

# Heritability estimates for 361 blood metabolites across 40 genome-wide association studies

Fiona A. Hagenbeek<sup>1,2\*</sup>, René Pool<sup>1,2</sup>, Jenny van Dongen<sup>1,2</sup>, Harmen H.M. Draisma<sup>1</sup>, Jouke Jan Hottenga<sup>1</sup>, Gonneke Willemsen<sup>1</sup>, Abdel Abdellaoui<sup>1,51</sup>, Iryna O. Fedko<sup>1</sup>, Anouk den Braber<sup>1,3,4</sup>, Pieter Jelle Visser<sup>3,5</sup>, Eco J.C.N. de Geus<sup>1,2,4</sup>, Ko Willems van Dijk<sup>6,7,8</sup>, Aswin Verhoeven<sup>9</sup>, H. Eka Suchiman<sup>10</sup>, Marian Beekman<sup>10</sup>, P. Eline Slagboom<sup>10</sup>, Cornelia M. van Duijn<sup>11</sup>, BBMRI Metabolomics Consortium, Amy C. Harms<sup>12</sup>, Thomas Hankemeier<sup>12</sup>, Meike Bartels<sup>1,2,4</sup>, Michel G. Nivard<sup>1,2,4,52\*</sup> & Dorret I. Boomsma<sup>1,2,4,52\*</sup>

Metabolomics examines the small molecules involved in cellular metabolism. Approximately 50% of total phenotypic differences in metabolite levels is due to genetic variance, but heritability estimates differ across metabolite classes. We perform a review of all genome-wide association and (exome-) sequencing studies published between November 2008 and October 2018, and identify >800 class-specific metabolite loci associated with metabolite levels. In a twin-family cohort ( $N = 5117$ ), these metabolite loci are leveraged to simultaneously estimate total heritability ( $h^2_{\text{total}}$ ), and the proportion of heritability captured by known metabolite loci ( $h^2_{\text{Metabolite-hits}}$ ) for 309 lipids and 52 organic acids. Our study reveals significant differences in  $h^2_{\text{Metabolite-hits}}$  among different classes of lipids and organic acids. Furthermore, phosphatidylcholines with a high degree of unsaturation have higher  $h^2_{\text{Metabolite-hits}}$  estimates than phosphatidylcholines with low degrees of unsaturation. This study highlights the importance of common genetic variants for metabolite levels, and elucidates the genetic architecture of metabolite classes.

<sup>1</sup>Department of Biological Psychology, Vrije Universiteit Amsterdam, Amsterdam, The Netherlands. <sup>2</sup>Amsterdam Public Health Research Institute, Amsterdam, The Netherlands. <sup>3</sup>Alzheimer Center Amsterdam, Department of Neurology, VU Amsterdam, Amsterdam UMC, Amsterdam, The Netherlands. <sup>4</sup>Amsterdam Neuroscience, Amsterdam, The Netherlands. <sup>5</sup>Department of Psychiatry and Neuropsychology, School of Mental Health and Neuroscience, Alzheimer Center Limburg, Maastricht University, Maastricht, The Netherlands. <sup>6</sup>Eindhoven Laboratory for Experimental Vascular Medicine, Leiden University Medical Center, Leiden, The Netherlands. <sup>7</sup>Department of Human Genetics, Leiden University Medical Center, Leiden, The Netherlands. <sup>8</sup>Department of Internal Medicine, Division of Endocrinology, Leiden University Medical Center, Leiden, The Netherlands. <sup>9</sup>Center for Proteomics and Metabolomics, Leiden University Medical Center, Leiden, The Netherlands. <sup>10</sup>Department of Biomedical Data Sciences, Section of Molecular Epidemiology, Leiden University Medical Center, Leiden, The Netherlands. <sup>11</sup>Department of Epidemiology, Erasmus Medical Center, Rotterdam, The Netherlands. <sup>12</sup>Division of Analytical Biosciences, Leiden Academic Center for Drug Research, Leiden University and The Netherlands Metabolomics Centre, Leiden, The Netherlands. <sup>51</sup>Department of Psychiatry, Amsterdam UMC, University of Amsterdam, Amsterdam, The Netherlands. <sup>52</sup>These authors contributed equally: Michel G. Nivard, Dorret I. Boomsma. A full list of consortium members appears at the end of the paper. \*email: [f.a.hagenbeek@vu.nl](mailto:f.a.hagenbeek@vu.nl); [m.g.nivard@vu.nl](mailto:m.g.nivard@vu.nl); [di.boomsma@vu.nl](mailto:di.boomsma@vu.nl)

The metabolome is defined as the collection of metabolites, i.e., small molecules involved in cellular metabolism, which are produced in cells<sup>1</sup> and can be categorized into many classes<sup>2</sup>. The overall aim of the field of metabolomics is to provide a holistic overview of the metabolome<sup>1</sup>, and its role in biological mechanisms and metabolic disturbances in diseases. Elucidating this role may offer new therapeutic targets or new biomarkers for disease diagnosis<sup>3</sup>. Variation in metabolite levels can arise due to gender<sup>4</sup>, and age<sup>5</sup>, as well as physiologic effects, behavior, and lifestyle factors, such as diet<sup>6</sup>. Genetic differences may be a source of direct variation in metabolomics profiles, or an indirect source of variation through genetic influences on physiology, behavior, and (or) lifestyle.

Genome- and metabolome-wide analysis of common genetic variants in human metabolism have successfully identified genetically influenced metabolites<sup>7</sup>. In 2008, the first genome-wide association study (GWAS;  $N=284$  participants) identified four genetic variants associated with metabolite levels<sup>8</sup>. Thereafter, GWAS with increasing sample sizes, and in diverse populations, identified hundreds of single nucleotide polymorphism (SNP) associations with metabolites from a wide range of metabolite classes<sup>7</sup>. Additional metabolite loci have been identified by leveraging low-frequency and rare-variant analyses using (exome-) sequencing. We conducted a comprehensive review of all quantitative trait loci (QTL) discovery for metabolites and supply the complete reference list in Supplementary Table 1.

Twin and family studies have established that the heritability ( $h^2$ ; proportion of phenotypic variance due to genetic factors) of metabolite levels is 50% on average, with a range from  $h^2 = 0\%$  to  $h^2 = 80\%$ <sup>6,9–16</sup>. Several studies reported differences in heritability estimates among different classes of lipid species<sup>13,15</sup> or lipoprotein subclasses<sup>14</sup>. For example, Rhee et al.<sup>12</sup> reported higher heritability estimates for amino acids than for lipids. Essential amino acids, which cannot be synthesized by an organism *de novo*<sup>17</sup>, had lower heritability than nonessential amino acids<sup>12</sup>, that are synthesized within the body<sup>17</sup>. Several techniques are available to estimate the contribution of measured SNPs to trait heritability<sup>18</sup>, and, given SNP data in family members, to simultaneously estimate SNP-associated ( $h^2_{\text{SNP}}$ ) and pedigree-associated genetic variance ( $h^2_{\text{ped}}$ )<sup>19</sup>. Together the SNP- and pedigree-associated genetic effects account for the narrow-sense heritability. However, when including data of family members, the variance explained by genetic effects ( $h^2_{\text{total}}$ ) may be biased upwards by shared environmental factors and/or nonadditive genetic effects<sup>19,20</sup>.

An improved understanding of the genetic background of the metabolome will benefit our understanding of the etiology of diseases and traits, such as cardiometabolic diseases<sup>21</sup>, migraine<sup>22</sup>, psychiatric disorders<sup>23</sup>, and cognition<sup>24</sup>. Here, we aim to further our understanding of the contribution of genetic factors to variation in fasting blood metabolic measures (henceforth referred to as metabolites for brevity) by the analysis of data from multiple metabolomics platforms in a large cohort of twins and family members ( $N=5117$ ). Specifically, we aim to estimate the total genetic variance of metabolite levels ( $h^2_{\text{total}}$ ), and to elucidate the contribution to metabolite levels of known metabolite class-specific and metabolite class-unspecific loci ( $h^2_{\text{Metabolite-hits}}$ ), on the basis of the results of a decade of GWA and (exome-) sequencing studies. To this end, we characterize all metabolite-SNP associations published between November 2008 and October 2018 by metabolite classification, and used linear-mixed models to estimate the  $h^2_{\text{total}}$ ,  $h^2_{\text{SNP}}$ , and  $h^2_{\text{Metabolite-hits}}$  simultaneously for 369 metabolites. In these models, the  $h^2_{\text{Metabolite-hits}}$  consists of two variance components, a component attributable to metabolite loci associated with metabolites of a specific superclass ( $h^2_{\text{Class-hits}}$ ) and a component attributable other metabolite loci ( $h^2_{\text{Notclass-hits}}$ ).

The median  $h^2_{\text{total}}$  for lipids is 0.47 and for organic acids 0.40, and the median lipid  $h^2_{\text{Metabolite-hits}}$  is 0.06 and 0.01 for organic acids, with most of the  $h^2_{\text{Metabolite-hits}}$  attributable to  $h^2_{\text{Class-hits}}$ . We further expand on the current knowledge of the genetic etiology of metabolite classes by employing mixed-effect meta-regression models to test differences in heritability estimates among metabolite classes and among lipid species. Although estimates of  $h^2_{\text{total}}$  do not differ significantly among metabolite classes, significant differences were observed among lipid and organic classes for  $h^2_{\text{Metabolite-hits}}$  and  $h^2_{\text{Class-hits}}$ .

Intriguingly, phosphatidylcholines<sup>11</sup> and triglycerides (TGs)<sup>16</sup> show increasing heritability with increasing number of carbon atoms and/or double bonds in their fatty acyl side chains. Draisma et al.<sup>11</sup> speculated this might be attributable to differences in the number of metabolic conversion rounds for phosphatidylcholines or TGs with a variable number of carbon atoms. To distinguish between the effects of the number of carbon atoms or number of double bonds in the fatty acyl side chains of phosphatidylcholines and TGs, we conduct additional univariate follow-up analyses. Our results indicate higher  $h^2_{\text{Metabolite-hits}}$  estimates for more complex phosphatidylcholines (i.e., with larger number of carbon atoms and/or double bonds). Univariate follow-up suggests this could be attributed to the number of double bonds in phosphatidylcholines (e.g., degree of unsaturation).

## Results

**Metabolite classification.** In the period of November 2008 to October 2018, 40 GWA and (exome-) sequencing studies identified 242,580 metabolite-SNP or metabolite ratio-SNP associations (see Supplementary Table 1). All 242,580 associations may be found in Supplementary Data 1, which lists the significant SNP-metabolite associations by study. These associations, included 1804 unique metabolites or ratios and 49,231 unique SNPs (43,830 after converting all SNPs to NCBI build 37; Supplementary Data 1). The human metabolome database (HMDB)<sup>2</sup> identifiers of each metabolite were retrieved in order to extract information concerning the metabolite's hydrophobicity and chemical classification (see Methods). Excluding the ratios and unidentified metabolites, we classified 953 metabolites into 12 super classes (Table 1), 43 classes, or 77 subclasses based on the HMDB classification (Supplementary Data 1). The majority of the metabolites were classified into the super classes lipids or organic acids. The lipids

**Table 1 Overview of the number of unique metabolites per super class.**

Super class	Number of unique metabolites
Lipids and lipid-like molecules (e.g., lipids)	662
Organic acids and derivatives (e.g., organic acids)	182
Organoheterocyclic compounds	45
Organic oxygen compounds	19
Nucleosides, nucleotides, and analogues	12
Benzenoids	12
Organic nitrogen compounds	11
Phenylpropanoids and polyketides	4
Proteins	3
Organic compounds	1
Trichlorophenols	1
Organoxygen compounds	1

For each Human Metabolome<sup>2</sup> super class the number of unique metabolites, for which significant SNP-metabolite associations have been published, is provided. See Supplementary Data 1 for an overview of the exact metabolites classified per super class, class, and subclass, as well as the SNPs associated with each metabolite

could be subdivided into 8 classes, with 1 to 95,795 metabolite-SNP associations per class (mean = 17,589; SD = 32,553), and in 32 subclasses, with the number of subclass metabolites-SNP associations ranging from 1 to 40,440 (mean = 4673; SD = 9124). The organic acids and derivatives were divided in 9 classes, with the number of metabolite-SNP associations ranging from 1 to 26,832 (mean = 3374; SD = 8832). The organic acids and derivatives were also divided into 17 organic acid subclasses, with the number of subclass metabolite-SNP associations ranging from 1 to 26,448 (mean = 1786; SD = 6371; Supplementary Data 1). Across all four platforms 427 metabolites were assessed. After excluding the ratios (17) and the metabolites of super classes not included in the curated metabolite-SNP association list (8), data were available for 402 metabolites. The full list of metabolites, with their classifications and the quartile values of the untransformed levels, is included in Supplementary Table 1. The 402 metabolites were classified as 336 lipids, 53 organic acids, 9 organic oxygen compounds, 3 proteins and one organic nitrogen compound. These super classes consisted of 12 classes (Supplementary Table 2). In this paper we mainly focus on the first two super classes. After quality control (QC), 369 metabolites from these two super classes were retained for analysis.

**Characterization of the heritable influences on metabolites.**

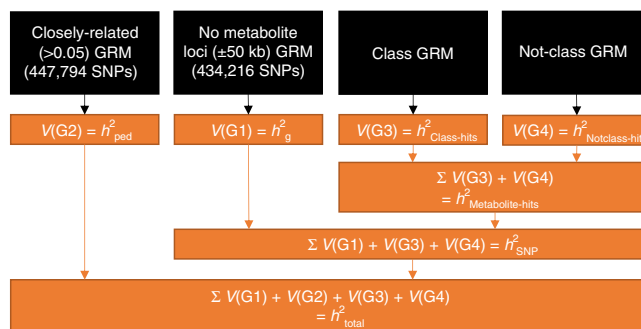
Data of 5117 participants were available from the following four metabolomics platforms: the Nightingale Health proton nuclear magnetic resonance (<sup>1</sup>H-NMR) platform, a ultra performance liquid chromatography mass spectrometry (UPLC-MS) lipidomics platform, the Leiden <sup>1</sup>H-NMR platform, and the Biocrates Absolute-IDQ™ p150 platform. The participants were registered with the Netherlands Twin Register (NTR)<sup>25</sup> and were clustered in 2445 nuclear families. Metabolomics and SNP data were available for all participants. Background and demographic characteristics for the sample can be found in Table 2.

We aimed to assess the variance explained by previously identified metabolite GWA and (exome-) sequencing genetic variants in our (independent) sample. Clearly, our results are conditional on the power of past the studies, as the list of metabolite genetic variants is based on previous GWA and (exome-) sequencing studies, which vary in power. We present the sample size of each past study in Supplementary Table 1, and the sample size per metabolite-SNP association in Supplementary Data 1.

Linear-mixed models including all loci for genetic variants associated with metabolites in a single genetic relatedness matrix (GRM) will contain SNPs that are associated with some metabolites, but not with others, or include many SNPs that are not associated with a given metabolite. We therefore created two GRMs for the loci associated with metabolite hits (see Methods): one class-specific and one nonclass specific (i.e., GRMs including metabolite loci for all metabolites, except for the target metabolite class). We explored models for the 12

class-specific and the corresponding not-class specific GRMs (Supplementary Note 2). These models displayed high degrees of non-convergence (37.9% total), with models including small class-specific GRMs displaying more non-convergence (Supplementary Table 2). Therefore, the results in the remainder of this paper were based on the metabolite super classes, i.e., lipids and organic acids.

For the 369 lipids and organic acids, we carried out unconstrained four-variance component analyses (Fig. 1). In genome-wide complex trait analysis (GCTA)<sup>18</sup> we specified a model in which we partition the metabolite variation into SNP-associated ( $h^2_{\text{SNP}}$ ), pedigree-associated ( $h^2_{\text{ped}}$ ), class-specific metabolite-loci-associated ( $h^2_{\text{class-hits}}$ ), and not-class metabolite-loci-associated ( $h^2_{\text{notclass-hits}}$ ) genetic variation (Fig. 1). We report the total heritability ( $h^2_{\text{total}}$ ), the proportion attributable to metabolite superclass-specific loci ( $h^2_{\text{Class-hits}}$ ), the proportion of variance attributable to non-superclass metabolite loci ( $h^2_{\text{Notclass-hits}}$ ) and the contribution of known metabolite loci to metabolite levels ( $h^2_{\text{Metabolite-hits}}$ ). The analyses were performed separately for lipids and organic acids, with class-specific and corresponding nonclass GRMs (created using the LDK program<sup>26,27</sup>) in both sets of analyses. The lipid analyses employed a class-specific GRM of 479 lipid loci and a corresponding nonclass GRM of 596 loci (Supplementary Fig. 1). The organic acid analyses included a class-specific GRM of 397 loci and a nonclass GRM of 683 loci (Supplementary Fig. 1). Before the analyses, the metabolite data were normalized (log-normal or inverse rank; see Methods). All models included age at blood draw, sex, the first ten principal



**Fig. 1 Overview of the four-variance component models.** Overview of the SNP-filtering and GRM construction can be found in Supplementary Fig. 1 and is explained in details in the Methods. This figure describes which GRMs (black boxes) are used to calculate which variance components (orange boxes) by drawing black arrows from the GRMs to the variance components. The variance components give rise to the four different heritability estimates:  $h^2_{\text{ped}}$ ,  $h^2_{\text{g}}$ ,  $h^2_{\text{Class-hits}}$ , and  $h^2_{\text{Notclass-hits}}$  (see Methods). The orange arrows indicate how the various variance components are summed to obtain estimates for  $h^2_{\text{metabolite-hits}}$ ,  $h^2_{\text{SNP}}$ , and  $h^2_{\text{total}}$  (see Methods).

**Table 2 Participant characteristics per metabolomics platform.**

Metabolomics platform	N	N families	Age <sup>a</sup> (mean ± SD)	Female (%)	Twins (%)	BMI (mean ± SD)	Cholesterol <sup>b</sup> (mean ± SD)	LDL <sup>b</sup> (mean ± SD)	HDL <sup>b</sup> (mean ± SD)
All participants	5117	2445	42.1 ± 14.2	62.8	63.4	24.8 ± 4.1	4.9 ± 1.2	3.0 ± 1.0	1.7 ± 1.0
Nightingale Health <sup>1</sup> H-NMR	4227	2179	40.7 ± 13.7	67.3	69.7	24.6 ± 4.0	4.9 ± 1.2	3.0 ± 1.0	1.7 ± 1.0
UPLC-MS lipidomics	2324	1251	39.0 ± 12.9	66.6	89.2	24.4 ± 4.1	5.0 ± 1.0	3.0 ± 0.9	1.4 ± 0.4
Leiden <sup>1</sup> H-NMR	2324	1323	37.6 ± 12.5	67.0	89.0	24.2 ± 4.1	4.6 ± 1.3	2.7 ± 1.0	2.0 ± 1.4
Biocrates	1448	946	45.7 ± 15.3	43.8	39.6	25.2 ± 3.9	4.6 ± 1.5	2.8 ± 1.1	2.3 ± 1.7

This table gives an overview of the number of individuals (N) per platform, specifies the number of families these individuals belong to and the percentage of females and twins in each dataset. In addition, for each platform the mean and standard deviation (SD) of the age at blood draw in years, the body mass index (BMI), the cholesterol level in mmol/l, the low-density lipoprotein cholesterol (LDL) levels in mmol/l, and the high-density lipoprotein cholesterol (HDL) levels in mmol/l are given. All participant characteristics are given after preprocessing, which was done separately for each metabolomics platform (see Methods)

<sup>a</sup>Age at blood draw in years  
<sup>b</sup>Levels in mmol/l

**Table 3 Summary of the heritability estimates of the four-variance component models.**

		Mean	Median	Range
Lipids and lipid-like molecules	$h^2_{\text{total}}$ estimate	0.47	0.47	(0.11–0.66)
	$h^2_{\text{total}}$ S.e.	0.04	0.03	(0.02–0.07)
	$h^2_{\text{Metabolite-hits}}$ estimate	0.06	0.06	(–0.05–0.16)
	$h^2_{\text{Metabolite-hits}}$ S.e.	0.03	0.03	(0.01–0.04)
	$h^2_{\text{Class-hits}}$ estimate	0.06	0.06	(–0.02–0.16)
	$h^2_{\text{Class-hits}}$ S.e.	0.02	0.02	(0.01–0.03)
	$h^2_{\text{Notclass-hits}}$ estimate	0.00	0.01	(–0.06–0.12)
	$h^2_{\text{Notclass-hits}}$ S.e.	0.02	0.02	(0.01–0.03)
Organic acids and derivatives	$h^2_{\text{total}}$ estimate	0.41	0.40	(0.14–0.72)
	$h^2_{\text{total}}$ S.e.	0.04	0.03	(0.02–0.07)
	$h^2_{\text{Metabolite-hits}}$ estimate	0.01	0.02	(–0.08–0.11)
	$h^2_{\text{Metabolite-hits}}$ S.e.	0.02	0.02	(0.01–0.04)
	$h^2_{\text{Class-hits}}$ estimate	0.01	0.01	(–0.04–0.14)
	$h^2_{\text{Class-hits}}$ S.e.	0.02	0.02	(0.01–0.03)
	$h^2_{\text{Notclass-hits}}$ estimate	0.00	0.00	(–0.06–0.05)
	$h^2_{\text{Notclass-hits}}$ S.e.	0.02	0.02	(0.01–0.03)

The mean, median, and range of the total heritability ( $h^2_{\text{total}}$ ), heritability based on the 479 significant metabolite loci for the 309 lipids or the 397 significant metabolite loci for the 52 organic acids ( $h^2_{\text{Class-hits}}$ ), the 596–683 significant metabolite loci not belonging to these classes ( $h^2_{\text{Notclass-hits}}$ ) and the total heritability explained by metabolite loci (e.g., sum of  $h^2_{\text{Class-hits}}$  and  $h^2_{\text{Notclass-hits}}$ :  $h^2_{\text{Metabolite-hits}}$ ), as well as their standard errors (s.e.'s), are depicted for all 361 successfully analyzed metabolites as included on all platforms. Supplementary Data 2 denotes which metabolites belong to each class and Supplementary Data 3 provides the estimates for each of the individual metabolites

components (PCs) from SNP genotype data, genotyping chip, and metabolomics measurement batch as covariates.

Supplementary Data 3 includes the estimates from the four-variance genetic component models for all 369 metabolites. The genomic relatedness matrix residual maximum likelihood (GREML) algorithm converged for 361 (97.8%) of the 53 organic acids and 316 lipids (Supplementary Table 3). Non-convergence of the GREML algorithm was observed for 6 metabolites (1.6%). The analyses of 2 metabolites (0.5%) were not completed due to non-invertible variance-covariance matrices. The estimates for  $h^2_{\text{total}}$  of the 309 lipids ranged from 0.11 to 0.66 (mean = 0.47; mean s.e. = 0.04). The estimates for  $h^2_{\text{Metabolite-hits}}$  ranged from –0.05 to 0.16 (mean = 0.06; mean s.e. = 0.03; Table 3). The 52 organic acids had  $h^2_{\text{total}}$  estimates ranging from 0.14 to 0.72 (mean = 0.41; mean s.e. = 0.04). The estimates for  $h^2_{\text{Metabolite-hits}}$  ranged from –0.08 to 0.11 (mean = 0.01; mean s.e. = 0.02; Table 3). On average, for both lipids and organic acids the  $h^2_{\text{class}}$  was higher than the  $h^2_{\text{Notclass}}$ , with  $h^2_{\text{Class-hits}}$  ranging from –0.02 to 0.16 (0.06; mean s.e. = 0.02) for lipids and from –0.04 to 0.14 for organic acids (mean = 0.01; mean s.e. = 0.02). For both lipids and organic acids  $h^2_{\text{Notclass-hits}}$  was zero (mean s.e. = 0.02), ranging from –0.06 to 0.12 for lipids and from –0.06 to 0.05 for organic acids (Table 3).

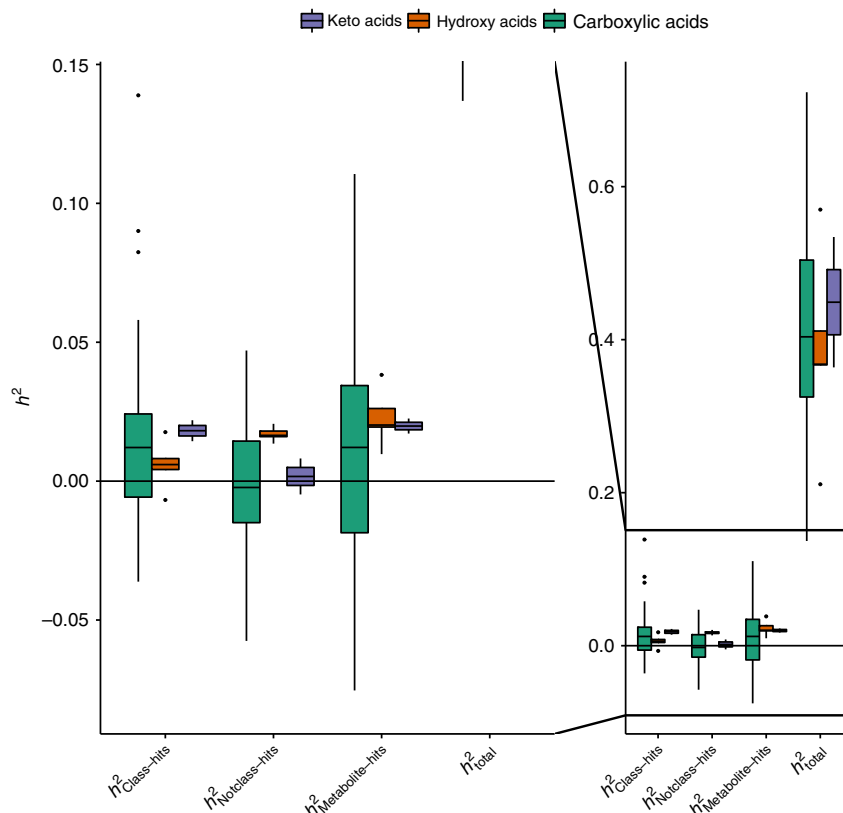
Including multiple metabolomics platforms allowed for a comparison of metabolites as measured on multiple platforms. An earlier study showed that 29 out of 43 metabolites present on two platforms to exhibit moderate heritability on both platforms<sup>28</sup>. In the current study, 61 metabolites were measured on multiple platforms (phenotypic correlations provided in Supplementary Data 4), with moderate  $h^2_{\text{total}}$  on each of the platforms and on average a positive correlation of 0.36 between the  $h^2_{\text{total}}$  of the same metabolite assessed on different platforms (Supplementary Data 4).

**Differential heritability among metabolite classes.** Figure 2 shows variation in median heritability among the following classes of organic acids: keto acids, hydroxy acids, and carboxylic acids (see Supplementary Data 2 for metabolites per class). Keto acids, followed by carboxylic acids, had the highest median  $h^2_{\text{total}}$  and  $h^2_{\text{Class-hits}}$  estimates (Fig. 2). While hydroxy acids had the highest median  $h^2_{\text{Notclass-hits}}$  and  $h^2_{\text{Metabolite-hits}}$  estimates, the lowest median  $h^2_{\text{total}}$  and  $h^2_{\text{Class-hits}}$  estimates were observed for these metabolites

(Fig. 2). To investigate whether heritability differs significantly among classes of organic acids, we applied multivariate mixed-effect meta-regression, corrected for metabolite platform effects (see Methods). The multivariate mixed-effect meta-regression models showed that  $h^2_{\text{total}}$  and  $h^2_{\text{Class-hits}}$  for the organic acid classes did not differ significantly. However, significant differences among the organic acid classes were observed with multivariate mixed-effect meta-regression models with respect to the  $h^2_{\text{Metabolite-hits}}$  estimates ( $F(4, 47) = 3.44$ , false discovery rate (FDR)-adjusted  $p$  value = 0.03), and the  $h^2_{\text{Notclass-hits}}$  estimates ( $F(4, 47) = 19.95$ , FDR-adjusted  $p$  value =  $1.25 \times 10^{-08}$ , Supplementary Data 5).

The multivariate mixed-effect meta-regressions were also applied to assess the significance of heritability differences among essential and non-essential amino acids (subdivision of carboxylic acids; see Supplementary Table 4) and among lipid classes (see Supplementary Data 2 for metabolites per lipid class). The meta-regression analyses revealed no significant mean differences among essential and non-essential amino acids (Table 4; Supplementary Data 6). Small but significant mean heritability differences were observed with multivariate mixed-effect meta-regression models among the different classes of lipids (Fig. 3). For lipid classes the  $h^2_{\text{Metabolite-hits}}$  estimates differed significantly ( $F(8, 300) = 8.47$ ; FDR-adjusted  $p$  value = 0.004; Supplementary Data 5).

Finally, we explored whether heritability of phosphatidylcholines and TGs increases with a larger number of carbon atoms and/or double bonds in their fatty acyl side chains. To this end we employed both uni- and multivariate mixed-effect meta-regression models separately for the TGs, diacyl phosphatidylcholines (PCaa) and acyl-alkyl phosphatidylcholines (PCae; see Methods). The platform specific heritability estimates for each of these lipid species are depicted in Supplementary Fig. 2. Multivariate mixed-effect meta-regression models showed that variation in the number of carbon atoms and double bonds was significantly associated with  $h^2_{\text{Metabolite-hits}}$  estimates for PCaa's ( $F(3, 52) = 7.05$ ; FDR-adjusted  $p$  value = 0.009) and PCae's ( $F(3, 45) = 3.41$ ; FDR-adjusted  $p$  value = 0.05; Supplementary Data 5). Phosphatidylcholines with a larger number of carbon atoms showed lower heritability estimates and phosphatidylcholines with a larger number of double bonds had higher heritability estimates (Supplementary Data 5). The differences among the phosphatidylcholines with a variable number of carbon atoms and/or double



**Fig. 2 Heritability of all 52 carboxylic acids by class.** Box- and dotplots of the  $h^2_{total}$ , and  $h^2_{Metabolite-hits}$  for all 52 successfully analyzed carboxylic acids and derivatives across all metabolomics platforms by class. The left-hand side of the figure is a close-up of the  $-0.08$  to  $0.15$  part of the heritability range, focusing on the  $h^2_{Class-hits}$  and  $h^2_{Notclass-hits}$  estimates. The boxes denote the 25th and 75th percentile (bottom and top of box), and median value (horizontal band inside box). The whiskers indicate the values observed within up to 1.5 times the interquartile range above and below the box. The purple, orange and green boxes denote the keto acid, hydroxyl acid and carboxylic acid classes, respectively. Supplementary Data 3 provides the estimates for each of the individual metabolites.

**Table 4 Summary of the heritability estimates for the essential and nonessential amino acids.**

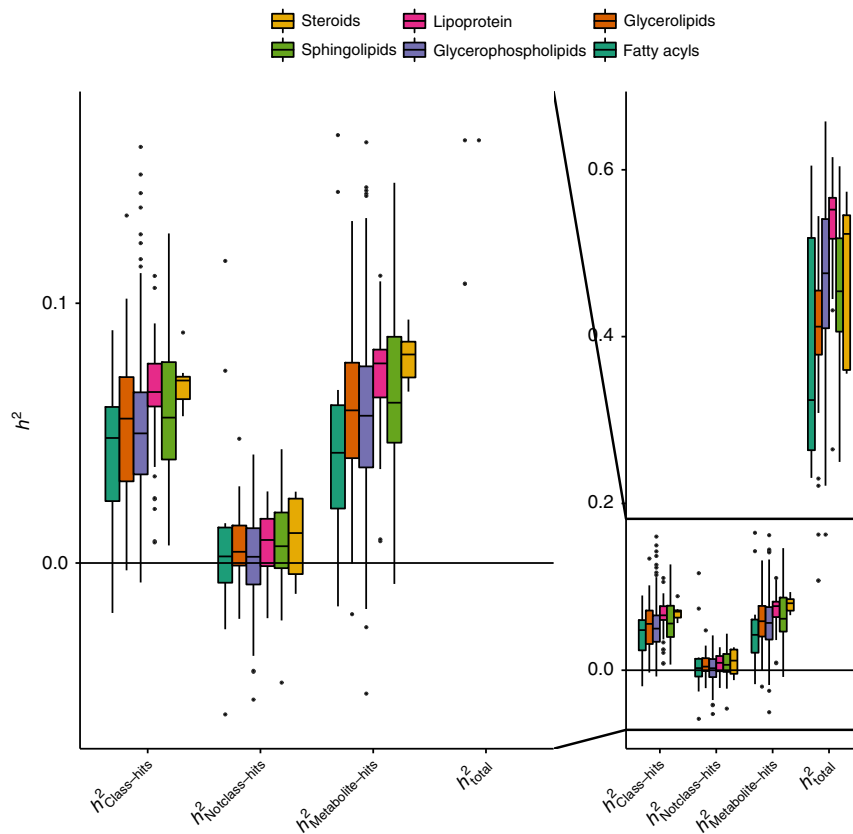
		Mean	Median	Range
Essential amino acids	$h^2_{total}$ estimate	0.42	0.40	(0.23-0.64)
	$h^2_{total}$ s.e.	0.04	0.03	(0.02-0.07)
	$h^2_{Metabolite-hits}$ estimate	0.00	0.00	(-0.05-0.05)
	$h^2_{Metabolite-hits}$ s.e.	0.02	0.02	(0.01-0.03)
	$h^2_{Class-hits}$ estimate	0.01	0.00	(-0.03-0.05)
	$h^2_{Class-hits}$ s.e.	0.02	0.02	(0.01-0.02)
	$h^2_{Notclass-hits}$ estimate	-0.01	-0.01	(-0.06-0.04)
	$h^2_{Notclass-hits}$ s.e.	0.02	0.02	(0.01-0.03)
Non-essential amino acids	$h^2_{total}$ estimate	0.39	0.39	(0.22-0.69)
	$h^2_{total}$ s.e.	0.04	0.04	(0.03-0.07)
	$h^2_{Metabolite-hits}$ estimate	0.02	0.01	(-0.07-0.11)
	$h^2_{Metabolite-hits}$ s.e.	0.03	0.03	(0.01-0.04)
	$h^2_{Class-hits}$ estimate	0.03	0.01	(-0.03-0.14)
	$h^2_{Class-hits}$ s.e.	0.02	0.02	(0.01-0.03)
	$h^2_{Notclass-hits}$ estimate	0.00	0.00	(-0.04-0.03)
	$h^2_{Notclass-hits}$ s.e.	0.02	0.02	(0.01-0.03)

The mean, median, and range of the total heritability ( $h^2_{total}$ ), and heritability based on the 397 significant metabolite loci for the organic acids ( $h^2_{Class-hits}$ ), the 683 significant metabolite loci not belonging to this class ( $h^2_{Notclass-hits}$ ) and the total heritability explained by metabolite loci (e.g., sum of  $h^2_{Class-hits}$  and  $h^2_{Notclass-hits}$ ;  $h^2_{Metabolite-hits}$ ), as well as their standard errors (s.e.'s), are depicted for all 31 successfully analyzed essential (17) and nonessential (14) amino acids as included on all platforms. Supplementary Data 2 denotes which metabolites belong to each class and Supplementary Data 3 provides the estimates for each of the individual metabolites

bonds may have contributed to differential  $h^2_{Class}$  estimates. Univariate models confirmed the results for the number of double bonds in PCaa's and PCae, though they were not significant after correction for multiple testing (Supplementary Data 6).

**Discussion**

We carried out a comprehensive assessment of GWA-metabolomics studies, and created a repository of all studies reporting on associations of SNPs and blood metabolites in



**Fig. 3 Heritability of all 309 lipids by class.** Box- and dotplots of the  $h^2_{\text{total}}$ , and  $h^2_{\text{Metabolite-hits}}$  for all 309 successfully analyzed lipids and lipid-like molecules across all metabolomics platforms by class. The left-hand side of the figure is a close-up of the  $-0.06$  to  $0.17$  part of the heritability range, focusing on the  $h^2_{\text{Class-hits}}$  and  $h^2_{\text{Notclass-hits}}$  estimates. The boxes denote the 25th and 75th percentile (bottom and top of box), and median value (horizontal band inside box). The whiskers indicate the values observed within up to 1.5 times the interquartile range above and below the box. The yellow, pink, orange, light green, purple, and dark green boxes denote the steroids, lipoprotein, glycerolipid, sphingolipid, glycerophospholipid, and fatty acyl classes, respectively. Supplementary Data 3 provides the estimates for each of the individual metabolites.

European ancestry samples. We curated 241,965 genome-wide metabolite associations and we classified the associated metabolites into super classes, classes and subclasses. The complete overview of all blood metabolite-SNP associations is provided in Supplementary Data 1, with the complete list of references in Supplementary Table 1. The information from the repository was used to construct GRMs, which served to identify genetic variance components in the analysis of 369 metabolites. The metabolite data in our study came from a large cohort of twin-families ( $N = 5117$  clustered in 2445 families) measured on four metabolomics platforms. We focused on two metabolite super classes. By mapping all metabolites to the HMDB<sup>2</sup> we were able to classify both the measured metabolites and all previously published metabolites as either lipids or organic acids. In the current study, we sought to elucidate the contribution of known metabolite loci, based on a decade of GWA and (exome-) sequencing studies, to metabolite levels ( $h^2_{\text{Metabolite-hits}}$ ). A unique feature of our study was the ability to disentangle the role of class-specific ( $h^2_{\text{Class-hits}}$ ) and nonclass ( $h^2_{\text{Notclass-hits}}$ ) metabolite loci on heritability differences among metabolite classes and lipid species.

To evaluate differences among metabolite classes and lipid species in the estimates for  $h^2_{\text{total}}$ , we applied multivariate mixed-effect meta-regression models to the estimates of  $h^2_{\text{Metabolite-hits}}$ ,  $h^2_{\text{Class-hits}}$ , and  $h^2_{\text{Notclass-hits}}$ . We observed no significant differences in  $h^2_{\text{total}}$  estimates among the metabolite classes. Consistent with a previous twin-family study<sup>10</sup>, none of the heritability estimates differed significantly among essential and nonessential

amino acids. We observed significant  $h^2_{\text{Metabolite-hits}}$  differences among the different classes of organic acids. Keto acids had significantly lower  $h^2_{\text{Metabolite-hits}}$  estimates as compared with carboxylic acids. Class-specific metabolite loci heritability estimates for fatty acyls, lipoproteins and steroids were significantly higher. Similarly, significant heterogeneity in lipid class heritability, with lower  $h^2_{\text{total}}$  and  $h^2_{\text{SNP}}$  for phospholipids than for sphingolipids or glycerolipids has been reported<sup>13,15,29</sup>. Lastly, we assessed whether heritability increases with added complexity in lipid species<sup>11,16</sup>. We found that this was the case with respect to  $h^2_{\text{Metabolite-hits}}$  estimates in more complex diacyl and acyl-alkyl phosphatidylcholines, but not for more complex TGs. Previous research reported significant higher  $h^2_{\text{SNP}}$  estimates in polyunsaturated fatty acid containing lipids<sup>15</sup>. Furthermore, loci associated with traditional lipid measures explained 2–21% of the variance in lipid levels<sup>15</sup>. Together these results suggest that higher heritability in phosphatidylcholines is driven by a lower number of carbon atoms and higher number of double bonds, e.g., a larger degree of unsaturation.

Evaluating the mean heritability differences among lipids and organic acids, it appears that lipids have higher  $h^2_{\text{total}}$ ,  $h^2_{\text{Class-hits}}$ , and  $h^2_{\text{Metabolite-hits}}$  estimates than organic acids (Table 3). Previous twin-family studies indicates that the heritability difference among lipids and organic acid is rarely investigated<sup>9–12</sup>. This is possibly because most metabolomics platforms focus mainly on either lipids or organic acids. Lipid metabolite classes tend to be very well represented on metabolomics platforms, whereas

organic acids are unrepresented, and as a consequence, the analysis to obtain  $h^2_{\text{Class-hits}}$  and  $h^2_{\text{Metabolite-hits}}$  estimates of the organic acids will be underpowered due to this imbalance.

The current study has several limitations. First, the extent to which our findings generalize to populations of non-European ancestry is unknown. Loci of common human metabolism pathways are most likely to replicate over ethnicities<sup>30</sup>. Second, estimates of the total variance explained may show upward bias when based on data from closely related individuals (e.g., first cousins or closer)<sup>19,20</sup>. This bias is caused by the influence of shared environmental influences, epistatic interactions, or dominance<sup>19,20</sup>. While the results of the current study may suffer of such biases by the inclusion of twins, siblings, and parents, the sample also includes many unrelated individuals which will reduce the possible bias (Supplementary Fig. 3).

Kettunen et al.<sup>31</sup> investigated 217 metabolites of the Nightingale Health <sup>1</sup>H-NMR platform in a classical twin design and reported dominance effects for 6.45% of the metabolites. Tsepilov et al.<sup>32</sup> performed GWA study targeting nonadditive genetic effects and concluded that most genetic effects on metabolite levels and ratios were in fact additive. Together, these studies suggested that the bias due to dominance effects on metabolite levels will be minor.

Relatively few twin-family studies explicitly investigated the role of shared environmental influences on metabolite levels. Overall, shared environmental influences are reported for a small number of metabolites (e.g., 14.3% of all Nightingale Health <sup>1</sup>H-NMR metabolites<sup>31</sup>) and the influence of the shared environment is small-to-moderate (platform and metabolite class-dependent averages range from 0.03 to 0.45<sup>6,13,33–35</sup> with larger estimates deriving from small studies). For studies including parents and offspring, or adult twin and siblings pairs the question arises which effects are captured by the shared environment. Are these the lasting influences of the environment offspring shared with their parents and with each other before they started living independently? Additional research is necessary to elucidate the role of the shared environment on metabolite levels<sup>19</sup>.

Third, standard errors of  $h^2_{\text{SNP}}$  estimates were high. While we have included all  $h^2_{\text{SNP}}$  estimates in the supplements, we stress that the primary goal of our paper was to investigate the contribution of known metabolite loci in an independent sample rather than obtaining the  $h^2_{\text{SNP}}$  estimates for metabolites.

Finally, the estimates for  $h^2_{\text{Metabolite-hits}}$  are based on SNPs of 40 different studies from a decade of GWA and (exome-) sequencing studies. The sample size, and therefore the power, of these studies vary, with some studies conducted with as few as 211 individuals while others included over 24,000 individuals (Supplementary Table 1). For underrepresented metabolites the low power may result in downward biased heritability estimates. However, leveraging information from a decade of research in 40 studies and extracting loci for metabolite classes across multiple studies, the number of such metabolites is not large. New<sup>29,36–38</sup> and future studies will increase the number of variants identified as metabolite loci. The investment in UK Biobank<sup>39</sup> is expected to dramatically increase sample sizes for large-scale genomic investigations of the human metabolome and subsequently the number of metabolite loci.

Mendelian randomization may benefit from the comprehensive overview of metabolite loci that we identified. The identified loci can serve as instruments in metabolome-wide Mendelian randomization studies of complex traits. In addition, our work offers valuable insights into the role of common genetic variants in class specific heritability differences among metabolite classes and lipids species. Further research is required to elucidate the contribution of rare genetic variants to metabolite levels, and differences in the contribution of rare genetic variants among

metabolite classes. A reasonable approach would be to carry out a similar study in a large sample of whole-genome sequencing data. Such an approach, using minor allele frequency (MAF)- and linkage disequilibrium (LD)-stratified GREML analysis<sup>40</sup>, identified additional variance due to rare variants for height and body mass index<sup>41</sup>.

In conclusion, we contributed to our understanding of the genetic architecture of fasting blood metabolite levels, and of differences in the genetic architecture among metabolite classes. Extending the GREML framework with the inclusion of known metabolite loci allowed us to simultaneously estimate  $h^2_{\text{total}}$  and  $h^2_{\text{metabolite-hits}}$  (which consists of  $h^2_{\text{Class-hits}}$  and  $h^2_{\text{Notclass-hits}}$ ) for 361 metabolites. Significant differences in  $h^2_{\text{Metabolite-hits}}$  estimates were observed among different classes of lipids and organic acids and for more complex diacyl and acyl-alkyl phosphatidylcholines. Future studies should address the proportion of metabolite variation influenced by heritable and nonheritable lifestyle factors, as this will facilitate the development of personalized disease prevention and treatment of complex disorders.

## Methods

**Participants.** At the NTR<sup>42</sup> metabolomics data for twins and family members as measured in blood samples were available for 6011 individuals of whom 5667 were genotyped. The blood samples for the four metabolomics experiments described in this study were mainly collected in participants of the NTR biobank project<sup>25,43</sup>. Blood samples were collected after a minimum of two hours of fasting (1.3%), with the majority of the samples collected after overnight fasting (98.7%). Fertile women were bled in their pill-free week or on day 2–4 of their menstrual cycle. For the current paper, we excluded participants who were not of European ancestry, who were on lipid-lowering medication at the time of blood draw, and who failed to adhere to the fasting protocol. The exact number of exclusions per dataset is listed in Supplementary Data 7. After completing the preprocessing of the metabolomics data, the separate subsets (e.g., different collection and measurement waves; see Supplementary Data 7) of each platform were merged into a single per platform dataset, retaining a single (randomly chosen) observation per platform when multiple observations were available. Supplementary Data 8 gives an overview of the overlap in participants among the different platforms, with the overlap among each metabolite that survived QC for all four platforms available in Supplementary Data 9. The final number of participants included in the study was 5117, with platform specific sample size ranging from 1448 to 4227 individuals clustered in 946–2179 families. Characteristics for the individuals can be found in Table 2. Supplementary Fig. 3 depicts the distribution of the relatedness in the sample. Informed consent was obtained from all participants. Projects were approved by the Central Ethics Committee on Research Involving Human Subjects of the VU University Medical Centre, Amsterdam, an Institutional Review Board certified by the U.S. Office of Human Research Protections (IRB number IRB00002991 under Federal-wide Assurance- FWA00017598; IRB/institute codes, NTR 03–180 and EMIF-AD 2014.210).

**Metabolite profiling.** Plasma and serum samples have been profiled on four metabolomics platforms: two proton nuclear magnetic resonance spectroscopy (<sup>1</sup>H-NMR) platforms and two mass spectrometry (MS) platforms. Plasma samples have been analyzed on the Nightingale Health <sup>1</sup>H-NMR platform (Nightingale Health Ltd., Helsinki, Finland), an MS lipidomics platform, and the Leiden <sup>1</sup>H-NMR platform. Serum samples were analyzed with the Biocrates Absolute-IDQ<sup>TM</sup> p150 platform (Biocrates Life Sciences AG, Innsbruck, Austria). Details about each of the metabolomics platforms have been included in Supplementary Note 2.

**Metabolomics data preprocessing.** Preprocessing of the metabolomics data was done separately for each of the platforms and each measurement batch. Metabolites were excluded from analysis when the mean coefficient of variation exceeded 25% and the missing rate exceeded 5%. Metabolite measurements were set to missing if they were below the lower limit of detection or quantification or could be classified as an outlier (five standard deviations greater or smaller than the mean). Metabolite measurements, which were set to missing because they fell below the limit of detection/quantification were imputed with half of the value of this limit, or when this limit was unknown with half of the lowest observed level for this metabolite. All remaining missing values were imputed using multivariate imputation by chained equations (mice)<sup>44</sup>. On average, 9 values were imputed for each metabolite (SD = 12; range: 1–151). Data for each metabolite on both <sup>1</sup>H-NMR platforms were normalized by inverse normal rank transformation<sup>45,46</sup>, while the imputed values of the Biocrates metabolomics platform and the UPLC-MS lipidomics platform were normalized by natural logarithm transformation<sup>11,47</sup>, conform previous normalization strategies applied to the data obtained using these platforms. The complete lists with full names of all detected metabolites that survived

QC and preprocessing for all platforms can be found in Supplementary Data 2, these tables also include the quartile values of the untransformed metabolites.

**Genotyping, imputation, and ancestry outlier detection.** Genotype information was available for 21,001 NTR participants from 6 different genotyping arrays (Affymetrix 6.0 [ $N = 8640$ ], Perlegen-Affymetrix [ $N = 1238$ ], Illumina Human Quad Bead 660 [ $N = 1439$ ], Affymetrix Axiom [ $N = 3144$ ], Illumina GSA [ $N = 5938$ ] and Illumina Omni Express 1 M [ $N = 238$ ]), as well as sequence data from the Netherlands reference genome project GONL (BGI full sequence at  $12 \times (N = 364)$ <sup>48</sup>). For each genotyping array samples were removed if they had a genotype call rate above 90%, gender-mismatch occurred or if heterozygosity (Plink F statistic) fell outside the range of  $-0.10$  to  $0.10$ . SNPs were removed if they were palindromic AT/GC SNPs with a MAF range between  $0.4$  and  $0.5$ , if the MAF was below  $0.01$ , if Hardy Weinberg Equilibrium (HWE) had  $p < 10^{-5}$ , and if the number of Mendelian errors was greater than  $20$  and the genotype call rate was  $< 0.95$ . After QC the six genotyping arrays were aligned to the GONL reference set (V4) and SNPs were removed if the alleles mismatched with this reference panel or the allele frequency different more than  $0.10$  between the genotyping array and this reference set.

The data from the six genotyping chips were subsequently merged into a single dataset (1,781,526 SNPs). Identity-by-descent (IBD) was estimated with PLINK<sup>49</sup> and KING<sup>50</sup> for all individual pairs based on the  $\sim 10.6$  K SNPs in common across the arrays. Next IBD was compared to expected family relations and individuals were removed in the event of a mismatch. Prior to imputation to the GONL reference data<sup>51,52</sup> the duplicate monozygotic pairs ( $N = 3032$ ) or trios ( $N = 7$ ) and NTR GONL samples ( $N = 364$ ) were removed and the data was cross-array phased using MACH-ADMIX<sup>53</sup>. Post-imputation the NTR GONL samples and the duplicated MZ pairs and trios were returned to the dataset. Filtering of the imputed dataset included the removal of SNPs that were significantly associated with a single genotyping chip ( $p < 10^{-5}$ ), had HWE  $p < 10^{-5}$ , the Mendelian error rate  $> \text{mean} + 3 \text{ SD}$ , or imputation quality ( $R^2$ ) below  $0.90$ . The final cross-platform imputed dataset included 1,314,639 SNPs, including 20,792 SNPs on the X-chromosome.

The cross-platform imputed data was aligned with PERL based HRC or 1000G Imputation preparation and checking tool (version 4.2.5; <https://www.well.ox.ac.uk/~wrayner/tools>). The remaining 1,302,481 SNPs were phased with EAGLE<sup>54</sup> for the autosomes, and SHAPEIT<sup>55</sup> for chromosome X and then imputed to 1000 Genomes Phase 3 (1000GP3 version 5)<sup>56</sup> on the Michigan Imputation server using Minimac3 following the standard imputation procedures of the server<sup>57</sup>. PC analysis (PCA) was used to project the first 10 PCs of the 1000 genomes references set population on the NTR cross-platform imputed data using SMARTPCA<sup>58</sup>. Ancestry outliers (non-Dutch ancestry;  $N = 1823$ ) were defined as individuals with PC values outside the European/British population range<sup>59</sup>. After ancestry outlier removal the first 10 PCs were recalculated.

**Curation of metabolite loci.** In October 2018 PubMed and Google Scholar were searched to identify published GWA and (exome-) sequencing studies on metabolomics or fatty acid metabolism in blood samples using <sup>1</sup>H-NMR, mass spectrometry or gas chromatography-based methods. In the period of November 2008 to October 2018 40 GWA or (exome-) sequencing studies on blood metabolomics in European samples were published (Supplementary Table 1). The genome-wide significant ( $p < 5 \times 10^{-8}$ ) metabolite-SNP associations of all studies were extracted, including only those observations for autosomal SNPs and reporting SNP effect sizes and  $p$  values based on the summary statistics excluding NTR samples<sup>46,47</sup>. In the 40 studies, 242,580 metabolite-SNP or metabolite ratio-SNP associations were reported. These associations included 1804 unique metabolites or ratios and 49,231 unique SNPs (Supplementary Data 1). For all metabolites their Human Metabolome Database (HMDB)<sup>2</sup>, PubChem<sup>60</sup>, Chemical Entities of Biological Interest<sup>61</sup> and International Chemical Identifier<sup>62</sup> identifiers were retrieved. Information with regards to the super class, class and subclass of metabolites was extracted from HMDB. If no HMDB identifier was available and categorization information could not be extracted, super class, class and subclass were provided based on expert opinion. Excluding the ratios and unidentified metabolites, 953 metabolites were classified into 12 super classes, 43 classes or 77 subclasses (Supplementary Data 1). Based on the metabolite identifiers we also extracted the  $\log(S)$  value for each metabolite to assess the hydrophobicity of the metabolites. The  $\log(S)$  value represents the log of the partition coefficient between 1-octanol and water, two fluids that hardly mix. The partition coefficient is the ratio of concentrations in water and in octanol when a substance is added to an octanol-water mixture and hence indicates the hydrophobicity of a compound. Thus, we classified a metabolite as hydrophobic if it is more hydrophobic than 1-octanol, and as hydrophilic otherwise (Supplementary Data 1).

The rsIDs or chromosome-base pair positions of the 49,231 unique SNPs were reported by different genome builds or dbSNP maps<sup>63</sup>, therefore we lifted all SNPs to HG19 build 37<sup>64</sup>, after which 43,830 unique SNPs remained (Supplementary Fig. 1; Supplementary Data 1). All biallelic metabolite SNPs were extracted from our 1000GP3 data, which excluded 295 triallelic SNPs, and 4256 SNPs that could not be retrieved from 1000GP3. Next,  $\text{MAF} > 1\%$  (2067 SNPs removed),  $R^2 > 0.70$  (2002 SNPs) and HWE  $p < 10^{-4}$  (72 SNPs) filtering was performed, resulting in 35,138 metabolite SNPs for NTR participants (Supplementary Fig. 1). Next, we

created two super class-specific lists of metabolite loci and two not-superclass lists of metabolite loci. To create a list of loci associated with the 652 unique metabolites classified as lipids and lipid-like molecules (e.g., lipids), we clumped (PLINK version 1.9) all 112,760 lipid-SNP associations using an LD-threshold ( $r^2$ ) of  $0.10$  in a 500 kb radius in 2500 unrelated individuals (Supplementary Fig. 1). Clumping identified 482 lead SNPs, or loci for lipids. An additional 12,169 SNPs were identified as LD-proxies for the lipid-loci (Supplementary Fig. 1). To obtain the not-superclass list of lipid loci the 12,651 lipid loci and proxies were removed from the list of all metabolite-SNP associations and the resulting list was clumped to obtain the 598 non-superclass loci (Supplementary Fig. 1). The same clumping procedure was applied to the 26,352 organic acid-SNP associations, identifying 398 organic acids loci, 10,781 organic acid LD-proxies, and 687 non-superclass loci (Supplementary Fig. 1).

**Construction of genetic relationship matrices.** In total six weighted GRMs were constructed, which were corrected for uneven and long-range LD between the SNPs (LDAK version 4.9<sup>26,27</sup>). In Supplementary Note 3, the use of weighted versus unweighted GRMs is compared using simulations. Two of the GRMs used the cross-platform imputed dataset as backbone and the other four GRMs were based on SNPs extracted from the 1000GP3 imputed data. Before calculating the first GRM, the autosomal SNPs of the cross-platform imputed dataset were filtered on  $\text{MAF} (< 1\%)$  and all lipid and organic acid loci, their LD-proxies and 50 kb surrounding both types of SNPs were removed (see curation of metabolite loci; Supplementary Fig. 1). The LDAK GRM was created after removal of the 50 kb surrounding the lipid and organic acid loci and their LD-proxies (as obtained by the clumping procedure as described above) and included 434,216 SNPs (Supplementary Fig. 1). The  $V(G1)$  variance component in the GREML analyses is based on this GRM (see heritability analyses; Fig. 1). The  $V(G2)$  variance component in the GREML analyses is based on the LDAK GRM including all autosomal SNPs with a MAF greater than  $1\%$  included on the cross-platform imputed dataset (447,794 SNPs), where ancestry outliers were removed, and genome sharing was set to zero for all individual pairs sharing less than  $0.05$  of their genome<sup>19</sup> (Fig. 1). Depending on the metabolite the  $V(G3)$  variance component in the GREML analyses was either based on an LDAK GRM of the 1000GP3 extracted lipid loci (479 SNPs) or the organic acid loci (397 SNPs), as obtained after the clumping procedure as described above (Supplementary Fig. 1; Fig. 1). Finally, depending on the metabolite either the not-lipid LDAK GRM (596 SNPs) or the not-organic acid LDAK GRM (683 SNPs) provided the  $V(G4)$  variance component in the GREML analyses (Supplementary Fig. 1; Fig. 1). The not-class metabolite loci on which the LDAK GRMs were built were obtained by the clumping procedure as described above (Supplementary Fig. 1). Supplementary Data 1 indicates for each listed SNP if it was included in any of the class-specific or not-class LDAK GRMs.

**Heritability analyses.** Mixed linear models<sup>19</sup>, implemented in the GCTA software package (version 1.91.7)<sup>18</sup>, were applied to compare three models including a variable number of covariates. Supplementary Table 5 gives the three different models, full descriptions of the covariates and model comparison have been given in Supplementary Note 4. The most parsimonious model was chosen for further analyses (full results in Supplementary Table 6). This final model included the first ten genetic PCs for the Dutch population, genotyping chip, sex, and age at blood draw as covariates. For metabolites of the Nightingale Health <sup>1</sup>H-NMR and Biocrates platform, measurement batch was included as covariate.

The final four-variance component model, including four GRMs, allows for the estimation of the proportion of variation explained by superclass-specific significant metabolite loci and non-superclass significant metabolite loci. The first two-variance components in the four-variance component model (Fig. 1),  $V(G1)$  and  $V(G2)$  allow for the estimation of the additive genetic variance effects captured by genome-wide SNPs ( $h^2_g$ ) and the additive genetic effects associated with pedigree ( $h^2_{ped}$ )<sup>19,65</sup>, and  $V(G3)$  and  $V(G4)$  capture the additive genetic effect associated with class-specific ( $h^2_{\text{class-hits}}$ ) and not-class ( $h^2_{\text{Notclass-hits}}$ ) metabolite loci. Based on the four-variance component model, three additional heritability estimates can be calculated: the total variance explained by significant metabolite loci ( $h^2_{\text{Metabolite-hits}}$ ) consists of the sum of  $\frac{V(G3)}{V_p}$  and  $\frac{V(G4)}{V_p}$ , where  $V_p$  is the phenotypic variance,  $h^2_{\text{SNP}}$  is defined as the sum of  $\frac{V(G1)}{V_p}$ ,  $\frac{V(G3)}{V_p}$  and  $\frac{V(G4)}{V_p}$ , and the total variance explained ( $h^2_{\text{total}}$ ) is defined as the sum of  $\frac{V(G1)}{V_p}$ ,  $\frac{V(G2)}{V_p}$ ,  $\frac{V(G3)}{V_p}$ , and  $\frac{V(G4)}{V_p}$  (Fig. 1). We note that the total variance explained by genetic factors may also include influences of the shared environment, dominance and epistasis, which may result in upward bias of the  $h^2_{\text{total}}$  estimates<sup>19,20</sup>. This bias is expected to arise by the presence of closely related participants, who may share these effects, in addition to the additive genetic effects. To calculate the standard errors (s.e.'s) for the composite variance estimates, we have randomly sampled 10,000 new variances from the parameter variance-covariance matrices of the  $V(G1)$ ,  $V(G3)$ , and  $V(G4)$  GRMs for each metabolite. Random sampling was performed in R by creating 10,000 multivariate normal distributions (mvrnorm function in MASS package version 7.3-50<sup>66</sup>) based on the original means and variance/covariance matrices. The s.e.'s of the specific ratio of interest were then based on the standard deviation of the ratio of interest across 10,000 samples. The four-variance component models included variance components that were not constrained to be positive, thus



allowing for negative  $h^2_{\text{SNP}}$  and  $h^2_{\text{Metabolite-hits}}$  estimates. All four-variance component models applied the --reml-bendV flag where necessary to invert the variance-covariance matrix  $V$  if  $V$  was not positive definite, which may occur when variance components are negative<sup>67</sup>. Finally, we calculated the log likelihood of a reduced model with either  $V(G3)$ ,  $V(G4)$ , or both dropped from the full model and calculated the LRT and  $p$  value (Supplementary Data 3).

**Mixed-effect meta-regression analyses.** To investigate differences in heritability estimates among metabolites of different classes we applied mixed-effect meta-regression models as implemented in the metafor package (version 2.0-0) in R (version 3.5.1)<sup>68</sup>. Here, we tested for the moderation of heritability estimates by metabolite class and metabolomics platform on all 361 successfully analyzed metabolites. We included a matrix combining the phenotypic correlations (Supplementary Data 10) and the sample overlap (Supplementary Data 9) between the metabolites as random factor to correct for dependence among the metabolites and participants. This matrix includes the sample size of the metabolite on the diagonal, with the off-diagonal computed by  $\frac{N_{12}}{\sqrt{n_1 n_2}} * r$  (Supplementary Data 11), where  $N_{1,2}$  is the sample overlap between the metabolites,  $n_1$  is the sample size of metabolite one,  $n_2$  is the sample size of metabolite two and  $r$  is the phenotypic (Spearman's rho) correlation between the metabolites. In all mixed-effect meta-regression analyses we obtained the robust estimates based on a sandwich-type estimator, clustered by the metabolites included in the models to correct for the sample overlap among the different metabolites<sup>69</sup>. First, we used multivariate mixed-effect meta-regression models to simultaneously estimate the effect of metabolite class and metabolomics platform on the  $h^2_{\text{total}}$ ,  $h^2_{\text{SNP}}$ , and the  $h^2_{\text{Metabolite-hits}}$ , as well as the  $h^2_{\text{Class-hits}}$  and  $h^2_{\text{Notclass-hits}}$  estimates. Subsequently, to separately assess the effect of the number of carbon atoms or double bonds in the fatty acyls chains of phosphatidylcholines and TGs univariate models were fitted, as follow-up. To account for multiple testing the  $p$ -values were adjusted with the with the FDR<sup>70</sup> using the p.adjust function in R. Multiple testing correction was done separately for the univariate and the multivariate models.

**Reporting summary.** Further information on research design is available in the Nature Research Reporting Summary linked to this article.

## Data availability

The curated list of all published metabolite-SNP associations is included in Supplementary Data 1 and is publicly available through the BBMRI—omics atlas (<http://bbmri.researchlumc.nl/atlas/#data>). All information on the metabolites in this study are in Supplementary Data 2; with full summary statistics for the four-variance component models included in Supplementary Data 3. The Nightingale Health metabolomics data may be requested through BBMRI-NL (<https://www.bbmri.nl/Omics-metabolomics>). All (other) data may be accessed, upon approval of the data access committee, through the NTR ([ntr.fgb@vu.nl](mailto:ntr.fgb@vu.nl)).

Received: 22 August 2018; Accepted: 25 November 2019;

Published online: 07 January 2020

## References

- Patti, G. J., Yanes, O. & Siuzdak, G. Innovation: metabolomics: the apogee of the omics trilogy. *Nat. Rev. Mol. Cell Biol.* **13**, 263–269 (2012).
- Wishart, D. S. et al. HMDB 4.0: the human metabolome database for 2018. *Nucleic Acids Res.* **46**, D608–D617 (2018).
- Kuehnbaum, N. L. & Britz-McKibbin, P. New advances in separation science for metabolomics: resolving chemical diversity in a post-genomic era. *Chem. Rev.* **113**, 2437–2468 (2013).
- Mittelstrass, K. et al. Discovery of sexual dimorphisms in metabolic and genetic biomarkers. *PLoS Genet.* **7**, e1002215 (2011).
- Chaleckis, R., Murakami, I., Takada, J., Kondoh, H. & Yanagida, M. Individual variability in human blood metabolites identifies age-related differences. *Proc. Natl Acad. Sci. USA* **113**, 4252–4259 (2016).
- Menni, C. et al. Targeted metabolomics profiles are strongly correlated with nutritional patterns in women. *Metabolomics* **9**, 506–514 (2013).
- Kastenmüller, G., Raffler, J., Gieger, C. & Suhre, K. Genetics of human metabolism: an update. *Hum. Mol. Genet.* **24**, R93–R101 (2015).
- Gieger, C. et al. Genetics meets metabolomics: a genome-wide association study of metabolite profiles in human serum. *PLoS Genet.* **4**, e1000282 (2008).
- Nicholson, G. et al. Human metabolic profiles are stably controlled by genetic and environmental variation. *Mol. Syst. Biol.* **7**, 525 (2011).
- Shah, S. H. et al. High heritability of metabolomic profiles in families burdened with premature cardiovascular disease. *Mol. Syst. Biol.* **5**, 258 (2009).
- Draisma, H. H. M. et al. Familial resemblance for serum metabolite concentrations. *Twin Res. Hum. Genet.* **16**, 948–961 (2013).
- Rhee, E. P. et al. A genome-wide association study of the human metabolome in a community-based cohort. *Cell Metab.* **18**, 130–143 (2013).
- Frahnow, T. et al. Heritability and responses to high fat diet of plasma lipids in a twin study. *Sci. Rep.* **7**, 1–11 (2017).
- Kaess, B. et al. The lipoprotein subfraction profile: heritability and identification of quantitative trait loci. *J. Lipid Res.* **49**, 715–723 (2008).
- Bellis, C. et al. Human plasma lipidome is pleiotropically associated with cardiovascular risk factors and death. *Circ. Cardiovasc. Genet.* **7**, 854–863 (2014).
- Draisma, H. H. M. *Analysis of Metabolomics Data from Twin Families* (Leiden, 2011).
- Reeds, P. J. Dispensable and indispensable amino acids for humans. *J. Nutr.* **130**, 1874S–1876S (2000).
- Yang, J., Lee, S. H., Goddard, M. E. & Visscher, P. M. GCTA: a tool for genome-wide complex trait analysis. *Am. J. Hum. Genet.* **88**, 76–82 (2011).
- Zaitlen, N. et al. Using extended genealogy to estimate components of heritability for 23 quantitative and dichotomous traits. *PLoS Genet.* **9**, e1003520 (2013).
- Zuk, O., Hechter, E., Sunyaev, S. R. & Lander, E. S. The mystery of missing heritability: genetic interactions create phantom heritability. *Proc. Natl Acad. Sci. USA* **109**, 1193–1198 (2012).
- Newgard, C. B. Metabolomics and metabolic diseases: where do we stand? *Cell Metab.* **25**, 43–56 (2017).
- Onderwater, G. L. J. et al. Large-scale plasma metabolome analysis reveals alterations in HDL metabolism in migraine. *Neurology* **92**, e1899–e1911 (2019).
- Nedic Erjavec, G. et al. Short overview on metabolomic approach and redox changes in psychiatric disorders. *Redox Biol.* **14**, 178–186 (2018).
- van der Lee, S. J. et al. Circulating metabolites and general cognitive ability and dementia: Evidence from 11 cohort studies. *Alzheimer's Dement.* **14**, 707–722 (2018).
- Willemsen, G. et al. The Netherlands Twin Register biobank: a resource for genetic epidemiological studies. *Twin Res. Hum. Genet.* **13**, 231–245 (2010).
- Speed, D., Hemani, G., Johnson, M. R. & Balding, D. J. Improved heritability estimation from genome-wide SNPs. *Am. J. Hum. Genet.* **91**, 1011–1021 (2012).
- Speed, D., Cai, N., Johnson, M. R., Nejentsev, S. & Balding, D. J. Reevaluation of SNP heritability in complex human traits. *Nat. Genet.* **49**, 986–992 (2017).
- Yet, I. et al. Genetic influences on metabolite levels: a comparison across metabolomic platforms. *PLoS One* **11**, e0153672 (2016).
- Tabassum, R. et al. Genetic architecture of human plasma lipidome and its link to cardiovascular disease. *Nat. Commun.* **10**, 4329 (2019).
- Yousri, N. A. et al. Whole-exome sequencing identifies common and rare variant metabolic QTLs in a Middle Eastern population. *Nat. Commun.* **9**, 1–13 (2018).
- Kettunen, J. et al. Genome-wide association study identifies multiple loci influencing human serum metabolite levels. *Nat. Genet.* **44**, 269–276 (2012).
- Tsepilov, Y. A. et al. Nonadditive effects of genes in human metabolomics. *Genetics* **200**, 707–718 (2015).
- Tukiaainen, T. et al. Detailed metabolic and genetic characterization reveals new associations for 30 known lipid loci. *Hum. Mol. Genet.* **21**, 1444–1455 (2012).
- Yet, I. et al. Genetic influences on metabolite levels: a comparison across metabolomic platforms. *PLoS One* **11**, e0153672 (2016).
- Tremblay, B. L., Guénard, F., Lamarche, B., Pérusse, L. & Vohl, M. C. Familial resemblances in human plasma metabolites are attributable to both genetic and common environmental effects. *Nutr. Res.* **61**, 22–30 (2019).
- Gallois, A. et al. A comprehensive study of metabolite genetics reveals strong pleiotropy and heterogeneity across time and context. *Nat. Commun.* **10**, 4788 (2019).
- Witteboms, L. B. L. et al. Assessing the causal association of glycine with risk of cardio-metabolic diseases. *Nat. Commun.* **10**, 1–13 (2019).
- Demirkan, A. et al. Genome-wide association study of plasma lipids. Preprint at: <https://doi.org/10.1101/621334> (2019).
- Sudlow, C. et al. UK Biobank: an open access resource for identifying the causes of a wide range of complex diseases of middle and old age. *PLoS Med.* **12**, 1–10 (2015).
- Yang, J. et al. Genetic variance estimation with imputed variants finds negligible missing heritability for human height and body mass index. *Nat. Genet.* **47**, 1114–1120 (2015).
- Wainschtein, P. et al. Recovery of trait heritability from whole genome sequence data. Preprint at: <https://doi.org/10.1101/588020> (2019).
- Boomsma, D. I. et al. Netherlands Twin Register: from twins to twin families. *Twin Res. Hum. Genet.* **9**, 849–857 (2006).
- Willemsen, G. et al. The Adult Netherlands Twin Register: twenty-five years of survey and biological data collection. *Twin Res. Hum. Genet.* **16**, 271–281 (2013).

44. Buuren, S. van & Groothuis-Oudshoorn, K. Mice: multivariate imputation by chained equations in R. *J. Stat. Softw.* **45**, 1–67 (2011).
45. Demirkan, A. et al. Insight in genome-wide association of metabolite quantitative traits by exome sequence analyses. *PLoS Genet.* **11**, e1004835 (2015).
46. Kettunen, J. et al. Genome-wide study for circulating metabolites identifies 62 loci and reveals novel systemic effects of LPA. *Nat. Commun.* **7**, 11122 (2016).
47. Draisma, H. H. M. et al. Genome-wide association study identifies novel genetic variants contributing to variation in blood metabolite levels. *Nat. Commun.* **6**, 7208 (2015).
48. Boomsma, D. I. et al. The Genome of the Netherlands: design, and project goals. *Eur. J. Hum. Genet.* **22**, 221–227 (2014).
49. Purcell, S. et al. PLINK: a tool set for whole-genome association and population-based linkage analyses. *Am. J. Hum. Genet.* **81**, 559–575 (2007).
50. Manichaikul, A. et al. Robust relationship inference in genome-wide association studies. *Bioinformatics* **26**, 2867–2873 (2010).
51. Fedko, I. O. et al. Estimation of Genetic Relationships between Individuals Across Cohorts and Platforms: application to childhood height. *Behav. Genet.* **45**, 514–528 (2015).
52. Deelen, P. et al. Improved imputation quality of low-frequency and rare variants in European samples using the ‘Genome of the Netherlands’. *Eur. J. Hum. Genet.* **22**, 1321–1326 (2014).
53. Liu, E. Y., Li, M., Wang, W. & Li, Y. MaCH-Admix: genotype Imputation for admixed populations. *Genet. Epidemiol.* **37**, 25–37 (2013).
54. Loh, P. R., Palamara, P. F. & Price, A. L. Fast and accurate long-range phasing in a UK Biobank cohort. *Nat. Genet.* **48**, 811–816 (2016).
55. Delaneau, O., Marchini, J. & Zagury, J.-F. A linear complexity phasing method for thousands of genomes. *Nat. Methods* **9**, 179–181 (2012).
56. Auton, A. et al. A global reference for human genetic variation. *Nature* **526**, 68–74 (2015).
57. Das, S. et al. Next-generation genotype imputation service and methods. *Nat. Genet.* **48**, 1284–1287 (2016).
58. Price, A. L. et al. Principal components analysis corrects for stratification in genome-wide association studies. *Nat. Genet.* **38**, 904–909 (2006).
59. Abdellaoui, A. et al. Population structure, migration, and diversifying selection in the Netherlands. *Eur. J. Hum. Genet.* **21**, 1277–1285 (2013).
60. Kim, S. et al. PubChem 2019 update: Improved access to chemical data. *Nucleic Acids Res.* **47**, D1102–D1109 (2019).
61. Hastings, J. et al. ChEBI in 2016: improved services and an expanding collection of metabolites. *Nucleic Acids Res.* **44**, D1214–D1219 (2016).
62. Heller, S. R., McNaught, A., Pletnev, I., Stein, S. & Tchekhovskoi, D. InChI, the IUPAC International Chemical Identifier. *J. Cheminform.* **7**, 1–34 (2015).
63. Sherry, S. T. et al. dbSNP: the NCBI database of genetic variation. *Nucleic Acids Res.* **29**, 308–311 (2001).
64. Haussler, M. et al. The UCSC Genome Browser database: 2019 update. *Nucleic Acids Res.* **47**, D853–D858 (2019).
65. Xia, C. et al. Pedigree- and SNP-associated genetics and recent environment are the major contributors to anthropometric and cardiometabolic trait variation. *PLoS Genet.* **12**, 1–25 (2016).
66. Venables, W. N. & Ripley, B. D. *Modern Applied Statistics With S* (Springer, 2002).
67. Hayes, J. F. & Hill, W. G. Modification of estimates of parameters in the construction of genetic selection indices (‘Bending’). *Biometrics* **37**, 483–493 (1981).
68. Viechtbauer, W. Conducting meta-analyses in R with the metafor Package. *J. Stat. Softw.* **36**, 1–48 (2010).
69. Hedges, L. V., Tipton, E. & Johnson, M. C. Robust variance estimation in meta-regression with dependent effect size estimates. *Res. Synth. Methods* **1**, 39–65 (2010).
70. Benjamini, Y. & Hochberg, Y. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *J. R. Stat. Soc. B* **57**, 289–300 (1995).

## Acknowledgements

We thank all twins and family members for their participation. We thank P.M. Visscher (University of Queensland) for his helpful comments and C.V. Dolan (Vrije Universiteit Amsterdam) for critically reading and commenting on the final version of the paper. Preliminary analyses of this paper were included in a presentation at the 46th Annual

Meeting of the Behavioral Genetics Association (BGA) in June 2016 (abstract in *Behav. Genet.* (2016) 46:785–786), and a presentation at the 49th Annual Meeting of the BGA in June 2019 (abstract forthcoming). This work was performed within the framework of the BBMRI Metabolomics Consortium funded by BBMRI-NL, a research infrastructure financed by the Dutch government (NWO, nos. 184.021.007 and 184.033.111). The European Network of Genomic and Genetic Epidemiology (ENGAGE) contributed to funding to perform the Biocrates Absolute-IDQ™ p150 metabolomics measurements (European Union Seventh Framework Program: FP7/2007–2013, grant number 201413). Analyses were supported by the Netherlands Organization for Scientific Research: Netherlands Twin Registry Repository: researching the interplay between genome and environment (480-15-001/674); the European Union Seventh Framework Program (FP7/2007–2013): ACTION Consortium (Aggression in Children: Unraveling gene–environment interplay to inform Treatment and InterventiON strategies; Grant number 602768). Genotyping was made possible by grants from NWO/SPI 56-464-14192, Genetic Association Information Network (GAIN) of the Foundation for the National Institutes of Health, Rutgers University Cell and DNA Repository (NIMH U24 MH068457-06), the Avera Institute, Sioux Falls (USA) and the National Institutes of Health (NIH R01 HD042157-01A1, MH081802, Grand Opportunity grants 1RC2 MH089951 and 1RC2 MH089995) and European Research Council (ERC-230374). EMIF-AD has received support from the EU/EFPIA Innovative Medicines Initiative Joint Undertaking EMIF grant agreement no. 115372. DIB acknowledges her KNAW Academy Professor Award (PAH/6635). M. Bartels is supported by an ERC consolidator grant (WELL-BEING 771057 PI Bartels). Jv.D. is supported by the NWO-funded X-omics project (184.034.019).

## Author contributions

Nightingale Health metabolomics data: H.E.S., M. Beekman, P.E.S. and C.Mv.D. Leiden <sup>1</sup>H-NMR metabolomics data: K.Wv.D. and A.V. UPLC-MS lipidomics data: A.C.H. and T.H. EMIF-AD data: Ad.B. and P.J.V. Genotype data: J.J.H., A.A., and I.O.F. NTR Biobank data: G.W. and E.J.Cd.G. Metabolomics preprocessing: R.P., H.H.M.D. and F.A.H. Statistical analyses: F.A.H. and M.G.N. Wrote the paper: F.A.H., Jv.D., M. Bartels, M.G.N. and D.I.B. All authors critically read and commented on the paper.

## Competing interests

The authors declare no competing interests.

## Additional information

**Supplementary information** is available for this paper at <https://doi.org/10.1038/s41467-019-13770-6>.

**Correspondence** and requests for materials should be addressed to F.A.H., M.G.N. or D.I.B.

**Peer review information** *Nature Communications* thanks Doug Speed and the other, anonymous, reviewer(s) for their contribution to the peer review of this work. Peer reviewer reports are available.

**Reprints and permission information** is available at <http://www.nature.com/reprints>

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2020

**BBMRI Metabolomics Consortium**

J.J.H. Barkey Wolf<sup>10</sup>, D. Cats<sup>10</sup>, N. Amin<sup>11</sup>, J.W. Beulens<sup>13,14</sup>, J.A. van der Bom<sup>15,16,17,18</sup>, N. Bomer<sup>19</sup>, A. Demirkan<sup>11</sup>, J.A. van Hilten<sup>20</sup>, J.M.T.A. Meessen<sup>21</sup>, M.H. Moed<sup>10</sup>, J. Fu<sup>22,23</sup>, G.L.J. Onderwater<sup>24</sup>, F. Rutters<sup>13</sup>, C. So-Osman<sup>20</sup>, W.M. van der Flier<sup>3,13</sup>, A.A.W.A. van der Heijden<sup>25</sup>, A. van der Spek<sup>11</sup>, F.W. Asselbergs<sup>26</sup>, E. Boersma<sup>27</sup>, P.M. Elders<sup>2,28</sup>, J.M. Geleijnse<sup>29</sup>, M.A. Ikram<sup>11,30,31</sup>, M. Kloppenburg<sup>18,32</sup>, I. Meulenbelt<sup>10</sup>, S.P. Mooijaart<sup>33</sup>, R.G.H.H. Nelissen<sup>34</sup>, M.G. Netea<sup>35,36</sup>, B.W.J.H. Penninx<sup>2,37</sup>, C.D.A. Stehouwer<sup>38,39</sup>, C.E. Teunissen<sup>40</sup>, G.M. Terwindt<sup>24</sup>, L.M. 't Hart<sup>2,10,13,41,42</sup>, A.M.J.M. van den Maagdenberg<sup>43,7</sup>, P. van der Harst<sup>18</sup>, I.C.C. van der Horst<sup>43</sup>, C.J.H. van der Kallen<sup>38,39</sup>, M.M.J. van Greevenbroek<sup>38,39</sup>, W.E. van Spil<sup>44</sup>, C. Wijmenga<sup>22</sup>, A.H. Zwinderman<sup>45</sup>, A. Zhernikova<sup>22</sup>, J.W. Jukema<sup>46</sup>, H. Mei<sup>10,47</sup>, M. Slofstra<sup>22</sup>, M. Swertz<sup>22</sup>, E.B. van den Akker<sup>10,48,49</sup>, J. Deelen<sup>10,50</sup> & M.J.T. Reinders<sup>48,49</sup>

<sup>13</sup>Department of Epidemiology and Biostatistics, Amsterdam University Medical Center, Amsterdam, The Netherlands. <sup>14</sup>Julius Center for Health Sciences and Primary Care, University Medical Center Utrecht, Utrecht, The Netherlands. <sup>15</sup>Centre for Clinical Transfusion Research, Sanquin Research, Leiden, The Netherlands. <sup>16</sup>Jon J van Rood Centre for Clinical Transfusion Research, Leiden University Medical Centre, Leiden, The Netherlands. <sup>17</sup>TIAS, Tilburg University, Tilburg, The Netherlands. <sup>18</sup>Department of Clinical Epidemiology, Leiden University Medical Centre, Leiden, The Netherlands. <sup>19</sup>Department of Cardiology, University Medical Center Groningen, University of Groningen, Groningen, The Netherlands. <sup>20</sup>Center for Clinical Transfusion Research, Sanquin Research, Leiden, The Netherlands. <sup>21</sup>Department of Orthopedics, Leiden University Medical Centre, Leiden, The Netherlands. <sup>22</sup>Department of Genetics, University Medical Center Groningen, University of Groningen, Groningen, The Netherlands. <sup>23</sup>Department of Pediatrics, University Medical Center Groningen, University of Groningen, Groningen, The Netherlands. <sup>24</sup>Department of Neurology, Leiden University Medical Center, Leiden, The Netherlands. <sup>25</sup>Department of General Practice, The EMGO Institute for Health and Care Research, VU University Medical Center, Amsterdam, The Netherlands. <sup>26</sup>Department of Cardiology, Division Heart and Lungs, University Medical Center Utrecht and the Julius Center for Health Sciences and Primary Care, University Medical Center Utrecht, Utrecht, The Netherlands. <sup>27</sup>Thorax Centre, Erasmus Medical Centre, Rotterdam, The Netherlands. <sup>28</sup>Department of General Practice and Elderly Care Medicine, VU University Medical Center, Amsterdam, The Netherlands. <sup>29</sup>Division of Human Nutrition and Health, Wageningen University, Wageningen, The Netherlands. <sup>30</sup>Department of Radiology, Erasmus University Medical Center Rotterdam, Rotterdam, The Netherlands. <sup>31</sup>Department of Neurology, Erasmus University Medical Center Rotterdam, Rotterdam, The Netherlands. <sup>32</sup>Department of Rheumatology, Leiden University Medical Center, Leiden, The Netherlands. <sup>33</sup>Department of Internal Medicine, Division of Gerontology and Geriatrics, Leiden University Medical Centre, Leiden, The Netherlands. <sup>34</sup>Department of Orthopaedics, Leiden University Medical Center, Leiden, The Netherlands. <sup>35</sup>Department of Internal Medicine, Radboud Center for Infectious Diseases, Radboud University Medical Center, Nijmegen, Netherlands. <sup>36</sup>Department for Genomics & Immunoregulation, Life and Medical Sciences Institute (LIMES), University of Bonn, Bonn, Germany. <sup>37</sup>Department of Psychiatry, VU University Medical Center, Amsterdam, The Netherlands. <sup>38</sup>Department of Internal Medicine, Maastricht University Medical Center (MUMC+), Maastricht, The Netherlands. <sup>39</sup>School for Cardiovascular Diseases (CARIM), Maastricht University, Maastricht, The Netherlands. <sup>40</sup>Neurochemistry Laboratory, Clinical Chemistry Department, Amsterdam University Medical Center, Amsterdam Neuroscience, Amsterdam, The Netherlands. <sup>41</sup>Department of Cell and Chemical Biology, Leiden University Medical Center, Leiden, The Netherlands. <sup>42</sup>Department of General practice, Amsterdam University Medical Center, Amsterdam, The Netherlands. <sup>43</sup>Department of Critical Care, University Medical Center Groningen, Groningen, The Netherlands. <sup>44</sup>UMC Utrecht, Department of Rheumatology & Clinical Immunology, Utrecht, The Netherlands. <sup>45</sup>Department of Clinical Epidemiology, Biostatistics, and Bioinformatics, Academic Medical Centre, University of Amsterdam, Amsterdam, The Netherlands. <sup>46</sup>Department of Cardiology, Leiden University Medical Center, Leiden, The Netherlands. <sup>47</sup>Sequencing Analysis Support Core, Leiden University Medical Center, Leiden, The Netherlands. <sup>48</sup>Leiden Computational Biology Center, Leiden University Medical Center, Leiden, The Netherlands. <sup>49</sup>Department of Pattern Recognition and Bioinformatics, Delft University of Technology, Delft, The Netherlands. <sup>50</sup>Max Planck Institute for Biology of Ageing, Cologne, Germany