1    **A novel framework for analysis of the shared genetic background of**

2    **correlated traits**

3

4

5    Gulnara R. Svishcheva[1,2], Evgeny S. Tiys[1,3], Elizaveta E. Elgaeva[1,3], Sofia G.

6    Feoktistova[1,3], Paul R. H. J. Timmers[4,5], Sodbo Zh. Sharapov[1,3], Tatiana I. Axenovich[1,3],

7    Yakov A. Tsepilov[1,3*]

8

9    [1] Institute of Cytology and Genetics, Siberian Branch of the Russian Academy of Sciences,

10    Novosibirsk, Russia

11    [2] Vavilov Institute of General Genetics, Russian Academy of Sciences, Moscow, Russia

12    [3] Novosibirsk State University, Novosibirsk, Russia

13    [4] MRC Human Genetics Unit, MRC Institute of Genetics and Cancer, University of

14    Edinburgh, Edinburgh, UK

15    [5] Centre for Global Health Research, Usher Institute, University of Edinburgh, Edinburgh,

16    UK

17

18    *correspondence to: tsepilov@bionet.nsc.ru

19    Keywords: GWAS, shared genetic component, linear combination of traits, shared

20    heritability, proportion of heritability explained by SGF

21

22 **Abstract**

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

35

## 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 this shared genetics between traits can help discover pleiotropic interactions, common genes and pathways, and identify genetic effects that are unique 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. Write in 1921 [6]. There are widely used classic twin designs to have this problem solved. They are based on structural equation modelling, in particular, multivariate pathway models assuming the existence of the genetic influences common for all traits and unique 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 impacts.

There are several terms for these common and unique genetic impacts. Hereafter 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.

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). Also, if a locus belonging in fact to SGF was not identified as having pleiotropic effects on all traits due to a limited statistical power of the analysis, then the locus can be erroneously assigned to UGF. 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,

SHAHER framework, 2021

68  there is no specific method that could be good for both variance component decomposition

69  and identification of SGF and UGF.

70      We had previously developed a method for obtaining genetically independent

71  phenotypes (GIPs) [2]. This method is based on the calculation of the principal components

72  using genetic rather than phenotypic correlations. We applied this method to genetically

73  correlated pain phenotypes and aging related phenotypes and showed that the first GIP

74  component, GIP1, that explains the largest proportion of the genetic variance probably could

75  be interpreted as SGI [2, 15]. This makes GIP promising for identification of loci attributed to

76  SGF. However, this method was not designed specifically for SGI analysis. In addition, no

77  specific experiments have been performed to validate the approach or to estimate its

78  statistical properties.

79      Here, we present a novel general framework for the estimation of shared and unshared

80  heritability and identification of the shared and unshared genetic factors using the summary

81  statistics of original traits. The essence of our approach is to find the optimum linear

82  combination of traits which has the maximum proportion of its genetic variance explained

83  by the SGF. We validated our framework using simulation studies under different scenarios,

84  by comparing it with the developed GIP approach, and assessed its performance using three

85  real datasets: anthropometric indices, psychiatric disorders and conditions, and lipid

86  concentrations.

87

88  **Results**

89      **Abbreviations and terms**

90  SHAHER: a framework for the estimation of the shared and unshared heritability of studied

91  traits and identification of the shared and unshared genetic factors using the summary

92  statistics of original traits.

93  SGI: shared genetic impact.

94  UGI: unshared genetic impact.

95  SGF (shared genetic factors): genetic factors involved in the control of all studied traits and

96  whose effects are collinear between all studied traits; SGI is due to SGF.

97  Shared heritability: the proportion of the trait variance explained by SGF.

98  SGIT (shared genetic impact trait): a trait defined as a linear combination of original traits

99  maximizing its shared heritability.

100  $\alpha$: the coefficients of an optimum linear combination of original traits for building the SGIT.

101     UGF (unshared genetic factors): the residual genetic factors of an original trait after

102     exclusion of the SGF; UGI is due to UGF.

103     Unshared heritability: the proportion of the trait variance explained by UGFs.

104     UGIT (unshared genetic impact trait): an original trait after adjustment for the SGIT.

105     $\gamma$: the coefficients of a linear combination of original traits for building the UGIT.

106     MaxSH (MAXimization of Shared Heritability): a method for estimating the shared and

107     unshared heritability of each trait and calculating the coefficients of the linear combination

108     of the original traits: $\alpha$, to build the SGIT, and $\gamma$, to build the UGITs.

109     sumCOT (summary-level GWAS for linear Combination of Traits): a method to compute

110     GWAS summary statistics for the linear combination of the original traits using their

111     summary statistics.

112

113          **Shared heredity model**

114          We adopted a commonly used multivariate pathway model [7] in terms of SGF and

115     UGF. We call it the 'shared heredity model'. For simplicity, we consider SGF and UGF as

116     biallelic SNPs and consider a sample of $N$ unrelated individuals measured for $K$ traits and

117     genotyped for $M$ SNPs. For a standardized normal trait, $y$ ($N \times 1$), the traditional polygenic

118     (null) model takes the form: $y = G\beta + \varepsilon$, where $G$ is an ($N \times M$) matrix of standardized

119     genotypes; $\beta$ ($M \times 1$) and $\varepsilon$ ($N \times 1$) are genetic and non-genetic random effects, respectively;

120     $\beta \sim N(\mathbf{0}, h^2 I_M)$ and $\varepsilon \sim N(\mathbf{0}, (1-h^2)I_N)$, where $\mathbf{0}$ is a null mean vector, $h^2$ is the trait

121     heritability, and $I$ is an identity matrix of the given dimension. For unrelated individuals, we

122     expect $y \sim N(\mathbf{0}, I_N)$.

123          We propose to divide $M$ SNPs into two non-overlapping SNP sets with sizes $M_0$ and

124     $M_1$ ($M_0 + M_1 = M$). The set of $M_0$ SNPs called 'SGF' includes only those SNPs whose effects

125     on all traits are collinear. The set of $M_1$ SNPs consists of the other SNPs, which do not have

126     shared joint influence on all traits at once, this set being called 'UGF'. In accordance with $M$,

127     $G$ is divided into two matrices, $G_0$ ($N \times M_0$) and $G_1$ ($N \times M_1$). To decompose every trait into

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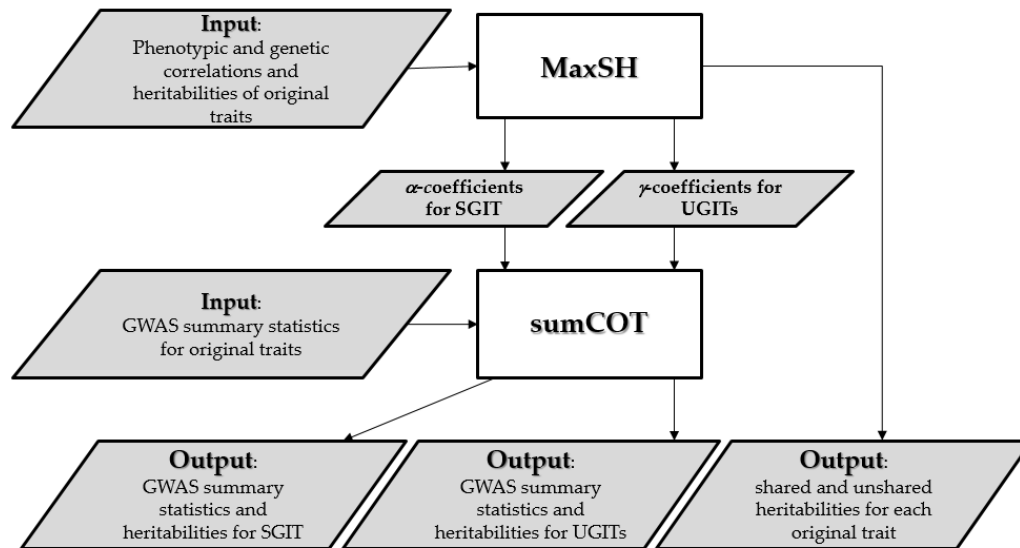
129     terms of $G_0$ and $G_1$

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

SHAHER framework, 2021

130  Here, the first and second terms are genetic components explained by SGF and UGF,

131  respectively, which are assumed independent. In the first term, $b_{0_i}$ is an ($M_0 \times 1$) vector of

132  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

133  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

134  trait explained by SGF. Here $w_i$ is a non-zero trait-specific multiplier: $w_i^2$ denotes the

135  proportion of $h_i^2$ explained by SGF; the value of $w_i$ can be positive and negative, indicating

136  the direction of the SGF effect on the $i$-th trait. $G_0 \beta_0$ is the so-called shared genetic impact

137  or SGI. In the second term of Model (1), $b_{1_i}$ is an ($M_1 \times 1$) vector of UGF effects, which can

138  be presented as $b_{1i} = \beta_{1i} \sqrt{(1 - w_i^2) h_i^2}$, $\beta_{1i} \sim N(0, I_{M_1})$. In contrast to $\beta_0$, $\beta_{1_i}$ are different

139  for different traits, moreover they are not collinear. For illustrative purposes, we rewrote

140  Equation (1) as:

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

141

142  **Overview of the SHAHER framework**



143
144  **Figure 1. Flowchart of the SHAHER framework.** Details are given in the text.
145

146  For analyses of the SGI and UGI on a set of correlated traits, we propose an effective

147  multi-stage framework named SHAHER (see Figure 1). The concept of the framework is

148 first to partition the genetic basis of each original trait into two components: one shared by
149 all the original traits and one shared not by all the original traits, and then to identify the
150 SNPs that contribute to these genetic components. To do this, we propose to construct new
151 traits: (1) an SGIT as a linear combination of original traits, which has the maximum
152 possible heritability explained by the SGF, and (2) UGITs as linear combinations of the
153 original traits, which are obtained by adjusting the original traits for the SGIT. This means
154 that the genetic basis of the UGITs is predominantly determined by the UGF.

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

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

166 *The MaxSH stage.* To determine the $\alpha$ and $\gamma$ coefficients for the linear combinations of
167 the original traits to build the SGIT and UGITs, we developed the MaxSH method, which is
168 based on the correlation component model given below. This model partitions the
169 phenotypic correlation matrix, $U_{phen}$, into environmental and genetic components, $U_{env}$ and
170 $U_{gen}$, respectively, the latter being further subdivided into two components caused by the
171 SGF and UGF:

$$U_{phen} = \underbrace{\sqrt{H^2}U_{gen}\sqrt{H^2}}_{genetic\ component} + \underbrace{\sqrt{I-H^2}U_{env}\sqrt{I-H^2}}_{environmental\ component}$$

$$U_{gen} = \underbrace{W\mathbf{1}\mathbf{1}^{\mathrm{T}}W}_{due\ to\ SGF} + \underbrace{\sqrt{I-W^2}U_{unsh}\sqrt{I-W^2}}_{due\ to\ UGF}$$

(2)

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

7
SHAHER framework, 2021

175        First of all, using only the genetic correlation matrix, $U_{gen}$, we estimate the proportion

176     of heritability of every trait explained by SGF ($W$). To do this, we minimize the difference

177     between $U_{gen}$ and the auxiliary matrix $V$. This matrix is built using formula (2), with the

178     identity matrix used instead of $U_{unsh,}$. The second task is to determine the $\alpha$-coefficients,

179     which is solved by maximizing the shared heritability of the SGIT. This task is analytically

180     solved as

$$a = \frac{U_{phen}^{-\frac{1}{2}} HW\mathbf{1}}{\sqrt{\mathbf{1}^T W H U_{phen}^{-1} HW\mathbf{1}}}.$$

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

182        After determining the $\alpha$-coefficients and building the SGIT, we build a UGIT for every

183     trait using the residual regression equation $UGIT_i = y_i - SGIT*c_i$, where $c_i$ is the impact of

184     the SGIT on the $i$-th original trait, defined as

185     $$c_i = cov_{gen}(y_i, SGIT)/h_{SGIT}^2.$$

186     Here $cov_{gen}$ denotes a genetic covariance. Note that we should use genetic rather than

187     phenotypic covariances, as our goal is to adjust only the genetic components of the original

188     traits.

189        Since the SGIT is the linear combination of the original traits, the UGITs are linear

190     combinations of the original traits, too. The coefficients of these linear combinations called

191     the $\gamma$-coefficients form the matrix of the $\gamma$-coefficients $\Gamma = (I_K - \alpha c^T)$, where the $i$-th

192     column of $\Gamma$ corresponds to linear combination coefficients for building the $i$-th UGIT.

193

194        *The sumCOT stage.* This stage is aimed at obtaining GWAS summary statistics for

195     the SGIT and UGITs using the previously determined $\alpha$ and $\gamma$ coefficients, GWAS summary

196     statistics (Z-scores, allele frequencies and sample sizes for each SNP) for the original traits

197     and the matrix of phenotypic correlations. The method can use Z scores obtained from any

198     regression model and allows for varying sample sizes and sample overlap between traits.

199     This sample overlap is incorporated into the estimation of the matrix of phenotypic

200     correlations. In short, the SNP effects for combined trait are calculated by summing effect

201     estimates from the individual trait GWASes, each multiplied by their corresponding linear

202     coefficient ($\alpha$ or $\gamma$), and standardized by the expected variance. The standard errors of the

203     SNP effect are calculated using variance-covariance arithmetic, taking into account the

204     phenotypic covariance between GWAS results to adjust for the sample overlap. Effective

SHAHER framework, 2021

205    sample sizes are then estimated based on the median Z statistic and allele frequencies by

206    solving Equation (1) in [17].

207    At the final stage, SHAHER checks for the correctness of the output. In particular, we

208    anticipate that UGITs do not have a shared genetic basis. This is verified by applying

209    MaxSH to the matrix of correlations between UGITs.

210    To summarize, our framework estimates shared and unshared heritabilities for each of

211    the studied original traits and produces GWAS summary statistics for the SGIT and UGITs

212    as outputs.

213    The full details and mathematical formulae of SHAHER are in *Supplementary*

214    *Methods*.

## Simulation study

216    To assess the MaxSH performance, we conducted simulation studies. We (1) assessed

217    the accuracy of $w$ estimates (using $\Delta W$ metrics estimated as $\left(\frac{w_0 - w_{est}}{w_o}\right)^2$, where $w_0$ and $w_{est}$

218    are modeled and estimated $w$, respectively) with respect to the loss function given in Fig. 3,

219    (2) assessed the proportion of the shared heritability to the total heritability of the SGIT (the

220    $Q$-value) with respect to the loss function, and (3) compared the analytically predicted

221    total/shared heritabilities of two traits: SGIT and the first component, GIP1, obtained by GIP

222    method [2]. The $Q$-value can be interpreted as the specificity metrics of the SGIT: the closer

223    the $Q$-value to 1, the lower the share of unshared heritability in the total heritability of the

224    SGIT. The simulation scenarios were based on six varying parameters that describe the

225    properties of the genetic and phenotypic correlation matrices. Under each scenario, we

226    considered two situations, where all traits have the same $w^2$ and different $w^2$'s. To distinguish

227    between these situations, we will hereinafter write either '$w^2$' or 'different $w^2$'s'. In total, we

228    performed 10,000 iterations of simulations for each of 288 scenarios.

229    The results are presented in Supplementary Figures S1-18. For all scenarios, there are

230    few general patterns: (1) the higher simulated $w$ values, the higher the accuracy of the $w$

231    estimates, (2) the accuracy of the $w$ estimates and the $Q$-value increase with an increasing in

232    the number, $K$, of traits, (3) for all scenarios with $w^2 > 0.8$, $\Delta W$ was very low (<0.025) and the

233    $Q$-value was more than 90%.

234    For all scenarios with three traits, the accuracy of the $w$ estimates was in general low:

235    $\Delta W$ was not higher than 0.7 for scenarios with $w^2 = 0.2$ and 0.3, although at $w^2$ equal to or

236    higher than 0.4 $\Delta W$ was less than 0.2. The $Q$-value was higher than 60% for almost all

SHAHER framework, 2021

237    scenarios with $w^2 \geq 0.4$. In almost all cases, the total and shared heritabilities of the SGIT

238    were higher than the corresponding heritabilities of GIP1, except for the scenarios with

239    $h^2 = 0.8$.

240      For the scenarios with four and five traits, the accuracy of $w$ estimates was higher:

241    $\Delta W < 0.15$ 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

242    was more than 70% for four traits and more than 80% for five traits. Again, the total and

243    shared heritabilities of the SGIT were higher than the corresponding heritabilities of GIP1

244    under all scenarios, except for the scenarios with $h^2 = 0.8$. In the scenarios with $h^2 = 0.8$, the

245    total and shared heritabilities of the SGIT were higher than those of the GIP1 at $w^2 \geq 0.5$.
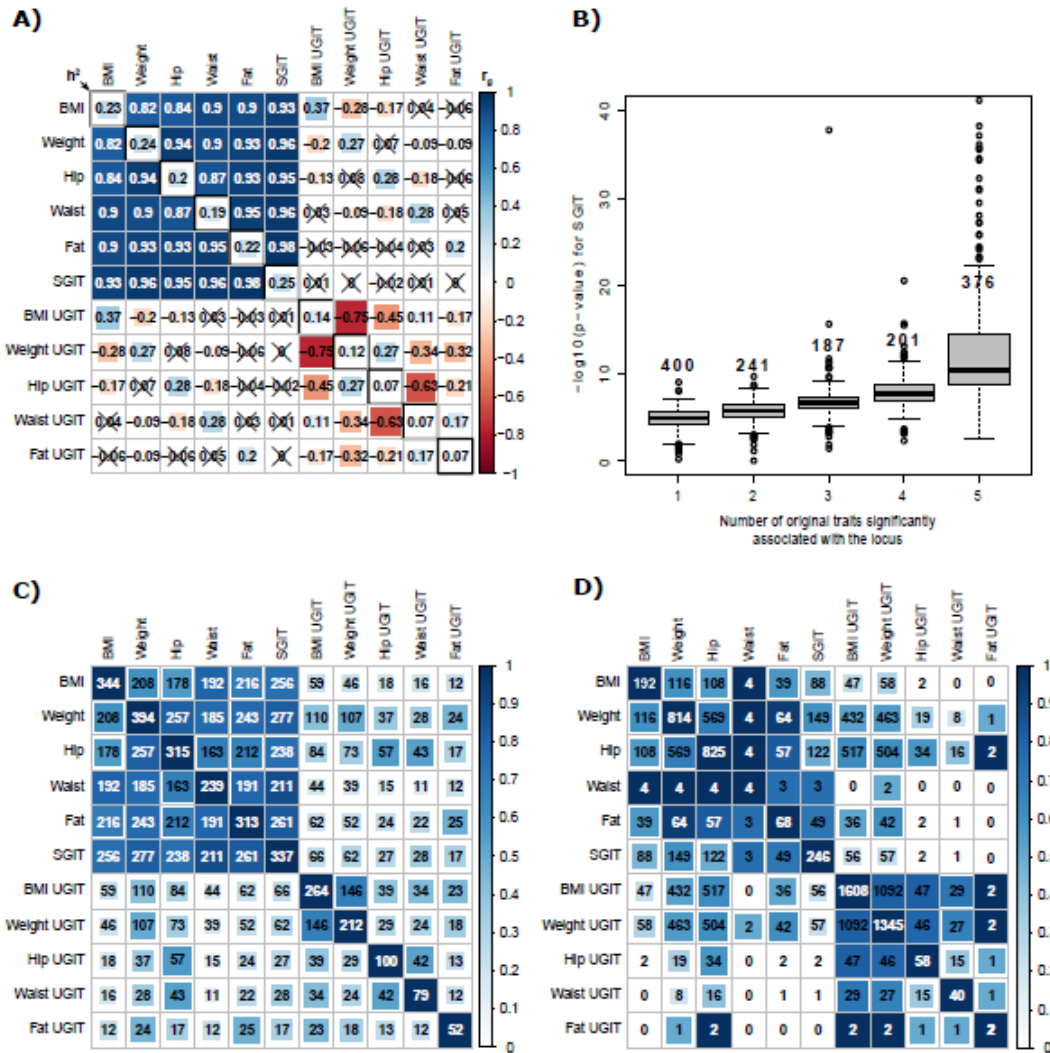
246      In summary, the performance of MaxSH was suitable at $w^2 \geq 0.5$ and when the number

247    of traits being higher than or equal to four. In the case of small $w$ or three traits, the results of

248    MaxSH should be interpreted with caution.

249    **Real data assessment**

250      We applied SHAHER to three datasets: anthropometric (five traits), psychometric

251    (four traits) and lipid traits (three traits). We should note that the performance of SHAHER

252    applied to three traits is limited (see simulation results), yet still passable, although the

253    results should be interpreted with caution. We present SHAHER results for anthropometric

254    traits in the main text as an example. The full results for the psychometric and lipid traits are

255    presented in Supplementary Results.

256      At the first step, we confirmed that SGI exists for five traits. At the second step, we

257    determined the $\alpha$ and $\gamma$ coefficients and their CI (see Supplementary Table 1a). At the third

258    step, we applied sumCOT and obtained GWAS results for the SGIT and UGITs (see

259    Supplementary Table 2a for heritability estimates and LD score regression intercepts).

260    SHAHER results are presented in Figure 2.

261

**Figure 2. 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-value) for the SGIT with respect to the number of the original traits significantly associated with the locus. Two outliers for loci with $-\log_{10}$(p-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 the 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 the scale from 0 to 1). The diagonal elements represent the number of gene sets found for every trait.

11
SHAHER framework, 2021

280    Figure 2A demonstrates genetic correlations between all pairs of the original
281    anthropometric traits, SGIT and UGITs. All the original traits were positively correlated with
282    r > 0.82. We did not observe any significant genetic correlation between the SGIT and the
283    UGITs. Moreover, we did not observe additional SGI among UGITs, which was up to
284    expectation. The heritabilities of the UGITs varied from 0.07 to 0.14.

285    We revealed a dependence of the SGIT p-value from the number of the original traits
286    significantly associated with the locus (Figure 2B). It clearly shows that the loci associated
287    with all the original traits have lower SGIT p-values than the other loci.

288    Joint clumping of 11 traits (five original traits, five UGITs and SGIT) resulted in 820
289    genome-wide significantly associated loci (p-value $< 5{\times}10^{-8}$, Supplementary Table 3a). If a
290    locus was not significantly associated with any of the original traits, it was considered new.
291    SGIT was significantly associated with 337 SNPs. We detected no new loci among SGIT
292    loci. The clumping of UGITs revealed 422 loci, of which 199 were new. At the same time,
293    the clumping of only original traits allowed 621 loci to be detected, of which 161 could not
294    be detected by analyzing SGIT or UGITs. Thus, the joint analysis of SGIT and UGITs
295    increased the number of associated loci by more than 32%. Figure 2C reflects the
296    overlapping between significantly associated loci for 11 analyzed traits. There is a weak
297    albeit non-zero overlap between loci for UGITs and SGIT, although the genetic correlation
298    between them is zero. It could be due to the conservative settings of the clumping procedure,
299    which tends to clump together closely located loci, and due to some level of unspecificity of
300    the SHAHER.

301    Next, we checked how enriched gene sets overlap between the SGIT, UGITs and
302    original traits (see Figure 2D). Significant results (FDR<5%) of enriched gene sets and tissue
303    enrichment analyses are presented in Supplementary Table 4. As expected, the heatmap of
304    the overlapping gene sets looks similar to the heatmap of genetic correlations and the
305    heatmap of the overlapping loci. Moreover, there were almost no overlap between SGIT and
306    any UGIT. For the original traits, the number of enriched gene sets varied a lot: from 4 for
307    the waist to 825 for the hip circumference. It should be noted that we observed a high
308    number of enriched gene sets for the BMI UGIT, 1608, which was almost ten times the value
309    for BMI (192).

310    Finally, we obtained GIP1 GWAS statistics and calculated the genetic correlations
311    between the SGIT and GIP1. The genetic correlation was higher than 0.97.
312

313 **Discussion**

314   We developed a new fast and efficient framework, which allows us to decompose the

315 heritability of each trait from a given set of traits into two components. One of them is

316 explained by shared genetic factors common to all traits. Another one is explained by

317 unshared genetic factors specific for each trait. The framework not only decomposes

318 heritability, but also identifies SNPs associated with the shared and unshared genetic

319 impacts. To our knowledge, this framework is unparalleled. It has an additional advantage: it

320 uses GWAS summary statistics obtained for original traits and does not require raw

321 genotype or phenotype data.

322   We compared the performances of MaxSH and GIP in identifying the shared genetic

323 components. GIP calculates the linear combination coefficients via the eigenvalues of the

324 genetic covariance matrix and can be considered a close approximation to MaxSH. In our

325 simulations, GIP and MaxSH were similar in almost all scenarios, with MaxSH being

326 somewhat superior in terms of the power (total heritability) and quality (shared heritability).

327 If obtaining genetically independent phenotypes is not the aim, we suggest using SHAHER,

328 because it is more robust and gives additional metrics like SGI contributions to the

329 heritability of the original traits.

330   The framework is computationally effective. The stage using sumCOT is the most

331 time consuming. However, it takes only several minutes for an average computer to conduct

332 a GWAS of a linear combination of traits with 6M SNPs using a C++ implementation of the

333 sumCOT. MaxSH, based on numerical optimization procedures, and the other parts of the

334 framework take seconds.

335   The proposed sumCOT method can be applied as an independent tool to address

336 additional tasks. One of them is making a summary-level adjustment of traits by other traits

337 using the same scheme as was used for obtaining the UGIT GWAS statistics. This can be

338 helpful, for example, for ridding the studied trait's genetic component of the genetic

339 component that was caused by the confounding or unaccounted effects of assortative mating

340 or family effects, which is quite a problem in GWAS at the biobank scale [15, 18]. Another task

341 is a GWAS for the trait that appears as a linear combination of the original traits. The

342 sumCOT method is robust to differences in sample sizes used for GWASs of original traits

343 and is applicable to different GWAS models (Cox, linear or logistic).

344   The main interest in the application of the SHAHER framework lies in the possibility

345 of obtaining novel biological insights into a trait's heritability composition. This can be

SHAHER framework, 2021

346 achieved by the application of a huge variety of *in-silico* follow-up techniques to the SGIT
347 and UGITs. The SGIT is of interest by itself, but we also emphasize the importance of the
348 comparison of shared and unshared impacts for each trait. In our real data application, the
349 most remarkable case is BMI in the set of anthropometric traits (see Figure 3C). We found
350 246 and 1608 significantly enriched gene sets for the SGIT and UGIT of BMI, respectively,
351 with negligible overlapping between them of size 56. By analyzing only BMI, we would
352 have detected only 192 enriched gene sets. By analyzing each of the impacts separately, we
353 dramatically increased the number of observed unique gene sets (1798 in total for both SGI
354 and UGI). It means that each sub-phenotype controlled by the SGF and UGF is less
355 heterogeneous than the original trait. According to the significant gene sets, the UGIT of
356 BMI (see Supplementary Tables 4) controls some structural changes in body compositions
357 and bone formation, while the SGIT is involved in some general signaling pathways and
358 pathways related to nervous system development and probably to general psycho-social
359 aspects of BMI, obesity and other anthropometric traits [19].

360     Although SHAHER is effective, it has several limitations. First, when trait-trait
361 genetic correlations are weak, it is expected that the contributions of these traits to the shared
362 heritability will be small, too. In this case, MaxSH may overestimate these contributions.
363 Secondly, the framework is applicable only if the number of traits is no less than three. In the
364 case of three traits, the performance is limited and the SHAHER results should be interpreted
365 with caution. We have shown in simulations and real dataset examples that MaxSH works
366 better at higher numbers of genetically correlated traits being analyzed. However, an
367 increase in the number of weakly correlated traits leads to a decrease in the proportion of
368 SNPs associated with all traits simultaneously and to a decrease in the efficiency of the
369 framework. Thirdly, although the set of SNPs identified by the SGIT GWAS is enriched for
370 the SGF, each SNP should be interpreted with caution for whether it is shared or not,
371 because SHAHER has some level of unspecificity. Finally, if any confounding effects were
372 included in the GWAS of the original traits, these effects are amplified in the SGIT [15]. The
373 confounding effects can be controlled easily using special methods like LD score regression
374 [20], although this method fails to distinguish a polygenic component if the trait was measured
375 in the sample with the assortative mating or family effects. Thus, we suggest a thorough
376 check of the original GWAS for the presence of any effects of possible confounders before
377 proceeding to SHAHER.

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

386

## Materials and Methods

387

### Simulation study

388

389    Under different scenarios, we designed simulations to assess the performance of
390    MaxSH. We (1) assessed the accuracy of $w$ estimates, (2) assessed the proportion of SGIT
391    heritability explained by the SGF to the total heritability of the SGIT (the $Q$-value), and (3)
392    compared the analytically predicted total and shared heritabilities of the SGIT and GIP1 with
393    respect to the loss function. The design of our simulation experiment is shown in Figure 3.
394    To generate the input for the MaxSH and GIP approaches, we used a six-parameter
395    simulation model, in which $K$ is the number of traits; $W_0^2$ is a ($K \times K$) diagonal matrix, where
396    the $i$-th diagonal element is $w_i^2$ (the proportion of heritability explained by SGF); $s$ is the
397    proportion of zeros in the matrix $U_{unsh}$; $d_1$ is the amplitude of the uniform distribution for
398    non-zero values of $U_{unsh}$ and $d_2$ is the amplitude of the uniform distribution for $U_{env}$; $H_2$ is the
399    diagonal matrix with diagonal elements equal to the trait heritabilities. The parameters
400    values used are given in Figure 3.

401



402

**Figure 3. A schematic depicting the overall workflow of a simulation study.** All details are given in the text.

405

406     For each fixed number, $K$, of the original traits and fixed heritability, $h_i^2$ ($i=1,\ldots,K$), of

407    each trait, we simulated $U_{gen}$. To do this, we separately modelled two its components caused

408    by SGF and UGF as $W\mathbf{11}^T W$ and $\sqrt{I-W^2}U_{unsh}\sqrt{I-W^2}$, respectively (see the 'Model'

409    box in Figure M1 of Supplementary Methods). Here $\mathbf{1}$ is a ($K\times1$) vector of units, and $U_{unsh}$ is

410    a ($K\times K$) matrix randomly generated using the parameters $s$ and $d_1$ (see Supplementary

411    Methods). Then we randomly generated the trait-trait correlation matrix $U_{env}$ explained by

412    the environmental factors, by giving the parameter $d_2$ (see Supplementary Methods). Finally,

413    we modeled a matrix of phenotypic correlations by using Model (2) with regard to simulated

414    values $W_0$.

415     Using simulation data, $U_{phen}$, $U_{gen}$ and $H^2$, we estimated $W_{est}$ and calculated its squared

416    relative difference with the simulated values of $W_0$ ($\Delta W$). We revealed a dependence of $\Delta W$

417    on the loss function (*Loss*). The *Loss* value characterizes the difference between $U_{gen}$ and the

418    auxiliary matrix $V$.

419    Then we estimated $\alpha$ in three ways: (1) using MaxSH and $W_0$, (2) using MaxSH and

420    $W_{est}$, and (3) using the GIP method [2]. On the basis of these estimates, we formed three traits

421    being the linear combinations of the original traits. For these combined traits, we calculated

422    the total heritability and the heritability explained by SGF.

423    The simulated experiments were repeated 10,000 times for each set of parameters.

424    The model parameters and formulas for all calculated values are shown in Figure 3.

425

**Application to real data**

**Data sets**

428    We used three publicly available real data sets: anthropometric traits, psychiatric

429    conditions and lipid concentrations, which contain five, four and three traits respectively.

430    The group of anthropometric traits consisted of UK Biobank GWAS summary

431    statistics obtained from the Neale lab (http://www.nealelab.is/blog/2017/7/19/rapid-gwas-of-

432    thousands-of-phenotypes-for-337000-samples-in-the-uk-biobank) for people of European

433    ancestry: BMI (N = 336,107), weight (N = 336,227), hip (N = 336,601), waist circumference

434    (N = 336,639) and whole body fat mass (N = 330,762).

435    The second dataset reflecting psychometric traits was constructed from GWAS

436    results    provided    by    the    Psychiatric    Genomics    Consortium

437    (https://www.med.unc.edu/pgc/download-results/) for bipolar disorder, BIP (N cases =

438    20,352; N controls = 31,358) [21], major depressive disorder, MDD (N cases = 43,204; N

439    controls = 95,680; without UK Biobank and 23andMe data) [22] and schizophrenia, SCZ (N

440    cases = 36,989; N controls = 113,075). Summary statistics for the fourth trait – subjective

441    well-being (N = 110,935) – were derived from UK Biobank data from the Neale lab. All the

442    psychometric trait GWASs were conducted using samples of white Europeans.

443    The last dataset corresponding to lipid traits was formed using GWAS data for

444    European    participants    from    the    Global    Lipid    Genetics    Consortium

445    (http://csg.sph.umich.edu/willer/public/lipids2013/) for LDL cholesterol (N = 173,082),

446    triglycerides (N = 177,860) and total cholesterol (N = 187,365).

447    Summary statistics for the three data sets were integrated and quality controlled by

448    the GWAS-MAP platform developed by our group [23]. The GWAS-MAP database contains

449    implemented software for quality control of GWAS results, estimation of phenotypic

450    correlations and LD Score regression (LDSC) [20].

451     We conducted the quality control of all data and unified them within the GWAS-

452     MAP platform [23]. We filtered all summary statistics by minor allele frequencies ≥ 0.01.

453     Additionally, we filtered GWAS results for BIP by imputation qualities ≥ 0.9. We did not

454     apply this filter to the other traits due to the absence of imputation quality in summary

455     statistics data. Finally, using GWAS-MAP, we performed a correction for genomic control

456     for all traits (including the original traits, SGIT and UGITs) with an LDSC intercept greater

457     than 1 [20]. Thus, we corrected all traits from the psychometric dataset apart from MDD, all

458     original anthropometric traits and their SGIT and lipid SGIT as their LDSC intercept

459     exceeded 1 (see Supplementary Tables 2a-c). Moreover, all SNPs with the p-value equal to 0

460     were excluded from analysis.

**Genetic analysis**

462     Pairwise phenotypic correlations between traits were computed using the GWAS-

463     MAP platform described above. The used method is based on correlations between

464     insignificant $z$-statistics for independent SNPs as previously described in [9]. SNP-based

465     heritability and genetic correlation coefficients were estimated using the LD Score

466     regression software [16] embedded in the GWAS-MAP platform. The significance threshold

467     for genetic correlations was set at $4.5 \times 10^{-4}$ (0.05/112, where 112 is the number of pairwise

468     combinations between all original traits, their SGIT and UGITs in each dataset - between 11,

469     9 and 7 traits for anthropometry, psychometric and lipid traits respectively).

470     SHAHER analysis included checking if there was an SGI or not, the application of

471     MaxSH and conducting SGIT and UGIT GWASs. The threshold for confirming the

472     existence of an SGI at the first stage was empirically set to 0.2.

473     For each dataset, we visualized the full genetic correlation matrices using the

474     *corrplot*() function from the corrplot R package (v.0.84) [24]. We also placed the SNP-based

475     heritability estimates on the diagonal and crossed out non-significant values.

476     Finally, we compared the GWAS results obtained for the SGIT by MaxSH and GIP

477     (the principal component analysis on the matrix of genetic covariances)[2].

**Gene set and tissue/cell type enrichment analyses**

479     We performed a gene set enrichment analysis and a tissue/cell type enrichment

480     analysis combined with a gene prioritization using the Data-driven Expression Prioritized

481     Integration for Complex Traits (DEPICT) tool v.1.1, release 194 [25]. We selected genome-

482     wide significant SNPs (p-value < $5 \times 10^{-8}$) from summary statistics before the genomic

SHAHER framework, 2021

483    control and applied the DEPICT software with default parameters

484    (https://data.broadinstitute.org/mpg/depict/). The MHC region was excluded from analyses.

485    Next, for the gene set enrichment results, we calculated the number of significant

486    enriched gene sets (FDR < 5%) and constructed an overlapping matrix, in which each cell

487    represents the number of overlapping gene sets for each pair of traits. For each pair of traits,

488    we scaled the number of overlapping gene sets by the minimum number of significant gene

489    sets for this pair of traits. The resulting matrix was visualized using the corrplot R-package

490    as descried above.

491    **The number of original traits associated with SGIT loci**

492    We performed a clumping procedure to search for loci associated with each of the

493    original traits, SGIT and UGITs at a genome-wide significance level of $5\times10^{-8}$. The

494    associated locus was defined as a genomic region spanning 500 kb in either direction of the

495    lead SNP. Those loci that were significantly associated with SGIT, but not with the original

496    traits, were assumed to be new loci.

497    We expected that the loci associated with all the original traits used to obtain SGIT

498    are likely to be SGF. To test this expectation, for each dataset we selected all independent

499    loci that were significantly associated with at least one of the original traits and calculated

500    the number of the original traits significantly associated with these loci. For the original

501    anthropometric and lipid traits, we empirically set the significance threshold at p-value =

502    $1\times10^{-5}$. For the psychometric traits, it was set at $1\times10^{-3}$. We then analyzed the SGIT p-values

503    for the selected loci and constructed boxplots of $-\log_{10}$ for them with regard to the number of

504    the original traits significantly associated with these loci.

505    ## Data Availability

506    The SHAHER framework is implemented as a set of R/C++ scripts and is freely

507    available at https://github.com/Sodbo/shared_heredity.

508    ## Acknowledgements

509    ## Funding

## Author contributions

518 YAT, GRS and TIA conceived and oversaw the study. YAT, GRS, SZS and TIA

519 contributed to the design of the study and interpretation of the results. GRS developed the

520 MaxSH method, including the algorithm and program, and conducted simulation studies. PT

521 developed the C++ version of sumCOMB and tested the developed framework. EST, EEE,

522 SGF, SZS and YAT wrote the source code for the framework and performed real data

523 analyses. All co-authors discussed the results and contributed to preparing the draft and final

524 version of the manuscript.

## Conflict of interest

526 P.R.H.J.T. is an employee of BioAge Labs. The remaining authors declare no conflict

527 of interest.

## Supplementary Information legend

529 Supplementary Tables 1a-c: linear combination coefficients and CIs of SGIT for real

530 data sets

531 Supplementary Tables 2a-c: results of LD score regression for original traits from

532 real data sets

533 Supplementary Tables 3a-c: clumping results for real data sets

534 Supplementary Tables 4-6: DEPICT results for real data sets

## References

536 1. Jiang X, Finucane HK, Schumacher FR, Schmit SL, Tyrer JP, Han Y, et al. Shared
537 heritability and functional enrichment across six solid cancers. Nature communications.
538 2019;10(1):1-23.

539  2.      Tsepilov YA, Freidin MB, Shadrina AS, Sharapov SZ, Elgaeva EE, van Zundert J, et
540  al. Analysis of genetically independent phenotypes identifies shared genetic factors
541  associated with chronic musculoskeletal pain conditions. Communications biology.
542  2020;3(1):1-13.
543  3.      Sampson JN, Wheeler WA, Yeager M, Panagiotou O, Wang Z, Berndt SI, et al.
544  Analysis of heritability and shared heritability based on genome-wide association studies for
545  13 cancer types. JNCI: Journal of the National Cancer Institute. 2015;107(12).
546  4.      Brainstorm C, Anttila V, Bulik-Sullivan B, Finucane HK, Walters RK, Bras J, et al.
547  Analysis of shared heritability in common disorders of the brain. Science. 2018 Jun
548  22;360(6395). PubMed PMID: 29930110. Pubmed Central PMCID: PMC6097237.
549  5.      Yang Y, Zhao H, Heath AC, Madden PA, Martin NG, Nyholt DR. Shared genetic
550  factors underlie migraine and depression. Twin Research and Human Genetics.
551  2016;19(4):341-50.
552  6.      Wright S. Correlation and Causation. JouMal of. Agricultural Research. 1921.
553  7.      Rijsdijk FV, Sham PC. Analytic approaches to twin data using structural equation
554  models. Briefings in bioinformatics. 2002;3(2):119-33.
555  8.      Galesloot TE, Van Steen K, Kiemeney LA, Janss LL, Vermeulen SH. A comparison
556  of multivariate genome-wide association methods. PloS one. 2014;9(4):e95923.
557  9.      Stephens M. A unified framework for association analysis with multiple related
558  phenotypes. PloS one. 2013;8(7):e65245.
559  10.     Yang X, Zhang S, Sha Q. Joint analysis of multiple phenotypes in association studies
560  based on cross-validation prediction error. Scientific reports. 2019;9(1):1-10.
561  11.     Turley P, Walters RK, Maghzian O, Okbay A, Lee JJ, Fontana MA, et al. Multi-trait
562  analysis of genome-wide association summary statistics using MTAG. Nature genetics.
563  2018;50(2):229-37.
564  12.     Fatumo S, Carstensen T, Nashiru O, Gurdasani D, Sandhu M, Kaleebu P.
565  Complimentary methods for multivariate genome-wide association study identify new
566  susceptibility genes for blood cell traits. Frontiers in genetics. 2019;10:334.
567  13.     Ning Z, Tsepilov YA, Sharapov SZ, Grishenko AK, Feng X, Shirali M, et al. Beyond
568  power: multivariate discovery, replication, and interpretation of pleiotropic loci using
569  summary association statistics. bioRxiv. 2019:022269.
570  14.     Pasaniuc B, Price AL. Dissecting the genetics of complex traits using summary
571  association statistics. Nature Reviews Genetics. 2017;18(2):117-27.
572  15.     Timmers PRHJ, Tiys ES, Sakaue S, Akiyama M, Kiiskinen TTJ, Zhou W, et al.
573  Genetically independent phenotype analysis identifies LPA and VCAM1 as drug targets for
574  human ageing. bioRxiv. 2021:2021.01.22.427837.
575  16.     Bulik-Sullivan B, Finucane HK, Anttila V, Gusev A, Day FR, Loh PR, et al. An atlas
576  of genetic correlations across human diseases and traits. Nat Genet. 2015 Nov;47(11):1236-
577  41. PubMed PMID: 26414676. Pubmed Central PMCID: PMC4797329.
578  17.     Winkler TW, Day FR, Croteau-Chonka DC, Wood AR, Locke AE, Mägi R, et al.
579  Quality control and conduct of genome-wide association meta-analyses. Nature protocols.
580  2014;9(5):1192-212.
581  18.     Howe LJ, Lawson DJ, Davies NM, Pourcain BS, Lewis SJ, Smith GD, et al. Genetic
582  evidence for assortative mating on alcohol consumption in the UK Biobank. Nature
583  communications. 2019;10(1):1-10.
584  19.     Marcellini F, Giuli C, Papa R, Tirabassi G, Faloia E, Boscaro M, et al. Obesity and
585  body mass index (BMI) in relation to life-style and psycho-social aspects. Archives of
586  gerontology and geriatrics. 2009;49:195-206.

587  20.    Bulik-Sullivan B, Finucane HK, Anttila V, Gusev A, Day FR, Loh P-R, et al. An
588  atlas of genetic correlations across human diseases and traits. Nature genetics.
589  2015;47(11):1236.
590  21.    Stahl EA, Breen G, Forstner AJ, McQuillin A, Ripke S, Trubetskoy V, et al.
591  Genome-wide association study identifies 30 loci associated with bipolar disorder. Nature
592  genetics. 2019;51(5):793-803.
593  22.    Wray NR, Ripke S, Mattheisen M, Trzaskowski M, Byrne EM, Abdellaoui A, et al.
594  Genome-wide association analyses identify 44 risk variants and refine the genetic
595  architecture of major depression. Nature genetics. 2018;50(5):668-81.
596  23.    Gorev D, Shashkova T, Pakhomov E, Torgasheva A, Klaric L, Severinov A, et al.,
597  editors. GWAS-MAP: a platform for storage and analysis of the results of thousands of
598  genome-wide association scans. Bioinformatics of Genome Regulation and
599  Structure\Systems Biology (BGRS\SB-2018); 2018.
600  24.    Wei T, Simko V, Levy M, Xie Y, Jin Y, Zemla J. corrplot: visualization of a
601  correlation matrix. R package v. 0.84. 2017.
602  25.    Pers TH, Karjalainen JM, Chan Y, Westra H-J, Wood AR, Yang J, et al. Biological
603  interpretation of genome-wide association studies using predicted gene functions. Nature
604  communications. 2015;6(1):1-9.

605

SHAHER framework, 2021