

Manuscript Number:	GIGA-D-17-00337	
Full Title:	A network-based conditional genetic association analysis of the human metabolome	
Article Type:	Technical Note	
Funding Information:	the European Union FP7 framework project Pain-Omics (602736)	Dr. Christian Gieger
	Ministry of Education and Science of the Russian Federation (the 5-100 Excellence Programme)	Mr. Sodbo Zh. Sharapov
	the Federal Agency of Scientific Organisations via the Institute of Cytology and Genetics (0324-2018-0017)	Dr. Yakov A. Tsepilov
	the Federal Agency of Scientific Organisations via the Institute of Cytology and Genetics (0324-2018-0017)	Prof. Yurii S. Aulchenko
Abstract:	<p>Background: Genome-wide association studies have identified hundreds of loci that influence a wide variety of complex human traits; however, little is known regarding the biological mechanism of action of these loci. The recent accumulation of functional genomics ("omics"), including metabolomics data, has created new opportunities for studying the functional role of specific changes in the genome. Functional genomic data are characterized by their high dimensionality, the presence of (strong) statistical dependency between traits, and—potentially—complex genetic control. Therefore, the analysis of such data requires specific statistical genetics methods.</p> <p>Results: To facilitate our understanding of the genetic control of omics phenotypes, we propose a trait-centered, network-based conditional genetic association (cGAS) approach for identifying the direct effects of genetic variants on omics-based traits. For each trait of interest, we selected from a biological network a set of other traits to be used as covariates in the cGAS. The network can be reconstructed either from biological pathway databases or directly from the data, using a Gaussian Graphical Model applied to the metabolome. We derived mathematical expressions which allow comparison of the power of univariate analyses with conditional genetic association analyses. We then tested our approach using data from a population-based KORA study (n=1784 subjects, 1.7 million SNPs) with measured data for 151 metabolites.</p> <p>Conclusions: We found that compared to single-trait analysis, performing a genetic association analysis that includes biologically relevant covariates increases the power of the analysis by providing more accurate estimates of genetic effects; however, analysis can either gain additional power or even lose power, depending on specific pleiotropic scenarios, for which we provide empirical examples. We also show the importance of properly selecting sets of traits to be entered in the multivariate analysis. In the context of analyzed metabolomics data, the knowledge-based network approach has increased power. Nevertheless, we believe that our analysis shows that neither a prior-knowledge-only approach nor a phenotypic-data-only approach is optimal, and we discuss possibilities for improvement.</p>	
Corresponding Author:	Yurii Aulchenko Institute of Cytology and Genetics SB RAS Novosibirsk, RUSSIAN FEDERATION	
Corresponding Author Secondary Information:		
Corresponding Author's Institution:	Institute of Cytology and Genetics SB RAS	
Corresponding Author's Secondary Institution:		

First Author:	Yakov A. Tsepilov, Ph.D.
First Author Secondary Information:	
Order of Authors:	Yakov A. Tsepilov, Ph.D.
	Sodbo Zh. Sharapov
	Olga O. Zaytseva, Ph.D.
	Jan Krumsek, Ph.D.
	Cornelia Prehn, Ph.D.
	Jerzy Adamski, Ph.D.
	Gabi Kastenmüller, Ph.D.
	Rui Wang-Sattler, Ph.D.
	Konstantin Strauch, Ph.D.
	Christian Gieger, Ph.D.
	Yurii S. Aulchenko, Ph.D.
Order of Authors Secondary Information:	
Opposed Reviewers:	
Additional Information:	
Question	Response
Are you submitting this manuscript to a special series or article collection?	No
Experimental design and statistics	Yes
<p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist. Information essential to interpreting the data presented should be made available in the figure legends.</p> <p>Have you included all the information requested in your manuscript?</p>	
Resources	Yes
<p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist?</p>	

<p>Availability of data and materials</p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in publicly available repositories (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.</p> <p>Have you have met the above requirement as detailed in our Minimum Standards Reporting Checklist?</p>	<p>Yes</p>

A network-based conditional genetic association analysis of the human metabolome

Y.A. Tsepilov^{1,2}, S.Z. Sharapov², O.O. Zaytseva^{1,2}, J. Krumsek³, C. Prehn⁴, J. Adamski^{4,5,6}, G.
Kastenmüller⁷, R. Wang-Sattler^{6,8,9}, K. Strauch^{10,11}, C. Gieger^{6,8,9}, Y.S. Aulchenko^{1,2,12*}

1 Institute of Cytology and Genetics SB RAS, Novosibirsk, Russia

2 Novosibirsk State University, Novosibirsk, Russia

3 Institute of Computational Biology, Helmholtz Center Munich - German Research Center
for Environmental Health, Neuherberg, Germany

4 Institute of Experimental Genetics, Genome Analysis Center, Helmholtz Center Munich -
German Research Center for Environmental Health, Neuherberg, Germany

5 Institute of Experimental Genetics, Life and Food Science Center Weihenstephan, Technical
University of Munich, Freising-Weihenstephan, Germany

6 German Center for Diabetes Research, Neuherberg, Germany

7 Institute of Bioinformatics and Systems Biology, Helmholtz Center Munich - German
Research Center for Environmental Health, Neuherberg, Germany

8 Research Unit of Molecular Epidemiology, Helmholtz Center Munich - German Research
Center for Environmental Health, Neuherberg, Germany

9 Institute of Epidemiology II, Helmholtz Center Munich - German Research Center for
Environmental Health, Neuherberg, Germany

10 Institute of Genetic Epidemiology, Helmholtz Center Munich - German Research Center
for Environmental Health, Neuherberg, Germany

11 Institute of Medical Informatics, Biometry and Epidemiology, Chair of Genetic
Epidemiology, Ludwig Maximilian University of Munich, Munich, Germany

12 PolyOmica, 's-Hertogenbosch, The Netherlands

* Correspondence to

Yurii S. Aulchenko

Institute of Cytology and Genetics SB RAS, 630090 Novosibirsk, Russia

yurii@bionet.nsc.ru

*Keywords: genome-wide association study; multivariate model; metabolomics; conditional
analysis; pleiotropy*

38 **Abstract**

1 39 **Background:** Genome-wide association studies have identified hundreds of loci that influence a
2 wide variety of complex human traits; however, little is known regarding the biological mechanism
3 of action of these loci. The recent accumulation of functional genomics (“omics”), including
4 metabolomics data, has created new opportunities for studying the functional role of specific
5 changes in the genome. Functional genomic data are characterized by their high dimensionality,
6 the presence of (strong) statistical dependency between traits, and—potentially—complex genetic
7 control. Therefore, the analysis of such data requires specific statistical genetics methods.
8

9 46 **Results:** To facilitate our understanding of the genetic control of omics phenotypes, we propose a
10 trait-centered, network-based conditional genetic association (cGAS) approach for identifying the
11 direct effects of genetic variants on omics-based traits. For each trait of interest, we selected from
12 a biological network a set of other traits to be used as covariates in the cGAS. The network can be
13 reconstructed either from biological pathway databases or directly from the data, using a Gaussian
14 Graphical Model applied to the metabolome. We derived mathematical expressions which allow
15 comparison of the power of univariate analyses with conditional genetic association analyses. We
16 then tested our approach using data from a population-based KORA study (n=1784 subjects,
17 1.7 million SNPs) with measured data for 151 metabolites.
18

19 55 **Conclusions:** We found that compared to single-trait analysis, performing a genetic association
20 analysis that includes biologically relevant covariates increases the power of the analysis by
21 providing more accurate estimates of genetic effects; however, analysis can either gain additional
22 power or even lose power, depending on specific pleiotropic scenarios, for which we provide
23 empirical examples. We also show the importance of properly selecting sets of traits to be entered
24 in the multivariate analysis. In the context of analyzed metabolomics data, the knowledge-based
25 network approach has increased power. Nevertheless, we believe that our analysis shows that
26 neither a prior-knowledge-only approach nor a phenotypic-data-only approach is optimal, and we
27 discuss possibilities for improvement.
28

29 64

65 **Short abstract**

1 66 We propose a trait-centric network-based conditional approach for performing a genetic
2
3 67 association analysis of multivariate omics phenotypes. This approach can incorporate existing
4
5 68 biological knowledge regarding biological pathways obtained from external sources and is
6
7 69 designed to specifically test for direct genetic effects. We applied this approach to existing
8
9 70 metabolomics data and found that it increases power by having increased accuracy of genetic effect
10
11 71 estimates in the presence of specific “counterintuitive” pleiotropic scenarios in which genetically
12
13 72 induced and residual covariance are opposite. We provide examples of different pleiotropic
14
15 73 scenarios, and we discuss possible additional applications for this approach.
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

74 **Background**

1 75 Genome-wide association studies (GWASs) are a highly popular method for identifying alleles
2
3 76 that affect complex traits in humans, including the risk of common diseases. In the past decade,
4
5 77 GWASs have enabled the identification of thousands of loci, significantly increasing our
6
7 78 understanding of the genetic basis underlying the control of complex human traits [1]. On the other
8
9 79 hand, this has had only a limited impact on the development of biomarkers and therapeutic agents;
10
11 80 in most cases, any association found using GWAS approach can only serve as a starting point for
12
13 81 future research, rather than providing a direct answer to the question of the genetic region's precise
14
15 82 biological function. The recent accumulation of functional genomics (or “omics” for short) data—
16
17 83 including information regarding the levels of gene expression (the transcriptome), metabolites (the
18
19 84 metabolome), proteins (the proteome), and glycosylation (the glycome)—can provide new insight
20
21 85 into the functional role of specific changes in the genome [2,3].

22 86 Metabolomics is an emerging field that has been studied extensively in the past decade. A
23
24 87 number of GWASs of metabolites have been performed using various platforms [4–8], revealing
25
26 88 literally dozens of loci associated with variations in various lipid species, amino acids, and other
27
28 89 small molecules. Linking the variants that underlie these variations in metabolomics with various
29
30 90 diseases can provide functional insight into the many disease-related associations that were
31
32 91 reported in previous studies, including cardiovascular and kidney disease, type 2 diabetes, cancer,
33
34 92 gout, venous thromboembolism, and Crohn's disease [5].

35 93 However, analyzing metabolomics data requires specialized statistical methods due to their
36
37 94 characteristically high dimensionality and the presence of statistical dependencies that reflect
38
39 95 biological relationships between different variables. Conventional univariate GWAS (uGAS)
40
41 96 approaches ignore any possible dependencies between different omics traits, which can confound
42
43 97 the biological interpretation of the results and may lead to a loss of statistical power. On the other
44
45 98 hand, utilizing multivariate phenotype information increases the statistical power of the association
46
47 99 tests compared to univariate analysis [9–12]. Despite a large number of methodological studies,
48
49 100 however, only a few empirical multivariate GWASs have been published using data for humans.
50
51 101 We recently demonstrated [13] that using a multivariate analysis can substantially increase the
52
53 102 power of locus identification in the context of human *N*-glycomics; indeed, not only did our
54
55 103 multivariate analysis double the number of loci identified in the analysis sample, but also all five
56
57 104 novel loci were strongly replicated. With respect to metabolomics, Inouye et al. [6] performed a
58
59 105 multivariate GWAS on 130 metabolites (grouped in 11 sets) measured in approximately 6600
60
61 106 individuals. They found that multivariate analysis doubled the number of loci detected in this
62
63 107 sample; seven of these additional loci discovered were novel loci that had not been identified
64
65

108 previously in other GWAS analyses of related traits. While no replication of novel loci was
109 performed by Inouye et al., we compared the authors' results with a recently published univariate
110 GWAS of metabolomics derived from a cohort containing nearly 25,000 individuals [8]. We found
111 that three of the seven SNPs reported by Inouye et al. have a p -value $< 5 \times 10^{-11}$ for at least one
112 metabolite (i.e., are significant at the genome-wide level after Bonferroni correction for 130
113 analyses). These findings provide empirical evidence supporting the value of using multivariate
114 methods to analyze the genomics of metabolic traits, at least in the context of locus discovery.

115 It should be noted that these multivariate methods and tests were developed by statistical
116 geneticists to specifically increase the power of gene identification. In such "gene-centric" tests,
117 the model that includes the effects of genotype on multiple traits is contrasted with the null model
118 in which the gene has no effect on any trait analyzed. Although useful and powerful for genetic
119 mapping, this approach may have limited interpretability in a context in which one is interested in
120 the genetic control and biology of specific trait or a subset of traits (the "trait-centered" view).
121 Several statistical methods have been suggested to address the question of which specific traits are
122 affected in an analyzed ensemble (see for example [10,14]). One such method is based on
123 conditional analysis [15], in which a "target trait" is analyzed as a genotype-dependent variable
124 and related traits are included in the regression model as covariates. Such a modeling approach
125 allows—at least in theory—one to rule out indirect genetic effects (e.g., effects that are in fact
126 solely mediated through some other trait) and study only the genetic effects that directly affect the
127 trait of interest.

128 Here, we present a statistical model in which a given trait depends on a genetic
129 polymorphism and in which a number of related traits are included in the model as covariates. In
130 this model, the relationship between the genotype and the trait of interest is our primary focus.
131 Analyzing such a model allows us to identify the direct effect of genetics on the trait of interest.
132 We first contrast our conditional genetic association (cGAS) approach with the standard model in
133 which a trait of interest depends solely on genotype, without other traits used as covariates (i.e.,
134 the univariate genetic association—or uGAS—model). We do so by mathematically deriving
135 expressions that allow us to examine the relative power of the uGAS and cGAS approaches, and
136 we identify the situations in which these models are expected to yield different results.

137 As might be expected—and as demonstrated here—the choice of covariates plays a critical
138 role in conditional analyses. First, we used the assumption that the covariates (i.e., biologically
139 relevant traits) are known. Second, we explored the problem of selecting appropriate covariates,
140 and we tested the approaches by performing a proof-of-principle study using metabolomics data
141 consisting of 151 metabolites (Biocrates assay) obtained from the KORA F4 study (n=1785
142 individuals). Specifically, we selected covariates based on existing knowledge from metabolite

143 biochemical networks (BN-cGAS) and using a data-driven approach based on Gaussian Graphical
144 Modeling (GGM-cGAS). Finally, we compare and discuss the obtained results, and we discuss
145 possible applications for this analysis based on biologically and/or statistically relevant traits.

146
147

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

148 Results

149 The power of performing a conditional analysis of genetic associations

150 We start with the theoretical substantiation and identification of specific scenarios in which
151 adjusting for biologically relevant covariates can modify the power of an association analysis.

152 Let us consider a trait of interest, y , covariate c , and genotype g . We can formulate this
153 problem in terms of a linear regression as follows: $y = \beta_g * g + \beta_c * c + e$, where β_g and β_c are
154 the effects of the genotype and covariate, respectively, and e is the residual noise. Without a loss
155 of generality, we assume that all random variables in this equation are distributed with a mean of
156 zero and a standard deviation of 1; if this is not the case, one would need to use partial regression
157 coefficients and covariance instead of the correlation parameters that we use here. Given these
158 assumptions made, the joint distribution of y , g , and c can be specified using a set of three
159 correlation coefficients, ρ . Given specific parameters, the value of the “univariate” score test
160 statistic for the association between y and g is calculated as follows: $T_u^2 = n \rho_{yg}^2 / \sigma_u^2$, where n is
161 the sample size and $\sigma_u^2 = 1 - \rho_{yg}^2$ is the residual variance of y . For the conditional test, $T_c^2 =$
162 $n \beta_{yg}^2 / \sigma_c^2 = n(\rho_{yg} - \beta_{yc} \rho_{cg})^2 / \sigma_c^2$, where β represents the partial correlation coefficients from
163 the conditional model and σ_c^2 is the residual variance of y . Consequently, the log-ratio of the
164 conditional and univariate test statistics can be partitioned into two components:

$$165 \quad \log\left(\frac{T_c^2}{T_u^2}\right) = \log\left(\frac{\sigma_u^2}{\sigma_c^2}\right) + \log\left(\left[1 - \frac{\beta_{yc} \rho_{cg}}{\rho_{yg}}\right]^2\right) \quad (1)$$

166 Because the first term in Eq. (1) is dependent only upon residual variances of the two
167 models, we call this term the “noise” component. The second term depends upon the correlations
168 between traits and between the traits and the genotype; we call this term the “pleiotropic”
169 component. Because the noise component (σ_u^2 / σ_c^2) is always ≥ 1 , any possible decrease in the ratio
170 between univariate and conditional tests is determined by the sign and the magnitude of the term
171 $\beta_{yc} \rho_{cg} / \rho_{yg}$. If this term is negative, there will always be an increase in the power of the
172 conditional analysis.

173 We can re-write $\beta_{yc} \rho_{cg} / \rho_{yg}$ as $\beta_{yc} \rho_{yc}^*$, where $\rho_{yc}^* = \rho_{cg} / \rho_{yg}$ is the quantity interpreted
174 in a Mendelian randomization analysis as the effect of the covariate on the trait independent of
175 non-genetic confounders [16]. Note that whereas ρ_{yc}^* reflects the covariance between the trait and
176 the covariate (which is induced by the effect of the genotype), β_{yc} is related to the residual (in
177 most cases, environmental) sources of covariance between y and c . Thus, we conclude somewhat
178 surprisingly that when genotype-induced and environmental correlations are similar in sign (i.e.,
179 both are positive or both are negative), the product $\beta_{yc} \rho_{yc}^*$ is positive and the contribution of the

180 second term in Eq. (1) to the relative power is negative. Note that the contribution of the first term
181 in Eq. (1) is always positive; therefore, even if $\beta_{yc}\rho_{yc}^*$ is positive, the power of a conditional
182 analysis may still be higher than the power of a univariate analysis. In contrast, an “unexpected”
183 product (in which the signs are different and hence $\beta_{yc}\rho_{yc}^*$ is negative) contributes positively to
184 the relative power of the conditional model. Note that in such a situation, the power of a conditional
185 analysis will always be higher than the power of a univariate analysis.

186 We can readily extend Eq. (1) to a situation in which k covariates are included in the
187 conditional model. Denoting the coefficients of correlation between g and covariate i as ρ_{gi} and
188 the partial coefficients of regression of y on covariate i as β_i yields the following equation:

$$\log\left(\frac{T_c^2}{T_u^2}\right) = \log\left(\frac{\sigma_u^2}{\sigma_c^2}\right) + \log\left(\left[1 - \frac{1}{\rho_{yg}} \sum_{i=1}^k \beta_i \rho_{gi}\right]^2\right) \quad (2)$$

190 When appropriate covariates are selected, performing cGAS using individual-level data
191 becomes rather trivial and can be achieved using standard statistical and software tools in which
192 one estimates the effects of a SNP and covariates. However, cGAS becomes somewhat less trivial
193 if one chooses to use summary-level univariate GWAS data such as data available from previously
194 published studies. The formalization of cGAS in terms of summary univariate GWAS statistics is
195 described in **Supplementary Note 1**. Here, we used methods based on analyzing summary-level
196 data.

197 **Network-based selection of covariates**

199 The ability to select appropriate covariates is extremely important, as it can have direct
200 implications regarding the outcome of the analysis. If the biological/biochemical relationships
201 between traits of interest are known and are summarized in a database(s), this knowledge can be
202 used directly, for example by using all direct neighbors as covariates. We refer to this approach as
203 a biochemical-network driven cGAS (BN-cGAS). Alternatively, the network can be reconstructed
204 in a hypothesis-free, empirical manner from the data, for example using a Gaussian Graphical
205 Model (GGM) [17]. We refer to this approach as a GGM-cGAS.

206 We compared cGAS and uGAS by performing a genome-wide analysis of genetic effects
207 using summary-level data obtained from the KORA F4 study. This study included 151 metabolites
208 measured in 1784 individuals using the Biocrates assay and imputed at 1,717,498 SNPs.

209 First, we examined the potential of using cGAS when the covariates are selected based on
210 a known biochemical network (i.e., BN-cGAS). Thus, our analysis was restricted to a subset of
211 105 metabolites for which at least the one-reaction-step immediate biochemical neighbors are
212 known [17]. This biochemical network incorporates only lipid metabolites, and the pathway

213 reactions cover two groups of pathways: (1) fatty acid biosynthesis reactions, which apply to the
214 metabolite classes lyso-PC, diacyl-PC, acyl-alkyl-PC, and sphingomyelins; and (2) β -oxidation
215 reactions that reflect fatty acid degradation and apply to acylcarnitines. The β -oxidation model
216 consists of a linear chain of C2 degradation steps (C10 to C8 to C6, etc.). The number of covariates
217 ranged from 1 to 4, with mean and median values of 2.48 covariates and 2 covariates, respectively.

218 **Table 1** lists the 11 loci that were significant in either BN-cGAS or uGAS and fell into
219 known associated regions (see **Supplementary Note 2**). Of these 11 loci, ten and nine loci could
220 be identified by BN-cGAS and uGAS, respectively. Compared to uGAS, BN-cGAS identified one
221 fewer locus (*ETFDH*), but identified two more (*ACSL1* for PC ae C42:5 and *PKD2L1* for lyso-
222 PC a C16:1). It is interesting to note that for *ACSL1*, the effect of SNP rs4862429 on PC ae C42:5
223 was highly significant ($p=7e-11$) with BN-cGAS, but was not significant ($p=0.7$) with uGAS; this
224 outcome is to be expected under the model of unexpected pleiotropy.

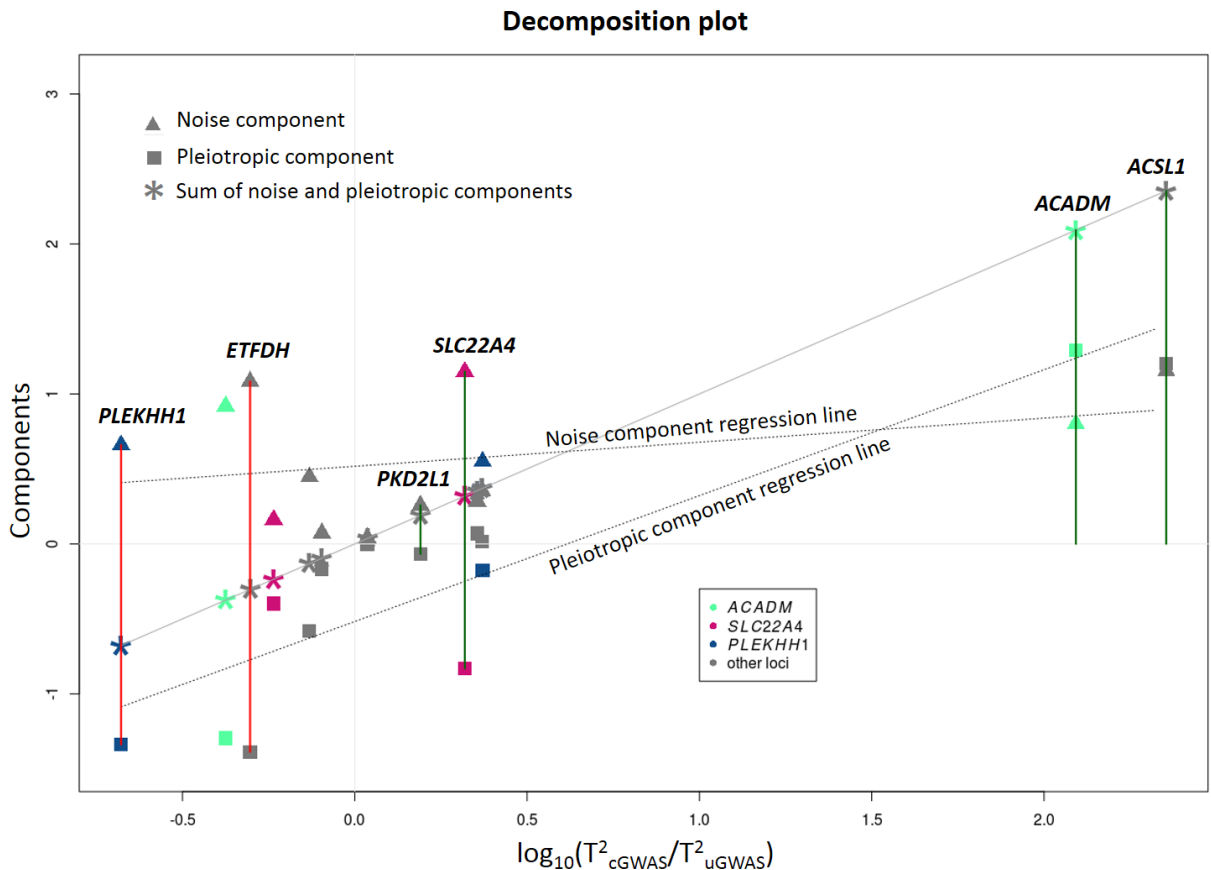
225 Next, to test whether using BN-cGAS increases the average power of the association
226 analysis, we compared the BN-cGAS and uGAS chi-square test results for the loci listed in **Table**
227 **1**. Within a given locus, we compared the maximum test value. The ratio of the average maximum
228 test statistic between BN-cGAS and uGAS was 1.59; however, a paired-sample Wilcoxon test
229 comparing the best chi-square test results between BN-cGAS and uGAS was not significant
230 ($p=0.067$).

231 For the SNPs listed in **Table 1**, we then used Eq. (2) to partition the log-ratio of the BN-
232 cGAS and uGAS statistics values into “noise” and “pleiotropic” components. As shown in **Figure**
233 **1**, the ratio is determined primarily by the second (i.e., “pleiotropic”) term in Eq. (2). Moreover,
234 with the exception of the *SLC22A4* locus, the SNP-trait pairs for which BN-cGAS had increased
235 power are the pairs in which the second term in Eq. (2) is either positive or close to zero. In
236 contrast, in the SNP-trait pairs that were not identified using BN-cGAS, the “pleiotropic” term in
237 Eq. (2) had a strong negative contribution.

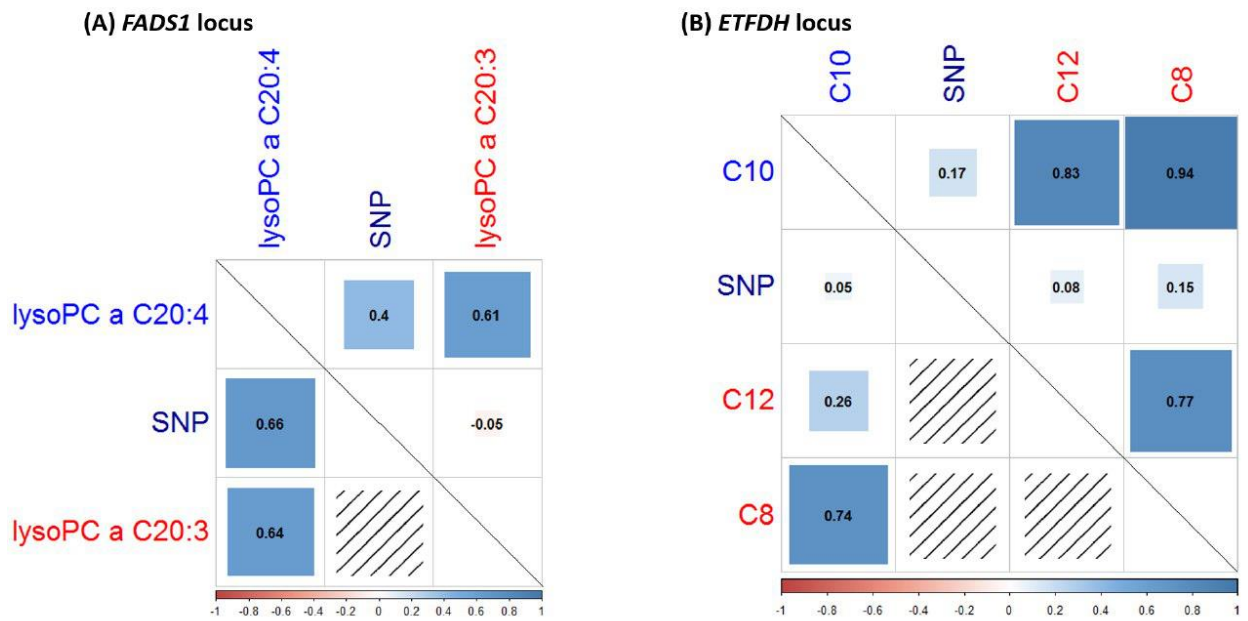
238 Next, we investigated the variance-covariance structure of the loci with positive and
239 negative pleiotropic terms. We therefore selected a locus in which the pleiotropic component’s
240 contribution to power was positive (rs174547 at *FADS1*) and a locus in which the pleiotropic
241 component’s contribution to power was negative (rs8396 at *ETFDH*). **Figure 2** shows the
242 corresponding correlations between the SNP, the trait, and the covariates involved, together with
243 the partial coefficients for the conditional regression of the trait on the SNP and the covariates.
244 With respect to *FADS1* (**Figure 2A**), the correlations between the SNP and the trait (lyso-
245 PC a C20:4) and between the SNP and the covariate (lyso-PC a C20:3) are in opposite directions,
246 generating negative genetically induced covariance between lyso-PC a C20:4 and lyso-
247 PC a C20:3. In contrast, the residual correlation between the trait and the covariate is positive.

248 Therefore, the value of the partial coefficient of regression between the SNP and lyso-PC a C20:4,
 249 conditional on lyso-PC a C20:3, is greater than that of the coefficient of regression without
 250 covariates.

251 With respect to the second example, *ETFDH* (**Figure 2B**), we found that the conditional
 252 regression of C10 on rs8396 and two covariates (C8 and C12, two medium-chain acylcarnitines)
 253 led to a smaller SNP partial regression coefficient compared to an unconditional regression; this
 254 is because all of the terms in $\sum_{i=1}^k \beta_i \rho_{gi} / \rho_{yg}$ are positive.
 255



256
 257 **Figure 1. Decomposition of the log-T² ratio for cGAS and uGAS into pleiotropic and noise**
 258 **components.** Vertically grouped trios (each composed of a square, triangle, and asterisk)
 259 correspond to one of fourteen associations in Table 1. The position of a trio on the x-axis
 260 corresponds to the log-ratio between conditional and univariate test statistic. On the y-axis,
 261 the asterisk corresponds to the log-ratio. The value of the pleiotropic component is depicted by a
 262 square, and the value of the noise component is depicted by a triangle. Each trio is shown in gray,
 263 except the trios representing the *ACADM*, *SLC22A4*, and *PLEKHH1* loci, for which we have two
 264 different associations. The two dotted lines correspond to the regression lines for the two
 265 components. The three dark-green vertical lines indicate the associations that were significant in
 266 the cGAS analysis but not in the uGAS analysis, and the two dark-red lines indicates the
 267 associations that were significant only in the uGAS analysis.



269

270 **Figure 2.** Matrix of correlations (above diagonal line) and the partial coefficients of regression of
 271 the trait of interest on the SNP genotype and covariate(s) (the first column below diagonal line)
 272 for the *FADS1* (A) and *ETFDH* (B) loci. Names of traits used as covariates are in red. The number
 273 in a cell indicates the value of correlation (partial regression coefficient). The area of a square is
 274 proportional to the absolute value of correlation (partial regression coefficient); the effect
 275 magnitude is also reflected by square's color (the scale provided at the bottom of the graph). The
 276 *FADS1* locus represents scenario in which the pleiotropic term in Eq. (2) is strongly positive, while
 277 for *ETFDH* this term is negative.

278

279 Although using a known biochemical network to select covariates has many advantages, it
 280 may be somewhat unpractical and perhaps even harmful, as our biochemical knowledge is still
 281 relatively incomplete. Therefore, we explored the potential of performing a cGAS in which the
 282 covariates are selected using a data-driven approach (GGM-cGAS). The network of metabolites
 283 was reconstructed using Gaussian Graphical Models based on partial correlations. For a given
 284 metabolite, we selected covariates based on significant partial correlations. Specifically, we used
 285 the following threshold as proposed previously [17]: a p -value \leq
 286 $(0.01/\text{number of calculated partial correlations})$, which corresponds to a cut-off at $p \leq 8.83 \times 10^{-7}$.
 287 The network used in our analysis is shown in **Supplementary Figure S1**.

288 To compare GGM-cGAS with BN-cGAS, we used the same set of metabolites that we used
 289 for BN-cGAS to run our GGM-cGAS analysis; these results are presented in **Supplementary**
 290 **Table S1**. We found 16 SNP-trait pairs clustered at 10 loci that were detected by either GGM-
 291 cGAS or BN-cGAS. More covariates were included in the GGM-cGAS analysis (ranging from 1

292 to 18, with mean and median values of 7.6 covariates and 7 covariates, respectively) than in the
293 BN-cGAS analysis. Thus, we predicted that GGM-cGAS would have relatively more power than
294 BN-cGAS due to reduced noise (term 1 in Eq. (2)); on the other hand, GGM-cGAS might lose
295 power because of reduced occurrence of unexpected pleiotropy (term 2 in Eq. (2)).

296 For the best SNP-trait pairs detected by GGM-cGAS or BN-cGAS, we computed the
297 components in Eq. (2) and compared these components using a paired-sample Wilcoxon test. We
298 found that the noise component in Eq. (2) was always larger for GGM-cGAS, with a mean
299 difference of 0.66 ($p=3 \times 10^{-5}$). Moreover, the second “pleiotropic” component in Eq. (2) was on
300 generally smaller in GGM-cGAS than for BN-cGAS, with a mean difference of -0.54 ($p=0.013$);
301 nevertheless, for three out of 16 GGM-cGAS SNP-trait pairs, the pleiotropic component was
302 positive. The average chi-square value was 33% smaller for GGM-cGAS than for BN-cGAS,
303 indicating an average loss of power (although this loss was not significant; $p=0.5$ based on a paired
304 Wilcoxon test), yet 22% larger for GGM-cGAS than for uGAS ($p=0.8$ based on a paired Wilcoxon
305 test). We therefore conclude that although GGM-cGAS may not serve as an ideal proxy for
306 analyzing a bona fide biochemical network, it can still have increased power due to reduced target
307 trait residual variance and the potential to detect unexpected pleiotropy.

308 Next, we investigated further the potential of using cGAS under realistic conditions to a
309 full extent by analyzing all 151 available metabolites using GGM-cGAS and comparing these
310 results with the results of uGAS (**Table 2** and **Supplementary Figure S2**). In total, uGAS detected
311 15 loci at the genome-wide significance level $p \leq 5 \times 10^{-8}/151$ (i.e., $p < 3.3 \times 10^{-10}$). On the other
312 hand, GGM-cGAS identified 19 significant loci using the same threshold. As expected, the
313 standard errors of the genetic effect estimates were smaller for GGM-cGAS than for uGAS (**Table**
314 **2** and **Supplementary Figure S3**). A total of 14 loci were detected by both uGAS and GGM-
315 cGAS. GGM-cGAS failed to identify one locus that was identified by uGAS (C5:1-DC at
316 rs2943644), but identified five loci that were missed by uGAS. Three of the five loci identified
317 solely by GGM-cGAS affect amino acids, and the remaining two loci affect acylcarnitines. It is
318 important to note that the loci identified by BN-cGAS (when we analyzed 105 metabolites) are a
319 subset of the 19 loci that were identified by GGM-cGAS (when we used all 151 metabolites).

320 Finally, we searched the available literature for the loci listed in **Table 2** (see
321 **Supplementary Note 2** for details). From the 20 loci that we report here, 15 were found to be
322 significant at the genome-wide level in a recent large ($n=7478$) meta-analysis of Biocrates
323 metabolomics data reported by Draisma et al. [7]. Some of the metabolites analyzed in our study
324 were not analyzed by Draisma et al. [7]; nevertheless, for 11 out of these 15 loci, we observed a
325 significant association for the same SNP-metabolite pair; for three loci, the strongest association
326 was with a metabolite in the same class, and for one locus the strongest association was with a

1 327 metabolite from a different lipid class (see **Supplementary Table S2**). For the other five loci that
2 328 were not significant in the study by Draisma et al. [7], we determined whether these five loci were
3 329 significant and replicated in a study by Tsepilov et al. [18]. It should be noted that Tsepilov et al.
4 330 analyzed the ratios of metabolites and also used the KORA F4 data set in their discovery stage,
5 331 although they used another cohort (TwinsUK) for replication. Of these five loci, two were also
6 332 significant in the study by Tsepilov et al. [18]; moreover, for both of these loci the metabolite
7 333 analyzed in our study was included in the ratios analyzed by Tsepilov et al. One of the five loci
8 334 was associated with the same trait in two other studies [19,20]. Finally, we found no prior
9 335 published evidence of any association with metabolites for rs2943644 (*LOC646736*) or
10 336 rs17112944 (*LOC728755*); therefore, we conclude the observed associations with rs17112944 and
11 337 rs2943644 as likely false positives, and these two loci were excluded from further consideration.
12 338
13 339
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Table 1. Eleven loci identified by BN-cGAS and uGAS on metabolites for which at least one one-reaction-step neighbor was available.

Locus	SNP	Metabolite	chr:pos	Gene	effA/refA	EAF	uGAS		cGAS		
							beta (se)	P-value	Beta (se)	P-value	N _{cov}
uGAS & cGAS											
1	rs211718	C8	1:75879263	<i>ACADM</i>	T/C	0.30	-0.45 (0.034)	3.26E-37	-0.10 (0.012)	4.83E-17	1
1	rs211718	C12	1:75879263	<i>ACADM</i>	T/C	0.30	-0.04 (0.036)	2.19E-01	0.20 (0.014)	1.67E-40	3
2	rs7705189	PC ae C42:5	5:131651257	<i>SLC22A4</i>	G/A	0.47	0.15 (0.034)	8.65E-06	0.06 (0.009)	1.49E-10	3
2	rs419291	C5	5:131661254	<i>SLC22A4</i>	T/C	0.38	0.26 (0.035)	7.03E-14	0.17 (0.029)	1.01E-08	1
3	rs9368564	PC aa C42:5	6:11168269	<i>ELOVL2</i>	G/A	0.25	-0.29 (0.039)	1.14E-13	-0.15 (0.024)	1.63E-10	3
4	rs12356193	C0	10:61083359	<i>SLC16A9</i>	G/A	0.17	-0.51 (0.046)	1.84E-27	-0.42 (0.042)	1.67E-22	1
5	rs174547	lyso-PC a C20:4	11:61327359	<i>FADS1</i>	C/T	0.70	0.61 (0.033)	1.24E-69	0.66 (0.024)	2.96E-141	1
6	rs2066938	C4	12:119644998	<i>ACADS</i>	G/A	0.27	0.73 (0.033)	2.42E-93	0.72 (0.032)	2.13E-100	1
7	rs10873201	PC ae C36:5	14:67036352	<i>PLEKHH1</i>	T/C	0.45	-0.26 (0.034)	4.37E-14	-0.21 (0.018)	2.38E-30	2
7	rs1077989	PC ae C32:2	14:67045575	<i>PLEKHH1</i>	C/A	0.46	-0.30 (0.034)	2.23E-18	-0.06 (0.016)	5.33E-05	3
8	rs4814176	PC ae C40:2	20:12907398	<i>SPTLC3</i>	T/C	0.36	0.24 (0.035)	5.74E-12	0.25 (0.023)	1.58E-25	4
Only uGAS											
9	rs8396	C10	4:159850267	<i>ETFDH</i>	C/T	0.71	0.26(0.037)	2.11E-12	0.05 (0.011)	6.67E-07	2
Only cGAS											
10	rs4862429	PC ae C42:5	4:186006834	<i>ACSL1</i>	T/C	0.31	0.02(0.037)	6.62E-01	-0.06 (0.010)	6.57E-11	3
11	rs603424	Lyso-PC a C16:1	10:102065469	<i>PKD2LI</i>	A/G	0.80	0.23(0.042)	5.34E-08	0.21 (0.031)	1.39E-11	1

Notes: The best SNP-metabolite pair is shown for each locus. chr:pos refers to the physical position of the SNP; EAF, effect allele frequency; beta (se), the estimated effect and standard error of the SNP; effA/refA, effect allele/reference allele; P-value, the p-value of the additive model; Gene, the most likely (according to DEPICT) associated gene in the region; N_{cov}, the number of covariates used in cGAS.

16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Table 2. Twenty loci identified by GGM-cGAS and uGAS.

Locus	SNP	Metabolite	chr:pos	Gene	effA/refA	EAF	uGAS		cGAS		
							beta (se)	P-value	beta (se)	P-value	N _{cov}
uGAS & cGAS											
1	rs211718	C6 (C4:1-DC)	1:75,879,263	<i>ACADM</i>	T/C	0.30	-0.48 (0.034)	4.64E-42	-0.13 (0.017)	2.00E-13	7
1	rs7552404	C6 (C4:1-DC)	1:75,908,534	<i>ACADM</i>	G/A	0.30	-0.48 (0.034)	3.10E-42	-0.12 (0.017)	3.25E-13	7
2	rs483180	Ser	1:120,069,028	<i>PHGDH</i>	G/C	0.30	-0.24 (0.037)	3.34E-11	-0.24 (0.028)	1.50E-17	2
2	rs477992	Ser	1:120,059,099	<i>PHGDH</i>	A/G	0.70	0.24 (0.037)	5.15E-11	0.24 (0.028)	5.82E-18	2
3	rs2286963	C9	2:210,768,295	<i>ACADL</i>	G/T	0.63	-0.49 (0.032)	1.10E-49	-0.48 (0.027)	1.48E-67	3
4	rs8396	C10	4:159,850,267	<i>ETFDH</i>	C/T	0.71	0.26 (0.037)	2.02E-12	0.04 (0.010)	1.49E-05	8
4	rs8396	C7-DC	4:159,850,267	<i>ETFDH</i>	C/T	0.71	-0.09 (0.037)	1.67E-02	-0.13 (0.020)	3.29E-11	8
5	rs419291	C5	5:131,661,254	<i>SLC22A4</i>	T/C	0.38	0.26 (0.035)	7.03E-14	0.17 (0.026)	2.28E-10	3
5	rs270613	C5	5:131,668,482	<i>SLC22A4</i>	A/G	0.61	-0.26 (0.035)	7.93E-14	-0.17 (0.026)	8.48E-11	3
6	rs9393903	PC aa C42:5	6:11,150,895	<i>ELOVL2</i>	A/G	0.75	0.29 (0.039)	2.19E-13	0.18 (0.020)	4.51E-19	6
6	rs9368564	PC aa C42:5	6:11,168,269	<i>ELOVL2</i>	G/A	0.25	-0.29 (0.039)	1.14E-13	-0.19 (0.021)	7.84E-19	6
7	rs816411	Ser	7:56,138,983	<i>PHKG1</i>	C/T	0.51	-0.22 (0.034)	2.15E-10	-0.19 (0.026)	5.16E-13	2
7	rs1894832	Ser	7:56,144,740	<i>PHKG1</i>	C/T	0.51	0.21 (0.034)	3.23E-10	0.19 (0.026)	1.69E-13	2
8	rs12356193	C0	10:61,083,359	<i>SLC16A9</i>	G/A	0.17	-0.51 (0.046)	1.84E-27	-0.27 (0.034)	9.72E-16	3
9	rs174547	lyso-PC a C20:4	11:61,327,359	<i>FADS1</i>	C/T	0.70	0.61 (0.033)	1.44E-69	0.07 (0.011)	1.41E-10	9
9	rs174556	PC ae C44:4	11:61,337,211	<i>FADS1</i>	T/C	0.27	0.09 (0.038)	1.55E-02	0.21 (0.014)	3.16E-46	3
10	rs2066938	C4	12:119,644,998	<i>ACADS</i>	G/A	0.27	0.73 (0.033)	5.87E-94	0.71 (0.025)	1.31E-151	2
11	rs12879147	PC aa C28:1	14:63,297,349	<i>SYNE2</i>	A/G	0.85	-0.46 (0.050)	2.07E-19	-0.12 (0.019)	5.94E-11	14
11	rs17101394	SM(OH) C14:1	14:63,302,139	<i>SYNE2</i>	A/G	0.83	-0.32 (0.050)	1.00E-10	-0.10 (0.011)	1.17E-17	7
12	rs1077989	PC ae C36:5	14:67,045,575	<i>PLEKHH1</i>	C/A	0.46	-0.26 (0.034)	3.42E-14	-0.08 (0.010)	8.25E-16	10
12	rs1077989	PC ae C32:2	14:67,045,575	<i>PLEKHH1</i>	C/A	0.46	-0.30 (0.034)	2.23E-18	-0.05 (0.016)	1.31E-03	6
13	rs4814176	SM(OH).C22:1	20:12,907,398	<i>SPTLC3</i>	T/C	0.36	0.03 (0.035)	4.51E-01	-0.07 (0.009)	1.10E-16	10

13	rs4814176	SM(OH) C24:1	20:12,907,398	<i>SPTLC3</i>	T/C	0.36	0.24 (0.035)	5.40E-12	0.09 (0.013)	3.04E-11	9
14	rs5746636	Pro	22:17,276,301	<i>PRODH</i>	T/G	0.24	-0.31 (0.039)	3.00E-15	-0.32 (0.034)	1.91E-20	2
Only uGAS											
15	rs2943644	C5:1-DC	2:226,754,586	<i>LOC646736</i>	C/T	0.68	0.32 (0.042)	5.14E-14	0.09 (0.022)	3.58E-05	5
Only cGAS											
16	rs1374804	Gly	3:127,391,188	<i>ALDH1L1</i>	A/G	0.64	0.20 (0.036)	1.88E-08	0.21 (0.030)	8.08E-13	3
17	rs4862429	PC ae C42:5	4:186,006,834	<i>ACSL1</i>	T/C	0.31	0.02 (0.037)	6.62E-01	-0.06 (0.008)	1.25E-12	8
18	rs603424	C16:1	10:102,065,469	<i>PKD2L1</i>	A/G	0.80	0.16 (0.042)	9.51E-05	0.14 (0.018)	1.32E-13	9
19	rs2657879	Gln	12:55,151,605	<i>GLS2</i>	G/A	0.21	-0.24 (0.042)	2.82E-08	-0.27 (0.031)	9.37E-18	5
20	rs17112944	C6:1	14:27,179,297	<i>LOC728755</i>	A/G	0.90	-0.28 (0.059)	2.09E-06	-0.21 (0.032)	1.38E-10	9

Notes: The best SNP-metabolite pair is shown for each locus. chr:pos refers to the physical position of the SNP; EAF, effect allele frequency; beta (se), the estimated effect and standard error of the SNP; effA/refA, effect allele/reference allele; *P*-value, the *p*-value of the additive model; Gene, the most likely (according to DEPICT) associated gene in the region; N_{cov}, the number of covariates used in cGAS.

1 351 **Discussion**

2
3
4 352 We report a new “trait-centric” approach for analyzing genetic determinants of multivariate
5 353 “omics” traits by performing a network-based conditional genetic association analysis (cGAS). In
6
7 354 the context of metabolomics, for each trait we selected a set of other metabolites to be used as
8
9 355 covariates in our genetic association analysis. The selection of covariates can be either mechanistic
10
11 356 (e.g., based on known biological relationships between traits of interest) or data-driven (e.g., based
12
13 357 on partial correlations). Importantly, this approach can use either individual-level or summary-
14
15 358 level data. We first mathematically compared the power of conditional and standard single-trait
16
17 359 genetic association analyses (univariate genetic association, uGAS), and we identified scenarios
18
19 360 in which these analyses are expected to produce different results; next, we applied cGAS to 151
20
21 361 metabolomics traits (Biocrates panel) in a large (n=1784 individuals) population-based KORA
22
23 362 cohort.

24 363 We found that the log-ratio between the cGAS and uGAS test statistic can be decomposed
25
26 364 in a “noise” component (which depends on residual variance of the trait and is always positive)
27
28 365 and a “pleiotropic” component. The pleiotropic component is negative in cases in which
29
30 366 genetically induced covariance (between the trait of interest and the trait used as the covariate) and
31
32 367 the residual covariance have the same sign (i.e., act in the same direction). The pleiotropic
33
34 368 component is positive in cases in which the genetically induced covariance and residual covariance
35
36 369 act in opposite directions.

37 370 Should one expect that genetically induced and residual covariance act in the same or
38
39 371 opposite directions? In the context of complex polygenic traits, one would expect that genetic and
40
41 372 environmental correlations have the same sign. This is reflected in animal breeding studies; for a
42
43 373 recent example in humans, see [21]. It should be noted that a negative sign for the pleiotropic
44
45 374 component does not necessarily indicate higher power of the uGAS, as the noise component may
46
47 375 still dominate the relative non-centrality parameter. This will happen, for example, when ρ_{cg} (the
48
49 376 effect of the genotype on the covariate) is small while β_{yc} (partial residual regression between the
50
51 377 trait and covariate) is relatively large, thereby reducing σ_c^2 .

52 378 Nevertheless, in the case of metabolomic traits, genetic and environmental sources do not
53
54 379 necessarily generate consistent covariance. Moreover, for a given locus that affects the activity of
55
56 380 an enzyme involved in a biochemical reaction, the unexpected inconsistency between genetically
57
58 381 induced covariance and residual covariance may not be so unexpected after all. Indeed, consider
59
60 382 an allele associated with an increased activity of an enzyme that converts substrate A into product
61
62 383 B. One would expect that the levels of A and B are positively correlated; one would also expect

384 that the allele is positively correlated with the level of product B and negatively correlated with
385 the level of substrate A. This is precisely the scenario that yields a positive value for the second
386 term in Eq. (1), thus providing an additional increase in power above and beyond the power
387 provided by the first term in Eq. (1) (noise reduction).

388 Our empirical investigation of real data on the genetic association between the genome and
389 metabolites confirmed the existence of both scenarios. An extreme example of concordance
390 between genetic covariance and residual covariance is provided by the effects of rs8396 on C10,
391 with C8 and C12 used as covariates (see Figure 2B). The *ETFDH* gene, which was prioritized by
392 DEPICT as the best candidate in this region (with a false-discover rate <5%), encodes the enzyme
393 electron transfer flavoprotein (ETF) dehydrogenase, which plays a role in mitochondrial fatty acid
394 oxidation. During this process, the acyl group is transferred from a long chain acylcarnitine to a
395 long-chain acetyl-CoA, which is then catabolized. ETF dehydrogenase participates in the catabolic
396 process by transferring electrons from acyl-CoA dehydrogenase to the oxidative phosphorylation
397 pathway. Thus, the *ETFDH* gene should affect all forms of long-chain acylcarnitines in the same
398 way, and we can expect that the pleiotropic effect of this gene on the acylcarnitines in our example
399 (C8, C10, C12, etc.) will be unidirectional. The presence of unidirectional genetic effects and the
400 positive correlation between these acylcarnitines makes the second term in Eq. (2) negative, which
401 determines that—in this situation—univariate GAS has more power than cGAS.

402 An empirical example of discordance between genetically induced covariance and residual
403 covariance is provided by the effects of the SNP rs174547 on lyso-PC a C20:4, with lyso-PC a
404 C20:3 used as a covariate. This SNP exhibits opposite correlations with lyso-PC a C20:4 and lyso-
405 PC a C20:3, resulting in negative genetically induced covariance between these traits. At the same
406 time, the residual correlation between these traits is positive, resulting in steep increase in the
407 power of conditional analysis. In this region, the *FADS1/2/3* gene cluster is an attractive candidate,
408 providing the detected model with biological relevance. The *FADS1* gene encodes the enzyme
409 fatty acid desaturase 1, whereas the two traits differ by only one double bond. Thus, this example
410 mimics perfectly the biochemical scenario in which we would expect a conditional analysis to
411 have increased power.

412 The trait-centric methods considered here provide an attractive framework to identify and
413 study direct genetic effects on a trait of interest. Conditional analysis is an attractive option in cases
414 in which we wish to clearly interpret the results in terms of the effect of the genotype on a particular
415 trait. Such specific interpretation may be important when comparing genetic association results
416 obtained for our trait of interest with results obtained for other traits (e.g., using the methods
417 described in [22–24]). It should be noted, though, that a trait-centric approach is not intended to
418 maximize the power of identifying genes that affect metabolomics as a whole. Such a gene-centric

1
2
3
4
5
6
7
8
9
10
419 view would favor analysis using joint—and not conditional—modeling of sets of traits. Such an
420 approach can maintain power across a wide range of scenarios, including the scenario of
421 concordance between genetically induced and residual covariance [13]. In this gene-centric
422 framework, other formulations of conditional analysis have also been proposed [25] in order to
423 specifically increase power of gene identification by selecting covariates that—using our
424 terminology—affect the “noise reduction” component of the model while avoiding the problems
425 associated with the pleiotropic component.

11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
426 The proper selection of sets of biologically related traits is extremely important for the
427 conditional genetic association analysis method described here, as well as for multivariate methods
428 that model the joint effects of genotype on an ensemble of traits. Here, we considered two
429 alternative approaches—knowledge based and data-driven—to finding the networks of related
430 traits, with a subnetwork centered around a trait of interest used as the analyzed set. In principle,
431 in the context of analyzed metabolomics data, the knowledge-based network approach has slightly
432 higher power in the context of trait-centric genetic association analysis. However, we believe that
433 our analysis revealed that both approaches are suboptimal. The knowledge-based network
434 reconstruction has many advantages, but it may be somewhat unpractical, as our biochemical
435 knowledge is still relatively incomplete. Secondly, by reconstructing the network while relying
436 only on current knowledge, we may be missing new knowledge that may be revealed by the data.
437 Finally, by including neighbors that are based only on biochemical information, we may miss
438 covariance induced by technical confounders; adjusting for this may increase the power of analysis
439 [25]. Learning the network from the same data that were used for genetic analysis has the
440 disadvantages of potentially ignoring existing knowledge and being sensitive to sample size.
441 Finally, we note that the total observed correlation between metabolites is determined by the
442 balance between genetic and environmental sources of covariance; it is possible to imagine a
443 situation in which total correlation is smaller than one or more of its components, and our analysis
444 provides examples of such a situation. We may speculate that—ideally—one should use a method
445 that allows one to combine prior knowledge and new information obtained from the data, thereby
446 allowing the simultaneous learning of the structure of dependencies between different metabolites
447 and between the metabolites and the genome. Such learning from the data while allowing for the
448 incorporation of previous knowledge (e.g., biochemical relations between traits) might be achieved
449 by applying, e.g. a machine-learning approach that allows for differential shrinkage. It is also
450 important to note that the proper application of such an approach would require the availability of
451 vast samples of data, thereby allowing for separate training, validation, verification, testing, and
452 replication of detected dependencies and associations.

454 **Materials and Methods**

455 **KORA study**

456 The KORA study (Cooperative Health Research in the region of Augsburg) is a series of
457 population-based studies in the region of Augsburg in Southern Germany [26]. KORA F4 is a
458 follow-up survey (conducted from 2006 through 2008) of the baseline KORA S4 survey, which
459 was conducted from 1999 through 2001. All study protocols were approved by the ethics
460 committee of the Bavarian Medical Chamber, and all participants provided written informed
461 consent.

462 The concentration of 163 metabolites were measured in 3061 serum samples obtained from
463 KORA F4 participants using flow injection electrospray ionization tandem mass spectrometry and
464 the AbsoluteIDQ p150 Kit (Biocrates Life Sciences AG, Innsbruck, Austria) [27]. After applying
465 quality control screening, a total of 151 metabolite measurements were used in our analysis.
466 Details regarding the methods and quality control of the metabolite measurements, as well as
467 details regarding the metabolite nomenclature, have been published previously [27]. The
468 nomenclature for the metabolites in this study is provided in **Supplementary Table S3**.

469 Genotyping was performed using the Affymetrix 6.0 SNP array (534,174 SNP markers
470 after quality control), with further imputation using HapMap2 (release 22) as a reference panel,
471 resulting in a total of 1,717,498 SNPs (for details, see ref. [28]). Both the metabolite concentrations
472 and genotype were available for 1785 participants in the KORA F4 study.

473

474 **Statistical analysis**

475 Partial correlation coefficients and their p -values were calculated using the “ppcor” package [29]
476 in R. Graphical representations were generated using the “ggm” [30] package in R. Consistent
477 with previous studies [17], we considered a partial correlation coefficient to be significant at p
478 $< 0.01/(151*150/2)$ (i.e., $p < 8.83 \times 10^{-7}$).

479 For the GWAS analysis, we used OmicABEL software [31]. Prior to GWAS, all traits were
480 first adjusted for the participant’s sex, age, and batch effect; subsequently, the residual traits were
481 transformed using an inverse-normal transformation [32]. The genotypes from the KORA F4
482 cohort were used. Only SNPs that had a call rate ≥ 0.95 , $R^2 \geq 0.3$, Hardy–Weinberg equilibrium
483 (HWE) $p \geq 10^{-6}$, and MAF ≥ 0.1 (1,717,498 SNPs in total) were included in the analysis. The
484 genomic control method was used to correct for any possible inflation of the test statistics. The
485 genomic control [33] lambda value for all traits was between 1.00 and 1.03.

486 In a specific analysis (i.e., cGAS or uGAS), we defined independent loci as groups of
487 genome-wide significant associations that were separated by at least 500 kb or were located on

488 different chromosomes. The strongest association (i.e., the association with the lowest p -value)
489 was selected to represent this locus. The cGAS and uGAS results were considered to reflect
490 different loci if the strongest associations were in loci that were separated by at least 500. The
491 threshold for the genome-wide significance for 151 traits was set to $p=5 \times 10^{-8}/151$ (i.e.,
492 $p=3.31 \times 10^{-10}$).

493 When partitioning the log(cGAS/uGAS) test statistics into the noise and pleiotropic
494 components (see Eq. (2) and **Figure 1**), we used all known loci that were significant in either the
495 cGAS or uGAS analysis (see Table 1). If a locus included two SNPs associated with different
496 traits, we included both associations during partitioning. If a locus included two SNPs associated
497 with the same trait, to be conservative we included only the locus with the lower uGAS p -value
498 during partitioning. After partitioning, we compared the pleiotropic and noise components using
499 the paired-samples Wilcoxon test. For comparing the chi-square test results for the two methods,
500 for each locus we first selected the method that yielded the strongest association (and hence the
501 largest chi-square value). We compared that chi-square value with the maximal chi-square value
502 observed for the second method within a 500-kb region centered around the strongest association
503 observed using the first method.

504 ***In silico* functional annotation**

506 We conducted functional annotation for our findings. To prioritize genes in associated regions,
507 gene set enrichment, and tissue/cell-type enrichment analyses, we used DEPICT software [34]
508 (release 140721) with the following settings: flag_loci = 1; flag_genes = 1; flag_genesets = 1;
509 flag_tissues = 1; param_ncores = 2; and further manual annotation (h37 assembly). All 27 SNPs
510 (clustered in 20 loci) identified by cGAS or uGAS (see **Table 2**) were included in the analysis. If
511 more than one gene was annotated for a SNP by DEPICT, we selected the gene with the lowest
512 nominal DEPICT P-value. In most cases, the results of manual annotation matched the annotation
513 results using DEPICT annotation (see **Supplementary Note 2**). In addition, we looked up each
514 SNP using the Phenoscanner [35] database to check whether it was previously reported to be
515 associated with metabolic traits at $p < 5 \times 10^{-8}$ and proxy $r^2 < 0.7$.

516

517 **Additional Files**

- 1 518 Supplementary Note 1 – cGAS using summary level data
2
3 519 Supplementary Note 2 – Literature search for loci identified by cGAS and uGAS
4
5 520 Supplementary Tables
6
7 521 Supplementary Table S1 – BN- cGAS and GGM- cGAS for 105 metabolites
8
9 522 Supplementary Table S2 – GGM-cGAS and uGAS for 151 metabolites
10
11 523 Supplementary Table S3 - List of metabolites measured using the AbsoluteIDQ p150 Kit
12
13 524 Supplementary Figures
14
15 525 Supplementary Figure S1 – Partial correlations network
16
17 526 Supplementary Figure S2 – Manhattan plots for cGAS and uGAS for 151 metabolites
18
19 527 Supplementary Figure S3 – Comparison of effect estimates and their standard errors for
20
21 528 SNPs from Table 2
22

23 529

24 530 **Abbreviations**

- 25
26 531 GWAS – genome-wide association study
27
28 532 cGAS – conditional GWAS
29
30 533 uGAS – univariate GWAS (trait-by-trait)
31
32 534 BN-cGAS – cGAS based on biochemical networks
33
34 535 GGM-cGAS – Gaussian Graphical Modeling cGAS based on partial correlations network
35

36 536

37 537 **Acknowledgments**

38
39
40 538 We thank Athina Spilopoulou and Felix Agakov for helpful discussions. We also thank Alexander
41
42 539 Zlobin and Alexander Grishenko for help preparing the tables and figures in the manuscript, and
43
44 540 we thank Sophie Molnos for help with data management.
45

46 541

47 542 **Funding**

48
49
50 543 The KORA study was initiated and financed by the Helmholtz Center Munich – German
51
52 544 Research Center for Environmental Health, which is funded by the German Federal Ministry of
53
54 545 Education and Research (BMBF) and by the State of Bavaria. Furthermore, the KORA study
55
56 546 was supported by the Munich Center of Health Sciences (MC-Health), Ludwig Maximilian
57
58 547 University of Munich, as part of the LMUinnovativ project.

59 548 This work was supported by the European Union FP7 framework project Pain-Omics (grant
60
61 549 number 602736).
62
63
64
65

1 550 SS was supported by the Russian Ministry of Science and Education under the 5-100 Excellence
2 551 Programme. YA and YT were supported by the Federal Agency of Scientific Organisations via
3
4 552 the Institute of Cytology and Genetics (project number 0324-2018-0017)
5

6 553

7 **Authors Contributions**

8
9 554 YT, CG, and YA designed and supervised the study; PC, CP, JA, KG, and RW-S collected the
10 555 data; CG and KS contributed data for the analysis; YT, OZ, and SS analyzed the data; YT, YA,
11 556 CG, OZ, JK, and KS discussed and interpreted the results; YT, OZ, CG, and YA wrote the
12 557 manuscript. All authors contributed to and approve the final version of the manuscript.
13
14 558
15

16 559

17 **Availability of Data and Materials**

18
19 560 The code produced in relation to this work is distributed under the LGPL license and is available
20 561 from [WHILE MANUSCRIPT IS UNDER REVIEW, THE LINK TO THE DATA IS ONLY
21 562 AVAILABLE TO THE REVIEWERS FROM THE EDITOR]. All summary statistics and
22 563 association data that are necessary to reproduce our results are licensed under CC0 and are
23 564 accessible online at [WHILE MANUSCRIPT IS UNDER REVIEW, THE LINK TO THE
24 565 DATA IS ONLY AVAILABLE TO THE REVIEWERS FROM THE EDITOR]. The informed
25 566 consent given by the KORA study participants does not cover the posting of participant-level
26 567 phenotype or genotype data in public databases. However, the KORA data are available upon
27 568 request from KORA-gen (<https://www.helmholtz-muenchen.de/en/kora/index.html>). Requests
28 569 can be submitted online and are subject to approval by the KORA board.
29
30 570
31
32
33
34
35
36
37
38
39
40
41

42 571

43 **Competing Interests**

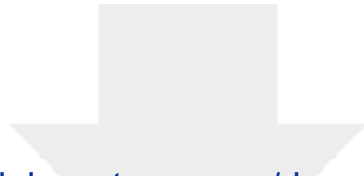
44
45 572 Y. Aulchenko is the founder and co-owner of PolyOmica, a private research organization that
46 573 specializes in computational and statistical (gen)omics.
47 574
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

- 1 577 1. Visscher PM, Brown MA, McCarthy MI, Yang J. Five Years of GWAS Discovery. *Am. J.*
2 578 *Hum. Genet.* [Internet]. 2012;90:7–24. Available from:
3 579 <http://linkinghub.elsevier.com/retrieve/pii/S0002929711005337>
- 4 580 2. Ritchie MD, Holzinger ER, Li R, Pendergrass SA, Kim D. Methods of integrating data to
5 581 uncover genotype–phenotype interactions. *Nat. Rev. Genet.* [Internet]. 2015;16:85–97. Available
6 582 from: <http://www.nature.com/doi/10.1038/nrg3868>
- 7 583 3. van der Sijde MR, Ng A, Fu J. Systems genetics: From GWAS to disease pathways. *Biochim.*
8 584 *Biophys. Acta - Mol. Basis Dis.* [Internet]. 2014;1842:1903–9. Available from:
9 585 <http://linkinghub.elsevier.com/retrieve/pii/S0925443914001124>
- 10 586 4. Hicks AA, Pramstaller PP, Johansson A, Vitart V, Rudan I, Ugocsai P, et al. Genetic
11 587 determinants of circulating sphingolipid concentrations in European populations. *PLoS Genet.*
12 588 [Internet]. 2009 [cited 2013 Dec 19];5:e1000672. Available from:
13 589 <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2745562&tool=pmcentrez&renderty>
14 590 [pe=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2745562&tool=pmcentrez&renderty)
- 15 591 5. Suhre K, Shin S-Y, Petersen A-K, Mohny RP, Meredith D, Wägele B, et al. Human
16 592 metabolic individuality in biomedical and pharmaceutical research. *Nature* [Internet]. 2011 [cited
17 593 2013 Dec 19];477:54–60. Available from:
18 594 <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3832838&tool=pmcentrez&renderty>
19 595 [pe=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3832838&tool=pmcentrez&renderty)
- 20 596 6. Inouye M, Ripatti S, Kettunen J, Lyttikäinen L-P, Oksala N, Laurila P-P, et al. Novel Loci for
21 597 metabolic networks and multi-tissue expression studies reveal genes for atherosclerosis. Visscher
22 598 PM, editor. *PLoS Genet.* [Internet]. 2012;8:e1002907. Available from:
23 599 <http://dx.plos.org/10.1371/journal.pgen.1002907>
- 24 600 7. Draisma HHM, Pool R, Kobl M, Jansen R, Petersen A-K, Vaarhorst AAM, et al. Genome-
25 601 wide association study identifies novel genetic variants contributing to variation in blood
26 602 metabolite levels. *Nat. Commun.* [Internet]. England; 2015;6:7208. Available from:
27 603 <http://www.ncbi.nlm.nih.gov/pubmed/26068415>
- 28 604 8. Kettunen J, Demirkan A, Würtz P, Draisma HHMM, Haller T, Rawal R, et al. Genome-wide
29 605 study for circulating metabolites identifies 62 loci and reveals novel systemic effects of LPA.
30 606 *Nat. Commun.* [Internet]. 2016;7:11122. Available from:
31 607 <http://www.ncbi.nlm.nih.gov/pubmed/27005778> <http://www.nature.com/doi/10.1038/ncomms11122>
- 32 608 9. Cichonska A, Rousu J, Martinen P, Kangas AJ, Soininen P, Lehtimäki T, et al. metaCCA:
33 609 summary statistics-based multivariate meta-analysis of genome-wide association studies using
34 610 canonical correlation analysis. *Bioinformatics* [Internet]. 2016;32:1981–9. Available from:
35 611 <http://www.ncbi.nlm.nih.gov/pubmed/27153689>
- 36 612 10. Stephens M. A unified framework for association analysis with multiple related phenotypes.
37 613 *PLoS One* [Internet]. 2013;8:e65245. Available from:
38 614 <http://www.ncbi.nlm.nih.gov/pubmed/23861737>
- 39 615 11. O'Reilly PF, Hoggart CJ, Pomyen Y, Calboli FCF, Elliott P, Jarvelin M-R, et al. MultiPhen:
40 616 joint model of multiple phenotypes can increase discovery in GWAS. *PLoS One* [Internet]. 2012
41 617 [cited 2014 Sep 20];7:e34861. Available from:
42 618 <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3342314&tool=pmcentrez&renderty>
43 619 [pe=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3342314&tool=pmcentrez&renderty)
- 44 620 12. Galesloot TE, van Steen K, Kiemeny LALM, Janss LL, Vermeulen SH. A comparison of
45 621 multivariate genome-wide association methods. *PLoS One* [Internet]. 2014 [cited 2014 Sep
46 622 20];9:e95923. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/24763738>
- 47 623 13. Shen X, Klarić L, Sharapov S, Mangino M, Ning Z, Wu D, et al. Multivariate discovery and
48 624 replication of five novel loci associated with Immunoglobulin G N-glycosylation. *Nat. Commun.*
49 625 [Internet]. 2017;8:447. Available from: <http://www.nature.com/articles/s41467-017-00453-3>
- 50 626

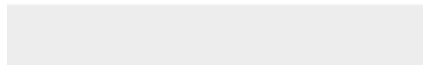
- 627 14. Schaid DJ, Tong X, Larrabee B, Kennedy RB, Poland GA, Sinnwell JP. Statistical Methods
628 for Testing Genetic Pleiotropy. *Genetics*. 2016;204:483–97.
- 1 629 15. Deng Y, Pan W. Conditional analysis of multiple quantitative traits based on marginal
2 630 GWAS summary statistics. *Genet. Epidemiol.* [Internet]. 2017;41:427–36. Available from:
3 631 <http://www.ncbi.nlm.nih.gov/pubmed/28464407>
- 4 632 16. Smith GD, Ebrahim S. “Mendelian randomization”: can genetic epidemiology contribute to
5 633 understanding environmental determinants of disease? *Int. J. Epidemiol.* [Internet]. 2003;32:1–
6 634 22. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/12689998>
- 7 635 17. Krumsiek J, Suhre K, Illig T, Adamski J, Theis FJ. Gaussian graphical modeling reconstructs
8 636 pathway reactions from high-throughput metabolomics data. *BMC Syst. Biol.* [Internet]. BioMed
9 637 Central Ltd; 2011 [cited 2013 May 23];5:21. Available from:
10 638 <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3224437&tool=pmcentrez&renderty>
11 639 [pe=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3224437&tool=pmcentrez&renderty)
- 12 640 18. Tsepilov YA, Shin S-Y, Soranzo N, Spector TD, Prehn C, Adamski J, et al. Nonadditive
13 641 Effects of Genes in Human Metabolomics. *Genetics* [Internet]. 2015;200:707–18. Available
14 642 from: <http://www.genetics.org/cgi/doi/10.1534/genetics.115.175760>
- 15 643 19. Xie W, Wood AR, Lyssenko V, Weedon MN, Knowles JW, Alkayyali S, et al. Genetic
16 644 variants associated with glycine metabolism and their role in insulin sensitivity and type 2
17 645 diabetes. *Diabetes* [Internet]. 2013;62:2141–50. Available from:
18 646 <http://www.ncbi.nlm.nih.gov/pubmed/23378610>
- 19 647 20. Shin S-Y, Fauman EB, Petersen A-K, Krumsiek J, Santos R, Huang J, et al. An atlas of
20 648 genetic influences on human blood metabolites. *Nat. Genet.* [Internet]. 2014 [cited 2014 May
21 649 12];46:543–50. Available from: <http://www.nature.com/doi/10.1038/ng.2982>
- 22 650 21. Bulik-Sullivan B, Finucane HK, Anttila V, Gusev A, Day FR, Loh P-R, et al. An atlas of
23 651 genetic correlations across human diseases and traits. *Nat. Genet.* Nature Publishing Group;
24 652 2015;47:1236–41.
- 25 653 22. Zhu Z, Zhang F, Hu H, Bakshi A, Robinson MR, Powell JE, et al. Integration of summary
26 654 data from GWAS and eQTL studies predicts complex trait gene targets. *Nat. Genet.* [Internet].
27 655 2016;48:481–7. Available from:
28 656 <http://www.nature.com/doi/10.1038/ng.3538>
29 657 <http://www.ncbi.nlm.nih.gov/pubmed/27019110>
- 30 658 23. Pickrell JK, Berisa T, Liu JZ, Séguérel L, Tung JY, Hinds DA. Detection and interpretation of
31 659 shared genetic influences on 42 human traits. *Nat. Genet.* 2016;19885.
- 32 660 24. Giambartolomei C, Vukcevic D, Schadt EE, Franke L, Hingorani AD, Wallace C, et al.
33 661 Bayesian Test for Colocalisation between Pairs of Genetic Association Studies Using Summary
34 662 Statistics. Williams SM, editor. *PLoS Genet.* [Internet]. 2014;10:e1004383. Available from:
35 663 <http://dx.plos.org/10.1371/journal.pgen.1004383>
- 36 664 25. Aschard H, Guillemot V, Vilhjalmsson B, Patel CJ, Skurnik D, Ye CJ, et al. Covariate
37 665 selection for association screening in multiphenotype genetic studies. *Nat. Genet.* [Internet].
38 666 2017; Available from: <http://www.ncbi.nlm.nih.gov/pubmed/29038595>
- 39 667 26. Wichmann H-E, Gieger C, Illig T. KORA-gen--resource for population genetics, controls
40 668 and a broad spectrum of disease phenotypes. *Gesundheitswesen* [Internet]. 2005 [cited 2013 Jun
41 669 6];67 Suppl 1:S26-30. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/16032514>
- 42 670 27. Illig T, Gieger C, Zhai G, Römisch-Margl W, Wang-Sattler R, Prehn C, et al. A genome-
43 671 wide perspective of genetic variation in human metabolism. *Nat. Genet.* [Internet]. Nature
44 672 Publishing Group; 2010 [cited 2013 May 23];42:137–41. Available from:
45 673 <http://www.ncbi.nlm.nih.gov/pubmed/20037589>
- 46 674 28. Kolz M, Johnson T, Sanna S, Teumer A, Vitart V, Perola M, et al. Meta-analysis of 28,141
47 675 individuals identifies common variants within five new loci that influence uric acid
48 676 concentrations. *PLoS Genet.* [Internet]. 2009 [cited 2013 May 30];5:e1000504. Available from:
49 677 <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2683940&tool=pmcentrez&renderty>
50 678 [pe=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2683940&tool=pmcentrez&renderty)

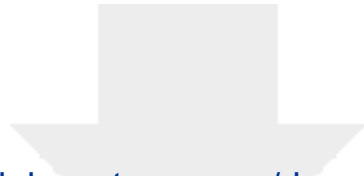
679 29. Kim S. ppcor: An R Package for a Fast Calculation to Semi-partial Correlation Coefficients.
680 Commun. Stat. Appl. Methods [Internet]. 2015;22:665–74. Available from:
681 <http://www.csam.or.kr/journal/view.html?doi=10.5351/CSAM.2015.22.6.665>
682 30. Marchetti GM. Independencies Induced from a Graphical Markov Model after
683 Marginalization and Conditioning: The R Package ggm. J. Stat. Softw. [Internet]. 2006;15.
684 Available from: <http://www.jstatsoft.org/v15/i06/>
685 31. Fabregat-Traver D, Sharapov SZ, Hayward C, Rudan I, Campbell H, Aulchenko Y, et al.
686 High-Performance Mixed Models Based Genome-Wide Association Analysis with omicABEL
687 software. F1000Research [Internet]. 2014;3:200. Available from:
688 <http://f1000research.com/articles/3-200/v1>
689 32. Beasley TM, Erickson S, Allison DB. Rank-based inverse normal transformations are
690 increasingly used, but are they merited? Behav. Genet. [Internet]. 2009 [cited 2013 Nov
691 7];39:580–95. Available from:
692 <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2921808&tool=pmcentrez&renderty>
693 [pe=abstract](#)
694 33. Devlin B, Roeder K. Genomic control for association studies. Biometrics [Internet]. 1999
695 [cited 2013 Jun 5];55:997–1004. Available from:
696 <http://www.ncbi.nlm.nih.gov/pubmed/11315092>
697 34. Pers TH, Karjalainen JM, Chan Y, Westra H-J, Wood AR, Yang J, et al. Biological
698 interpretation of genome-wide association studies using predicted gene functions. Nat. Commun.
699 [Internet]. 2015;6:5890. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/25597830>
700 35. Staley JR, Blackshaw J, Kamat MA, Ellis S, Surendran P, Sun BB, et al. PhenoScanner: a
701 database of human genotype–phenotype associations. Bioinformatics [Internet]. 2016;32:3207–
702 9. Available from:
703 <http://bioinformatics.oxfordjournals.org/lookup/doi/10.1093/bioinformatics/btw373>
704



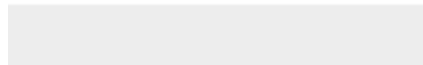


Click here to access/download
Supplementary Material
Supplementary Table 1.xlsx





Click here to access/download
Supplementary Material
Supplementary Table 2.xlsx

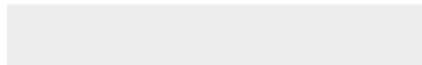




Click here to access/download
Supplementary Material
Supplementary Table 3.docx



Click here to access/download
Supplementary Material
Supplementary Note 1.docx





Click here to access/download
Supplementary Material
Supplementary Note 2.docx

