

METHODOLOGY ARTICLE

Open Access

Combined genotype and haplotype tests for region-based association studies

Sergii Zakharov^{1,2*}, Tien Yin Wong^{3,4}, Tin Aung^{3,4}, Eranga Nishanthie Vithana³, Chiea Chuen Khor^{1,2}, Agus Salim² and Anbupalam Thalamuthu^{1,5*}

Abstract

Background: Although single-SNP analysis has proven to be useful in identifying many disease-associated loci, region-based analysis has several advantages. Empirically, it has been shown that region-based genotype and haplotype approaches may possess much higher power than single-SNP statistical tests. Both high quality haplotypes and genotypes may be available for analysis given the development of next generation sequencing technologies and haplotype assembly algorithms.

Results: As generally it is unknown whether genotypes or haplotypes are more relevant for identifying an association, we propose to use both of them with the purpose of preserving high power under both genotype and haplotype disease scenarios. We suggest two approaches for a combined association test and investigate the performance of these two approaches based on a theoretical model, population genetics simulations and analysis of a real data set.

Conclusions: Based on a theoretical model, population genetics simulations and analysis of a central corneal thickness (CCT) Genome Wide Association Study (GWAS) data set we have shown that combined genotype and haplotype approach has a high potential utility for applications in association studies.

Keywords: Genotype-based tests, Haplotype-based tests, Association analysis, Test statistic combination

Background

The development of genotyping and sequencing technologies has enabled scientists to investigate the impact of genomic loci on complex disorders and traits. Indeed, genome-wide association studies (GWAS) and sequencing studies have identified many common single-nucleotide polymorphisms (SNPs) (for GWAS publication list, see <http://www.genome.gov/gwastudies/>) and rare variations [1-4] associated with common diseases. Although single-SNP analysis has proven to be useful in discovering many disease-associated loci, this strategy may be limited due to very stringent significance threshold and poor reproducibility [5]. Region-based association studies have the advantages of less stringent significance level and potentially higher power if multiple associated variants are found within a region. Indeed, several empirical studies

have demonstrated the superiority of genotype gene-based association analysis over single-SNP strategy [6,7]. Also, there is some theoretical [8,9] and empirical evidence that haplotype-based tests may possess higher power than SNP-based tests. When intending to use haplotypes in an association study, one faces a problem of phase inference. While several statistical algorithms have been developed to infer unknown haplotypes from genotype data [10-12], the improvements of sequencing technologies will enable researchers to assemble haplotypes from sequencing data with very high accuracy (for examples of existing assembly algorithms, see Bansal et al. [13], Bansal et al. [14], and Schatz et al. [15]). This opens up the opportunity to use high-quality haplotypes and genotypes in sequencing association studies.

Numerous studies have reported cases when haplotype-based analysis resulted in detection of an association, while SNP-based analysis either did not yield any significant results or yielded much higher p-values [16-19]. A haplotype-based test may be more powerful than a genotype-based test if haplotypes tag a true causal variant

* Correspondence: zakharovs99@gis.a-star.edu.sg; a.thalamuthu@unsw.edu.au

¹Human Genetics, Genome Institute of Singapore, Singapore, Singapore

⁵Centre for Healthy Brain Ageing (CHeBA), School of Psychiatry, University of New South Wales, Sydney, Australia

Full list of author information is available at the end of the article

better (although the imputation of untyped SNPs using publicly available reference panels may also be a powerful strategy), or if a SNP-SNP interaction is present within a region. In general, it is unknown whether haplotype- or genotype-based tests are more relevant for identifying an association of a genomic region with a phenotype. In this article we propose two statistical tests that explicitly combine both genotype and haplotype information for the purpose of preserving high power under both genotype and haplotype disease scenarios. We investigate two methods based on a combination of p-values from genotype- and haplotype-based association tests. The first method is a minimum of p-values (MinP-val), and the other is a sum test statistic based on inverse standard normal transformation of two p-values (SumP-val). Based on simulations, theoretical power calculations and application to a GWAS data set, we have highlighted the merits and the drawbacks of genotype- and haplotype-based tests, and those of our combined approaches. The major conclusions from our work are as follows:

1. Combination of haplotype- and genotype-based test statistics preserves power for both genotype and haplotype disease models;
2. In some of the considered scenarios, the performance of the MinP-val approach is comparable to those of the SumP-val method;
3. MinP-val is much more robust than SumP-val when one of the underlying tests has low power.

Methods

Genotype- and haplotype-based tests

Let us assume that we are interested in testing the joint association of all the variants within a genomic region with either a dichotomous phenotype or quantitative trait. Next, assume we have chosen the two statistical tests for a region-based association analysis: one genotype- and one haplotype-based test. For haplotype-based tests haplotypes can be inferred from genotypes [10-12] or assembled from sequencing data [13,15,20]. Several conventional genotype-based methods [21-23] are applicable for common variants testing, whereas for sequencing data numerous recently-developed rare variants approaches are available [24-28]. Haplotype-based methodologies have also been extensively published elsewhere [29-31], including rare haplotype tests [32-34].

The combined approaches

Let us denote p-values from a genotype- and a haplotype-based tests as p_1 and p_2 respectively. Our first approach is SumP-val [35]. Let us consider the inverse standard normal transformation of both p-values and which are distributed as standard normal random variables under the null hypothesis. Here, we assume that y_1, y_2 is bivariate

normal. The SumP-val test statistic is $P_{\text{sum}} = y_1 + y_2$. Under the null hypothesis, it is distributed as a normal random variable with zero mean and variance $\text{Var}(y_1 - y_2) = \text{Var}(y_1) + 2\text{Cor}(y_1, y_2) - \text{Var}(y_2) = 2 + 2p$, where p is a correlation coefficient between y_1 and y_2 , since two statistical tests for the same genomic region may not be independent. The correlation coefficient p may be estimated via permutation procedure. The rejection region is large values of the test statistic, which is equivalent to low values for p_1 and/or p_2 . The theoretical p-value for SumP-val test is calculated as $P(P_{\text{sum}} > x)$, where $\Phi(a, b, c)$ is a value of normal cumulative distribution function with mean b and variance c taken at the point a .

Our second approach (MinP-val) is to utilize the minimum of the two p-values as a test statistic, namely, $\min\{p_1, p_2\}$ [36]. Let us represent a test statistic as $\min\{p_1, p_2\} = \min\{1 - \Phi(y_1), 1 - \Phi(y_2)\}$, where y_1 and y_2 defined above are distributed as standard normal random variables under the null. Thus, the theoretical cumulative distribution function of MinP-val test statistic under the null hypothesis can be calculated as follows:

$$\begin{aligned} P(\min\{1 - \Phi(y_1), 1 - \Phi(y_2)\} < x) &= P(1 - \max\{\Phi(y_1), \Phi(y_2)\} < x) = \\ &= P((1-x) < \max\{\Phi(y_1), \Phi(y_2)\}) = 1 - P(1-x > \Phi(y_1), 1-x > \Phi(y_2)) = \\ &= 1 - P(\Phi^{-1}(1-x) > y_1, \Phi^{-1}(1-x) > y_2), \end{aligned} \quad (1)$$

where $0 < x < 1$. Given the rejection region is small values of the test statistic, theoretical p-value for MinP-val test is straightforward to compute using (1).

Theoretical power model

Within our theoretical framework the following model is adopted: the two test statistics S_g and S_h of the underlying genotype- and haplotype-based tests, respectively, are assumed to asymptotically follow central chi-squared distribution X_1^2 with 1 degree of freedom under the null hypothesis, and non-central chi-squared X_{1a}^2 and X_{1b}^2 with NCPs a and b , respectively, under the alternative hypothesis. One of the examples of the test which results in such null and alternative distributions is Rao's score test on, for example, genotype or haplotype scores described in Additional file 1. Since the two tests are applied to the same data, the chi-squared test statistics are likely to be positively correlated. The correlation between the two test statistics may vary from very low to high. For example, if within a region there are few SNPs in very high LD then we would expect the correlation between the tests to be high. Alternatively, we would expect the correlation to be low when variants within a region are independent. The correlation is modeled via underlying multivariate normal distribution, namely, to simulate the test statistics $S_g \sim X_{1a}^2$ and $S_h \sim X_{1b}^2$ a bivariate normal random vector $y = (y_1, y_2)$ with mean (\sqrt{a}, \sqrt{b}) ,

unit variances and correlation coefficient $p > 0$ is generated, and the squares of the coordinates are taken as the proxy for the test statistics: $S_g = y_1^2$, $S_h = y_2^2$. To estimate the power of MinP-val and SumP-val tests we simulated 500,000 independent pairs (S_g, S_h) under the alternative hypothesis, calculated the test statistics for the combined approaches, and noted the share of statistically significant pairs. This procedure was done for every theoretical scenario (see “Results” section).

Population genetics simulations

King et al. [37] provided the SFS_CODE (<http://sfscode.sourceforge.net>) implementation of population genetics simulation for ANGPTL4 gene exons (<http://home.uchicago.edu/~crk8e>). The authors assumed the demographic and distribution fitness effect parameters from Boyko et al. [38] and Gutenkunst et al. [39]. Using the SFS_CODE program 1000 haplotype pools each containing 20000 sampled “individuals” (40000 chromosomes) from a European population were generated. A data replicate was created from each haplotype pool by iterative random sampling of two haplotypes (thus, defining the genotype of an “individual”) and assigning a dichotomous phenotype conditional on a multi-site genotype or a pair of haplotypes. Each data replicate contains 500 cases and 500 controls. Let us assume that there are L variants within the genomic region of interest, and the genotype of “an individual” $\{g_1, \dots, g_L\}$ is constructed from the sampled haplotypes. To describe the genotype-based disease model let us, without loss of generality, denote the genotypes at rare (MAF < 1%) causal SNPs as $\{g_1, \dots, g_c\}$, causal common SNP g_{c-1} (if present depending on a model), and other SNPs as $\{g_{c-2}, \dots, g_L\}$. Let us also define the assigned odds ratios of causal variants $\{b_1, \dots, b_{c-1}\}$. The probability of a disease $P(A)$ for a genotype-based scenario is calculated from the following:

$$\log\left(\frac{P(A)}{1-P(A)}\right) = \log\left(\frac{0.01}{1-0.01}\right) + \sum_{l=1}^{c-1} g_l \log(b_l). \quad (2)$$

For the haplotype-based scenarios let us consider the two sampled haplotypes $\{h_1, h_2\}$. Also, denote H_r and H_c as the sets of rare (frequency in a haplotype pool < 1%) and common causal haplotypes, respectively (depending on the disease model H_c may be empty). The probability of a disease is defined as:

$$\log\left(\frac{P(A)}{1-P(A)}\right) = \log\left(\frac{0.01}{1-0.01}\right) + \sum_{i=1}^2 (I\{h_i \in H_r\} \log(d_r) + I\{h_i \in H_c\} \log(d_c)), \quad (3)$$

where $I\{A\}$ is an indicator of an event A , and $\{d_r, d_c\}$ are the odds ratios for causal rare and common haplotypes,

respectively. For our simulations, we considered three phenotype models: “Rare” (only rare variants or haplotypes are risk-contributing), “Common” (only common variants or haplotypes), and “Both” (both types of variants or haplotypes). Following the scenarios of exome-scale simulations of Wu et al. [40], we have assigned 50%, 20% and 10% of the observed rare variants (haplotypes) to be causal. Additionally, we chose one common causal SNP (haplotype) for “Both” and “Common” models. The odds ratios $\{b_l\}_{l=1}^{c-1}$ in (2) and $\{d_r, d_c\}$ in (3) were assigned as follows:

- for the “Rare” model: $b_l = d_r = 4, l = 1, \dots, c$ and $b_{c-1} = d_c = 0$
- for the “Both” model: $b_l = d_r = 3, l = 1, \dots, c$ and $b_{c-1} = d_c = 1.2$
- for the “Common” model: $b_l = d_r = 1.5, l = 1, \dots, c$ and $b_{c-1} = d_c = 2$

The average number of variants across data replicates is shown in Table 1.

Real data analysis

For the purpose of demonstrating the performance of the described methodologies we conducted a gene-based analysis of the central corneal thickness (CCT) GWAS data sets described in Vithana et al. [41]. Briefly, the Singapore Indian Eye Study (SINDI), which is part of the Singapore Indian Chinese Cohort Eye Study (SICC) [42], consists of 2538 Indian subjects aged 40 and above, and the Singapore Malay Eye Study (SiMES) [43-45] is a genome-wide association study of CCT phenotype which contains 2542 Malay subjects aged 40 and above. Both SiMES and SICC adhered to the Declaration of Helsinki. Ethics approval for the both studies was obtained from the Singapore Eye Research Institute Institutional Review Board [41]. The combined data set consists of 5080 individuals genotyped at 552318 SNPs after quality control. In total, 5049 individuals were analyzed after excluding those with missing phenotype. Also, we attempted to replicate all the genome-wide significant regions using Chinese samples

Table 1 The average number of variants within a region across 1000 data replicates in population genetics simulations

Phenotype model	Proportion of causal variants/haplotypes		
	50%	20%	10%
Haplotype common	32.4	31.6	31.6
Haplotype both	35.5	33.2	32.6
Haplotype rare	37.2	34.1	33.0
Genotype common	33.1	32.4	32.2
Genotype both	36.3	33.7	32.9
Genotype rare	37.6	34.4	33.0

from the SICC. This data set contains 2837 samples with non-missing phenotype and covariates (age and gender). SNPs were mapped to genes based on the method outlined by Zhao et al. [46]. Briefly, information on gene identifiers (IDs), names, start and end positions on a chromosome were downloaded from the NCBI Genome database (<http://www.ncbi.nlm.nih.gov/Genomes>). Gene regions included 10 kb upstream and downstream. Hierarchical mapping scheme (coding > intronic > 5'UTR > 3'UTR) was used if a variant was within 10 kb of multiple genes. The remaining inter-gene variants between two genes were grouped together. Haplotype inference was performed using the software Beagle [10] with reference panel from 1000Genomes Project (<http://www.1000genomes.org/>). In our analysis we adjusted for age, gender and the first ten principal components from Eigenstrat [47].

Statistical tests for population genetics simulations and real data application

Sequence Kernel Association Test (SKAT), introduced by Wu et al. [40], is a variance component score test derived from a semi-parametric regression model. It was initially proposed to test the association of phenotype with multi-site genotype; however, we also used SKAT to test an association of phenotype with haplotypes as described below. To show the consistency of empirical results we applied the same pair of underlying tests, namely, genotype SKAT and haplotype SKAT with linear kernel and uniform weights, to both population genetics simulations and real data. For genotype SKAT all rare variants (MAF < 1% in the sample) within a region were collapsed according to the method described by Thalamuthu et al. [48]. Briefly, the collapsed super-variant is the sum of minor alleles across rare variants within a region; if this sum is greater than 2 then the value of 2 is assigned. We did not apply weighting as described by Wu et al. [40] because that would substantially decrease power to identify an association with common SNPs. Haplotype-based SKAT is SKAT applied to a haplotype regression matrix R , which is constructed similar to those used by Zaykin et al. [31]. First, we pooled all rare haplotypes into one haplotype group, whereas each of the common haplotypes formed a separate haplotype group. Let us define the following notations: n is a number of individuals; $\{H_1, \dots, H_M\}$ is the haplotype groups with H_M being the most common group; $R = \{R_{ij}\}_{ij=1}^{n, M-1}$ is a haplotype regression matrix; and $\{h_{i,1}; h_{i,2}\}$ is a pair of haplotypes for i th individual. The haplotype matrix $\{R_{ij}\}$ is constructed as follows:

$$R_{ij} = I\{h_{i,1} \in H_j\} + I\{h_{i,2} \in H_j\}, i = 1, \dots, n; j = 1, \dots, M-1, \quad (4)$$

where $I\{A\}$ is an indicator of an event A . If there were no common haplotypes within a region, we formed three groups of haplotype: those with a frequency less than 0.05%,

those in between 0.05% and 0.1%, and those with a frequency greater than 0.1%. For both tests we used the R (<http://www.r-project.org/>) package SKAT (<http://www.hsph.harvard.edu/~xlin/software.html>). For population genetics simulations p-values for all the tests were estimated using 1000 permutations. In real data analysis for the underlying tests we used theoretical p-values as we believe that they reasonably approximate empirical p-values given large sample size and normally-distributed quantitative trait [41]. Then we tested an assumption of bivariate normality by applying the Shapiro-Wilk test (R package "mvnormtest" <http://cran.r-project.org/web/packages/mvnormtest/>). If the normality test was not significant on the genome-wide level, we used theoretical p-values for both SumP-val and MinP-val; otherwise we used permutations. The permutation procedure and estimation of correlation coefficient are described in the next section.

Permutation procedure and estimation of correlation coefficient

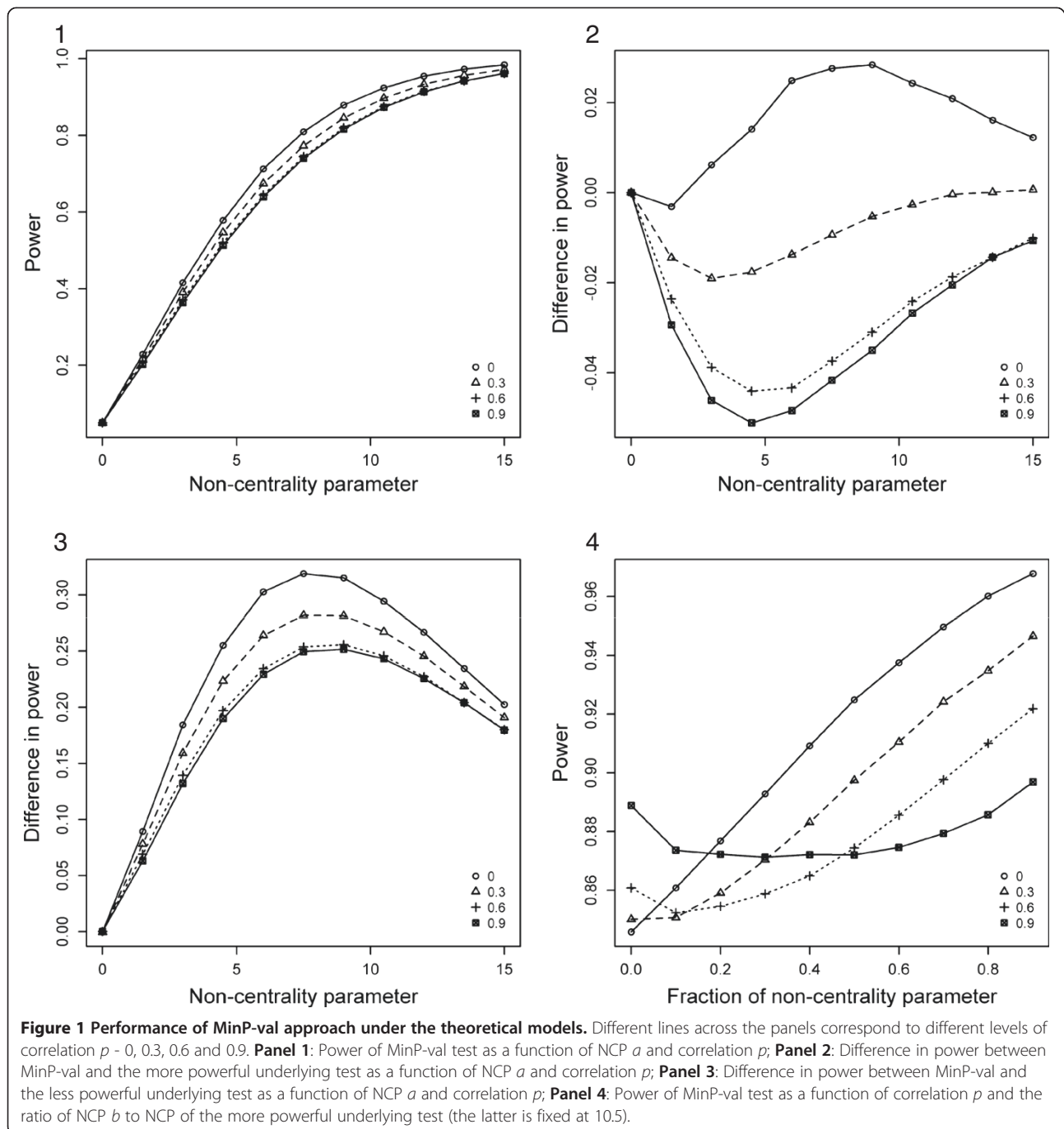
To calculate theoretical p-values for the proposed methods we estimated the correlation coefficient ρ using 500 permutations. The difficulty in applying permutations lies in the fact that the permutation procedure should preserve the relationship between all the covariates, and also between phenotype and covariates, but disrupt the relationship between phenotype and genotype. Several techniques have been developed for conducting permutation tests of partial coefficients in a multiple regression model [49-51]. Among them the permutation of residuals under the reduced model [49] was shown to preserve correct type-1 error for t-test [52] and was previously applied to microarray data analysis [53]. As the SKAT test can be obtained from a semi-parametric regression model [40], let us consider the following genotype and haplotype regression models: $Y = a_1 - f_1(P) - Cc - \varepsilon$ and $Y = a_2 - f_2(R) - Cc - \varepsilon$, where P is $n \times L$ collapsed genotype matrix, n is the sample size, L is the number of common SNPs within a region plus one for collapsed rare variants super-locus, Y is $n \times 1$ vector of quantitative phenotype (CCT), C is $n \times 12$ matrix of covariates which include age, gender and the first ten genotype principal components obtained from Eigenstrat [47], R is haplotype regression matrix, and f_1 and f_2 are unknown functions. To obtain the permutation values for the test statistics the reduced model $Y = a - Cc - \varepsilon$ is fitted, and a, c, ε are the estimated constant coefficient, regression coefficients and residuals, respectively. Next, the residuals ε are permuted to obtain ε , and $Y = a + Cc - \varepsilon$. The permuted statistic values for both genotype and haplotype SKAT tests are calculated as respective SKAT statistics from semi-parametric models $Y = a_1 - f_1(P) - Cc - \varepsilon$ and $Y = a_2 - f_2(R) - Cc - \varepsilon$. Each p-value obtained from permutations was transformed using the inverse standard normal transformation, and the value of ρ was estimated by a Pearson correlation coefficient.

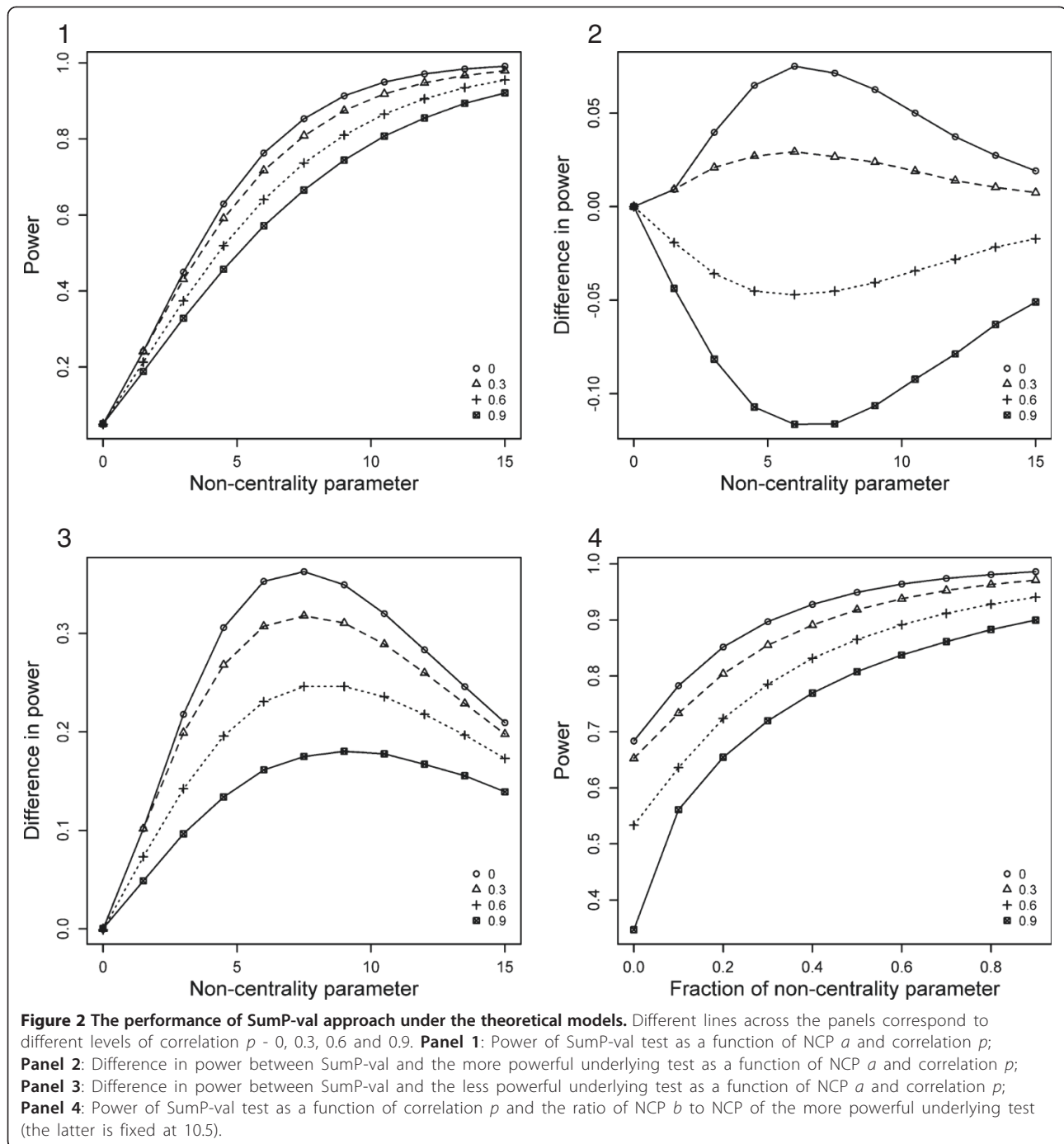
Results

Theoretical power results

Depending on the disease model, one of the underlying tests (genotype- or haplotype-based) is expected to be more powerful than the other underlying test. So, we assume that under the alternative hypothesis the non-centrality parameter (NCP) of the more powerful underlying test is a , and the NCP of the less powerful underlying test is $b = a/2$. Figures 1 and 2 (Panel 1) show the power of MinP-val and SumP-val strategies as a function of correlation

coefficient p and NCP a at the fixed type-1 error of 0.05. As can be seen in Panel 1, power in general decreases slightly with increasing correlation. Panel 2 depicts the difference in power between the combined approaches and the more powerful underlying test. It is notable that both MinP-val and SumP-val achieved greater power than the more powerful underlying test for lower correlation. Also, MinP-val approach lost a maximum of 5% power for high correlation and gained a maximum of 2% for low





correlation, whereas for SumP-val these values are more than 5% in both cases.

In Panel 3, where the difference in power between the combined approaches and the less powerful underlying test is shown, it can be seen that both MinP-val and SumP-val are consistently better than the less powerful test. This suggests that combination of statistical tests may prove beneficial when the underlying disease model

is unknown. To investigate the impact of change of NCP b on the performance of the proposed approaches we fixed NCP a to be equal to 10.5 (corresponding to 90% power of a chi-squared test with χ^2_{1a} distribution under the alternative hypothesis, the type-1 error is 0.05). Panel 4 of Figures 1 and 2 depicts the power of MinP-val and SumP-val as a function of correlation and a “fraction of NCP” – the ratio of b to 10.5. As can be seen in Panel 4,

MinP-val test achieved higher power than SumP-val in the majority of scenarios. It is notable that SumP-val lost much power when the value of b is low. Hence, MinP-val approach is more robust with respect to underperformance of one of the underlying tests.

Population genetics simulation results

Panel 4 of Figure 3 shows the empirical type-1 error estimate for the theoretical level of 0.05 for all the tests. The estimate of type-1 error is distributed as a binomial random variable with 1000 trials and the probability of success 0.05 under the hypothesis of no inflation. The one-sided 99% quantile of the described distribution is 0.067. As can be seen, in our simulations the type-1 error was well controlled for all the tests.

Panels 1–3 of Figure 3 depict the results of population genetics simulations analysis for all the phenotype models with 50%, 20% and 10% or rare causal variants/haplotypes, respectively, at the fixed 5% type-1 error. For all the tests 1000 permutations were performed to estimate p-values. Haplotypes were assumed to be known without ambiguity. Under the genotype-based disease scenarios, genotype SKAT is expected to be more powerful than haplotype SKAT, and vice versa under the haplotype-based scenarios. However, genotype SKAT was less powerful for many genotype-based phenotype models. A possible explanation of this observation is that when rare variants are strongly associated with phenotype, for some statistical tests pooling of rare haplotypes may be a better strategy than pooling of rare variants. Also, it should be noted that with the decrease in the percentage of causal rare variants/haplotypes, the power for “Rare” and “Both” phenotype models decrease substantially since for these models rare variants/haplotypes are the major carriers of an association signal. For “Common” phenotype model one common variant/haplotype has a significant impact on phenotype; so, the decrease in power with the lower proportion of causal rare variants/haplotypes is not as high as for other phenotype models.

As can be seen from the Panels 1–3 of Figure 3, for all the phenotype disease models, when both underlying tests were almost equally powerful (e.g. Panel 1 haplotype disease scenario “Common” model, and genotype disease scenario “Both” and “Common” models), the power of both MinP-val and SumP-val were on the same level or even higher than those of the underlying tests. However, when genotype-based SKAT significantly underperformed haplotype-based SKAT (e.g. Panel 1 haplotype disease scenario “Rare” and “Both” models), MinP-val approach showed slightly lower power than the more powerful underlying test and greater power compared with SumP-val approach. The maximum power loss of SumP-val and MinP-val compared with the more powerful underlying test across all phenotype models was

6.3% and 3.8% respectively (haplotype disease scenario “Both” model). These results are consistent with those obtained from the theoretical power considerations, and illustrate the great potential of the proposed methods in their application to real association studies.

To examine the effect of phasing on our results we repeated the analysis using the most probable haplotypes inferred by Beagle [10]. The reference panel consisted of 1094 simulated individuals to mimic the size of the publicly available reference panel from the 1000 Genomes Project (www.1000genomes.org). The results of this analysis were very similar to those described above (data not shown). In addition, we applied the proposed methods with a different pair of underlying tests. The results are similar to those described above. For more details, see Additional files 1 and 2.

Application to central corneal thickness GWAS data set

A total of 552318 SNPs were mapped using the hierarchical mapping algorithm described in the “Methods” section. As a result, we obtained 36146 genes and between-gene blocks. Regions that reached genome-wide significance ($1.38E-6$ after Bonferroni correction) for at least one of the four applied tests are presented in Table 2. We identified all of the significant regions reported by Vithana et al. [41]: COL8A2 gene, an interval between the genes RXRA and COL5A1 and a region near ZNF469 gene. As can be seen from Table 2, genotype-based SKAT and MinP-val tests achieved genome-wide significance for all the five regions listed, whereas SumP-val failed to reach genome-wide significance for both RXRA-COL5A1 and C7orf42 regions. This highlights that MinP-val approach performed better than the SumP-val approach. Haplotype-based SKAT failed to identify any association signal for four out of the five regions. From our results it is clear that for this data set genotypes were more relevant for identifying associated regions. Judging from the performance of the Haplotype SKAT there is no evidence of association of haplotypes with a phenotype except for COL8A2 gene. It is also of interest to compare our results to those reported by Vithana et al. [41] for single-SNP analysis. For the two out of three regions genotype-based SKAT and MinP-val methods outperformed the single-SNP analysis and yielded lower p-values, which may be explained by the utilization of linkage disequilibrium (LD) within a region to boost power. Given the sensitivity of SumP-val approach to underperformance of one of the underlying tests, this method showed higher p-values. To justify our assumption of bivariate normality for calculation of p-values for the proposed tests, we used the Shapiro-Wilk test. The corresponding p-values for the significant regions are presented in the (Additional file 1: Table S1) in the first row. All the p-values, except those for COL8A2-TRAPPC3 region, are non-significant at the genome-wide level, which suggests

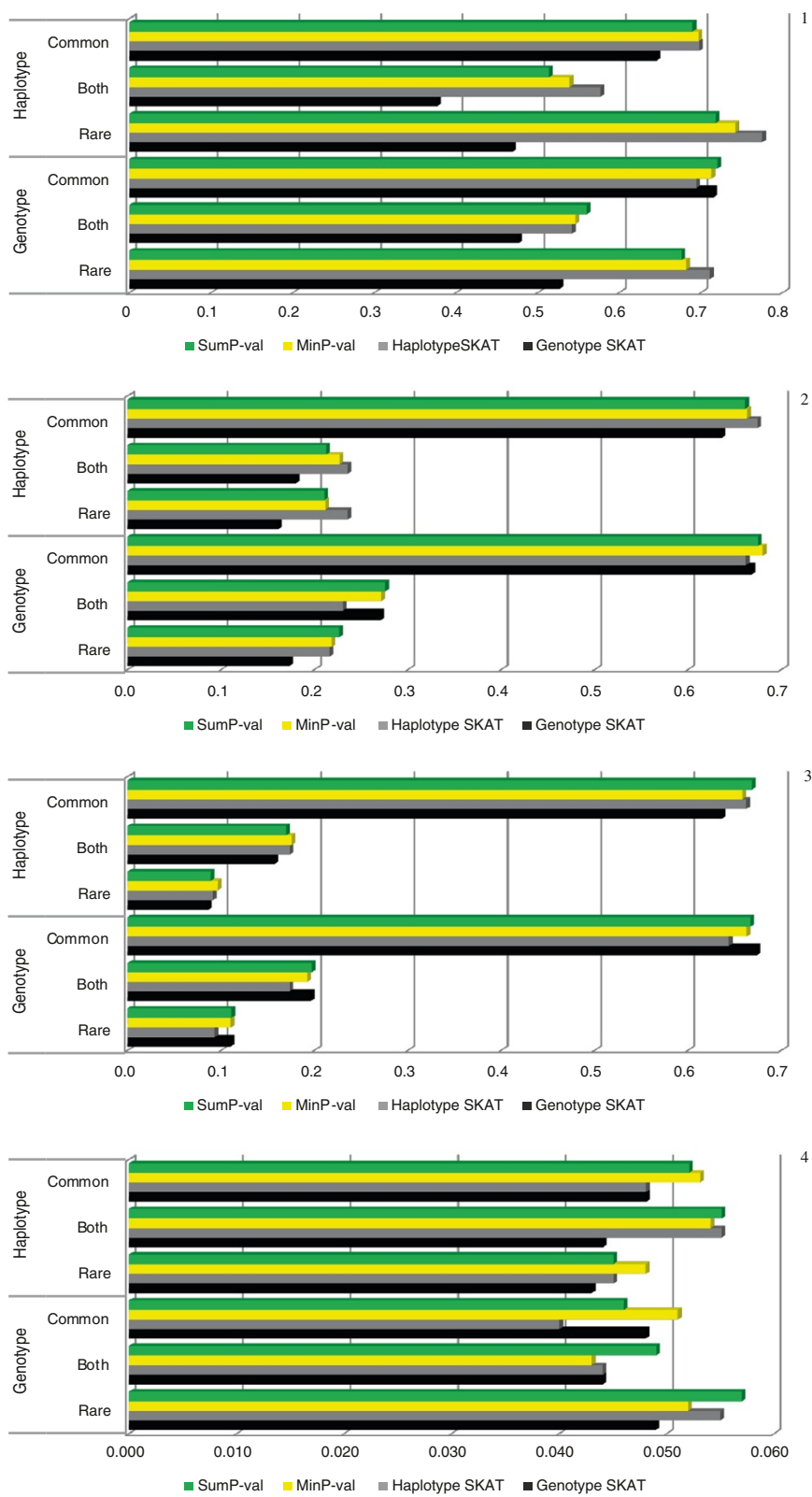


Figure 3 (See legend on next page.)

(See figure on previous page.)

Figure 3 Power comparison of genotype-based SKAT, haplotype-based SKAT, MinP-val and SumP-val tests for population genetics simulations, and an estimate of empirical type-1 error. In each panel the top three disease models correspond to the haplotype-based disease scenario, whereas the lower three correspond to the genotype-based scenario. Disease models "Rare", "Both" and "Common" are described in the section "Population genetics simulation". Type-1 error is set to 5%. **Panel 1:** 50% of rare variants/haplotypes were assumed to be causal; **Panel 2:** 20% of rare variants/haplotypes were assumed to be causal; **Panel 3:** 10% of rare variants/haplotypes were assumed to be causal; **Panel 4:** empirical type-1 error estimate for simulations under the null hypothesis.

there is no evidence against the bivariate normality assumption. The Shapiro-Wilk test of COL8A2-TRAPPC3 region yielded a marginally significant p-value on the genome-wide level. Hence, the p-value for this region in the Table 2 is based on permutations.

Table 3 shows the results of the replication analysis. For the region COL8A2-TRAPPC3 the reported replication p-value is based on permutations, whereas for other regions there was no evidence against bivariate normality assumption (Additional file 1: Table S1, second row). As can be seen from Table 3, only RXRA-COL5A1 region was significant after Bonferroni correction for all the tests except for the haplotype-based SKAT. Our replication results are consistent with those of Cornes et al. [54] who found strong evidence of association of multiple SNPs within the RXRA-COL5A1 region, and marginal significance of COL8A2 SNP rs96067. It is worth noting that in our analysis C7orf42 gene, which was not identified by Vithana et al. [41], reached genome-wide significance in SiMES + SINDI data set and had moderate p-value in the replication dataset. Cornes et al. [54] found this gene to be significant in a meta-analysis of SiMES, SINDI, 1883 samples from the Singapore Chinese Eye Study and 798 samples from the Beijing Eye Study [55]. The role of C7orf42 gene in central corneal thickness (CCT) phenotype requires further investigation. The results of our analysis suggest that RXRA-COL5A1 region may have an impact on CCT phenotype.

In addition to the gene-based analysis, we tried to replicate the four genome-wide significant SNPs found by Vithana et al. [41] in our Chinese samples using single-SNP analysis. Having tested an association of these SNPs with CCT trait using trend test within a

linear additive model adjusting for age, gender and the first ten principal components, we found that none of the SNPs was significant on the corrected type-1 error rate $0.0125 = 0.05/4$. This result suggests that gene-based replication may be a more powerful strategy than single-SNP replication.

In addition to the main genome-wide analysis of SiMES + SINDI data set, we applied the proposed methods with a different pair of underlying tests to the three regions reported by Vithana et al. [41]. Both MinP-val and SumP-val identified the three regions on genome-wide significance level. This result suggests that our combined approaches work as well with other underlying tests (for more details, see Additional file 1).

Discussion

When the underlying disease model is unknown, combining statistical tests tailored for different disease scenarios may be a much better strategy than application of a statistical test designed for one specific disease model. In this article we have described the two approaches of combining genotype- and haplotype-based statistical tests. The results of theoretical power considerations, population genetics simulations and real data analysis showed strong performance of MinP-val approach for different disease scenarios, whereas SumP-val method was shown to perform poorly when one of the underlying tests had low power. Our analysis of SiMES + SINDI identified the three regions found by Vithana et al. [41], and additionally, the C7orf42 gene. The replication analysis confirmed an association of RXRA-COL5A1 region, which is consistent with the results of Cornes et al. [54], and showed a moderate p-value for C7orf42 gene. The analysis of real data

Table 2 The results of the combined SiMES and SINDI data analysis and the single-SNP p-values from the original article

	COL8A2	ZNF469-LOC100128913	RXRA-COL5A1	COL8A2- TRAPPC3	C7orf42
Chromosome	1	16	9	1	7
Number of SNPs	4	27	73	3	6
Genotype SKAT	3.68E-13	2.13E-15	4.06E-12	2.63E-08	2.55E-07
Haplotype SKAT	5.58E-10	0.149394	0.79	2.78E-05	0.005
MinP-val	3.68E-13	4.22E-15	8.11E-12	4E-08	4.96E-07
SumP-val	1.77E-11	9.44E-10	7.90E-06	4E-08	2.60E-05
Single-SNP analysis*	rs96067: 5.4E-13	rs9938149: 1.63E-16 rs12447690: 1.92E-14	rs1536478: 3.5E-9	-	-

Genome-wide significant p-values for gene-based tests are shown in bold. * SiMES and SINDI meta analysis p-values from Vithana et al. [41].

Table 3 Replication results on Chinese samples from the Singapore Indian Chinese cohort eye study

	COL8A2	ZNF469-LOC100128913	RXRA-COL5A1	COL8A2- TRAPPC3	C7orf42
Genotype SKAT	0.019	0.117	0.001	1	0.014
Haplotype SKAT	0.599	0.479	0.1	1	0.27
MinP-val	0.037	0.223	0.002	0.989	0.028
SumP-val	0.089	0.186	0.001	0.788	0.027
Single-SNP analysis*	rs96067: 0.036	rs9938149: 0.4 rs12447690: 0.03	rs1536478: 0.016	-	-

Significant p-values are shown in bold. For single-SNP analysis the Bonferroni correction corresponds to the four tests. * Trend test within a linear regression model.

highlighted the applicability of our combined approaches to real association studies.

In our simulations the Haplotype SKAT was the most powerful test in many cases, but in real data analysis it performed the worst. It is not known beforehand whether a genotype- or a haplotype-based test would perform better; hence, our proposal to apply a combined approach is a robust choice. Indeed, MinP-val did well in both simulations and real data. This emphasizes the major point of the combined strategy: MinP-val may have slightly lower power when a disease model fits Haplotype SKAT and higher power when the disease model is closer to the second underlying tests. One of the possible reasons for the apparent inconsistency of Haplotype SKAT performance may be that for “Rare” and “Both” simulation models we assumed that rare variants bear the major association signal whereas in the real data only common SNPs were present. However, Haplotype SKAT performed well even for “Common” model when a common SNP was causal. We suppose that for this scenario genotype association translated into an association of haplotypes with a phenotype, which is possible if common SNPs within a region are in high LD with each other. On the other hand, if a causal common SNP within a region is in low LD with other common SNPs within a region then under a genotype-based disease scenario haplotype-based test may have much lower power than a genotype-based test which is observed in the results of the real data analysis.

The methods proposed in this study may be easily generalized to multiple statistical tests, namely, instead of two underlying tests it is possible to apply more tests and combine all of them via the described methodology. In this case the arguments for theoretical p-value calculation for the proposed approaches can be extended in a straightforward manner.

Recently Derkach et al. [56] investigated the performance of the combined approaches, namely, the minimum of p-values and the Fisher p-value combination, for rare variants association scenarios. Although the approaches we propose are similar, our major idea is different. We combine two test statistics for the purpose of widening the set of alternatives for which our test is powerful;

thus, we choose the underlying tests designed for very different phenotype models, whereas Derkach et al. [56] used linear and quadratic tests which are likely to be both powerful under many models. As a result, our conclusions are different from those of Derkach et al. [56]. For example, the authors stated that “hybrid test statistics provide much needed robustness in terms of power for association tests”, whereas we observed that only minimum p-value approach really preserves power when one of the underlying tests underperforms. Secondly, the authors found that in many cases Fisher method outperforms both of the underlying tests, and the minimum p-value approach. However, from our work it is clear that SumP-val (which is similar to the Fisher p-value combination) outperforms all the three tests only when both of the underlying tests have comparable power which is unlikely if the two underlying tests are deliberately chosen to fit very different phenotype models.

One of the limitations of the proposed approaches is the need to use permutations. For theoretical p-value calculation both SumP-val and MinP-val require a correlation coefficient to be estimated via permutations. Moreover, permutations need to be applied when asymptotic distributions of the underlying test statistics are unknown or inadequate to describe the empirical distributions.

The described methodologies may be extended to preserve power under other disease models. For example, the combination of rare-variants and common-variants statistical tests applied to a sequenced region may preserve high power when either only rare or only common variants are associated with a phenotype. However, it is not known how the combined approaches will perform if both common and rare variants are associated with phenotype.

Conclusions

In this study we have investigated the performance of combined haplotype- and genotype-based tests for the purpose of preserving high power under both genotype and haplotype disease scenarios. Based on theoretical power calculations, population genetics simulations and analysis of the real data set we have illustrated high performance and potential utility of combined approaches for association studies.

Additional files

Additional file 1: Table S1. Shapiro-Wilk bivariate normality test p-values for the genome-wide significant genes, additional simulation and real data analysis results. Shapiro-Wilk test was used in real data analysis to justify our assumption of bivariate normality for calculation of theoretical p-values for MinP-val and SumP-val tests. Additional simulations and real data analysis were performed using different pair of underlying tests.

Additional file 2: Power comparison of the gene score haplotype test, the gene score genotype test, MinP-val and SumP-val statistical tests for population genetics simulations, and an estimate of empirical type-1 error. In each panel the top three disease models correspond to the haplotype-based disease scenario, whereas the lower three correspond to the genotype-based scenario. Disease models "Rare", "Both" and "Common" are described in the section "Population genetics simulation". Type-1 error is set to 5%. Panel 1: 50% of rare variants/haplotypes were assumed to be causal; Panel 2: 20% of rare variants/haplotypes were assumed to be causal; Panel 3: 10% of rare variants/haplotypes were assumed to be causal; Panel 4: empirical type-1 error estimate for simulations under the null hypothesis.

Competing interests

The authors declare that they have no competing interests.

Authors' contributions

SZ, AS and AT conceived the study. SZ and AT designed the experiments. SZ conducted the experiments and performed the analysis. TYW, TA, ENV, and CCK provided the GWAS data. SZ and AT wrote the manuscript. SZ, TYW, TA, EV, KCC, AS and AT approved the manuscript.

Funding

This work was supported by the Agency for Science, Technology and Research (A*STAR), Singapore. The first author is a recipient of the Singapore International Graduate Award.

Author details

¹Human Genetics, Genome Institute of Singapore, Singapore, Singapore.
²Saw Swee Hock School of Public Health, National University of Singapore, Singapore, Singapore.
³Singapore Eye Research Institute, Singapore, Singapore.
⁴Department of Ophthalmology, National University Health System, Singapore, Singapore.
⁵Centre for Healthy Brain Ageing (CHeBA), School of Psychiatry, University of New South Wales, Sydney, Australia.

Received: 9 April 2013 Accepted: 13 August 2013

Published: 21 August 2013

References

- Green EK, Grozeva D, Sims R, Raybould R, Forty L, Gordon-Smith K, Russell E, St. Clair D, Young AH, Ferrier IN, et al: **DISC1 exon 11 rare variants found more commonly in schizoaffective spectrum cases than controls.** *Am J Med Genet B Neuropsychiatr Genet* 2011, **156**(4):490–492.
- Norton N, Li D, Rieder Mark J, Siegfried Jill D, Rampersaud E, Züchner S, Mangos S, Gonzalez-Quintana J, Wang L, McGee S, et al: **Genome-wide studies of copy number variation and exome sequencing identify rare variants in BAG3 as a cause of dilated cardiomyopathy.** *The American Journal of Human Genetics* 2011, **88**(3):273–282.
- Ramagopalan SV, Dymment DA, Cader MZ, Morrison KM, Disanto G, Morahan JM, Berlanga-Taylor AJ, Handel A, De Luca GC, Sadovnick AD, et al: **Rare variants in the CYP27B1 gene are associated with multiple sclerosis.** *Ann Neurol* 2011, **70**(6):881–886.
- Xie P, Kranzler HR, Krauthammer M, Cosgrove KP, Oslin D, Anton RF, Farrer LA, Picciotto MR, Krystal JH, Zhao H, et al: **Rare nonsynonymous variants in alpha-4 Nicotinic Acetylcholine receptor gene protect against nicotine dependence.** *Biol Psychiatry* 2011, **70**(6):528–536.
- Wu MC, Kraft P, Epstein MP, Taylor DM, Chanoock SJ, Hunter DJ, Lin X: **Powerful SNP-Set analysis for case-control genome-wide association studies.** *Am J Hum Genet* 2010, **86**(6):929–942.
- Caporaso N, Gu F, Chatterjee N, Sheng-Chih J, Yu K, Yeager M, Chen C, Jacobs K, Wheeler W, Landi MT, et al: **Genome-wide and candidate gene**

- association study of cigarette smoking behaviors. *PLoS ONE* 2009, **4**(2):e4653.
- Hong M-G, Reynolds CA, Feldman AL, Kallin M, Lambert J-C, Amouyel P, Ingelsson E, Pedersen NL, Prince JA: **Genome-wide and gene-based association implicates FRMD6 in alzheimer disease.** *Hum Mutat* 2012, **33**(3):521–529.
- Akey J, Jin L, Xiong M: **Haplotypes vs single marker linkage disequilibrium tests: what do we gain?** *Eur J Hum Genet* 2001, **9**:291–300.
- Schaid DJ: **Power and sample size for testing associations of haplotypes with complex traits.** *Ann Hum Genet* 2006, **70**(1):116–130.
- Browning SR, Browning BL: **Rapid and accurate Haplotype phasing and missing-data inference for whole-genome association studies by use of localized Haplotype clustering.** *The American Journal of Human Genetics* 2007, **81**(5):1084–1097.
- Howie BN, Donnelly P, Marchini J: **A flexible and accurate genotype imputation method for the next generation of genome-wide association studies.** *PLoS Genet* 2009, **5**(6):e1000529.
- Stephens M, Smith NJ, Donnelly P: **A New statistical method for haplotype reconstruction from population data.** *Am J Hum Genet* 2001, **68**(4):978–989.
- Bansal V, Halpern AL, Axelrod N, Bafna V: **An MCMC algorithm for haplotype assembly from whole-genome sequence data.** *Genome Res* 2008, **18**(8):1336–1346.
- Bansal V, Libiger O, Torkamani A, Shork JN: **Statistical analysis strategies for association studies involving rare variants.** *Nat Rev Genet* 2011, **11**:773–785.
- Schatz MC, Delcher AL, Salzberg SL: **Assembly of large genomes using second-generation sequencing.** *Genome Res* 2010, **20**(9):1165–1173.
- Begnini A, Tessari G, Turco A, Malerba G, Naldi L, Gotti E, Boschiero L, Forni A, Ruggi C, Piaserico S, et al: **PTCH1 gene haplotype association with basal cell carcinoma after transplantation.** *Brit J Dermatol* 2010, **163**(2):364–370.
- DIEUDE P, DAWIDOWICZ K, GUEJ M, LEGRAIN Y, WIPFF J, HACHULLA E, DIOT E, SIBILIA J, MOUTHON L, CABANE J, et al: **Phenotype-Haplotype Correlation of IRF5 in Systemic Sclerosis: role of 2 Haplotypes in Disease Severity.** *J Rheumatol* 2010, **37**(5):987–992.
- Lambert JC, Grenier-Boley B, Harold D, Zelenika D, Chouraki V, Kamatani Y, Sleegers K, Ikram MA, Hiltunen M, Reitz C, et al: **Genome-wide haplotype association study identifies the FRMD4A gene as a risk locus for Alzheimer's disease.** *Mol Psychiatry* 2013, **18**(4):461–470.
- Tregouet D-A, König IR, Erdmann J, Munteanu A, Braund PS, Hall AS, Groszhenig A, Linsel-Nitschke P, Perret C, DeSureau M, et al: **Genome-wide haplotype association study identifies the SLC22A3-LPAL2-LPA gene cluster as a risk locus for coronary artery disease.** *Nat Genet* 2009, **41**(3):283–285.
- Bansal V, Bafna V: **HapCUT: an efficient and accurate algorithm for the haplotype assembly problem.** *Bioinformatics* 2008, **24**(16):i153–i159.
- Gauderman WJ, Murcray C, Gilliland F, Conti DV: **Testing association between disease and multiple SNPs in a candidate gene.** *Genet Epidemiol* 2007, **31**(5):383–395.
- Li M, Fu W, Lu Q: **An aggregating U-Test for a genetic association study of quantitative traits.** *BMC Proceedings* 2011, **5**(Suppl 9):S23.
- Li YM, Xiang Y, Sun ZQ: **An entropy-based measure for QTL mapping using extreme samples of population.** *Hum Hered* 2008, **65**(3):121–128.
- Bansal V, Libiger O, Torkamani A, Schork N: **Statistical analysis strategies for association studies involving rare variants.** *Nature Review Genetics* 2010, **11**:773–785.
- Ionita-Laza I, Buxbaum JD, Laird NM, Lange C: **A New testing strategy to identify rare variants with either risk or protective effect on disease.** *PLoS Genet* 2011, **7**(2):e1001289.
- Madsen BE, Browning SR: **A groupwise association test for rare mutations using a weighted Sum statistic.** *PLoS Genet* 2009, **5**(2):e1000384.
- Neale BM, Rivas MA, Voight BF, Altshuler D, Devlin B, Orho-Melander M, Kathiresan S, Purcell SM, Roeder K, Daly MJ: **Testing for an unusual distribution of rare variants.** *PLoS Genet* 2011, **7**(3):e1001322.
- Price AL, Kryukov GV, de Bakker PIW, Purcell SM, Staples J, Wei L-J, Sunyaev SR: **Pooled association tests for rare variants in exon-resequencing studies.** *Am J Hum Genet* 2010, **86**(6):832–838.
- Jin L, Zhu W, Guo J: **Genome-wide association studies using haplotype clustering with a new haplotype similarity.** *Genet Epidemiol* 2010, **34**(6):633–641.
- Sha Q, Dong J, Jiang R, Zhang S: **Tests of association between quantitative traits and haplotypes in a reduced-dimensional space.** *Ann Hum Genet* 2005, **69**(6):715–732.

31. Zaykin DV, Westfall PH, Young SS, Karnoub MA, Wagner MJ, Ehm MG: **Testing association of statistically inferred haplotypes with discrete and continuous traits in samples of unrelated individuals.** *Hum Hered* 2002, **53**(2):79–91.
32. Guo W, Lin S: **Generalized linear modeling with regularization for detecting common disease rare haplotype association.** *Genet Epidemiol* 2009, **33**(4):308–316.
33. Li Y, Byrnes AE, Li M: **To identify associations with rare variants, just WHaT: weighted haplotype and imputation-based tests.** *Am J Hum Genet* 2010, **87**(5):728–735.
34. Zhu X, Feng T, Li Y, Lu Q, Elston RC: **Detecting rare variants for complex traits using family and unrelated data.** *Genet Epidemiol* 2010, **34**(2):171–187.
35. Stouffer SA, Suchman EA, Devinney LC, Star SA, Williams RMJ: **Adjustment during army life.** In *The American soldier*. 1st edition. Princeton, NJ: Princeton Univ; 1949.
36. Tippet LHC: *The method of statistics*. London: Williams and Northgate; 1931.
37. King CR, Rathouz PJ, Nicolae DL: **An evolutionary framework for association testing in resequencing studies.** *PLoS Genet* 2010, **6**(11):e1001202.
38. Boyko AR, Williamson SH, Indap AR, Degenhardt JD, Hernandez RD, Lohmueller KE, Adams MD, Schmidt S, Sninsky JJ, Sunyaev SR, et al: **Assessing the evolutionary impact of amino acid mutations in the human genome.** *PLoS Genet* 2008, **4**(5):e1000083.
39. Gutenkunst RN, Hernandez RD, Williamson SH, Bustamante CD: **Inferring the joint demographic history of multiple populations from multidimensional SNP frequency data.** *PLoS Genet* 2009, **5**(10):e1000695.
40. Wu Michael C, Lee S, Cai T, Li Y, Boehnke M, Lin X: **Rare-variant association testing for sequencing data with the sequence kernel association test.** *Am J Hum Genet* 2011, **89**(1):82–93.
41. Vithana EN, Aung T, Khor CC, Cornes BK, Tay W-T, Sim X, Lavanya R, Wu R, Zheng Y, Hibberd ML, et al: **Collagen-related genes influence the glaucoma risk factor, central corneal thickness.** *Hum Mol Genet* 2011, **20**(4):649–658.
42. Lavanya R, Jeganathan VSE, Zheng Y, Raju P, Cheung N, Tai ES, Wang JJ, Lamoureux E, Mitchell P, Young TL, et al: **Methodology of the Singapore Indian Chinese cohort (SICC) Eye study: quantifying ethnic variations in the epidemiology of eye diseases in Asians.** *Ophthalmic Epidemiol* 2009, **16**(6):325–336.
43. Foong AWP, Saw S-M, Loo J-L, Shen S, Loon S-C, Rosman M, Aung T, Tan DTH, Tai ES, Wong TY: **Rationale and methodology for a population-based study of Eye diseases in Malay people: the Singapore Malay Eye study (SiMES).** *Ophthalmic Epidemiol* 2007, **14**(1):25–35.
44. Su DHW, Wong TY, Foster PJ, Tay W-T, Saw S-M, Aung T: **Central corneal thickness and its associations with ocular and systemic factors: the Singapore Malay Eye study.** *Am J Ophthalmol* 2009, **147**(4):709–716.e701.
45. Wong TCE, et al: **Prevalence and causes of low vision and blindness in an urban malay population: the singapore malay eye study.** *Arch Ophthalmol* 2008, **126**(8):1091–1099.
46. Zhao J, Gupta S, Seielstad M, Liu J, Thalamuthu A: **Pathway-based analysis using reduced gene subsets in genome-wide association studies.** *BMC Bioinformatics* 2011, **12**:17.
47. Price AL, Patterson NJ, Plenge RM, Weinblatt ME, Shadick NA, Reich D: **Principal components analysis corrects for stratification in genome-wide association studies.** *Nat Genet* 2006, **38**(8):904–909.
48. Thalamuthu A, Zhao J, Keong G, Kondragunta V, Mukhopadhyay I: **Association tests for rare and common variants based on genotypic and phenotypic measures of similarity between individuals.** *BMC Proceedings* 2011, **5**(Suppl 9):S89.
49. Freedman D, Lane D: **A nonstochastic interpretation of reported significance levels.** *Journal of Business and Economic Statistics* 1983, **1**:292–298.
50. Kennedy PE: **Randomization tests in econometrics.** *Journal of Business and Economic Statistics* 1995, **13**:85–94.
51. Ter Braak CJF: **Permutation versus bootstrap significance tests in multiple regression and ANOVA.** In *Bootstrapping and related techniques*. Edited by Jockel KH, Rothe G, Sendler W. Berlin: Springer; 1992.
52. Anderson MJ, Legendre P: **An empirical comparison of permutation methods for tests of partial regression coefficients in a linear model.** *J Stat Comput Sim* 1999, **62**(3):271–303.
53. Wagner BD, Zerbe GO, Mexal S, Leonard SS: **Permutation-based adjustments for the significance of partial regression coefficients in microarray data analysis.** *Genet Epidemiol* 2008, **32**(1):1–8.
54. Cornes BK, Khor CC, Nongpiur ME, Xu L, Tay W-T, Zheng Y, Lavanya R, Li Y, Wu R, Sim X, et al: **Identification of four novel variants that influence central corneal thickness in multi-ethnic Asian populations.** *Hum Mol Genet* 2012, **21**(2):437–445.
55. Zhang H, Xu L, Chen C, Jonas J: **Central corneal thickness in adult Chinese. Association with ocular and general parameters. The Beijing Eye Study.** *Graefes Archive for Clinical and Experimental Ophthalmology* 2008, **246**(4):587–592.
56. Derkach A, Lawless JF, Sun L: **Robust and powerful tests for rare variants using Fisher's method to combine evidence of association from two or more complementary tests.** *Genet Epidemiol* 2013, **37**(1):110–121.

doi:10.1186/1471-2164-14-569

Cite this article as: Zakharov et al.: Combined genotype and haplotype tests for region-based association studies. *BMC Genomics* 2013 **14**:569.

Submit your next manuscript to BioMed Central and take full advantage of:

- Convenient online submission
- Thorough peer review
- No space constraints or color figure charges
- Immediate publication on acceptance
- Inclusion in PubMed, CAS, Scopus and Google Scholar
- Research which is freely available for redistribution

Submit your manuscript at
www.biomedcentral.com/submit

