villoslada<sup>114</sup>, Hugh S. Markus<sup>115</sup>, Kai Wang<sup>36,66</sup>, Christopher G. Mathew<sup>116</sup>, James Wason<sup>117</sup>, Colin N. A. Palmer<sup>118</sup>, H-Erich Wichmann<sup>7,6,119,120</sup>, Robert Plomin<sup>121</sup>, Ernest Willoughby<sup>122</sup>, Anna Rautanen<sup>2</sup>, Juliane Winkelmann<sup>1,3,123,124</sup>, Michael Wittig<sup>58,125</sup>, Richard C. Trembath<sup>116</sup>, Jacqueline Yao<sup>126</sup>, Ananth C. Viswanathan<sup>127</sup>, Haitao Zhang<sup>36,66</sup>, Nicholas W. Wood<sup>128</sup>, Rebecca Zuvich<sup>103</sup>, Panos Deloukas<sup>7</sup>, Cordelia Langford<sup>7</sup>, Audrey Duncanson<sup>129</sup>, Jorge R. Oksenberg<sup>33</sup>, Margaret A. Pericak-Vance<sup>87</sup>, Jonathan L. Haines<sup>103</sup>, Tomas Olsson<sup>26</sup>, Jan Hillert<sup>26</sup>, Adrian J. Iverson<sup>51,130</sup>, Philip L. De Jager<sup>4,5,51</sup>, Leena Peltonen-Toussainen<sup>1</sup>, Graeme J. Stewart<sup>8</sup>, David A. Hafler<sup>4,131</sup>, Stephen L. Hauser<sup>33</sup>, Gil McVean<sup>2</sup>, Peter Donnelly<sup>2,6\*</sup> & Alastair Compston<sup>1\*</sup>

# Collaboration

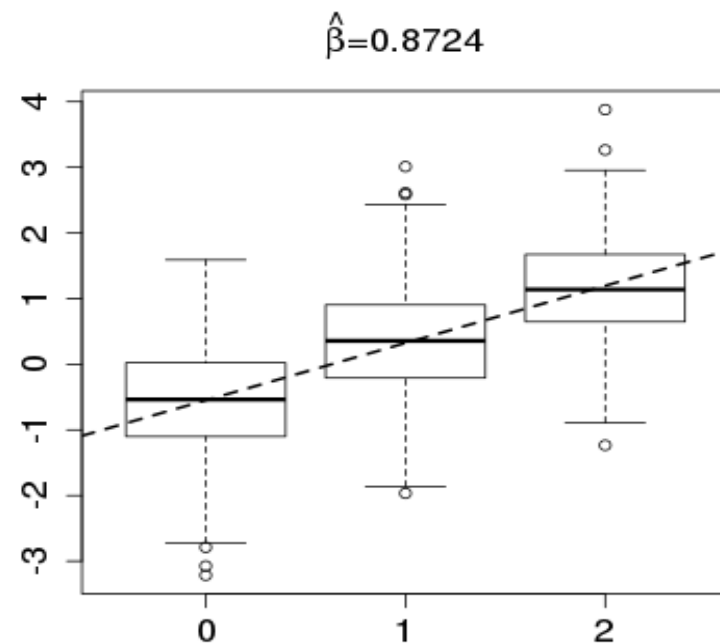
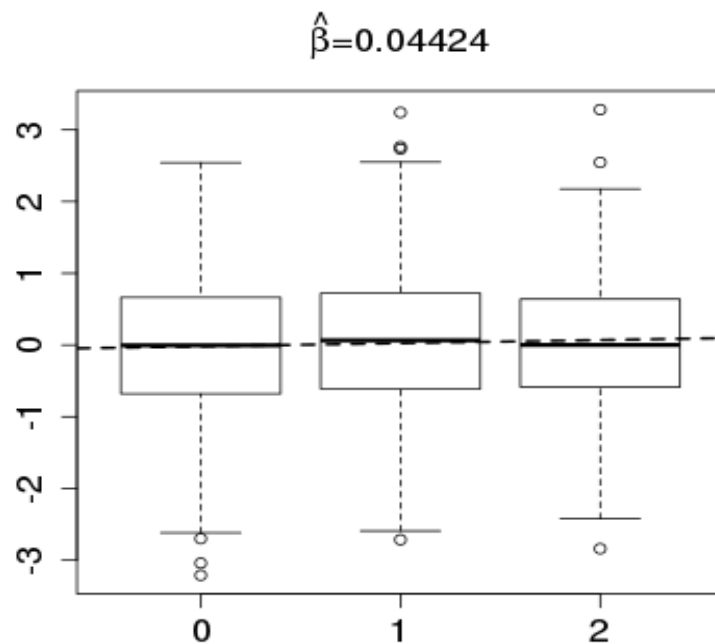
- 245 authors with 140 affiliations around the world

# Part I: Common variants

- 1. Motivation
  - What is a “genetic association” ?
  - Why is this important ?
- 2. How to read the genome (for SNPs) ?
- 3. An Example GWAS of Multiple Sclerosis
- 4. Interpreting a GWAS

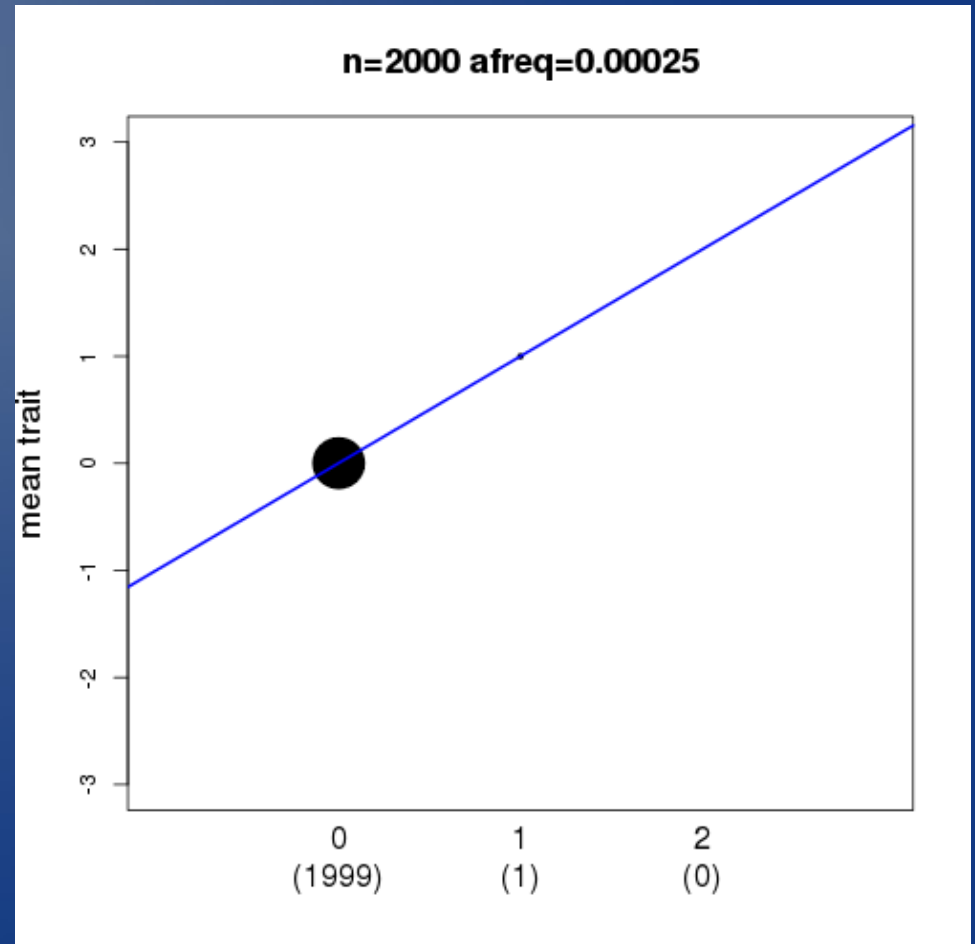
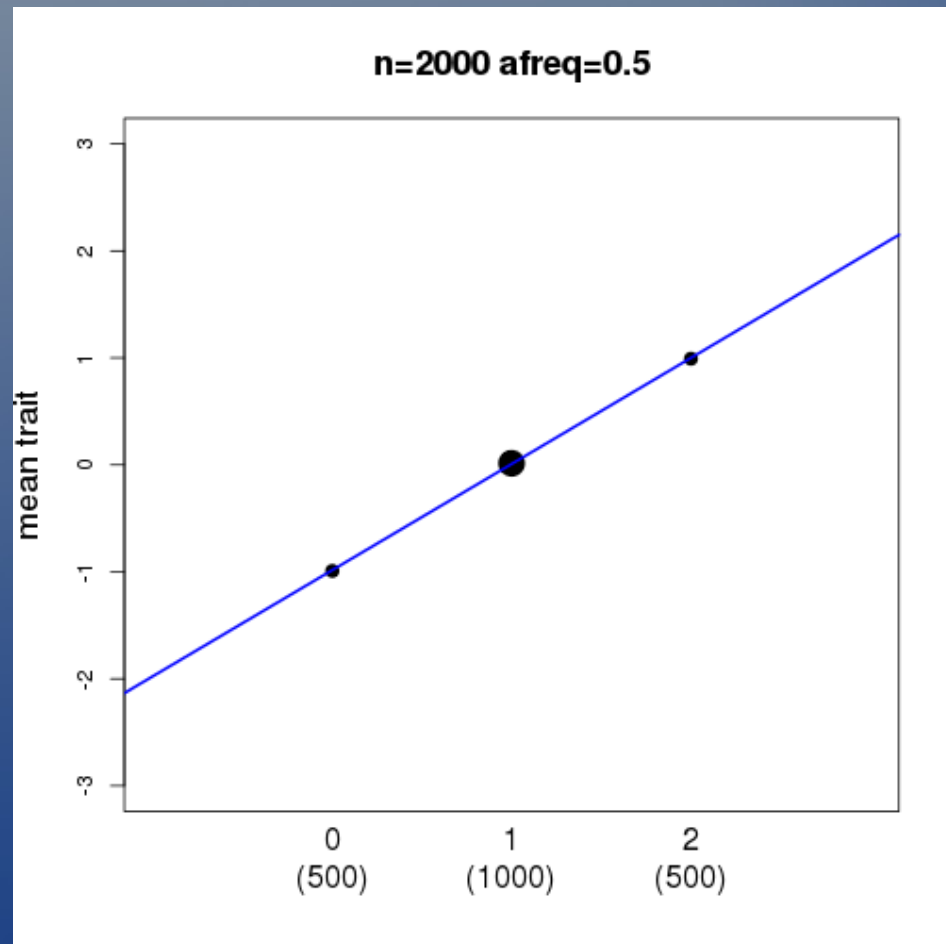
# Linear model to measure association

- Fit a line through 3 genotypes
  - Large slope = strong association (any prob's?)
  - Why is Manhattan plot not about |slope| but instead about “p-value”?



# Why is slope not everything?

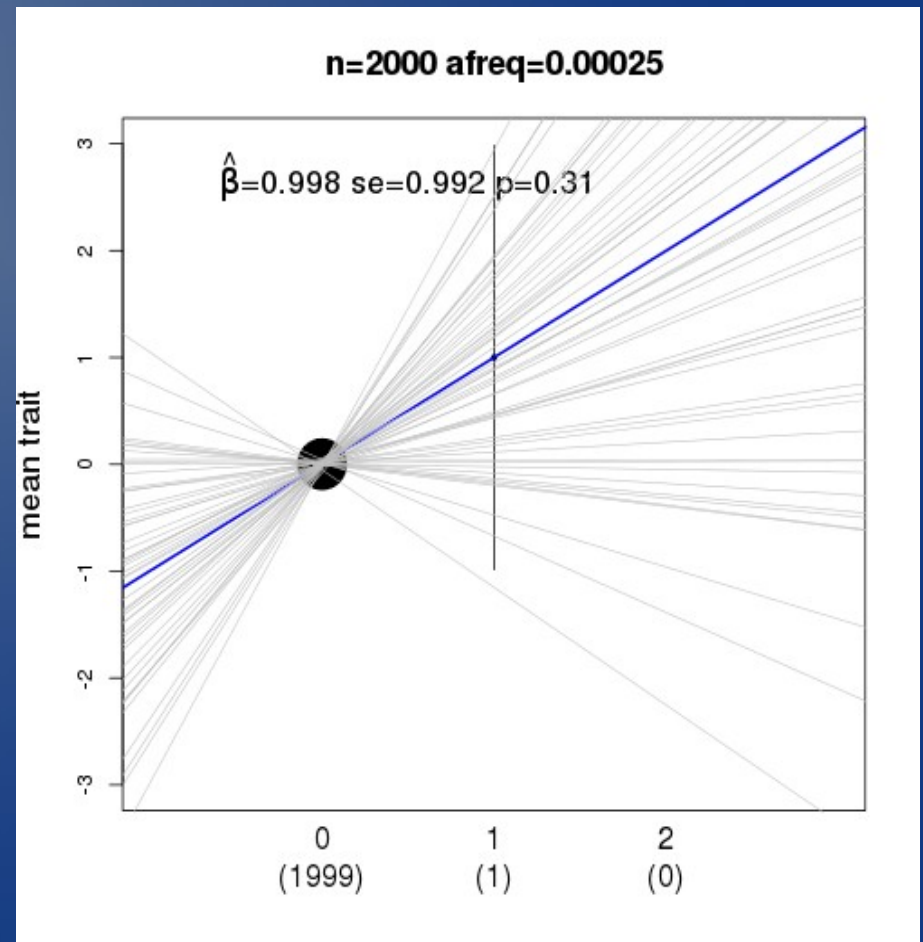
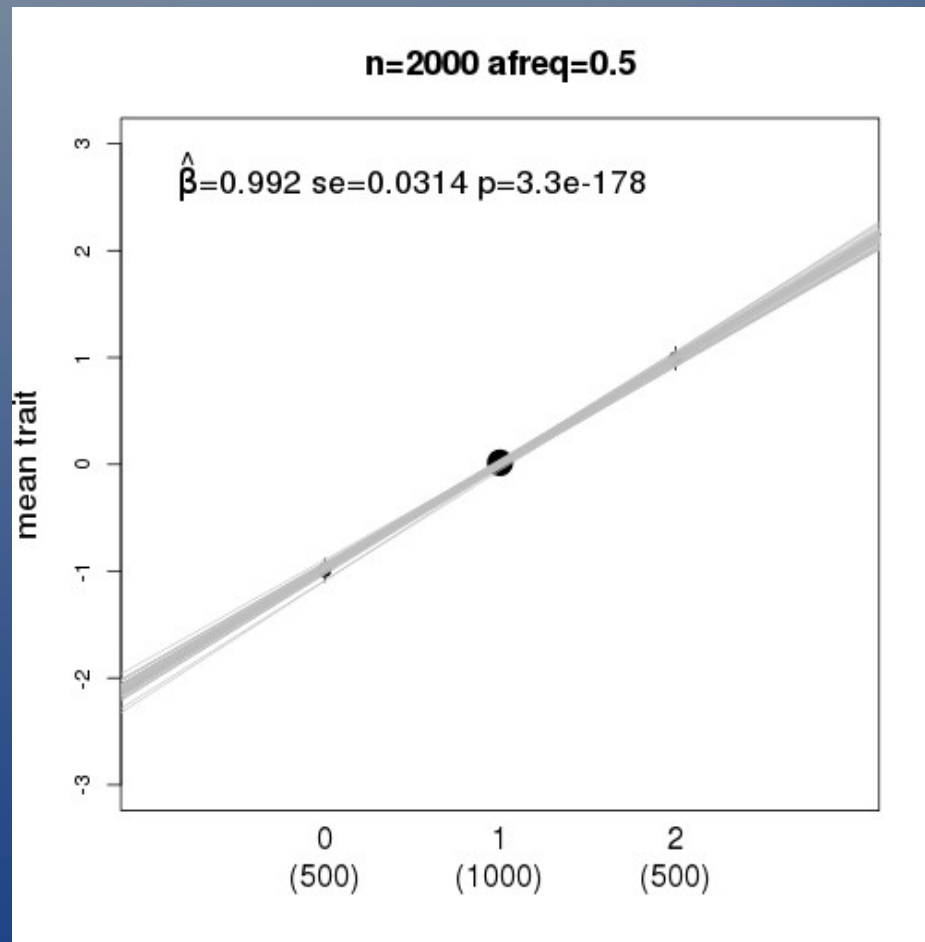
- Two SNPs with a slope of  $\sim 1.0$  for  $n=2000$





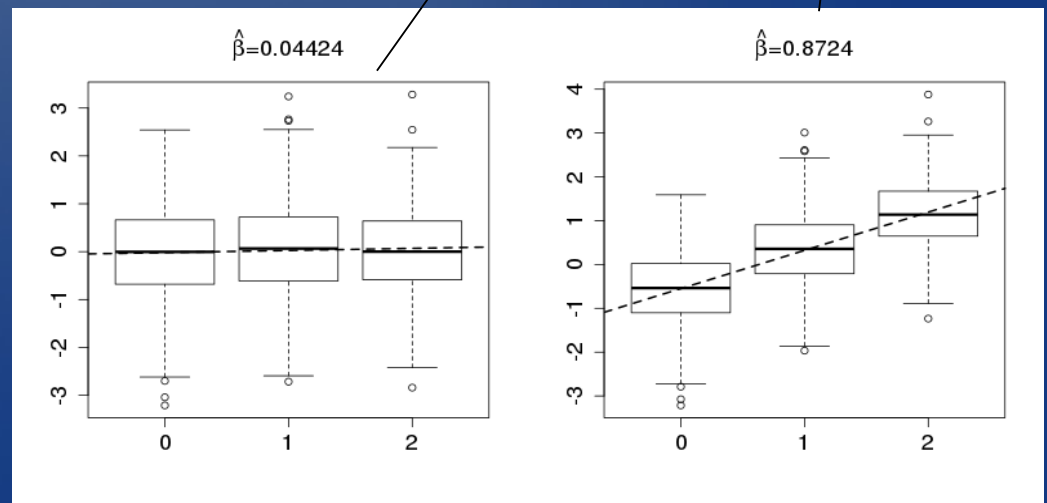
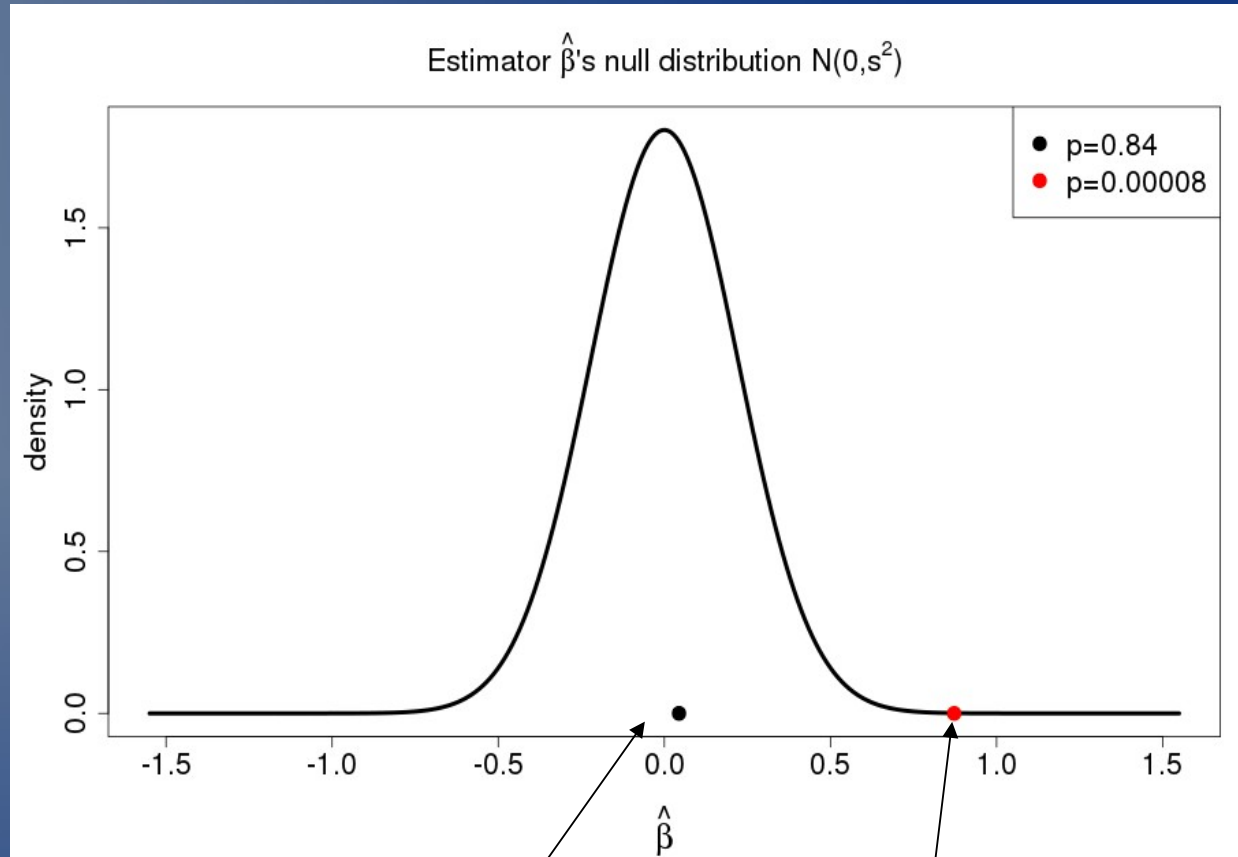
# Why is slope not everything?

- Uncertainty about the slope
- Left 1.0 (0.97 ... 1.03); Right 1.0 (-1.0 ... 3.0)




# P-value

- Is the observed slope plausible if true slope=0?
- P-value: Probability that “by chance” we get as extreme value as observed
- $P=0.84$ : No evidence for deviation from null
- $P=8e-5$ : Unlikely under the null  $\rightarrow$  maybe not null



# Crude inference procedure based on “statistical significance”

P-value cutoff 

	NULL SNPs	TRUE EFFECTS
NON-SIGNIFICANT	MANY	?
SIGNIFICANT	(almost) none	?

- Label SNPs as “significant” if their p-value is small enough
- Use a stringent p-value threshold in order that there is (almost) no false positives
- Hope that there will be some true positives

# Genome-wide significance threshold

- There are  $\sim 10^6$  independent regions in the genome
  - Genome has block structure due to recombination process (linkage disequilibrium)
- If we use threshold  $p = 0.05/10^6 = \underline{5 \times 10^{-8}}$  then, on average, 1 out of 20 GWAS reports a false positive association
- Very small p-value protects from false positives

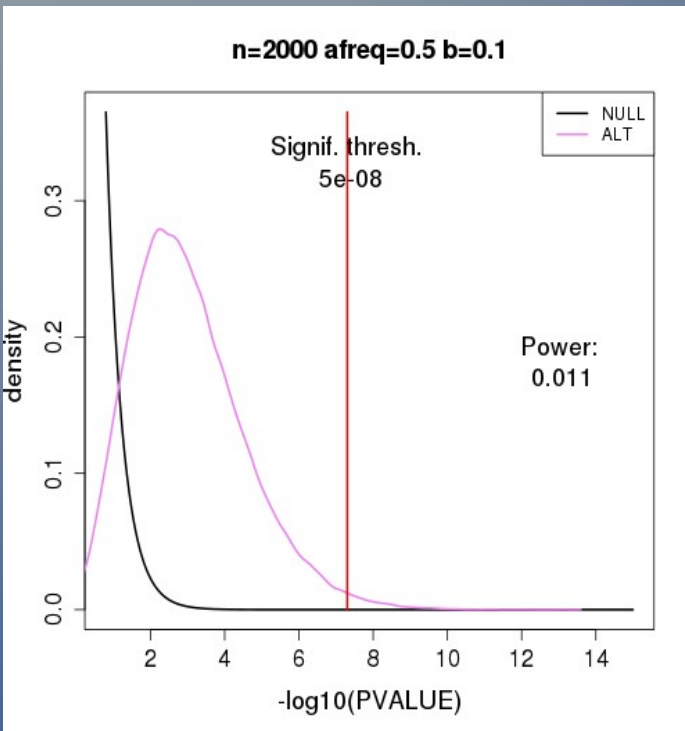
# Genome-wide significance threshold

- What if I have data only on one SNP
  - Can I now use p-value threshold 0.05?
  - What if my SNP is known to truncate a protein? Is the same p-value threshold used as for a random SNP?

# Genome-wide significance threshold

- What if I have data only on one SNP
  - Can I now use p-value threshold 0.05?
  - What if my SNP is known to truncate a protein? Is the same p-value threshold used as for a random SNP?
- Number of tests is NOT the general rule for defining consistent thresholds
  - We look at this more soon ... after statistical power

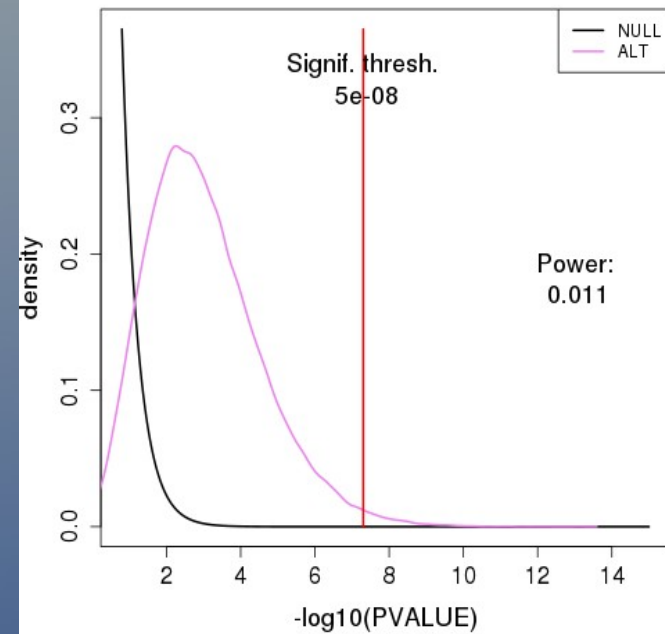
# Power



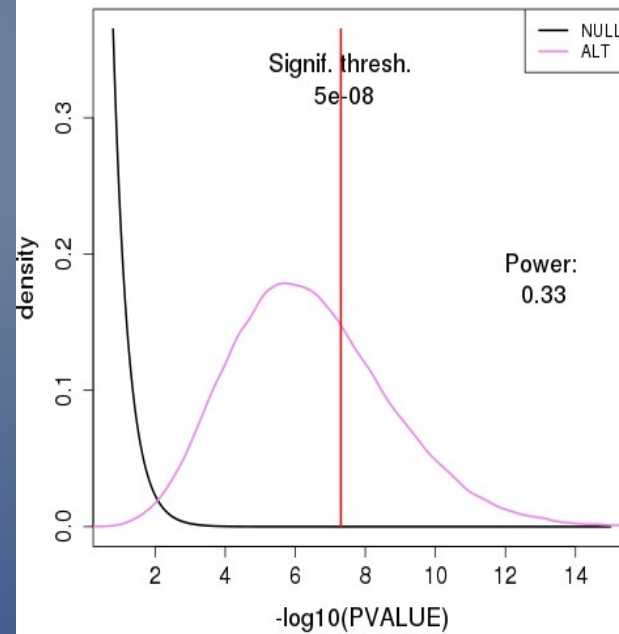
- Power = Probability that a SNP will reach a given significance threshold
  - Depends on sample size, allele freq. and effect size

# Power

n=2000 afreq=0.5 b=0.1



n=5000 afreq=0.5 b=0.1

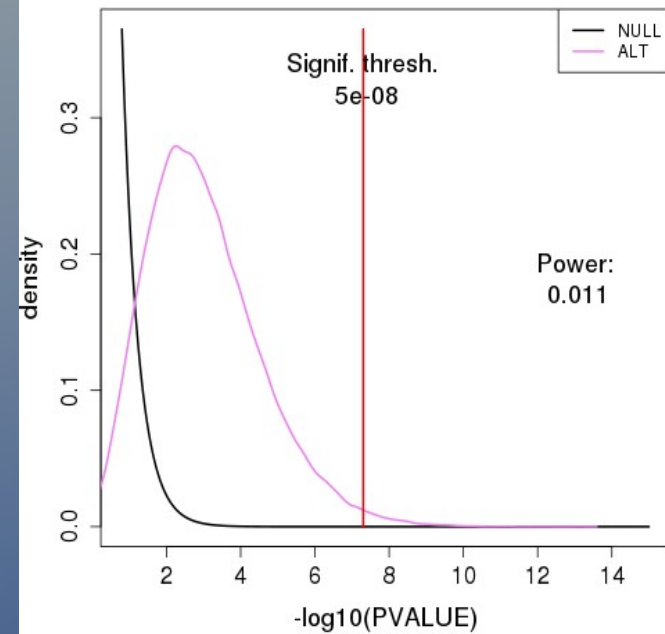


- Power = Probability that a SNP will reach a given significance threshold
  - Depends on sample size, allele freq. and effect size

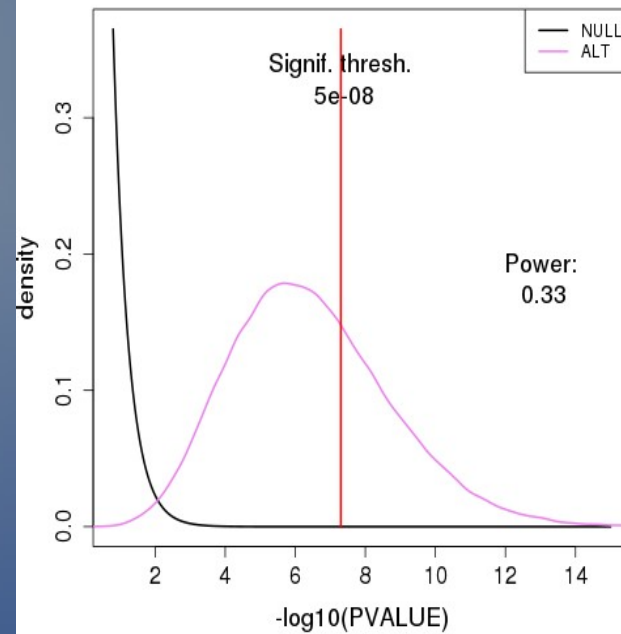


# Power

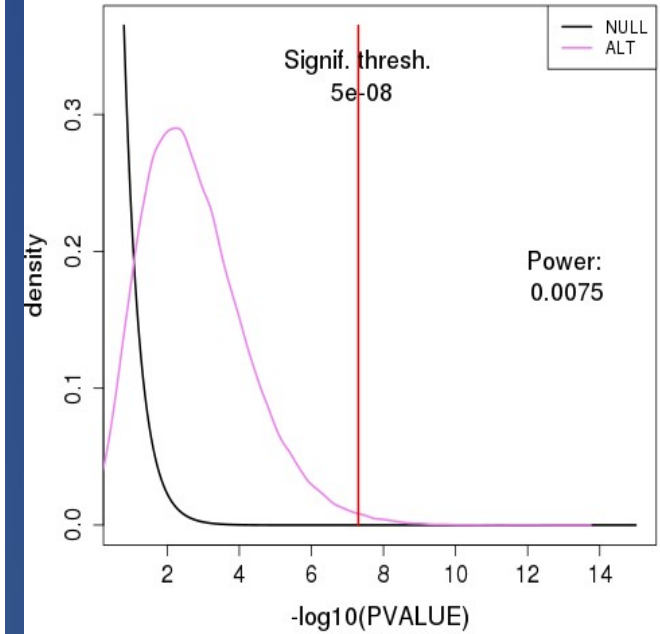
n=2000 afreq=0.5 b=0.1



n=5000 afreq=0.5 b=0.1

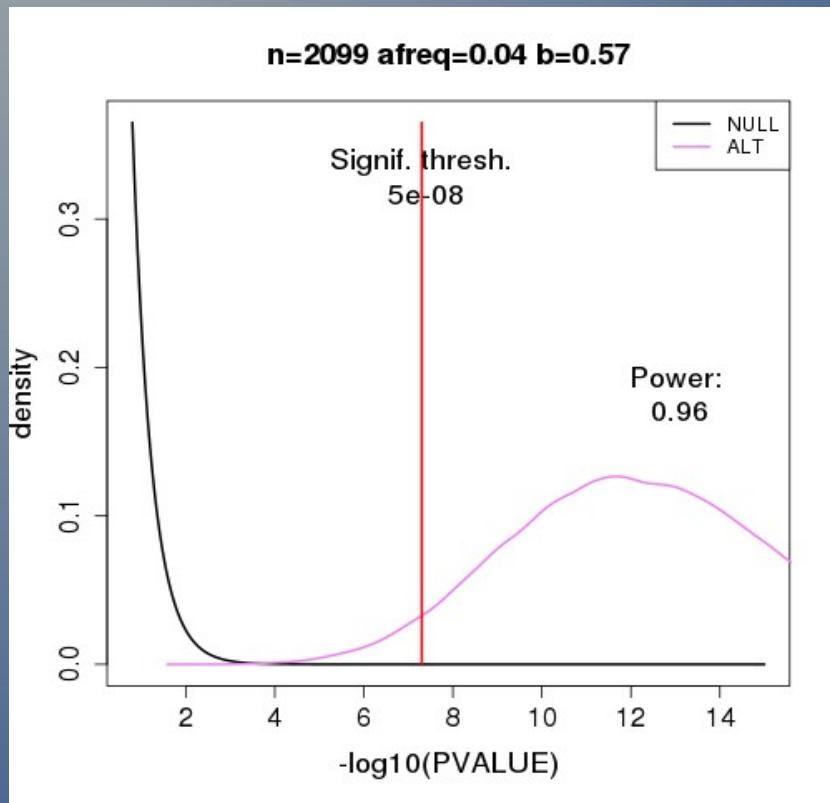


n=5000 afreq=0.1 b=0.1



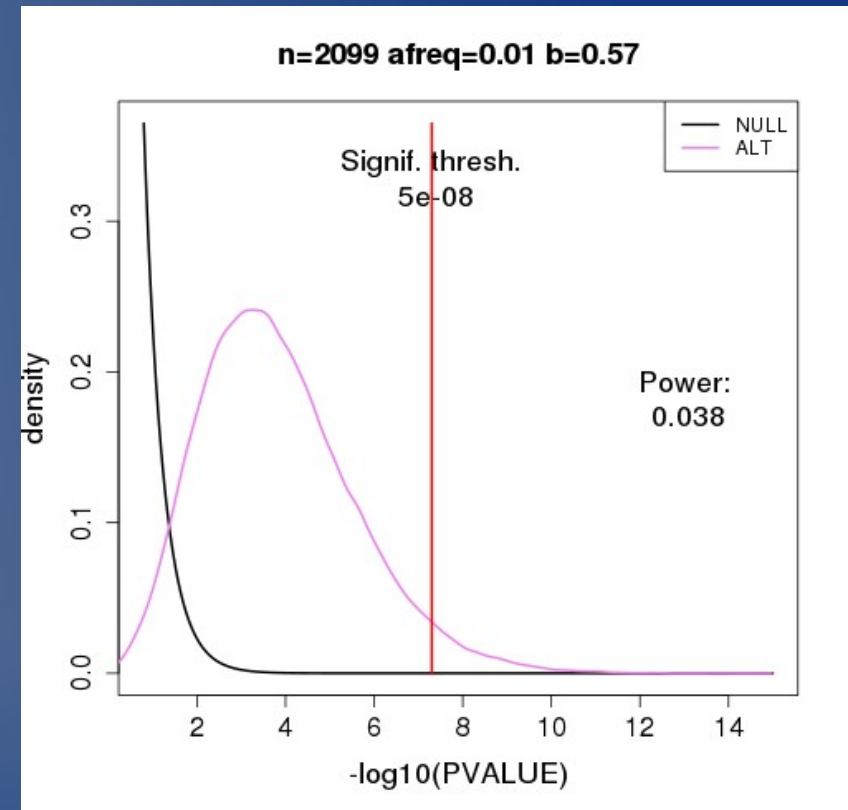
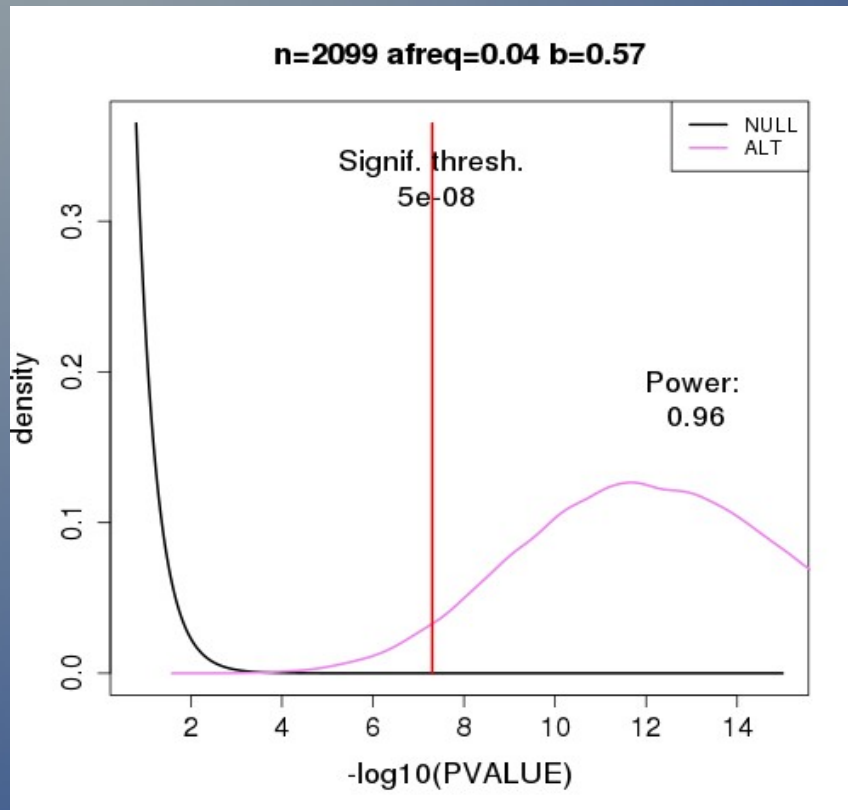
- Power = Probability that a SNP will reach a given significance threshold
  - Depends on sample size, allele freq. and effect size

# Power



- Our earlier PCSK9 variant was almost destined to be found from our Finnish data
- Power calculations needed in study design
  - Tell what kind of effects have we found / have not found

# Power



- Our earlier PCSK9 variant was almost destined to be found from our Finnish data
- But would almost certainly not been found from European data where frequency is 0.01 (compared to 0.04 in FIN)

# Ingredients of statistical power

- For quantitative traits power increases with

$$N f (1-f) b^2$$

- $N$  = sample size
- $f$  = minor allele frequency
- $b$  = effect size (“slope”) per one allele

- For case-control study power increases with

$$N t (1-t) f (1-f) b^2$$

- $t$  = proportion of cases in the study

# Properties of power

- If, for a given SNP, pop1 has MAF=4% and pop2 has MAF=1%, how large a sample from pop2 is needed for same power as a sample of  $n=2,000$  from pop1 ?

$$N f (1-f) b^2$$

# Properties of power

- If, for a given SNP, pop1 has MAF=4% and pop2 has MAF=1%, how large a sample from pop2 is needed for same power as a sample of  $n=2,000$  from pop1 ?

$$N f (1-f) b^2$$

$$N \times 0.01 \times (1-0.01) = 2000 \times 0.04 \times (1-0.04)$$
$$N=7758$$

# Properties of power

- Which one of the following two options for collecting a case-control study should you choose?
  - 2,000 inds: 1,000 cases and 1,000 controls
  - 2,500 inds: 500 cases and 2,000 controls

$$N t (1-t) f (1-f) b^2$$

# Properties of power

- Which one of the following two options for collecting a case-control study should you choose?
  - 2,000 inds: 1,000 cases and 1,000 controls
  - 2,500 inds: 500 cases and 2,000 controls

$$N t (1-t) f (1-f) b^2$$

$$2000 \times 1000/2000 \times (1-1000/2000) = 500$$

$$2500 \times 500/2500 \times (1-500/2500) = 400$$

Option 1 has more power



# From “significance” to probability of a true effect

P-value cutoff



	NULL SNPs	TRUE EFFECTS
NON- SIGNIF.	MANY	?
SIGNIF.	(almost) none	?

T = "true effect"

N = "null SNP"

S = "significant p-value"

# From “significance” to probability of a true effect

P-value cutoff



	NULL SNPs	TRUE EFFECTS
NON-SIGNIF.	MANY	?
SIGNIF.	(almost) none	?

T = “true effect”

N = “null SNP”

S = “significant p-value”

$$\frac{P(T|S)}{P(N|S)} = \frac{P(S|T)}{P(S|N)} \times \frac{P(T)}{P(N)} = \frac{\text{power}}{\text{signif. cutoff}} \times \text{prior-odds of assoc.}$$

# From “significance” to probability of a true effect

P-value cutoff



	NULL SNPs	TRUE EFFECTS
NON-SIGNIF.	MANY	?
SIGNIF.	(almost) none	?

T = “true effect”

N = “null SNP”

S = “significant p-value”

$$\frac{P(T|S)}{P(N|S)} = \frac{P(S|T)}{P(S|N)} \times \frac{P(T)}{P(N)} = \frac{\text{power}}{\text{signif. cutoff}} \times \text{prior-odds of assoc.}$$

- Small p-value threshold is because of LOW PRIOR of association NOT because of the number of tests done
  - Often prior is not known and number of tests works OK in practice
- A “signif.” finding from a well-powered study is more likely to be true than that from a study with low power (WTCCC. Nature 2007)

# Part I: Common variants

- 1. Motivation
  - What is a “genetic association” ?
  - Why is this important ?
- 2. How to read the genome (for SNPs) ?
- 3. An Example GWAS of Multiple Sclerosis
- 4. Interpreting a GWAS
- 5. Current state of GWAS

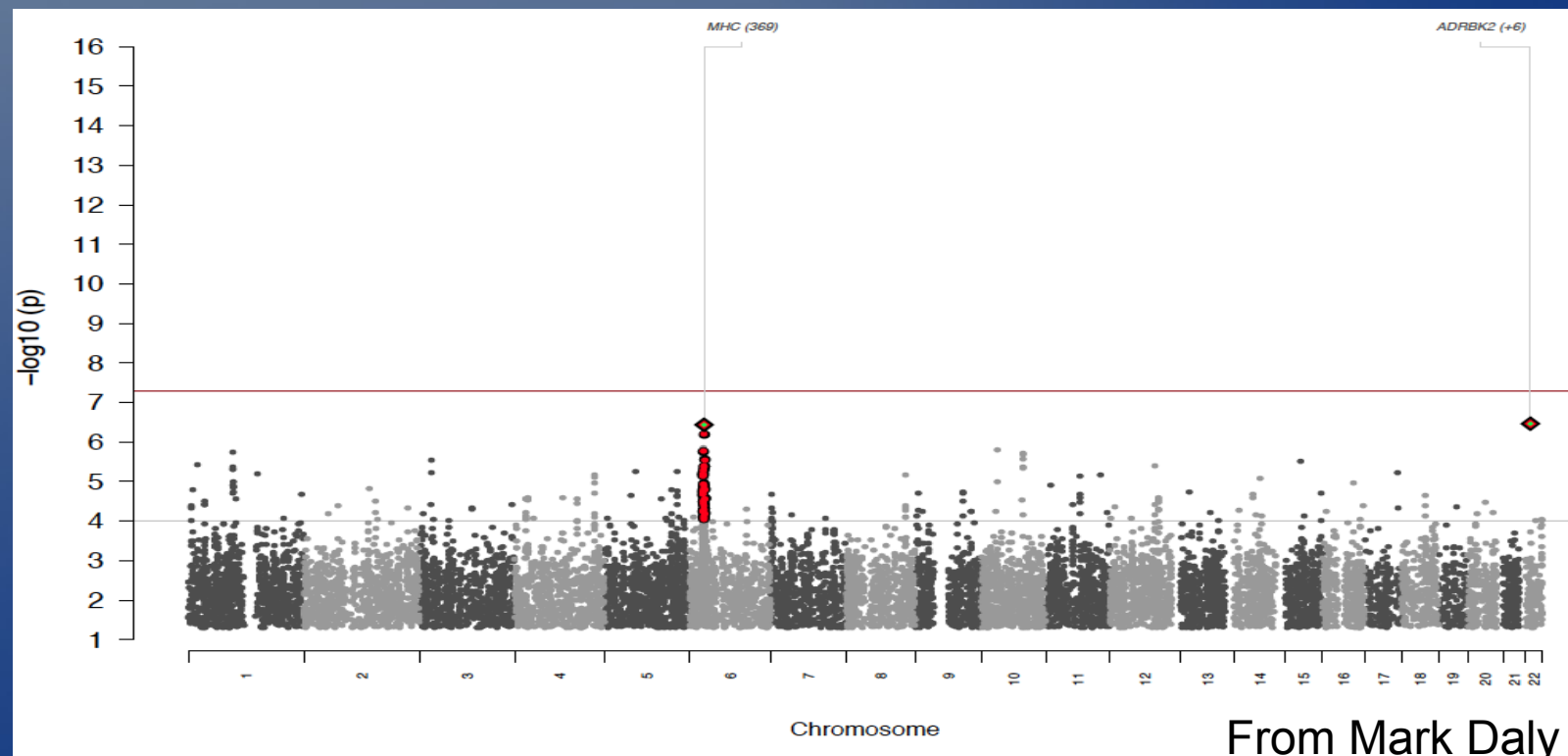
# Schizophrenia

## (as an example of GWAS evolution)

- Mental disorder with abnormal social behavior and failure to recognize what is real
- Onset as young adults, prevalence 0.5%-1%
- High heritability, estimates up to 80%
- Linkage studies with families in 1980s-1990s were not successful
  - Unlikely to exist only a few “SZ-genes” which would account for most heritability

# Int'l SZ Consortium, 2009, Nature

- 3,332 SZ cases and 3,587 controls at 1M SNPs
- Suggestive evidence for HLA-region on chr 6
- Evidence for highly polygenic basis for SZ
- But unable to identify “SZ genes”
  - GWAS is a failure ?



# PGC 2011

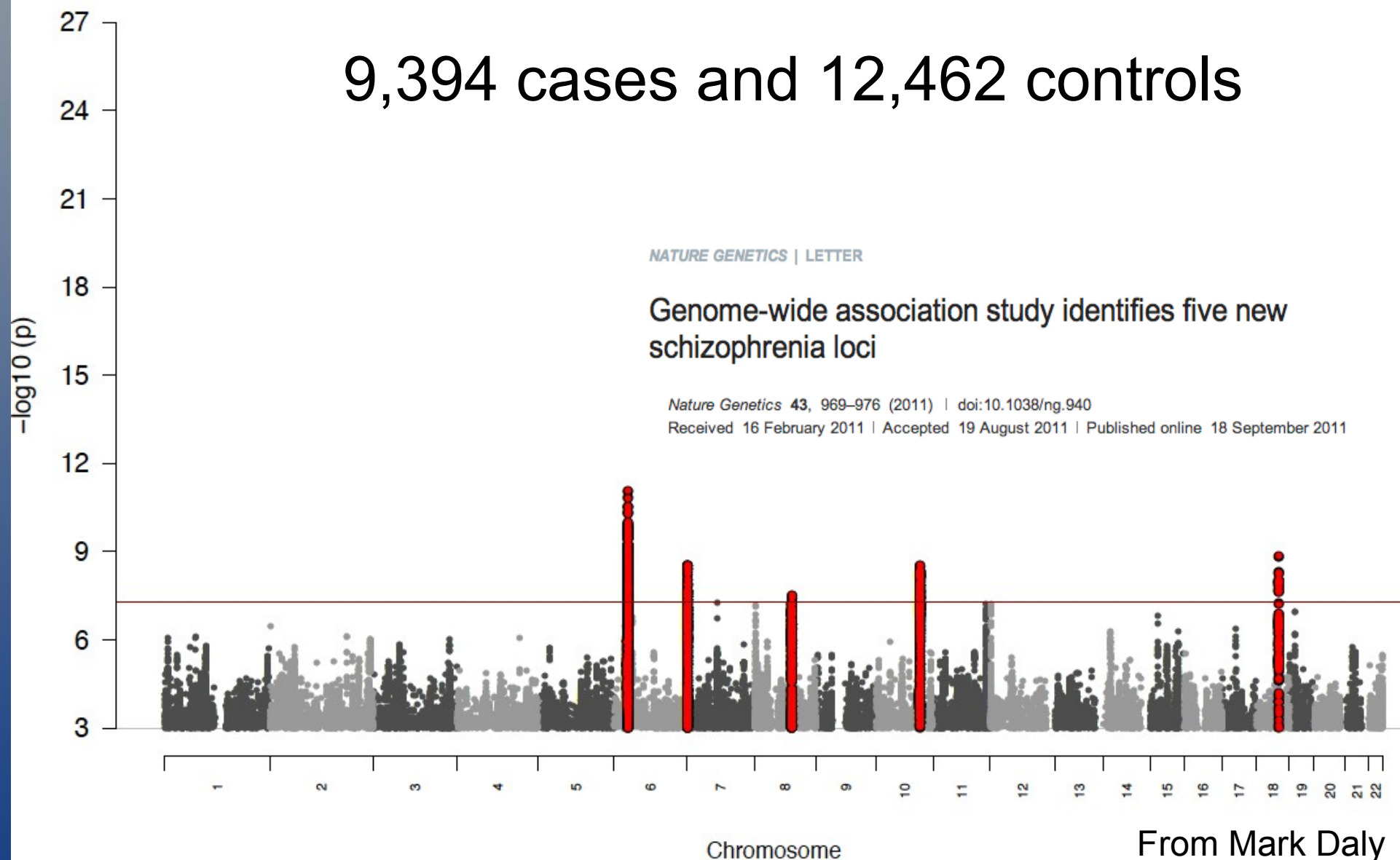
9,394 cases and 12,462 controls

NATURE GENETICS | LETTER

Genome-wide association study identifies five new schizophrenia loci

*Nature Genetics* 43, 969–976 (2011) | doi:10.1038/ng.940

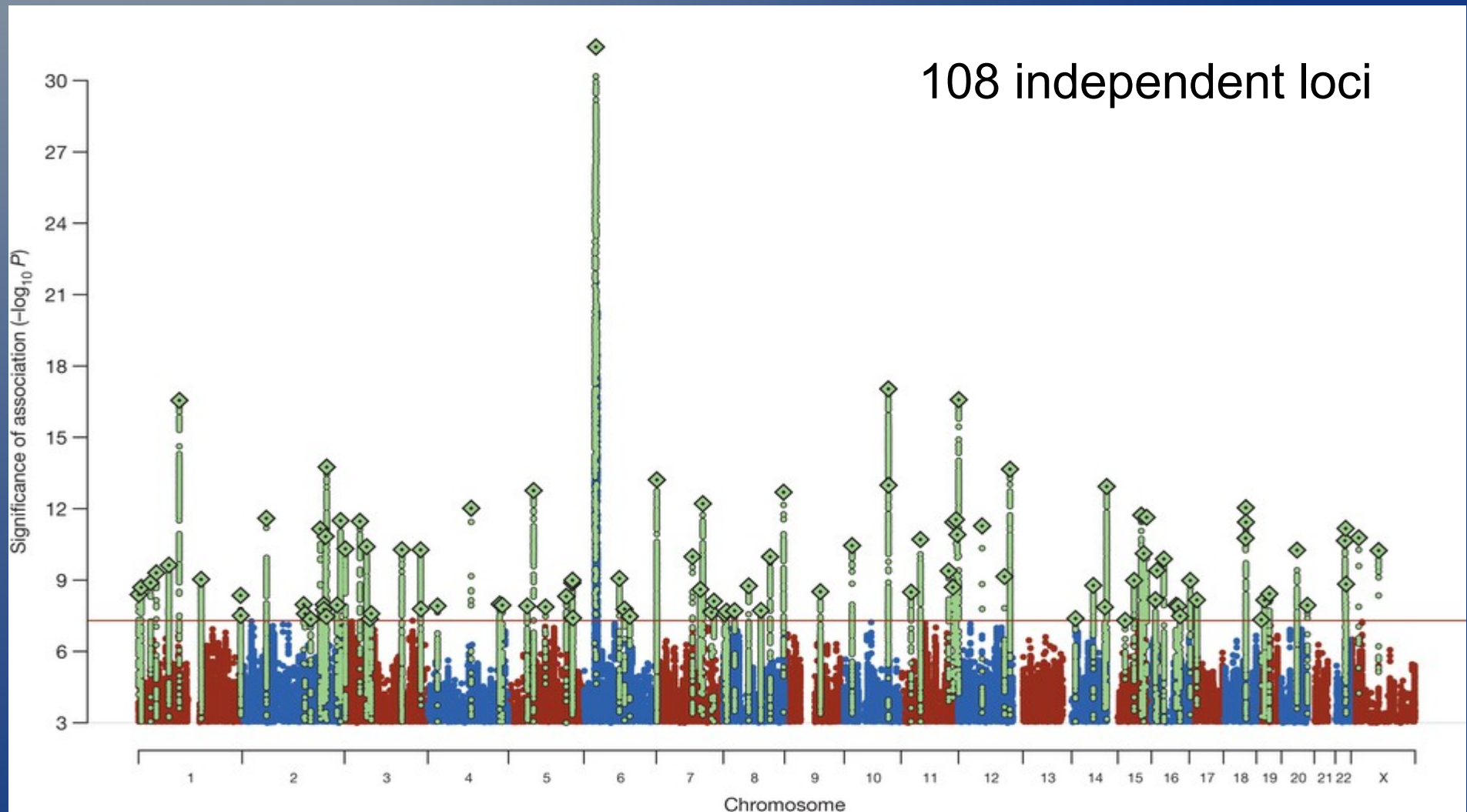
Received 16 February 2011 | Accepted 19 August 2011 | Published online 18 September 2011



From Mark Daly

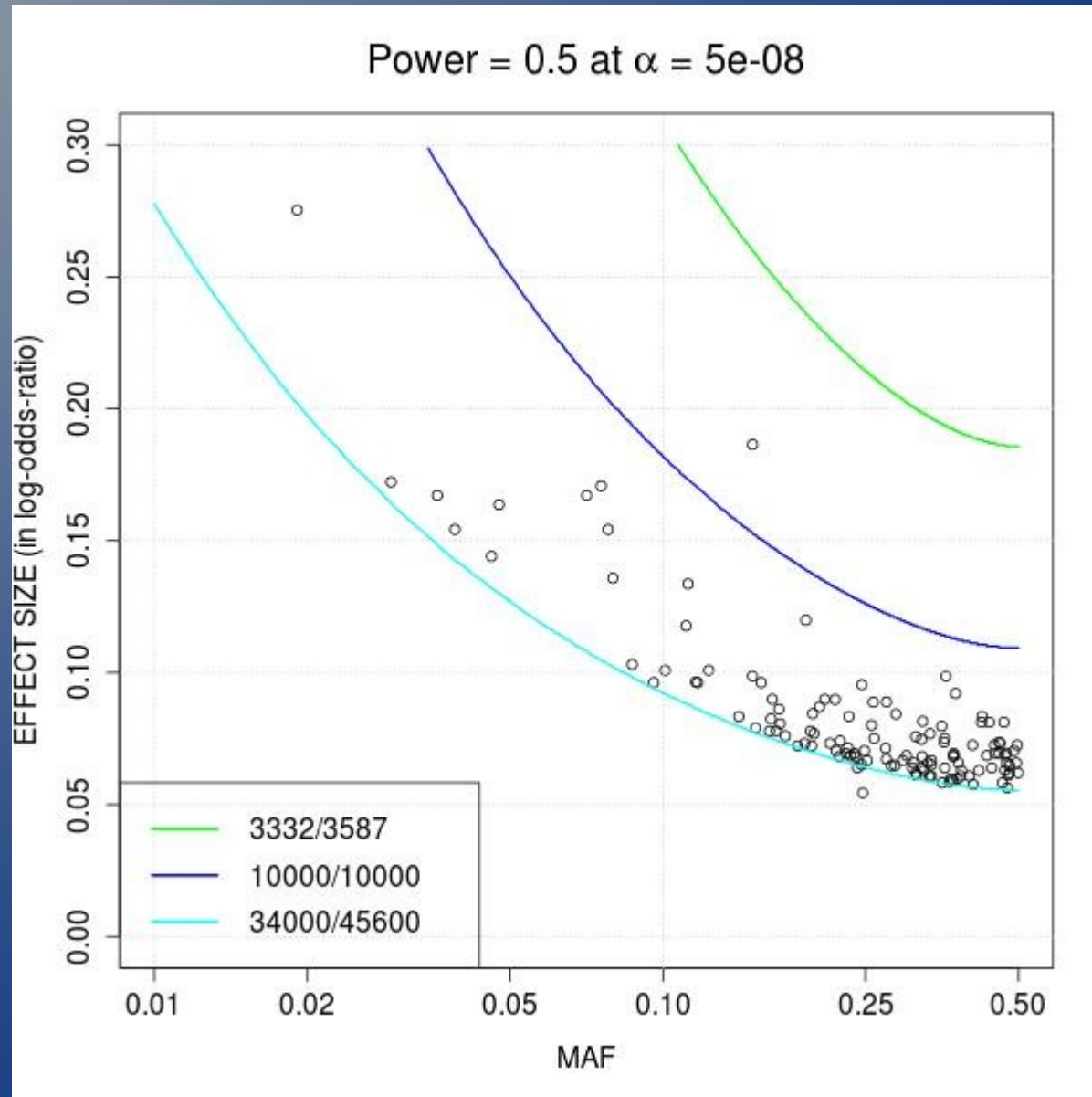
# PGC 2014, Nature

- 34,000 SZ cases and 45,600 controls at 9.5M SNPs





# Power for schizophrenia GWAS



# Biological hints emerging...

- Finding many genes implies pathways and processes underlying schizophrenia
  - 4 voltage-gated calcium channels (CACNA1C, CACNA1D, CACNA1I, CACNB2)
  - Enrichment of protein interacting with FMRP
  - Pathways highlighting postsynaptic density and dendritic spine heads
  - Enrichment of enhancer elements in specific brain regions (angular gyrus, inferior temporal lobe) and immune system

# A Pattern across common diseases and traits

	Adult height	Crohn's	Schizophrenia
Per N	10,000	1,000/1,000	3,000/3,000
1x	0	2	1
2x	2	4	2
3x	7	5	6
9x	68	51	62
18x	180	-	-

Regions with  $p < 5e-08$

Schizophrenia is a heritable, medical disorder with a genetic architecture similar to non-brain diseases and traits

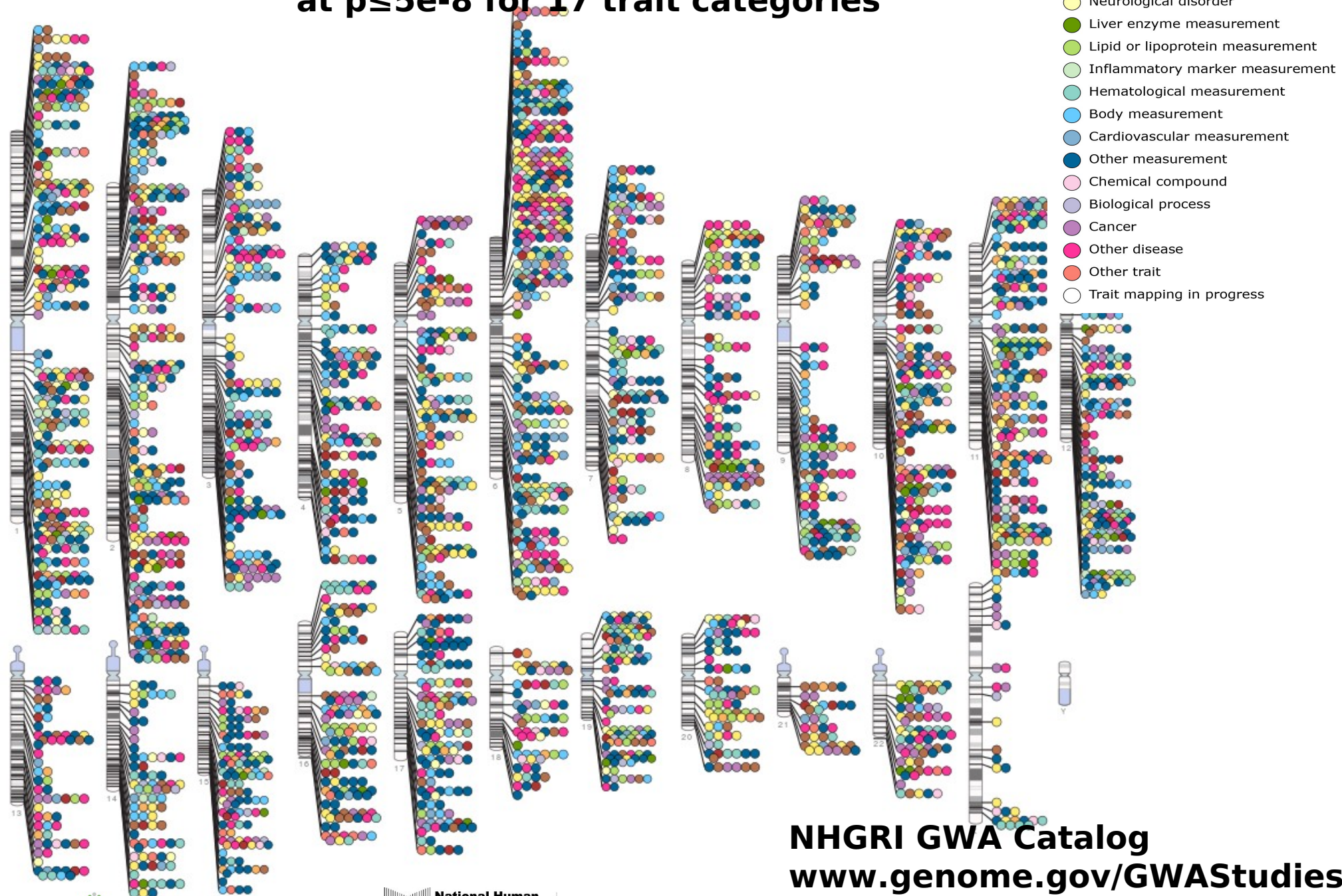
From Mark Daly

# An era of Meta-analysis

- Meta-analysis is to combine results from several individual studies on the same question e.g. a particular disease
- Increases the sample size without re-analysing everything from the very beginning
- Large consortia have been formed to carry out meta-analyses for each disease



# Published Genome-Wide Associations through 12/2013 at $p \leq 5e-8$ for 17 trait categories



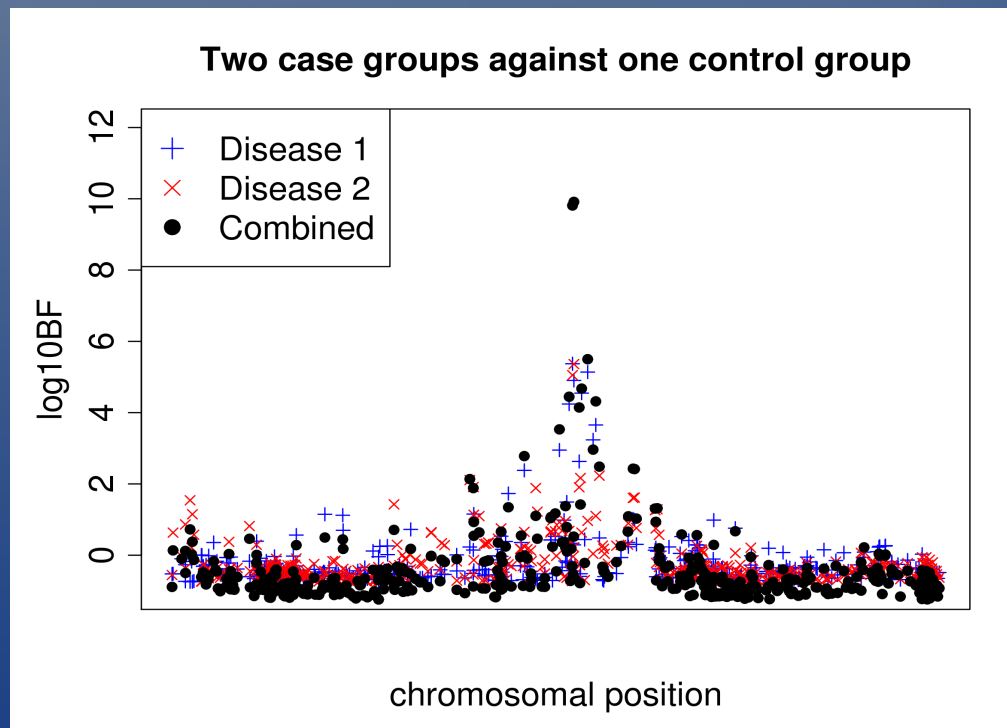
**NHGRI GWA Catalog**  
[www.genome.gov/GWASStudies](http://www.genome.gov/GWASStudies)  
[www.ebi.ac.uk/fgpt/gwas/](http://www.ebi.ac.uk/fgpt/gwas/)

# Picture emerging from GWAS

- A lot of common variants with small effects
  - Some are tagging rare variants

# Picture emerging from GWAS

- A lot of common variants with small effects
  - How many tagging rare variants?
- Many shared effects across traits
  - Need joint analyses & phenotype refinement



Psoriasis and  
Ankylosing spondylitis  
around *IL23R*

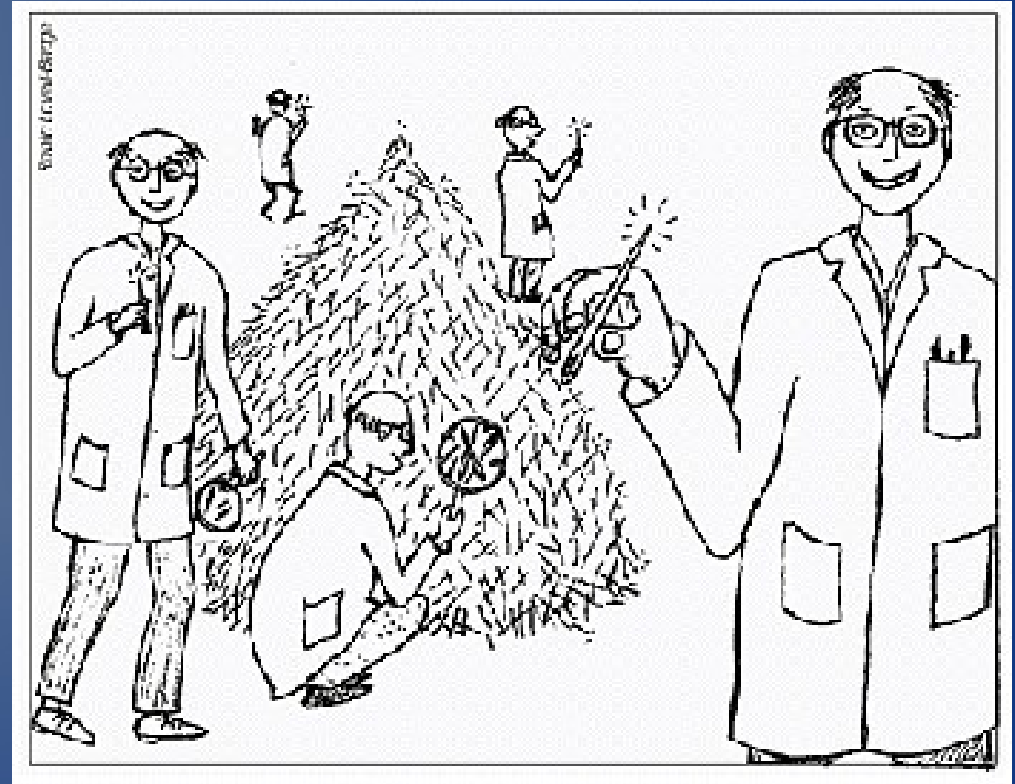
# Picture emerging from GWAS

- A lot of common variants with small effects
  - How many tagging rare variants?
- Many shared effects across traits
  - Need joint analyses & phenotype refinement
- Much to do on the biological side
  - Pathways
  - From association to function



# GWAS criticism

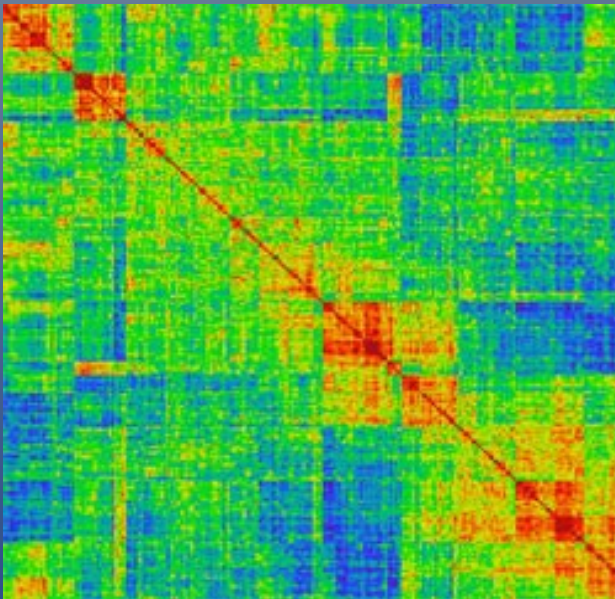
- Missing heritability
  - GWAS hits explain only a little of total genetic variance
- Missing mechanisms
- Small effect sizes
- Methodological flaws
  - Population structure
  - Separate handling of cases and controls



Weiss & Terwilliger  
Nat Gen 2000

# Criticism 1: Missing heritability

- GWAS SNPs typically explain only 10-20% of the estimated genetic variance (“heritability”)



For height in 14,500 Finnish samples, genetic component of common SNPs explains 52% of variation.

GWAS meta-analysis of 250,000 inds reports over 400 regions with top SNPs together explaining ~ 15% of the variance.

Where is the rest?

# Criticism 1: Missing heritability

- GWAS SNPs typically explain only 10-20% of the estimated genetic variance (“heritability”)
  - We don't have power to pick out (myriad of?) still smaller effects (Yang et al. 2010)

# Criticism 1: Missing heritability

- GWAS SNPs typically explain only 10-20% of the estimated genetic variance (“heritability”)
  - We don't have power to pick out (myriad of?) still smaller effects (Yang et al. 2010)
  - We haven't covered rare variants well (Gibson 2012)

# Criticism 1: Missing heritability

- GWAS SNPs typically explain only 10-20% of the estimated genetic variance (“heritability”)
  - We don't have power to pick out (myriad of?) still smaller effects (Yang et al. 2010)
  - We haven't covered rare variants well (Gibson 2012)
  - Estimates of heritability may be biased (Zuk et al. 2012)

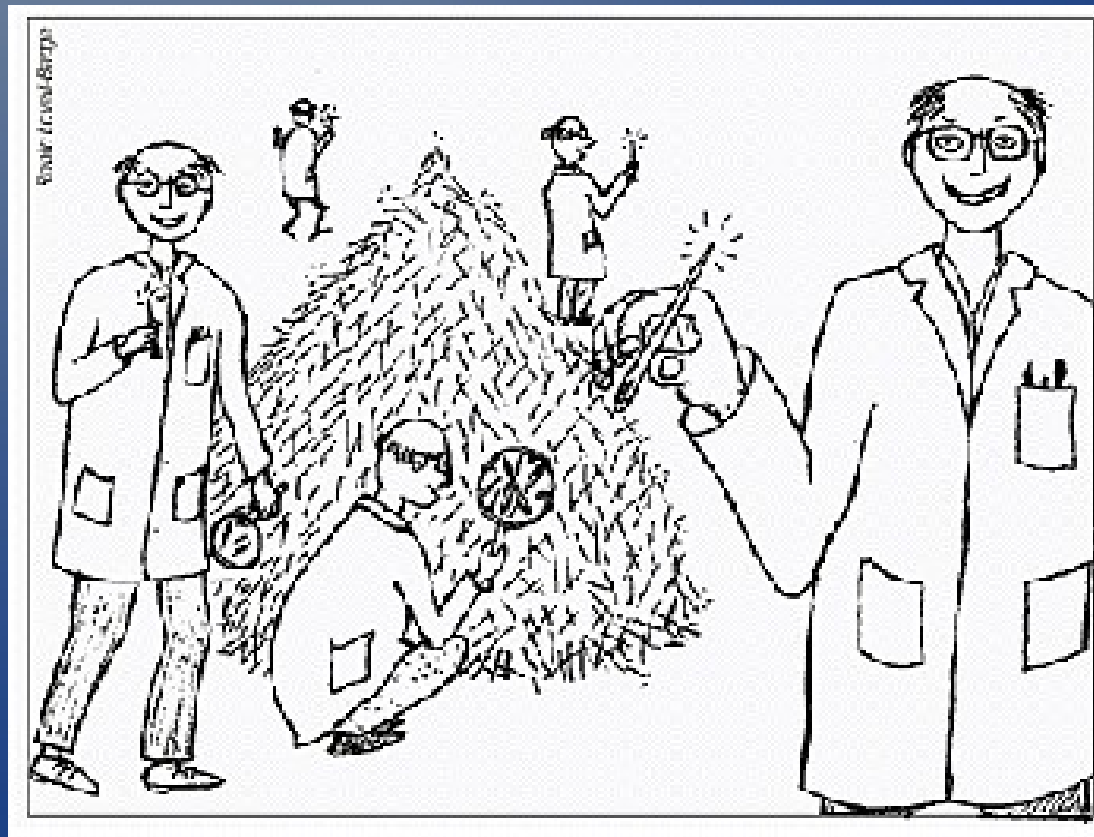
# Criticism 2: Missing mechanisms

- GWAS hits pointers but we need time to sort out
- There are examples of mechanisms
  - “From noncoding variant to phenotype via SORT1 at the 1p13 cholesterol locus” Nature 2010

Recent genome-wide association studies (GWASs) have identified a locus on chromosome 1p13 strongly associated with both plasma low-density lipoprotein cholesterol (LDL-C) and myocardial infarction (MI) in humans. Here we show through a series of studies in human cohorts and human-derived hepatocytes that a common noncoding polymorphism at the 1p13 locus, rs12740374, creates a C/EBP (CCAAT/enhancer binding protein) transcription factor binding site and alters the hepatic expression of the SORT1 gene. With small interfering RNA (siRNA) knockdown and viral overexpression in mouse liver, we demonstrate that Sort1 alters plasma LDL-C and very low-density lipoprotein (VLDL) particle levels by modulating hepatic VLDL secretion. Thus, we provide functional evidence for a novel regulatory pathway for lipoprotein metabolism and suggest that modulation of this pathway may alter risk for MI in humans. **We also demonstrate that common noncoding DNA variants identified by GWASs can directly contribute to clinical phenotypes.**

# Criticism 3: Small effect sizes

- If an effect of an allele is 1 mm in height or relative risk of 1.1 for multiple sclerosis is that of any use?
- Can we make everything affect everything when sample size is large enough?



# Small effect sizes

- Maybe individual common variants have small effects because nature does NOT tolerate large effects at those loci
- By therapies we can perturb those pathways with much larger effects than present in nature
- “The HMGCR locus has a common variant at 40% frequency that changes LDL by a modest 2.8 mg/dl and no known rare mutations of large effect, presumably because they would be lethal. Yet, the encoded protein is the target of statins, drugs taken by tens of millions of patients that can significantly reduce both LDL levels and myocardial infarction risk.” E Lander 2011



# Criticism 4: Flawed design of GWAS

- Population structure
- Different genotyping errors and other properties in cases and controls if genotyped separately

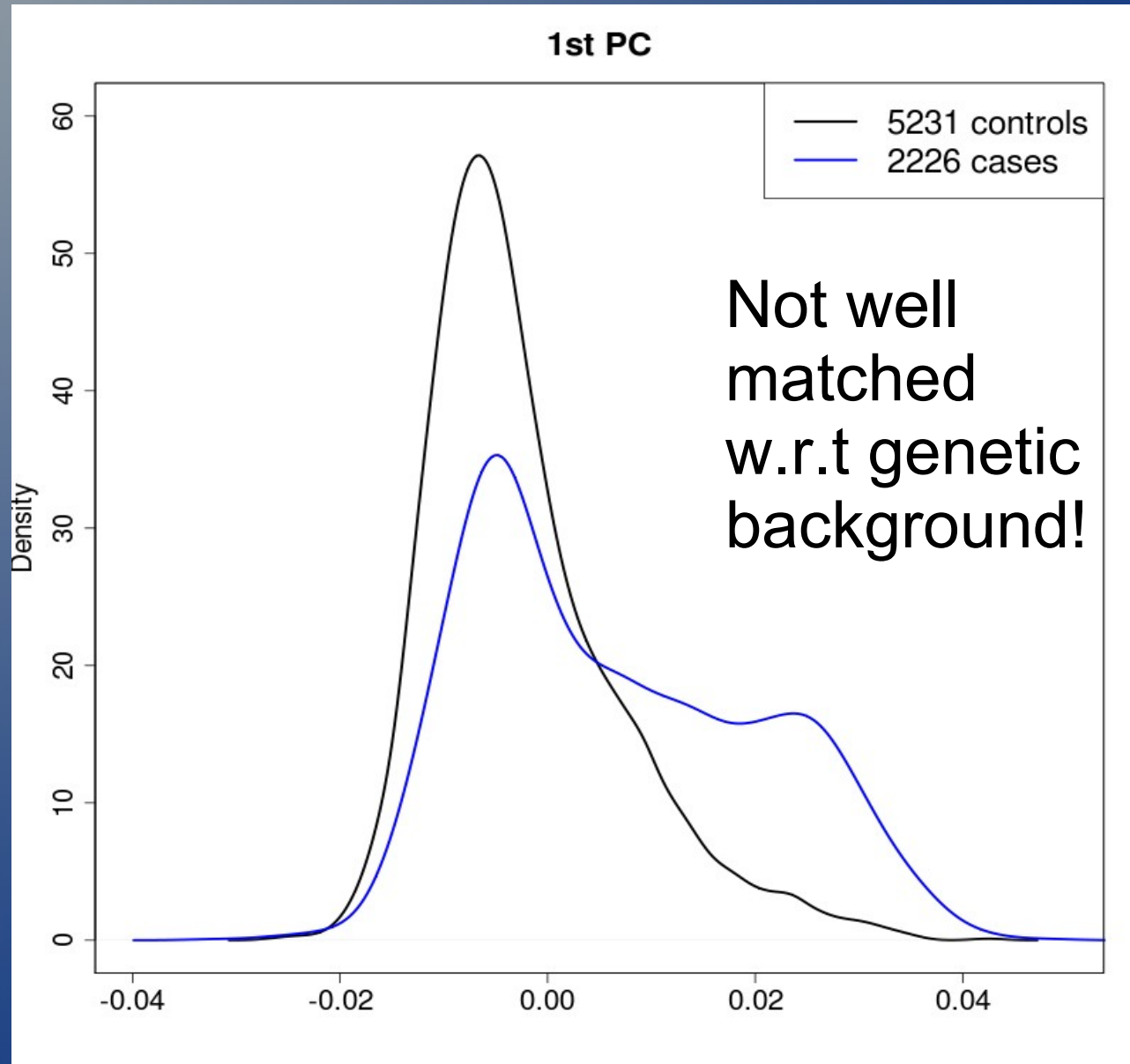
# To tackle confounding in GWAS

- Good methods to warn about and correct for systematic biases: mixed models, covariates, careful quality control, qq-plots
- In large meta-analyses unlikely that biases of individual studies would multiply
- Replication in several cohorts provides convincing evidence
- Stringent statistical thresholds are applied (p-value  $5e-8$ )

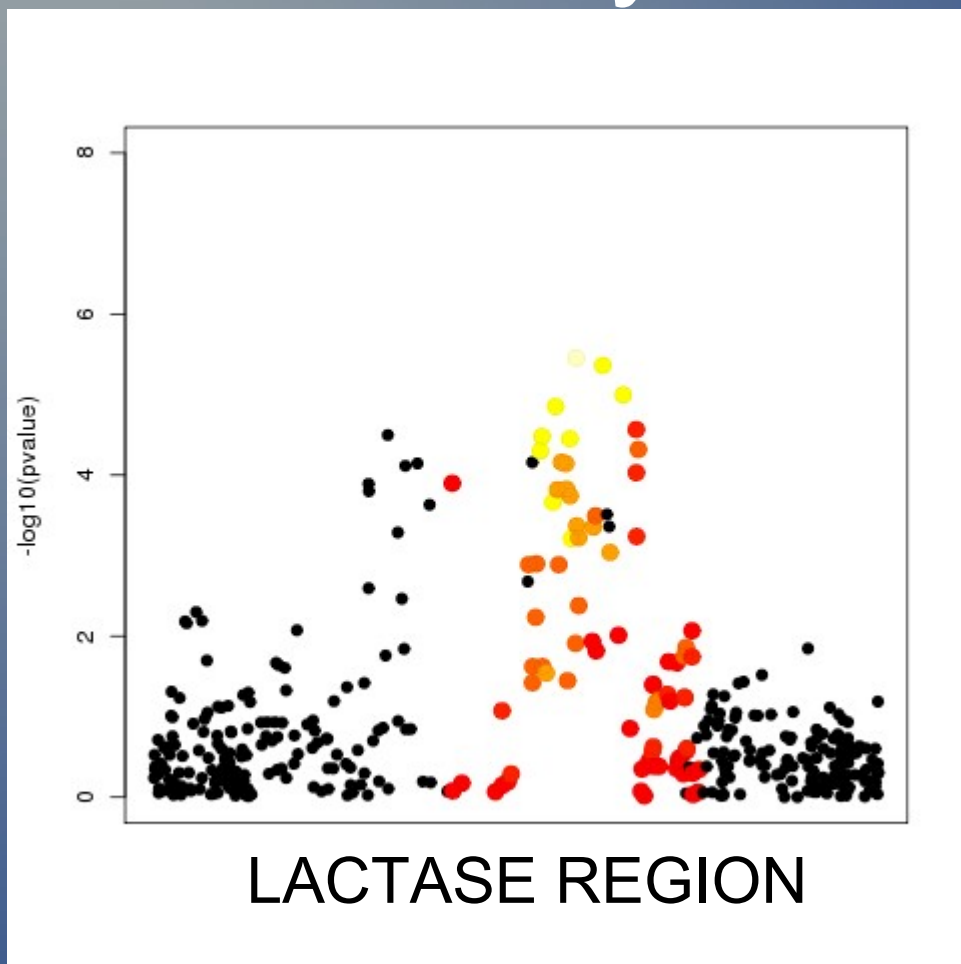
# Population structure in case-control association studies

- Are there differences in genotype frequencies between cases and controls?
  - If yes, then locus is possibly interesting
  - But could also reflect ascertainment scheme if cases and controls are not well matched w.r.t. genetic background!
  - Example: Let's look at analysis where 5,000 UK controls are compared against 1,800 UK + 500 Irish Psoriasis cases

# UK+Irish Psoriasis study



# Basic analysis



# First PC as covariate

