

A minicourse on  
Genomewide association studies  
(GWAS)  
Part II: Analysis

Matti Pirinen  
FIMM, University of Helsinki

December 5th 2012



Contents:  
1. Concepts  
And Rationale

2. Technology  
and Data

**3. Statistics  
and Analysis**

4. Results and  
Hot topics

- Digestive system disorder
- Cardiovascular disorder
- Metabolic disorder
- Immune system disorder
- Neurological disorder
- Liver enzyme measurement
- Lipid or lipoprotein measurement
- Inflammatory marker measurement
- Hematological measurement
- Body measurement
- Cardiovascular measurement
- Other measurement
- Chemical compound
- Biological process
- Cancer
- Other disease
- Other trait
- Trait mapping in progress

# GWAS Statistics

- Linear model, ML-estimates, Wald's test, p-value, power (R EXAMPLE), odds (WTCCC 2007 BOX on significance)
- Bayesian model, marginal likelihood, Bayes factor (BF) (Wakefield 2009, Vukcevic 2009 Ch. 5&6) (R EXAMPLE)

# Effect size

- Quantitative: phenotypic scale (e.g., cm, std dev)
- Disease: odds-ratio (~relative risk for rare diseases)
  - Logistic regression
- Meta-analysis: Combining results over studies
  - Inverse variance weighting
  - Forest plots
- Effect size is not (the only) measure of importance of a GWAS finding! (Lander 2011, BOX 1: HMGCR and statins)

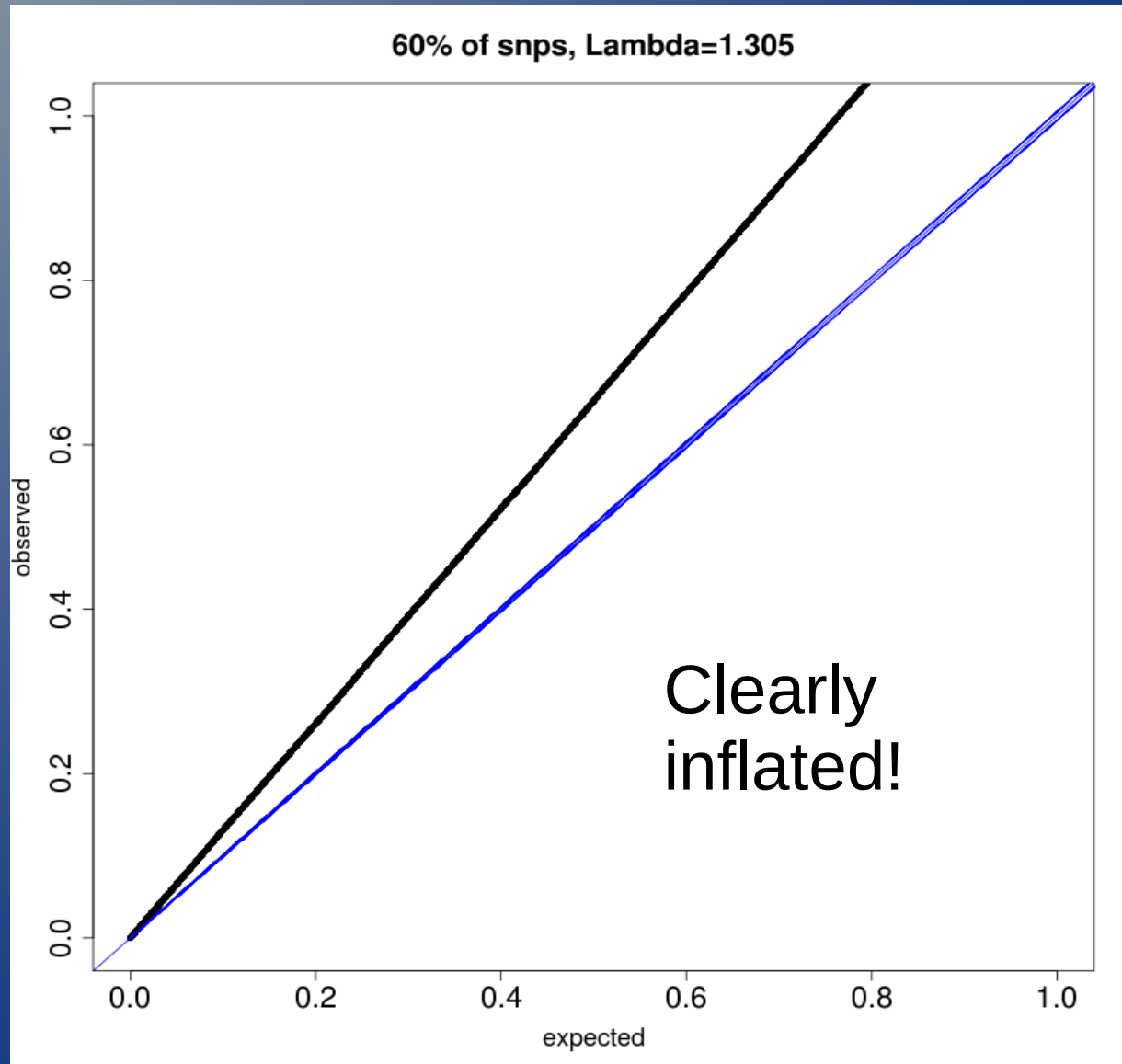
# QQ-plot

- R EXAMPLE: 3 QQ-plots
  - NULL, inflated and exciting
- Inflation (measure lambda) arises because of
  - Polygenic effects AND / OR
  - Ascertainment bias (in case-control studies)

# 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

# Plot association statistics over the genome and compare to the expectation under the null



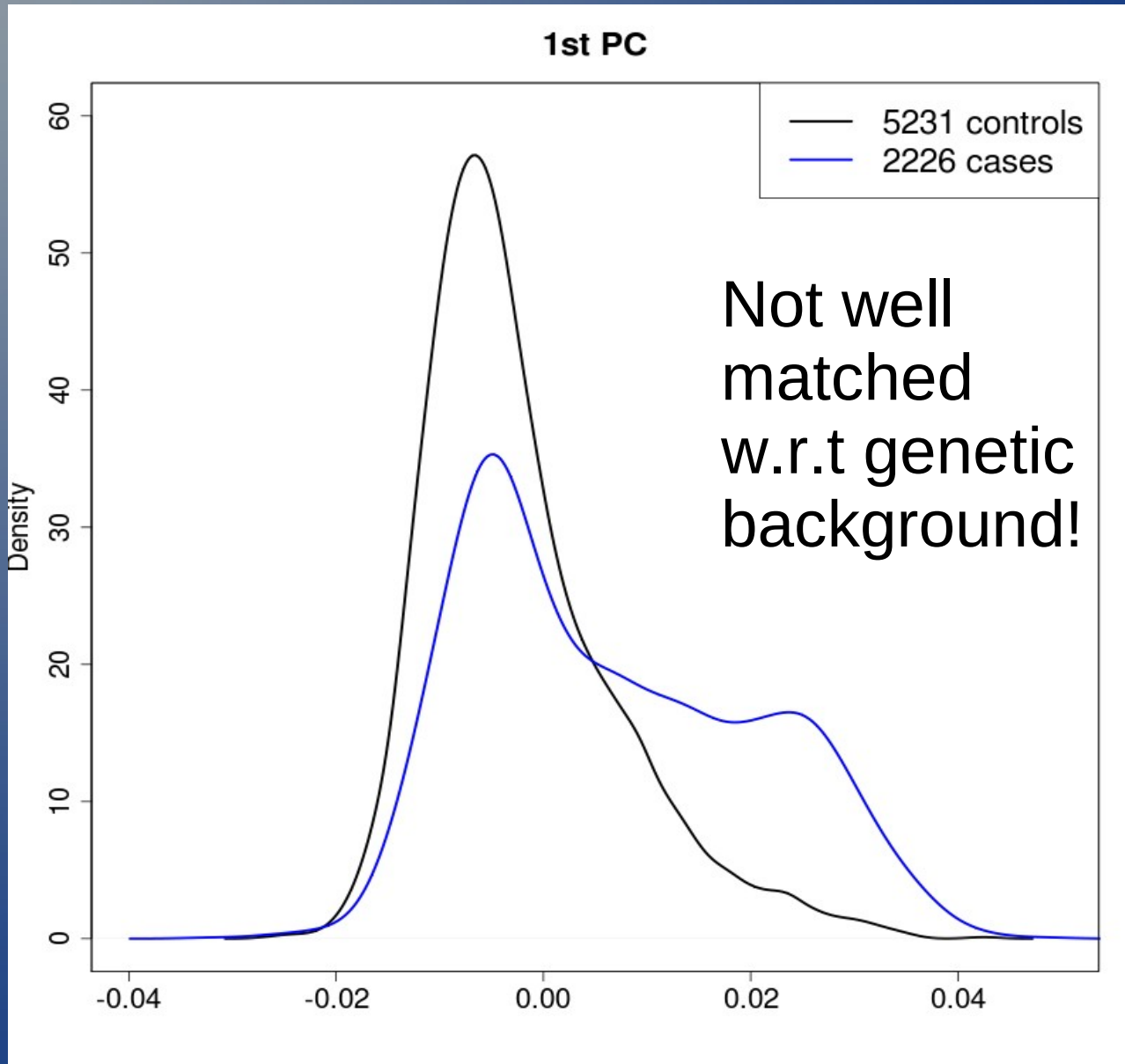


# Principal components of population structure

- Find linear combination of SNPs that has the largest sample variance (1st PC)
  - Sequentially, find further PCs ORTHOGONAL to the previous ones
- Novembre 2008: European PC plot
- R EXAMPLE

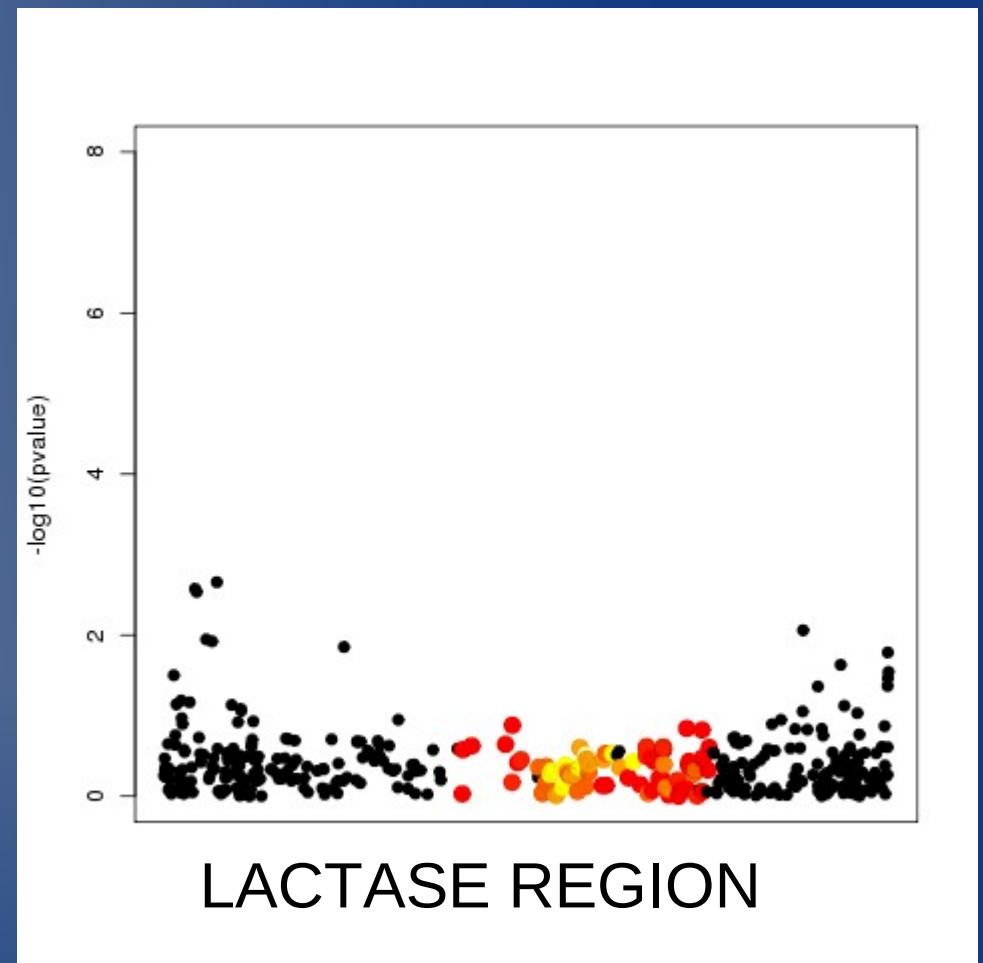
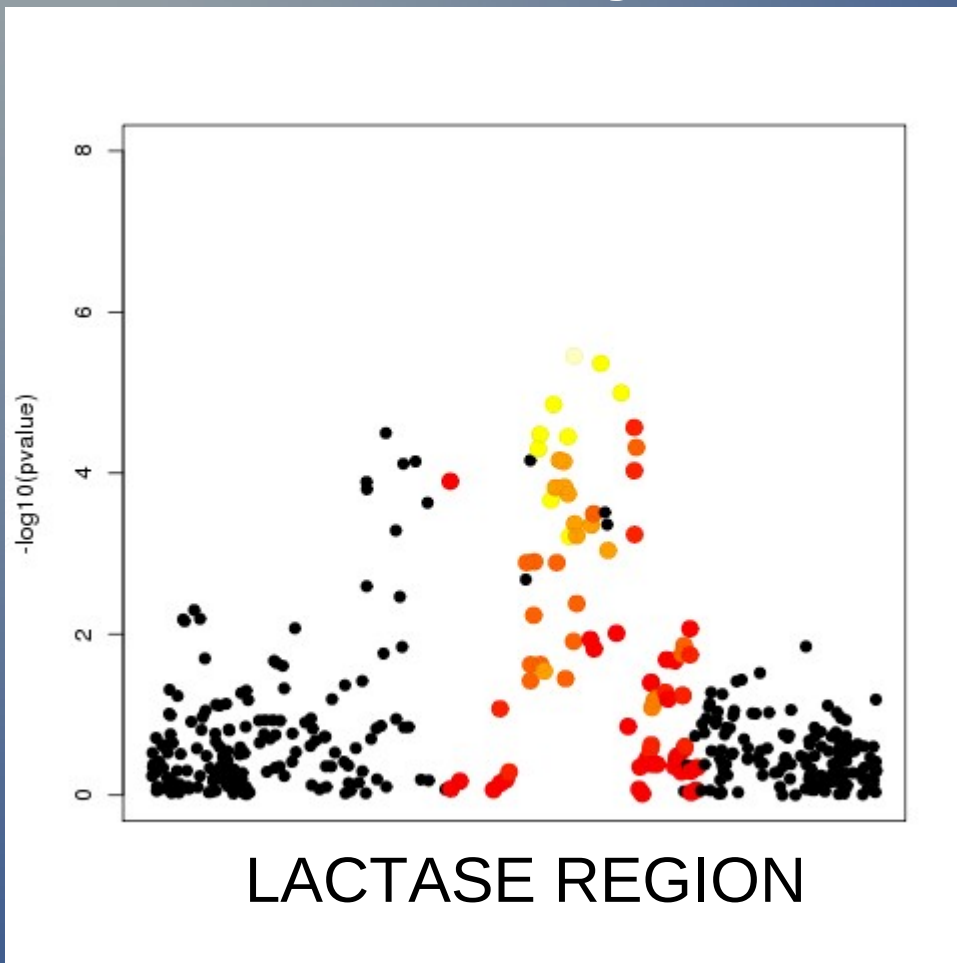


# PCA on UK+Irish Psoriasis study

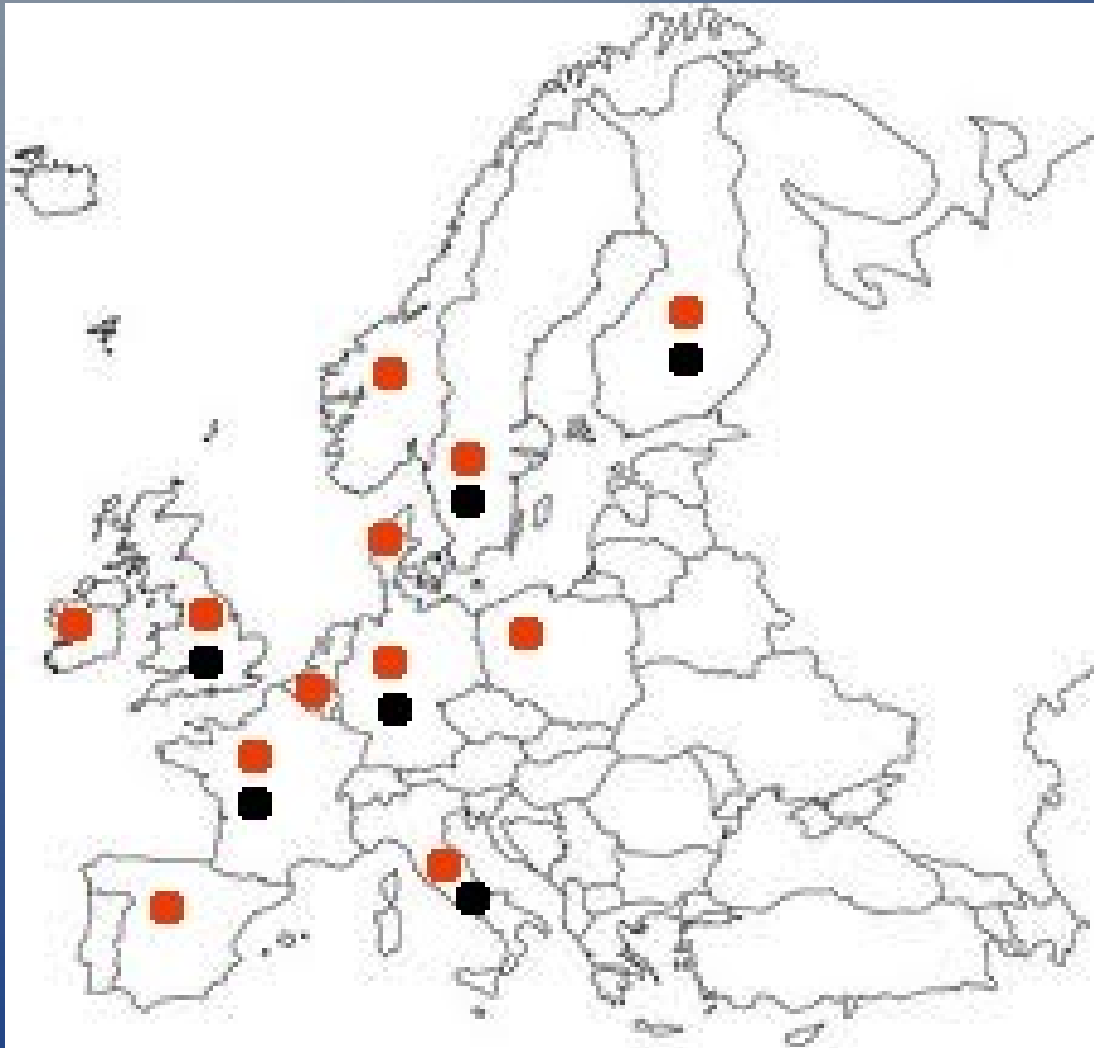


# Basic analysis

# First PC as covariate



Multiple sclerosis:  
10,000 cases/  
17,000 controls

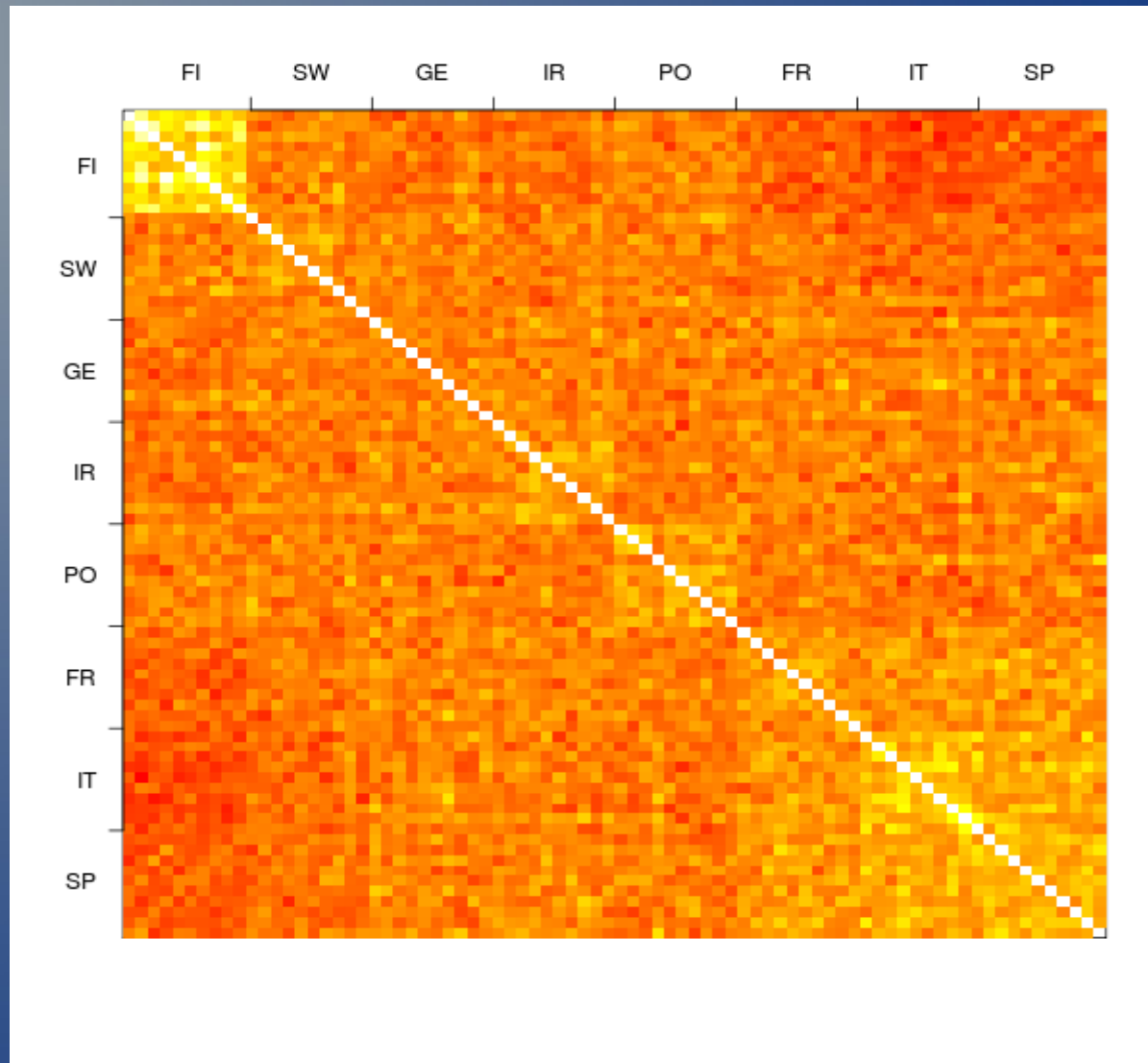


| Pop       | Cases | Controls |
|-----------|-------|----------|
| Spain     | 205   | --       |
| Norway    | 953   | 121      |
| UK        | 1854  | 5175     |
| Sweden    | 685   | 1928     |
| France    | 479   | 347      |
| Poland    | 58    | --       |
| Denmark   | 332   | --       |
| Belgium   | 544   | --       |
| Italy     | 745   | 571      |
| Germany   | 1100  | 1699     |
| Ireland   | 61    | --       |
| Finland   | 581   | 2165     |
| Australia | 648   | --       |
| NZ        | 146   | --       |
| USA       | 1383  | 5370     |

# Possible strategies

- Each country analysed separately and results combined (loses samples and power due to differing case-control ratios)
- Clustering individuals into genetically homogeneous groups which are analysed separately and results combined
- Regression models on whole data with population structure included in the model
  - A few leading principal components of the genetic structure
  - Linear mixed model

# Genome-wide correlation matrix (R)



LMM:  $Y \sim \mu + G \beta + z$ , where  
 $E(z) = 0$  and  $\text{var}(z) \propto R$

PC:  $Y \sim \mu + G \beta + P \gamma$ , where  
P is the leading eigenvector of R

# Linear Mixed model

- In the standard linear model it is computationally possible to introduce individual specific effects that follow empirical genetic correlations (R matrix)
- $Y = \alpha + G\beta + Z + \varepsilon$ , where
  - $Z \sim N(0, h\sigma^2 R)$   $\varepsilon \sim N(0, (1-h)\sigma^2 I)$
- Is this appropriate for case-control data?

# Linear model with case-control data

- Likelihood ratio statistic approaches the Armitage trend test
- Least-squares applied to binary data approximates the ML-estimates of logistic regression when effects are small
  - Relative error in log-odds estimates is less than 0.5% in typical GWAS application



# QQ-plots

$\lambda = 1.04$

Linear mixed model

$\lambda = 1.22$

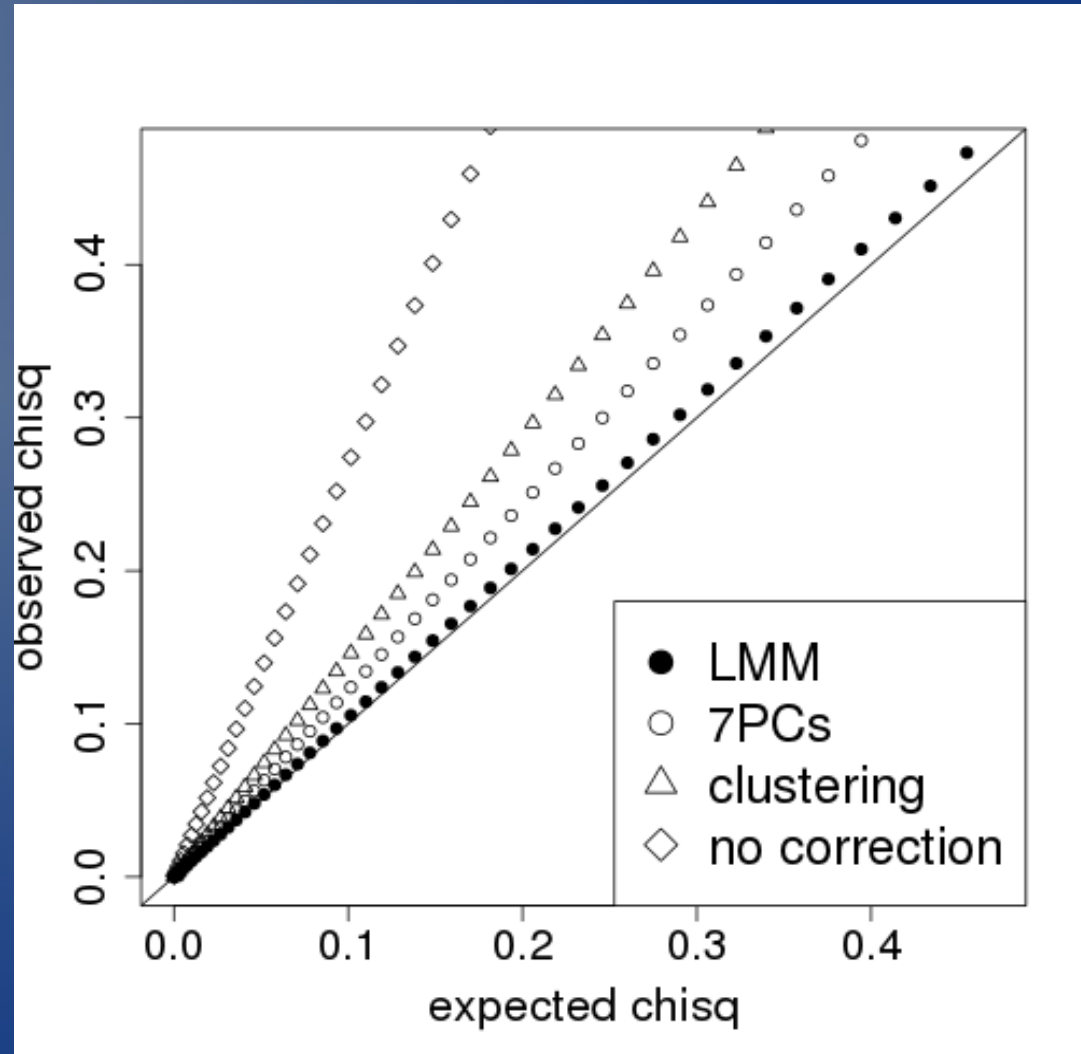
Logistic regression with 7  
PCs as covariates

$\lambda = 1.44$

Clustering individuals to  
homogeneous groups

$\lambda = 2.7$

No structure corrections

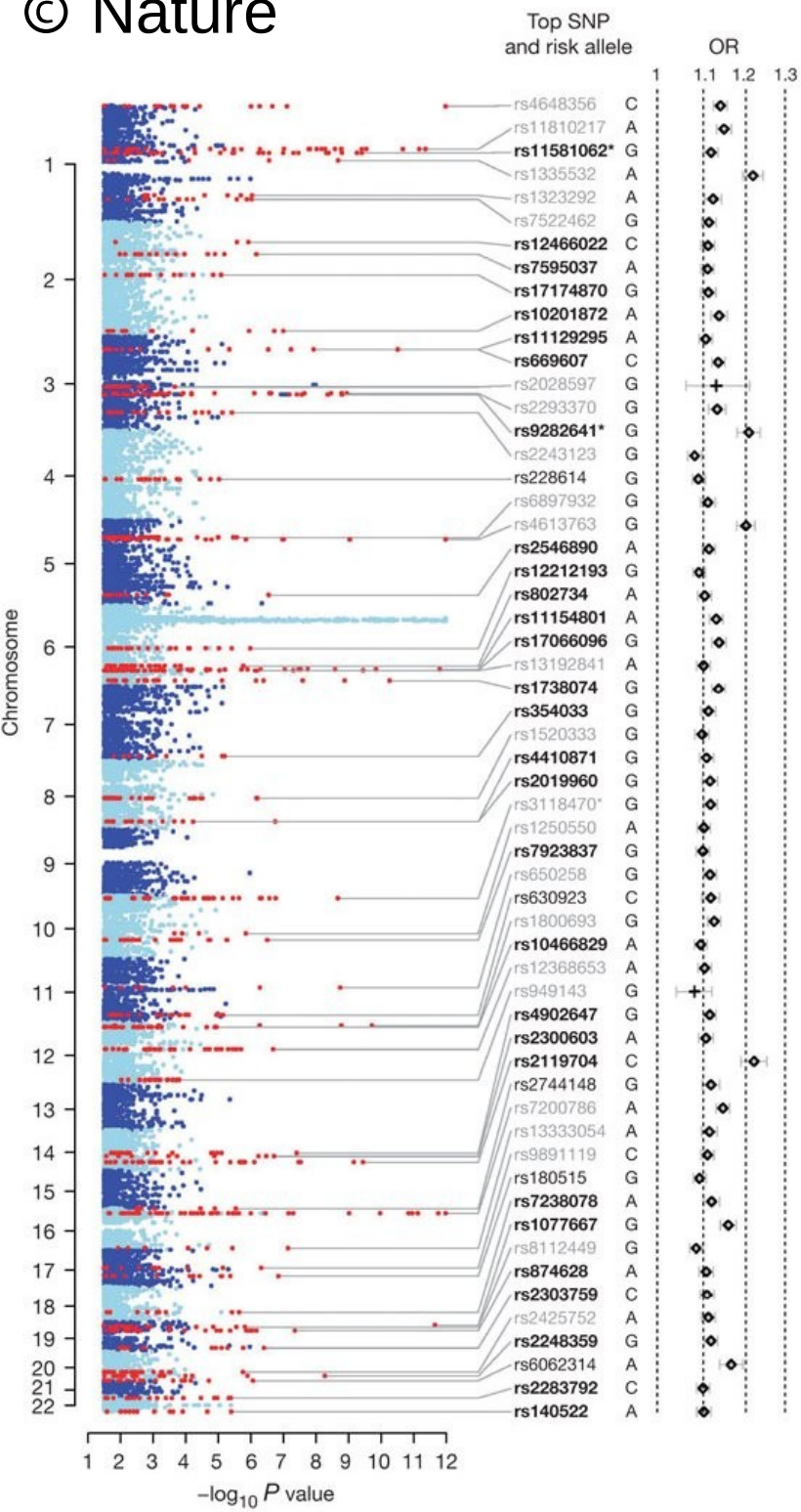


# We also applied

- A beta-binomial model to exclude variants that differ among control data sets more than expected under no selection
- Empirical allele variation scores
- 1000 Genomes database to see whether our hits tag strongly differentiated variants
- Bayesian hierarchical model to assess whether the effect sizes differ between the populations

# Replication

- We took 102 SNPs to replication based on LMM analysis
- Replication attempted in silico from meta-analysis of 6 previous MS-studies having together 4,200 cases and 7,300 controls
- Logistic regression with 10 PCs
- Carried out at the Broad Institute



- 98 /102 SNPs have consistent effects in replication data
- 29 novel:  $p < 5e-8$  (combined),  $p < 1e-4.5$  (discovery),  $p < 0.05$  (replication, 1-sided)
- plus 5 novel:  $p < 5e-7$  (combined),  $p < 1e-4.5$  (discovery),  $p < 0.05$  (replication, 1-sided)
- 23 / 26 previously reported MS loci had  $p < 1e-3$  in our discovery data
- Immunological genes are overrepresented among the hits; in particular, T-helper cell differentiation pathway

# EXAMPLE GWAS: Multiple sclerosis

- WTCCC2 2011
- <http://wattle.well.ox.ac.uk/wtccc2/external/ms/>
- Over 50 associations now known
  - Primary role for immune mechanism
  - T-helper cell differentiation
  - 2 GWAS hits close to known drug targets, what about remaining ~50 loci!?

# Further study on one SNP

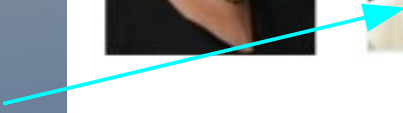
- Reported before: Tumour necrosis factor (TNF) blocking drugs are bad for MS patients
- Gregory et al. (July 2012 NATURE) “TNF receptor 1 genetic risk mirrors outcome of anti-TNF therapy in multiple sclerosis”
  - MS risk variant at rs1800693 (one of MS GWAS hits) promote TNF-blocking effect
  - Consistent effect between GWAS findings and empirical experience with a drug



# WTCCC2 Oxford



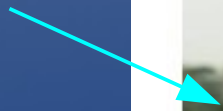
Chris  
Spencer



Garrett  
Hellenthal



Peter  
Donnelly





# Thank you

- 245 authors with 140 affiliations around the world
  - MS disease groups (including Finland)
  - Genotyping at Sanger institute
  - Replication analyses at the Broad institute
- Stephen Sawcer, and Alastair Compston, Clinical Neuroscience, Cambridge, UK

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>2</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 Little<sup>7</sup>, Hanne F. Harbo<sup>2,122</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>3</sup>, Cristin Aubin<sup>4</sup>, Amie Baker<sup>1</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. Broadley<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>, Katleen Clysters<sup>9</sup>, Giancarlo Comi<sup>20</sup>, Mark Cossburn<sup>29</sup>, Isabelle Cournu-Rebeix<sup>12</sup>, Matthew 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 Dobos<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>, Refujia 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>6</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>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>25</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 Marriot<sup>89,40</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>2,183</sup>, Tania Mihalova<sup>27</sup>, Xavier Montalban<sup>25</sup>, John Mottershead<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 Pickersgil<sup>29</sup>, Fredrik Piehli<sup>26</sup>, Susan Pobywajlo<sup>5</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>99</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 Sellebjerg<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 Stankovic<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 Bramon<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>76,119,120</sup>, Robert Plomin<sup>121</sup>, Ernest Willoughby<sup>122</sup>, Anna Rautanen<sup>2</sup>, Juliane Winkelmann<sup>13,123,124</sup>, Michael Wittig<sup>58,125</sup>, Richard C. Trembath<sup>116</sup>, Jacqueline Yao uanq<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. Ivinson<sup>51,130</sup>, Philip L. De Jager<sup>4,5,51</sup>, Leena Peltonen†, Graeme J. Steffer<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>

# References

- **IMSGC & WTCCC2.** (Aug 2011) “Genetic risk and a primary role for cell-mediated immune mechanisms in multiple sclerosis.” **Nature.**
- **Pirinen, Donnelly, Spencer.** (2012) “Efficient computation with a linear mixed model on large-scale data sets with applications to genetic studies.” (available online) **Annals of Applied Statistics.**
  - **Binary data**
  - **Algorithmic improvement from  $O(n^3)$  to  $O(n^2)$**
  - **Bayesian computation**

# Summaries of GWAS

- YOUTUBE: “Genome-Wide Association Studies - Karen Mohlke (2012)” ~90 minutes
- Paper: Pearson, Manolio (2008) “How to Interpret a Genome-wide Association Study?”  
JAMA 299(11):1335-1344