

## Article

# A Novel Framework for Analysis of the Shared Genetic Background of Correlated Traits

Gulnara R. Svishcheva <sup>1,2</sup>, Evgeny S. Tiys <sup>1</sup>, Elizaveta E. Elgaeva <sup>1,3</sup>, Sofia G. Feoktistova <sup>1</sup>, Paul R. H. J. Timmers <sup>4,5</sup>, Sodbo Zh. Sharapov <sup>1,3</sup>, Tatiana I. Axenovich <sup>1</sup> and Yakov A. Tsepilov <sup>1,3,\*</sup>

<sup>1</sup> Institute of Cytology and Genetics, Siberian Branch of the Russian Academy of Sciences, 630090 Novosibirsk, Russia

<sup>2</sup> Vavilov Institute of General Genetics, Russian Academy of Sciences, 117971 Moscow, Russia

<sup>3</sup> Novosibirsk State University, 630090 Novosibirsk, Russia

<sup>4</sup> MRC Human Genetics Unit, MRC Institute of Genetics and Cancer, University of Edinburgh, Edinburgh EH8 9YL, UK

<sup>5</sup> Centre for Global Health Research, Usher Institute, University of Edinburgh, Edinburgh EH8 9YL, UK

\* Correspondence: tsepilov@bionet.nsc.ru

**Abstract:** We propose a novel effective framework for the analysis of the shared genetic background for a set of genetically correlated traits using SNP-level GWAS summary statistics. This framework called SHAHER is based on the construction of a linear combination of traits by maximizing the proportion of its genetic variance explained by the shared genetic factors. SHAHER requires only full GWAS summary statistics and matrices of genetic and phenotypic correlations between traits as inputs. Our framework allows both shared and unshared genetic factors to be effectively analyzed. We tested our framework using simulation studies, compared it with previous developments, and assessed its performance using three real datasets: anthropometric traits, psychiatric conditions and lipid concentrations. SHAHER is versatile and applicable to summary statistics from GWASs with arbitrary sample sizes and sample overlaps, allows for the incorporation of different GWAS models (Cox, linear and logistic), and is computationally fast.

**Keywords:** GWAS; shared genetic component; linear combination of traits; shared heritability; proportion of heritability explained by SGF



**Citation:** Svishcheva, G.R.; Tiys, E.S.; Elgaeva, E.E.; Feoktistova, S.G.; Timmers, P.R.H.J.; Sharapov, S.Z.; Axenovich, T.I.; Tsepilov, Y.A. A Novel Framework for Analysis of the Shared Genetic Background of Correlated Traits. *Genes* **2022**, *13*, 1694. <https://doi.org/10.3390/genes13101694>

Academic Editors: Xiaolei Liu, Yunlong Ma, Yuhua Fu and Lilin Yin

Received: 11 August 2022

Accepted: 16 September 2022

Published: 21 September 2022

**Publisher's Note:** MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Copyright:** © 2022 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

## 1. Introduction

There is a growing interest in studying the shared genetic background between genetically correlated traits [1–5] (see, for example, the number of PubMed search results by year for keywords related to “shared genetic background”). Studying the shared genetics between traits can help with the discovery of pleiotropic interactions, common genes and pathways, and identify genetic effects that are specific for each trait.

The problem of the decomposition of the variance of several traits into the shared/unshared genetic and environment components were first formulated by S. Wright in 1921 [6]. There are widely used classic twin designs to address this problem. They are based on structural equation modelling; in particular, multivariate pathway models assuming the existence of the genetic influences common for all traits and specific for each trait [7]. These designs are implemented only for the variance decomposition, but not for the identification of the genetic factors that determine these genetic effects.

There are several terms for these common and specific genetic impacts. We will call them the “shared genetic impact” (SGI) and “unshared genetic impacts” (UGI). The genetic factors that determine these impacts will be called “shared genetic factors” (SGF) and “unshared genetic factors” (UGF), respectively. The heritability of each trait explained by SGF and UGF will be called “shared heritability” and “unshared heritability”, respectively.

Note that the term “unshared” for every trait means the rest other than the “shared”. For any incomplete set of traits, UGFs can be partially overlapping.

The application of different methods of multivariate analysis in genome-wide association studies (GWAS) allows the problem of SGF and UGF identification to be partially solved [8–13]. The multivariate methods involve complicated genetic or/and phenotypic correlation structures of traits in the analysis. In most cases, this increases the power of detection of the loci associated with several traits due to pleiotropic effects. If the detected locus has a pleiotropic effect on all studied traits, the locus could potentially be attributed to SGF, and if not, to UGF. However, a pleiotropic effect of the locus on all studied traits is necessary but insufficient for inclusion of this locus in SGF (at least effects should be also collinear between traits; see the model description below). Moreover, this approach of SGF identification assumes a manual classification of loci, which prevents the use of more sophisticated modern in-silico approaches for genetic analysis, for example, the ones that rely on GWAS summary statistics [14]. To our knowledge, there is no specific method that could be good for both variance component decomposition and identification of SGF and UGF.

We had previously developed a method for obtaining genetically independent phenotypes (GIPs) [2]. This method is based on the calculation of the principal components using genetic rather than phenotypic correlations. We applied this method to genetically correlated pain phenotypes and aging related phenotypes and showed that the first GIP component, GIP1, that explains the largest proportion of the genetic variance probably could be interpreted as SGI [2,15]. This makes GIP promising for the identification of loci attributed to SGF. However, this method was not designed specifically for SGI analysis. In addition, no specific experiments have been performed to validate the approach or to estimate its statistical properties.

Here, we present a novel general framework for the estimation of shared and unshared heritability and identification of the shared and unshared genetic factors using the summary statistics of original traits. The essence of our approach is to find the optimum linear combination of traits which has the maximum proportion of its genetic variance explained by the SGF. We validated our framework using simulation studies under different scenarios by comparing it with the developed GIP approach, and assessed its performance using three real datasets: anthropometric indices, psychiatric disorders and conditions, and lipid concentrations.

## 2. Materials and Methods

### 2.1. Shared Heredity Model

We adopted a commonly used multivariate pathway model [7] in terms of SGF and UGF. We call it the “shared heredity model”. For simplicity, we consider SGF and UGF as biallelic SNPs and consider a sample of  $N$  unrelated individuals measured for  $K$  traits and genotyped for  $M$  SNPs. For a standardized normal trait,  $y$  ( $N \times 1$ ), the traditional polygenic (null) model takes the form:  $y = G\beta + \varepsilon$ , where  $G$  is an ( $N \times M$ ) matrix of standardized genotypes;  $\beta$  ( $M \times 1$ ) and  $\varepsilon$  ( $N \times 1$ ) are genetic and non-genetic random effects, respectively;  $\beta \sim N(\mathbf{0}, h^2 I_M)$  and  $\varepsilon \sim N(\mathbf{0}, (1 - h^2) I_N)$ , where  $\mathbf{0}$  is a null mean vector,  $h^2$  is the trait heritability, and  $I$  is an identity matrix of the given dimension. For unrelated individuals, we expect  $y \sim N(\mathbf{0}, I_N)$ .

We propose to divide  $M$  SNPs into two non-overlapping SNP sets with sizes  $M_0$  and  $M_1$  ( $M_0 + M_1 = M$ ). The set of  $M_0$  SNPs called “SGF” includes only those SNPs whose effects on all traits are collinear. The set of  $M_1$  SNPs consists of the other SNPs which do not have shared joint influence on all traits at once, this set being called “UGF”. In accordance with  $M$ ,  $G$  is divided into two matrices,  $G_0$  ( $N \times M_0$ ) and  $G_1$  ( $N \times M_1$ ). To decompose every

trait into components explained by SGF and UGF, we rewrote the traditional polygenic model in terms of  $G_0$  and  $G_1$

$$y_i = \underbrace{G_0 b_{0_i}}_{\text{due to SGF}} + \underbrace{G_1 b_{1_i}}_{\text{due to UGF}} + \varepsilon_i. \tag{1}$$

Here, the first and second terms are genetic components explained by SGF and UGF, respectively, which are assumed to be independent. In the first term,  $b_{0_i}$  is an  $(M_0 \times 1)$  vector of non-zero SGF effects, which can be presented as  $\beta_0 w_i \sqrt{h_i^2}$ , where  $\beta_0$  is an  $(M_0 \times 1)$  non-zero vector that is the same for all traits,  $\beta_0 \sim N(0, I_{M_0})$ , and  $w_i^2 h_i^2$  is the heritability of the  $i$ -th trait explained by SGF. Here  $w_i$  is a non-zero trait-specific multiplier:  $w_i^2$  denotes the proportion of  $h_i^2$  explained by SGF; and the value of  $w_i$  can be positive and negative, indicating the direction of the SGF effect on the  $i$ -th trait.  $G_0 \beta_0$  is the so-called shared genetic impact or SGI. In the second term of Model (1),  $b_{1_i}$  is an  $(M_1 \times 1)$  vector of UGF effects, which can be presented as  $b_{1_i} = \beta_{1i} \sqrt{(1 - w_i^2) h_i^2}$ ,  $\beta_{1i} \sim N(0, I_{M_1})$ . In contrast to  $\beta_0$ ,  $\beta_{1i}$  are different for different traits, and they are not collinear. For illustrative purposes, we rewrote Equation (1) as:

$$y_i = \underbrace{G_0 \beta_0 w_i \sqrt{h_i^2}}_{\substack{\text{SGI} \\ \text{due to SGF}}} + \underbrace{G_1 \beta_{1i} \sqrt{1 - w_i^2} \sqrt{h_i^2}}_{\text{due to UGF}} + \varepsilon_i. \tag{2}$$

2.2. Overview of the SHAHER Framework

For analyses of SGI and UGI on a set of correlated traits, we propose an effective multi-stage framework named SHAHER (see Figure 1). The concept of the framework is first to partition the genetic basis of each original trait into two components: one shared by all the original traits and one not shared by all the original traits, and then to identify the SNPs that contribute to these genetic components. To do this, we propose to construct new traits: (1) an SGIT as a linear combination of original traits, which has the maximum possible heritability explained by SGF, and (2) UGITs as linear combinations of the original traits, which are obtained by adjusting the original traits for SGIT. This means that the genetic basis of UGITs is predominantly determined by UGF.

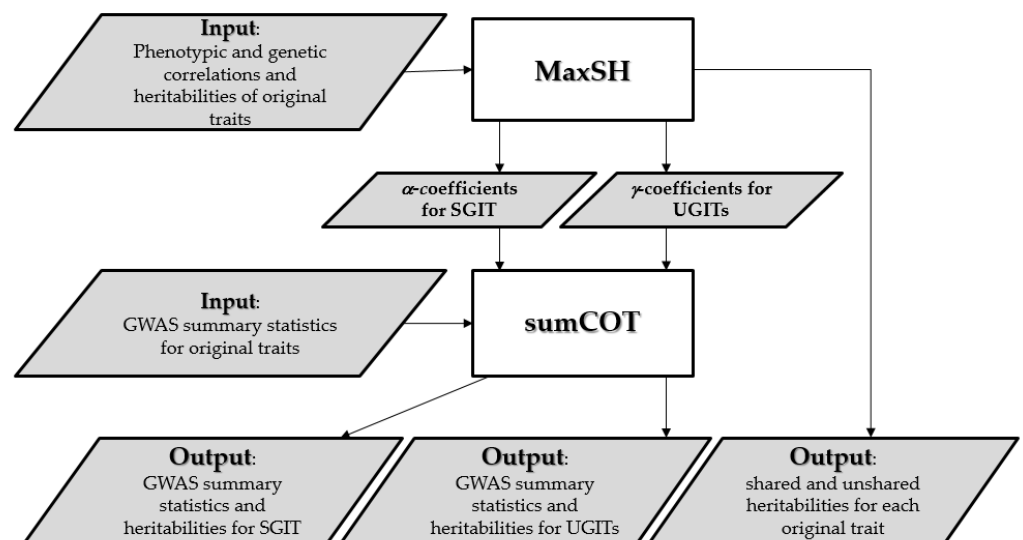


Figure 1. Flowchart of the SHAHER framework. Details are given in the text.

Our framework requires matrices of phenotypic correlations ( $U_{phen}$ ) between the original traits, the matrices of genetic correlations ( $U_{gen}$ ) between the original traits, the heritabilities of the original traits, and GWAS summary statistics of the original traits as inputs. It is worth noting that  $U_{phen}$ ,  $U_{gen}$  and heritabilities could be estimated using GWAS summary statistics of the original traits, for example, by the LD score regression method [16].

SHAHER starts with a preliminary stage, which verifies the presence of SGI in a given set of traits. This is achieved by checking the following requirements for  $U_{gen}$ : it must be positive definite; the absolute values of its elements must be significantly more than a given threshold, and the rank of the correlation matrix derived from  $U_{gen}$  by rounding its elements to extremal correlation values, either  $-1$  or  $1$ , must be equal to one. If the requirements are met, we turn to the basic stages of SHAHER.

*The MaxSH stage.* To determine the  $\alpha$  and  $\gamma$  coefficients for the linear combinations of the original traits to build SGIT and UGITs, we developed the MaxSH method, which is based on the correlation component model given below. This model partitions the phenotypic correlation matrix,  $U_{phen}$ , into environmental and genetic components,  $U_{env}$  and  $U_{gen}$ , respectively, the latter being further subdivided into two components caused by SGF and UGF:

$$\begin{aligned} U_{phen} &= \underbrace{\sqrt{H^2}U_{gen}\sqrt{H^2}}_{\text{genetic component}} + \underbrace{\sqrt{I-H^2}U_{env}\sqrt{I-H^2}}_{\text{environmental component}} \\ U_{gen} &= \underbrace{W\mathbf{1}\mathbf{1}^T W}_{\text{due to SGF}} + \underbrace{\sqrt{I-W^2}U_{unsh}\sqrt{I-W^2}}_{\text{due to UGF}} \end{aligned} \quad (3)$$

Here  $W$  is a diagonal matrix, whose  $i$ -th diagonal element is  $w_i$ ;  $U_{unsh}$  is a matrix of genetic correlations explained by UGF;  $H^2$  is a diagonal matrix, whose  $i$ -th diagonal element is  $h_i^2$ , and  $\mathbf{1}$  is a  $(k \times 1)$  vector of units. Using this model, MaxSH solves several tasks.

First of all, using only the genetic correlation matrix,  $U_{gen}$ , we estimate the proportion of heritability of every trait explained by SGF ( $W$ ). To do this, we minimize the difference between  $U_{gen}$  and the auxiliary matrix  $V$ . This matrix is built using Formula (2), with the identity matrix used instead of  $U_{unsh}$ . The second task is to determine the  $\alpha$ -coefficients, which is solved by maximizing the shared heritability of SGIT. This task is analytically solved as

$$a = \frac{U_{phen}^{-\frac{1}{2}}HW\mathbf{1}}{\sqrt{\mathbf{1}^TWHU_{phen}^{-1}HW\mathbf{1}}} \quad (4)$$

It requires  $U_{phen}$ ,  $H^2$  and  $W$  as input data.

After determining the  $\alpha$ -coefficients and building SGIT, we build a UGIT for every trait using the residual regression equation  $UGIT_i = y_i - SGIT * c_i$ , where  $c_i$  is the impact of SGIT on the  $i$ -th original trait, defined as

$$c_i = \frac{cov_{gen}(y_i, SGIT)}{h_{SGIT}^2}. \quad (5)$$

Here  $cov_{gen}$  denotes a genetic covariance. Note that we should use genetic rather than phenotypic covariances, as our goal is to adjust only the genetic components of the original traits.

Since SGIT is the linear combination of the original traits, UGITs are linear combinations of the original traits, as well. The coefficients of these linear combinations called the  $\gamma$ -coefficients form the matrix of the  $\gamma$ -coefficients  $\Gamma = (I_K - \alpha c^T)$ , where the  $i$ -th column of  $\Gamma$  corresponds to the linear combination coefficients for building the  $i$ -th UGIT.

*The sumCOT stage.* This stage is aimed at obtaining GWAS summary statistics for SGIT and UGITs using the previously determined  $\alpha$  and  $\gamma$  coefficients, GWAS summary statistics (Z-scores, allele frequencies and sample sizes for each SNP) for the original traits and the matrix of phenotypic correlations. The method can use Z scores obtained from

any regression model and allows for varying sample sizes and sample overlap between traits. This sample overlap is incorporated into the estimation of the matrix of phenotypic correlations. In short, the SNP effects for combined traits are calculated by summing effect estimates from the individual trait GWASs, each multiplied by their corresponding linear coefficient ( $\alpha$  or  $\gamma$ ), and standardized by the expected variance. The standard errors of the SNP effect are calculated using variance-covariance arithmetic, taking into account the phenotypic covariance between GWAS results to adjust for the sample overlap. Effective sample sizes are then estimated based on the median Z statistic and allele frequencies by solving Equation (1) in [17].

At the final stage, SHAHER checks for the correctness of the output. In particular, we anticipate that UGITs do not have a shared genetic basis. This is verified by applying MaxSH to the matrix of correlations between UGITs.

To summarize, our framework estimates shared and unshared heritabilities for each of the studied original traits and produces GWAS summary statistics for SGIT and UGITs as outputs.

The full details and mathematical formulae of SHAHER are in Supplementary Methods.

### 2.3. Simulation Study

Under different scenarios, we designed simulations to assess the performance of MaxSH. We (1) assessed the accuracy of  $w$  estimates, (2) assessed the proportion of SGIT heritability explained by SGF to the total heritability of SGIT (the  $Q$ -value), and (3) compared the analytically predicted total and shared heritabilities of SGIT and GIP1 with respect to the loss function. The design of our simulation experiment is shown in Figure 2. To generate the input for the MaxSH and GIP approaches, we used a six-parameter simulation model, in which  $K$  is the number of traits;  $W_0^2$  is a ( $K \times K$ ) diagonal matrix, where the  $i$ -th diagonal element is  $w_i^2$  (the proportion of heritability explained by SGF);  $s$  is the proportion of zeros in the matrix  $U_{unsh}$ ;  $d_1$  is the amplitude of the uniform distribution for non-zero values of  $U_{unsh}$  and  $d_2$  is the amplitude of the uniform distribution for  $U_{env}$ ;  $H^2$  is the diagonal matrix with diagonal elements equal to the trait heritabilities. The parameter values used are given in Figure 2.

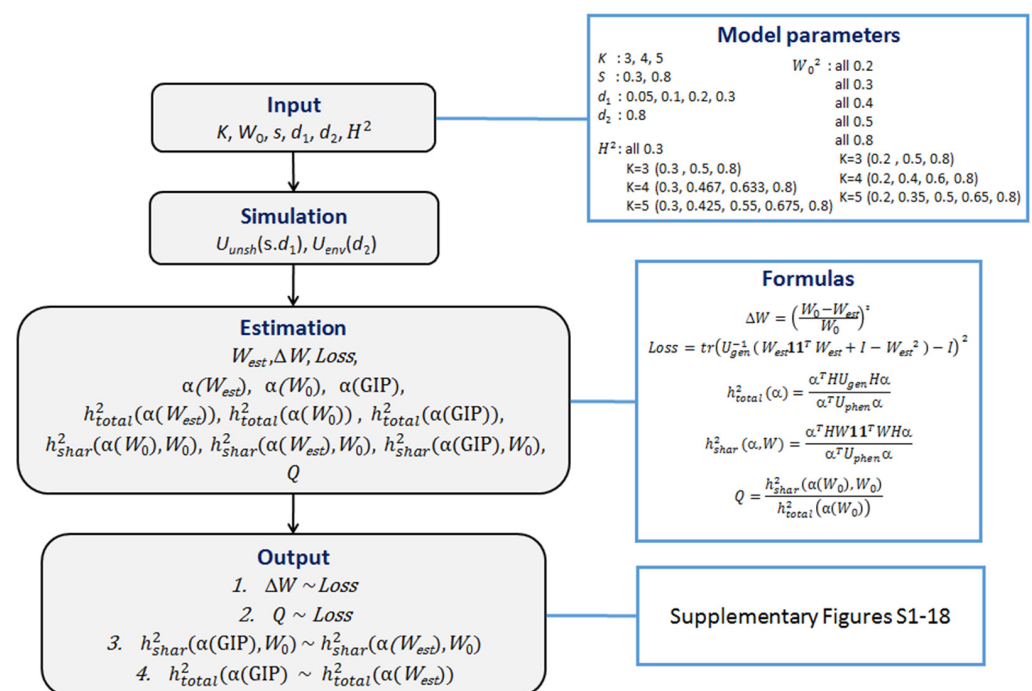


Figure 2. A schematic depicting the overall workflow of a simulation study. All details are given in the text.

For each fixed number,  $K$ , of the original traits and fixed heritability,  $h_i^2$  ( $i = 1, \dots, K$ ), of each trait, we simulated  $U_{gen}$ . To do this, we separately modelled two of its components caused by SGF and UGF as  $W\mathbf{1}^T W$  and  $\sqrt{I - W^2}U_{unsh}\sqrt{I - W^2}$ , respectively (see the “Model” box in Figure S1 in Supplementary Methods). Here  $\mathbf{1}$  is a  $(K \times 1)$  vector of units, and  $U_{unsh}$  is a  $(K \times K)$  matrix randomly generated using the parameters  $s$  and  $d_1$  (see Supplementary Methods). We then randomly generated the trait-trait correlation matrix  $U_{env}$  explained by the environmental factors, by giving the parameter  $d_2$  (see Supplementary Methods). Finally, we modeled a matrix of phenotypic correlations by using Model (2) with regard to simulated values  $W_0$ .

Using simulation data  $U_{phen}$ ,  $U_{gen}$  and  $H^2$ , we estimated  $W_{est}$  and calculated its squared relative difference with the simulated values of  $W_0$  ( $\Delta W$ ). We revealed a dependence of  $\Delta W$  on the loss function ( $Loss$ ). The  $Loss$  value characterizes the difference between  $U_{gen}$  and the auxiliary matrix  $V$ .

We then estimated  $\alpha$  in three ways: (1) using MaxSH and  $W_0$ , (2) using MaxSH and  $W_{est}$ , and (3) using the GIP method [2]. Based on these estimates, we formed three traits being the linear combinations of the original traits. For these combined traits, we calculated the total heritability and the heritability explained by SGF.

The simulated experiments were repeated 10,000 times for each set of parameters. The model parameters and formulas for all calculated values are shown in Figure 2.

## 2.4. Application to Real Data

### 2.4.1. Data Sets

We used three publicly available real data sets: anthropometric traits, psychiatric conditions and lipid concentrations, which contain five, four and three traits, respectively.

The group of anthropometric traits consisted of UK Biobank GWAS summary statistics obtained from the Neale lab (<http://www.nealelab.is/blog/2017/7/19/rapid-gwas-of-thousands-of-phenotypes-for-337000-samples-in-the-uk-biobank>, accessed on 1 September 2020) for people of European ancestry: BMI ( $N = 336,107$ ), weight ( $N = 336,227$ ), hip ( $N = 336,601$ ), waist circumference ( $N = 336,639$ ) and whole body fat mass ( $N = 330,762$ ).

The second dataset reflecting psychometric traits was constructed from GWAS results provided by the Psychiatric Genomics Consortium (<https://www.med.unc.edu/pgc/download-results/>, accessed on 1 September 2020) for bipolar disorder, BIP ( $N$  cases = 20,352;  $N$  controls = 31,358) [18], major depressive disorder, MDD ( $N$  cases = 43,204;  $N$  controls = 95,680; without UK Biobank and 23andMe data) [19] and schizophrenia, SCZ ( $N$  cases = 36,989;  $N$  controls = 113,075). Summary statistics for the fourth trait—subjective well-being ( $N = 110,935$ )—were derived from UK Biobank data from the Neale lab. All psychometric trait GWASs were conducted using samples from white Europeans.

The last dataset corresponding to lipid traits was formed using GWAS data for European participants from the Global Lipid Genetics Consortium (<http://csg.sph.umich.edu/willer/public/lipids2013/>, accessed on 1 September 2020) for LDL cholesterol ( $N = 173,082$ ), triglycerides ( $N = 177,860$ ), and total cholesterol ( $N = 187,365$ ).

Summary statistics for the three data sets were integrated and quality controlled by the GWAS-MAP platform developed by our group [20]. The GWAS-MAP database contains implemented software for quality control of GWAS results, the estimation of phenotypic correlations, and LD Score regression (LDSC) [16].

We conducted the quality control of all data and unified them within the GWAS-MAP platform [20]. We filtered all summary statistics by minor allele frequencies  $\geq 0.01$ . Additionally, we filtered GWAS results for BIP by imputation qualities  $\geq 0.9$ . We did not apply this filter to the other traits due to the absence of imputation quality in summary statistics data. Finally, using GWAS-MAP, we performed a correction for genomic control for all traits (including the original traits, SGIT and UGITs) with an LDSC intercept greater than 1 [16]. Thus, we corrected all traits from the psychometric dataset apart from MDD, all original anthropometric traits and their SGIT and lipid SGIT, as their LDSC intercept

exceeded 1 (see Supplementary Table S2a–c). Moreover, all SNPs with a  $p$ -value equal to 0 were excluded from analysis.

#### 2.4.2. Genetic Analysis

Pairwise phenotypic correlations between traits were computed using the GWAS-MAP platform described above. The used method is based on correlations between insignificant  $z$ -statistics for independent SNPs as previously described in [9]. SNP-based heritability and genetic correlation coefficients were estimated using the LD Score regression software [16] embedded in the GWAS-MAP platform. The significance threshold for genetic correlations was set at  $4.5 \times 10^{-4}$  ( $0.05/112$ , where 112 is the number of pairwise combinations between all original traits, their SGIT and UGITs in each dataset—between 11, nine and seven traits for anthropometry, psychometric and lipid traits, respectively).

The SHAHER analysis included checking if there was an SGI or not, the application of MaxSH, and conducting SGIT and UGIT GWASs. The threshold for confirming the existence of an SGI at the first stage was empirically set to 0.2.

For each dataset, we visualized the full genetic correlation matrices using the *corrplot()* function from the *corrplot* R package (v.0.84) [21]. We also placed the SNP-based heritability estimates on the diagonal and crossed out non-significant values.

Finally, we compared the GWAS results obtained for SGIT by MaxSH and GIP (the principal component analysis on the matrix of genetic covariances) [2].

#### 2.4.3. Gene set and Tissue/Cell Type Enrichment Analyses

We performed a gene set enrichment analysis and a tissue/cell type enrichment analysis combined with a gene prioritization using the Data-driven Expression Prioritized Integration for Complex Traits (DEPICT) tool v.1.1, release 194 [22]. We selected genome-wide significant SNPs ( $p$ -value  $< 5 \times 10^{-8}$ ) from summary statistics before the genomic control and applied the DEPICT software with default parameters (<https://data.broadinstitute.org/mpg/depict/>, accessed on 1 September 2020). The MHC region was excluded from analysis.

Next, for the gene set enrichment results, we calculated the number of significant enriched gene sets ( $FDR < 5\%$ ) and constructed an overlapping matrix in which each cell represents the number of overlapping gene sets for each pair of traits. For each pair of traits, we scaled the number of overlapping gene sets by the minimum number of significant gene sets for this pair of traits. The resulting matrix was visualized using the *corrplot* R-package, as described above.

#### 2.4.4. The Number of Original Traits Associated with SGIT Loci

We performed a clumping procedure to search for loci associated with each of the original traits, SGIT and UGITs at a genome-wide significance level of  $5 \times 10^{-8}$ . The associated locus was defined as a genomic region spanning 500 kb in either direction of the lead SNP. Those loci that were significantly associated with SGIT, but not with the original traits, were assumed to be new loci.

We expected that the loci associated with all the original traits used to obtain SGIT were likely to be SGF. To test this expectation, for each dataset we selected all independent loci that were significantly associated with at least one of the original traits and calculated the number of the original traits significantly associated with these loci. For the original anthropometric and lipid traits, we empirically set the significance threshold at  $p$ -value =  $1 \times 10^{-5}$ . For the psychometric traits, it was set at  $1 \times 10^{-3}$ . We then analyzed the SGIT  $p$ -values for the selected loci and constructed boxplots of  $-\log_{10}$  for them with regard to the number of the original traits significantly associated with these loci.

### 3. Results

#### 3.1. Simulation Study

To assess the MaxSH performance, we conducted simulation studies. We (1) assessed the accuracy of  $w$  estimates (using  $\Delta W$  metrics estimated as  $\left(\frac{w_0 - w_{est}}{w_0}\right)^2$ , where  $w_0$  and  $w_{est}$  are modeled and estimated  $w$ , respectively) with respect to the loss function given in Figure 2; (2) assessed the proportion of the shared heritability to the total heritability of SGIT (the  $Q$ -value) with respect to the loss function; and (3) compared the analytically predicted total/shared heritabilities of two traits: SGIT and the first component, GIP1, obtained by the GIP method [2]. The  $Q$ -value can be interpreted as the specificity metrics of SGIT: the closer the  $Q$ -value to 1, the lower the share of unshared heritability in the total heritability of SGIT. The simulation scenarios were based on six varying parameters that describe the properties of the genetic and phenotypic correlation matrices. Under each scenario, we considered two situations where all traits have the same  $w^2$  and different  $w^2$ s. To distinguish between these situations, we will hereinafter write either “ $w^2$ ” or “different  $w^2$ s”. In total, we performed 10,000 iterations of simulations for each of 288 scenarios.

The full results are presented in Figures S1–S18 in Supplementary Results. For all scenarios, there are few general patterns: (1) the higher simulated  $w$  values, the higher the accuracy of the  $w$  estimates, (2) the accuracy of the  $w$  estimates and the  $Q$ -value increase with an increasing in the number,  $K$ , of traits, (3) for all scenarios with  $w^2 > 0.8$ ,  $\Delta W$  was very low ( $<0.025$ ) and the  $Q$ -value was more than 90%.

For all scenarios with three traits, the accuracy of the  $w$  estimates was generally low:  $\Delta W$  was not higher than 0.7 for scenarios with  $w^2 = 0.2$  and 0.3, although at  $w^2$  equal to or higher than 0.4,  $\Delta W$  was less than 0.2. The  $Q$ -value was higher than 60% for almost all scenarios with  $w^2 \geq 0.4$ .

For the scenarios with four and five traits, the accuracy of  $w$  estimates was higher:  $\Delta W < 0.15$  for  $w^2 \geq 0.4$  and  $\Delta W < 0.05$  for  $w^2 \geq 0.5$ . For the scenarios with  $w^2 \geq 0.5$ , the  $Q$ -value was more than 70% for four traits and more than 80% for five traits.

The selected results of comparison of the shared heritability of SGIT and the shared heritability of GIP1 for scenarios with  $s = 0.3$  are presented in Figure 3. The results with  $s = 0.8$  were similar to those with  $s = 0.3$ . They are presented in Figures S1–S18 in Supplementary Results. For three traits in almost all cases, the shared heritabilities of SGIT were higher than the corresponding heritabilities of GIP1, except for the scenarios with  $h^2 = 0.8$ . For four and five traits, the shared heritabilities of SGIT were higher than the corresponding heritabilities of GIP1 under all scenarios, except for the scenarios with  $h^2 = 0.8$ . In the scenarios with  $h^2 = 0.8$ , the shared heritabilities of SGIT were higher than those of the GIP1 at  $w^2 \geq 0.5$ . The patterns of the total heritabilities for all scenarios reproduced the corresponding patterns of the shared heritabilities (Figures S1–S18 in Supplementary Results).

In summary, the performance of MaxSH was suitable at  $w^2 \geq 0.5$  and when the number of traits was higher than or equal to four. In the case of small  $w$  or three traits, the results of MaxSH should be interpreted with caution.

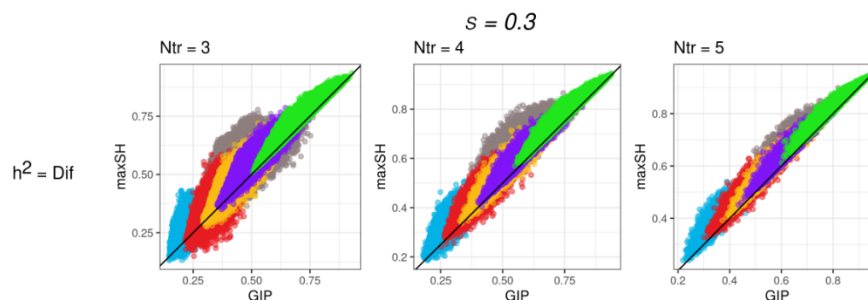
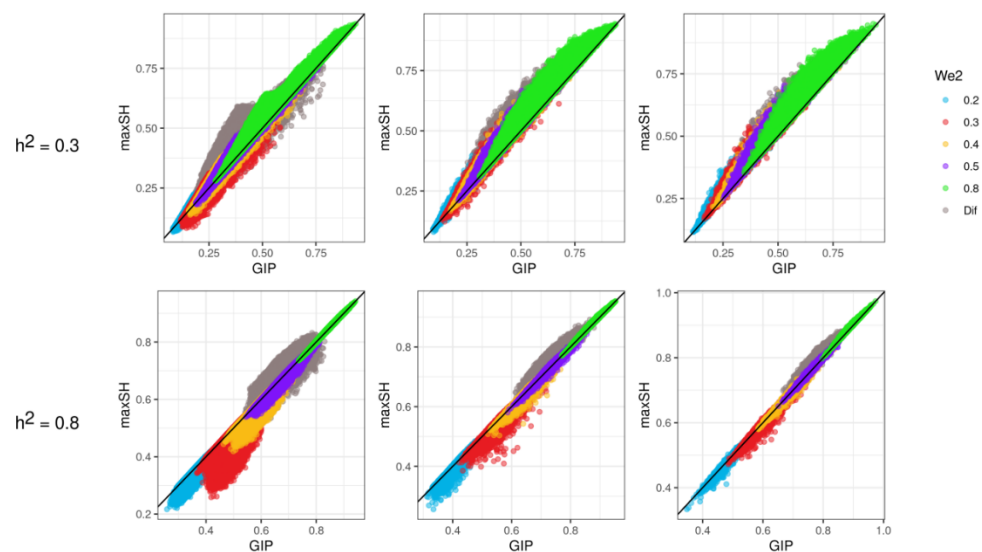


Figure 3. Cont.





**Figure 3.** The plot of the shared heritability of SGIT (the maxSH method) versus the shared heritability of GIP1 (the GIP method). For different  $h^2$  and different number of traits (Ntr) for  $s = 0.3$ .

### 3.2. Real Data Assessment

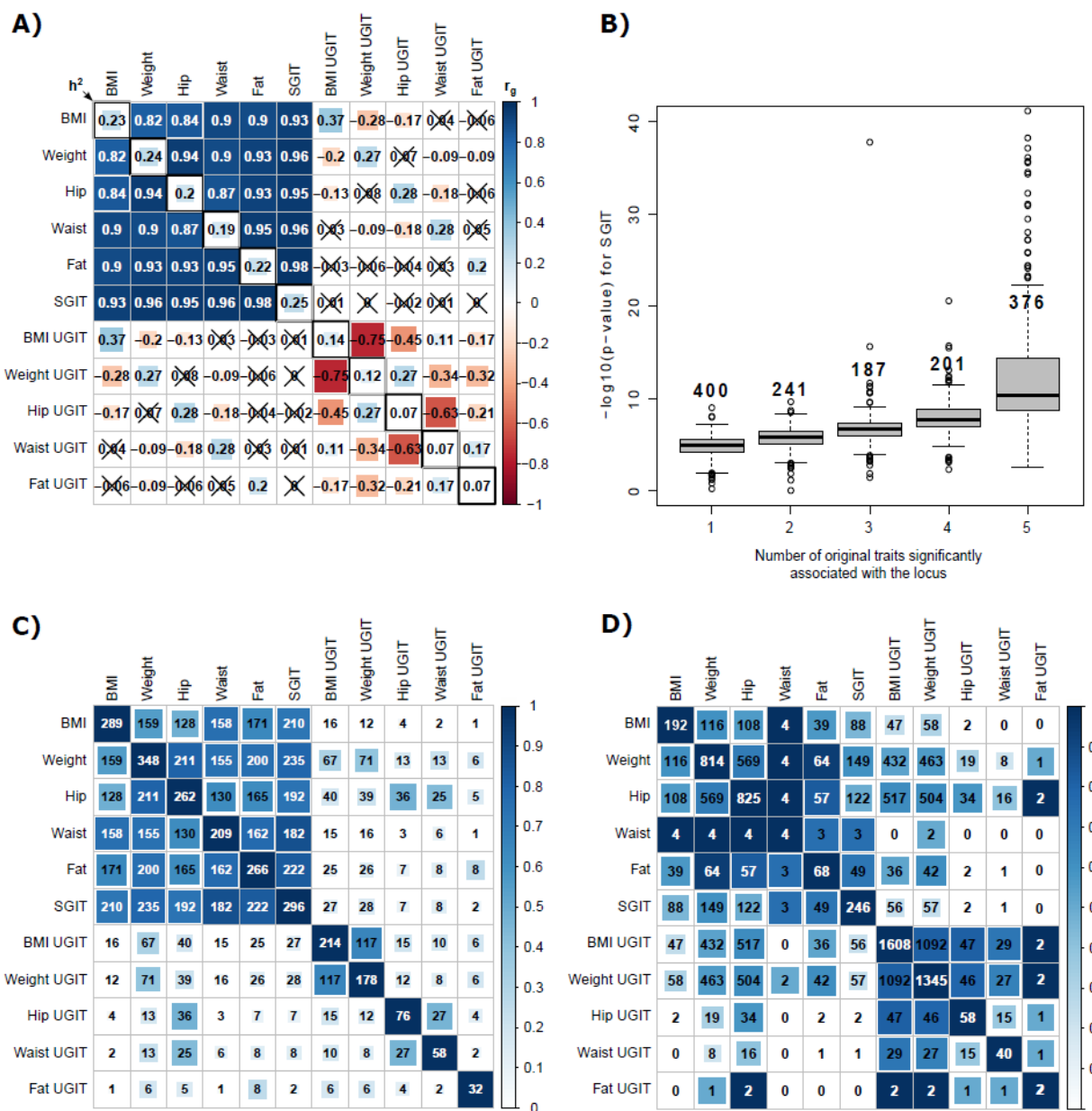
We applied SHAHER to three datasets: anthropometric (five traits), psychometric (four traits) and lipid traits (three traits). We should note that the performance of SHAHER applied to three traits is limited (see simulation results), yet still passable, although the results should be interpreted with caution. The number of identified loci for each trait for each data set is given in Table 1. We present SHAHER results for anthropometric traits in the main text as an example. The full results for the psychometric and lipid traits are presented in Supplementary Results.

**Table 1.** Number of significant loci ( $p$ -values  $< 5 \times 10^{-8}$ ) identified for each trait applying SHAHER for three data sets.

Trait Name	Number of Significant Loci		
	Real Trait	SGIT *	UGIT *
<b>Anthropometric Traits</b>			
<b>BMI</b>	289	296 (210)	214 (16)
<b>Weight</b>	348	296 (235)	178 (71)
<b>Hip</b>	262	296 (192)	76 (36)
<b>Waist</b>	209	296 (182)	58 (6)
<b>Fat</b>	266	296 (222)	32 (8)
<b>Psychometric Traits</b>			
<b>BIP</b>	12	57 (8)	2 (0)
<b>MDD</b>	3	57 (0)	2 (1)
<b>SCZ</b>	92	57 (26)	2 (0)
<b>Happiness</b>	0	57 (0)	1 (0)
<b>Lipid Traits</b>			
<b>LDL</b>	85	97 (69)	43 (31)
<b>Triglycerides</b>	71	97 (26)	59 (30)
<b>Cholesterol</b>	101	97 (84)	51 (21)

\* number of loci overlapping with those identified using the original trait is given in parentheses.

At the first step, we confirmed that SGI exists for five traits. At the second step, we determined the  $\alpha$  and  $\gamma$  coefficients and their CI (see Supplementary Table S1a). At the third step, we applied sumCOT and obtained GWAS results for SGIT and UGITs (see Supplementary Table S2a for heritability estimates and LD score regression intercepts). SHAHER results are presented in Figure 4.



**Figure 4.** Results of the application of SHAHER to anthropometric traits. (A) The heatmap of genetic correlations between the original, SGI and UGI traits. The number, color strength and size of the squares in the matrix show the values of the correlation coefficients between the traits. The diagonal elements represent heritabilities. Crossed out values indicate insignificant correlations. (B) Boxplots of  $-\log_{10}(p\text{-value})$  for SGIT with respect to the number of the original traits significantly associated with the locus. Two outliers for loci with  $-\log_{10}(p\text{-value}) > 40$  are omitted. The number at the top of the boxplot corresponds to the number of significant SNPs. (C) The heatmap of the numbers of overlapping loci between traits. The numbers in the cells represent the absolute numbers of overlapping loci. The color strength and size of the squares in the cells show the relative scaled number of overlapping loci (on a scale from 0 to 1). The diagonal elements represent the number of loci found for every trait. (D) The heatmap of the numbers of overlapping gene sets between traits. The color strength and size of the squares in the cells show the relative scaled number of overlapping gene sets (on a scale from 0 to 1). The diagonal elements represent the number of gene sets found for every trait.

Figure 4A demonstrates genetic correlations between all pairs of the original anthropometric traits, SGIT and UGITs. All the original traits were positively correlated with  $r > 0.82$ . We did not observe any significant genetic correlation between SGIT and UGITs. Moreover, we did not observe additional SGI among UGITs, which was expected. The heritabilities of UGITs varied from 0.07 to 0.14.

We revealed a dependence of the SGIT  $p$ -value from the number of the original traits significantly associated with the locus (Figure 4B). It clearly shows that the loci associated with all the original traits have lower SGIT  $p$ -values than the other loci.

Joint clumping of 11 traits (five original traits, five UGITs and SGIT) resulted in 820 genome-wide significantly associated loci ( $p$ -value  $< 5 \times 10^{-8}$ , Supplementary Table S3a). If a locus was not significantly associated with any of the original traits, it was considered new. SGIT was significantly associated with 296 SNPs. We detected no new loci among SGIT loci. The clumping of UGITs revealed 379 loci, of which 246 were new. At the same time, the clumping of only original traits allowed 574 loci to be detected, of which 187 could not be detected by analyzing SGIT or UGITs. Thus, the joint analysis of SGIT and UGITs increased the number of associated loci by more than 42.8%. Figure 4C reflects the overlapping between significantly associated loci for 11 analyzed traits. There is a weak albeit non-zero overlap between loci for UGITs and SGIT, although the genetic correlation between them is zero. It could be due to the conservative settings of the clumping procedure, which tends to clump together closely located loci, and is due to some level of unspecificity of the SHAHER.

Next, we checked how enriched gene sets overlap between SGIT, UGITs and the original traits (see Figure 4D). Significant results (FDR  $< 5\%$ ) of enriched gene sets and tissue enrichment analyses are presented in Supplementary Table S4. As expected, the heatmap of the overlapping gene sets looks similar to the heatmap of genetic correlations and the heatmap of the overlapping loci. Moreover, there was almost no overlap between SGIT and any UGIT. For the original traits, the number of enriched gene sets varied a lot: from four for the waist to 825 for the hip circumference. For BMI UGIT, the number of enriched gene sets was 1608, which was almost ten times the value for BMI (192).

Finally, we obtained GIP1 GWAS statistics and calculated the genetic correlations between SGIT and GIP1. The genetic correlation was higher than 0.97.

#### 4. Discussion

We developed a new fast and efficient framework which allows us to decompose the heritability of each trait from a given set of traits into two components. One of them is explained by shared genetic factors common to all traits. Another one is explained by unshared genetic factors specific for each trait. The framework not only decomposes heritability but also identifies SNPs associated with the shared and unshared genetic effects. To our knowledge, this framework is unparalleled. It has an additional advantage: it uses GWAS summary statistics obtained for original traits and does not require raw genotype or phenotype data.

We compared the performances of MaxSH and GIP in identifying the shared genetic components. GIP calculates the linear combination coefficients via the eigenvalues of the genetic covariance matrix and can be considered a close approximation to MaxSH. In our simulations, GIP and MaxSH were similar in almost all scenarios, with MaxSH being somewhat superior in terms of the power (total heritability) and quality (shared heritability). If obtaining genetically independent phenotypes is not the aim, we suggest using SHAHER, because it is more robust and gives additional metrics like SGI contributions to the heritability of the original traits.

The framework is computationally effective. The stage using sumCOT is the most time consuming. However, it only takes several minutes for an average computer to conduct a GWAS of a linear combination of traits with 6M SNPs using a C++ implementation of the sumCOT. MaxSH, based on numerical optimization procedures, and the other parts of the framework take seconds.

The proposed sumCOT method can be applied as an independent tool to address additional tasks. One of them is making a summary-level adjustment of traits by other traits using the same scheme as was used for obtaining the UGIT GWAS statistics. This can be helpful, for example, for ridding the studied trait's genetic component of the genetic component that was caused by the confounding or unaccounted effects of assortative mating or family effects, which is quite a problem in GWAS at the biobank scale [15,23]. Another task is a GWAS for the trait that appears as a linear combination of the original traits. The sumCOT method is robust to differences in sample sizes used for GWASs of original traits and is applicable to different GWAS models (Cox, linear or logistic).

The main interest in the application of the SHAHER framework lies in the possibility of obtaining novel biological insights into a trait's heritability composition. This can be achieved by the application of a huge variety of in-silico follow-up techniques to SGIT and UGITs. SGIT is of interest by itself, but we also emphasize the importance of the comparison of shared and unshared effects for each trait. In our real data application, the most remarkable case is BMI in the set of anthropometric traits (see Figure 4C). We found 246 and 1608 significantly enriched gene sets for SGIT and the UGIT of BMI, respectively, with negligible overlapping between them of size 56. By analyzing BMI only, we would have detected only 192 enriched gene sets. By analyzing each of the impacts separately, we dramatically increased the number of observed unique gene sets (1798 in total for both SGI and UGI). This means that each sub-phenotype controlled by SGF and UGF is less heterogeneous than the original trait. According to the significant gene sets, UGIT of BMI (see Supplementary Table S4) controls some structural changes in body compositions and bone formation, while SGIT is involved in some general signaling pathways and pathways related to nervous system development and probably to general psycho-social aspects of BMI, obesity and other anthropometric traits [24]. Note that all new loci were associated with UGITs. We can speculate that these new SNPs were detected due to the decreased genetic heterogeneity of UGITs compared to the original traits.

To validate the findings of SHAHER, we have compared the association results for anthropometric and psychometric traits with the biggest publicly available GWAS results for BMI and MDD. The idea is as follows: if the locus was not significant on the original trait but was detected on SGIT or corresponding UGIT, it will be detected on the original trait if the GWAS sample size is increased. For BMI we have used the largest meta-analysis of the UK Biobank and GIANT GWAS ( $N = 806,834$ ) to date [25] (see Supplementary Table S3a). Among 264 loci significant on SGIT and BMI UGIT but not on BMI, 57 loci (22%) became significantly associated with BMI in the biggest GWAS. If we consider only loci associated with SGIT, the validation ratio is higher: out of 86 loci, 49 (57%) were significant in the biggest GWAS. For MDD we have used the biggest meta-analysis of the UK Biobank and PGC GWAS ( $N = 500,199$ ) to date [26] (see Supplementary Table S3b). Among 58 loci significantly associated with SGIT and MDD UGIT but not with MDD, seven (12%) became significant in the biggest MDD GWAS. The similar validation of lipid loci was not performed since there was no bigger GWAS available in the open access. The validation of the loci on the biggest GWAS is not a proper replication, however it still greatly increases the confidence that detected by SHAHER loci are true positives.

Although SHAHER is effective, it has several limitations. First, when trait-trait genetic correlations are weak, it is expected that the contributions of these traits to the shared heritability will be small as well. In this case, MaxSH may overestimate these contributions. Secondly, the framework is applicable only if the number of traits is no less than three. In the case of three traits, the performance is limited and the SHAHER results should be interpreted with caution. We have shown in simulations and real dataset examples that MaxSH works better at higher numbers of genetically correlated traits being analyzed. However, an increase in the number of weakly correlated traits leads to a decrease in the proportion of SNPs associated with all traits simultaneously and to a decrease in the efficiency of the framework. Thirdly, although the set of SNPs identified by the SGIT GWAS is enriched for the SGF, each SNP should be interpreted with caution for whether

it is shared or not, because SHAHER has some level of nonspecificity. Finally, if any confounding effects were included in the GWAS of the original traits, these effects are amplified in the SGIT [15]. The confounding effects can be controlled easily using special methods like LD score regression [16], although this method fails to distinguish a polygenic component if the trait was measured in the sample with the assortative mating or family effects. Thus, we suggest a thorough check of the original GWAS for the presence of any effects of possible confounders before proceeding to SHAHER. In principle, if the LD score regression intercept was estimated, it is possible to correct for residual inflation by adjusting the standard errors of the effects by multiplying them by the square root of the intercept.

We should highlight the distinctive specificity of SHAHER which distinguishes it from existing approaches for multivariate analysis. There are a lot of frameworks that allow for the incorporation of several correlated traits in one analysis to increase the power of mapping [8–13]. Our framework is not aimed to increase the power of mapping itself (although empirically we showed that SHAHER has higher power compared to univariate analyses). Our framework is aimed to estimate the shared and unshared heritability and to identify the shared and unshared genetic factors. Therefore, we did not compare the power of SHAHER with the powers of existing approaches and do not expect it to have the highest power among them. Moreover, our definition of shared genetic factors is stricter than just the pleiotropy of all analyzed traits. This is why using multivariate approaches aimed to increase the power of mapping is not the optimal way to identify shared genetic factors.

In conclusion, we propose a novel effective framework for analysis of the shared genetic background for a set of genetically correlated traits using GWAS summary statistics. The framework allows us to obtain novel biological insights into the trait's genetic impact composition. By analyzing shared and unshared genetic impacts separately, we increased the number of identified loci and observed unique gene sets, identified genetic mechanisms that are common for all traits or specific for every trait. Of note, sumCOT can be used as a stand-alone method for obtaining GWAS results of the linear combination of the traits using their summary statistics.

**Supplementary Materials:** The following are available online at <https://www.mdpi.com/article/10.3390/genes13101694/s1>, Supplementary Tables S1–S6, Supplementary Results, Supplementary Methods. References [27–31] are cited in the supplementary materials.

**Author Contributions:** Y.A.T., G.R.S. and T.I.A. conceived and oversaw the study. Y.A.T., G.R.S., S.Z.S. and T.I.A. contributed to the design of the study and interpretation of the results. G.R.S. developed the MaxSH method, including the algorithm and program, and conducted simulation studies. P.R.H.J.T. developed the C++ version of sumCOMB and tested the developed framework. E.S.T., E.E.E., S.G.F., S.Z.S. and Y.A.T. wrote the source code for the framework and performed real data analyses. All authors have read and agreed to the published version of the manuscript.

**Funding:** The work of GRS was supported by the Russian Foundation for Basic Research (project 20-04-00464). The work of Y.A.T. and T.I.A. was supported by the Russian Science Foundation (RSF) grant and Government of the Novosibirsk region No. 22-15-20037. The work of E.E.E. was supported by the grant for the implementation of the strategic academic leadership program “Priority 2030” in Novosibirsk State University. The work of S.Z.S. was supported by budget project No. FWNR-2022-0020. The work of PRHJT was supported by the Medical Research Council Human Genetics Unit (MC\_UU\_00007/10).

**Institutional Review Board Statement:** Not applicable.

**Informed Consent Statement:** Not applicable.

**Data Availability Statement:** UK Biobank GWAS summary statistics for anthropometric traits and for subjective well-being: <http://www.nealelab.is/blog/2017/7/19/rapid-gwas-of-thousands-of-phenotypes-for-337000-samples-in-the-uk-biobank>, accessed on 1 September 2020. GWAS summary statistics for psychometric trait: the Psychiatric Genomics Consortium, <https://www.med.unc.edu/pgc/download-results/>, accessed date September 2020. GWAS for lipid traits: the Global Lipid Genetics Consortium, <http://csg.sph.umich.edu/willer/public/lipids2013/>, accessed on 1 September 2020. The biggest to date MDD GWAS meta-analysis of the UK Biobank and PGC: <https://www.nature.com/articles/s41467-021-25888-8>, accessed on 1 September 2020.

[//datashare.ed.ac.uk/handle/10283/3203](https://datashare.ed.ac.uk/handle/10283/3203), accessed on 1 August 2022. The biggest to date BMI GWAS meta-analysis of the UK Biobank and GIANT: [https://portals.broadinstitute.org/collaboration/giant/index.php/GIANT\\_consortium\\_data\\_files](https://portals.broadinstitute.org/collaboration/giant/index.php/GIANT_consortium_data_files), accessed on 1 August 2022.

**Conflicts of Interest:** P.R.H.J.T. is an employee of BioAge Labs. The remaining authors declare that they have no conflict of interest.

**Code Availability:** The SHAHER framework is implemented as a set of R/C++ scripts and is freely available at [https://github.com/Sodbo/shared\\_heridity](https://github.com/Sodbo/shared_heridity) (accessed on 1 September 2022).

## References

- Jiang, X.; Finucane, H.K.; Schumacher, F.R.; Schmit, S.L.; Tyrer, J.P.; Han, Y.; Michailidou, K.; Lesueur, C.; Kuchenbaecker, K.B.; Dennis, J. Shared heritability and functional enrichment across six solid cancers. *Nat. Commun.* **2019**, *10*, 431. [CrossRef] [PubMed]
- Tsepilov, Y.A.; Freidin, M.B.; Shadrina, A.S.; Sharapov, S.Z.; Elgaeva, E.E.; van Zundert, J.; Karssen, L.C.; Suri, P.; Williams, F.M.; Aulchenko, Y.S. Analysis of genetically independent phenotypes identifies shared genetic factors associated with chronic musculoskeletal pain conditions. *Commun. Biol.* **2020**, *3*, 329. [CrossRef]
- Sampson, J.N.; Wheeler, W.A.; Yeager, M.; Panagiotou, O.; Wang, Z.; Berndt, S.I.; Lan, Q.; Abnet, C.C.; Amundadottir, L.T.; Figueroa, J.D. Analysis of heritability and shared heritability based on genome-wide association studies for 13 cancer types. *JNCI J. Natl. Cancer Inst.* **2015**, *107*, djv279. [CrossRef] [PubMed]
- Brainstorm, C.; Anttila, V.; Bulik-Sullivan, B.; Finucane, H.K.; Walters, R.K.; Bras, J.; Duncan, L.; Escott-Price, V.; Falcone, G.J.; Gormley, P.; et al. Analysis of shared heritability in common disorders of the brain. *Science* **2018**, *360*, aap8757. [CrossRef]
- Yang, Y.; Zhao, H.; Heath, A.C.; Madden, P.A.; Martin, N.G.; Nyholt, D.R. Shared genetic factors underlie migraine and depression. *Twin Res. Hum. Genet.* **2016**, *19*, 341–350. [CrossRef] [PubMed]
- Wright, S. Correlation and Causation. *J. Agric. Res.* **1921**, *XX*, 557–585.
- Rijsdijk, F.V.; Sham, P.C. Analytic approaches to twin data using structural equation models. *Brief. Bioinform.* **2002**, *3*, 119–133. [CrossRef]
- Galesloot, T.E.; Van Steen, K.; Kiemeny, L.A.; Janss, L.L.; Vermeulen, S.H. A comparison of multivariate genome-wide association methods. *PLoS ONE* **2014**, *9*, e95923. [CrossRef] [PubMed]
- Stephens, M. A unified framework for association analysis with multiple related phenotypes. *PLoS ONE* **2013**, *8*, e65245. [CrossRef] [PubMed]
- Yang, X.; Zhang, S.; Sha, Q. Joint analysis of multiple phenotypes in association studies based on cross-validation prediction error. *Sci. Rep.* **2019**, *9*, 1073. [CrossRef]
- Turley, P.; Walters, R.K.; Maghziyan, O.; Okbay, A.; Lee, J.J.; Fontana, M.A.; Nguyen-Viet, T.A.; Wedow, R.; Zacher, M.; Furlotte, N.A. Multi-trait analysis of genome-wide association summary statistics using MTAG. *Nat. Genet.* **2018**, *50*, 229–237. [CrossRef] [PubMed]
- Fatumo, S.; Carstensen, T.; Nashiru, O.; Gurdasani, D.; Sandhu, M.; Kaleebu, P. Complimentary methods for multivariate genome-wide association study identify new susceptibility genes for blood cell traits. *Front. Genet.* **2019**, *10*, 334. [CrossRef]
- Ning, Z.; Tsepilov, Y.A.; Sharapov, S.Z.; Wang, Z.; Grishenko, A.K.; Feng, X.; Shirali, M.; Joshi, P.K.; Wilson, J.F.; Pawitan, Y.; et al. Nontrivial Replication of Loci Detected by Multi-Trait Methods. *Front. Genet.* **2021**, *12*, 627989. [CrossRef] [PubMed]
- Pasaniuc, B.; Price, A.L. Dissecting the genetics of complex traits using summary association statistics. *Nat. Rev. Genet.* **2017**, *18*, 117–127. [CrossRef] [PubMed]
- Timmers, P.R.H.J.; Tiys, E.S.; Sakaue, S.; Akiyama, M.; Kiiskinen, T.T.J.; Zhou, W.; Hwang, S.-J.; Yao, C.; Kamatani, Y.; Deelen, J.; et al. Mendelian randomization of genetically independent aging phenotypes identifies LPA and VCAM1 as biological targets for human aging. *Nat Aging* **2022**, *2*, 19–30. [CrossRef]
- Bulik-Sullivan, B.; Finucane, H.K.; Anttila, V.; Gusev, A.; Day, F.R.; Loh, P.R.; ReproGen Consortium; Psychiatric Genomics Consortium; Genetic Consortium for Anorexia Nervosa of the Wellcome Trust Case Control Consortium 3; Duncan, L.; et al. An atlas of genetic correlations across human diseases and traits. *Nat. Genet.* **2015**, *47*, 1236–1241. [CrossRef]
- Winkler, T.W.; Day, F.R.; Croteau-Chonka, D.C.; Wood, A.R.; Locke, A.E.; Mägi, R.; Ferreira, T.; Fall, T.; Graff, M.; Justice, A.E. Quality control and conduct of genome-wide association meta-analyses. *Nat. Protoc.* **2014**, *9*, 1192–1212. [CrossRef]
- Stahl, E.A.; Breen, G.; Forstner, A.J.; McQuillin, A.; Ripke, S.; Trubetskoy, V.; Mattheisen, M.; Wang, Y.; Coleman, J.R.; Gaspar, H.A. Genome-wide association study identifies 30 loci associated with bipolar disorder. *Nat. Genet.* **2019**, *51*, 793–803. [CrossRef] [PubMed]
- Wray, N.R.; Ripke, S.; Mattheisen, M.; Trzaskowski, M.; Byrne, E.M.; Abdellaoui, A.; Adams, M.J.; Agerbo, E.; Air, T.M.; Andlauer, T.M. Genome-wide association analyses identify 44 risk variants and refine the genetic architecture of major depression. *Nat. Genet.* **2018**, *50*, 668–681. [CrossRef] [PubMed]
- Gorev, D.; Shashkova, T.; Pakhomov, E.; Torgasheva, A.; Klaric, L.; Severinov, A.; Sharapov, S.; Alexeev, D.; Aulchenko, Y. GWAS-MAP: A platform for storage and analysis of the results of thousands of genome-wide association scans. In Proceedings of the Bioinformatics of Genome Regulation and Structure Systems Biology (BGRS/SB-2018), Novosibirsk, Russia, 20–25 August 2018; p. 43.
- Wei, T.; Simko, V. R Package ‘Corrplot’: Visualization of a Correlation Matrix. (Version 0.92). 2021. Available online: <https://github.com/taiyun/corrplot> (accessed on 1 September 2022).

22. Pers, T.H.; Karjalainen, J.M.; Chan, Y.; Westra, H.-J.; Wood, A.R.; Yang, J.; Lui, J.C.; Vedantam, S.; Gustafsson, S.; Esko, T. Biological interpretation of genome-wide association studies using predicted gene functions. *Nat. Commun.* **2015**, *6*, 5890. [[CrossRef](#)] [[PubMed](#)]
23. Howe, L.J.; Lawson, D.J.; Davies, N.M.; Pourcain, B.S.; Lewis, S.J.; Smith, G.D.; Hemani, G. Genetic evidence for assortative mating on alcohol consumption in the UK Biobank. *Nat. Commun.* **2019**, *10*, 5039. [[CrossRef](#)] [[PubMed](#)]
24. Marcellini, F.; Giuli, C.; Papa, R.; Tirabassi, G.; Faloia, E.; Boscaro, M.; Polito, A.; Ciarapica, D.; Zaccaria, M.; Mocchegiani, E. Obesity and body mass index (BMI) in relation to life-style and psycho-social aspects. *Arch. Gerontol. Geriatr.* **2009**, *49*, 195–206. [[CrossRef](#)] [[PubMed](#)]
25. Pulit, S.L.; Stoneman, C.; Morris, A.P.; Wood, A.R.; Glastonbury, C.A.; Tyrrell, J.; Yengo, L.; Ferreira, T.; Marouli, E.; Ji, Y. Meta-analysis of genome-wide association studies for body fat distribution in 694 649 individuals of European ancestry. *Hum. Mol. Genet.* **2019**, *28*, 166–174. [[CrossRef](#)] [[PubMed](#)]
26. Howard, D.M.; Adams, M.J.; Clarke, T.-K.; Hafferty, J.D.; Gibson, J.; Shiri, M.; McIntosh, A. Genome-wide meta-analysis of depression in 807,553 individuals identifies 102 independent variants with replication in a further 1,507,153 individuals. *BioRxiv* **2018**, 6288, 433367.
27. Falconer, D.; Mackay, T. *Introduction to Quantitative Genetics*, 2nd ed.; Longman: New York, NY, USA, 1981; p. 315.
28. Khodadadi, Z.; Tarami, B. Robust Empirical Bayes Estimation of the Elliptically Countoured Covariance Matrix. *J. Math. Ext.* **2011**, *5*, 31–46.
29. Konno, Y. Estimation of Multivariate Complex Normal Covariance Matrices Under an Invariant Quadratic Loss. *Commun. Stat. Theory Methods* **2010**, *39*, 1490–1497. [[CrossRef](#)]
30. Oualkacha, K.; Labbe, A.; Ciampi, A.; Roy, M.A.; Maziade, M. Principal components of heritability for high dimension quantitative traits and general pedigrees. *Stat. Appl. Genet. Mol. Biol.* **2012**, *11*. [[CrossRef](#)]
31. Mardia, K.; Kent, J.; Bibby, J. *Multivariate Analysis*; Academic Press Inc.: London, UK, 1979; Volume 15, p. 518.