

1 **Genome-wide analysis of blood lipid metabolites in over 5,000 South Asians**
2 **reveals biological insights at cardiometabolic disease loci**

3
4 Eric L. Harshfield^{1,2}, Eric B. Fauman³, David Stacey¹, Dirk S. Paul^{1,4,5,6,7,8}, Daniel Ziemek⁹,
5 Rachel M. Y. Ong¹, John Danesh^{1,4,5,6,7,8}, Adam S. Butterworth^{1,4,5,6,7,8}, Asif Rasheed¹⁰,
6 Taniya Sattar¹⁰, Zameer-ul-Asar¹⁰, Imran Saleem¹⁰, Zoubia Hina¹⁰, Unzila Ishtiaq¹⁰,
7 Nadeem Qamar¹¹, Nadeem Hayat Mallick¹², Zia Yaqub¹¹, Tahir Saghir¹¹, Syed Nadeem
8 Hasan Rizvi¹¹, Anis Memon¹¹, Mohammad Ishaq¹³, Syed Zahed Rasheed¹³, Fazal-ur-
9 Rehman Memon¹⁴, Anjum Jalal¹⁵, Shahid Abbas¹⁵, Philippe Frossard¹⁰, Danish
10 Saleheen^{10,16,*}, Angela M. Wood^{1,4,5,6,7,8,*}, Julian L. Griffin^{17,18,*}, & Albert Koulman^{19,*}.

- 11
12 1. British Heart Foundation Cardiovascular Epidemiology Unit, Department of Public
13 Health and Primary Care, University of Cambridge, Cambridge, CB1 8RN, UK.
14 2. Stroke Research Group, Department of Clinical Neurosciences, University of
15 Cambridge, Cambridge, CB2 0QQ, UK.
16 3. Internal Medicine Research Unit, Pfizer Worldwide Research, Development and
17 Medical, Cambridge, Massachusetts 02139, USA.
18 4. British Heart Foundation Centre of Research Excellence, University of Cambridge,
19 Cambridge, CB2 0QQ, UK.
20 5. National Institute for Health Research Blood and Transplant Research Unit in Donor
21 Health and Genomics, University of Cambridge, Cambridge, CB1 8RN, UK.
22 6. National Institute for Health Research Cambridge Biomedical Research Centre,
23 University of Cambridge and Cambridge University Hospitals, Cambridge, CB2 0QQ,
24 UK.
25 7. Health Data Research UK Cambridge, Wellcome Genome Campus and University of
26 Cambridge, Cambridge, CB10 1SA, UK.
27 8. Department of Human Genetics, Wellcome Sanger Institute, Hinxton, CB10 1SA,
28 UK.
29 9. Inflammation and Immunology, Pfizer Worldwide Research, Development and
30 Medical, 10785 Berlin, Germany.
31 10. Center for Non-Communicable Diseases, Karachi 75300, Pakistan.
32 11. National Institute of Cardiovascular Diseases 75510, Karachi, Pakistan.
33 12. Punjab Institute of Cardiology, Lahore 42000, Pakistan.
34 13. Karachi Institute of Heart Diseases, Karachi 75950, Pakistan.

NOTE: This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.

- 1 14. Red Crescent Institute of Cardiology, Hyderabad 71500, Pakistan.
2 15. Faisalabad Institute of Cardiology, Faisalabad 38000, Pakistan.
3 16. Department of Biostatistics & Epidemiology, University of Pennsylvania,
4 Philadelphia, Pennsylvania 19104, USA.
5 17. Department of Biochemistry and Cambridge Systems Biology Centre, University of
6 Cambridge, Cambridge, CB2 1GA, UK.
7 18. Section of Biomolecular Medicine, Division of Systems Medicine, Department of
8 Metabolism, Digestion, and Reproduction, Imperial College London, SW7 2AZ, UK.
9 19. Core Metabolomics and Lipidomics Laboratory, National Institute for Health
10 Research, Cambridge Biomedical Research Centre, Cambridge, CB2 0QQ, UK.

11

12 * These authors contributed equally.

13

14 Correspondence and requests for materials should be addressed to E.L.H. (email:
15 eh457@medschl.cam.ac.uk) or J.L.G. (julian.griffin@imperial.ac.uk) or A.K.
16 (ak675@medschl.cam.ac.uk).

17

18 Short title: Genetic determinants of lipid metabolites in South Asians

19 Word count: Abstract: 218; Main text: 6,463

20 5 main figures, 15 supplementary tables, and 6 supplementary figures

21

1 **ABSTRACT**

2 **Background:** Genetic, lifestyle, and environmental factors can lead to perturbations in
3 circulating lipid levels and increase risk of cardiovascular and metabolic diseases. However,
4 how changes in individual lipid species contribute to disease risk is often unclear.
5 Moreover, little is known about the role of lipids on cardiovascular disease in Pakistan, a
6 population historically underrepresented in cardiovascular studies.

7
8 **Methods:** We characterised the genetic architecture of the human blood lipidome in 5,662
9 hospital controls from the Pakistan Risk of Myocardial Infarction Study (PROMIS) and
10 13,814 healthy British blood donors from the INTERVAL study. We applied a candidate
11 causal gene prioritisation tool to link the genetic variants associated with each lipid to the
12 most likely causal genes, and Gaussian Graphical Modelling network analysis to identify
13 and illustrate relationships between lipids and genetic loci.

14
15 **Results:** We identified 359 genetic associations with 255 lipids measured using direct
16 infusion high-resolution mass spectrometry in PROMIS, and 616 genetic associations with
17 326 lipids in INTERVAL. Our analyses revealed new biological insights at genetic loci
18 associated with cardiometabolic diseases, including novel lipid associations at the *LPL*,
19 *MBOAT7*, *LIPC*, *APOE-C1-C2-C4*, *SGPP1*, and *SPTLC3* loci.

20
21 **Conclusions:** Our findings, generated using a distinctive lipidomics platform in an
22 understudied South Asian population, strengthen and expand the knowledge base of the
23 genetic determinants of lipids and their association with cardiometabolic disease-related
24 loci.

25
26 **Keywords:** lipidomics, genetics, Gaussian Graphical Modelling, network analysis, South
27 Asian

28

1 **BACKGROUND**

2 Mass spectrometry-based lipidomics, which aims to capture information on the full
3 complement of lipid metabolites in a given biological sample [1], holds the potential to
4 identify novel insights leading to lipid regulation and dyslipidaemia, potentially providing
5 new mechanisms that link lipid perturbances with cardiometabolic disorders. While
6 pathways underlying dyslipidaemia have been widely studied, we still do not understand
7 how individual lipid species are regulated or contribute to disease. With increasing rates
8 of cardiometabolic diseases in low- and middle-income countries, there is a need for well-
9 powered studies to understand the mechanisms that lead to such disorders in these
10 settings. This need is especially acute for genetic studies where the overrepresentation of
11 individuals of European ancestry amongst genotyped cohorts has led to ancestral bias in
12 effect size estimates at both the genotype and polygenic score levels [2].

13
14 In this study, we aimed to identify novel genetic associations with lipid metabolites in an
15 understudied South Asian population and determine plausible metabolic pathways for the
16 significantly associated lipid metabolites. We performed a comprehensive interrogation of
17 genetic influences on the human blood serum lipidome using direct infusion high resolution
18 mass spectrometry (DIHRMS). We quantified 360 lipid metabolites in 5,662 individuals
19 from Pakistan, from which we identified 359 genotype–lipid associations (lipid quantitative
20 trait loci, or lipid QTLs [3, 4]) at 24 independent loci, providing new insights into lipid
21 metabolism and its impact on cardiovascular and metabolic diseases.

22
23 To help disentangle which of these findings are specific to the Pakistani population and
24 which are unique to the lipid platform itself, we also carried out a parallel set of analyses
25 using the same lipidomics platform in a much larger cohort of individuals from the UK. We
26 measured 432 lipid metabolites in 13,814 healthy British blood donors, from which we
27 identified 616 lipid QTLs at 38 independent loci.

28

1 **METHODS**

2 **Study description**

3 Our primary analyses involved a subset of participants from the Pakistan Risk of Myocardial
4 Infarction Study (PROMIS), a case-control study of first-ever acute myocardial infarction
5 (MI) in nine urban centres in Pakistan consisting of approximately 16,700 cases and
6 18,600 controls. Details of PROMIS have been described previously [5]. In this analysis
7 we analysed controls (individuals free from MI at baseline), who were identified and
8 recruited at the same hospitals as cases according to the following order of priority: (1)
9 visitors of patients attending the outpatient department, (2) patients attending outpatient
10 clinics for non-cardiac-related symptoms, and (3) non-first-degree relative visitors of MI
11 cases. The present analysis involved serum samples from 5,662 PROMIS controls for which
12 genetic and lipid-profiling data were available. Ethical approval was obtained from the
13 relevant ethics committee of each of the institutions involved in participant recruitment
14 and the Center for Non-Communicable Diseases in Karachi, Pakistan, and informed
15 consent was obtained from each participant recruited into the study, including for use of
16 samples in genetic, biochemical, and other analyses.

17
18 Comparative, parallel analyses were performed in INTERVAL, a prospective cohort study
19 of approximately 50,000 healthy blood donors from the UK. The present analyses involved
20 13,814 participants from INTERVAL with both genetic and lipid-profiling data. Details
21 concerning the INTERVAL study, including DNA extraction and genotyping, lipid profiling,
22 and genome-wide association analyses are provided in the Supplementary Methods
23 (Additional file 1).

24 25 **Lipid profiling**

26 Lipid levels in human serum were quantified using direct infusion high-resolution mass
27 spectrometry (DIHRMS) using an Exactive Orbitrap (Thermo, Hemel Hempstead, UK).
28 Data processing, peak-picking, normalisation, cleaning, and quality control were
29 performed to identify and record signals for 360 known lipids in 5,662 PROMIS

1 participants. The 360 lipids corresponded to five broad lipid categories (fatty acyls and
2 derivatives, glycerolipids, glycerophospholipids, sphingolipids, and sterol lipids), which are
3 further subdivided into fourteen lipid subclasses (Supplementary Table 1 in Additional file
4 2). We have previously described all the details of our lipid profiling, data processing,
5 quality control, and peak-picking process [6]. In brief, lipid profiles were obtained using
6 an open-profiling technique that measured all lipid species across a spectrum. We
7 developed a novel peak-picking algorithm [6] to select all lipids within an m/z window of
8 185-1000, with a time window of 20-70 seconds for lipids in positive ionisation mode and
9 95-145 seconds for lipids in negative ionisation mode. A lipid list containing all known
10 lipids within this m/z range was used to extract information on the lipid concentrations at
11 specific peaks of interest, consisting of 1,305 lipids in positive ionisation mode and 3,772
12 lipids in negative ionisation mode. Quality control samples and blanks were used to remove
13 lipids that were not able to be detected or had poor quality of assessment, resulting in a
14 final list of 360 distinct lipid annotations across both ionisation modes.

15

16 **Genotyping and imputation**

17 DNA from PROMIS participants was extracted from leukocytes in Pakistan and genotyped
18 at the Wellcome Sanger Institute in Cambridge, UK on either (1) the Illumina 660-Quad
19 GWAS platform, which consisted of 527,925 genotyped autosomal variants after quality
20 control (QC) steps were performed, or (2) the Illumina HumanOmniExpress GWAS
21 platform, which consisted of 643,333 genotyped autosomal variants after QC. Genetic
22 samples were removed if (1) they were heterozygosity outliers (heterozygosity $>$ mean \pm
23 3 SD), (2) the sample call rate was less than 97%, (3) there was discordant sex between
24 genetically-inferred and self-reported sex, or (4) they were duplicate or related pairs
25 (kinship coefficient $>$ 0.375). Single nucleotide polymorphisms (SNPs) were excluded if
26 (1) the SNP call rate was less than 97%, (2) there was evidence of departure from Hardy-
27 Weinberg Equilibrium (HWE) at a P -value of less than 1×10^{-7} , or (3) the minor allele
28 frequency (MAF) was less than 1%. Imputation was applied to the cleaned PROMIS
29 datasets using the 1000 Genomes Project March 2012 (v3) release [7] as the reference

1 panel. Imputation was conducted using IMPUTE v2.1.0 [8] using 5-Mb non-overlapping
2 intervals for the whole genome. Once imputation had been performed for the samples on
3 both genotyping platforms separately, there were over 7.2 million imputed SNPs available
4 for analyses in either dataset before further QC. SNPs were removed if they were poorly
5 imputed, i.e. if they had an information score (an assessment of the level of accuracy of
6 imputation) < 80%. The results were then extracted from the output files, and once the
7 final QC filters were reapplied, 6,720,657 SNPs were available for analyses of the
8 lipidomics data. In total, 5,662 individuals from PROMIS had concomitant information on
9 lipidomics data and imputed SNPs.

10

11 **Primary genome-wide association analyses**

12 In PROMIS, linear regression was used to determine the association of each lipid with each
13 SNP using SNPTEST v2.4.1 [9], which was performed separately for the samples
14 genotyped on each of the two genetic platforms. Residuals were calculated from the null
15 model for each lipid, which included adjustment for age group, sex, date of survey, plate
16 (batch), and fasting status. To account for population stratification and genetic
17 substructure in the data, principal component analysis was conducted on the multi-
18 dimensional scaling matrix created from autosomal SNPs as implemented in PLINK; the
19 first six principal components were subsequently added to each model. A missing data
20 likelihood score test was used when testing for association at imputed SNPs to account for
21 genotype uncertainty. Beta estimates and standard errors from the association results for
22 the two genetic platforms were combined in a fixed-effect inverse-variance-weighted
23 meta-analysis using METAL version 2011-03-25 [10]. The threshold for genome-wide
24 significance level was set to $P < 8.929 \times 10^{-10}$, which corrected for multiple testing by
25 dividing the standard genome-wide significance level (5×10^{-8}) by the number of principal
26 components (56) that explained over 95% of the variance in the levels of the lipids. All
27 traits gave genomic inflation factors (λ) in the meta-analysis less than 1.05 [mean (SD)
28 1.0139 (0.0129); range 0.9741-1.0455], indicating that there was little evidence of
29 systematic bias in the test statistics.

1

2 To verify the robustness and validity of the results, post-analysis quality control (QC) was
3 performed by comparing the meta-analysis results with the results on each GWAS
4 platform. The lead SNPs from the meta-analysis were only kept if they (1) passed QC in
5 the raw SNPTTEST results from both GWAS platforms (i.e. HWE $P < 1 \times 10^{-7}$, call rate $<$
6 0.97 , MAF < 0.01 , and info score < 0.80); (2) had beta (β) estimates in the same direction
7 on both platforms (i.e. betas were both negative or both positive); and (3) had $P < 0.01$
8 on both platforms (with $P < 8.9 \times 10^{-10}$ in the meta-analysis).

9

10 **Genome-wide analysis of ratios of lipids**

11 A second discovery step was carried out in PROMIS by testing genome-wide associations
12 on 26 pairwise ratios of lipid concentrations. Ratios were identified based on those that
13 had strong biological rationales and that acted through thoroughly understood metabolic
14 pathways (Supplementary Table 5). Meta-analysis was performed to combine results from
15 the two genotyping platforms using a fixed-effect inverse-variance weighted meta-
16 analysis. Since there were fewer statistical tests for the ratios than for the individual lipids,
17 the combined results file for each ratio was filtered using the standard threshold for
18 genome-wide significance of $P < 5 \times 10^{-8}$.

19

20 **Conditional analyses**

21 We conducted conditional analyses on the significant loci from the meta-analysis results
22 of the univariate GWAS for each lipid in PROMIS. All SNPs were selected where $P < 8.9 \times$
23 10^{-10} , the 5-Mb chunks were identified where each of these SNPs were located, and the
24 lead SNPs were selected within each chunk that had the strongest P -value. On an individual
25 lipid basis, for each 5-Mb chunk that was identified, SNPTTEST was run on the imputed data
26 for each genotyping platform using the same null model as before, except also conditioning
27 on the lead SNP in the identified chunk. The results from the samples analysed on each
28 genotyping platform were combined in a meta-analysis using METAL as described above,
29 and any SNPs where $P < 8.9 \times 10^{-10}$ were identified. The lead SNP from the meta-analysed

1 results of the first conditional analysis (i.e. the SNP with the strongest *P*-value) was
2 identified, and this process was repeated for each chunk. Additional SNPs to be conditioned
3 on were repeatedly added to the model on each chunk for each lipid until there were no
4 more significant SNPs left within that chunk. The final set of SNPs that were “conditionally
5 independent” for each lipid were combined into a single list across all lipids, resulting in
6 359 SNP-lipid associations (lipid QTLs) for 255 lipids, or 90 unique lead SNPs. These
7 variants were grouped into 24 loci using a distance measure of ± 500 -Kb.

8
9 We identified the proportion of variation in the lipidome explained by inherited genetic
10 variants by regressing each lipid on the number of copies of each allele held by each
11 participant for each of the conditional analysis sentinel SNPs.

12

13 **Candidate gene annotation**

14 In order to prioritise candidate genes that might underpin the genotype—lipid associations,
15 we applied the ProGeM framework (Supplementary Figure 5 in Additional file 1) to both
16 PROMIS and INTERVAL [11]. In addition to reporting the nearest gene to the sentinel
17 variant, ProGeM combines information from complementary “bottom-up” and “top-down”
18 approaches to assess the credibility of potential candidate genes [11] (Supplementary
19 Table 7). In the bottom-up approach, we annotated SNPs according to their putative
20 effects on proximal gene function by examining whether these SNPs influence protein
21 sequencing, gene splicing, and/or mRNA levels of a local gene (Supplementary Table 8).
22 Conversely, in the top-down approach, we annotated SNPs according to previous
23 knowledge concerning local gene function by examining whether proximal genes have
24 been previously implicated in lipid metabolism (Supplementary Table 9). In cases where
25 (1) SNPs were purported to exert effects on more than one local gene and/or (2) more
26 than one local gene was previously implicated in lipid metabolism, we assigned SNPs to
27 multiple genes rather than force-assigning each to a single gene. In cases where it was
28 not possible to annotate SNPs using either the bottom-up or top-down approach, we

1 assigned the SNPs to their nearest gene. Further details of the candidate gene annotation
2 approach that we followed are described in the Supplementary Methods (Additional file 1).

3

4 After performing comprehensive annotation of SNPs as per the bottom-up and top-down
5 procedures, we then integrated this information to try to predict the most likely causal
6 gene(s) using a hierarchical approach as follows: (1) For those lead SNPs where the same
7 gene was highlighted by both the bottom-up and the top-down approach, we selected this
8 gene as the putative causal gene; (2) If both the SNP (from this study) and the proximal
9 gene (from IPA) were associated with the same lipid subclass, we made further SNP-gene
10 assignments accordingly; (3) Finally, for each of the remaining lead SNPs, we assigned
11 the highest scoring top-down gene and any bottom-up genes as the likely causal gene(s).

12
13 Separately, we assigned an expertly-curated causal gene to each variant and compared
14 the predicted causal genes identified by the functional annotation pipeline to assess
15 concordance and validate the pipeline.

16

17 **Gaussian Graphical Modelling**

18 As described previously [6], we estimated a Gaussian Graphical Model (GGM) on the
19 normalised relative intensities of the lipids in PROMIS to better resolve lipid cross-
20 correlations. The GGM resulted in a set of edges in which each edge connected two
21 detected lipids if their cross-correlation conditioned on all other lipids was significantly
22 different from zero. Subjects with more than 10% missing lipids as well as lipids with more
23 than 20% missing subjects were removed from the analysis. The “genenet” R package
24 was used to infer the GGM [12]. A similar approach for metabolomics data has been
25 suggested previously [13]. To focus on strong effects we retained only edges in the model
26 that met an FDR cutoff of 0.05 and had a partial correlation coefficient greater than 0.2.

27

1 **Fatty acid chain enrichment analysis**

2 We manually annotated detected lipids in PROMIS with their constituent fatty acid chains.
3 For each combination of fatty acid chains, we counted the number of GGM edges
4 connecting lipids with that specific combination, which we used to directly estimate *P*-
5 values of enrichment and depletion. To test whether edges from the GGM were enriched
6 for any combination of fatty acid chains, we permuted the annotation 1000 times using
7 the R package "BiRewire" [14], keeping the number of annotations per lipid and fatty acid
8 chain constant.

10 **Network of genetic and metabolic associations**

11 We used Cytoscape v3.2.1 [15] to generate a network of associations between genes and
12 lipid subclasses in PROMIS (Figure 3). Using a previously described approach [16], we
13 constructed a GGM to connect lipids to each other based on partial correlation coefficients,
14 and we also connected lipids with genetic loci using the conditional analysis results, with
15 one link for each genome-wide significant association. The full network facilitates
16 visualisation of the genetic determinants of human metabolism and the relationships
17 between genetic loci and lipid subclasses.

18
19 The network diagrams were created by combining two parts to integrate different sources
20 of information. The first part was created by loading the reported associations between
21 lipids and genes into Cytoscape. Lipid species were clustered according to the lipid subclass
22 they belong to, resulting in fourteen distinct lipid subclass nodes in the network. The 90
23 identified lead SNPs from the conditional analyses were clustered according to their
24 corresponding predicted causal gene(s), which was determined using the ProGeM
25 framework [11]. In cases where it was not possible to confidently identify a single
26 predicted causal gene, loci were entered into the network instead. For the second part, a
27 functional interaction network consisting solely of our list of predicted causal genes/loci
28 was created in Cytoscape using interaction network data downloaded from Ingenuity
29 Pathway Analysis (IPA) that had been merged using in-house R scripts to create a .sif file.

1 For loci with multiple potential causal genes, interaction networks for all genes were
2 extracted from IPA and an edge was drawn if at least one gene at that locus functionally
3 interacts with another of our lipid-associated genes according to IPA. Finally, these two
4 parts were merged together by node names (i.e. gene symbols). No enrichment statistics
5 (e.g. KEGG pathways or GO terms) or other statistical information was used to produce
6 the network, since this information was already incorporated to inform the predictions of
7 the most likely “causal” genes, and would therefore invalidate the conclusions if it was also
8 used to inform the network.

9
10 A second network diagram was created containing a subset of the first network containing
11 only the triglyceride species (Figure 4). It also provides more detail as it shows the
12 individual triglycerides rather than the lipid subclass as a whole. Thus, it portrays the
13 partial correlations of the triglycerides with each other and the association of each
14 triglyceride with genetic loci.

15

1 RESULTS

2 Genetic architecture of the lipidome in South Asians and in the UK

3 We performed a genome-wide association study (GWAS) on the levels of 360 lipid
4 metabolites using 6.7 million imputed autosomal variants in 5,662 hospital-based controls
5 from PROMIS. We applied DIHRMS to quantify serum lipid metabolites across five broad
6 lipid categories, i.e. fatty acyls and derivatives, glycerolipids, glycerophospholipids,
7 sphingolipids, and sterol lipids [6]. We demonstrated the robustness of these lipid
8 measurements in several ways, including validation of lipid signals against blanks, pooled
9 samples, and internal standards, as we described previously [6]. Additionally, we
10 replicated known associations of lipid metabolites with previously reported major lipid loci
11 (Supplementary Table 15). After Bonferroni correction for multiple testing of variants and
12 lipid metabolites ($P < 8.929 \times 10^{-10}$), we found 359 significant associations between 255
13 lipid metabolites and 24 genomic regions (Figure 1, Figure 2, Supplementary Figure 1,
14 Supplementary Table 2). The majority of these lipid metabolites (67%; $n = 171$) were
15 associated with variation at a single locus, while 26% of lipid metabolites were associated
16 with two loci and 7% were associated with three or more loci (Supplementary Figures 2a
17 and 3). To detect multiple independent associations at the same locus, we used stepwise
18 conditional analysis, identifying 90 conditionally independent variants associated with lipid
19 metabolites (Supplementary Table 3). 335 (93%) of the lipid QTLs had multiple
20 conditionally significant associations (Supplementary Figure 2b).

21
22 Using the same DIHRMS platform, we also performed a GWAS on levels of 432 lipid
23 metabolites using 87.7 million imputed autosomal variants in 13,814 British blood donors
24 from INTERVAL. We identified significant associations with lipids at 38 independent loci
25 (Figure 1, Supplementary Table 4). There was considerable consistency in the genomic
26 regions identified in each study, with 18 (75%) of the significant genetic loci from PROMIS
27 also found in INTERVAL (Figure 1). Six genetic loci were specific to lipid levels in the
28 Pakistani population: *ANGPTL3*, *UGT8*, *PCTP*, *C19orf80*, *XBP1*, and *GAL3ST1*. There were

1 also twenty genetic loci associated with lipids in the British population that were not
2 significantly associated with lipids in the Pakistani population.

3

4 In PROMIS, the median proportion of variation in the lipidome explained by the genome-
5 wide significant conditionally independent variants was 1.7% (interquartile range: 1.5-
6 1.9%) (Supplementary Figure 2c), which is slightly less than that reported in
7 metabolomics studies [16–19] but similar to the reported variation explained in previous
8 lipidomics studies [20, 21]. There was a strong inverse relationship between effect size
9 and minor allele frequency (MAF) (Supplementary Figure 2d), consistent with previous
10 GWAS of quantitative traits [22, 23]. Approximately 70% of the analysed genetic variants
11 in this analysis were common (MAF >5%) and 30% were low-frequency (MAF: 1-5%) with
12 a median MAF of 8%. To help identify candidate causal genes through which genetic loci
13 may influence lipid levels and thereby impact disease risk, we applied the ProGeM
14 framework [11] (Supplementary Tables 7-13, Supplementary Figure 5). We identified a
15 plausible or established link to biochemical function for 16 of the 24 loci (including *GCKR*,
16 *LPL*, *FADS1-2-3*, and *APOA5-C3*), involving 34 unique genes. In cases where it was not
17 possible to annotate SNPs using our systematic approach, we assigned them to their
18 nearest protein-coding gene.

19
20 Previous studies have shown that the ratios of metabolites can strengthen association
21 signals and lead to a better understanding of possible mechanisms [16]. Thus, in addition
22 to the individual lipid metabolites, we selected twenty-six ratios of lipid metabolites that
23 act through well-understood metabolic pathways. These included ratios associated with
24 lipase activity, elongases, docosahexaenoic acid (DHA) levels, dairy fat intake, insulin
25 production, glucose control, *de novo* lipogenesis, and cardiovascular disease risk
26 (Supplementary Table 5). Genome-wide association analyses of these ratios in PROMIS
27 resulted in the identification of four additional loci that were not detected in the GWAS of
28 individual lipid metabolites (*MYCL1-MFSD2A*, *LPGAT1*, *LOC100507470*, and *HAPLN4-
29 TM6SF1*) (Supplementary Table 6).

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

Network of genetic and metabolic associations

To identify and visualise the connectivity between lipid subclasses, we generated a network of genetic and metabolic associations in PROMIS by summarising within each subclass the pairwise partial correlations between lipid metabolites and their genetic associations (Figure 3). This network diagram highlights that the number of connections between diglycerides and triglycerides was strongly over-represented in the Gaussian Graphical Model (GGM), indicating that there were more significant partial correlations between lipids from these subclasses than would be expected due to chance alone, whereas the number of connections between sphingomyelins and triglycerides was strongly under-represented in the GGM. In addition to being associated with variants from the *SPTLC3* and *FADS1-2-3* loci, we found that sphingomyelins were associated exclusively with four loci that were not associated with any other lipid subclasses: *GCKR*, *SGPP1*, *MLXIPL*, and *XBP1*.

Given the striking findings for triglycerides in the overall network diagram, we also generated a network in PROMIS for a subset of the triglyceride species showing the partial correlations of individual triglycerides and their detailed associations with genetic loci (Figure 4). This network diagram shows that variants in the *APOA5-C3* locus are associated with a wide range of triglycerides, consistent with previous associations of Apolipoprotein A-V (ApoA5) with plasma triglyceride levels. ApoA5 is a component of a number of lipoprotein fractions including HDL, VLDL, and chylomicrons, and it may regulate the catabolism of triglyceride-rich lipoprotein particles by *LPL* and/or play a role in the assembly of VLDL particles [24–28]. The network mainly shows links with triglycerides containing polyunsaturated fatty acids (PUFAs), suggesting that variants in the *APOA5-C3* locus mainly affect the catabolism of lipoproteins containing triglycerides derived from adipose tissues that are relatively enriched in more unsaturated fatty acids. In contrast, we did not see direct links of fully saturated triglycerides with the *APOA5-C3* locus, suggesting that genetic variation at this locus is not particularly involved in the assembly of VLDL particles in the liver as part of *de novo* lipogenesis (see Supplementary Figure 4).

1
2 Fatty acid desaturase is key in the production of PUFAs; therefore, differences in *FADS1*-
3 2-3 activity are expected to be observed in triglycerides with a large number of double
4 bonds and carbon atoms. Indeed, the GGM concurs with established biochemistry since
5 this locus is associated with triglycerides (TG) 56:6, 56:7, and 58:9 but is not associated
6 with triglycerides with fewer double bonds or carbon atoms. In contrast, it is unclear why
7 variants in the *PNPLA3* locus also have the strongest associations with triglycerides with a
8 relatively larger number of carbon atoms and double bonds, namely TG(56:5) and
9 TG(56:6) (see also Figure 5). One possible explanation is that significantly associated
10 variants in the *PNPLA3* locus are changing the substrate specificity so that there is a shift
11 in the relative amounts of triglycerides that are exported from the liver.

12
13 Additionally, the network diagram confirms that *LPL* is mainly active on MUFAs in
14 triglyceride species. Variants in the *LPL* locus are significantly associated with TG(52:2),
15 TG(52:3), TG(53:2), and TG(53:3), which have a high probability of containing one or
16 more MUFAs within their fatty acid side chains. Figure 2 also shows that triglycerides and
17 diglycerides are predominantly inversely associated with *LPL* variants, while triglycerides
18 are positively associated with *PNPLA3* variants. Variants in the *LPL* locus are also positively
19 associated with phosphocholines, sphingomyelins, and cholesterol esters, although the
20 associations for the majority of the lipids in these subclasses did not reach genome-wide
21 significance.

22

23 **New biological insights into lipid metabolism**

24 Our analysis replicated known associations between lipids and genetic loci while also
25 further extending what is known about these loci. We found significant associations of a
26 wide range of lipids, including phosphatic acid (PA) 39:1, phosphatidylcholine (PC) 35:4,
27 and phosphatidylethanolamines (PE) 36:4, 36:5, and 38:6, with variants in the *LIPC* locus
28 (Supplementary Figure 6i); and significant associations of six specific sphingomyelins
29 [SM(34:0), SM(40:0), SM(40:1), SM(40:2), SM(42:0)+AcO⁻, and SM(42:1)] and two

1 phosphatidylcholines [PC-O(37:1) and PC-O(39:1)] with variants in the *APOE-C1-C2-C4*
2 locus (Supplementary Figure 6b). We also identified significant associations of four further
3 sphingomyelins [SM(31:1)-H⁻, SM(32:1), SM(32:1)+AcO⁻, and SM(39:1)] with variants in
4 the *SGPP1* locus (Supplementary Figure 6s). Additionally, we found significant associations
5 of nine ceramides [Cer(40:0)-H⁻, Cer(40:1)-H⁻, Cer(40:2)-H⁻, Cer(41:0)-H⁻, Cer(41:1)-H⁻
6 , Cer(41:2)-H⁻, Cer(42:0)-H⁻, Cer(42:1)-H⁻, and Cer(42:2)-H⁻] with variants in the *SPTLC3*
7 locus, which have not previously been reported in relation to this locus, as well as
8 significant associations with three phosphatidylcholines and fifteen sphingomyelins
9 (Supplementary Figure 6t).

10

11 We also discovered genetic associations with lipids at the patatin-like phospholipase
12 domain containing protein 3 (*PNPLA3*) and membrane bound *O*-acyltransferase domain
13 containing 7 (*MBOAT7*) loci that may have important biological and clinical implications.
14 We found significant associations of two triglycerides—TG(56:6) (*m/z* 924.801) and
15 TG(56:5) (*m/z* 926.817)—with rs12484809, an intronic variant in the *PNPLA3* locus
16 (Supplementary Figure 6q). We also we found that the lead SNP in the *MBOAT7* locus,
17 rs8736 (chr19:54677189), was associated with a wide range of phosphatic acids [e.g.
18 PA(40:5) and PA(44:6)], phosphatidylcholines [e.g. PC(36:6) and PC(42:11)],
19 phosphatidylethanolamines [e.g. PE(39:7)], and phosphoinositols [e.g. PI(34:1) and
20 PI(36:1)] (Supplementary Figure 6k).

21

22 We undertook further investigation of a related nonsynonymous *PNPLA3* variant that is in
23 moderate LD ($r^2 = 0.695$), rs738409 (p.Ile148Met), to study the associations of lipids with
24 *PNPLA3* in greater detail, including those that did not reach genome-wide significance. We
25 focused on this variant rather than rs12484809 because I148M is already known to be
26 associated with total triglycerides [29] and has been extensively characterised in previous
27 genetic and functional analyses, and therefore is more likely to have potential clinical
28 applications. As shown in Figure 5a, the *PNPLA3* I148M allele was associated with
29 increased levels of lipids of higher carbon number and double-bond content, and

1 consistently, with decreased levels of lipids of lower carbon number and double-bond
2 content. There were also significant differences between the mean levels of the
3 triglycerides TG(57:10), TG(46:0), and TG(56:6) between individuals stratified by *PNPLA3*
4 I148M genotype (Figures 5b, 5c, and 5d).

1 **DISCUSSION**

2 Based on a comprehensive analysis of genetic influences on 360 human blood lipids
3 assayed in 5,662 individuals from Pakistan, we identified 359 significant associations
4 between 255 lipids and 24 genetic loci. Additionally, in our analysis of 432 lipids in 13,814
5 British blood donors, we identified significant associations between 326 lipids and 38
6 independent loci. The majority of genetic regions associated with lipids in PROMIS were
7 also found in INTERVAL; those that did not replicate may be due to the increased sample
8 size in INTERVAL which gave a substantial boost in power. These findings suggest that
9 genetically determined aspects of lipid metabolism are broadly similar in individuals of
10 South Asian and European ancestry, and that DIHRMS can reliably capture differences in
11 lipid levels across diverse populations.

12
13 There were six genetic loci specific to lipid levels in PROMIS: *ANGPTL3*, *UGT8*, *PCTP*,
14 *C19orf80*, *XBP1*, and *GAL3ST1*. Angiopoietin-like 3 (*ANGPTL3*) is involved in regulation of
15 lipid and glucose metabolism. SNPs in the *ANGPTL3* region have previously been shown to
16 be associated with major lipids, including LDL-C and total cholesterol [30, 31]. In PROMIS,
17 rs6657050, an intronic variant in the *ANGPTL3* locus, was significantly associated with
18 PC(38:7)+AcO⁻ (*m/z* 862.5603) and PI(36:2)-H⁻ (*m/z* 861.5498) (Supplementary Figure
19 6a).

20
21 UDP glycosyltransferase 8 (*UGT8*) catalyses the transfer of galactose to ceramide, a key
22 enzymatic step in the biosynthesis of galactocerebrosides, which are abundant
23 sphingolipids of the myelin membrane of the central and peripheral nervous system. In
24 PROMIS, rs28870381, an intergenic variant in *UGT8*, was associated with PG(32:1) (*m/z*
25 779.5078) (Supplementary Figure 6u).

26
27 Phosphatidylcholine transfer protein (*PCTP*) catalyses the transfer of phosphatidylcholines
28 between membranes and is involved in lipid binding. Through regulation of plasma lipid
29 concentrations it may also modulate the development of atherosclerosis [32]. In PROMIS,

1 rs11079173, an intronic variant in the *PCTP* locus, was associated with PA(40:5)+AcO⁻ or
2 PG(39:5)-H⁻ (*m/z* 809.5337) (Supplementary Figure 6n).

3
4 *C19orf80*, also known as angiopoietin-like 8 (*ANGPTL8*), is involved in the regulation of
5 serum triglyceride levels, and is associated with major lipids including HDL-C and
6 triglycerides [31]. In PROMIS, rs8101801, an intronic variant in the *C19orf80* locus, was
7 significantly associated with PC(40:9)+AcO⁻ (*m/z* 886.5603) and PI(38:4)-H⁻ (*m/z*
8 885.5498) (Supplementary Figure 6d).

9
10 Galactose-3-*O*-sulfotransferase 1 (*GAL3ST1*) catalyses the sulfation of membrane
11 glycolipids and the synthesis of galactosylceramide sulfate, a major lipid component of the
12 myelin sheath. In PROMIS, rs2267161, a missense variant in the *GAL3ST1* locus, was
13 associated with PG(32:1) (*m/z* 779.5078) (Supplementary Figure 6g).

14
15 X-box binding protein 1 (*XBP1*) functions as a transcription factor during endoplasmic
16 reticulum stress by regulating the unfolded protein response. It is also a major regulator
17 of the unfolded protein response in obesity-induced insulin resistance and T2D for the
18 management of obesity and diabetes prevention. Recent studies have shown that
19 compounds targeting the *XBP1* pathway are a potential approach for the treatment of
20 metabolic diseases [33]. In addition, *XBP1* protein expression, which is induced in the liver
21 by a high carbohydrate diet, is directly involved in fatty acid synthesis through *de novo*
22 lipogenesis. Therefore, compounds that inhibit *XBP1* activation may also be useful for
23 treatment of NAFLD [34]. In PROMIS, rs71661463, an intronic variant for which *XBP1* is
24 the candidate causal gene, was associated with SM(37:1) (*m/z* 745.6216) (Supplementary
25 Figure 6v). Recent research across many species has shown that *XBP1* is a transcription
26 factor regulating hepatic lipogenesis. In mice, hepatic *XBP1* expression is regulated by
27 proopiomelanocortin (POMC) during sensory food perception and coincides with changes
28 in the lipid composition of the liver with increases in PCs and PEs [35]. Although previous
29 studies have shown direct links between *XBP1* and overall lipid metabolism, this is the first

1 time a genetic association has been reported between *XBP1* and lipid metabolites in
2 humans, affecting sphingomyelins, PCs, and PEs (Supplementary Figure 6v).

3
4 Our findings for the *PNPLA3* and *MBOAT7* loci were also notable. *PNPLA3* is a
5 multifunctional enzyme that encodes a triacylglycerol lipase, which mediates
6 triacylglycerol hydrolysis in adipocytes and has acylglycerol *O*-acyltransferase activity. The
7 relationship between rs738409, a nonsynonymous variant (p.Ile148Met) in the *PNPLA3*
8 gene, and non-alcoholic fatty liver disease (NAFLD) has been well established [36]. This
9 variant has been shown to impair triglyceride hydrolysis in the liver and secretion of
10 triglyceride-rich very low density lipoproteins, leading to altered fatty acid composition of
11 liver triglycerides, and is also associated with reduced risk of CHD [37] and increased risk
12 of type 2 diabetes (T2D) [38]. This suggests that targeting hepatic pathways to reduce
13 cardiovascular risk may be complex, despite the clustering of cardiovascular and hepatic
14 diseases in people with metabolic syndrome. Our analysis offers granularity to the
15 previously identified total triglyceride associations with *PNPLA3* by identifying two specific
16 triglyceride species that may have a role in *PNPLA3* function.

17
18 *MBOAT7*, which contributes to the regulation of free arachidonic acid in the cell through
19 the remodelling of phospholipids, was reported as being associated with the metabolite 1-
20 arachidonoylglycerophosphoinositol in a previous mGWAS [16] [known as PI(36:4) in our
21 study], but we found that the lead SNP in this locus, rs8736 (chr19:54677189), was also
22 associated with a wide range of phosphatic acids, phosphatidylcholines,
23 phosphatidylethanolamines, and phosphoinositols (Supplementary Figure 6k). Several
24 studies have shown that *MBOAT7* (also known as lysophosphatidylinositol-acyltransferase
25 1 [*LPIAT1*]) is responsible for the transfer of arachidonoyl-CoA to lysophosphoinositides
26 [39]. The creation of *MBOAT7*-deficient macrophages show a decreased level of PI(38:4)
27 and an increase of PI(34:1) as well as PI(40:5) [40]. The T allele of rs8736, a 3' UTR SNP,
28 shows a similar shift in the phosphatidylinositol metabolism. Our work shows that this SNP
29 is also strongly associated with PI(38:3), which is likely to be the dihomogamma linoleic

1 acid (20:3n6)-containing phosphoinositol. None of the papers testing the substrate
2 specificity of *MBOAT7* have included dihomo-gamma linoleic acid or PI(38:3) in their
3 analysis. Thus, we provide novel evidence in humans that there is an association between
4 *MBOAT7* activity and circulating phosphatidylinositols, a finding that requires further
5 replication.

6
7 Our network diagram helped identify sphingomyelins that were associated exclusively with
8 four loci that were not associated with any other lipid subclasses: *GCKR*, *SGPP1*, *MLXIPL*,
9 and *XBP1*. Sphingomyelins have previously been shown to be associated with *SGPP1* [41],
10 but the associations of sphingomyelins with these other three loci are reported here for
11 the first time. *GCKR* has been shown to be associated with total cholesterol and
12 triglycerides (see Figure 2), and has also been associated with the plasma phospholipid
13 fraction fatty acids 16:0 and 16:1 [42, 43]; most lipids that we found to be associated
14 with *GCKR* (Supplementary Figure 6g) are likely to contain these particular fatty acids. It
15 has been suggested that the glucokinase receptor, encoded by *GCKR*, affects the
16 production of malonyl-CoA, an important substrate for *de novo* lipogenesis [42]. To a
17 similar extent there is a known relation between *MLXIPL* and carbohydrate and lipid
18 metabolism. *MLXIPL* is a transcription factor affecting carbohydrate response element
19 binding protein (CREBP) and therefore also plays a role in lipogenesis. Although both these
20 genes have previously been linked to lipogenesis, we discovered that genetic variation at
21 genes involved in the regulation of lipogenesis have been implicated in altering
22 sphingomyelin concentrations.

23
24 The network diagram also helped recapitulate known biological relationships between
25 lipids. As we established in our previous analysis [6], the number of significant partial
26 correlations between lipids of different subclasses was significantly higher than would be
27 expected due to chance alone. This analysis further showed that genes that were
28 significantly associated with lipids of a particular subclass regulated all of the lipids within
29 the subclass in a similar manner. Therefore, the total concentrations of a given lipid class

1 associated with a genetic locus are less affected by the proportion of fatty acids present
2 in those lipid species.

3

4 In summary, our analyses resulted in the following new insights in an understudied South
5 Asian population: (1) we established that decreased levels of sphingomyelins are
6 associated with genetically lower *LPL* activity; (2) we revealed a wide range of
7 glycerophospholipids that are associated with variants in the *MBOAT7* locus; (3) we
8 identified several new associations of phosphatic acids, phosphocholines, and
9 phosphoethanolamines with variants in the *LIPC* region; (4) we found several novel
10 associations of sphingomyelins and phosphocholines with variants in the *APOE-C1-C2-C4*
11 cluster; (5) we discovered four new associations of sphingomyelins with variants in the
12 *SGPP1* locus; and (6) we found several previously unreported associations of
13 phosphocholines, sphingomyelins, and ceramides with variants in the *SPTLC3* locus. These
14 findings can help further the identification of novel therapeutic targets for prevention and
15 treatment.

16

17 Our investigation into the genetic influences of lipids has several strengths. First, the
18 research involved participants from a population cohort in Pakistan, thereby enhancing
19 scientific understanding of lipid associations in this understudied population, and we
20 compared the findings with a typical Western population of British blood donors using the
21 same lipid-profiling platform. Second, the analysis was based on a relatively large dataset
22 of 5,662 participants from Pakistan and an even larger cohort of 13,814 individuals from
23 the UK, thereby increasing statistical power to detect associations. Third, our mGWAS was
24 performed in individuals free from established MI at baseline in PROMIS and healthy blood
25 donors in INTERVAL, which reduces spurious associations due to the disease state or
26 potential treatments. Finally, our newly developed open-profiling lipidomics platform was
27 utilised to provide detailed lipid profiles, with a wider coverage of lipids than most other
28 high-throughput profiling methods [6], which improved our ability to detect novel

1 associations and our understanding of the detailed effects of known lipid loci at the level
2 of individual lipid species.

3

4 Nevertheless, our study has several potential limitations. First, possible selection biases
5 arise from the case-control design of PROMIS, although this was minimised by the
6 recruitment of controls from patients, visitors of patients attending out-patient clinics, and
7 unrelated visitors of cardiac patients. Second, serum samples in PROMIS were stored in
8 freezers at -80 °C for between two to eight years before aliquots were taken for the
9 lipidomics measurements, which we accounted for by adjusting the analyses by the
10 number of years that the samples had been stored. Although residual confounding and
11 deterioration of lipid profiles may still exist, such deterioration is unlikely to have been
12 related to genotype. Third, a majority (76%) of PROMIS participants had not fasted prior
13 to blood draw, and a small proportion of participants (7%) had reportedly fasted for an
14 unknown duration. Recent food consumption may have had significant effects on lipid
15 levels and influenced the results. Our analyses adjusted for fasting status although we
16 lacked statistical power to stratify by fasting status. Fourth, PROMIS participants were
17 recruited from multiple centres in urban Pakistan [6], but it is unclear whether the findings
18 from this study would be generalizable to individuals living in rural villages and other parts
19 of Pakistan, or in other countries in South Asia. However, the confirmatory analysis in
20 INTERVAL, in which we identified significant associations with lipids for the majority of the
21 genetic loci found in PROMIS, helps strengthen the argument that these findings are
22 generalizable. Additionally, many of the lipids were associated with known genetic regions
23 such as *APOA5-C3* and *FADS1-2-3*, which have already been shown to be associated with
24 multiple lipids in other Western populations, further strengthening the validity of the
25 findings from this analysis. Finally, although two-sample Mendelian randomization
26 approaches to make causal inferences about the association of lipids with CHD risk factors
27 and disease outcomes holds great promise in the lipidomics arena [44], extensive
28 pleiotropy made it too difficult to disentangle the findings and we chose not to pursue this
29 avenue. Therefore, although especially stringent procedures were followed, highly

1 conservative cut-offs were used to determine statistical significance, and rigorous pre-
2 analysis and post-analysis quality control steps were performed, there is still a possibility
3 that some of the findings were false positives that arose due to artefacts rather than being
4 true signals. Additional analyses in other populations using the DIHRMS lipidomics platform
5 would be helpful to further replicate our findings. Moreover, the identified pathways and
6 proposed molecular mechanisms require validation through functional analyses in model
7 organisms and humans.

8

9 Further research will be able to leverage these lipidomics results in combination with
10 whole-genome and whole-exome sequencing performed in PROMIS and INTERVAL to help
11 understand the consequences of loss-of-function mutations identified in these participants
12 [45].

13

14 **CONCLUSIONS**

15 In conclusion, this article presents the results from a comprehensive analysis of genetic
16 influences on human blood lipids in South Asians with a comparative analysis in the UK.
17 Our findings strengthen and expand the knowledge base for understanding the genetic
18 determinants of lipids and their association with cardiometabolic disease-related loci.
19 These findings have important implications for the identification of novel therapeutic
20 targets and advancement of mechanistic understanding of metabolic pathways that may
21 lead to the onset of chronic diseases and lipid-related abnormalities.

1 **LIST OF ABBREVIATIONS**

- 2 CHD: Coronary heart disease
- 3 CVD: Cardiovascular disease
- 4 DHA: Docosahexaenoic acid
- 5 DIHRMS: Direct infusion high resolution mass spectrometry
- 6 FDR: False discovery rate
- 7 GGM: Gaussian Graphical Model
- 8 HWE: Hardy-Weinberg Equilibrium
- 9 MAF: Minor allele frequency
- 10 MI: Myocardial infarction
- 11 m/z : Mass-charge ratio
- 12 NAFLD: Non-alcoholic fatty liver disease
- 13 PROMIS: Pakistan Risk of Myocardial Infarction Study
- 14 PUFA: Polyunsaturated fatty acid
- 15 SD: Standard deviation
- 16 SNP: Single nucleotide polymorphism
- 17 T2D: Type 2 diabetes
- 18 QC: Quality Control
- 19 QTL: Quantitative trait loci
- 20

1 **DECLARATIONS**

2 **Ethics approval and consent to participate:**

3 *PROMIS*: The institutional review board at the Center for Non-Communicable Diseases in
4 Karachi, Pakistan approved the study (IRB: 00007048, IORG0005843, FWAS00014490)
5 and all participants gave informed consent, including for use of samples in genetic,
6 biochemical, and other analyses.

7
8 *INTERVAL*: The National Research Ethics Service approved this study (11/EE/0538) and
9 all participants gave electronic informed consent.

10

11 **Consent for publication:** Not applicable.

12

13 **Availability of data and materials:** The datasets used and/or analysed during the
14 current study are available from the corresponding author on reasonable request.

15

16 **Competing interests:** E.B.F. and D.Z. are employees and shareholders of Pfizer, Inc.
17 J.D. has received research funding from the British Heart Foundation, the National Institute
18 for Health Research Cambridge Comprehensive Biomedical Research Centre, the Bupa
19 Foundation, diaDexus, the European Research Council, the European Union, the Evelyn
20 Trust, the Fogarty International Centre, GlaxoSmithKline, Merck, the National Heart, Lung,
21 and Blood Institute, the National Institute for Health Research [Senior Investigator
22 Award], the National Institute of Neurological Disorders and Stroke, NHS Blood and
23 Transplant, Novartis, Pfizer, the UK Medical Research Council, and the Wellcome Trust [*].
24 J.L.G. has received funding from Agilent, Waters, GlaxoSmithKline, Medimmune, Unilever,
25 AstraZeneca, the Medical Research Council, the Biotechnology and Biological Sciences
26 Research Council, the National Institutes of Health, the British Heart Foundation, and the
27 Wellcome Trust. D.Sa. has received funding from Pfizer, Regeneron Pharmaceuticals,
28 Genentech, and Eli Lilly. All other authors declare no competing interests. [*] The views

1 expressed are those of the authors and not necessarily those of the NHS, the NIHR, or the
2 Department of Health and Social Care.

3

4 **Funding:**

5 *PROMIS*: Fieldwork, genotyping, and standard clinical chemistry assays in *PROMIS* were
6 principally supported by grants awarded to the University of Cambridge from the British
7 Heart Foundation (SP/09/002; RG/13/13/30194), the UK Medical Research Council
8 (G0800270; MR/L003120/1), the Wellcome Trust, the EU Framework 6-funded
9 Bloodomics Integrated Project, Pfizer, Novartis, and Merck.

10

11 *INTERVAL*: Participants in the *INTERVAL* randomised controlled trial were recruited with
12 the active collaboration of NHS Blood and Transplant England (<http://www.nhsbt.nhs.uk>),
13 which has supported field work and other elements of the trial. DNA extraction and
14 genotyping was co-funded by the National Institute for Health Research (NIHR), the NIHR
15 BioResource (<http://bioresource.nihr.ac.uk>), and the NIHR [Cambridge Biomedical
16 Research Centre at the Cambridge University Hospitals NHS Foundation Trust] [*]. The
17 academic coordinating centre for *INTERVAL* was supported by core funding from: NIHR
18 Blood and Transplant Research Unit in Donor Health and Genomics (NIHR BTRU-2014-
19 10024), UK Medical Research Council (MR/L003120/1), British Heart Foundation
20 (SP/09/002, RG/13/13/30194; RG/18/13/33946) and the NIHR [Cambridge Biomedical
21 Research Centre at the Cambridge University Hospitals NHS Foundation Trust] [*]. A
22 complete list of the investigators and contributors to the *INTERVAL* trial is provided in
23 reference [46]. The academic coordinating centre would like to thank blood donor staff
24 and blood donors for participating in the *INTERVAL* trial. [*] The views expressed are those
25 of the authors and not necessarily those of the NHS, the NIHR, or the Department of
26 Health and Social Care.

27

28 This work was supported by Health Data Research UK, which is funded by the UK Medical
29 Research Council, Engineering and Physical Sciences Research Council, Economic and

1 Social Research Council, Department of Health and Social Care (England), Chief Scientist
2 Office of the Scottish Government Health and Social Care Directorates, Health and Social
3 Care Research and Development Division (Welsh Government), Public Health Agency
4 (Northern Ireland), British Heart Foundation, and Wellcome.

5
6 J.L.G. and A.K. are funded by the UK Medical Research Council under the Lipid Dynamics
7 and Regulation supplementary grant (MC_PC_13030) and Lipid Programming and
8 Signalling program grant (MC_UP_A090_1006) and Cambridge Lipidomics Biomarker
9 Research Initiative (G0800783). D.S.P. and D.St. are funded by the Wellcome Trust
10 (105602/Z/14/Z).

11
12 **Authors' contributions:** E.L.H., J.D., D.Sa., J.L.G., and A.K. conceived and designed the
13 study. J.D. and D.Sa. are principal investigators of PROMIS. A.M.W., J.L.G., and A.K. jointly
14 supervised the research. A.K. generated the lipidomics data. E.L.H. and A.K. processed
15 the lipidomics data. E.L.H. performed the bioinformatics and statistical analyses. E.B.F.,
16 D.St., D.S.P., D.Z., R.M.Y.O., A.S.B., A.M.W., J.L.G., and A.K. contributed important
17 intellectual content to the study and manuscript. E.L.H., A.S.B., A.M.W., J.L.G., and A.K.
18 were involved in drafting the manuscript. All authors read and approved the final
19 manuscript.

20
21 **Acknowledgements:** The authors would like to thank Michael Inouye for his helpful
22 comments on an earlier version of the manuscript.

23

1 **SUPPLEMENTARY INFORMATION**

2 **Additional file 1.** Supplementary Methods; Supplementary Figures 1-6; References for
3 Supplementary Material

4 **Additional file 2.** Supplementary Tables 1-15

5 **Additional file 3.** Supplementary Figure 1 (high resolution)

6

7

1 REFERENCES

- 2 1. Griffin JL, Atherton H, Shockcor J, Atzori L. Metabolomics as a tool for cardiac research.
3 Nat Rev Cardiol. 2011;8:630–43. doi:10.1038/nrcardio.2011.138.
- 4 2. Martin AR, Gignoux CR, Walters RK, Wojcik GL, Neale BM, Gravel S, et al. Human
5 demographic history impacts genetic risk prediction across diverse populations. Am
6 J Hum Genet. 2017;100:635–49. doi:10.1016/J.AJHG.2017.03.004.
- 7 3. Jha P, McDevitt MT, Halilbasic E, Williams EG, Quiros PM, Gariani K, et al. Genetic
8 regulation of plasma lipid species and their association with metabolic phenotypes.
9 Cell Syst. 2018;6:709–21. doi:10.1016/j.cels.2018.05.009.
- 10 4. Jha P, McDevitt MT, Gupta R, Quiros PM, Williams EG, Gariani K, et al. Systems analyses
11 reveal physiological roles and genetic regulators of liver lipid species. Cell Syst.
12 2018;6:722–33. doi:10.1016/J.CELS.2018.05.016.
- 13 5. Saleheen D, Zaidi M, Rasheed A, Ahmad U, Hakeem A, Murtaza M, et al. The Pakistan
14 Risk of Myocardial Infarction Study: a resource for the study of genetic, lifestyle and
15 other determinants of myocardial infarction in South Asia. Eur J Epidemiol.
16 2009;24:329–38. doi:10.1007/s10654-009-9334-y.
- 17 6. Harshfield EL, Koulman A, Ziemek D, Marney L, Fauman EB, Paul DS, et al. An unbiased
18 lipid phenotyping approach to study the genetic determinants of lipids and their
19 association with coronary heart disease risk factors. J Proteome Res. 2019;18:2397–
20 410. doi:10.1021/acs.jproteome.8b00786.
- 21 7. 1000 Genomes Project Consortium, Abecasis GR, Auton A, Brooks LD, DePristo MA,
22 Durbin RM, et al. An integrated map of genetic variation from 1,092 human genomes.
23 Nature. 2012;491:56–65. doi:10.1038/nature11632.
- 24 8. Howie BN, Donnelly P, Marchini J. A flexible and accurate genotype imputation method
25 for the next generation of genome-wide association studies. PLoS Genet.
26 2009;5:e1000529. doi:10.1371/journal.pgen.1000529.
- 27 9. Marchini J, Howie B. Genotype imputation for genome-wide association studies. Nat Rev
28 Genet. 2010;11:499–511. doi:10.1038/nrg2796.
- 29 10. Willer CJ, Li Y, Abecasis GR. METAL: fast and efficient meta-analysis of genomewide
30 association scans. Bioinformatics. 2010;26:2190–1.
31 doi:10.1093/bioinformatics/btq340.
- 32 11. Stacey D, Fauman EB, Ziemek D, Sun BB, Harshfield EL, Wood AM, et al. ProGeM: a
33 framework for the prioritization of candidate causal genes at molecular quantitative
34 trait loci. Nucleic Acids Res. 2019;47:e3. doi:10.1093/nar/gky837.
- 35 12. Opgen-Rhein R, Strimmer K. From correlation to causation networks: a simple
36 approximate learning algorithm and its application to high-dimensional plant gene
37 expression data. BMC Syst Biol. 2007;1:37. doi:10.1186/1752-0509-1-37.
- 38 13. Krumsiek J, Suhre K, Illig T, Adamski J, Theis FJ. Gaussian graphical modeling

- 1 reconstructs pathway reactions from high-throughput metabolomics data. *BMC Syst*
2 *Biol.* 2011;5:21. doi:10.1186/1752-0509-5-21.
- 3 14. Gobbi A, Iorio F, Dawson KJ, Wedge DC, Tamborero D, Alexandrov LB, et al. Fast
4 randomization of large genomic datasets while preserving alteration counts.
5 *Bioinformatics.* 2014;30:i617-23. doi:10.1093/bioinformatics/btu474.
- 6 15. Cline MS, Smoot M, Cerami E, Kuchinsky A, Landys N, Workman C, et al. Integration
7 of biological networks and gene expression data using Cytoscape. *Nat Protoc.*
8 2007;2:2366–82. doi:10.1038/nprot.2007.324.
- 9 16. Shin S-Y, Fauman EB, Petersen A-K, Krumsiek J, Santos R, Huang J, et al. An atlas of
10 genetic influences on human blood metabolites. *Nat Genet.* 2014;46:543–50.
11 doi:10.1038/ng.2982.
- 12 17. Chasman DI, Pare G, Mora S, Hopewell JC, Peloso G, Clarke R, et al. Forty-three loci
13 associated with plasma lipoprotein size, concentration, and cholesterol content in
14 genome-wide analysis. *PLoS Genet.* 2009;5:e1000730.
15 doi:10.1371/journal.pgen.1000730.
- 16 18. Rueedi R, Ledda M, Nicholls AW, Salek RM, Marques-Vidal P, Morya E, et al. Genome-
17 wide association study of metabolic traits reveals novel gene-metabolite-disease
18 links. *PLoS Genet.* 2014;10:e1004132. doi:10.1371/journal.pgen.1004132.
- 19 19. Teslovich TM, Kim DS, Yin X, Stančáková A, Jackson AU, Wielscher M, et al.
20 Identification of seven novel loci associated with amino acid levels using single-
21 variant and gene-based tests in 8545 Finnish men from the METSIM study. *Hum Mol*
22 *Genet.* 2018;27:1664–74. doi:10.1093/hmg/ddy067.
- 23 20. Demirkan A, van Duijn CM, Ugocsai P, Isaacs A, Pramstaller PP, Liebisch G, et al.
24 Genome-wide association study identifies novel loci associated with circulating
25 phospho- and sphingolipid concentrations. *PLoS Genet.* 2012;8:e1002490.
26 doi:10.1371/journal.pgen.1002490.
- 27 21. Tabassum R, Rämö JT, Ripatti P, Koskela JT, Kurki M, Karjalainen J, et al. Genetic
28 architecture of human plasma lipidome and its link to cardiovascular disease. *Nat*
29 *Commun.* 2019;10:4329. doi:10.1038/s41467-019-11954-8.
- 30 22. Astle WJ, Elding H, Jiang T, Allen D, Ruklisa D, Mann AL, et al. The allelic landscape of
31 human blood cell trait variation and links to common complex disease. *Cell.*
32 2016;167:1415–29. doi:10.1016/j.cell.2016.10.042.
- 33 23. Sun BB, Maranville JC, Peters JE, Stacey D, Staley JR, Blackshaw J, et al. Genomic
34 atlas of the human plasma proteome. *Nature.* 2018;558:73–9. doi:10.1038/s41586-
35 018-0175-2.
- 36 24. Schaap FG, Rensen PC, Voshol PJ, Vrins C, van der Vliet HN, Chamuleau RA, et al.
37 ApoAV reduces plasma triglycerides by inhibiting very low density lipoprotein-
38 triglyceride (VLDL-TG) production and stimulating lipoprotein lipase-mediated VLDL-

- 1 TG hydrolysis. *J Biol Chem*. 2004;279:27941–7. doi:10.1074/jbc.M403240200.
- 2 25. Ariza MJ, Sanchez-Chaparro MA, Baron FJ, Hornos AM, Calvo-Bonacho E, Rioja J, et
3 al. Additive effects of LPL, APOA5 and APOE variant combinations on triglyceride
4 levels and hypertriglyceridemia: results of the ICARIA genetic sub-study. *BMC Med*
5 *Genet*. 2010;11:66. doi:10.1186/1471-2350-11-66.
- 6 26. Johansen CT, Wang J, Lanktree MB, Cao H, McIntyre AD, Ban MR, et al. Excess of rare
7 variants in genes identified by genome-wide association study of
8 hypertriglyceridemia. *Nat Genet*. 2010;42:684–7. doi:10.1038/ng.628.
- 9 27. Weissglas-Volkov D, Aguilar-Salinas CA, Nikkola E, Deere KA, Cruz-Bautista I,
10 Arellano-Campos O, et al. Genomic study in Mexicans identifies a new locus for
11 triglycerides and refines European lipid loci. *J Med Genet*. 2013;50:298–308.
12 doi:10.1136/jmedgenet-2012-101461.
- 13 28. De Castro-Orós I, Cenarro A, Tejedor MT, Baila-Rueda L, Mateo-Gallego R, Lamiquiz-
14 Moneo I, et al. Common genetic variants contribute to primary hypertriglyceridemia
15 without differences between familial combined hyperlipidemia and isolated
16 hypertriglyceridemia. *Circ Cardiovasc Genet*. 2014;7:814–21.
17 doi:10.1161/CIRCGENETICS.114.000522.
- 18 29. Tang CS, Zhang H, Cheung CY, Xu M, Ho JC, Zhou W, et al. Exome-wide association
19 analysis reveals novel coding sequence variants associated with lipid traits in
20 Chinese. *Nat Commun*. 2015;6:10206. doi:10.1038/ncomms10206.
- 21 30. Global Lipids Genetics Consortium, Willer CJ, Schmidt EM, Sengupta S, Peloso GM,
22 Gustafsson S, et al. Discovery and refinement of loci associated with lipid levels. *Nat*
23 *Genet*. 2013;45:1274–83. doi:10.1038/ng.2797.
- 24 31. Klarin D, Damrauer SM, Cho K, Sun Y V., Teslovich TM, Honerlaw J, et al. Genetics of
25 blood lipids among ~300,000 multi-ethnic participants of the Million Veteran
26 Program. *Nat Genet*. 2018;50:1514–23.
- 27 32. Wang WJ, Baez JM, Maurer R, Dansky HM, Cohen DE. Homozygous disruption of Pctp
28 modulates atherosclerosis in apolipoprotein E-deficient mice. *J Lipid Res*.
29 2006;47:2400–7. doi:10.1194/jlr.M600277-JLR200.
- 30 33. Piperi C, Adamopoulos C, Papavassiliou AG. XBP1: a pivotal transcriptional regulator
31 of glucose and lipid metabolism. *Trends Endocrinol Metab*. 2016;27:119–22.
32 doi:10.1016/j.tem.2016.01.001.
- 33 34. Glimcher LH, Lee AH. From sugar to fat: how the transcription factor XBP1 regulates
34 hepatic lipogenesis. *Ann N Y Acad Sci*. 2009;1173 Suppl:E2-9. doi:10.1111/j.1749-
35 6632.2009.04956.x.
- 36 35. Brandt C, Nolte H, Henschke S, Engström Ruud L, Awazawa M, Morgan DA, et al. Food
37 perception primes hepatic ER homeostasis via melanocortin-dependent control of
38 mTOR activation. *Cell*. 2018;175:1321-1335.e20. doi:10.1016/J.CELL.2018.10.015.

- 1 36. Macaluso FS, Maida M, Petta S. Genetic background in nonalcoholic fatty liver disease:
2 a comprehensive review. *World J Gastroenterol.* 2015;21:11088.
3 doi:10.3748/wjg.v21.i39.11088.
- 4 37. Simons N, Isaacs A, Koek GH, Kuc S, Schaper NC, Brouwers MC. PNPLA3, TM6SF2,
5 and MBOAT7 genotypes and coronary artery disease. *Gastroenterology.*
6 2017;152:912–3. doi:10.1053/j.gastro.2016.12.020.
- 7 38. Mahajan A, Wessel J, Willems SM, Zhao W, Robertson NR, Chu AY, et al. Refining the
8 accuracy of validated target identification through coding variant fine-mapping in
9 type 2 diabetes. *Nat Genet.* 2018;50:559–71. doi:10.1038/s41588-018-0084-1.
- 10 39. Gijón MA, Riekhof WR, Zarini S, Murphy RC, Voelker DR. Lysophospholipid
11 acyltransferases and arachidonate recycling in human neutrophils. *J Biol Chem.*
12 2008;283:30235–45. doi:10.1074/jbc.M806194200.
- 13 40. Takemasu S, Ito M, Morioka S, Nigorikawa K, Kofuji S, Takasuga S, et al.
14 Lysophosphatidylinositol-acyltransferase-1 is involved in cytosolic Ca²⁺ oscillations
15 in macrophages. *Genes to Cells.* 2019;24:366–76. doi:10.1111/gtc.12681.
- 16 41. Draisma HH, Pool R, Kobl M, Jansen R, Petersen AK, Vaarhorst AA, et al. Genome-wide
17 association study identifies novel genetic variants contributing to variation in blood
18 metabolite levels. *Nat Commun.* 2015;6:7208. doi:10.1038/ncomms8208.
- 19 42. Wu JH, Lemaitre RN, Manichaikul A, Guan W, Tanaka T, Foy M, et al. Genome-wide
20 association study identifies novel loci associated with concentrations of four plasma
21 phospholipid fatty acids in the de novo lipogenesis pathway: results from the Cohorts
22 for Heart and Aging Research in Genomic Epidemiology (CHARGE) consortium. *Circ*
23 *Cardiovasc Genet.* 2013;6:171–83. doi:10.1161/CIRCGENETICS.112.964619.
- 24 43. Hu Y, Tanaka T, Zhu J, Guan W, Wu JHY, Psaty BM, et al. Discovery and fine-mapping
25 of loci associated with MUFAs through trans-ethnic meta-analysis in Chinese and
26 European populations. *J Lipid Res.* 2017;58:974–81. doi:10.1194/jlr.P071860.
- 27 44. Burgess S, Harshfield E. Mendelian randomization to assess causal effects of blood
28 lipids on coronary heart disease: lessons from the past and applications to the future.
29 *Curr Opin Endocrinol Diabetes Obes.* 2016;23:124–30.
30 doi:10.1097/MED.000000000000230.
- 31 45. Saleheen D, Natarajan P, Armean IM, Zhao W, Rasheed A, Khetarpal SA, et al. Human
32 knockouts and phenotypic analysis in a cohort with a high rate of consanguinity.
33 *Nature.* 2017;544:235–9. doi:10.1038/nature22034.
- 34 46. Di Angelantonio E, Thompson S, Kaptoge S, Moore C, Walker M, Armitage J, et al.
35 Efficiency and safety of varying the frequency of whole blood donation (INTERVAL):
36 a randomised trial of 45 000 donors. *Lancet.* 2017;390:2360–71.
37 doi:10.1016/S0140-6736(17)31928-1.

38

1 **FIGURES**

2

3 Figure 1. Miami plot of combined association results from genome-wide association
4 analysis for all lipids in PROMIS and INTERVAL

5 Figure 2. Heat map showing associations of significant loci from conditional analyses
6 with selected lipid metabolites in PROMIS

7 Figure 3. Combined network graph summarising genetic associations and a Gaussian
8 graphical model (GGM) relating to levels of individual lipid species in
9 PROMIS

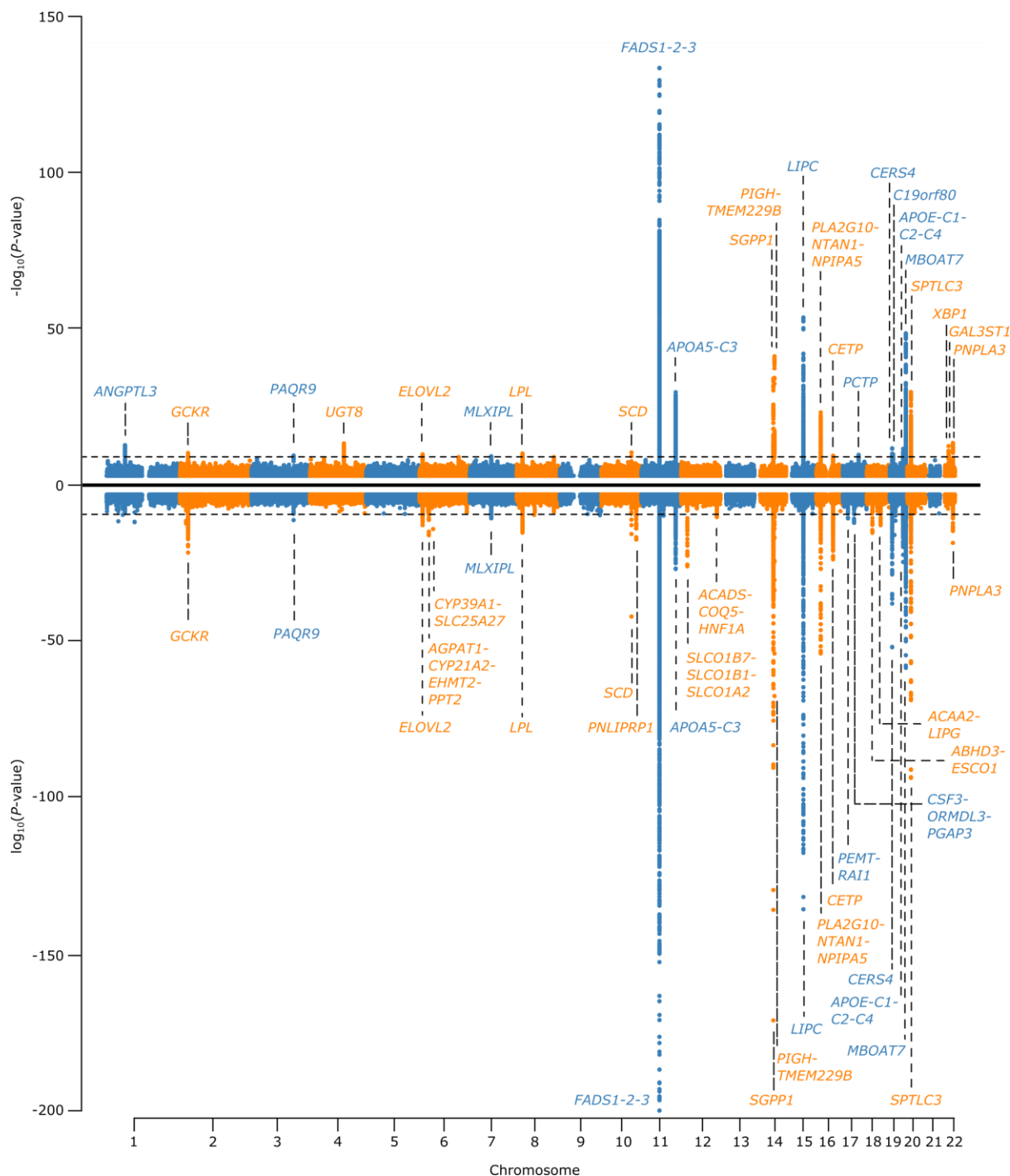
10 Figure 4. Combined network graph summarising genetic associations and a Gaussian
11 graphical model (GGM) relating to levels of individual triglycerides in
12 PROMIS

13 Figure 5. Association of lipids in PROMIS with *PNPLA3* and differences in levels of
14 triglycerides by genotype

15

16

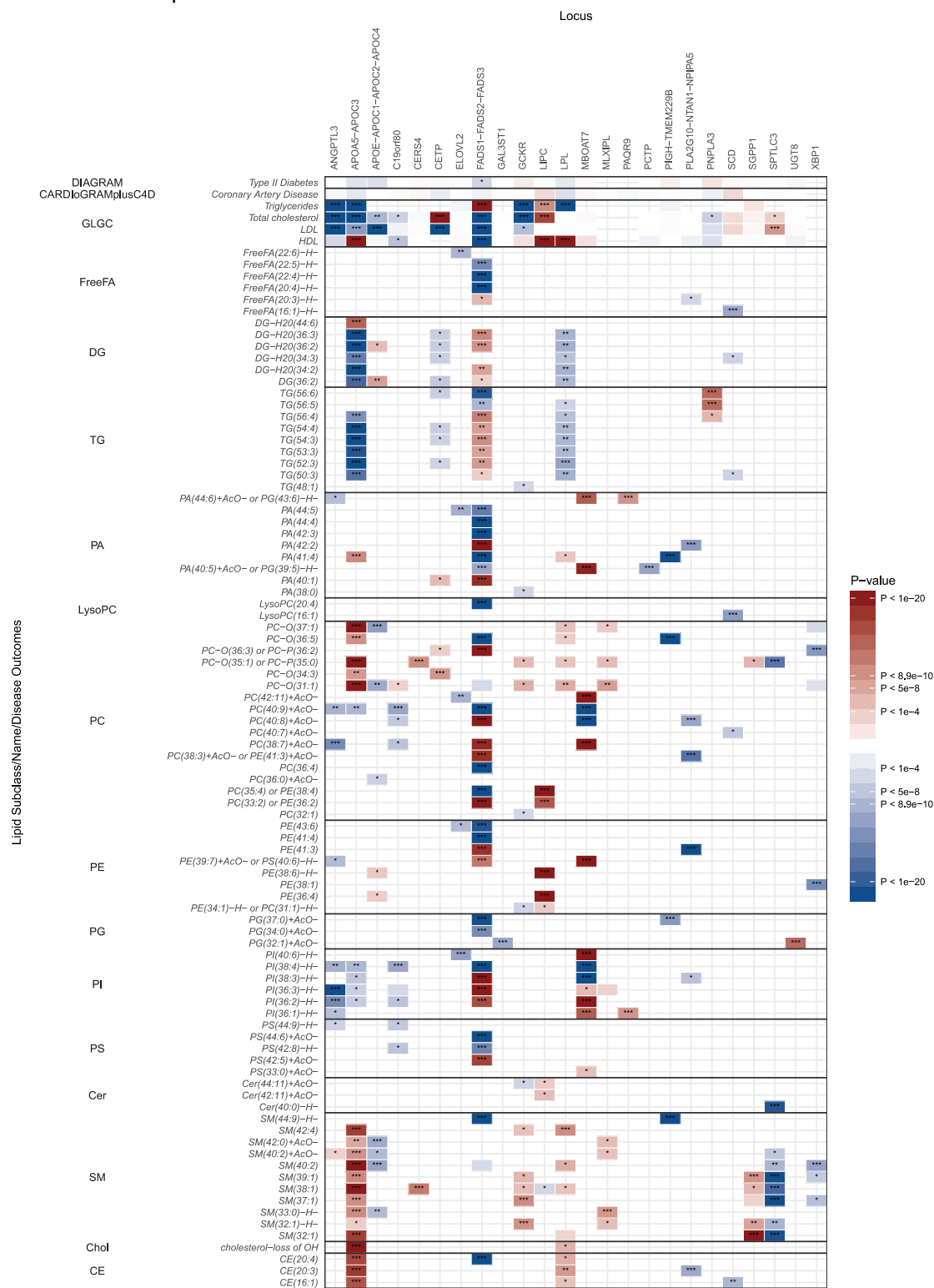
1 **Figure 1.** Miami plot of combined association results from genome-wide association
 2 analysis for all lipids in PROMIS and INTERVAL
 3



4
 5
 6 The combined association results are shown for all lipids with each variant in PROMIS (top)
 7 and INTERVAL (bottom). P -values $> 1 \times 10^{-3}$ have been truncated at 1×10^{-3} , and P -values
 8 $< 1 \times 10^{-200}$ have been truncated at 1×10^{-200} . Actual P -value for lead SNP in *FADS1-2-3*
 9 locus in INTERVAL is 1.6×10^{-286} .

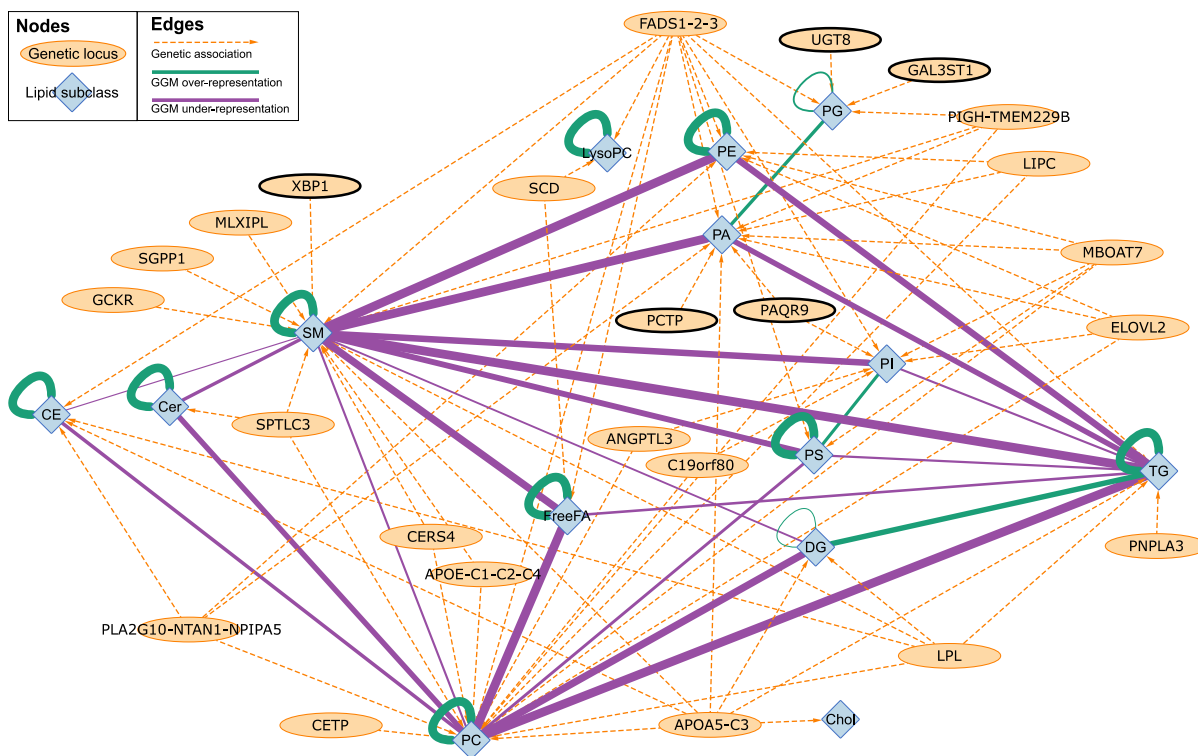
10
 11

1 **Figure 2.** Heat map showing associations of significant loci from conditional analyses
 2 with selected lipid metabolites in PROMIS



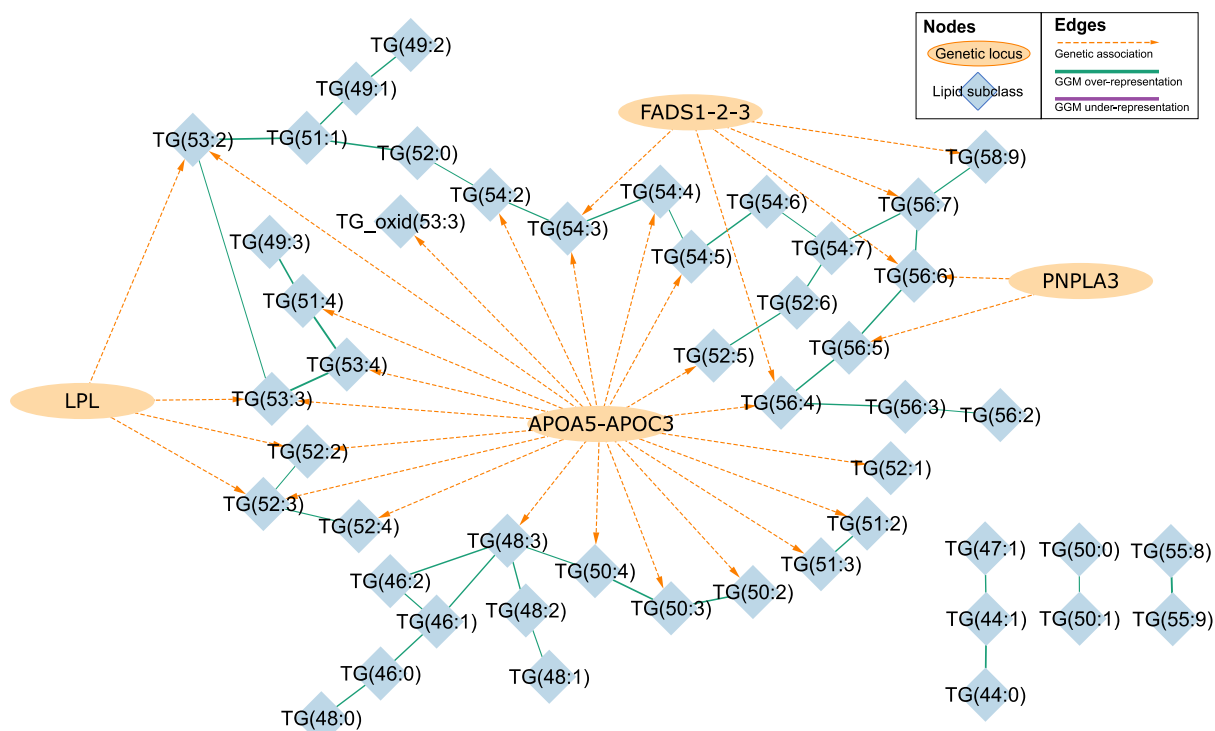
3 The effect estimates of the associations between significant variants and selected lipids are plotted as a heat
 4 map. Results are shown for selected top lipids with the strongest associations within each subclass (rows)
 5 against the most strongly associated genetic variant within each locus (columns). The associations with major lipids from
 6 the GLGC (total cholesterol, HDL-C, LDL-C, and triglycerides), DIAGRAM Consortium (type 2 diabetes), and
 7 CARDIoGRAMplusC4D Consortium (coronary artery disease) are also shown. The magnitude and direction of the
 8 effect estimates (standardised per 1-SD) are indicated by a colour scale, with blue indicating a negative
 9 association and red indicating a positive association with respect to the SNP effect on the trait. Asterisks indicate
 10 the degree of significance of the P-values of association. * = $P < 1 \times 10^{-4}$; ** = $P < 5 \times 10^{-8}$; *** = $P < 8.9 \times$
 11 10^{-10} .
 12

1 **Figure 3.** Combined network graph summarising genetic associations and a
 2 Gaussian graphical model (GGM) relating to levels of individual lipid
 3 species in PROMIS
 4



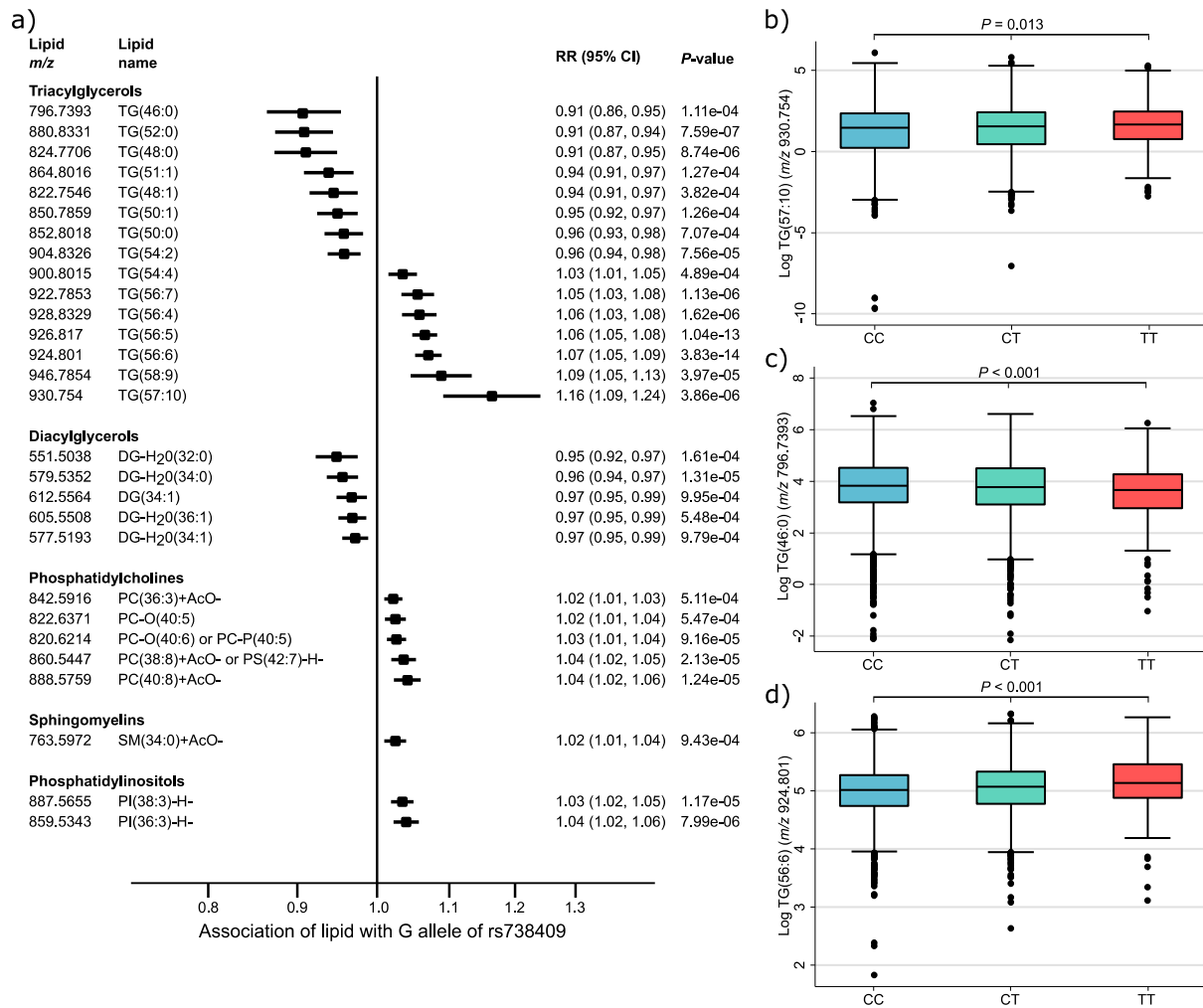
5
 6
 7 Nodes representing genetic loci are each labelled with the most likely “causal” gene at that
 8 locus according to our functional annotation (see Methods). In order for an edge to be
 9 drawn between a genetic locus and a lipid subclass, there must have been a minimum of
 10 one variant at that locus significantly ($P < 8.9 \times 10^{-10}$) associated with a minimum of one
 11 lipid species belonging to that lipid subclass. Edges between lipid subclasses indicate
 12 whether there was either a significant over- (green) or under- (purple) representation (the
 13 magnitude is indicated in the thickness of the edges) of GGM connections between lipid
 14 species belonging to different lipid subclasses.
 15
 16
 17

1 **Figure 4.** Combined network graph summarising genetic associations and a
 2 Gaussian graphical model (GGM) relating to levels of individual
 3 triglycerides in PROMIS
 4
 5



6
 7
 8 Nodes representing genetic loci are each labelled with the most likely “causal” gene at that
 9 locus according to our functional annotation (see Methods). In order for an edge to be
 10 drawn between a genetic locus and a triglyceride, there must have been a minimum of
 11 one variant at that locus significantly ($P < 8.9 \times 10^{-10}$) associated with at least one
 12 triglyceride. Edges between triglycerides indicate whether there was either a significant
 13 over- (green) or under- (purple) representation, with the magnitude indicated by the
 14 thickness of the edges.
 15

1 **Figure 5.** Association of lipids in PROMIS with *PNPLA3* and differences in levels of
 2 triglycerides by genotype
 3



4
 5
 6 (a) Association of G allele of rs738409 in *PNPLA3* locus with levels of various lipids in PROMIS. The
 7 black lines denote 95% confidence intervals. Difference in levels of triglycerides in PROMIS by
 8 genotype: (b) TG(57:10) (*m/z* 930.754), (c) TG(46:0) (*m/z* 796.7393), and (d) TG(56:6) (*m/z*
 9 924.801). *P*-values are for ANOVA test of difference in mean levels of triglycerides by genotype.
 10