



THE UNIVERSITY *of* EDINBURGH

Edinburgh Research Explorer

SNP and Haplotype Regional Heritability Mapping (SNHap-RHM): Joint Mapping of Common and Rare Variation Affecting Complex Traits

Citation for published version:

Oppong, RF, Boutin, T, Campbell, A, McIntosh, AM, Porteous, DJ, Hayward, C, Haley, CS, Navarro, P & Knott, S 2022, 'SNP and Haplotype Regional Heritability Mapping (SNHap-RHM): Joint Mapping of Common and Rare Variation Affecting Complex Traits', *Frontiers in genetics*.
<https://doi.org/10.3389/fgene.2021.791712>

Digital Object Identifier (DOI):

<https://doi.org/10.3389/fgene.2021.791712>

Link:

[Link to publication record in Edinburgh Research Explorer](#)

Document Version:

Publisher's PDF, also known as Version of record

Published In:

Frontiers in genetics

Publisher Rights Statement:

This is an open-access article distributed under the terms of the Creative Commons Attribution License (CC BY).

General rights

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact openaccess@ed.ac.uk providing details, and we will remove access to the work immediately and investigate your claim.





SNP and Haplotype Regional Heritability Mapping (SNHap-RHM): Joint Mapping of Common and Rare Variation Affecting Complex Traits

Richard F. Opong^{1,2}, Thibaud Boutin³, Archie Campbell⁴, Andrew M. McIntosh⁵, David Porteous⁴, Caroline Hayward³, Chris S. Haley^{3,6}, Pau Navarro^{3*} and Sara Knott^{2*}

¹Longitudinal Studies Section, Translational Gerontology Branch, National Institute on Aging, National Institutes of Health, Baltimore, MD, United States, ²Institute of Evolutionary Biology, School of Biological Sciences, The University of Edinburgh, Edinburgh, United Kingdom, ³MRC Human Genetics Unit, Institute of Genetics and Cancer, The University of Edinburgh, Edinburgh, United Kingdom, ⁴Centre for Genomic and Experimental Medicine, Institute of Genetics and Cancer, The University of Edinburgh, Edinburgh, United Kingdom, ⁵Division of Psychiatry, The University of Edinburgh, Edinburgh, United Kingdom, ⁶The Roslin Institute and Royal (Dick) School of Veterinary Studies, The University of Edinburgh, Edinburgh, United Kingdom

OPEN ACCESS

Edited by:

Simon Charles Heath,
Center for Genomic Regulation (CRG),
Spain

Reviewed by:

David Duffy,
The University of Queensland,
Australia
Doug Speed,
Aarhus University, Denmark

*Correspondence:

Pau Navarro
Pau.Navarro@ed.ac.uk
Sara Knott
s.knott@ed.ac.uk

Specialty section:

This article was submitted to
Statistical Genetics and Methodology,
a section of the journal
Frontiers in Genetics

Received: 08 October 2021

Accepted: 14 December 2021

Published: 06 January 2022

Citation:

Opong RF, Boutin T, Campbell A,
McIntosh AM, Porteous D, Hayward C,
Haley CS, Navarro P and Knott S
(2022) SNP and Haplotype Regional
Heritability Mapping (SNHap-RHM):
Joint Mapping of Common and Rare
Variation Affecting Complex Traits.
Front. Genet. 12:791712.
doi: 10.3389/fgene.2021.791712

We describe a genome-wide analytical approach, SNP and Haplotype Regional Heritability Mapping (SNHap-RHM), that provides regional estimates of the heritability across locally defined regions in the genome. This approach utilises relationship matrices that are based on sharing of SNP and haplotype alleles at local haplotype blocks delimited by recombination boundaries in the genome. We implemented the approach on simulated data and show that the haplotype-based regional GRMs capture variation that is complementary to that captured by SNP-based regional GRMs, and thus justifying the fitting of the two GRMs jointly in a single analysis (SNHap-RHM). SNHap-RHM captures regions in the genome contributing to the phenotypic variation that existing genome-wide analysis methods may fail to capture. We further demonstrate that there are real benefits to be gained from this approach by applying it to real data from about 20,000 individuals from the Generation Scotland: Scottish Family Health Study. We analysed height and major depressive disorder (MDD). We identified seven genomic regions that are genome-wide significant for height, and three regions significant at a suggestive threshold (p -value $< 1 \times 10^{-5}$) for MDD. These significant regions have genes mapped to within 400 kb of them. The genes mapped for height have been reported to be associated with height in humans. Similarly, those mapped for MDD have been reported to be associated with major depressive disorder and other psychiatry phenotypes. The results show that SNHap-RHM presents an exciting new opportunity to analyse complex traits by allowing the joint mapping of novel genomic regions tagged by either SNPs or haplotypes, potentially leading to the recovery of some of the “missing” heritability.

Keywords: MDD, height, haplotypes, regional heritability mapping, missing heritability, rare variation, genome-wide analysis

1 INTRODUCTION

Estimates of the genetic component of complex trait variation using genotyped SNPs led to the conclusion that a proportion of the heritability of complex traits is still unexplained or “missing” (Maher, 2008; Manolio et al., 2009). Full sequence data will contain all the variants that account for all the heritability of complex traits (Wainschtein et al., 2019). Moreover, some of these true causal variants may be rare (Pritchard, 2001) and therefore may be in incomplete linkage disequilibrium (LD) with genotyped SNPs (Yang et al., 2010). Thus, some of the “missing” heritability may be “hidden” in rare variants whose effects are difficult to capture because of lack of statistical power. There is, therefore, some benefit to be gained in terms of improving the heritability estimates and uncovering gene variants involved in the control of traits by fitting genome-wide analytical models that adequately capture the combined effects of rare genetic variants (Cirulli and Goldstein, 2010; Gonzalez-Recio et al., 2015).

In light of this, we proposed a genome-wide analytical approach that draws its theoretical basis from the genome-based restricted maximum likelihood (GREML) approach (Maher, 2008; Manolio et al., 2009; Clarke and Cooper, 2010; Yang et al., 2011; Speed et al., 2012) which utilises both local and genome-wide relationship matrices to provide regional estimates of the heritability across locally defined regions in the genome (Nagamine et al., 2012; Uemoto et al., 2013). This regional heritability analysis can capture the combined effect of SNPs in a region, and thus small effect variants may be detectable. However, the analysis only captures effects associated with common SNPs present on genotyping chips.

Haplotypes may provide a better strategy to capture genomic relationships amongst individuals in the presence of causal rare variants. Although rare variants are not in LD with genotyped variants and thus are difficult to capture in conventional GWAS, these rare variants may be in LD with some haplotypes and thus can be captured using haplotype methods. Compared with genotyped SNPs, capturing haplotype effects may offer an advantage because haplotypes can be functional units (Vormfelde and Brockmüller, 2007). Therefore, haplotype effects may reflect the combined effects of closely linked cis-acting causal variants (Balding, 2006) and using haplotypes could provide real benefit over SNPs in recovering some of the “missing” heritability and identifying novel trait-associated variants. Therefore, we extended the SNP-based regional heritability analysis further by incorporating haplotypes in addition to SNPs in the calculation of the regional GRMs used in the analysis (Shirali et al., 2018). This approach includes two regional GRMs and divides the genome into windows based on local haplotype blocks delimited by recombination boundaries.

This paper further explores the properties of both the SNP-based and the haplotype-based regional heritability mapping (SNP-RHM and Hap-RHM respectively). We hypothesise and show by simulation that the Hap-RHM complements existing SNP-RHM analytical approaches by capturing regional effects in the genome that existing SNP-based methods fail to capture. This

leads us to propose a mapping strategy that jointly utilises SNP and haplotype GRMs in a single analysis called SNHap-RHM. We then confirm the utility of this approach by applying it to real data obtained from about 20,000 individuals from the Generation Scotland: Scottish Family Health Study (GS: SFHS) (Smith et al., 2012). We analysed two phenotypes: height and major depressive disorder (MDD). The aim was to uncover novel genetic loci that may affect these traits and improve the estimates of the genetic components of the variation in these traits.

2 MATERIALS AND METHODS

2.1 Materials

The data used in this study is from the Generation Scotland: Scottish Family Health Study (GS: SFHS), comprising of 23,960 participants recruited from Scotland (Smith et al., 2006; Smith et al., 2012). The DNA from about 20,032 of the participants had been genotyped using the Illumina HumanOmniExpressExome8v1-2_A chip (~700 K genome-wide SNP chip) (Smith et al., 2012). GRCh37 was used throughout. Quality control excluded SNPs and individuals with a call rate less than 98%, SNPs with minor allele frequency (MAF) less than 1% and SNPs that were out of Hardy-Weinberg equilibrium (p -value < 0.000001). A total of 555,091 autosomal SNPs passed quality control for downstream analysis. Ethical approval for the GS: SFHS study was obtained from the Tayside Committee on Medical Research Ethics (on behalf of the National Health Service).

2.2 Methods

We have shown previously that regional GREML analysis (Regional Heritability Mapping or RHM) using fixed region sizes in the genome is a suitable mapping method for finding local genetic effects (Nagamine et al., 2012). The conventional RHM model fits two genomic relationship matrices (GRMs) in the analyses to map genetic loci that affect trait variation: a local GRM (rGRM) calculated using SNPs located in the region and a genome wide GRM (gwGRM) calculated from SNPs outside the region. We have since extended this conventional regional heritability analysis to incorporate haplotypes in the calculation of the local GRM and have successfully implemented this in a simulation study (Shirali et al., 2018). This study, like our previous (Shirali et al., 2018), utilises a regional heritability model that breaks the genome into naturally defined regions by delimiting them by recombination hotspots. Two types of regional heritability models are then fitted in turn to the phenotypes. One model (SNP-RHM) uses SNPs to estimate local genetic relationships between study individuals, and the other model (Hap-RHM) estimates local genetic relationships amongst individuals using haplotypes. We go a step further in this study to perform a regional heritability analysis that jointly fits the SNP and the haplotype GRM in an approach that we termed SNP and Haplotype Regional Heritability Mapping (SNHap-RHM). An overview of SNHap-RHM is shown in **Figure 1**.

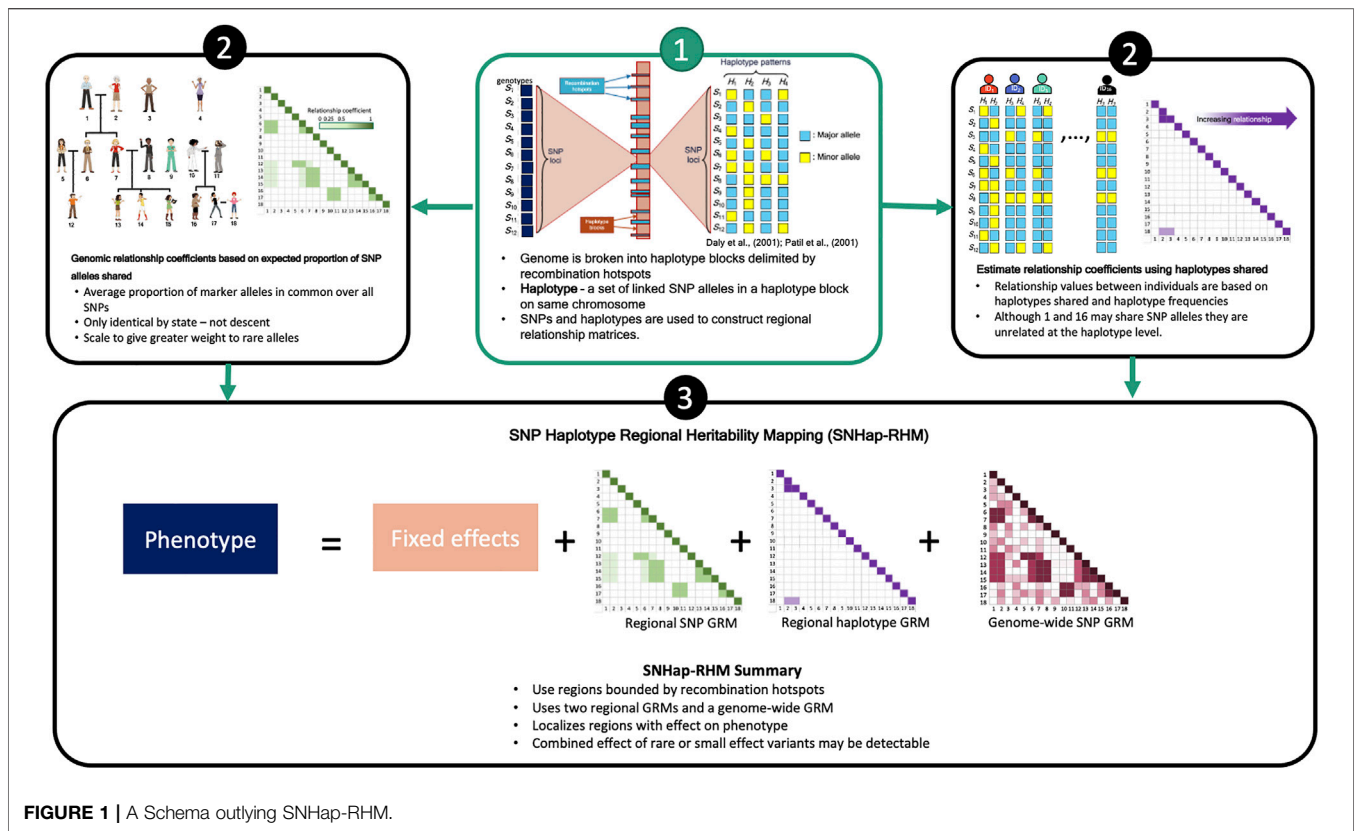


FIGURE 1 | A Schema outlying SNHap-RHM.

2.2.1 The General Statistical Setting of a Regional Heritability Analysis

Consider a vector y of phenotype values with length n , the linear mixed-effects model for fitting the effects of genomic region i and background polygenic markers is given as

$$y = X\beta + W_i u_i + Z u_b + e$$

where y is a vector of phenotypes, X is a design matrix of fixed effects, and β is a vector of fixed effects, W_i is a design matrix relating phenotype measures to genetic markers in region i and u_i is a vector of random genetic effects due to region i assumed to be multivariate normal, $MVN(0, \sigma_{u_i}^2 L_{u_i})$. L_{u_i} is a relationship matrix calculated using markers (SNPs or haplotypes) in region i : calculated in the subsequent sections as G for the SNP and H for the haplotype-based models. Z is a design matrix for background polygenic effects of markers outside the region i and u_b is a vector of random polygenic effect of genetic markers excluded from region i , assumed to be multivariate normal, $MVN(0, \sigma_{u_b}^2 B_{u_b})$. B_{u_b} is a relationship matrix calculated using the markers outside the region i : calculated in the subsequent section in the same way as G . And e is a vector of residual effects assumed to be multivariate normal, $MVN(0, \sigma_e^2 I)$. I is an identity matrix.

Under the model, the vector of phenotypes y is assumed to be normally distributed, $N(X\beta, V)$ where the variance is:

$$V = \sigma_{u_i}^2 W_i L_{u_i} W_i^T + \sigma_{u_b}^2 Z B_{u_b} Z^T + \sigma_e^2 I$$

2.2.1.1 SNP-RHM: SNP-Based Regional Heritability Model

A SNP-based regional heritability analysis was first reported by Nagamine et al. (2012). The regional heritability analysis approach we employ here differs from the analysis done by Nagamine et al. (2012) in the way the regions are defined. That analysis defined local regions by breaking the genome into smaller user-defined windows of r SNPs, which overlapped by s SNPs. Here, however, we define regions based on recombination boundaries in the genome.

The regional heritability model fits two genetic relationship matrices (GRMs): one local GRM for the region and a whole-genome GRM for the remaining SNPs in the genome that are outside the region. The GRMs are genomic relatedness matrices calculated as the weighted proportion of the local or genome-wide autosomal SNPs shared identity by state (IBS) between pairs of individuals. The SNP IBS matrices are calculated as follows, following the second scaling factor proposed by VanRaden (2008)

$$G = \frac{MM'}{m}$$

where m is the total number of r local or b background autosomal SNPs, and M is a matrix of genotype codes for the sampled individuals centred by loci means and normalised by the standard deviation of each locus. M is calculated as follows for individual i at locus j

$$M_{ij} = \frac{(x_{ij} - 2p_j)}{\sqrt{2p_j(1 - p_j)}}$$

where x_{ij} is the genotype code at locus j for individual i and takes the values 0, 1 and 2 for AA, Aa and aa genotypes respectively, p_j is the frequency of allele “a” at locus j . The SNP-based relationship for individuals i and k is therefore calculated as follows

$$G_{ik} = \frac{1}{m} \times \sum_{j=1}^m \frac{(x_{ij} - 2p_j)(x_{kj} - 2p_j)}{2p_j(1 - p_j)}$$

2.2.1.2 Hap-RHM: Haplotype-Based Regional Heritability Model

The haplotype-based regional heritability model follows theoretically from the SNP-based analysis and utilises haplotypes instead of SNPs as the genetic markers for the regional analysis. The analysis fits two GRMs, a haplotype-based regional GRM and a SNP-based background genome-wide GRM. The haplotype-based GRM is similar to the SNP-based GRM defined in the previous section. For a locally defined region (haplotype block) containing h haplotype variants, the haplotype-based kinship for individuals i and k is calculated as follows

$$H_{ik} = \frac{1}{h} \times \sum_{j=1}^h \frac{(d_{ij} - 2p_j)(d_{kj} - 2p_j)}{2p_j(1 - p_j)}$$

where d_{ij} is the diplotype code (coded as the number of copies of haplotype j) for individual i and takes the values 0, 1, and 2 for the $h_t h_t$, $h_t h_j$, $h_j h_j$ diplotypes respectively where haplotype t is any haplotype other than haplotype j , i.e., $t \neq j$, p_j is the haplotype frequency for haplotype j .

2.2.2 Simulation Study

Five phenotypes were simulated using available genotypic information of 20,032 individuals from the Generation Scotland: Scottish Family Health Study (Smith et al., 2012). The five phenotypes were simulated to have a total variance of 1. This total is composed of 0.6 environmental (residual) variance and genetic variance of 0.4. The genetic variance was partitioned into two components, a polygenic variance of 0.3 and a total QTL variance of 0.1 (20 QTLs, each explaining a variance of 0.005). A common polygenic variance was simulated for all five phenotypes from 20,000 markers randomly selected across the genome. The polygenic variance was simulated to be normally distributed with zero mean and variance of 0.3.

Phasing of the GS: SFHS data was done using SHAPEIT2 (Delaneau et al., 2013). Best guess haplotypes were used. Haplotypes variants within blocks were determined using the phased data. For each phenotype, 20 regions (haplotype blocks) were randomly selected, one on each autosome (except chromosomes 6 and 8 because of the unusually high LD in the MHC regions on chromosome 6 and a large inversion on chromosome 8 (Amador et al., 2015)), to simulate quantitative

trait loci (QTL). This gave a total of 20 QTLs for each phenotype. The haplotype blocks were delimited by natural boundaries: recombination hotspots where the estimated recombination frequency exceeds ten centiMorgans per Megabase (10 cM/Mb) with the estimated recombination frequency between boundaries being less than ten centiMorgans per Megabase (10 cM/Mb) based on the Genome Reference Consortium Human Build 37 (International Human Genome Sequencing Consortium, 2004). This recombination threshold resulted in a total of 48,772 regions across the genome. The number and type of marker used to simulate the QTL are what defined the five phenotypes. The five phenotypes are, a 1-SNP QTL within the haplotype block, a multiple-SNP (5 SNPs) QTL within the haplotype block, two types of 1-haplotype QTL within the haplotype block (taking either a common or a rare haplotype as causal) and multiple (5) haplotype QTL within the haplotype block. Details of these phenotypes are described below.

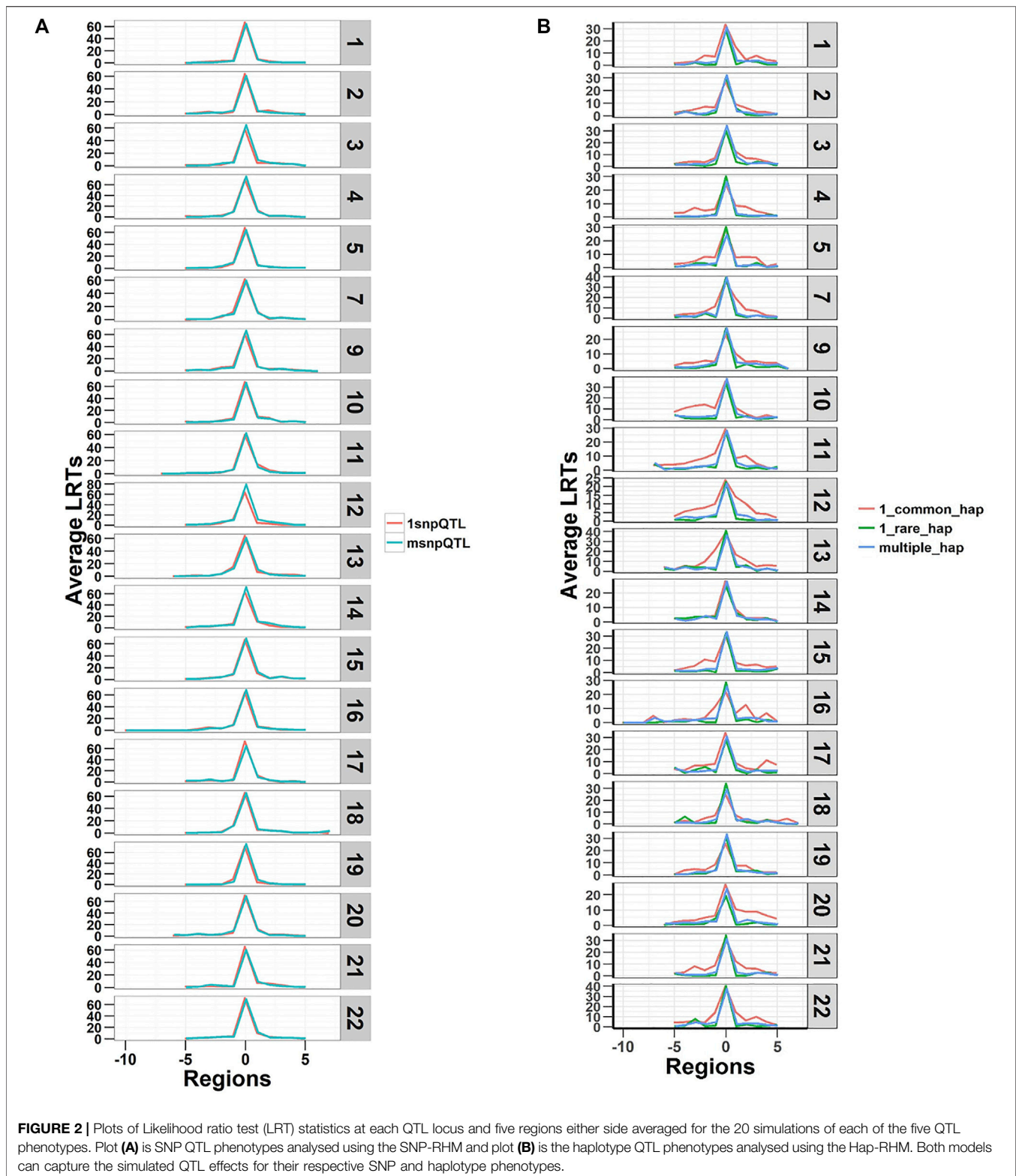
For the haplotype QTL phenotypes, a haplotype block is treated as a single genetic locus having multiple alleles. Each haplotype variant within a block is considered as an allele of that locus. Each study individual will carry two alleles, or have a diplotype, for each locus or haplotype block. The genotype data used to simulate the phenotypes were phased using SHAPEIT2 (Delaneau et al., 2013) to produce the haplotypes for study individuals. The multiple haplotype QTL phenotypes were simulated by randomly sampling two rare haplotypes and three common haplotypes within each haplotype block to give five haplotypes per block. The two types of 1-haplotype QTL phenotypes were simulated by randomly sampling a rare haplotype per haplotype block for one type and for the other type a common haplotype was randomly sampled within each haplotype block. **Supplementary Figure S10** gives an indication of the frequencies for the rare (0.00002–0.036) and common haplotype (0.008–0.906) randomly sampled to simulate the phenotypes. There is a slight overlap between the frequencies for rare and common haplotypes because the regions had already been randomly selected before proceeding to randomly select rare and common haplotypes in those regions. Which means what is rare in one region may be common in another.

The individual marker contribution to the polygenic effect and the QTL effects were calculated as follows

$$\sigma_j^2 = 2p_j(1 - p_j)g_j^2$$

$$g_j = \sqrt{\frac{\sigma_j^2}{2p_j(1 - p_j)}}$$

where σ_j^2 is the contribution of a marker to the QTL or polygenic variance, g_j is the effect of a SNP j or haplotype j randomly sampled to have polygenic or QTL effect, p_j is the frequency of haplotype j or the effect allele of the SNP j . For the single marker QTL phenotypes, each QTL explained a variance of 0.005. For the multiple marker QTL phenotypes, each causal variant explained the same variance, with the effects scaled to account for LD in the region so each QTL locus explained a variance of 0.005. For the multiple haplotype QTL effects, the haplotype effects were scaled



relative to the inverse of their frequency to give a total variance explained by the region of 0.005.

Common environmental effects were randomly sampled for the five phenotypes from a normal distribution $N(0, \sigma_e^2)$

where σ_e^2 is 0.6. This, together with a genetic variance of 0.4, gave a total variance of 1 for each phenotype. The final simulated phenotype for an individual i was then calculated as follows

$$y(\text{single markers per QTL region})_i = \sum_{j=1}^{20000} x_{ij}g_j + \sum_{j=1}^{20} x_{ij}g_j + e_i,$$

$$y(\text{multiple markers per QTL region})_i = \sum_{j=1}^{20000} x_{ij}g_j + \sum_{l=1}^{20} \sum_{j=1}^5 x_{ij}g_j + e_i,$$

where x_{ij} is the number of copies of the effect allele of SNP j for individual i (for haplotypes, this is defined as d_{ij} ; the number of copies of haplotype j for individual i) and g_j is the effect of haplotype j or SNP j . Twenty replicates were analysed for each of the five phenotypes with a different set of QTL markers sampled for each replicate.

2.2.2.1 Analysis of Simulated Data

In this simulation study, the five simulated phenotypes were analysed using the two models, the SNP-based regional heritability model (SNP-RHM for the SNP QTL phenotypes) and the haplotype-based regional heritability model (Hap-RHM for the haplotype QTL phenotypes). To test the analytical models' specificity, we applied Hap-RHM to SNP QTL phenotypes and SNP-RHM to the haplotype QTL phenotypes. We also performed a Hap-RHM analysis in which the units of analysis in the haplotype blocks were restricted to regions of 20 or fewer SNPs per haplotype block. This was because we observed that longer haplotype blocks had many SNPs (and hence many, many haplotypes, up to 14,000 in some blocks), and this impacted the estimation of the simulated regional effect. The hybrid Hap-RHM, therefore, investigates whether the regional effect is well captured by the haplotype-based model when shorter haplotypes are used.

We estimated the regional genetic variance and polygenic variance using restricted maximum likelihood (REML). For each simulated phenotype, we analysed 220 regions in total to map the 20 simulated QTLs. This involved analysing the region containing the QTL and ten adjacent regions (five in either direction). In this way, we limit the analysis to the regions in the genome with simulated effects, thereby reducing computation time considerably. Also, by analysing neighbouring regions, we are able to explore the precision of estimates of the location of regional effects. We assessed the significance of a region using the Likelihood Ratio Test (LRT). The genome-wide significance threshold was calculated to be $LRT = 23.9$ ($p\text{-value} < 1.02 \times 10^{-6}$) using a Bonferroni correction for testing 48,772 regions.

Also, we selected one replicate for each simulated phenotype and performed SNHap-RHM (SNP and Haplotype Regional Heritability Mapping), a regional heritability analysis that jointly fitted the SNP and the haplotype GRM.

2.2.3 SNHap-RHM of MDD and Height

MDD status for GS: SFHS participants was assigned following an initial mental health screening questionnaire with the questions: "Have you ever seen anybody for emotional or psychiatric problems?" or "Was there ever a time when you,

or someone else, thought you should see someone because of the way you were feeling or acting?" Participants who answered yes to one or both of the screening questions were further interviewed by the Structured Clinical Interview for DSM-IV (SCID) (First et al., 2002). A total of 18,725 participants (2,603 MDD cases and 16,122 controls) were retained for analysis for MDD. A total of 19,944 participants from the GS: SFHS were analysed for height.

SNHap-RHM fits jointly, the two types of regional GRMs, SNP-based and haplotype-based, in the analysis of phenotypes (Figure 1). We pre-corrected the phenotypes with the whole-genome GRM before performing SNHap-RHM to speed up the GREML analysis of each block. This pre-correction has previously been shown to speed the regional heritability analysis by Shirali et al. (2018). This is a leave-one-chromosome-out step (Yang et al., 2014), which involved 22 separate GREML analyses each fitting a whole-genome GRM that excluded SNPs from one chromosome. The residuals from the pre-correction step were then used in the SNHap-RHM analysis. The models adjusted for sex, age, age², and the first 20 principal components calculated from the study participants' genomic relationship matrix (calculated using 555,091 autosomal SNPs).

The significance of a region was tested with a likelihood ratio test (LRT) with two degrees of freedom which compared a model with three variance components fitted (the two regional variances together with the residual variance) against a model with only the residual variance component fitted. The individual regional variance components in all regions were subsequently tested with an LRT with one degree of freedom which compared a model with three variance components fitted against a model with two variance components fitted (one regional variance component dropped from the model). We assumed the appropriate null distribution that results from testing on the boundary of the parameter space and therefore calculated the p -values as $0.5 \times$ the p -value of a chi-squared distribution with one degree of freedom for the one degree of freedom test and as $0.5 \times$ the p -value of a chi-squared distribution with one degree of freedom plus $0.25 \times$ the p -value of a chi-squared distribution with two degrees of freedom for the two degrees of freedom test.

The p -values obtained from the LRTs were used to generate genome-wide association plots for each phenotype (equivalent to GWAS Manhattan plots). The genome-wide significance threshold was calculated to be $LRT = 23.9$ ($p\text{-value} < 1.02 \times 10^{-6}$) using a Bonferroni correction for testing 48,772 regions. The suggestive significance threshold of a region was set at an $LRT = 19.5$ ($p\text{-value} < 1.02 \times 10^{-5}$).

3 RESULTS

3.1 Simulation Study: SNP-RHM, Hap-RHM and SNHap-RHM

We performed a regional heritability analysis that fits two GRMs (one for the region and one for the rest of the genome) per region across multiple genomic regions delimited by recombination

hotspots (where the estimated recombination frequency exceeds ten centiMorgans per Megabase (10 cM/Mb)). This recombination threshold resulted in a total of 48,772 regions across the genome. We tested two types of regional heritability models, SNP-RHM and Hap-RHM, on 20 replicates of five simulated phenotypes. In SNP-RHM, the regional matrix is derived from SNP genotypes whereas in Hap-RHM the regional matrix is derived from haplotypes. The phenotypes were simulated to be determined by 20 regional QTL effects and genome-wide polygenic effects. The regional QTL effects of the five phenotypes were simulated using SNPs as causal variants for two of them and haplotypes for the remaining three as described in the methods section.

A likelihood ratio test (LRT) was used to test the null hypothesis, H_0 : that the genetic variance explained by the region is not significant, against the alternative hypothesis, H_1 : that the region accounts for a significant proportion of the phenotypic variance. A large LRT statistic is evidence against the null hypothesis, and therefore means the region explains a significant proportion of the phenotypic variance.

The LRTs averaged over the 20 replicates of the five phenotypes are shown in **Figure 2**. The figure shows plots of average LRT for the QTL regions and ten adjacent regions (five to each side). The results show that both models detected the simulated regional effects at the genome-wide significance level (LRT = 23.9) (p -value $< 1.02 \times 10^{-6}$, Bonferroni correction for testing 48,772 regions) and can capture true causal loci in traits with different genetic architectures. The LRTs were higher on average for the SNP-based model (SNP-RHM) than the haplotype-based model (Hap-RHM). This could be because for Hap-RHM, the genome-wide GRM which is a SNP-based GRM does not tag any of the background haplotype effects that are outside any one particular region being analysed, and thus the residual variance may be inflated by the other haplotype QTLs which downwardly impact the LRTs.

We provide further investigation of the results from the simulation in the supplementary material (**Supplementary Text**). For both analysis models, we have presented detailed results of the relationships between the LRT statistics, region size, variance estimates and allele frequencies (**Supplementary Figures S3–S10**). We observed that the longer haplotype blocks had many SNPs (and hence many, many haplotypes, up to 14,000 in some blocks), and this impacted the estimation of the simulated regional variance (**Supplementary Figure S8**). We, therefore, performed a hybrid-Hap-RHM analysis that restricted the natural haplotype block sizes to 20 or fewer SNPs per haplotype block. This hybrid-Hap-RHM was to investigate whether the regional variance is well captured by Hap-RHM when shorter haplotypes are used. The hybrid-Hap-RHM underestimated the regional variance for larger regions but did not offer any discernible improvement in the LRT statistics (**Supplementary Figure S9**). The relationship between region size and estimated variance was different between the Hap-RHM and hybrid-Hap-RHM, while we observed a similar relationship between LRTs and the region size.

Both SNP-RHM and Hap-RHM fail to capture the simulated regional effects when the simulated phenotype has a genetic architecture that does not match the analysis model, i.e., SNP or haplotype (**Figure 3; Supplementary Figure S1**). These figures show the results for the situation where the SNP QTL phenotypes were analysed with the haplotype-based model (Hap-RHM) and the haplotype QTL phenotypes were analysed with the SNP-based model (SNP-RHM). Both models fail to detect the simulated effects in such situations, therefore, showing that the models complement each other since they capture effects due to different types of genetic variants (i.e., tagged by SNPs or haplotypes).

To confirm that two models are complementary and independent of each other, we implemented SNHap-RHM that fits the regional SNP and haplotype GRMs jointly, on a replicate of each of the five simulated phenotypes. The significance of regional effects was tested with an LRT with two degrees of freedom. The results are shown in **Figure 4** and confirm that the two models are complementary since even when we fitted jointly the two regional matrices (SNP and Haplotype-based), we can still capture the simulated effects.

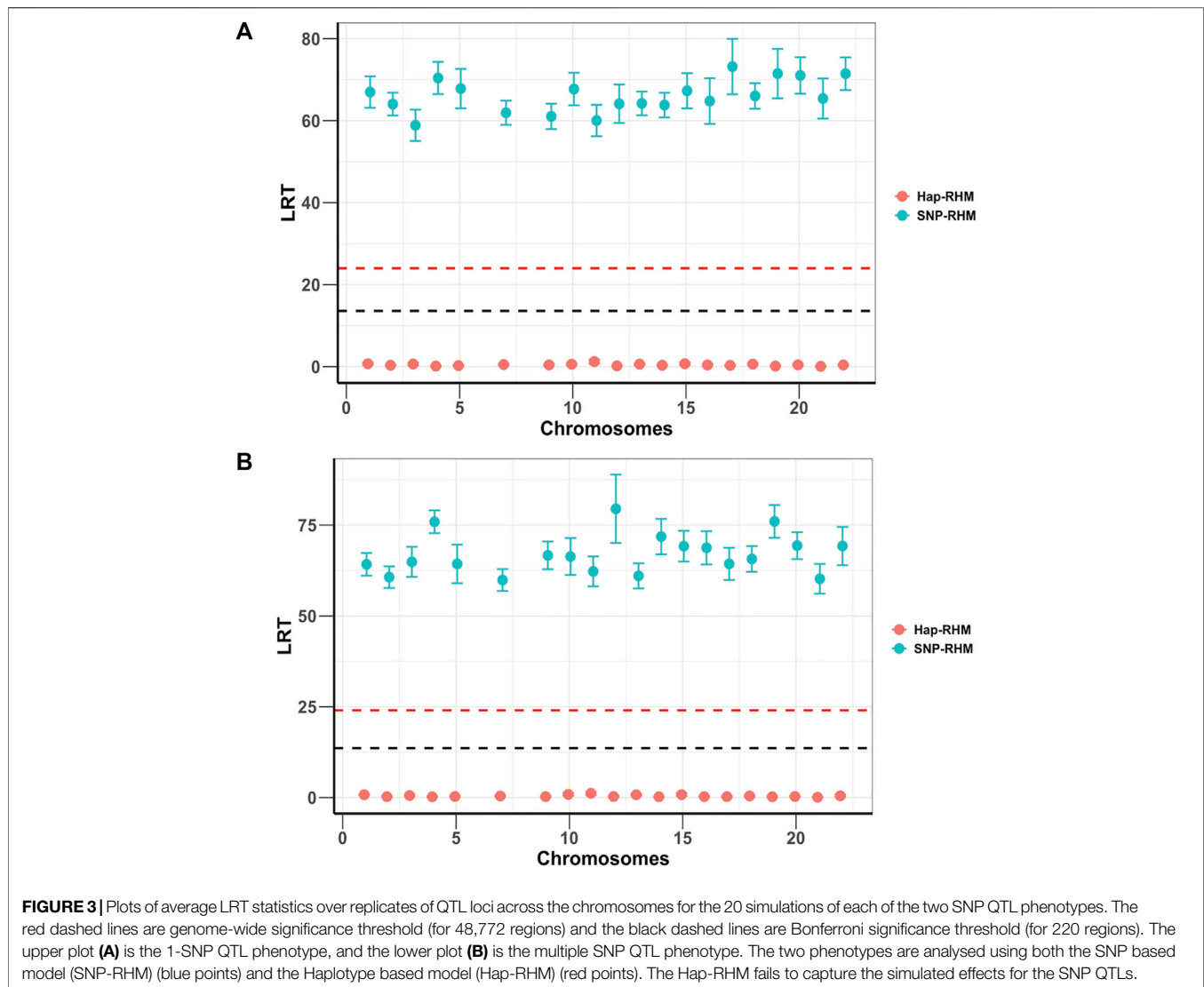
3.2 SNHap-RHM Analysis of Height and Major Depressive Disorder

The heritability estimates for height and MDD in the GS: SFHS dataset, calculated using the whole-genome GRM, were 81.4% (0.92) and 13.8% (1.35) respectively. There were no overlaps between regions identified as significant (tested with an LRT with one degree of freedom) by the haplotype and SNP-based models for either of the two traits (**Supplementary Figure S2**). This reaffirms our hypothesis tested by simulation that the Hap-RHM is complementary to SNP-RHM in mapping associated genomic loci.

The regional heritability results for height and MDD are presented as plots of minus-Log₁₀ of the LRT p -values (**Figures 5, 6**). The plots for the SNHap-RHM, SNP-RHM and Hap-RHM analyses are shown.

The results for height show that nine regions passed the Bonferroni-corrected genome-wide significance threshold in the analysis using SNP-RHM. No region was genome-wide significant for height when analysed with Hap-RHM. Furthermore, seven of the nine associated regions still come up as genome-wide significant when SNPs and haplotypes in those regions are analysed jointly using SNHap-RHM. There are GWAS reported genes that lie in or are within 400 kb of these regions (**Supplementary Table S1**).

For MDD, no region passed the Bonferroni-corrected genome-wide significance threshold for the analysis done with the SNP-based and haplotype-based regional heritability models (**Figure 6**). Three regions passed the suggestive significance threshold at p -value $< 1 \times 10^{-5}$ for Hap-RHM analysis of MDD. A further nine regions were significant at p -value $< 5 \times 10^{-5}$ for the haplotype-based analysis, and one region for the SNP-based analysis (**Supplementary Table S2**). **Figure 6** shows that when the two local GRMs are fitted jointly using SNHap-RHM, the genomic regions associated with MDD can still be mapped. The associated regions mapped by the haplotype-based



model for MDD contain genes reported by GWAS to be associated with several psychiatric phenotypes (**Figure 6**; **Supplementary Table S2**). The most strongly associated region was within 400 kb of the *DCC* gene. This gene is part of the *NETRIN1* pathway, which has been reported to be associated with major depressive disorder in two GWAS samples (GS: SFHS and Psychiatric Genomics Consortium) (Zeng et al., 2017). Zeng et al. (2017) used a SNP-RHM guided by pathway analysis (to first uncover pathway association and then localise *DCC* within the pathway) to show the *DCC* association with major depressive disorder. The second most strongly associated region was on chromosome 8, and this region had no gene mapped to it.

A linear mixed effects model was used to test for association of the SNPs within the suggestive significant region identified by the haplotype-based model on chromosome 3 for MDD. The model tested for association of SNPs by fitting their allelic dosages individually in a regression model and fitting a GRM to account for relatedness of individuals. The region on

chromosome 3 was chosen in this example because there is a psychiatric phenotype associated gene, *MYRIP* (Luciano et al., 2011), mapped to it, unlike the *DCC* region which has the gene outside the region. The results are shown in **Table 1**. Five SNPs within this region are nominally significant at p -value < 0.05 . Four out of these five SNPs confer about 2% increased risk of the disease each. These four SNPs lie within the *MYRIP* gene sequence. The *MYRIP* gene is expressed in the brain (Ganat et al., 2012). A SNP (rs9985399) in this gene is reported to be associated with brain processing speed in the Lothian birth cohort (Luciano et al., 2011). Brain processing speed is an important cognitive function that is compromised in psychiatric illness such as schizophrenia and depression, and old age. Also, a SNP (rs6599077) in the *MYRIP* gene region is associated with sleep duration (Gottlieb et al., 2007). Sleep durations outside the normal range (both short sleep and long sleep) is significantly associated with increased risk of depression (Roberts and Duong, 2014; Watson et al., 2014; Zhai et al., 2015; Mohan et al., 2017). The *MYRIP* gene is also reported to have a role in insulin

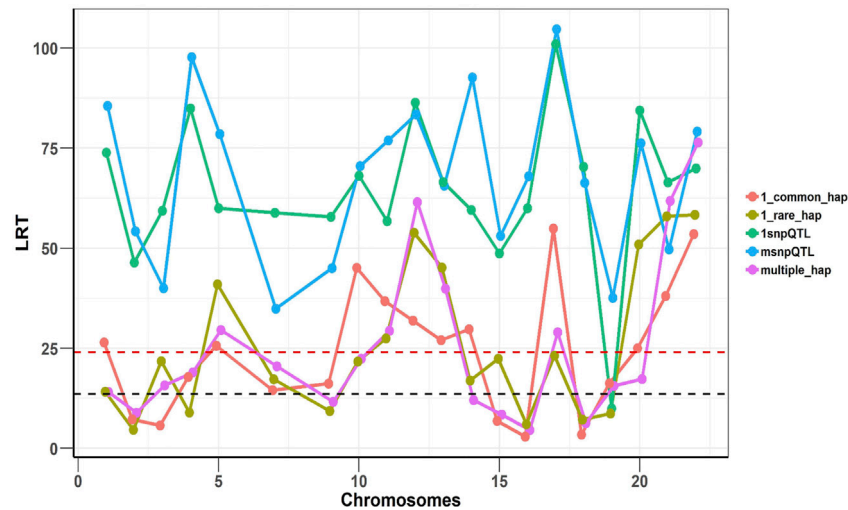


FIGURE 4 | Joint analysis of the SNP and haplotype phenotypes using SNHap-RHM. The plot is an analysis of one replicate of each of the simulated phenotypes. The LRT statistics are plotted over QTL loci across the chromosomes. The red dashed lines are genome-wide significance threshold (for 48,772 regions) and the black dashed lines are Bonferroni significance threshold (for 220 regions).

secretion (Waselle et al., 2003) and low insulin levels have been linked to depression (Pearson et al., 2010; Greenwood et al., 2015; Webb et al., 2017).

3.2.1 Comparison With Published GWAS SNPs

For both traits, the SNPs in the regions that were significant at p -value $< 5 \times 10^{-5}$ were compared to SNPs reported in the GWAS catalogue (MacArthur et al., 2017) to be significant for the two traits. The GWAS catalogue was accessed on the January 15, 2021. The results are presented in **Table 2**. The SNP-based and haplotype-based models identified 1,380 and 45 SNPs respectively for height, and 78 and 495 SNPs respectively for MDD taking all SNPs within haplotype blocks significant at p -value $< 5 \times 10^{-5}$. Out of the 1,380 SNPs identified for height by the SNP-based model, 57 SNPs spanning 20 haplotype regions were in common with published GWAS results for height. The number of SNPs found in common with published GWAS results are modest, and this could be because of the differences in genotyping chips used in this study and the published studies. Which means if we were to consider proxies of our SNPs ($LD > 0.8$) in the comparison, the numbers might increase. Also, our sample size compared to most of those published GWASs is quite small which means that we might not have enough power to detect all associations.

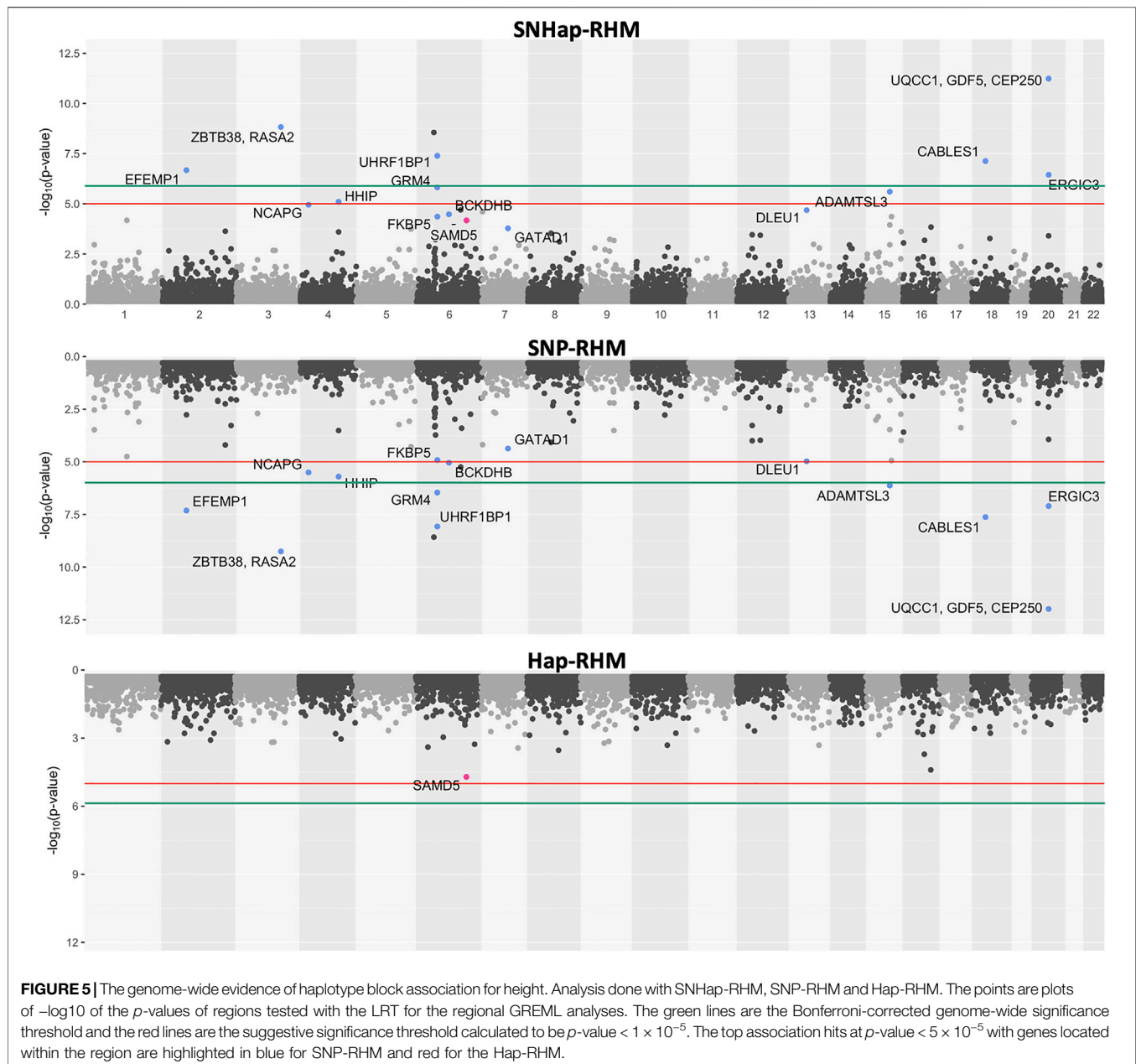
4 DISCUSSION

We have proposed and implemented a genome-wide analytical method that analyses genomic regions using a regional heritability model (Nagamine et al., 2012). We have since extended this method to include haplotypes by fitting a regional haplotype-based GRM (Hap-RHM) and redefined

genomic regions in our analysis to be delimited by recombination hotspots generated using HapMap Phase II (Frazer et al., 2007; Shirali et al., 2018). In this study, we build on our previous regional heritability methods by exploring the properties of the SNP and haplotype-based regional heritability mapping models by simulation and demonstrate that the two variance components fitted are largely independent of each other (**Supplementary Figure S2**). The novelty in this study is that we show that the two regional matrices fitted in SNP-RHM and Hap-RHM capture two different kinds of effects in terms of genetic architecture, and thus the two variance components can be fitted jointly (by fitting the SNP and haplotype regional matrices together) in a joint marker regional heritability mapping procedure that we call SNHap-RHM.

We hypothesised that the Hap-RHM would complement the SNP-RHM. We investigated this hypothesis in a simulation study in which we simulated 20 replicates each of two types of SNP QTL phenotypes and three types of haplotype QTL phenotypes. The results show that the two heritability models can capture the effects of causal variants within genomic loci associated with the phenotype analysed. The results also show that the two models are specific about the type of causal effect they can capture, therefore, providing support for the hypothesis that haplotype-based regional heritability models will complement SNP-based regional heritability models. We provide further support for this hypothesis by fitting the two GRMs jointly and showing (using an LRT with two degrees of freedom) that we can still capture the simulated effects and real effects from real data.

We applied SNHap-RHM to height and MDD phenotypes from the Generation Scotland: Scottish Family Health Study. Again, we draw comparisons between the effects captured by the SNP-RHM and the Hap-RHM. The SNP-RHM identified more Bonferroni-corrected genome-wide (GW) significant regions (p -value $< 1.02 \times 10^{-6}$) for height compared to MDD. Fifty-



seven of the SNPs identified for height by the SNP-RHM have been reported by other studies to be associated with height. These SNPs spanned 20 genomic regions in the GS: SFHS cohort. Height is a highly polygenic trait with many common genetic variants accounting for most of the additive genetic variation (Yang et al., 2015). These common genetic variants may be in LD with genotyped SNPs on SNP chips (these chips are disproportionately enriched for common SNPs). Therefore, the SNP-based regional heritability model is better suited for capturing SNP loci in height compared to MDD.

MDD is a very heterogeneous phenotype, and thus every MDD case could have a set of genetic and non-genetic risk factors exclusive to them (Levinson et al., 2014). These unique genetic

risk factors will mean that a lot of the genetic variants driving the disease will be rare at the population level. Three genomic regions were identified for MDD by the haplotype-based regional heritability model at the suggestive level, p -value $< 1 \times 10^{-5}$. The Hap-RHM works well for MDD because MDD is believed to be driven by rare genetic variants, and the model can capture rare genetic variants. The haplotype model can capture rare variants because of the LD between rare variants (both typed and untyped) and the flanking variants that aggregate to form the haplotypes within the genomic regions. There were no overlaps between regions identified by the Hap-RHM and SNP-RHM for each trait, which again supports the hypothesis that the two models complement each other in mapping associated loci.

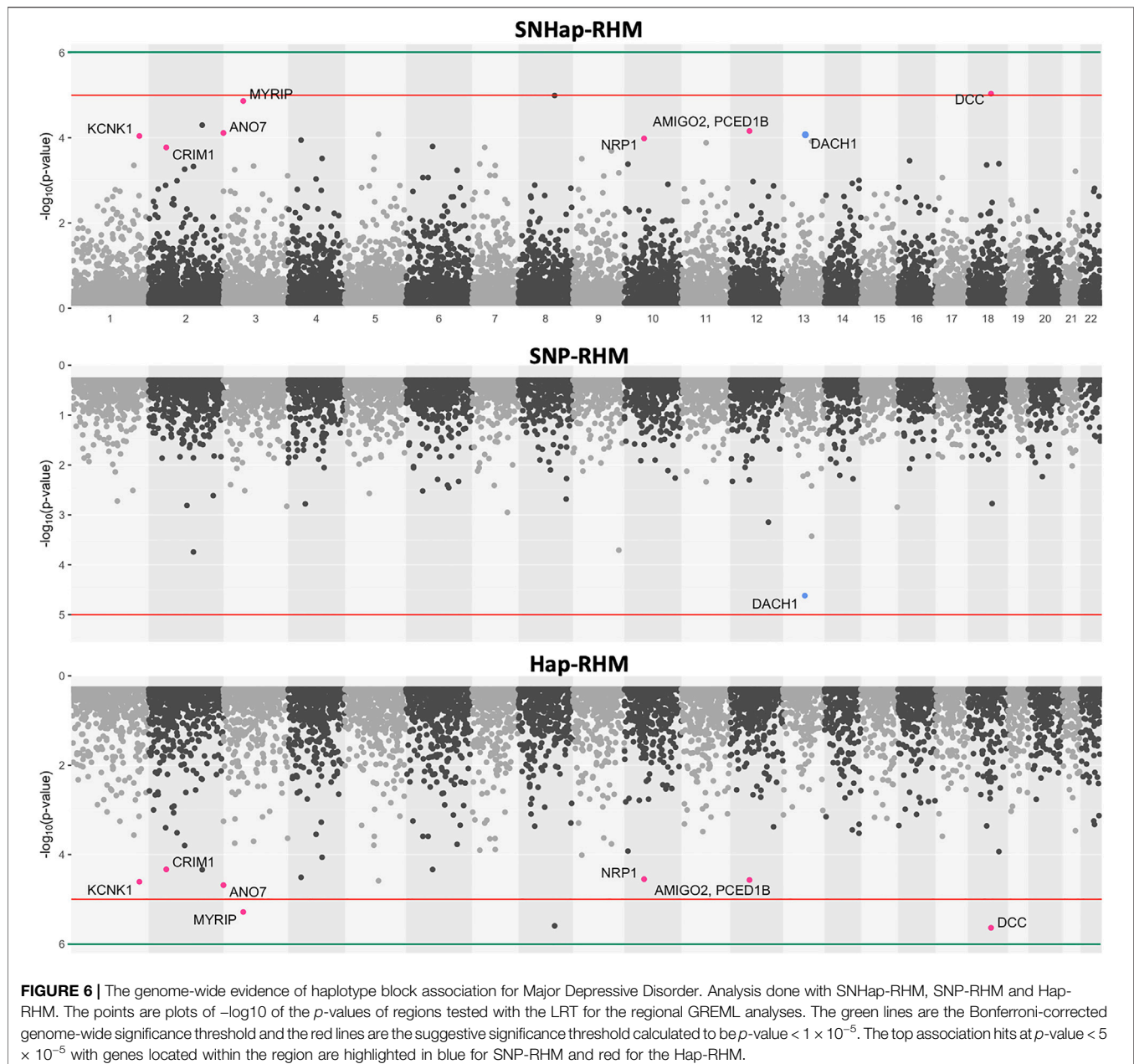


TABLE 1 | SNP-based association test of MDD in the *MYRIP* gene region.

SNP information				Major depressive disorder association			
SNP ID	Chr	Pos	MAF	OR	Log (OR)	SE (logOR)	p
rs9842160	3	39844703	0.14	0.97	-0.030	0.013	0.02
rs9858242	3	39847606	0.19	1.02	0.025	0.011	0.03
rs1599902	3	39954674	0.41	1.02	0.019	0.009	0.04
rs7618607	3	39947936	0.41	1.02	0.019	0.009	0.04
rs9860916	3	39944942	0.41	1.02	0.019	0.009	0.04

The columns are the SNP ID, chromosome, genome position of SNP, minor allele frequency, odds ratio, log of odds ratio, standard error of log odds ratio and association p -value.

TABLE 2 | Comparison of SNPs within significant regions identified by both models and published GWAS results for height and MDD.

Trait	Number of SNPS			Number of overlapping SNPS		
	SNP-RHM	Hap-RHM	pubGWAS	SNP-RHM & Hap-RHM	SNP-RHM & pubGWAS	Hap-RHM & pubGWAS
Height	1,380	45	4,960	0	57	0
MDD	78	495	1,815	0	0	0

The columns are the name of trait, number of SNPS in regions identified by SNP-RHM and HAP-RHM with p -value $< 5 \times 10^{-5}$ and SNPS in published GWAS (pubGWAS) for the traits, and the number of SNPS overlapping between the three.

In both traits, the top significant regions we mapped at p -value $< 5 \times 10^{-5}$ had genes mapped to those regions or within 400 kb of those regions. For height, these genes have been reported to be associated with height in humans (Gudbjartsson et al., 2008; Weedon et al., 2008; Lango Allen et al., 2010; Wood et al., 2014; Nagy et al., 2017; Tachmazidou et al., 2017; Kichaev et al., 2019). For MDD, these genes have been reported to be associated with major depressive disorder and other psychiatry phenotypes (Luciano et al., 2011; Zeng et al., 2017; Wray et al., 2018; Arnau-Soler et al., 2019; Howard et al., 2019; Liu et al., 2019). In one of such regions for MDD, five SNPs within the region are individually significantly associated with MDD at the nominal level (p -value < 0.05). Four of these SNPs lie within the gene sequence of *MYRIP*, and they each confer 2% disease risk. A conventional GWAS analysis would have missed these nominally associated SNPs because they will not reach the suggestive significance threshold, let alone genome-wide (GW) significance. However, analysing these SNPs within the region as haplotypes allowed us to detect the combined effect of these SNPs in the region at a suggestive-significance level even with our relatively small sample size compared to recent genome-wide association studies of MDD: 322,580 (Howard et al., 2018) and 480,359 (Wray et al., 2018).

The current study's primary strength is that we show the ability of SNHap-RHM to incorporate SNP and haplotype information jointly to map genomic regions that affect complex traits. This gives SNHap-RHM a uniquely useful role to play in the future of complex traits analysis. The plummeting costs of whole-genome resequencing (Caulfield et al., 2013) have shifted research focus in GWA studies towards sequence data analysis (Höglund et al., 2019). Although whole-genome sequence data analysis allows incorporating all the genetic variants that drive the phenotypic variation, there may still be some variants whose individual effects may be too small to be picked up in a conventional GWA analysis. However, regionally analysing sequence information can help overcome this because multiple small-effect variants in a region can add up to a substantial regional effect that can be captured by a regional SNP GRM or tagged by a haplotype GRM. Moreover, by defining haplotype blocks using recombination hotspots, whole-genome information can be summarised naturally without setting an arbitrary number of SNPs, and that facilitates integration and comparison across studies. More so, regional heritability analysis of sequence data would be an efficient way to deal with the burden of multiple testing, which has long been a problem of conventional GWAS.

One limitation of the current study is the computation burden of the analyses, which necessitates the pre-correction of the phenotypes with the whole-genome GRM before performing SNHap-RHM. This was a leave-one-chromosome-out step involving 22 separate GREML analyses, each fitting a whole-genome GRM that excluded SNPs from one chromosome (Yang et al., 2014). For our sample of about 20,000 individuals, the precorrection step reduced the computation time needed to perform GREML analysis at each region by approximately 33% (15 min) and used about 20% (16 gigabytes) less memory. Although this was done to speed up the analysis, the precorrection step was used as an approximation to account for the background polygenic effects of genetic markers outside each region; this would have been about 48,772 separate GREMLs to account for each region. One way to get around the computational burden of accounting for the background polygenic effect and speed up the analysis would be to sidestep the computation of the whole-genome GRM by using a decomposition step similar to the one used by FaST-LMM (Lippert et al., 2011). Additionally, it would be interesting to explore the incorporation of other GRMs that account for allele frequencies and LD (Speed et al., 2020) in the genomic background correction stage of SNHap-RHM, going forward. However, whether that will perform better than the commonly employed standard GRM proposed by VanRaden (2008) remains unclear (Rawlik et al., 2020). Moreover, for the regional matrices in SNHap-RHM, it is important to retain SNPs in LD as these determine the haplotype structure that we wish to explore. Also, due to the two degrees of freedom test applied in SNHap-RHM, we observed a slight drop in the significance of the associated regions in both height and MDD when SNHap-RHM was applied to those traits. One option would be to use a less stringent test for SNHap-RHM, effectively testing regions assuming only one degree of freedom so that if only one of the variance components significantly contributed to the phenotypic variance the region would be identified for subsequent formal testing of the individual variance components.

Finally, although this study thoroughly evaluates the robustness of SNP and Haplotype RHM using simulation and demonstrates the utility of SNHap-RHM in real phenotype analysis, seeking replication in other cohorts will improve our understanding and, more importantly, demonstrate that the analysis is portable across studies and genotyping platforms.

5 CONCLUSION

We have implemented a regional heritability analysis and undertaken analyses of regions in the genome delimited by recombination boundaries and shown by simulation that haplotype-based GRMs can capture genetic variance that may be missed by conventional SNP-based GRMs. We then applied this method in the analysis of real phenotype data from GS: SFHS. Again, we show that the haplotype-based regional heritability model uncovers associations in regions of the genome that explain genetic variance missed by the SNP-based heritability model. In light of this, we further showed that regional effects can still be captured when the two regional GRMs (SNP and haplotype-based) are fitted jointly: an analytical procedure we termed SNHap-RHM. This SNHap-RHM presents an exciting new opportunity to analyse complex traits by allowing the joint mapping of novel genomic regions tagged by either SNPs or haplotypes, potentially leading to the recovery of some of the “missing” heritability.

DATA AVAILABILITY STATEMENT

The data analyzed in this study is subject to the following licenses/restrictions: All relevant data supporting the conclusion of this article are included within the article and its **Supplementary Material**. Generation Scotland data are available from the MRC IGC Institutional Data Access/Ethics Committee (<https://www.ed.ac.uk/generation-scotland/for-researchers/access>) for researchers who meet the criteria for access to confidential data. The managed access process ensures that approval is granted only to research which comes under the terms of participant consent which does not allow making participant information publicly available. Requests to access these datasets should be directed to Archie Campbell, archie.campbell@igc.ed.ac.uk/ access@generationscotland.org.

ETHICS STATEMENT

The studies involving human participants were reviewed and approved by the Tayside Committee on Medical Research Ethics (on behalf of the National Health Service). Written informed consent to participate in this study was provided by the participants' legal guardian/next of kin.

REFERENCES

- Amador, C., Huffman, J., Trochet, H., Campbell, A., Porteous, D., Wilson, J. F., et al. (2015). Recent Genomic Heritage in Scotland. *BMC Genomics* 16, 1–17. doi:10.1186/s12864-015-1605-2
- Arnau-Soler, A., Macdonald-Dunlop, E., Macdonald-Dunlop, E., Adams, M. J., Clarke, T.-K., MacIntyre, D. J., et al. (2019). Genome-Wide by Environment Interaction Studies of Depressive Symptoms and Psychosocial Stress in UK Biobank and Generation Scotland. *Transl Psychiatry* 9, 14. doi:10.1038/s41398-018-0360-y

AUTHOR CONTRIBUTIONS

Conceived and designed the experiments: RO, PN, CH, SK. Provided data: TB, AC, AM, DP, CH. Performed the experiments: RO. Analysed the data: RO. Wrote the paper: RO, PN, CH, SK.

FUNDING

The first author (RO) was funded by the Darwin Trust of Edinburgh (<https://darwintrust.bio.ed.ac.uk/>) for his PhD study (no grant number). CH and PN acknowledge funding from the Medical Research Council UK (MRC, <https://mrc.ukri.org/funding/>): MC_UU_00007/10, MC_PC_U127592696, MC_PC_U127561128; the BBSRC (<https://bbsrc.ukri.org/funding/>): BBS/E/D/30002275, BBS/E/D/30002276 and a Wellcome Trust (<https://wellcome.org/grant-funding>) Investigator Award to AM: 220857/Z/20/Z. Generation Scotland received core support from the Chief Scientist Office of the Scottish Government Health Directorates (CZD/16/6) and the Scottish Funding Council (HR03006) and is currently supported by the Wellcome Trust (216767/Z/19/Z). Genotyping of the GS: SFHS samples was funded by the Medical Research Council UK and the Wellcome Trust (Wellcome Trust Strategic Award “STratifying Resilience and Depression Longitudinally” (STRADL) Reference 104036/Z/14/Z).

ACKNOWLEDGMENTS

We are grateful to all the families who took part, the general practitioners, and the Scottish School of Primary Care for their help in recruiting them, and the whole Generation Scotland team, which includes interviewers, computer and laboratory technicians, clerical workers, research scientists, volunteers, managers, receptionists, healthcare assistants and nurses. We also acknowledge Eilidh Fummey for coming up with the name for the joint mapping method.

SUPPLEMENTARY MATERIAL

The Supplementary Material for this article can be found online at: <https://www.frontiersin.org/articles/10.3389/fgene.2021.791712/full#supplementary-material>

- Balding, D. J. (2006). A Tutorial on Statistical Methods for Population Association Studies. *Nat. Rev. Genet.* 7, 781–791. doi:10.1038/nrg1916
- Caulfield, T., Evans, J., McGuire, A., McCabe, C., Bubela, T., Cook-Deegan, R., et al. (2013). Reflections on the Cost of “Low-Cost” Whole Genome Sequencing: Framing the Health Policy Debate. *PLoS Biol.* 11, e1001699. doi:10.1371/journal.pbio.1001699
- Cirulli, E. T., and Goldstein, D. B. (2010). Uncovering the Roles of Rare Variants in Common Disease Through Whole-Genome Sequencing. *Nat. Rev. Genet.* 11, 415–425. doi:10.1038/nrg2779
- Clarke, A. J., and Cooper, D. N. (2010). GWAS: Heritability Missing in Action? *Eur. J. Hum. Genet.* 18, 859–861. doi:10.1038/ejhg.2010.35

- Delaneau, O., Zagury, J.-F., and Marchini, J. (2013). Improved Whole-Chromosome Phasing for Disease and Population Genetic Studies. *Nat. Methods* 10, 5–6. doi:10.1038/nmeth.2307
- First, M. B., Spitzer, R. L., Gibbon, M., and Williams, J. B. W. (2002). *Structured Clinical Interview for DSM-IV-TR Axis I Disorders, Research Version, Non-Patient Edition*. New York, NY.
- Frazer, K. A., Ballinger, D. G., Cox, D. R., Hinds, D. A., Stuve, L. L., Gibbs, R. A., et al. (2007). A Second Generation Human Haplotype Map of over 3.1 Million SNPs. *Nature* 449, 851–861. doi:10.1038/nature06258
- Ganat, Y. M., Calder, E. L., Kriks, S., Nelander, J., Tu, E. Y., Jia, F., et al. (2012). Identification of Embryonic Stem Cell-Derived Midbrain Dopaminergic Neurons for Engraftment. *J. Clin. Invest.* 122, 2928–2939. doi:10.1172/JCI58767
- Gonzalez-Recio, O., Daetwyler, H. D., MacLeod, I. M., Pryce, J. E., Bowman, P. J., Hayes, B. J., et al. (2015). Rare Variants in Transcript and Potential Regulatory Regions Explain a Small Percentage of the Missing Heritability of Complex Traits in Cattle. *PLoS One* 10, e0143945. doi:10.1371/journal.pone.0143945
- Gottlieb, D. J., O'Connor, G. T., and Wilk, J. B. (2007). Genome-wide Association of Sleep and Circadian Phenotypes. *BMC Med. Genet.* 8, S9. doi:10.1186/1471-2350-8-S1-S9
- Greenwood, E. A., Pasch, L. A., Shinkai, K., Cedars, M. I., and Huddleston, H. G. (2015). Putative Role for Insulin Resistance in Depression Risk in Polycystic Ovary Syndrome. *Fertil. Sterility* 104, 707–714.e1. doi:10.1016/j.fertnstert.2015.05.019
- Gudbjartsson, D. F., Walters, G. B., Thorleifsson, G., Stefansson, H., Halldorsson, B. V., Zusmanovich, P., et al. (2008). Many Sequence Variants Affecting Diversity of Adult Human Height. *Nat. Genet.* 40, 609–615. doi:10.1038/ng.122
- Höglund, J., Rafati, N., Rask-Andersen, M., Enroth, S., Karlsson, T., Ek, W. E., et al. (2019). Improved Power and Precision with Whole Genome Sequencing Data in Genome-Wide Association Studies of Inflammatory Biomarkers. *Sci. Rep.* 9, 16844. doi:10.1038/s41598-019-53111-7
- Howard, D. M., Adams, M. J., Adams, M. J., Clarke, T.-K., Hafferty, J. D., Gibson, J., et al. (2019). Genome-Wide Meta-Analysis of Depression Identifies 102 Independent Variants and Highlights the Importance of the Prefrontal Brain Regions. *Nat. Neurosci.* 22, 343–352. doi:10.1038/s41593-018-0326-7
- Howard, D. M., Adams, M. J., Adams, M. J., Shirali, M., Clarke, T.-K., Marioni, R. E., et al. (2018). Genome-wide Association Study of Depression Phenotypes in UK Biobank Identifies Variants in Excitatory Synaptic Pathways. *Nat. Commun.* 9, 1470. doi:10.1038/s41467-018-03819-3
- International Human Genome Sequencing Consortium (2004). Finishing the Euchromatic Sequence of the Human Genome. *Nature* 431, 931–945. doi:10.1038/nature03001
- Kichaev, G., Bhatia, G., Loh, P.-R., Gazal, S., Burch, K., Freund, M. K., et al. (2019). Leveraging Polygenic Functional Enrichment to Improve GWAS Power. *Am. J. Hum. Genet.* 104, 65–75. doi:10.1016/j.ajhg.2018.11.008
- Lango Allen, H., Estrada, K., Lettre, G., Berndt, S. L., Weedon, M. N., Rivadeneira, F., et al. (2010). Hundreds of Variants Clustered in Genomic Loci and Biological Pathways Affect Human Height. *Nature* 467, 832–838. doi:10.1038/nature09410
- Levinson, D. F., Mostafavi, S., Milaneschi, Y., Rivera, M., Ripke, S., Wray, N. R., et al. (2014). Genetic Studies of Major Depressive Disorder: Why are There No Genome-wide Association Study Findings and What Can We do About it? *Biol. Psychiatry* 76, 510–512. doi:10.1016/j.biopsych.2014.07.029
- Lippert, C., Listgarten, J., Liu, Y., Kadie, C. M., Davidson, R. I., and Heckerman, D. (2011). FaST Linear Mixed Models for Genome-Wide Association Studies. *Nat. Methods* 8, 833–835. doi:10.1038/nmeth.1681
- Liu, M., Jiang, Y., Wedow, R., Li, Y., Brazier, D. M., Chen, F., et al. (2019). Association Studies of up to 1.2 Million Individuals Yield New Insights into the Genetic Etiology of Tobacco and Alcohol Use. *Nat. Genet.* 51, 237–244. doi:10.1038/s41588-018-0307-5
- Luciano, M., Hansell, N. K., Lahti, J., Davies, G., Medland, S. E., Rääkkönen, K., et al. (2011). Whole Genome Association Scan for Genetic Polymorphisms Influencing Information Processing Speed. *Biol. Psychol.* 86, 193–202. doi:10.1016/j.biopsycho.2010.11.008
- MacArthur, J., Bowler, E., Cerezo, M., Gil, L., Hall, P., Hastings, E., et al. (2017). The New NHGRI-EBI Catalog of Published Genome-wide Association Studies (GWAS Catalog). *Nucleic Acids Res.* 45, D896–D901. doi:10.1093/nar/gkx1133
- Maher, B. (2008). Personal Genomes: The Case of the Missing Heritability. *Nature* 456, 18–21. doi:10.1038/456018a
- Manolio, T. A., Collins, F. S., Cox, N. J., Goldstein, D. B., Hindorf, L. A., Hunter, D. J., et al. (2009). Finding the Missing Heritability of Complex Diseases. *Nature* 461, 747–753. doi:10.1038/nature08494
- Mohan, J., Xiaofan, G., and Yingxian, S. (2017). Association Between Sleep Time and Depression: a Cross-Sectional Study from Countries in Rural Northeastern China. *J. Int. Med. Res.* 45, 984–992. doi:10.1177/0300060517701034
- Nagamine, Y., Pong-Wong, R., Navarro, P., Vitart, V., Hayward, C., Rudan, I., et al. (2012). Localising Loci Underlying Complex Trait Variation Using Regional Genomic Relationship Mapping. *PLoS ONE* 7, e46501. doi:10.1371/journal.pone.0046501
- Nagy, R., Boutin, T. S., Marten, J., Huffman, J. E., Kerr, S. M., Campbell, A., et al. (2017). Exploration of Haplotype Research Consortium Imputation for Genome-wide Association Studies in 20,032 Generation Scotland Participants. *Genome Med.* 9, 23. doi:10.1186/s13073-017-0414-4
- Pearson, S., Schmidt, M., Patton, G., Dwyer, T., Blizzard, L., Otahal, P., et al. (2010). Depression and Insulin Resistance: Cross-Sectional Associations in Young Adults. *Diabetes Care* 33, 1128–1133. doi:10.2337/dc09-1940
- Pritchard, J. K. (2001). Are Rare Variants Responsible for Susceptibility to Complex Diseases. *Am. J. Hum. Genet.* 69, 124–137. doi:10.1086/321272
- Rawlik, K., Canela-Xandri, O., Woolliams, J., and Tenesa, A. (2020). SNP Heritability: What are we Estimating? *bioRxiv* doi:10.1101/2020.09.15.276121
- Roberts, R. E., and Duong, H. T. (2014). The Prospective Association Between Sleep Deprivation and Depression Among Adolescents. *Sleep* 37, 239–244. doi:10.5665/sleep.3388
- Shirali, M., Knott, S. A., Pong-Wong, R., Navarro, P., and Haley, C. S. (2018). Haplotype Heritability Mapping Method Uncovers Missing Heritability of Complex Traits. *Sci. Rep.* 8, 4982. doi:10.1038/s41598-018-23307-4
- Smith, B. H., Campbell, A., Linksted, P., Fitzpatrick, B., Jackson, C., Kerr, S. M., et al. (2012). Cohort Profile: Generation Scotland: Scottish Family Health Study (GS:SFHS). The Study, its Participants and Their Potential for Genetic Research on Health and Illness. *Int. J. Epidemiol.* 42, 689–700. doi:10.1093/ije/dys084
- Smith, B. H., Campbell, H., Blackwood, D., Connell, J., Connor, M., Deary, I. J., et al. (2006). Generation Scotland: The Scottish Family Health Study; A New Resource for Researching Genes and Heritability. *BMC Med. Genet.* 7. doi:10.1186/1471-2350-7-74
- Speed, D., Hemani, G., Johnson, M. R., and Balding, D. J. (2012). Improved Heritability Estimation from Genome-Wide SNPs. *Am. J. Hum. Genet.* 91, 1011–1021. doi:10.1016/j.ajhg.2012.10.010
- Speed, D., Holmes, J., and Balding, D. J. (2020). Evaluating and Improving Heritability Models Using Summary Statistics. *Nat. Genet.* 52, 458–462. doi:10.1038/s41588-020-0600-y
- Tachmazidou, I., Süveges, D., Min, J. L., Ritchie, G. R. S., Steinberg, J., Walter, K., et al. (2017). Whole-Genome Sequencing Coupled to Imputation Discovers Genetic Signals for Anthropometric Traits. *Am. J. Hum. Genet.* 100, 865–884. doi:10.1016/j.ajhg.2017.04.014
- Uemoto, Y., Pong-Wong, R., Navarro, P., Vitart, V., Hayward, C., Wilson, J. F., et al. (2013). The Power of Regional Heritability Analysis for Rare and Common Variant Detection: Simulations and Application to Eye Biometrical Traits. *Front. Genet.* 4, 232. doi:10.3389/fgene.2013.00232
- VanRaden, P. M. (2008). Efficient Methods to Compute Genomic Predictions. *J. Dairy Sci.* 91, 4414–4423. doi:10.3168/jds.2007-0980
- Vormfelde, S. V., and Brockmüller, J. (2007). On the Value of Haplotype-Based Genotype-Phenotype Analysis and on Data Transformation in Pharmacogenetics and -genomics. *Nat. Rev. Genet.* 8, 983. doi:10.1038/nrg1916-c1
- Wainschein, P., Jain, D., Zheng, Z., Cupples, L. A., Shadyab, A. H., McKnight, B., et al. (2019). Recovery of Trait Heritability from Whole Genome Sequence Data. *bioRxiv* 588020. doi:10.1101/588020
- Waselle, L., Coppola, T., Fukuda, M., Iezzi, M., El-Amraoui, A., Petit, C., et al. (2003). Involvement of the Rab27 Binding Protein Slac2c/MyRIP in Insulin Exocytosis. *Mol. Biol. Cell* 14, 4103–4113. doi:10.1091/mbc.E03-01-0022
- Watson, N. F., Harden, K. P., Buchwald, D., Vitiello, M. V., Pack, A. I., Strachan, E., et al. (2014). Sleep Duration and Depressive Symptoms: A Gene-Environment Interaction. *Sleep* 37, 351–358. doi:10.5665/sleep.3412
- Webb, M. B., Davies, M., Ashra, N., Bodicoat, D., Brady, E., Webb, D., et al. (2017). The Association Between Depressive Symptoms and Insulin Resistance,

- Inflammation and Adiposity in Men and Women. *PLoS One* 12, e0187448. doi:10.1371/journal.pone.0187448
- Weedon, M. N., Lango, H., Lango, H., Lindgren, C. M., Wallace, C., Evans, D. M., et al. (2008). Genome-Wide Association Analysis Identifies 20 Loci that Influence Adult Height. *Nat. Genet.* 40, 575–583. doi:10.1038/ng.121
- Wood, A. R., Esko, T., Yang, J., Vedantam, S., Pers, T. H., Gustafsson, S., et al. (2014). Defining the Role of Common Variation in the Genomic and Biological Architecture of Adult Human Height. *Nat. Genet.* 46, 1173–1186. doi:10.1038/ng.3097
- Wray, N. R., Ripke, S., Mattheisen, M., Trzaskowski, M., Byrne, E. M., Abdellaoui, A., et al. (2018). Genome-Wide Association Analyses Identify 44 Risk Variants and Refine the Genetic Architecture of Major Depression. *Nat. Genet.* 50, 668–681. doi:10.1038/s41588-018-0090-3
- Yang, J., Bakshi, A., Bakshi, A., Zhu, Z., Hemani, G., Vinkhuyzen, A. A. E., et al. (2015). Genetic Variance Estimation with Imputed Variants Finds Negligible Missing Heritability for Human Height and Body Mass Index. *Nat. Genet.* 47, 1114–1120. doi:10.1038/ng.3390
- Yang, J., Benyamin, B., McEvoy, B. P., Gordon, S., Henders, A. K., Nyholt, D. R., et al. (2010). Common SNPs Explain a Large Proportion of the Heritability for Human Height. *Nat. Genet.* 42, 565–569. doi:10.1038/ng.608
- Yang, J., Lee, S. H., Goddard, M. E., and Visscher, P. M. (2011). GCTA: A Tool for Genome-Wide Complex Trait Analysis. *Am. J. Hum. Genet.* 88, 76–82. doi:10.1016/j.ajhg.2010.11.011
- Yang, J., Zaitlen, N. A., Goddard, M. E., Visscher, P. M., and Price, A. L. (2014). Advantages and Pitfalls in the Application of Mixed-Model Association Methods. *Nat. Genet.* 46, 100–106. doi:10.1038/ng.2876
- Zeng, Y., Navarro, P., Fernandez-Pujals, A. M., Hall, L. S., Clarke, T.-K., Thomson, P. A., et al. (2017). A Combined Pathway and Regional Heritability Analysis Indicates NETRIN1 Pathway is Associated With Major Depressive Disorder. *Biol. Psychiatry* 81, 336–346. doi:10.1016/j.biopsych.2016.04.017
- Zhai, L., Zhang, H., and Zhang, D. (2015). Sleep Duration and Depression Among Adults: A Meta-Analysis of Prospective Studies. *Depress. Anxiety* 32, 664–670. doi:10.1002/da.22386
- Conflict of Interest:** AM has received research support from Eli Lilly and Company, Janssen and the Sackler Trust and speaker fees from Illumina and Janssen.
- The remaining authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.
- Publisher's Note:** All claims expressed in this article are solely those of the authors and do not necessarily represent those of their affiliated organizations, or those of the publisher, the editors and the reviewers. Any product that may be evaluated in this article, or claim that may be made by its manufacturer, is not guaranteed or endorsed by the publisher.
- Copyright © 2022 Oppong, Boutin, Campbell, McIntosh, Porteous, Hayward, Haley, Navarro and Knott. This is an open-access article distributed under the terms of the Creative Commons Attribution License (CC BY). The use, distribution or reproduction in other forums is permitted, provided the original author(s) and the copyright owner(s) are credited and that the original publication in this journal is cited, in accordance with accepted academic practice. No use, distribution or reproduction is permitted which does not comply with these terms.