

# SCIENTIFIC REPORTS



OPEN

## A method to estimate the contribution of regional genetic associations to complex traits from summary association statistics

Received: 29 February 2016

Accepted: 18 May 2016

Published: 08 June 2016

Guillaume Pare<sup>1,2,3,4</sup>, Shihong Mao<sup>3</sup> & Wei Q. Deng<sup>5</sup>

Despite considerable efforts, known genetic associations only explain a small fraction of predicted heritability. Regional associations combine information from multiple contiguous genetic variants and can improve variance explained at established association loci. However, regional associations are not easily amenable to estimation using summary association statistics because of sensitivity to linkage disequilibrium (LD). We now propose a novel method, LD Adjusted Regional Genetic Variance (LARGV), to estimate phenotypic variance explained by regional associations using summary statistics while accounting for LD. Our method is asymptotically equivalent to a multiple linear regression model when no interaction or haplotype effects are present. It has several applications, such as ranking of genetic regions according to variance explained or comparison of variance explained by two or more regions. Using height and BMI data from the Health Retirement Study ( $N = 7,776$ ), we show that most genetic variance lies in a small proportion of the genome and that previously identified linkage peaks have higher than expected regional variance.

Currently known genetic associations only explain a relatively small proportion of complex traits variance. In accordance with the widely accepted polygenic nature of complex traits, it has been proposed that weak, yet undetected, associations underlie complex trait heritability of a wide variety of phenotypes such as height<sup>1</sup>, cognitive function<sup>2</sup> or rheumatoid arthritis<sup>3</sup>, etc. We have recently shown that the joint association of multiple weakly associated variants over large chromosomal regions contributes to complex traits variance<sup>4</sup>. Such regional associations are not easily amenable to estimation using summary-level association statistics because of sensitivity to linkage disequilibrium (LD). Nonetheless, only large meta-analyses have the necessary power to identify weakly associated variants and results are typically reported in the form of summary association statistics. In this report, we propose a novel method to assess the contribution of regional associations to complex traits variance using summary association statistics. Estimation of regional associations with our novel method, LD Adjusted Regional Genetic Variance (LARGV), can help identify key genomic regions involved in regulation of complex traits.

Clustering of weak associations within defined chromosomal regions has been previously suggested<sup>5</sup> and can increase variance explained at established association loci as compared to genome-wide significant SNPs alone<sup>6</sup>. Such regional associations extended up to 433.0 Kb from genome-wide significant SNPs<sup>4</sup>, a distance compatible with long-range *cis* regulation of gene expression<sup>7,8</sup>. Furthermore, regional associations appeared to be the results of multiple weak associations rather than one or a few very significant univariate associations. These results point towards the existence of key regulatory regions where functional genetic variants aggregate, the identification of which can lead to novel biological insights and a better understanding of complex traits genetics.

Several methods have been described to estimate the overall contribution of common genetic variants to complex traits variance, but no method was specifically designed to estimate regional association using summary association statistics while accounting for LD (Table 1). For instance, a popular approach is based on variance

<sup>1</sup>Department of Pathology and Molecular Medicine, McMaster University, Hamilton, ON L8S 4L8, Canada.

<sup>2</sup>Population Genomics Program, Department of Clinical Epidemiology and Biostatistics, McMaster University, Hamilton, ON L8S 4L8, Canada. <sup>3</sup>Population Health Research Institute, Hamilton Health Sciences and McMaster University, Hamilton, ON L8L 2X2, Canada. <sup>4</sup>Thrombosis and Atherosclerosis Research Institute, Hamilton, ON L8L 2X2, Canada. <sup>5</sup>Department of Statistical Sciences, University of Toronto, Toronto, ON M5S 3G3, Canada. Correspondence and requests for materials should be addressed to G.P. (email: pareg@mcmaster.ca)

Method	Advantages	Disadvantages
<b>Multi-SNP locus association</b> <sup>13</sup>	-Can estimate genetic variance of medium to large regions (i.e. small regions might not have SNPs remaining after LD and $p$ -value pruning)	-Need for LD pruning as it assumes uncorrelated markers -Need for $p$ -value pruning -Assumes population structure has been entirely adjusted for
<b>LDscore</b> <sup>11</sup>	-Includes all SNPs irrespective of LD -Adjusts for potential bias such as population stratification -Can estimate genetic covariance between two traits	-Need for a reference LD structure -Best suited for genome-wide analysis or large regions
<b>Distribution of effect size</b> <sup>3,14,15</sup> (AVENGEME and ABPA)	-Estimates the proportion of markers affecting a trait -Can estimate genetic covariance between two traits	-Need for LD pruning as it assumes uncorrelated markers -Assumes population structure has been entirely adjusted for -Best suited for genome-wide analysis or large regions
<b>LD Adjusted Regional Genetic Variance (LARGV)</b>	-Includes all SNPs irrespective of LD -Can estimate genetic variance of any region, small or large	-Need for a reference LD structure -Assumes population structure has been entirely adjusted for

**Table 1. Advantages and disadvantages of existing methods to estimate regional genetic variance from summary association statistics.**

component models using genetic relatedness as the variance-covariance matrix of the random effect. An implementation of this approach uses REML<sup>1</sup> to estimate the genetic effect, and modifications have been reported to either take into account LD between SNPs<sup>9</sup> or to handle large datasets<sup>10</sup>. While very useful and informative, all of these approaches require individual-level data. This latter pitfall has been overcome by the development of LDscore<sup>11</sup>, which uses summary-level association statistics as inputs and has been shown to be equivalent to Elston regression<sup>12</sup>. However, LDscore is ill suited for estimation of regional heritability because it is based on regressing genetic effects on the sum of linkage disequilibrium and requires large number of SNPs for precise estimations. A multi-SNP locus-association method has been described that uses summary association statistics, but it necessitates to first prune SNPs for  $p$ -value ( $p < 0.01$ ) and LD ( $r^2 < 0.1$ )<sup>13</sup>. Finally, alternative methods<sup>14,15</sup> estimate genetic variance from the distribution of effect sizes and are not appropriate for regional genetic variance because the small number of SNPs will lead to an imprecise estimation and linkage equilibrium is assumed. There is thus a need for a method to estimate the regional contribution of common genetic variants using summary association statistics while simultaneously taking LD into account.

## Results

**Comparison of genetic variance estimated using summary statistics and variance component models.** Our method estimates the regional contribution to complex trait variance using summary data by adjusting the variance explained by each SNP for its LD with neighboring SNPs. Simulations using 1000 Genomes Project<sup>16</sup> (1000G) data showed that genetic regional variance is accurately estimated by our method when no haplotype or interaction effects are present (Figure S1), as predicted by theoretical derivations (see Methods). We sought to compare estimation of overall genetic variance by our novel method to variance component models. We divided the genome into SNP blocks of median size 250 Kb (85–95 contiguous SNPs) minimizing inter-block LD and applied our method to BMI and height in the Health Retirement Study<sup>17</sup> (HRS;  $N = 7,776$ ) for which individual-level genotypes were available. Using only summary association statistics for our method and corresponding individual-level data for variance component models, both methods provided consistent estimates of genome-wide genetic variance (Fig. 1). We explored the impact of adjustment for genetic principal components. The first 20 components provided adequate protection against population stratification while the inclusion of fewer components led to the inflation in genetic variance, especially for height. Using 20 components, genetic variance was estimated at 0.12 (95% CI 0.01–0.24) for BMI with our novel method and 0.14 (95% CI 0.05–0.23) with variance component models<sup>18</sup>. The corresponding estimates for height were 0.28 (95% CI 0.17–0.39) and 0.30 (95% CI 0.20–0.39).

**Using summary association statistics to identify regional associations.** We next sought to compare the ability of our method to identify regional associations with other approaches also using summary association statistics. We ranked SNP blocks (median size of 250 Kb) according to decreasing regional genetic variance by applying three different methods on Genetic Investigation of Anthropometric Traits consortium (GIANT) summary association statistics<sup>19,20</sup>: (1) our proposed approach (LARGV), (2) LDscore and (3) the multi-SNP locus-association proposed by Ehret and colleagues<sup>13</sup>. We then used SNP block ranks derived from GIANT by each of the three methods to estimate genetic variance in HRS (which is not part of GIANT) with variance component models, successively adding a higher proportion of SNP blocks. As illustrated in Fig. 2, the genetic variance estimated by LARGV increased more rapidly with the proportion of top SNP blocks included as compared

to other two methods. The multi-SNP locus-association method performed well, but a large proportion of SNP blocks did not have any SNP left after LD and  $p$ -value pruning (68.4% for BMI and 26.7% for height), highlighting the disadvantage of pruning instead of adjusting for LD. We obtained consistent results when using LD data from 1000G data instead of HRS (Figure S2), changing SNP block size (Figure S3), or changing the number of neighboring SNPs included in LD calculations (Figure S4).

A relatively small proportion of blocks contributed disproportionately to genetic variance. The top 25% SNP blocks explained 0.94 of BMI genetic variance (i.e. = 0.132/0.140) and 0.83 of height genetic variance when using LARGV on GIANT data to rank the genetic variance of each SNP block. These results could potentially be explained by the presence of one or more very strong associations in each of these top SNP blocks. To explore this possibility, we recorded the minimum univariate association SNP  $p$ -value for each block in both GIANT and HRS (Fig. 3). Median minimum univariate  $p$ -value was  $2.0 \times 10^{-2}$  for BMI in GIANT, with 1% of blocks having one or more genome-wide significant associations ( $p < 5 \times 10^{-8}$ ). On the other hand, the median minimum  $p$ -value was  $1.9 \times 10^{-3}$  for height in GIANT, with 9% of blocks having one or more genome-wide significant associations. The median minimum univariate  $p$ -values were  $1.8 \times 10^{-2}$  and  $1.7 \times 10^{-2}$  for BMI and height in HRS. No SNP reached genome-wide significance in HRS.

**Analysis of known linkage peaks.** A unique application of our method is the estimation of genetic variance over extended genomic regions using summary association statistics. We therefore tested the hypothesis that previously identified linkage peaks are enriched for regional associations. Based on results from the largest linkage study of height and BMI, three peaks with suggestive (LOD > 2.0) evidence of linkage in Europeans were identified, all for height<sup>21</sup>. The only peak with LOD > 3.0 showed a significant ( $p = 0.002$ ) enrichment in regional association within a distance of  $\pm 7.5$  Mb from the linkage marker, corresponding to an estimated excess regional genetic variance of 0.0044 over the genome-wide average (Table 2). Upon a closer inspection, the region encompasses several sub-regions with genome-wide significant associations.

## Discussion

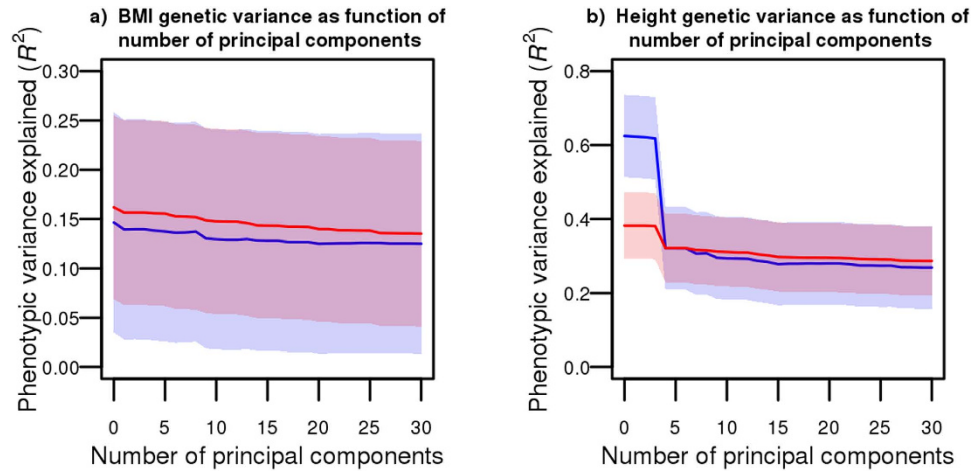
We propose a novel method to estimate regional genetic variance from summary association statistics. Using this method, we confirmed a major role of regional associations in complex trait heritability, whereby the aggregation of genetic associations contributes disproportionately to phenotypic variance. Selecting the top SNP blocks from the GIANT meta-analysis, we showed that 25% of the genome is responsible for up to 0.94 and 0.83 of BMI and height genetic variance. A large proportion of these blocks had unremarkable minimum univariate  $p$ -values, suggesting the presence of multiple weak associations underlies their impact on phenotypic variance, especially for BMI. The concentration of genetic associations within these regions supports the existence of critical nodes in the genetic regulation of complex traits such as height and BMI, with implications not only for association testing but also for population genetics and natural selection. These results also suggest that a combination of strong genetic associations and regional associations contribute to complex traits variance, with the relative proportions varying across traits. For instance, a higher proportion of genetic variance was found in the top 25% blocks for BMI yet these blocks had less significant minimum univariate  $p$ -value than height.

Our method can also be used to estimate the genetic variance explained by extended regions. We therefore tested the hypothesis that some of the previously identified linkage peaks are the result of regional associations. The only known linkage peak with LOD > 3.0 for height showed a marked and significant enrichment in regional association. This region had been previously identified in multiple linkage studies<sup>21,22</sup>. Genetic variance explained by the region was estimated at 0.0044, which is unlikely to explain the linkage peak by itself. Nonetheless, the juxtaposition of linkage and regional associations points towards concentration of functional variants as a potential explanation for the observed linkage. The lack of regional association at other peaks can be explained by false-positive linkage results, rare variant associations undetected by genome-wide association studies or by differences in genetic architecture between studied populations.

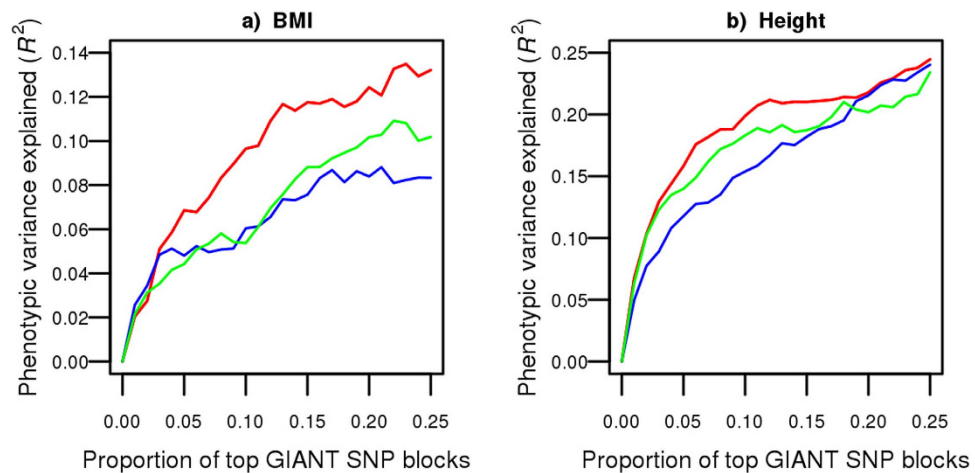
Our proposed method has several advantages and is complementary to other methods. First, it provides results that are highly consistent with the widely used variance component models requiring individual-level data. Second, it is computationally straightforward and can be applied to large datasets. Third, it is agnostic and therefore complementary to functional annotations of the genome. Fourth, since SNPs are not pruned by either LD or significance, every SNP block or region can be evaluated.

A few limitations are worth mentioning. First, the assumption that SNPs contribute to genetic variance without any interaction or haplotype effects can lead to an underestimation of genetic variance. While our estimates were consistent with variance component models, there is a need for statistical models that better capture genetic variance when strong haplotype effects are expected<sup>23</sup>. Second, estimation of overall genome-wide genetic variance using summary association statistics is dependent on both the accuracy of SNP effect size estimates and the correct specification of LD structure. For instance, differences in adjustment for population stratification (e.g. principal components) in individual studies that participate in a meta-analysis could potentially influence the results. Ideally, LD data would come from the same population used to derive the summary association statistics, which could be included as summary statistics for each SNP (i.e.  $\eta_{di}$ , see Methods). While this information is currently not available, consistent results were observed when using LD data from either HRS or 1000G, demonstrating the robustness of our approach to LD misspecification. Third, we have only tested continuous traits. Nevertheless, our method can be easily adapted to other outcome types through the use of generalized linear models.

In this report, we establish a novel method to estimate the regional contribution of common variants to complex traits variance using summary association statistics. Our method has several applications, such as ranking of genetic regions for genetic variance and identification of key regions contributing to genetic variance. Our



**Figure 1. Genome-wide genetic variance estimated by summary association statistics and variance component models.** The overall genetic variance estimated by summary association statistics and variance component models is illustrated as a function of number of principal components for BMI (a) and height (b). Blue lines represent estimates of genetic variance using summary association statistics with LARGV; with 95% confidence intervals illustrated as blue shaded area. Corresponding estimates for variance component models are in red.

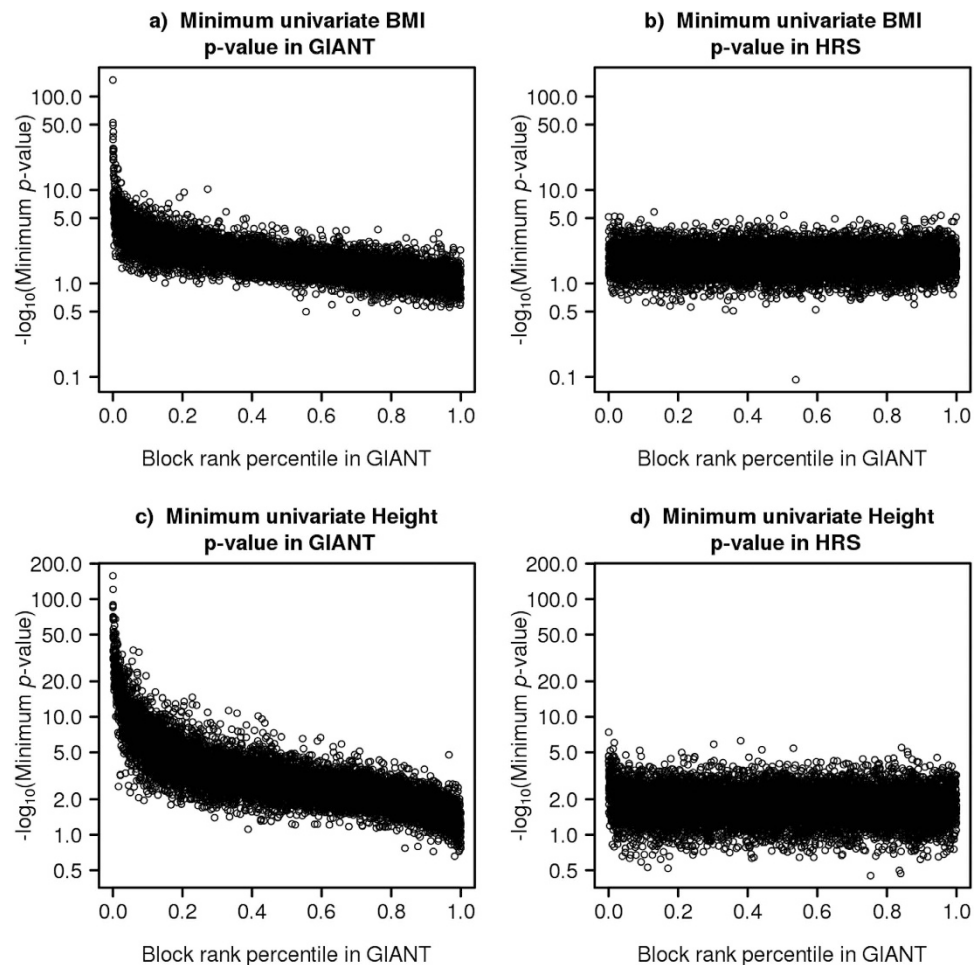


**Figure 2. Genetic variance as a function of proportion of top SNP blocks.** Genetic variance in HRS based on SNP block ranking derived from GIANT summary association statistics. Three methods were tested to rank SNP blocks: LARGV (red), LD Score (blue) and multi-SNP locus-association (green). The median SNP block size was 250 Kb (i.e. 85–95 SNPs). Genetic variance was calculated in HRS using variance component models for BMI (a) and height (b).

method can also be used to perform network analysis using summary association statistics, or to combine summary association statistics with other types of genetic annotations as we have shown with linkage results.

## Methods

**Methods overview.** We have previously shown that large-region joint association, where multiple genetic variants are included as independent variables in a linear model, is a simple and powerful way to test for regional associations when individual-level data are available<sup>4</sup>. However, such regional joint associations are not easily amenable to estimation using summary data because of sensitivity to linkage disequilibrium. To evaluate the contribution of large regions joint associations to the genetic variance of complex traits, we first devised an algorithm to divide the genome in blocks of SNPs in such a way to minimize inter-block linkage disequilibrium and thus “spillage” associations. We then derived a method to estimate regional associations using summary data and showed this method to be equivalent to a multiple linear regression model when genetic effects are strictly additive (i.e. no haplotype or interaction effect). Using regional variance estimates from summary association statistics for height and BMI from the Genetic Investigation of Anthropometric Traits (GIANT) consortium, we estimated the distribution of genetic effects across the genome and validated results in the Health Retirement



**Figure 3.** Minimum univariate SNP association  $p$ -value for each SNP block. SNP block ranks are based on GIANT data using our novel approach while the minimum univariate SNP association  $p$ -values were taken from GIANT (a,c) or calculated in HRS (b,d) for BMI (a,b) and height (c,d). The median SNP block size was 250 Kb (i.e. 85–95 SNPs).

Chr.	Peak Marker	LOD Score	Excess genetic variance	95% CI Upper limit	95% CI lower limit	$p$ -value
11	D11S2000	2.74	-0.0021	0.0002	-0.0044	0.07
12	D12S1301	2.07	-0.0005	0.0014	-0.0024	0.61
15	D15S655	3.00	0.0044	0.0072	0.0017	<b>0.002</b>

**Table 2.** Excess regional genetic variance at three suggestive (LOD > 2.0) linkage peaks for height using GIANT summary association statistics. Regional genetic variance was calculated within  $\pm 7.5$  Mb of each peak marker and compared to genome-wide average for regions of equivalent size.  $P$ -values are two-sided.

Study (HRS), which is not part of GIANT. A user-friendly software is available to produce SNP blocks and LD adjustment upon request at [pareg@mcmaster.ca](mailto:pareg@mcmaster.ca).

**Dividing the genome into SNP blocks.** We first divided the genome into regions of contiguous SNPs varying in size (e.g. from 195 SNPs to 205 SNPs), herein referred to as SNP blocks and used as units for regional associations. To minimize inter-block LD and thus “spillage” associations, we devised a greedy algorithm optimizing the choice of block boundary sequentially from one end of a chromosome to the other. Briefly, using a user-defined minimum and maximum block size (in number of SNPs) and starting at one end of a chromosome arm, each possible “cut-point” between the first and second block are tested and maximal LD ( $r^2$ ) between pairs of SNPs crossing block boundary is calculated. The cut-point that minimizes the maximal LD is chosen, thus defining the first block, and the procedure is repeated for each subsequent block until all SNPs on a chromosome arm have been assigned to a block. We empirically determined that SNP blocks of size 85 SNPs to 95 SNPs (median 90 SNPs) had a median physical size of 250 Kb.



**Estimating the contribution of regional genetic associations with individual-level genotypes.** The use of adjusted  $R^2$  lends itself nicely to estimation of regional variance explained when individual-level genotypes are available<sup>4</sup>. In this context, SNPs comprised in a given SNP block are included as independent variables in a multiple linear regression model and the goodness of fit statistic, adjusted  $R^2$  calculated. Because the adjusted  $R^2$  accounts for the number of SNPs included in each block, the expected adjusted  $R^2$  is zero under the null hypothesis of no association and the expected sum of adjusted  $R^2$  over all SNP blocks is also zero. The overall contribution of regional associations to complex traits variance can be estimated by simply summing the adjusted  $R^2$  over all (or selected) SNP blocks. Furthermore, the distribution of adjusted  $R^2$  under the null hypothesis of no association has been previously described<sup>24</sup> and can be used to derive the distribution of the sum of adjusted  $R^2$ .

**Estimating the contribution of regional genetic associations with summary association statistics.** Estimating the contribution of regional genetic associations from summary association statistics is challenging when the exact SNP linkage disequilibrium structure of source populations is unknown. While approaches have been described to perform joint or conditional associations<sup>23,25</sup> using estimated SNP covariance matrices, they do not perform well when estimating regional variance explained because of sensitivity to misspecification of linkage disequilibrium (data not shown) and ensuing an overestimation of regional associations. We therefore created a simple procedure to estimate regional variance explained from summary association statistics data, adjusting for linkage disequilibrium.

Without loss of generalizability, we assume a quantitative trait ( $Y$ ) standard normally distributed and genotypes normalized to have mean = 0 and standard deviation = 1 throughout. Given an  $n \times m$  genotype matrix  $X$  representing genotypes at  $m$  SNPs in  $n$  individuals and the pairwise linkage disequilibrium ( $r^2$ ) between two SNPs  $k$  and  $l$  as  $r_{k,l}^2$  for a SNP  $d$ , the following LD adjustment ( $\eta_d$ ) can be defined as the summation of LD between the  $d^{\text{th}}$  SNP and 100 SNPs upstream and downstream:

$$\eta_d = \sum_{e=d-100}^{e=d+100} r_{d,e}^2 \quad (1)$$

with a distance of 100 SNPs assumed sufficient to ensure linkage equilibrium (other values might be used). Only including SNPs with summary GWAS statistics in the sum, variance explained by each SNP  $d$  is given by:

$$R_d^2 = \frac{b_d^2}{\eta_d} \quad (2)$$

where  $b_d$  denotes the univariate regression coefficient commonly reported in GWAS results (assuming genotypes have been standardized to have mean zero and SD = 1). Regional variance explained is then given by the sum of  $R_d^2$  over SNPs in a given region. Assuming a strictly additive genetic model where each SNP contributes additively to a trait without any interaction or haplotype effects, we demonstrate the expected total variance explained over a region  $E(\sum_d R_d^2)$  is approximately equal to the expected value of the multiple linear regression variance explained  $E(R^2)$  when the sample size is sufficiently large.

To simplify the calculation, we define  $D$  such that  $X'X = D = [D_1 \ D_2 \ \dots \ D_m]$  is an  $m$  by  $m$  symmetric matrix, where  $D_k$  is a  $m \times 1$  vector whose entries represent the  $k^{\text{th}}$  column of  $D$ . We will make use of the following properties:

$$\begin{aligned} D'D_{(m \times m)} &= DD_{(m \times m)} = DD'_{(m \times m)} = D_1D_1' + D_2D_2' + \dots + D_mD_m' \\ \frac{D_k'D_k}{\eta_k} &= \frac{\text{tr}(D_kD_k')}{\eta_k} = \frac{\sum_{d=1}^m N^2 r_{k,d}^2}{\eta_k} = \frac{\sum_{d=1}^m N^2 r_{k,d}^2}{\sum_{d=k-100}^{k+100} r_{k,d}^2} \sim N^2 \\ \text{tr}(DD') &= \sum_{k=1}^m D_k'D_k = \sum_{k=1}^m \text{tr}(D_kD_k') = \sum_{k=1}^m \sum_{d=1}^m N^2 r_{k,d}^2 \sim \sum_{k=1}^m \eta_k N^2 \end{aligned} \quad (3)$$

where  $N$  is the total number of individuals in the sample.

*Estimation of regional genetic variance with multiple linear regression models.* Suppose the genotype matrix is fixed while the true genetic effect is a random vector  $\beta$ , whose individual components, i.e. the SNPs,  $i = 1, 2, \dots, m$ , have mean zero and variance  $\sigma^2$ . The size of the variance  $\sigma^2$  is on the scale of  $M^{-1}$ , where  $M$  is the number of genome-wide SNPs. The genetic model can be expressed as:

$$Y = X\beta + \varepsilon \quad (4)$$

where  $\varepsilon$  is a vector of standard normal error with identity variance covariance matrix. Then, the vector of estimated multiple linear regression coefficients  $B$  is given by:

$$B = (X'X)^{-1}X'Y = (X'X)^{-1}X'(X\beta + \varepsilon) = \beta + (X'X)^{-1}X'\varepsilon \quad (5)$$

The multivariate variance explained  $R^2$  can be written in terms of the true effect and the error term:

$$\begin{aligned}
 R^2 &= \frac{1}{N}(XB)'(XB) \\
 &= \frac{1}{N}(X\beta)'(X\beta) + \frac{1}{N}(X(X'X)^{-1}X'\varepsilon)'(X(X'X)^{-1}X'\varepsilon) \\
 &= \frac{1}{N}(X\beta)'(X\beta) + \frac{1}{N}\varepsilon'(X(X'X)^{-1}X')\varepsilon
 \end{aligned}
 \tag{6}$$

and since the error term has identity variance covariance and the true effect  $\beta$  has variance covariance matrix  $\sigma^2I$ , the expected variance explained is simplified to

$$\begin{aligned}
 E[R^2] &= \frac{1}{N}E(\beta'X'X\beta + \varepsilon'(X(X'X)^{-1}X')\varepsilon) \\
 &= \frac{\text{tr}(X'X\sigma^2I)}{N} + \frac{\text{tr}(X(X'X)^{-1}X')}{N} \\
 &= m\sigma^2 + \frac{m}{N}
 \end{aligned}
 \tag{7}$$

and the variance can be calculated accordingly:

$$\begin{aligned}
 \text{Var}[R^2] &= \frac{1}{N^2}\text{Var}(\beta'X'X\beta + \varepsilon'(X(X'X)^{-1}X')\varepsilon) \\
 &= \frac{2\text{tr}(X'X\sigma^2IX'X\sigma^2I)}{N^2} + \frac{2\text{tr}(X(X'X)^{-1}X'X(X'X)^{-1}X')}{N^2} \\
 &= \frac{2\sigma^4\text{tr}(DD)}{N^2} + \frac{2m}{N^2} \\
 &\sim 2\sigma^4\sum_{k=1}^m\eta_k + \frac{2m}{N^2} \\
 &\sim \frac{2m}{N^2}
 \end{aligned}
 \tag{8}$$

*Estimation of regional genetic variance using summary association statistics.* The univariate regression coefficients, denoted by lower case  $b$  and directly obtained from GWAS summary statistics, are given by

$$b = \frac{X'Y}{N} = \frac{1}{N}X'X\beta + \frac{1}{N}X'\varepsilon
 \tag{9}$$

The total variance explained over a region  $\sum_d R_d^2$  can be calculated using only the univariate regression coefficients from GWAS:

$$\begin{aligned}
 \sum_d R_d^2 &= \sum_d \frac{b_d^2}{\eta_d} \\
 &= b' \begin{bmatrix} 1/\eta_1 & 0 & \dots & 0 \\ 0 & 1/\eta_2 & \dots & \vdots \\ \vdots & \vdots & \ddots & 0 \\ 0 & \dots & 0 & 1/\eta_m \end{bmatrix} b \\
 &= \frac{1}{N^2}(X'X\beta)' \begin{bmatrix} 1/\eta_1 & 0 & \dots & 0 \\ 0 & 1/\eta_2 & \dots & \vdots \\ \vdots & \vdots & \ddots & 0 \\ 0 & \dots & 0 & 1/\eta_m \end{bmatrix} (X'X\beta) \\
 &\quad + \frac{1}{N^2}(X'\varepsilon)' \begin{bmatrix} 1/\eta_1 & 0 & \dots & 0 \\ 0 & 1/\eta_2 & \dots & \vdots \\ \vdots & \vdots & \ddots & 0 \\ 0 & \dots & 0 & 1/\eta_m \end{bmatrix} (X'\varepsilon)
 \end{aligned}$$

We can rewrite the expression by defining:

$$\Lambda = \begin{bmatrix} 1/\eta_1 & 0 & \cdots & 0 \\ 0 & 1/\eta_2 & \cdots & \vdots \\ \vdots & \vdots & \ddots & 0 \\ 0 & \cdots & 0 & 1/\eta_m \end{bmatrix}$$

Thus, we have (Equation 10):

$$\begin{aligned} \sum_d R_d^2 &= \frac{1}{N^2} \beta' \begin{bmatrix} D_1' \\ D_2' \\ \vdots \\ D_m' \end{bmatrix} \Lambda [D_1 \ D_2 \ \cdots \ D_m] \beta + \frac{1}{N^2} \sum_{d=1}^m (\epsilon' X_d X_d' \epsilon) / \eta_d \\ &= \frac{1}{N^2} \beta' \{D_1' D_1 / \eta_1 + D_2' D_2 / \eta_2 + \cdots + D_m' D_m / \eta_m\} \beta + \frac{1}{N^2} \sum_{d=1}^m (\epsilon' X_d X_d' \epsilon) / \eta_d \end{aligned}$$

Since  $\text{tr}(D_k D_k') / \eta_k = \frac{\sum_{d=1}^m N^2 r_{k,d}^2}{\eta_k} = \frac{\sum_{d=1}^m N^2 r_{k,d}^2}{\sum_{d=k-100}^{k+100} r_{k,d}^2} \sim N^2$  as LD tends to be weak 100 SNPs upstream and down-

stream away from the index SNP, we can simplify the expected total variance explained using the summary statistics to (Equation 11):

$$\begin{aligned} E(\sum_d R_d^2) &= E\left(\frac{1}{N^2} \beta' \{D_1' D_1 / \eta_1 + D_2' D_2 / \eta_2 + \cdots + D_m' D_m / \eta_m\} \beta + \frac{1}{N^2} \sum_{d=1}^m (\epsilon' X_d X_d' \epsilon) / \eta_d\right) \\ &= \frac{\sigma^2}{N^2} \text{tr}(D_1' D_1 / \eta_1 + D_2' D_2 / \eta_2 + \cdots + D_m' D_m / \eta_m) + \frac{1}{N^2} \sum_d \text{tr}(X_d X_d') / \eta_d \\ &\sim m\sigma^2 + \frac{1}{N} \sum_d 1 / \eta_d \end{aligned}$$

The variance of  $\sum_d R_d^2$  can be similarly derived. By using the cyclic property of trace, and the fact that  $D\Lambda D$  is a positive definite matrix, an upper bound for the variance can be expressed as (Equation 12):

$$\begin{aligned} \text{Var}(\sum_d R_d^2) &= \frac{1}{N^4} \text{Var}((X' X \beta)' \Lambda (X' X \beta)) + \frac{1}{N^4} \text{Var}(\epsilon' X \Lambda X' \epsilon) \\ &= \frac{1}{N^4} \text{Var}(\beta' D' \Lambda D \beta) + \frac{2}{N^4} \text{tr}(X \Lambda X' X \Lambda X') \\ &= \frac{2\sigma^4}{N^4} \text{tr}(D \Lambda D D \Lambda D) + \frac{2}{N^4} \text{tr}(D \Lambda D \Lambda) \\ &\quad - \frac{2\sigma^4}{N^4} \text{tr}(D \Lambda D)^2 + \frac{2}{N^4} \left\{ \sum_{d=1}^m \text{tr} \left( \frac{D_d D_d'}{\eta_d^2} \right) \right\} \\ &\sim 2m^2 \sigma^4 + \frac{2}{N^2} \left\{ \sum_{d=1}^m \frac{1}{\eta_d} \right\} \\ &\quad - \frac{2}{N^2} \left\{ \sum_{d=1}^m \frac{1}{\eta_d} \right\} \end{aligned}$$

where  $\frac{2\sigma^4}{N^4} \text{tr}(D \Lambda D D \Lambda D) \leq \frac{2\sigma^4}{N^4} \text{tr}(D \Lambda D)^2$  and both terms are small with respect to  $\frac{2}{N^2} \left\{ \sum_{d=1}^m \frac{1}{\eta_d} \right\}$ .

**Comparison of regional genetic variance estimated using multiple regression models and summary association statistics.** The expected values are equivalent between multiple linear regression model  $(m\sigma^2 + \frac{m}{N})$  and the summary statistics derived regional sum  $(m\sigma^2 + \frac{1}{N} \sum_d 1/\eta_d)$ , with the number of SNPs  $m$  replaced by the “effective” number of genetic markers  $\sum_{d=1}^m \frac{1}{\eta_d}$ . Variance is slightly bigger for multiple linear regression models  $(\sim \frac{2m}{N^2})$  as compared to regional sum  $(\sim \frac{2}{N^2} \left\{ \sum_{d=1}^m \frac{1}{\eta_d} \right\})$  because the “effective” number of genetic markers  $\sum_{d=1}^m \frac{1}{\eta_d}$  is always equal (no LD) or less than the number of markers ( $m$ ). In other words,  $\text{Var}(\sum_d R_d^2)$  is expected to be equal (no LD) or lower than the corresponding  $\text{Var}(R^2)$ .

**Health Retirement Study.** We conducted large region joint association analysis for height using genome-wide data from the publicly available Health Retirement Study (HRS; dbGaP Study Accession: phs000428.v1.p1). HRS quality control criteria were used for filtering of both genotype and phenotype data,



namely: (1) SNPs and individuals with missingness higher than 2% were excluded, (2) related individuals were excluded, (3) only participants with self-reported European ancestry genetically confirmed by principal component analysis were included, (4) SNPs with Hardy-Weinberg equilibrium  $p < 1 \times 10^{-6}$  were excluded, (5) individuals for whom the reported sex does not match their genetic sex were excluded, (6) SNPs with minor allele frequency lower than 0.02 were removed. The final dataset included 7,776 European participants genotyped for 740,748 SNPs. Height and BMI was adjusted for age and sex in all analyses. To mitigate the effect of outliers, we have removed values outside the 1<sup>st</sup> and 99<sup>th</sup> percentile range for each of height and BMI. All analyses are adjusted for the first 20 genetic principal components unless stated otherwise. All LD estimates used throughout the manuscript were derived from HRS genotypes. HRS was not part of the GIANT meta-analysis of height and BMI<sup>26,27</sup>.

## References

1. Yang, J. *et al.* Common SNPs explain a large proportion of the heritability for human height. *Nature Genetics* **42**, 565–569, doi: 10.1038/ng.608 (2010).
2. Davies, G. *et al.* Genetic contributions to variation in general cognitive function: a meta-analysis of genome-wide association studies in the CHARGE consortium (N = 53949). *Mol Psychiatry* **20**, 183–192, doi: 10.1038/mp.2014.188 (2015).
3. Stahl, E. A. *et al.* Bayesian inference analyses of the polygenic architecture of rheumatoid arthritis. *Nature genetics* **44**, 483–489, doi: 10.1038/ng.2232 (2012).
4. Pare, G., Asma, S. & Deng, W. Q. Contribution of large region joint associations to complex traits genetics. *PLoS Genet* **11**, e1005103, doi: 10.1371/journal.pgen.1005103 (2015).
5. Beyene, J., Tritchler, D., Asimit, J. L. & Hamid, J. S. Gene- or region-based analysis of genome-wide association studies. *Genet Epidemiol* **33** Suppl 1, S105–110, doi: 10.1002/gepi.20481 (2009).
6. Gusev, A. *et al.* Quantifying missing heritability at known GWAS loci. *PLoS Genet* **9**, e1003993, doi: 10.1371/journal.pgen.1003993 (2013).
7. Cheung, V. G. & Spielman, R. S. *Genetics of human gene expression: mapping DNA variants that influence gene expression*. **10**, 595–604, doi: 10.1038/nrg2630 (2009).
8. Consortium, T. E. P. *An integrated encyclopedia of DNA elements in the human genome*. **489**, 57–74, doi: 10.1038/nature11247 (2012).
9. Speed, D., Hemani, G., Johnson, M. R. & Balding, D. J. Improved heritability estimation from genome-wide SNPs. *Am J Hum Genet* **91**, 1011–1021, doi: 10.1016/j.ajhg.2012.10.010 (2012).
10. Loh, P.-R. *et al.* Contrasting regional architectures of schizophrenia and other complex diseases using fast variance components analysis. *bioRxiv* (2015).
11. Bulik-Sullivan, B. K. *et al.* LD Score regression distinguishes confounding from polygenicity in genome-wide association studies. *Nat Genet* **47**, 291–295, doi: 10.1038/ng.3211 (2015).
12. Bulik-Sullivan, B. Relationship between LD Score and Haseman-Elston Regression. *bioRxiv* (2015).
13. Ehret, G. B. *et al.* A multi-SNP locus-association method reveals a substantial fraction of the missing heritability. *Am J Hum Genet* **91**, 863–871, doi: 10.1016/j.ajhg.2012.09.013 (2012).
14. Palla, L. & Dudbridge, F. A Fast Method that Uses Polygenic Scores to Estimate the Variance Explained by Genome-wide Marker Panels and the Proportion of Variants Affecting a Trait. *Am J Hum Genet* **97**, 250–259, doi: 10.1016/j.ajhg.2015.06.005 (2015).
15. So, H. C., Li, M. & Sham, P. C. Uncovering the total heritability explained by all true susceptibility variants in a genome-wide association study. *Genet Epidemiol* **35**, 447–456, doi: 10.1002/gepi.20593 (2011).
16. Genomes Project, C. *et al.* A global reference for human genetic variation. *Nature* **526**, 68–74, doi: 10.1038/nature15393 (2015).
17. Sonnega, A. *et al.* Cohort Profile: the Health and Retirement Study (HRS). *Int J Epidemiol* **43**, 576–585, doi: 10.1093/ije/dyu067 (2014).
18. Yang, J., Lee, S. H., Goddard, M. E. & Visscher, P. M. GCTA: a tool for genome-wide complex trait analysis. *Am J Hum Genet* **88**, 76–82, doi: 10.1016/j.ajhg.2010.11.011 (2011).
19. Locke, A. E. *et al.* Genetic studies of body mass index yield new insights for obesity biology. *Nature* **518**, 197–206, doi: 10.1038/nature14177 (2015).
20. Wood, A. R. *et al.* Defining the role of common variation in the genomic and biological architecture of adult human height. *Nature genetics* **46**, 1173–1186, doi: 10.1038/ng.3097 (2014).
21. Sammalisto, S. *et al.* Genome-wide linkage screen for stature and body mass index in 3,032 families: evidence for sex- and population-specific genetic effects. *Eur J Hum Genet* **17**, 258–266, doi: 10.1038/ejhg.2008.152 (2009).
22. Perola, M. *et al.* Combined genome scans for body stature in 6,602 European twins: evidence for common Caucasian loci. *PLoS Genet* **3**, e97, doi: 10.1371/journal.pgen.0030097 (2007).
23. Vilhjalmsón, B. *et al.* Modeling Linkage Disequilibrium Increases Accuracy of Polygenic Risk Scores. *bioRxiv* (2015).
24. Ohtani, K. & Tanizaki, H. Exact Distributions of R<sup>2</sup> and Adjusted R<sup>2</sup> in a Linear Regression Model with Multivariate Error Terms. *Journal Of The Japan Statistical Society* **34**, 101–109, doi: 10.14490/jjss.34.101 (2004).
25. Yang, J. *et al.* Conditional and joint multiple-SNP analysis of GWAS summary statistics identifies additional variants influencing complex traits. *Nat Genet* **44**, 369–375, S361–363, doi: 10.1038/ng.2213 (2012).
26. Berndt, S. I. *et al.* Genome-wide meta-analysis identifies 11 new loci for anthropometric traits and provides insights into genetic architecture. *Nat Genet* **45**, 501–512, doi: 10.1038/ng.2606 (2013).
27. Lango Allen, H. *et al.* Hundreds of variants clustered in genomic loci and biological pathways affect human height. *Nature* **467**, 832–838, doi: 10.1038/nature09410 (2010).

## Author Contributions

G.P. designed the experiment; G.P. and W.Q.D. wrote the manuscript; S.M. analyzed the data and prepared tables and figures; All authors reviewed the manuscript.

## Additional Information

**Supplementary information** accompanies this paper at <http://www.nature.com/srep>

**Competing financial interests:** The authors declare no competing financial interests.

**How to cite this article:** Pare, G. *et al.* A method to estimate the contribution of regional genetic associations to complex traits from summary association statistics. *Sci. Rep.* **6**, 27644; doi: 10.1038/srep27644 (2016).



This work is licensed under a Creative Commons Attribution 4.0 International License. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in the credit line; if the material is not included under the Creative Commons license, users will need to obtain permission from the license holder to reproduce the material. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>