

# A genetic test for differential causative pathology in disease subgroups

James Liley<sup>1,2</sup>, John A Todd<sup>1</sup>, and Chris Wallace<sup>1,2,3</sup>

<sup>1</sup>*JDRF/Wellcome Trust Diabetes and Inflammation Laboratory, Department of Medical Genetics, NIHR Cambridge Biomedical Research Centre, Cambridge Institute for Medical Research, University of Cambridge, Cambridge, UK*

<sup>2</sup>*Department of Medicine, University of Cambridge, Addenbrookes Hospital, Cambridge, CB2 0SP, UK*

<sup>3</sup>*MRC Biostatistics Unit, Institute of Public Health, University Forvie Site, Robinson Way, CB2 0SR, Cambridge, UK*

## Abstract

Many common diseases show wide phenotypic variation. We present a statistical method for determining whether phenotypically defined subgroups of disease cases represent different genetic pathophysiologies, in which disease-associated variants have different effect sizes in the two subgroups. Our method models the genome-wide distributions of genetic association statistics with mixture Gaussians. We apply a global test without requiring explicit identification of disease-associated variants, thus maximising power in comparison to a standard variant by variant subgroup analysis. Where evidence for genetic subgrouping is found, we present methods for post-hoc identification of the contributing genetic variants.

We demonstrate the method on a range of simulated and test datasets where expected results are already known. We investigate subgroups of type 1 diabetes (T1D) cases defined by autoantibody positivity, establishing evidence for differential genetic basis with thyroid peroxidase antibody positivity, driven generally by variants in known T1D associated regions.

## Introduction

Analysis of genetic data in human disease typically uses a binary disease model of cases and controls. However, many common human diseases show extensive clinical and phenotypic diversity which may represent multiple causative pathophysiological processes. Because therapeutic approaches often target disease-causative pathways, understanding this phenotypic complexity is valuable for further development of treatment, and the progression towards personalised medicine. Indeed, identification of patient subgroups characterised by different clinical features can aid directed therapy [1] and accounting for phenotypic substructures can improve ability to detect causative variants by refining phenotypes into subgroups in which causative variants have larger effect sizes [2].

Such subgroups may arise from environmental effects, reflect population variation in non-disease related anatomy or physiology, correspond to partitions of the population in which disease heritability differs, or represent different forms of the underlying disease. In the latter two cases, we expect that the genetic architecture of the disease will differ between subgroups. We present a statistical method for assessing whether this is the case.

Our test is for a stronger assertion than the question of whether subgroups of a disease group exhibit any genetic differences at all, as these may be entirely disease-independent: for example, although there will be systematic genetic differences between Asian and European patient cohorts with type 1 diabetes (T1D) these differences will not in general relate

to the pathogenesis of disease.

We proceed by joint consideration of allelic differences between the combined case group and controls, and allelic differences between case subgroups independent of controls. We test for the presence of a category of SNPs with allelic differences between subgroups which additionally show evidence for association with the disease as a whole. This assumption has been used previously for genetic discovery [3].

Rather than attempting to analyse SNPs individually, a task for which GWAS are typically underpowered, we model allelic differences across all SNPs using multivariate normal models, and analyse the overall distribution. This can give insight into the structure of the genetic basis for disease. Given evidence that there exists some subset of SNPs that both differentiate controls and cases and differentiate subgroups, we can then reassess test statistics to search for single-SNP effects with more confidence, in an approach similar to the empirical Bayes method.

We establish properties of the method by a range of simulations, and demonstrate the method by establishing differentiation between T1D, type 2 diabetes (T2D), and RA datasets, considered as subgroups of a case group of generic disease phenotypes. We further demonstrate the method by analysing differences between Graves' disease (GD) and Hashimoto's thyroiditis (HT). Finally, we apply the method to a T1D dataset partitioned by positivity of disease-associated autoantibodies. Throughout, we use the term 'subgrouping' to mean a division of a case group into subgroups.

## Results

### Summary of proposed method

Our overall protocol is to fit two bivariate Gaussian mixture models, corresponding to null and alternative hypotheses, to summary statistics derived from SNP data. We assume a group of controls and two non-intersecting case subgroups.

For a SNP  $i$  we denote by  $\mu_1(i)$ ,  $\mu_2(i)$ ,  $\mu_{12}(i)$  and  $\mu_c(i)$  the population minor allele frequencies for the two case subgroups, the whole case group and the control group respectively, and  $P_d(i)$ ,  $P_a(i)$  GWAS p-values for comparisons of allelic frequency between case subgroups and between cases and controls, under the null hypotheses  $\mu_1(i) = \mu_2(i)$  and  $\mu_{12}(i) = \mu_c(i)$  respectively. We then derive absolute  $Z$  scores  $|Z_d(i)|$  and  $|Z_a(i)|$  from these p-values (see figure 1). We consider the values  $|Z_d(i)|, |Z_a(i)|$  as observations of random variables  $(Z_a, Z_d)$  which are samples from a mixture of three bivariate Gaussians. Further details are given in the supplementary material, section 1.

We consider each SNP to fall into one of three categories, each corresponding to a different joint distribution of  $Z_d, Z_a$ :

1. SNPs which do not differentiate subgroups and are not associated with the phenotype as a whole ( $\mu_c = \mu_1 = \mu_2$ )
2. SNPs which are associated with the phenotype as a whole but which are not differentially associated with the subgroups ( $\mu_c \neq \mu_{12}; \mu_1 = \mu_2 = \mu_{12}$ )
3. SNPs which have different population allele frequencies in subgroups, and may or may not be associated with the phenotype as a whole ( $\mu_1 \neq \mu_2$ )

If the SNPs in category 3 are not associated with the disease as a whole (null hypothesis,  $H_0$ ), we expect  $Z_d Z_a$  to be independent and the variance of  $Z_a$  to be 1. If SNPs in category

3 are also associated with the disease as a whole (alternative hypothesis,  $H_1$ ), the joint distribution of  $(Z_d Z_a)'$  will have both marginal variances greater than 1, and  $Z_a, Z_d$  may co-vary. Our test is therefore focussed on the form of the joint distribution of  $(Z_d, Z_a)$  in category 3. Importantly, we allow that the correlation between  $Z_d$  and  $Z_a$  may be simultaneously positive at some SNPs and negative at others. This allows for a subset of SNPs to specifically alter risk of one subgroup, and another subset alter risk for the other subgroup. To accommodate this, we only consider absolute  $Z$  scores and model the distribution of SNPs in category 3 with two mirror-image bivariate Gaussians.

Amongst SNPs with the same frequency in disease subgroups (categories 1 and 2),  $Z_a$  and  $Z_d$  are independent and the expected standard deviation of  $Z_d$  is 1. We therefore model the overall joint distribution of  $(Z_d, Z_a)$  as a Gaussian mixture in which the *pdf* of each observation  $(Z_d, Z_a)$  is given by

$$\begin{aligned}
 PDF_{Z_d, Z_a | \Theta}(d, a) &= \pi_0 N_{\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}}(d, a) && \text{(category 1)} \\
 &+ \pi_1 N_{\begin{pmatrix} 1 & 0 \\ 0 & \sigma_1^2 \end{pmatrix}}(d, a) && \text{(category 2)} \\
 &+ \pi_2 \left( \frac{1}{2} N_{\begin{pmatrix} \tau^2 & \rho \\ \rho & \sigma_2^2 \end{pmatrix}}(d, a) + \frac{1}{2} N_{\begin{pmatrix} \tau^2 & -\rho \\ -\rho & \sigma_2^2 \end{pmatrix}}(d, a) \right) && \text{(category 3)} \quad (1)
 \end{aligned}$$

where  $N_{\Sigma}(d, a)$  denotes the height of the bivariate normal *pdf* centered at  $\begin{pmatrix} 0 \\ 0 \end{pmatrix}$  with covariance matrix  $\Sigma$  at  $(d, a)$ .  $\Theta$  is the vector of values  $(\pi_0, \pi_1, \tau, \sigma_1, \sigma_2, \rho)$ . Under  $H_0$ , we have  $\rho = 0$  and  $\sigma_2 = 1$ . The values  $(\pi_0, \pi_1, \pi_2)$  represent the proportion of SNPs in each category, with  $\sum \pi_i = 1$ .

We use the product of values of the above *pdf* for a set of observed  $(Z_d Z_a)'$  as an objective function ('pseudo-likelihood', PL) to estimate the values of parameters. This is not a true likelihood as observations are not independent due to linkage disequilibrium (LD), although because we minimise the degree of LD between SNPs using the LDAK

method [4], the PL is close to a true likelihood.

## Model fitting and significance testing

We fit parameters  $\pi_0$ ,  $\pi_1$ ,  $\pi_2 (= 1 - \pi_0 - \pi_1)$ ,  $\sigma_1$ ,  $\sigma_2$ ,  $\tau$  and  $\rho$  under  $H_1$  and  $H_0$ . Under  $H_0$ ,  $(\rho, \sigma_2) = (0, 1)$ .

We then compare the fit of the two models using the log-ratio of PLs, giving an unadjusted pseudo-likelihood ratio (uPLR). We subtract a term to minimise the influence of the  $Z_a$  score distribution, and add a term to ensure the model is identifiable [5]. We term the resultant test statistic the pseudo-likelihood ratio (PLR).

$$PLR \sim \begin{cases} \gamma\chi_1^2 & \text{prob} = \kappa \\ \gamma\chi_2^2 & \text{prob} = 1 - \kappa \end{cases}. \quad (2)$$

The value  $\gamma$  arises from the weighting derived from the LDAK procedure causing a scale change in the observed  $PLR$ . The mixing parameter  $\kappa$  corresponds to the probability that  $\rho = 0$ , which is approximately  $\frac{1}{2}$ .

We estimate  $\gamma$  and  $\kappa$  by sampling random subgroups of the case group. Such subgroups only cover the subspace of  $H_0$  with  $\tau = 1$ , for which the asymptotic approximation of PLR by 2 is poor. We thus estimate the distribution of a similar alternative test statistic, the cPLR (see methods section and supplementary material, section 1.5.1), which majorises the distribution of PLR. Testing an observed PLR against this null distributions leads to a slightly conservative test statistic. The cPLR is effectively a pseudo-likelihood ratio conditioning on observed  $Z_a$  and asymptotically has the distribution 2.

If evidence for differential genetic basis can be established using the whole SNP set, as above, a natural next step is to search for the specific variants leading to the differ-

ence. An effective test statistic for testing subgroup differentiation for single SNPs is the Bayesian conditional false discovery rate (cFDR) [6, 7] applied to  $Z_d$  scores ‘conditioned’ on  $Z_a$  scores. This enables the association threshold for  $Z_d$  scores to be adjusted based on observed  $Z_a$  scores. However, this statistic alone cannot capture all the means by which the joint distribution of  $(Z_a, Z_d)$  can deviate from  $H_0$ , and we also propose three other test statistics, each with different advantages and compare their performance (Supplementary Note 4.1; Supplementary Figure 6).

## Power calculations and validation of method

We tested our method by application to a range of datasets, using simulated and resampled GWAS data. First, to confirm appropriate control of type 1 error rates across  $H_0$ , we simulated genotypes of case and control groups under  $H_0$  for a set of 50,000 autosomal SNPs in linkage equilibrium (supplementary note 2). Quantiles of the empirical PLR distribution were smaller than those for the empirical cPLR distribution and the asymptotic mixture- $\chi^2$ , leading to a conservative test when  $\tau > 1$  (estimated type 1 error rate 0.048, 95% CI 0.039-0.059) and when  $\tau \approx 1$  (estimated type 1 error rate 0.033, 95% CI 0.022-0.045); see figure 2. The distribution of cPLR closely approximated the asymptotic mixture- $\chi^2$  distribution across all values of  $\tau$  (supplementary material 2.1).

We then established the suitability of the test when SNPs are in LD and when there exist genetic differences between subgroups that are independent of disease status overall. First, we used a dataset of controls and autoimmune thyroid disease (ATD) cases. We repeatedly choose subgroups such that several SNPs had large allelic differences between subgroups, and estimated the null distribution of the PLR using the distribution of cPLR values from random subgroupings. The resultant p-values demonstrated good FDR control at all cutoffs (supplementary note, figure 2.2) and the overall type 1 error rate was 0.041

(95% CI 0.034-0.050). Second, we analysed a dataset of T1D cases with subgroups defined by geographical origin. Within the British Isles, there is clear genetic diversity associated with region [8]. As expected,  $Z_d$  scores for geographic subgroups showed inflation compared to for random subgroups (supplementary figure 3). None of the derived test statistics reached significance at a Bonferroni-corrected  $p < 0.05$  threshold (max. corrected p value  $> 0.8$ , supplementary figure 1).

To examine the power of our method, we used published GWAS data from the Wellcome Trust Case Control Consortium [9] comprising 1994 cases of Type 1 diabetes (T1D), 1903 cases of rheumatoid arthritis (RA), 1922 cases of type 2 diabetes (T2D) and 2953 common controls. We established that our test could differentiate between any pair of diseases, considered as subgroups of a general disease case group (all  $< 1 \times 10^{-8}$ , table 1).

		$\pi_0$	$\pi_1$	$\pi_2$	$\sigma_1$	$\sigma_2$	$\tau$	$\rho$	p-val
T1D/RA	$H_1$	0.997	$5.69 \times 10^{-4}$	$2.06 \times 10^{-3}$	2.76	1.39	1.74	1.815	$3.2 \times 10^{-12}$
	$H_0$	0.997	$6.26 \times 10^{-4}$	$2.48 \times 10^{-3}$	2.71	-	1.67	-	
T1D/T2D	$H_1$	0.573	0.426	$9.63 \times 10^{-4}$	1.00	2.03	2.25	1.68	$1.6 \times 10^{-9}$
	$H_0$	0.578	0.421	$8.91 \times 10^{-4}$	1.00	-	2.21	-	
T2D/RA	$H_1$	0.573	0.426	$8.71 \times 10^{-4}$	1.00	2.23	1.75	1.69	$5.1 \times 10^{-9}$
	$H_0$	0.91	$8.05 \times 10^{-4}$	0.0892	2.25	-	0.97	-	
GD/HT	$H_1$	0.506	0.487	0.007	1.12	2.90	1.65	2.61	$2.2 \times 10^{-15}$
	$H_0$	0.493	0.079	0.428	1.68	-	1.03	-	

Table 1: Fitted parameter values for models of T1D/RA, T1D/T2D, T2D/RA, and GD/HT.  $H_1$  is the null hypothesis (under which  $\sigma_2 = 1$ ,  $\rho = 0$ ) that SNPs differentiating the subgroups are not associated with the overall phenotype;  $H_1$  is the alternative (full model).  $p$  values for pseudo-likelihood ratio tests are also shown.

T1D and RA have considerable overlap in genetic basis [9, 10, 7], as well as many regions associated with only one of the two conditions. T1D and T2D have much less overlap [10] and T2D and RA less still. This was reflected in the fitted values (table 1). The fitted values parametrising category 2 in the full model for T1D/RA ( $\pi_1$ ,  $\sigma_1$ ) were consistent



with a subset of SNPs associated with case/control status (T1D+RA vs control) but not differentiating T1D/RA. By contrast, the parametrisation of category 2 for T1D/T2D and T2D/RA had marginal variance  $\sigma_1$  approximately 1, indicating that SNPs in categories 1 and 2 were essentially indistinguishable, and a subset of SNPs associated with case/control status but not with heterogeneity did not exist in these latter cases.

To explore the power of our test more generally, we first showed that power depends on the number of SNPs in category 3 and on the underlying parameters of the true model, depending on the number of samples only through the fitted model parameters (Supplementary Note 2.3). We therefore estimated the power of the test for varying numbers of SNPs in category 3 and for varying values of the parameters  $\sigma_2$ ,  $\tau$ , and  $\rho$ . (Figure 3; Supplementary Figure 2). As expected, power increases markedly with an increasing number of SNPs in category 3, reflecting the proportion of SNPs which differentiate case subgroups and are associated with the phenotype as a whole. Power also increases with increasing  $\tau$ ,  $\sigma_2$ , and correlation ( $\rho/(\sigma_2\tau)$ ) as high values make it easier to distinguish SNPs in the third category from those in the first and second.

We explored the dependence of power on sample size using the WTCCC data for RA and T1D, as above. We repeatedly drew subgroups of cases of varying size (originally c. 2000 cases for each), and computed the PLR using these subgroups. We then estimated the power of the PLR-based test to reject the null hypothesis at  $p \leq 0.05$  at each of these effect sizes (Figure 3) and compared it with the power to find any single SNP which differentiated the two diseases in several ways (see figure legend). We then repeated the analysis removing the known T1D- and RA- associated SNP rs17696736. Although the power of the PLR-based test was limited at reduced sample sizes, it remained consistently higher than the power to detect any single SNP which differentiated the two diseases. The power to detect a SNP with significant  $Z_d$  score (Bonferroni-corrected) amongst SNPs with

GW-significant  $Z_a$  score (sub-BCS) dropped dramatically, though the power of PLR was only slightly reduced. This indicated the non-reliance of the PLR test on large single-SNP effect sizes.

Estimating power requires an estimate of the underlying values of several parameters: the expected total number of SNPs in the pruned dataset with different population MAF in case subgroups, and the distribution of odds-ratios such SNPs between subgroups and between cases/controls. With sparse genome-wide cover, such as that in the WTCCC study, and discounting genomic regions with very large effect sizes (MHC),  $> 1250$  cases per subgroup are necessary for 90% power. If SNPs with greater coverage for the disease of interest are used (such as the ImmunoChip for autoimmune diseases) values of  $\pi_2$ ,  $\sigma_2$  and  $\tau$  are correspondingly higher, and around 500-700 cases per subgroup may be sufficient.

### **Properties and assumptions of the PLR test**

Our assumption that  $(Z_a, Z_d)$  follows a multivariate mixture Gaussian distribution is generally reasonable for complex phenotypes with a large number of associated variants [11] although our adjustment of the PL for the distribution of  $Z_a$ , similar to performing the analysis conditional on observed  $Z_a$ , reduces our reliance on this assumption. If subgroup prevalence is unequal between the study group and population, the expected distribution changes, but our method is still effective (supplementary material, section 1.4).

Our test is robust to confounders arising from differential sampling to the same extent as conventional GWAS. For example, if subgroups were defined based on population structure, and population structure also varied between the case and control group, SNPs which differed by ancestry would also appear associated with the disease, which could lead to an inappropriately high type 1 error rate. However, the same study design would also lead to identification of spurious association of ancestry-associated SNPs with the phenotype in a

conventional GWAS analysis. As for GWAS, this effect can be alleviated by accounting for the confounding trait (ancestry) as a covariate when computing p-values. Indeed, the only assumption we make in the form of the test statistic generating  $Z_d$  and  $Z_a$  is that they are linear functions of the allelic frequency difference between cases and controls or between subgroups (Supplementary Note 1). Therefore any covariates may be included without violating the assumptions of the test.

### **Application to autoimmune thyroid disease and type 1 diabetes**

Autoimmune thyroid disease (ATD) takes two major forms: Graves' disease (GD; clinically hyperthyroid) and Hashimoto's Thyroiditis (HT; clinically hypothyroid). The differential genetics of these two conditions have been investigated [12]. By contrast, T1D is relatively clinically homogenous with no major recognised subtypes, although heterogeneity arises between patients in levels of disease-associated autoantibodies, and disease course differs with age at diagnosis [3]. We analysed both of these diseases using our method.

For ATD, we were able to confidently detect evidence for differential genetic bases for GD and HT ( $p = 2.2 \times 10^{-15}$ ). Fitted values are shown in table 1. The distribution of cPLR statistics from random subgroups agreed well with the proposed mixture  $\chi^2$  (supplementary figure 10).

For T1D, we considered four subgroupings defined by plasma levels of the T1D-associated autoantibodies thyroid peroxidase antibody (TPO-Ab, n=5780), insulinoma-associated antigen 2 antibody (IA2-Ab, n=3197), glutamate decarboxylase antibody (GAD-Ab, n=3208) and gastric parietal cell antibodies (PCA-Ab, n=2240). A previous GWAS study on autoantibody positivity in T1D identified only two non-MHC loci at genome-wide significance: 1q23/*FCRL3* with IA2-Ab and 9q34/*ABO* with PCA-Ab - and 11 autoantibody positivity associations amongst known T1D-associated loci at  $FDR \approx 16\%$  [3].

We tested each of the four subgroupings both retaining and excluding the MHC region. Fitted values for models with and without the MHC region are shown in supplementary table 2, and plots of the relevant  $Z_a$  and  $Z_d$  scores are shown in supplementary figure 4. When retaining the MHC region, we were able to confidently reject  $H_0$  for subgroupings based on TPO-Ab, IA-2Ab and GAD-Ab ( $p$  values all  $< 1.0 \times 10^{-20}$ ). Although there was evidence that SNPs in the dataset were associated with PCA-Ab level ( $\tau \approx 2.5$ , null model), the improvement in fit in the full model was not significant, and we conclude that such SNPs determining PCA-Ab status are not in general associated with T1D. This pattern can be seen by in the plot of  $Z_a$  against  $Z_d$  (supplementary figure 4) in which SNPs with a high  $Z_d$  value do not have higher than expected  $Z_a$  values.

When the MHC region was removed, the subgrouping based on TPO-Ab remained significantly better-fit by the full model ( $p = 1.5 \times 10^{-4}$ ). There was weaker evidence to reject the null model for GAD-Ab ( $p = 0.002$ ) and IA2-Ab ( $p = 0.008$ ) (Bonferroni-corrected threshold: 0.006). The fitted values of  $\tau$  in both the full and null models for GAD-Ab were  $\approx 1$ , indicating a lack of evidence for a category of non-MHC T1D-associated SNPs additionally associated with GAD-Ab positivity. Collectively, this indicates that differential genetic basis for T1D with GAD-Ab and IA2-Ab positivity is driven almost entirely by the MHC region, and although PCA-Ab status may be genetically determined, the set of causative variants is independent of T1D causative pathways.

Finally, we applied the method to the same dataset with subgroups defined by age at diagnosis. Fitted values are shown in supplementary table 3 and a plot of  $Z_a$  and  $Z_d$  scores in supplementary figure 5. Rather than explicitly splitting cases into two groups, we considered age as an interval variable to compute  $Z_d$ . The hypothesis  $H_0$  could be rejected confidently when retaining the MHC region, and less so after removing it ( $p$  values  $< 1.0 \times 10^{-20}$  and 0.007 respectively). Signed  $Z_d$  and  $Z_a$  scores for age at diagnosis

showed a visible negative correlation ( $p = 0.002$ ; SNPs weighted by LDAK and probability of category 2 membership in null model) amongst  $Z_d$  and  $Z_a$  scores for disease-associated SNPs (figure 5). This is consistent with a higher genetic liability with lower age at diagnosis; the basis for differential diagnosis age may be differing heritability in the context of the same relative variant effect sizes.

### Assessment of individual SNPs

Many SNPs which discriminated between subgroups were in known disease-associated regions (Supplementary Tables 4, 5, and 6). In several cases, our method identified disease-associated SNPs which have reached genome-wide significance in subsequent larger studies but for which the  $Z_a$  score in the WTCCC study was not near significance. For example, the SNP rs3811019, in the *PTPN22* region, was identified as likely to discriminate T1D and T2D ( $p = 3.046 \times 10^{-6}$ ; see supplementary table 5), despite a p value of  $3 \times 10^{-4}$  for joint T1D/T2D association.

For GD and HT, SNPs were near the known ATD-associated loci *PTPN22* (rs7554023), *CTLA4* (rs58716662), and *CEP128* (rs55957493) as likely to be contributing to the difference (see supplementary table 7). The SNPs rs34244025 and rs34775390 are not known to be associated with ATD, but are in known loci for inflammatory bowel disease and ankylosing spondylitis, and our data suggest they may differentiate GD and HT (FDR  $3 \times 10^{-3}$ ).

We then sought to find the SNPs outside of the MHC region with differential effect sizes with TPOA positivity in T1D, the only subgrouping of T1D for which we could confidently reject  $H_0$ .

Previous work [3] identified several loci potentially associated with TPO-Ab positivity by restricting attention to known T1D loci, enabling use of a larger dataset than was

available to us. We list the top ten SNPs for each summary statistic for TPO-Ab positivity in supplementary table 8. Subgroup-differentiating SNPs included several near known T1D loci: *CTLA4* (rs7596727), *BACH2* (rs11755527), *RASGRP1* (rs16967120) and *UBASH3A* (rs2839511). These loci agreed with those found by Plagnol et al, but our analysis used only the available genotype data, without external information on confirmed T1D loci. Note that we were not able to replicate the same p-values due to a reduced availability of cases.

Finally, we analysed non-MHC SNPs with varying effect sizes with age at diagnosis in T1D (supplementary table 9). This implicated SNPs in or near *CTLA4* (rs2352551), *IL2RA* (rs706781), and *IKZF3* (rs11078927).

## Genetic correlation testing

Given the linear correlation between  $Z_d$  defined by T1D age at diagnosis and  $Z_a$  defined by T1D case/control comparison, it appears that recent methods to test for genetic correlation between diseases [13, 14] could be adapted to the subgrouping question. One potential approach is to estimate the narrow-sense genetic correlation ( $r_g$ ) across a set of SNPs between case/control traits of interest, either between  $Z$  scores derived from comparing the control group to each case subgroup, testing under the null hypothesis  $r_g = 1$  (method 1); or between the familiar  $Z_a$  and  $Z_d$ , under the null hypothesis  $r_g = 0$  (method 2). This approach should have the advantage of characterising heterogeneity using a single widely-interpretable metric.

We explored such adaptation in detail in supplementary note 3. We find that naive application of method 1 could lead to systematic false positives because subgroupings that are independent of the overall disease process still lead to rejections of the null hypothesis (e.g. hair colour in T2D). For method 2, we find that it is considerably less powerful

because it tests a rather narrower definition of  $H_1$ , in which the correlation between  $Z_d$  and  $Z_a$  is assumed to be either always positive or always negative for all SNPs in category 3, in contrast to the symmetric model for correlation that we allow (Figure 1). We note that neither of these methods were explicitly designed for subgroup testing, and our naive and direct application of them in this context does not suggest that they are less than optimal for the context in which they were designed - testing genetic correlation between two traits. The theoretical comparison did, however, usefully place our method in the context of these established methodologies, demonstrating that our parameter  $\rho$  is effectively an estimator of  $|r_g|$  (allowing for symmetry in correlation), and demonstrating the necessity of considering both variance parameters ( $\tau$ ,  $\sigma_2$ ) and covariance parameters ( $\rho$ ) in testing a subgrouping of interest.

## Discussion

The problem we address is part of a wider aim of adapting GWAS to complex disease phenotypes. As the body of GWAS data grows the analysis of between-disease similarity and within-disease heterogeneity has led to substantial insight into shared and distinct disease pathology [6, 7, 2, 15, 16]. We seek in this paper to use genomic data to infer whether such disease subtypes exist. Our problem is related to the question of whether two different diseases share any genetic basis [13] but differs in that the implicit null hypothesis relates to genetic homogeneity between subgroups rather than genetic independence of separate diseases.

We use ‘absolute covariance’  $\rho$  preferentially to standard covariance (see supplementary table 1) because we expect that  $Z_a$  and  $Z_d$  will frequently co-vary positively and negatively at different SNPs in the same analysis; for instance, if some variants are deleterious only for subgroup 1 and others deleterious only for subgroup 2. Our method tests whether

a subgrouping in question corresponds more closely to the third and fourth rows of supplementary table 1 than the first and second, but may be adaptable to other families of effect size distributions. A potential advantage of our method using unsigned  $Z$  scores we have not yet explored is the potential to generate  $Z_d$  scores from ANOVA-style tests for genetic homogeneity between three or more subgroups, rather than two, in which case the reconstructed  $Z$  scores would not have a direction.

Aetiologically and genetically heterogeneous subgroups within a case group correspond to substructures in the genotype matrix. Information about such substructures are lost in a standard GWAS, which only uses the column-sums (MAFs) of the matrix (linear-order information). Data-driven selection of appropriate case subgroups and corresponding analyses of these subgroups can use more of the remaining quadratic-order information the matrix contains. Indeed a 'two-dimensional' GWAS approach (using  $Z_a$  and  $Z_d$ ) instead of a standard GWAS (using only  $Z_a$ ) may improve SNP discovery, as we found for *PTPN22* in RA/T2D. However, this can only be the case if the subgroups genuinely correspond to different variant effect sizes; for other subgroupings, a two-dimensional GWAS will only add noise. While it might seem appealing to use this method to search for optimal partitions of patients in an unsupervised manner, we prefer to focus on testing subgroupings derived from independent clinical or phenotypic data. First, it is difficult to characterise subgroupings as "better" or "worse" than each other. Our proposed model necessarily complex, in order to account for the range of processes leading to disease heterogeneity (see supplementary table 1). No one parameter can parametrise the degree to which two subgroups differ; parameters  $\pi_2$ ,  $\tau$ , and  $\rho$  all contribute, and attempts to test the hypothesis using a single measure such as genetic correlation have serious shortcomings (see supplementary material, 3). Second, even if subgroups could be ranked by some metric, the search space of potential subgroupings of a case group is prohibitively large ( $2^N$  for  $N$  cases) so it is generally not



a tractable problem to search for the subgrouping which maximises some metric.

We demonstrated that effect sizes of T1D-causative SNPs differ with age at disease diagnosis. The strong negative correlation observed (figure 5) was consistent with an increased total genetic liability in samples with earlier age of diagnosis, a finding supported by candidate gene studies [17, 18, 19] and epidemiological data [20]. In this instance, the hypothesis  $H_0$  was most effectively tested by searching directly for correlation between  $Z_a$  and  $Z_d$ , as opposed to the PLR method. However, more generally, asserting correlation gives little insight into the structure of the two dimensional effect size distribution, and correlation need not necessarily be present to reject  $H_0$  (for example,  $\rho = 0$ ,  $\sigma_2 > 1$ , model 1). Whether such an age dependent effect on total liability corresponds to a sub-grouping is a somewhat philosophical point. Such a pattern arises naturally from a liability threshold model where total liability depends on both genetic effects and environmental influences which accumulate with age. Retrospective likelihood methods [21] have been proposed for such cases, and may provide a more powerful test of association, while our method focuses on identifying the pattern at the genomewide level.

Our method adds to the current body of knowledge by extracting additional information from a disease dataset over a standard GWAS analysis, and determines if further analysis of disease pathogenesis in subgroups is justified. Our approach is analogous to the intuitive and well-applied method of searching for between-subgroup differences in SNPs with known disease associations [3] but we do so without restricting attention to strong disease associations, enabling use of information from disease-associated SNPs which do not reach significance. Our parametrisation of effect size distributions allows insight into the structure of the genetic basis of the disease and potential subtypes, improving understanding of the allelic spectrum of diseases and genotype-phenotype relationships.

## Methods

### Ethics Statement

This paper re-analyses previously published datasets. All patient data were handled in accordance with the policies and procedures of the participating organisations.

### Joint distribution of variables $Z_a, Z_d$

We assume that SNPs may be divided into three categories, as described in equation 1 (figure 1). Under these assumptions,  $Z_a$  and  $Z_d$  scores have the joint *pdf*

$$\begin{aligned}
 PDF_{Z_d, Z_a | \Theta}(d, a) &= \pi_0 N_{\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}}(d, a) && \text{(category 1)} \\
 &+ \pi_1 N_{\begin{pmatrix} 1 & 0 \\ 0 & \sigma_1^2 \end{pmatrix}}(d, a) && \text{(category 2)} \\
 &+ \pi_2 \left( \frac{1}{2} N_{\begin{pmatrix} \tau^2 & \rho \\ \rho & \sigma_2^2 \end{pmatrix}}(d, a) + \frac{1}{2} N_{\begin{pmatrix} \tau^2 & -\rho \\ -\rho & \sigma_2^2 \end{pmatrix}}(d, a) \right) && \text{(category 3)}
 \end{aligned}$$

where  $\Theta$  is the vector of values  $(\pi_0, \pi_1, \pi_2, \tau, \sigma_1, \sigma_2, \rho)$ . Z scores  $Z_a$  and  $Z_d$  are reconstructed from GWAS p-values for SNP associations. In practice, since our model is symmetric, we only use require absolute Z scores, without considering effect direction.

For sample sizes  $n_1, n_2$  and 97.5% odds-ratio quantile  $\alpha$ , the expected observed standard deviation of Z scores (that is,  $\sigma_1, \sigma_2$ , and  $\tau$ ) is given by

$$E\{SD(Z)\} = \sqrt{1 + \frac{\log(\alpha)^2 n_1 n_2}{12(n_1 + n_2)}} \quad (3)$$

(supplementary note 2.3)

## Definition and distribution of PLR statistics

For a set of observed  $Z$  scores  $Z = (Z_a, Z_d)$  we define the joint unadjusted pseudo-likelihood  $PL_{da}(Z|\Theta)$  as

$$PL_{da}(Z|\Theta) = \prod_{Z_d^{(i)} \in Z_d, Z_a^{(i)} \in Z_a} PDF_{Z_d, Z_a|\Theta}(Z_d^{(i)}, Z_a^{(i)}) + C \log(\pi_0 \pi_1 \pi_2) \quad (4)$$

where the term  $C \log(\pi_0 \pi_1 \pi_2)$  is added to ensure identifiability of the model [5].

We now set

$$\begin{aligned} \hat{\theta}_1 &= \arg \max_{\Theta \in H_1} PL_{da}(Z_d, Z_a|\theta) \\ \hat{\theta}_0 &= \arg \max_{\Theta \in H_0} PL_{da}(Z_d, Z_a|\theta) \\ uPLR(Z) &= \log \left( \frac{PL_{da}(Z|\hat{\theta}_1)}{PL_{da}(Z|\hat{\theta}_0)} \right) \end{aligned} \quad (5)$$

recalling that  $H_0$  is the subspace of the parameter space  $H_1$  satisfying  $\sigma_2 = 1$  and  $\rho = 0$ .

If data observations are independent,  $uPLR$  reduces to a likelihood ratio. Under  $H_0$ , the asymptotic distribution of  $uPLR$  is then

$$uPLR \sim \frac{1}{2} \begin{cases} \chi_1^2 & p = 1/2 \\ \chi_2^2 & p = 1/2 \end{cases} \quad (6)$$

according to Wilk's theorem extended to the case where the null value of a parameter lies on the boundary of  $H_1$  (since  $\rho = 0$  under  $H_0$ ) [22].

The empirical distribution of  $uPLR$  may substantially majorise the asymptotic distribution when  $\tau \approx 1$ . In the full model, the marginal distribution of  $Z_a$  has more degrees of freedom (four;  $\pi_0, \pi_1, \sigma_1, \sigma_2$ ) than it does under the null model (two;  $\pi_1, \sigma_1; \sigma_2 \equiv 1$ ).

This can mean that certain distributions of  $Z_a$  can drive high values of  $uPLR$  independent of the values of  $Z_d$  (supplementary material, section 2), which is unwanted as the values  $Z_a$  reflect only case/control association and carry no information about case subgroups. If observed  $uPLR$ s from random subgroups (for which  $\tau = 1$  by definition) are used to approximate the null  $uPLR$  distribution, this effect would lead to serious loss of power when  $\tau \gg 1$ .

This effect can be largely alleviated by subtracting a correcting factor based on the pseudo-likelihood of  $Z_a$  alone, which reflects the contribution of  $Z_a$  values to the uPLR. We define

$$PL_a(Z_a|\Theta) = \prod_{Z_a^{(i)} \in Z_a} \left( \pi_0 N_{0,1}(Z_a^{(i)}) + \pi_1 N_{0,\sigma_1^2}(Z_a^{(i)}) + \pi_2 N_{0,\sigma_2^2}(Z_a^{(i)}) \right) \quad (7)$$

that is, the marginal likelihood of  $Z_a$ . Given  $\hat{\theta}_1, \hat{\theta}_0$  as defined above, we define

$$f(Z_a) = \min \left( \log \frac{PL_a(Z_a|\hat{\theta}_1)}{PL_a(Z_a|\hat{\theta}_0)}, 0 \right) \quad (8)$$

We now define the PLR as

$$PLR = uPLR - f(Z_a) \quad (9)$$

The action of  $f(Z_a)$  leads to the asymptotic distribution of PLR slightly minorising the asymptotic mixture- $\chi^2$  distribution of uPLR, to differential degrees dependent on the value of  $\tau$  (see supplementary note 2).

We define the similar test statistic  $cPLR$ :

$$\begin{aligned}
 cPL(Z_a|Z_d, \theta) &= \frac{PL_{da}(Z_a, Z_d|\theta)}{PL_a(Z_a|\theta)} \\
 \hat{\theta}_1^c &= \arg \max_{\theta \in H_1} cPL(Z_a|Z_d, \theta) \\
 \hat{\theta}_0^c &= \arg \max_{\theta \in H_0} cPL(Z_a|Z_d, \theta) \\
 cPLR &= \log \left( \frac{cPL(Z_d|Z_a, \hat{\theta}_1^c)}{cPL(Z_d|Z_a, \hat{\theta}_0^c)} \right)
 \end{aligned} \tag{10}$$

noting that the expression  $\frac{PL_{da}(Z_a, Z_d|\theta)}{PL_a(Z_a|\theta)}$  can be considered as a likelihood conditioned on the observed values of  $Z_a$ . Now

$$\begin{aligned}
 PLR &= \log \left( \frac{PL_{da}(Z|\hat{\theta}_1)}{PL_{da}(Z|\hat{\theta}_0)} \right) - \log \left( \frac{PL_a(Z_a|\hat{\theta}_1)}{PL_a(Z_a|\hat{\theta}_0)} \right) \\
 &= \log \left( \frac{cPL(Z_d|Z_a, \hat{\theta}_1)}{cPL(Z_d|Z_a, \hat{\theta}_0)} \right)
 \end{aligned} \tag{11}$$

The empirical distribution of cPLR for random subgroups majorises the empirical distribution of PLR (supplementary note, section 2). Furthermore, the approximation of the empirical distribution of cPLR by its asymptotic distribution is good, across all values of  $\tau$ ; that is, across the whole null hypothesis space.

Our approach is to compare the PLR of a test subgroup to the cPLR of random subgroups, which constitutes a slightly conservative test under the null hypothesis (see supplementary material 2).

## Allowance for linkage disequilibrium

The asymptotic approximation of the pseudo likelihood-ratio distribution breaks down when values of  $Z_a$ ,  $Z_d$  are correlated due to LD. One way to overcome this is to ‘prune’

SNPs by hierarchical clustering until only those with negligible correlation remain. A disadvantage with this approach is that it is difficult to control which SNPs are retained in an unbiased way without risking removal of SNPs which contribute greatly to the difference between subgroups.

We opted to use the LDAK algorithm [4], which assigns weights to SNPs approximately corresponding to their ‘unique’ contribution. Denoting by  $\rho_{ij}$  the correlation between SNPs  $i, j$ , and  $d(i, j)$  their chromosomal distance, the weights  $w_i$  are computed so that

$$w_i + \sum_{i \neq j} w_j \rho_{ij}^2 e^{-\lambda d(i,j)} \quad (12)$$

is close to constant for all  $i$ , and  $w_i > 0$  for all  $i$ . The motivation for this approach is that  $\sum_{i \neq j} \rho_{ij}^2$  represents the replication of the signal of SNP  $i$  from all other SNPs.

This approach has the advantage that if  $n$  SNPs are in perfect LD, and not in LD with any other SNPs, each will be weighted  $1/n$ , reducing the overall contribution to the likelihood to that of one SNP. In practice, the linear programming approach results in many SNP weights being 0. Using the LDAK algorithm therefore allows more SNPs to be retained and contribute to the model than would be retained in a pruning approach.

A second advantage of LDAK is that it homogenises the contribution of each genome region to the overall pseudo-likelihood. Many modern microarrays fine-map areas of the genome known or suspected to be associated with traits of interest [23] which could theoretically lead to peaks in the distribution of SNP effect sizes, disrupting the assumption of normality. LD pruning and LDAK both reduce this effect by homogenising the number of tags in each genomic region.

We adapted the pseudo-likelihood function to the weights by multiplying the contribution of each SNP to the log-likelihood by its weight; that is, if SNP  $i$  had weight  $w_i$  and  $Z$

scores  $Z_{ai}$ ,  $Z_{di}$ , we computed the final log pseudo-likelihood for parameters  $\Pi$ ,  $\Theta$  as:

$$\log\{PL_{da}(Z_a, Z_d|\Theta)\} = C \log(\pi_0\pi_1\pi_2) + \sum_{Z_d^{(i)} \in Z_d, Z_a^{(i)} \in Z_a} w_i \log(PDF_{Z_d, Z_a|\Theta}(Z_d^{(i)}, Z_a^{(i)})) \quad (13)$$

essentially counting the  $i$ th SNP  $w_i$  times over.

Adjusting using LDAK was effective in enabling the distributions of PLR to be well-approximated by mixture- $\chi^2$  distributions of the form 2 (supplementary plots 9, 10, 11).

### E-M algorithm to estimate model parameters

We use an expectation-maximisation algorithm [24, 25] to fit maximum-PL parameters. Given an initial estimate of parameters  $\Theta_0 = (\pi_0^i, \pi_1^i, \tau^i, \sigma_1^i, \sigma_2^i, \rho^i)$  we iterate three main steps:

1. Define for SNP  $s$

$$\zeta_g^{(s)} = Pr(s \in \text{category } g | \Theta_i) \quad (14)$$

2. For  $g \in (1, 2, 3)$  set

$$\pi_g^{i+1} = \overline{\zeta_g^{(s)}} \quad (15)$$

3. Set

$$(\tau^{i+1}, \sigma_1^{i+1}, \sigma_2^{i+1}, \rho^{i+1}) = \arg \max_{(\tau, \sigma_1, \sigma_2, \rho)} PL(Z_d, Z_a | \pi_0^{i+1}, \pi_1^{i+1}, \tau, \sigma_1, \sigma_2, \rho) \quad (16)$$

Step 3 is complicated by the lack of closed form expression for the maximum likelihood estimator of  $\rho$  (because of the symmetric two-Gaussian distribution of category 3), requiring a bisection method for computation. The algorithm is continued until  $|PLR(Z_d, Z_a|\Theta_i) - PLR(Z_d, Z_a|\Theta_{i-1})| < \epsilon$ ; we use  $\epsilon = 1 \times 10^{-5}$ .

The algorithm can converge to local rather than global minima of the likelihood. We overcome this by computing the pseudo-likelihood of the data at 1000 of points throughout the parameter space, retaining the top 100, and dividing these into 5 maximally-separated clusters. The full algorithm is then run on the best (highest-PL) point in each cluster

The speed of convergence of the algorithm depends on appropriate choice of  $\Theta_0$  can speed up the algorithm considerably; for simulations, we begin the model at previous maximum-PL estimates of parameters for earlier simulations.

Maximum-cPL estimations of parameters were made using generic numerical optimisation with the *optim* function in R. Prior to applying the algorithm, parameters  $\pi_1$  and  $\sigma_1$  are estimated as maximum-PL estimators of the objective function

$$f(Z_a|\pi_1, \sigma_1) = \sum w_i \log\{(1 - \pi_1)N_{0,1}(Z_a(i)) + \pi_1 N_{0,\sigma_1^2}(Z_a(i))\} \quad (17)$$

where  $w_i$  is the weight for SNP  $i$  (see supplementary material, section 2 for rationale). The conditional pseudo-likelihood was maximised over the remaining parameters.

The algorithm and other processing functions are implemented in an R package available at <https://github.com/jamesliley/subtest>

## Prioritisation of single SNPs

An important secondary problem to testing  $H_0$  is the determination of which SNPs are likely to be associated with disease heterogeneity. Ideally, we seek a way to test the association of a SNP with subgroup status (ie,  $Z_d$ ), which gives greater priority to SNPs



potentially associated with case/control status (ie, high  $Z_a$ ).

An effective test statistic meeting these requirements is the Bayesian conditional false discovery rate (cFDR) [6]. It tests against the null hypothesis  $H'_0$  that the population minor allele frequencies of the SNP in both case subgroups are equal (ie, that the SNP does not differentiate subgroups), but responds to association with case/control status in a natural way by relaxing the effective significance threshold on  $|Z_d|$ . This relaxation of threshold only occurs if there is systematic evidence that high  $|Z_d|$  scores and high  $|Z_a|$  scores typically co-occur. The test statistic is direction-independent.

Given a set of observed  $Z_a$  and  $Z_d$  values  $z_{ai}, z_{di}, i \in \text{SNPs}$  with corresponding two-sided p values  $p_{ai}, p_{di}$ , the value of  $X_4$  for SNP  $j$  is defined as

$$X_4 = p_{dj} \frac{|\{i : p_{ai} \leq p_{aj} \wedge p_{di} \leq p_{dj}\}|}{|\{i : p_{di} \leq p_{dj}\}|} \quad (18)$$

$$\approx Pr(H'_0 | P_a \leq p_{aj}, P_d \leq p_{dj})$$

The value gives the false-discovery rate for SNPs whose p-values fall in the region  $[0, p_{dj}] \times [0, p_{aj}]$ ; this can be converted into a false-discovery rate amongst all SNPs for whom  $X_4$  passes some threshold [7].

We discuss three other single-SNP test statistics in supplementary note 4.1, which test against different null hypotheses. If the hypothesis  $H'_0$  is to be tested, then we consider the cFDR the best of these.

Contour plots of the test statistics for several datasets are shown in supplementary figures 6,7.

## Description of GWAS datasets

ATD samples were genotyped on the ImmunoChip [23] a custom array targeting putative autoimmune-associated regions. Data were collected for GWAS-like analyses of dense SNP data [12]. The dataset comprised 2282 cases of Graves' disease, 451 cases of Hashimoto's thyroiditis, and 9365 controls.

T1D samples were genotyped on either the Illumina 550K or Affymetrix 500K platforms, gathered for a GWAS on T1D [26]. We imputed between platforms in the same way as the original GWAS. The dataset comprised genotypes from 5908 T1D cases and 8825 controls, of which all had measured values of TPO-Ab, 3197 had measured IA2-Ab, 3208 had measured GAD-Ab, and 2240 had measured PCA-Ab. Comparisons for each autoantibody were made between cases positive for that autoantibody, and cases not positive for it. We did not attempt to perform comparisons of individuals positive for different autoantibodies (for instance, TPO-Ab positive vs IA2-Ab positive) because many individuals were positive for both.

To generate summary statistics corresponding to geographic subgroups, we considered the subgroup of cases from each of twelve regions and each pair of regions against all other cases (78 subgroupings in total). To maximise sample sizes, we considered T1D cases as 'controls' and split the control group into subgroups.

## Quality control

Particular care had to be taken with quality control, as Z-scores had to be relatively reliable for all SNPs assessed, rather than just those putatively reaching genome-wide significance.. For the T1D/T2D/RA comparison, which we re-used from the WTCCC, a critical part of the original quality control procedure was visual analysis of cluster plots for SNPs reaching significance, and systematic quality control measures based on differential call rates and

deviance from Hardy-Weinberg equilibrium (HWE) were correspondingly loose [9]. Given that we were not searching for individual SNPs, this was clearly not appropriate for our method.

We retained the original call rate (CR) and MAF thresholds (MAF  $\geq$  1%, CR  $\geq$  95% if MAF  $\geq$  5%, CR  $\geq$  99% if MAF  $<$  5%) but employed a stricter control on Hardy-Weinberg equilibrium, requiring  $p \geq 1 \times 10^{-5}$  for deviation from HWE in controls. We also required that deviance from HWE in cases satisfied  $p \geq 1.91 \times 10^{-7}$ , corresponding to  $|z| \leq 5$ . The looser threshold for HWE in cases was chosen because deviance from HWE can arise due to true SNP effects [27]. We also required that call rate difference not be significant ( $p \geq 1 \times 10^{-5}$ ) between any two groups, included case-case and case-control differences. Geographic data was collected by the WTCCC and consisted of assignment of samples to one of twelve geographic regions (Scotland, Northern, Northwestern, East and West Ridings, North Midlands, Midlands, Wales, Eastern, Southern, Southeastern, and London [9]). In analysing differences between autoimmune diseases, we stratified by geographic location; when assessing subgroups based on geographic location, we did not.

For the ATD and T1D data, we used identical quality control procedures to those employed in the original paper [12, 26]. We applied genomic control [28] to computation of  $Z_a$  and  $Z_d$  scores except for our analysis of ATD (following the original authors [12]) and our geographic analyses (as discussed above). In all analyses except where otherwise indicated we removed the MHC region with a wide margin ( $\approx 5Mb$  either side).

## Acknowledgments

We acknowledge the help of the Diabetes and Inflammation Laboratory Data Service for access and quality control procedures on the datasets used in this study. The JDRF/Wellcome Trust Diabetes and Inflammation Laboratory is in receipt of a Wellcome Trust Strategic

Award (107212) and receives funding from the JDRF (5-SRA-2015-130-A-N) and the NIHR Cambridge Biomedical Research Centre. The research leading to these results has received funding from the European Union's 7th Framework Programme (FP7/2007-2013) under grant agreement no. 241447 (NAIMIT). JL is funded by the NIHR Cambridge Biomedical Research Centre and is on the Wellcome Trust PhD programme in Mathematical Genomics and Medicine at the University of Cambridge. CW is funded by the Wellcome Trust (089989,107881) and the MRC (MC\_UP\_1302/5). The Cambridge Institute for Medical Research (CIMR) is in receipt of a Wellcome Trust Strategic Award (100140). The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

### **Conflicts of Interest**

The JDRF/Wellcome Trust Diabetes and Inflammation Laboratory receives funding from Hoffmann La Roche and Eli-Lilly and Company.

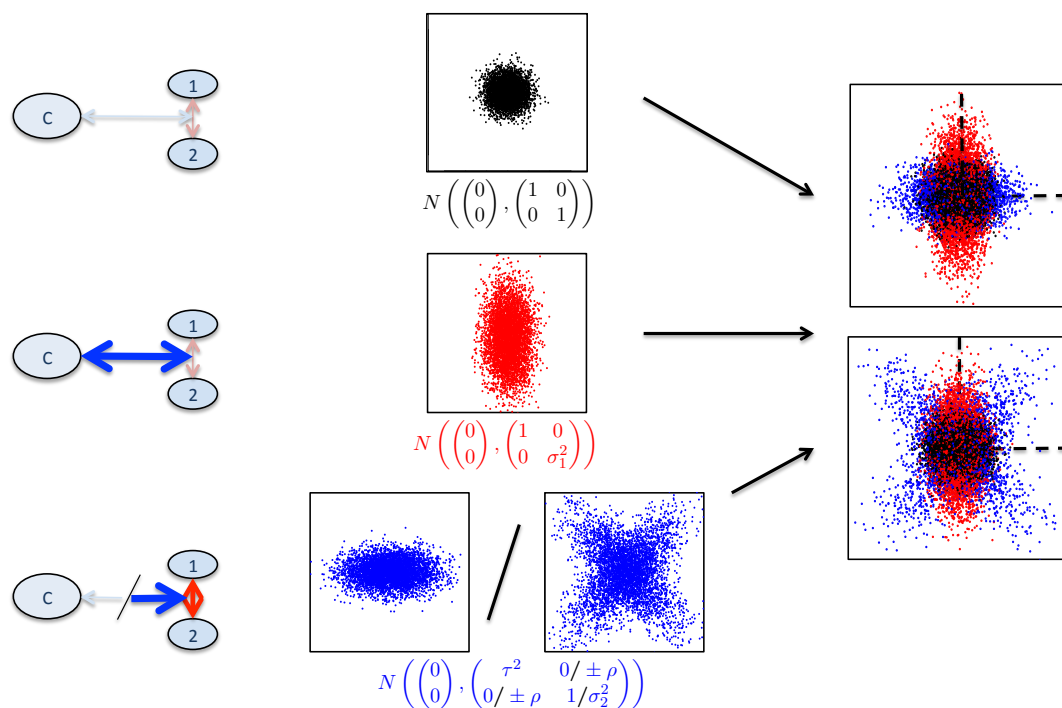


Figure 1: Overview of three-categories model.  $Z_d$  and  $Z_a$  are  $Z$  scores derived from GWAS p-values for allelic differences between case subgroups (1 vs 2), and between cases and controls (1 + 2 vs  $C$ ) respectively (top left). Within each category of SNPs, the joint distribution of  $(Z_d, Z_a)$  has a different characteristic form. In category 1,  $Z$  scores have a unit normal distribution; in category 2, the marginal variance of  $Z_a$  can vary. The distribution of SNPs in category 3 depends on the main hypothesis. Under  $H_0$  (that all disease-associated SNPs have the same effect size in both subgroups), only the marginal variance of  $Z_d$  may vary; under  $H_1$  (that subgroups correspond to differential effect sizes for disease-associated SNPs), any covariance matrix is allowed. The overall SNP distribution is then a mixture of Gaussians (rightmost panel). Visually, our test determines whether the observed overall  $Z_d, Z_a$  distribution more closely resembles the bottom rightmost panel than the top.

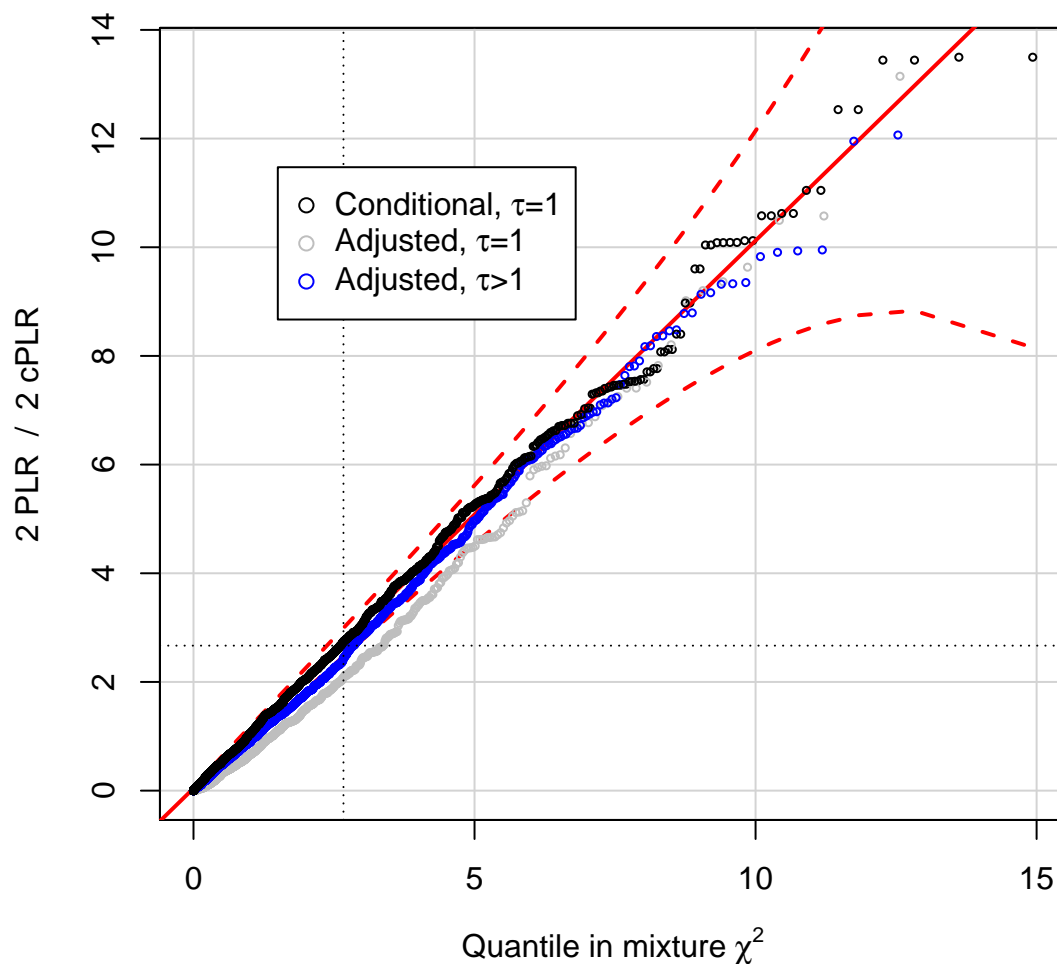


Figure 2: QQ plot from simulations demonstrating performance of PLR test. PLR values for test subgroups under  $H_0$  with either  $\tau = 1$  (random subgroups; grey) or  $\tau > 1$  (genetic difference between subgroups, but independent of main phenotype; blue) with cPLR values for random subgroups (black) and against proposed asymptotic distribution under simulation ( $\frac{1}{2}(\chi_1^2 + \chi_2^2)$ ); solid red line; 99% confidence limits dashed red line). The distribution of cPLR for random subgroups majorises the distribution of PLR, meaning the test statistic is conservative. Further details are shown in supplementary material, section 2.

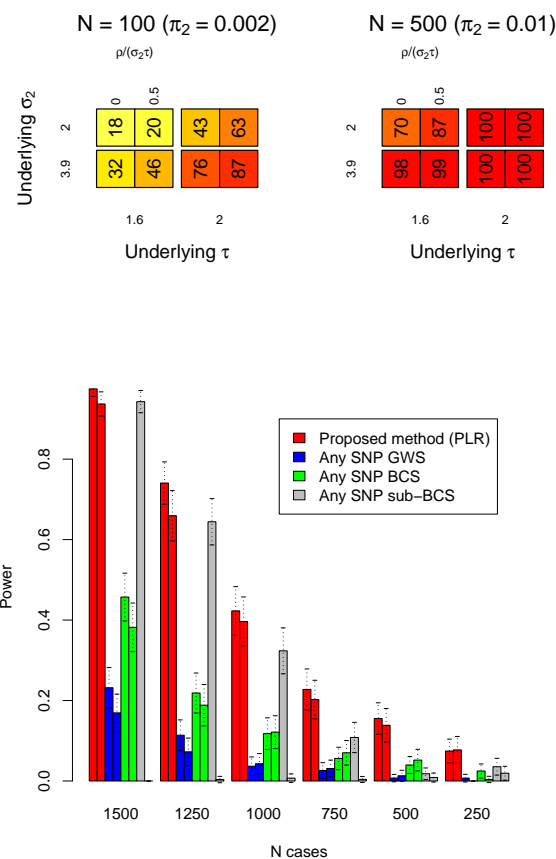


Figure 3: The power of the PLR to reject the null hypothesis (genetic homogeneity between subgroups) depends on the number of SNPs in category 3 and the true underlying values of the model parameters  $\sigma_1$ ,  $\sigma_2$ ,  $\tau$ , and  $\rho$ . Dependence on the number of case/control samples arises through the magnitudes of  $\sigma_2$  and  $\tau$  (supplementary material, section 2.3). The top figure shows estimates of power for various values of  $\pi_2$ ,  $\sigma_2$ ,  $\tau$ , and  $\rho$ . The value  $N$  is the approximate number of SNPs in category 3, corresponding to  $\pi_2$ . In total, each simulation was on  $5 \times 10^4$  simulated autosomal SNPs in linkage equilibrium. The value  $\rho/(\sigma_2\tau)$  is the correlation (rather than covariance) between  $Z_a$  and  $Z_d$  in category 3. More extensive plots of power are shown in supplementary figure 2. The lower figure shows the power of the PLR-based test to detect differences in genetic basis of T1D and RA subgroups of an autoimmune disease dataset, downsampling to varying numbers of cases (columns). Power is compared with: power to find  $\geq 1$  SNP with  $Z_d$  score reaching genome-wide significance (GWS, blue;  $p \leq 5 \times 10^{-8}$ ) or Bonferroni-corrected significance (BCS, green;  $p \leq 0.05/(\text{total \# of SNPs})$ ); and power to detect any SNP with  $Z_a$  score reaching genome-wide significance and  $Z_d$  score reaching Bonferroni-corrected significance (sub-BCS, grey;  $p \leq 0.05/(\text{total \# of SNPs with } Z_a \text{ reaching GWS})$ ). Error bars show 95% CIs. Left bars for each colour show power for all SNPs, rightmost for all SNPs except rs17696736. Power for sub-BCS drops dramatically but power for PLR is not markedly affected, indicating the robustness of PLR to single-SNP effects.

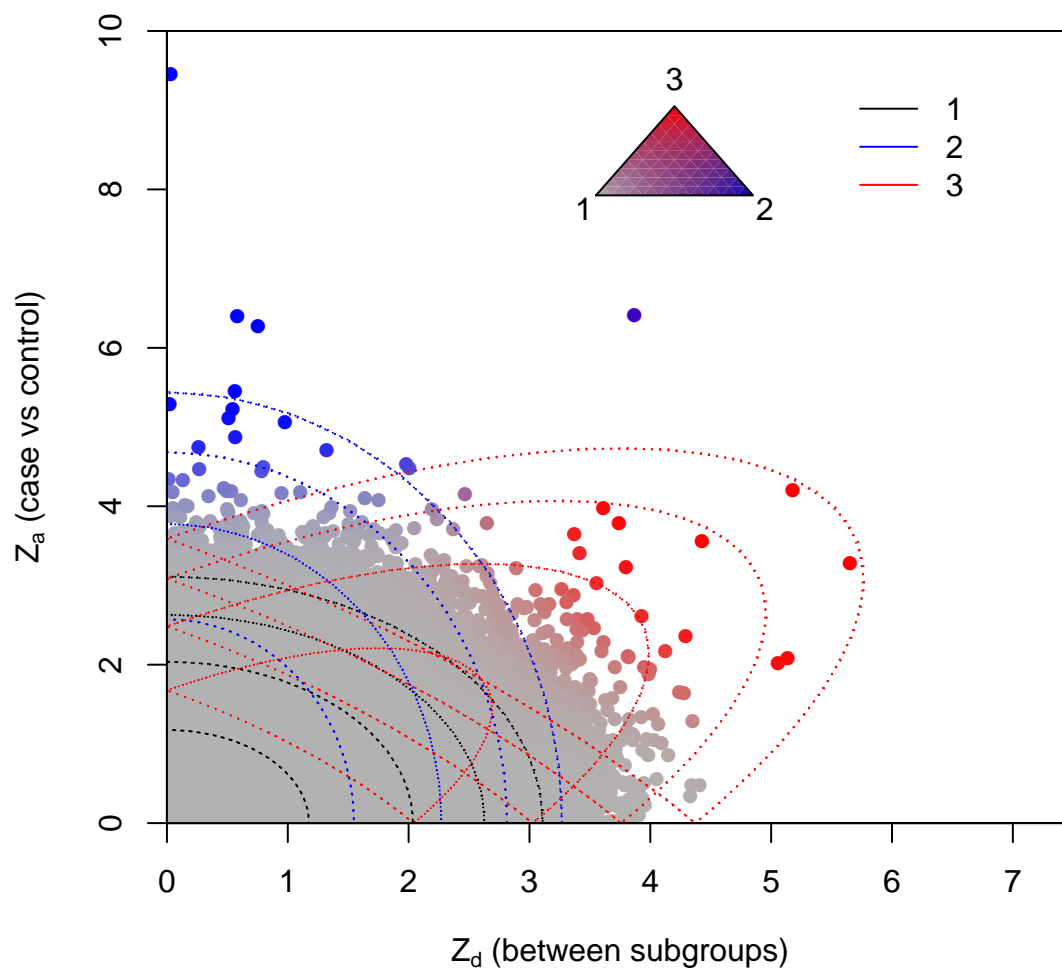


Figure 4: Observed absolute  $Z_a$  and  $Z_d$  for T1D/RA. Colourings correspond to posterior probability of category membership under full model (see triangle): grey - category 1, blue - category 2, red -category 3. Contours of the component Gaussians in the fitted full model (at are shown by dotted lines).



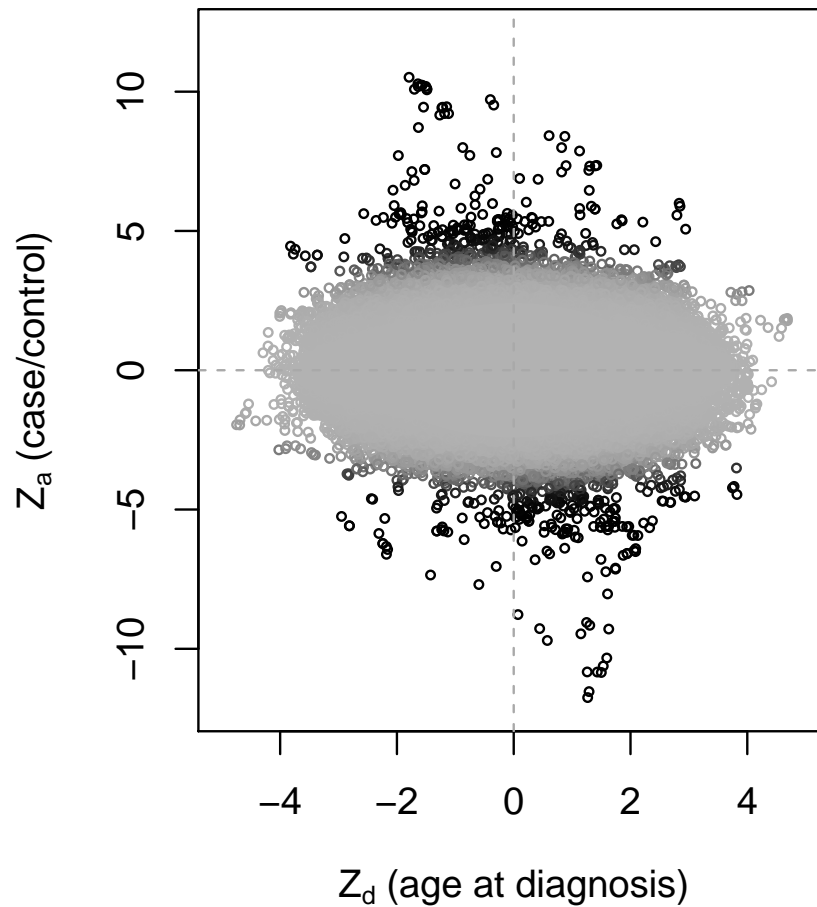


Figure 5:  $Z_a$  and  $Z_d$  scores for age at diagnosis in T1D, excluding MHC region. Colour corresponds to posterior probability of category 2 membership in null model (since categories in full model are assigned on the basis of correlation), with black representing a high probability. Data show a negatively correlated ( $p = 8.7 \times 10^{-5}$  with MHC included,  $p = 0.002$  with MHC removed)

## References

- [1] Li L, Cheng WY, Glicksberg BS, Gottesman O, Tamler R, et al. (2015) Identification of type 2 diabetes subgroups through topological analysis of patient similarity. *Science translational medicine* 7: 311ra174–311ra174.
- [2] Morris AP, Lindgren CM, Zeggini E, Timpson NJ, Frayling TM, et al. (2009) A powerful approach to sub-phenotype analysis in population-based genetic association studies. *Genetic Epidemiology* 34: 335-343.
- [3] Plagnol V, Howson JMM, Smyth DJ, Walker N, Hafler JP, et al. (2011) Genome-wide association analysis of autoantibody positivity in type 1 diabetes cases. *PLOS Genetics* 7.
- [4] Speed D, Hemani G, Johnson MR, Balding DJ (2012) Improved heritability estimation from genome-wide SNPs. *American Journal of Human Genetics* 91: 1011-1021.
- [5] Chen H, Chen J, Kalbfleisch JD (2001) A modified likelihood ratio test for homogeneity in finite mixture models. *Journal of the Royal Statistical Society, series B (methodological)* 63: 19-29.
- [6] Andreassen OA, Thompson WK, Schork AJ, Ripke S, Mattingsdal M, et al. (2013) Improved detection of common variants associated with schizophrenia and bipolar disorder using pleiotropy-informed conditional false discovery rate. *PLOS Genetics* 9(4).
- [7] Liley J, Wallace C (2015) A pleiotropy-informed bayesian false discovery rate adapted to a shared control design finds new disease associations from gwas summary statistics. *PLOS Genetics* .

- [8] Leslie S, Winney B, Hellenthal G, Davison D, Boumertit A, et al. (2015) The fine-scale genetic structure of the british population. *Nature* 519: 309-314.
- [9] The Wellcome trust case control consortium (2007) Genome-wide association study of 14000 cases of seven common diseases and 3000 shared controls. *Nature* 447: 661-678.
- [10] Fortune MD, Guo H, Burren O, Schofield E, Walker NM, et al. (2015) Statistical colocalization of genetic risk variants for related autoimmune diseases in the context of common controls. *Nature Genetics* 47: 839-846.
- [11] Lo PR, Tucker G, Bulik-Sullivan BK, Vilhjalmsson BJ, Finucane HK, et al. (2015) Efficient bayesian mixed model analysis increases association power in large cohorts. *Nature Genetics* 47: 284-90.
- [12] Cooper JD, Simmonds MJ, Walker NM, Burren O, Brand OJ, et al. (2012) Seven newly identified loci for autoimmune thyroid disease. *Human Molecular Genetics* 21: 5202-5208.
- [13] Bulik-Sullivan B, Finucane HK, Anttila V, Gusev A, Day FR, et al. (2015) An atlas of genetic correlations across human diseases and traits. *bioRxiv* .
- [14] Lee SH, Yang J, Goddard ME, Visscher PM, Wray NR (2012) Estimation of pleiotropy between complex diseases using single-nucleotide polymorphism-derived genomic relationships and restricted maximum likelihood. *Bioinformatics* 28: 2540–2542.
- [15] Traylor M, Bevan S, Rothwell PM, Sudlow C, 2 WTCCC, et al. (2013) Using phenotypic heterogeneity to increase the power of genome-wide association studies: Application to age at onset of ischaemic stroke subphenotypes. *Genetic Epidemiology* 37: 495-503.

- [16] Wen Y, Lu Q (2013) A multiclass likelihood ratio approach for genetic risk prediction allowing for phenotypic heterogeneity. *Genetic epidemiology* 37: 715–725.
- [17] Howson JMM, Walker NM, Smyth DJ, Todd JA (2009) Analysis of 19 genes for association with type 1 diabetes in the type 1 diabetes genetics consortium families. *Genes and Immunity* 10: S74-S84.
- [18] Howson JM, Rosinger S, Smyth DJ, Boehm BO, Todd JA, et al. (2011) Genetic analysis of adult-onset autoimmune diabetes. *Diabetes* 60: 2645–2653.
- [19] Howson JM, Cooper JD, Smyth DJ, Walker NM, Stevens H, et al. (2012) Evidence of gene-gene interaction and age-at-diagnosis effects in type 1 diabetes. *Diabetes* 61: 3012–3017.
- [20] Hyttinen V, Kaprio J, Kinnunen L, Koskenvuo M, Tuomilehto J (2003) Genetic liability of type 1 diabetes and the onset age among 22, 650 young finnish twin pairs in a nationwide follow up study. *Diabetes* 52: 1052-1055.
- [21] Chatterjee N, Carroll RJ (2005) Semiparametric maximum likelihood estimation exploiting gene-environment independence in case-control studies. *Biometrika* 92: 399–418.
- [22] Self SG, Liang KY (1987) Asymptotic properties of maximum likelihood estimators and likelihood ratio tests under nonstandard conditions. *Journal of the American Statistical Association* 82: 605-610.
- [23] Cortes A, Brown MA (2011) Promise and pitfalls of the ImmunoChip. *Arthritis Research and Therapy* 13.

- [24] Dempster AP, Laird NM, Rubin DB (1977) Maximum likelihood from incomplete data via the em algorithm. *Journal of the Royal Statistical Society, series B (methodological)* 39: 1-38.
- [25] Hastie T, Tibshirani R, Friedman J (2001) *The Elements of Statistical Learning*. Springer Series in Statistics. Springer.
- [26] Barrett JC, Clayton DG, Concannon P, Akolkar B, Cooper JD, et al. (2009) Genome-wide association study and meta-analysis find that over 40 loci affect risk of type 1 diabetes. *Nature genetics* 41: 703–707.
- [27] Anderson CA, Pettersson FH, Clarke GM, Cardon LR, Morris AP, et al. (2010) Data quality control in genetic case-control association studies. *Nature protocols* 5: 1564-1573.
- [28] Devlin B, Roeder K, Wasserman L (2001) Genomic control, a new approach to genetic-based association studies. *Theoretical Population Biology* 60: 155-166.