



OPEN

## Genetic determinants of plasma protein levels in the Estonian population

Anette Kalnapenkis<sup>1,2,7</sup>✉, Maarja Jöeloo<sup>1,2,7</sup>, Kaido Lepik<sup>1,3,4,5</sup>, Viktorija Kukuškina<sup>1</sup>, Mart Kals<sup>1</sup>, Kaur Alasoo<sup>6</sup>, Estonian Biobank Research Team\*, Reedik Mägi<sup>1</sup>, Tõnu Esko<sup>1,7</sup>✉ & Urmo Vösa<sup>1,7</sup>✉

The proteome holds great potential as an intermediate layer between the genome and phenotype. Previous protein quantitative trait locus studies have focused mainly on describing the effects of common genetic variations on the proteome. Here, we assessed the impact of the common and rare genetic variations as well as the copy number variants (CNVs) on 326 plasma proteins measured in up to 500 individuals. We identified 184 *cis* and 94 *trans* signals for 157 protein traits, which were further fine-mapped to credible sets for 101 *cis* and 87 *trans* signals for 151 proteins. Rare genetic variation contributed to the levels of 7 proteins, with 5 *cis* and 14 *trans* associations. CNVs were associated with the levels of 11 proteins (7 *cis* and 5 *trans*), examples including a 3q12.1 deletion acting as a hub for multiple *trans* associations; and a CNV overlapping *NAIP*, a sensor component of the NAIP-NLRC4 inflammasome which is affecting pro-inflammatory cytokine interleukin 18 levels. In summary, this work presents a comprehensive resource of genetic variation affecting the plasma protein levels and provides the interpretation of identified effects.

During the last decade, genome-wide association studies (GWASs) have successfully linked genetic variants to complex traits<sup>1</sup>. However, the mechanisms underlying many of these associations often remain unknown, as most of the associated genetic variants are located in non-coding regions of the genome, suggesting that they have regulatory effects on phenotypes<sup>2</sup>. To fill this knowledge gap, molecular traits are routinely used as intermediate phenotypes in association studies. The study of molecular phenotypes enables the assessment of the direct effects of genetic variants on, for example, the alteration of protein levels, and the potential underlying molecular mechanisms and links to endpoint phenotypes.

Proteins are functional products of the genome that provide insight about the normal processes of organisms; in addition, alterations in their levels are indicators of changes in disease status<sup>3</sup>. Recent technological advancements, including the development of multiplex immunoassays and aptamer assays, have provided opportunities for the measurement of thousands of plasma- and serum-based protein levels<sup>4–8</sup>.

The genetic backgrounds of protein levels are uncovered through the linking of these levels to genetic variability via protein quantitative trait locus (pQTL) analysis. Many recent pQTL studies have been large-scale<sup>4–8</sup>, with the largest of them including 54,306 individuals from the UK Biobank<sup>9</sup>. Their primary focus has been the identification of common [minor allele frequency (MAF) > 0.01] variants affecting inter-individual protein variability, but Sun et al.<sup>9</sup> reported that approximately 5.6% (570/10,248) and 1.5% (155/10,248) of the variants with primary associations had MAFs < 0.01 and < 0.005, respectively. In addition, the focus has been shifting toward the identification of associations with rare (MAF < 0.01) variants, using gene-based methods<sup>10–14</sup>. For example, a recent landmark study conducted on the Icelandic population revealed 18,084 genetic associations with protein levels, 19% of which were with rare variants<sup>8</sup>. Investigation of the effects of other structural variants, such as copy number variants (CNVs), on protein levels has thus far been limited<sup>15</sup>.

The combined examination of pQTL and GWAS results for disease phenotypes can lead to the validation and prioritisation of new and existing drug targets, and the identification of clinically relevant biomarkers.

<sup>1</sup>Estonian Genome Centre, Institute of Genomics, University of Tartu, Tartu, Estonia. <sup>2</sup>Institute of Molecular and Cell Biology, University of Tartu, Tartu, Estonia. <sup>3</sup>Department of Computational Biology, University of Lausanne, Lausanne, Switzerland. <sup>4</sup>Swiss Institute of Bioinformatics, Lausanne, Switzerland. <sup>5</sup>University Center for Primary Care and Public Health, Lausanne, Switzerland. <sup>6</sup>Institute of Computer Science, University of Tartu, Tartu, Estonia. <sup>7</sup>These authors contributed equally: Anette Kalnapenkis, Maarja Jöeloo, Tõnu Esko and Urmo Vösa. \*A list of authors and their affiliations appears at the end of the paper. ✉email: anette.kalnapenkis@ut.ee; tonu.esko@ut.ee; urmo.vosa@ut.ee

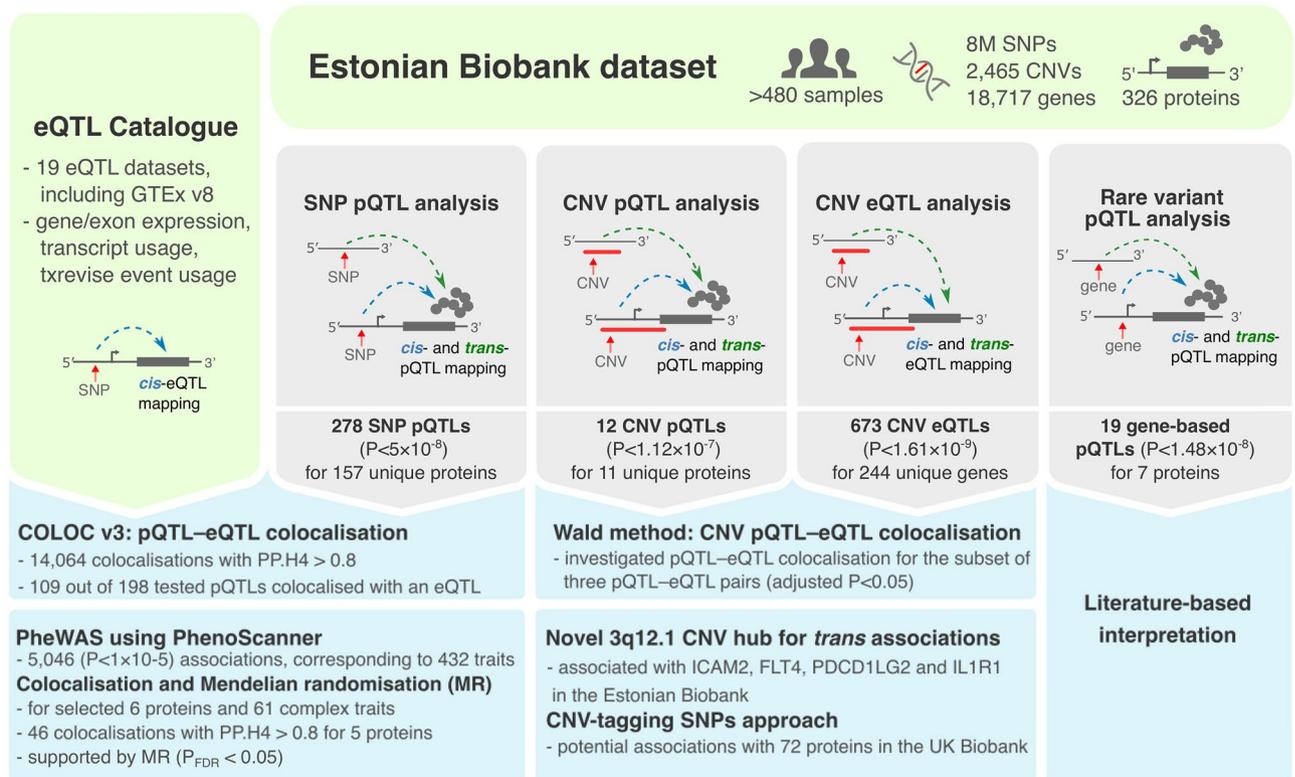
Ferkingstad et al.<sup>8</sup> found that 12% of 45,334 lead associations in the GWAS Catalog were with variants in high linkage disequilibrium (LD) with pQTLs. The application of Mendelian randomisation (MR) and colocalisation analysis to biomedical data for the identification of links between pQTLs and diseases enables the evaluation of the causality between protein levels and disease risk and the identification of potential disease pathways, respectively. Zheng et al.<sup>16</sup> used MR and colocalisation analysis to examine associations of 1002 plasma proteins with 153 diseases and 72 disease-related risk factors, and identified 413 protein–trait associations supported by MR, 130 (31.5%) of which were not supported by the colocalisation analysis. This example highlights the importance of intersecting the results from both analyses<sup>17</sup>.

Here, we integrated dense whole-genome sequencing (WGS) data to study the genetic contributions of rare and common variants to 326 plasma protein levels in the Estonian Biobank cohort (Fig. 1). We examined the effects of single nucleotide polymorphisms (SNPs) and common CNVs on the inter-individual protein variability, and identified several proteins that were affected by the latter. To assess the overlap of local (*cis*) and distal (*trans*) pQTL effects with gene expression levels, we conducted comprehensive colocalisation analyses with expression quantitative trait loci (eQTLs) and splicing QTLs using data from various tissues from the eQTL Catalogue<sup>18</sup>.

## Material and methods

### Study samples

The Estonian Biobank (EstBB) cohort consists of more than 200,000 Estonian volunteers aged  $\geq 18$  years, representing about 20% of the Estonian adult population, detailed information on the enrolment process and data collection is described in the Leitsalu et al. study<sup>19</sup>. Genotype data are available for all gene donors in this cohort. For a subcohort of 500 individuals [52.8% females and 47.2% males, mean age 54 (standard deviation 14.0) years], WGS, RNA sequencing and Olink proteomics data from the same timepoint are available. The WGS dataset was previously generated in 2015. Sample collection for RNA sequencing and Olink proteomics was conducted previously in years 2011–2012. RNA sequencing was previously performed in years 2015–2016 and protein levels were previously measured in year 2017. All the analyses in the current study were conducted using already existing data and no new data were collected during this study. The data were first accessed in year 2017. The research activities of the EstBB are regulated by § 6 section 1, § 16 section 1 and § 22 section 1 of the Human Genes Research Act (<https://www.riigiteataja.ee/en/eli/ee/531102013003/consolide/current>), which was adopted in year 2000 specifically for the operations of the EstBB. During enrolment, all the participants have signed a broad informed consent form which allows researchers to use their genomics and health data for scientific studies upon approval by the Estonian Committee on Bioethics and Human Research and previously by the Research Ethics Committee of the University of Tartu. Individual level data analysis for this project was carried out under approvals 234/T-12 from the Research Ethics Committee of the University of Tartu and 1.1-12/624 from the Estonian Committee on Bioethics and Human Research (Estonian Ministry of Social Affairs)



**Figure 1.** Overview of the main analyses conducted in this study.

and data extraction no. K29 from the Estonian Biobank. The current study was conducted using pseudonymised data. All methods were carried out in accordance with relevant guidelines and regulations.

### WGS data processing, variant calling and quality control

The 2284 EstBB WGS samples were sequenced at the Genomics Platform of the Broad Institute (Cambridge, MA, USA). Sequenced data were jointly variant called and quality controlled as described by Mitt et al.<sup>20</sup>; and the final WGS sample set was derived from 2244 individuals. We excluded multiallelic sites and genetic variants, based on quality/depth < 2, Hardy–Weinberg equilibrium test failure ( $P > 1 \times 10^{-9}$ ), and call rate < 90%. Data from individuals with available proteomics data ( $n = 500$ ) were retained for further analyses.

### CNV detection and quality control

The Genome STRiP pipeline (version 2.00.1611)<sup>21</sup> was applied to detect CNVs from aligned sequencing reads (in BAM format) for 2284 samples as described by Lepamets et al.<sup>22</sup>. In brief, CNV sites were identified and genotyped in five batches. After the exclusion of samples with excessive numbers of calls, the batches were combined and duplicate calls were merged. Low-quality calls and sites with call rates < 90% were excluded. We restricted the final dataset to deletions longer than 1000 bp and duplications longer than 2000 bp. The final sample set contained 51,026 CNV sites from 2230 individuals. Data from individuals with available proteomics data ( $n = 500$ ) were retained for further analyses.

### Measurement of plasma protein levels

Plasma concentrations in EDTA plasma samples from 500 Estonian Biobank donors were measured using four arrays with 92 protein targets each [ProSeek Cardiovascular Disease (CVD) II and III, Inflammation and Oncology II; Olink Biosciences, Uppsala, Sweden; Supplementary Table S1]. The procedure is described in detail elsewhere<sup>23</sup>, and a technical white paper with additional information is available at the manufacturer's website (<https://www.olink.com>). The native Olink data consisted of qPCR cycle threshold values corrected for extension control, followed by inter-plate control and the application of a correction factor predetermined by a negative control signal. The measurements were given at a natural logarithmic scale as normalised protein expression levels, a relative quantification unit. As part of the quality control, we excluded individual samples that did not pass the Olink internal quality control system. Final sample sizes per array ranged from 488 to 497, and the samples were measured in six batches. For arrays in which < 20% of samples had values below the limit of detection (LOD), protein level correction was performed by dividing the Olink-assigned LOD value by 2, as done in the SCALLOP CVD-I project<sup>6</sup>. A total of 341 protein traits (326 unique proteins; 13 proteins were measured by two arrays and one protein was measured by three arrays) passed quality control and were retained for further analyses (Supplementary Table S1).

### RNA sequencing data

RNA was extracted from samples in thawed Tempus tubes using TRIzol reagent (Invitrogen, Waltham, MA, USA) and further purified using an RNeasy Mini Kit (Qiagen, Hilden, Germany). Globin mRNA was depleted using GLOBINclear Kit (Invitrogen, Waltham, MA, USA). RNA quality was checked using electrophoresis (Agilent 2200 TapeStation; Agilent Technologies, Santa Clara, CA, USA). Sequencing libraries were prepared using 200 ng RNA according to the Illumina TruSeq stranded mRNA protocol. RNA sequencing was performed at the Estonian Genome Centre Core Facility using paired-end 50-bp sequencing technology (Illumina, San Diego, CA, USA), according to the manufacturer's specifications.

Adapters and leading and trailing bases with a quality score were removed using Trimmomatic (version 0.36)<sup>24</sup>. Quality control was done with FastQC (version 0.11.2)<sup>25</sup>. Reads were mapped to human genome reference version GRCh37.p13 with STAR (version 2.4.2a)<sup>26</sup>. Reads that mapped to each genomic feature were counted with STAR using the same algorithm as default htseq-count. Raw RNA sequencing counts were normalised with the weighted trimmed mean of M-values<sup>27</sup> method from the edgeR R package (version 3.12.1)<sup>28</sup>. Detailed information regarding RNA sequencing data pre-processing is described in Lepik et al.<sup>29</sup>. The final gene expression measure was in logarithmed count per million. In total, 486 RNA sequencing samples overlapped with available proteomics data and were used for eQTL mapping.

### Genome-wide SNP pQTL discovery

Protein trait levels were rank-based inverse normal transformed. We regressed out the effects of age, sex, the season of sample collection, smoking status, blood sample processing time (days), plasma sample storage time (in days) and protein analysis batch using a custom R script. The residuals were used in a single-variant pQTL analysis performed with the EMMAX linear mixed model<sup>30</sup> and the EPACTS software (version 3.3.0, *q.emmax* function; <https://genome.sph.umich.edu/wiki/EPACTS>). To account for population structure, a kinship matrix was generated in EPACTS using genetic variants with MAF > 0.01 and call rate > 95%. Depending on the panel, we tested between 8,856,032 and 8,891,303 autosomal genetic variants against 341 plasma protein traits.

We classified associated variants into two categories based on their positions in relation to the protein-coding genes. We defined *cis*-pQTLs as SNPs located within 1 Mb upstream or downstream of the transcription start sites (TSSs) of the corresponding protein-coding genes, and *trans*-pQTLs as SNPs located > 1 Mb upstream or downstream of the TSS or on a different chromosome. Heterodimers were classified based on the protein subunit gene closest to the associated variant. In the case of proteins that were present on multiple panels, weaker signals were omitted from the analyses.

To retain independent signals, associated variants were clumped in PLINK (version 1.9)<sup>31</sup>, using a 1 Mb window with the LD thresholds of  $R^2 = 0.1$  and  $P < 5 \times 10^{-8}$ . To flag potential 'pseudo-pQTL' signals caused by the

epitope effect, i.e. altered assay binding affinity due to a change in protein structure instead of an actual change in protein expression level, we followed the strategy described by Folkersen et al.<sup>6</sup>. Briefly, we determined whether any lead *cis* variant was a protein-altering variant (PAV) or in high LD ( $R^2 \geq 0.8$ ) with one, by using 2230 WGS samples as the reference for the LD calculations (Supplementary Table S2). Missense, frameshift, splice donor region and stop gain variants were flagged as PAVs. Lead pQTL variants were queried for evidence of location in a regulatory region using RegulomeDB<sup>32</sup>.

### Power calculation

For the common variant pQTL analysis power calculations with  $n = 500$ , we assumed linear regression with an additive genetic model with no gene–gene interactions using *genpwr* (version 1.0.4) R package<sup>33</sup>, following the strategy previously used by Yao et al.<sup>34</sup>. As each protein in the analysis was inverse normal transformed, we set the standard deviation = 1 and used  $\alpha = 5 \times 10^{-8}$  as well as  $\alpha = 2.76 \times 10^{-10}$ . To construct power curves, the power was estimated for pQTLs with effect sizes of 0.25, 0.5, 1, 1.25 and 1.5 with MAF values ranging from 0.001 to 0.5. The same package was also used to estimate the size of the pQTL effect required to detect the pQTL with 80% power.

### Corresponding eQTL discovery

In order to overlap the genome-wide significant ( $P < 5 \times 10^{-8}$ ) pQTLs with eQTLs, we used the RNA sequencing data from the overlapping samples of the same cohort<sup>29</sup>. We tested the eQTL effects on the genes encoding corresponding proteins by using a linear mixed model from EPACTS software (version 3.2.2)<sup>30</sup> with the same settings as for pQTL analysis. We included age, sex, body mass index, blood components (neutrophils, eosinophils, basophils, lymphocytes, monocytes, erythrocytes and thrombocytes) and RNA sequencing batch as covariates. To account for hidden batch effects on the gene expression, the first two principal components of the gene expression data were also included as covariates, as described in detail in Lepik et al.<sup>29</sup>. To correct for multiple testing, we adjusted *P*-values using false discovery rate (FDR) correction; eQTLs were considered as replicated at Benjamini–Hochberg  $FDR \leq 0.05$  and with concordant allelic direction with the pQTLs.

### Multiple testing correction for the pQTL analysis

From primary analyses, effects reaching per-protein genome-wide significance ( $P < 5 \times 10^{-8}$ ) were interpreted. To also provide the more conservative results accounting for the number of tested proteins, we used a strategy described by Gao et al.<sup>35</sup> and Kettunen et al.<sup>36</sup>, which accommodates the correlation between protein levels. Four matrices corresponding to inverse normal transformed and covariate-adjusted protein levels from the Olink panels were merged. Only samples that passed quality control on every panel ( $n = 478$ ) were included. The resulting matrix of standardised residuals was used in a principal components analysis implemented with the FactoMiner (version 1.41)<sup>37</sup> R package. As 181 principal components cumulatively explained > 95% of the total variance in the proteomics data, the stricter significance threshold was set to  $2.76 \times 10^{-10}$  ( $5 \times 10^{-8}/181$ ).

### Gene-based analysis of rare SNPs

Variants were annotated using the EPACTS ‘anno’ module (version 3.3.0; <https://genome.sph.umich.edu/wiki/EPACTS>) and GENCODE (version 14)<sup>38</sup> to ascertain their effects on protein sequences. A gene-based group file was generated with the inclusion of all nonsynonymous (missense and nonsense) variants in assigned genes. Only genes with more than two nonsynonymous variants were retained. We performed the gene-based SKAT test using the EMMAX *mmskat* function with adjustment for small sample size in EPACTS, using all variants with  $0.000001\% < \text{MAF} < 1\%$ . Covariates included in the rare variant pQTL analysis were the same as described in the Methods section for Genome-wide SNP pQTL discovery. The results were corrected for multiple testing based on Bonferroni-corrected threshold of  $P < 1.48 \times 10^{-8}$  [ $0.05 / (18,717 \text{ genes} \times 181 \text{ protein traits})$ ]. Associations between genes and levels of proteins encoded on the same gene were classified as *cis*, and all other associations were classified as *trans*. Using the GeneMANIA<sup>39,40</sup>, STRING (version 12.0)<sup>41</sup> and BioGRID (version 4.4.230)<sup>42,43</sup> databases, we investigated whether the associated genes also had gene–gene functional interactions with corresponding protein-coding genes. For overlapping the rare variant pQTL associations with eQTL data, we performed an eQTL mapping with EPACTS software (version 3.2.2) using the same gene-based SKAT test as in rare variant pQTL mapping. Covariates included in the rare variant eQTL analysis were the same as described in the Methods section for Corresponding eQTL discovery. Similar to single variant eQTL analysis, to account for multiple testing, we adjusted *P*-values using false discovery rate (FDR) correction; rare variant eQTLs were considered as replicated at Benjamini–Hochberg  $FDR \leq 0.05$  and directionally concordant with the rare variant pQTLs.

### Fine-mapping analysis

We conducted a fine-mapping analysis to pinpoint causal variants for protein level–significant ( $P < 5 \times 10^{-8}$ ) SNV–pQTL associations. We excluded the LTA and MICA–MICB proteins associated with variants in the major histocompatibility complex region on chromosome 6, due to the complexity of the associated *HLA* region. The fine-mapping procedure was based on the SuSiE ‘sum of single effects’ model<sup>44,45</sup> and was implemented using the *susie\_suff\_stat* function from *susieR* package (version 0.11.42). Fine-mapping pipeline <https://github.com/urmovosa/EstBBfinemap> was implemented in Nextflow<sup>46</sup> and some scripts were modified from the FINNGEN fine-mapping pipeline (<https://github.com/FINNGEN/finemapping-pipeline>). The SuSiE output contains single effect components, i.e., credible sets (CSs), with a > 95% probability of including a variant with a non-zero causal effect. We used a default setting of 10 for the maximum number of causal variants regulating a protein, because Wang et al. has demonstrated it to be the optimal choice for the number of causal variants<sup>44</sup>. LDstore (version 2)<sup>47</sup> was used to generate an LD matrix for each locus.

## Replication of pQTLs

All significant lead variants from the pQTL discovery analyses were queried for previously published associations with protein levels in the PhenoScanner database (version 2)<sup>48,49</sup> using the Python application (<https://github.com/phenoscanner/phenoscannerpy>, query date 4 October 2021). This database contains results from large pQTL studies<sup>4,50,51</sup>. For variant matching between datasets, we created variant names that were concatenations of the corresponding chromosome, chromosome position (hg19), and alphabetically ordered alleles. To match UniProt IDs from the discovery analyses to PhenoScanner trait names, the IDs were converted to recommended HUGO Gene Nomenclature Committee gene names using the UniProt conversion tool (<https://www.uniprot.org/uploadlists/>, latest query date 11 October 2021). We performed additional replication analysis using Pietzner et al. dataset by querying their publicly available results with  $P < 0.05$ <sup>7</sup>. The largest pQTL meta-analysis published to date ( $n = 30,931$ )<sup>6</sup> was conducted through the SCALLOP consortium and was not usable due to sample overlap with the current study. In order to ensure that each protein was represented by a single association, we restricted our comparisons to instances where either one subunit or the entire heterodimer complex was available. For instances where one protein was available multiple times, we conducted comparison with the association with the smallest  $P$ -value. To account for multiple testing, we adjusted  $P$ -values using false discovery rate (FDR) correction; pQTLs were considered as replicated at Benjamini–Hochberg  $FDR \leq 0.05$  and concordant allelic direction with the discovery pQTLs.

## Identification of relevant disease traits and molecular QTLs

To identify complex traits and diseases associated with the top pQTLs, we conducted a phenome-wide association analysis (PheWAS) by querying the lead variants from primary pQTL mapping and their proxies against the PhenoScanner database (version 2)<sup>48,49</sup>. Duplicate associations happening due to data resource overlap were removed. We considered only PhenoScanner associations with  $P < 1 \times 10^{-5}$ . Specifically, we sought to identify pQTLs associated with disease traits, methylation quantitative trait loci (meQTLs), histone modifications and metabolite quantitative trait loci (mQTLs), as well as percent-spliced-in (PSI) associations. We also searched for significant protein genes on a druggable genome list<sup>52</sup> and the drugs that interact with them<sup>53</sup>. For a subset of pQTLs we selected for in-depth analyses by coloc and Mendelian randomisation, an additional PheWAS was conducted with the Medical Research Council (MRC) Integrative Epidemiology Unit (IEU) OpenGWAS database<sup>54</sup>. This was done to extract region-wide associations, irrespective of association  $P$ -value.

## Colocalisation analysis

The colocalisation analyses between pQTLs and eQTLs, as well as between pQTLs and complex traits were carried out using coloc (version 3.2.1) R package<sup>55</sup>, which assumes that each locus has a single causal variant. Priors used for the colocalisation analysis were  $P_1 = 10^{-4}$ ,  $P_2 = 10^{-4}$  and  $P_{12} = 5 \times 10^{-6}$ , as suggested by Wallace et al.<sup>56</sup>. For each protein-level genome-wide-significant ( $P < 5 \times 10^{-8}$ ) pQTL locus, we extracted regions in a 1-Mb radius of its primary lead variant to test for colocalisation. The results were considered significant when the posterior probability for colocalisation ( $PP_4$ ) exceeded 0.8.

In an pQTL–eQTL colocalisation analysis, we compared our significant pQTL loci to all eQTL Catalogue datasets<sup>18</sup>, excluding those of Kasela et al.<sup>57</sup> and Lepik et al.<sup>29</sup> due to sample overlap, containing gene expression, exon expression, transcript usage and txrevise event usage data, and GTEx (version 8)<sup>58</sup> datasets containing gene expression data (<https://www.ebi.ac.uk/eqtl/Methods/>; Supplementary Table S3). We lifted the pQTL summary statistics over to an hg38 build to match with the eQTL Catalogue.

The region-wide associations for GWAS traits enrolled into the colocalisation analyses were extracted from the MRC IEU OpenGWAS database and were examined using the ieugwasr (version 0.1.5) R package (<https://github.com/MRCIEU/ieugwasr>; Supplementary Table S4). Since proteins were selected based on associated traits from the PheWAS, they were all associated with clinical traits (i.e. drugs, surgeries, diseases/conditions). In addition, all selected proteins except IL6R had primary pQTLs that did not include nonsynonymous variants, to minimise the possibility of association due to the epitope effect. IL6R was selected because it has been widely reported by previous pQTL studies as an example of the successful linking of molecular traits and diseases to discover drug targets<sup>50,59</sup>. The input data consisted of region-based summary statistics for six protein traits and 61 complex clinical traits.

## Two-sample MR

We conducted a two-sample MR analysis using protein levels with significant colocalisation ( $PP_4 \geq 0.8$ ) as exposures and complex traits as outcomes, using the TwoSampleMR (version 0.5.6) R package<sup>60,61</sup>. Independent variants obtained previously by clumping served as instrumental variables. We conducted the analysis using an inverse variance weighted fixed-effects method and a single instrument-based Wald ratio test. To correct for multiple testing, we adjusted  $P$ -values using false discovery rate (FDR) correction; results were considered significant at Benjamini–Hochberg  $FDR \leq 0.05$ .

## CNV pQTLs, eQTLs and colocalisation

To determine whether any of the examined proteins are genetically regulated by larger structural variants, we conducted a pQTL mapping using CNV data. Description of the used CNV data is in the Methods section for CNV detection and quality control. Associations between previously described standardised protein measure residuals and CNV sites were assessed by using the MatrxieQTL R package<sup>62</sup>. The post-quality control sample sizes for the Inflammation, Oncology II, CVD II and CVD III panels were 481, 480, 489 and 488 unrelated ( $PI\_HAT < 0.2$ ) individuals, respectively. To discard rare CNV events, all CNV sites with in-sample frequencies of the most frequent copy number  $> 0.95$  were excluded. Additionally, unique non-overlapping CNVs were

included. The final set used in the pQTL analyses comprised of 2465 CNV sites [1375 deletions ( $CN < 2$ ), 482 duplications ( $CN > 2$ ) and 608 combined deletions and duplications]. The genome-wide significance threshold was set to  $1.12 \times 10^{-7}$  ( $0.05/2465/181$ ).

For each significantly associated CNV, all SNP markers within a 500-kbp proximity were tested for potential tagging effects. For this purpose, the SNP pQTL analysis using EPACTS was repeated for these regions with the CNVs included as covariates.

The same CNVs were tested against the expression levels of 12,619 genes<sup>29</sup>, and the CNV pQTL results were then cross-referenced with eQTLs identified from the same set of individuals. The eQTL results were corrected for multiple testing and a Bonferroni-corrected threshold of  $P < 1.61 \times 10^{-9}$  [ $0.05/(2465 \text{ CNVs} \times 12,619 \text{ genes})$ ] was applied. Overlapping eQTL–pQTL pairings were tested in an MR framework using the summary statistics–based ratio estimate (Wald test)<sup>63</sup>, and Spearman's rank correlation coefficient was calculated for gene expression vs protein expression in the same individuals. We hypothesised that CNVs in gene regions would be considerably more likely than other causal variants to modulate the expression of those genes; thus, non-zero ratio estimates were taken to indicate shared causal CNVs of gene expression and protein traits.

### PheWAS of CNV pQTLs

CNV pQTLs from primary mapping that reached genome-wide significance ( $P < 1.12 \times 10^{-7}$ ) or the suggestive significance threshold ( $P < 2 \times 10^{-5}$ ) were included in a PheWAS, resulting in the inclusion of 38 CNV regions. All data included in the PheWAS were obtained using the *lm* function with custom R scripts from 2115 unrelated Estonian Biobank samples for which WGS data were available, and were corrected for age, sex and six genotype principal components (PCs; calculated from common SNPs). The 744 phenotypes examined were anthropometric traits (height, weight, body mass index, hip circumference, waist circumference, waist–hip circumferences ratio), cell counts from RNA-sequencing data (white blood cells, red blood cells, platelets, neutrophils, monocytes, lymphocytes, eosinophils, basophils), nuclear magnetic resonance spectroscopy–detected metabolites ( $n = 225$ ) and International Classification of Diseases, 10th revision (ICD-10) diagnoses with at least 20 carriers in the sample ( $n = 505$ ). Self-reported diagnoses not reported elsewhere were set to not available. Sex-specific diagnoses (ICD-10 codes F52, N4\* and N5\* for men, D25, D26, D27, E28, N7\*, N8\*, N9\*, O\* and Z3\* for women) were analysed using only samples of the relevant sex as controls. The PheWAS significance threshold was set to  $P < 0.05/420$ , as 420 PCs calculated on all included phenotypes explained 95% of the variability.

### Identification of CNV-tagging SNPs for pQTLs

To aid the interpretation of the CNV–pQTL results, we examined additional pQTLs not detected in this study due to the small sample size or the lack of protein measurements, by using a CNV-tagging proxy SNP approach. To detect additional CNV–protein associations, we extracted all SNPs with MAFs  $> 0.01$  from each common (major allele frequency  $< 0.95$ ) CNV and its 500-kb flanking region, as identified in 2230 Estonian WGS samples. We calculated Pearson correlation coefficients between the CNVs and SNPs using custom R scripts. SNPs with  $R^2 > 0.8$  were defined as CNV-tagging proxy SNPs. The proxy SNPs were then compared with a published set of SNP pQTLs in two larger sets of unique proteins<sup>4,9</sup> to determine the degree of overlap. We used data on 1021 independent autosomal lead pQTL variants for 1478 proteins from the large-scale pQTL study conducted by Sun et al.<sup>4</sup>; 824 (80.7%) of these variants were present in the EstBB WGS dataset. We extended the analysis to include data from the largest pQTL study to date, conducted with 35,571 samples and resulting in the detection of 10,248 independent autosomal pQTLs for 1463 proteins<sup>9</sup>. The two studies encompassed 2438 unique proteins, enabling broader investigation. The resulting loci were reported as potential cases in which the underlying CNVs might be the causal variants. Figure depicting tagged–CNV pQTLs was done by using the RIdeogram v02.2.2 R package<sup>64</sup>.

## Results

### Discovery of pQTLs

We identified 278 (184 *cis* and 94 *trans*) independent pQTLs for 157 (48.2%) of the 326 proteins examined, using a protein-level genome-wide significance threshold of  $P < 5 \times 10^{-8}$  (Supplementary Table S2). When using a strict multiple testing correction threshold of  $P < 2.76 \times 10^{-10}$ , 151 pQTLs (131 *cis* and 20 *trans*) for 99 proteins remained significant (Supplementary Table S2). The power calculations suggested that our dataset had 80% power to detect the pQTL effects that explained at least 7.6% and 9.7% of the variance in protein expression (when using  $P < 5 \times 10^{-8}$  and  $P < 2.76 \times 10^{-10}$  respectively, Supplementary Fig. S1, Supplementary Table S5). All interpretative analyses were conducted using protein-level genome-wide-significant results.

To provide a comparison with previous research, we compared our results with previously published data. From the Pietzner et al. study<sup>7</sup>, 147 pQTLs (52.88%) were nominally significant ( $P < 0.05$ ) and accessible for comparisons. After correcting for multiple testing, 147 pQTLs remained significant (Benjamini–Hochberg FDR  $< 0.05$ ) and 91.84% (135/147) of pQTLs were directionally concordant with the current study (Supplementary Table S2). 66.19% (184/278) of pQTLs were tested in the Sun et al. study<sup>4</sup>. Of them, 55.98% (103/184) were significant (Benjamini–Hochberg FDR  $< 0.05$ ) and 89.32% (92/103) were directionally concordant (Supplementary Table S2). 7.55% (21/278) pQTLs were also tested in the Suhre et al. study<sup>51</sup> and 57.14% (12/21) were significant (Benjamini–Hochberg FDR  $< 0.05$ ), and all the significant pQTLs were directionally concordant with the current study (Supplementary Table S2). 12.23% (34/278) pQTLs were tested in the Folkersen et al. study<sup>50</sup> and 85.29% (29/34) of the pQTLs were significant (Benjamini–Hochberg FDR  $< 0.05$ ) and all the significant pQTLs were also directionally concordant with the current study (Supplementary Table S2). Concordance with previous studies demonstrates the robustness of our results.

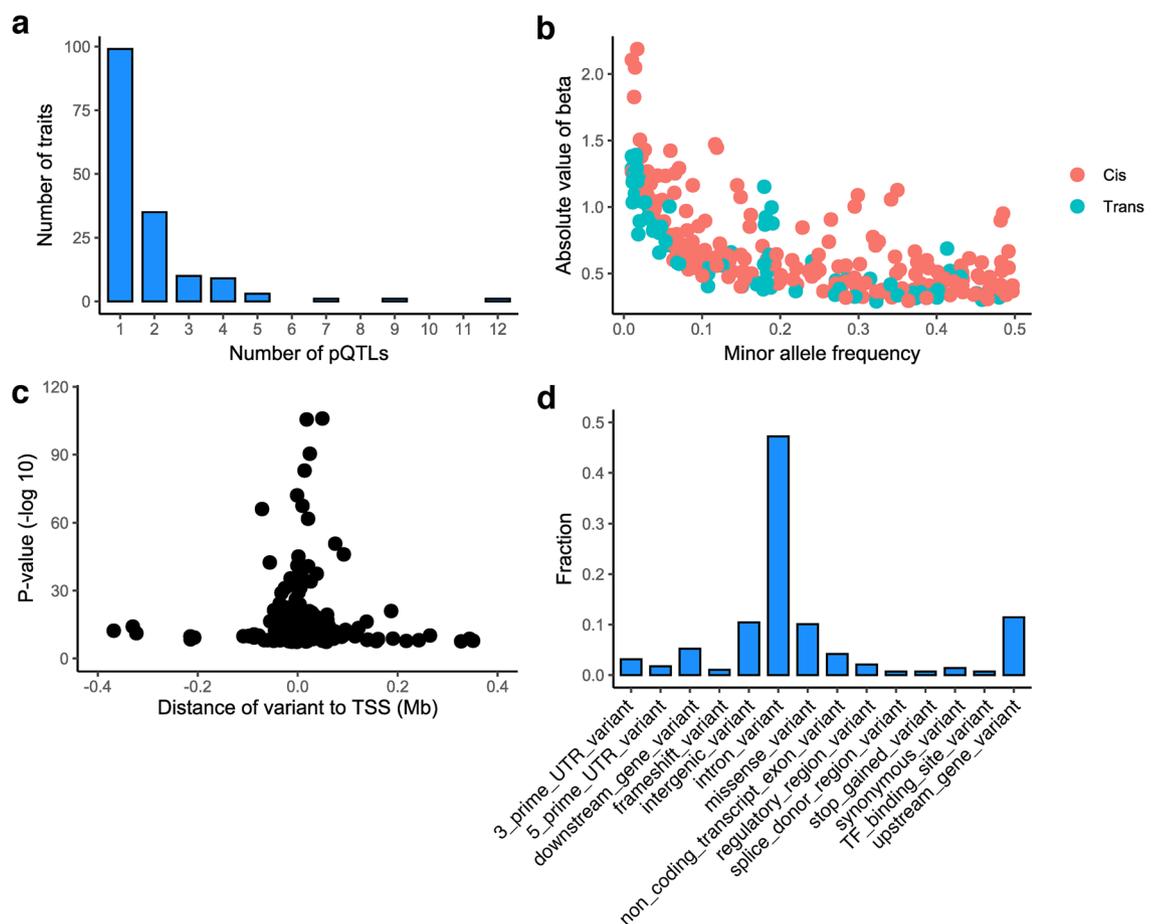
Fourteen (4.3%) of the proteins were measured in multiple arrays. Associations for the CXCL1, CCL3 and VEGFA proteins were validated by multiple independent arrays, in which the same genetic regions reached

genome-wide significance and showed concordant effect directions. The total numbers of associated proteins were similar for all panels and ranged from 33 to 43 (Supplementary Table S2). The detected associations included 278 independent pQTL variants [184 (66.2%) *cis* and 94 (33.8%) *trans*], 9.71% of which were indels. Of the 157 associated proteins, 61 (38.9%) had more than one independent pQTL. Twenty-one proteins had both *cis* and *trans* associations. A MICA-MICB heterodimer coded from the chromosome 6 *HLA* region had the largest number of independent associations ( $n = 12$ ; Fig. 2a). In concordance with previous studies<sup>4,9,65</sup>, there was an inverse relationship between the effect size and MAF (Fig. 2b), and the associations were the strongest for significant *cis*-pQTL variants located nearest to the TSSs of the relevant protein genes (Fig. 2c). The largest proportion of these *cis*-pQTLs [ $n = 78$  (42.4%)] was located in intronic regions (Fig. 2d). Of the 184 *cis* associations detected for 104 proteins, 31 (16.85%) were with protein-altering primary lead *cis*-pQTL variants and an additional 5 were with *cis*-pQTL variants in high LD with PAVs (Supplementary Table S2). These 36 (12.5%) pQTLs were designated as potential pseudo-pQTLs because currently it is difficult to exclude the possibility of technical signal happening due to the difference in antibody binding affinity.

The strongest *cis* association was between the missense variant rs2228145 (p.Asp358Ala) and the IL6RA level (MAF = 0.35,  $P = 1.04 \times 10^{-106}$ ). Additional strong *cis* associations included the rs1569960 and the SIRPA level (MAF = 0.34,  $P = 2.67 \times 10^{-106}$ ) association, with four independent signals in the SIRPA *cis* region; and a frameshift-causing insertion rs139130389 and the FOLR3 level (MAF = 0.12,  $P = 3.91 \times 10^{-91}$ ) association, with four independent signals in the FOLR3 *cis* region.

The most significant *trans* association was that of the *PLAUR* missense variant rs4760, located on chromosome 19, affecting the level of TNFRSF10C (8p21.3; MAF = 0.18,  $P = 4.60 \times 10^{-56}$ ). Strong *trans* associations were between the rs8176671 and the CDH5 level (16q21; MAF = 0.19,  $P = 8.83 \times 10^{-40}$ ) as well as between the deletion rs8176643 and the SELE level (1q24.2; MAF = 0.18,  $P = 7.98 \times 10^{-36}$ ); both of these variants are intronic variants for the 9q34.2 locus of the *ABO* gene. This locus was a *trans*-signal hotspot, with intronic variants additionally associated with the ICAM2, galectin-4 (LGALS4), PODXL and LIFR protein levels. Additional *ABO* variant rs12216891 was associated with the CTCR level (MAF = 0.19,  $P = 8.39 \times 10^{-30}$ ).

Two of the proteins examined (MICA/B and IL27) are heterodimers, made up of multiple subunits that are translated from two different genes at distinct loci. For IL27, we identified one independent *trans* signal



**Figure 2.** (a) Numbers of genome-wide significant associations of variants with protein traits. (b) Absolute beta values according to minor allele frequencies (MAFs). (c) Significance of primary pQTL mapping *cis* associations according to distances from transcription start sites (TSSs). (d) Functional annotation classes for the top *cis* variants from pQTL mapping, expressed as fractions.

for an intronic variant for CCDC94 (rs56075200; MAF = 0.32,  $P = 8.62 \times 10^{-35}$ ). For MICA/B, we identified ten independent signals in the *cis* region of one subunit on chromosome 6 (the strongest signal was for an intronic variant of MICA: rs3132467; MAF = 0.30,  $P = 3.04 \times 10^{-68}$ ) and two *trans* associations.

To determine if there were any corresponding eQTLs for pQTLs, we conducted an eQTL analysis, using the whole blood gene expression data from the same individuals and the same time point. Gene expression data was available for 109 proteins with 201 pQTLs, including two heterodimers with two subunits encoding the protein. In total, we detected 62 significant (Benjamini–Hochberg FDR < 0.05) eQTLs (59 *cis*, 3 *trans*) (Supplementary Table S6). 77% (48/62) of them were directionally concordant with corresponding pQTLs.

We found that 95% CSs for 151 proteins were linked to 131 independent genomic loci (Supplementary Table S7). LDLR, TNFRSF11B, TNFRSF6B, WISP1, CXCL1 and PLA2G1B proteins showed significant pQTL effects but yielded no CS. Signals for CCL3, CXCL1 and VEGFA from multiple assays were also validated by fine mapping to the same genetic regions. The 95% CSs contained an average of 15.7 variants (*cis* sets, 15.76; *trans* sets, 15.6). Fifty-five (36.4%) proteins had single-variant CSs. Of the 31 proteins with single-variant CSs in *cis* regions, 13 were fine-mapped to lead PAVs from primary pQTL mapping. Thirty-three (32.7%) out of 101 *cis* regions were fine-mapped to more than one signal (mean, 1.4 signals/region), with the CCL24 *cis* region having the largest number of independent CSs ( $n = 5$ ). In contrast, all associated regions for pQTL *trans* signals were fine-mapped to a single CS.

Since a large proportion (169/278) of primary pQTLs were located in intergenic and intronic regions, we queried RegulomeDB<sup>32</sup> to establish the variants' potential regulatory function. We obtained regulatory information for 260 of 278 pQTLs corresponding to 251 unique lead variants (Supplementary Table S8). Eleven variants (all *cis*) were previously established eQTLs and had evidence for transcription factor binding- and/or DNase peak-related functions. Seventeen lead variants (12 *cis* and 5 *trans*) had chromatin immunoprecipitation sequencing- and DNase-based evidence for regulatory functions, but were not eQTLs.

### pQTL–eQTL colocalisation

The pQTL–eQTL colocalisation analysis was performed with 198 pQTL loci (that contained 278 independent pQTL signals for 157 unique proteins;  $P < 5 \times 10^{-8}$ ), 18 eQTL Catalogue datasets<sup>18</sup> and GTEx tissue eQTL data<sup>58</sup>. We identified 14,064 cases of pQTL–eQTL colocalisation (PP4 > 0.8), involving 105 proteins [7936 (56.4%) *cis*- and 6128 (43.6%) *trans*-pQTLs; Table 1, Supplementary Table S9]. Colocalisations classified as *cis* consisted of 2021 (25.5%) cases in which colocalising eQTLs and pQTLs affected the same gene product and 5915 (74.5%) cases in which the colocalising loci affected different gene products in the *cis* regions. *Cis* and *trans* pairs were specific to 73 and 26 proteins, respectively, and 6 proteins (IL1R2, TEK, MIA, FCRLB, PDCD1LG2 and MICA-MICB) had colocalisations for both *cis* and *trans* associations. The largest number of colocalisations was found for pQTLs of the MICA-MICB heterodimer ( $n = 6583$ ), followed by OSCAR ( $n = 1207$ ) and ACP5 ( $n = 1105$ ) pQTLs.

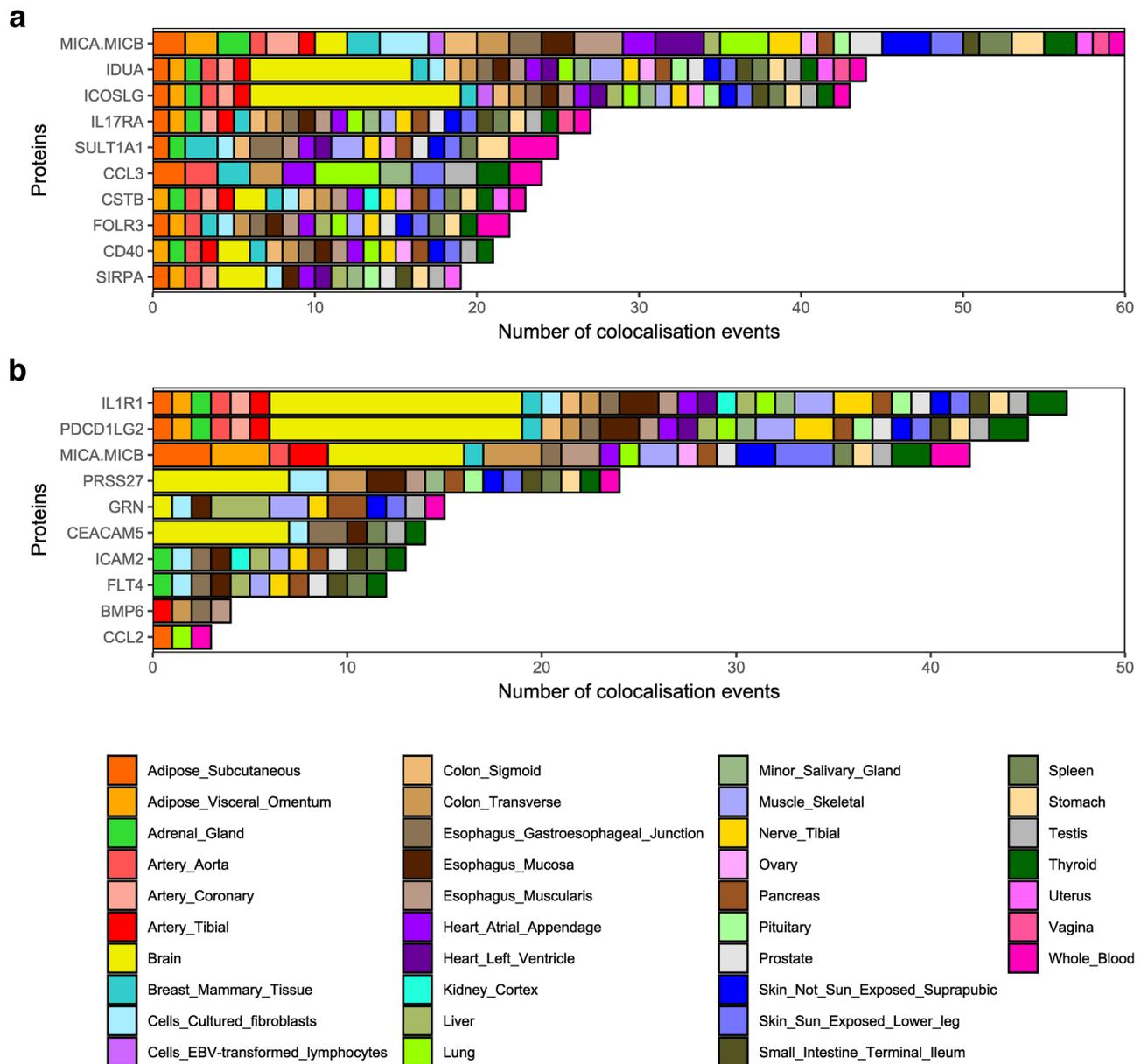
Since the protein measurements originated from blood, the most widely studied tissue, the largest fraction of pQTLs colocalised with blood eQTLs. However, while using the GTEx dataset<sup>58</sup>, we also found 739 cases of pQTL–eQTL colocalisation in multiple tissues (Fig. 3, Supplementary Table S9). For 55 proteins with *cis*-pQTLs, 503 (68.1%) colocalising eQTLs were identified; for 22 proteins with *trans*-pQTLs, 236 (31.9%) colocalising eQTLs were identified. *Cis*-pQTLs colocalising with eQTLs were detected in 49 tissues, and *trans*-pQTLs colocalising with eQTLs were identified in 46 tissues (not in Epstein-Barr virus-transformed lymphocytes or uterine or vaginal tissue).

### PheWAS on metabolite and epigenetic QTLs

Queries for the 268 unique variants corresponding to 278 significant pQTL lead variants led to the identification of 18 variants (from 7 *cis* and 13 *trans* associations for 19 proteins) associated with 160 metabolite traits (Supplementary Table S10). The majority [ $n = 158$  (52.3%)] of the mQTLs discovered were for the APOE missense variant rs7412, which had a *trans* association with the level of LDLR. Four metabolic traits [apolipoprotein B, the concentration of very small very-low-density lipoprotein (VLDL) particles, and phospholipids and total lipids in very small VLDL] had seven associations each.

Dataset	<i>Cis</i> -pQTL colocalising with eQTL (eQTL-pQTL same gene)	<i>Trans</i> -pQTL colocalising with eQTL
Gene expression (RNAseq)	710 (393)	398
Gene expression (microarray)	79 (51)	25
Exon expression	3899 (777)	2750
Txrevis	2533 (547)	2338
Transcript usage	715 (253)	617
Total	7936 (2021)	6128

**Table 1.** Overview of significant colocalisation events for eQTLs from eQTL Catalogue datasets and pQTLs. The numbers of colocalisation with genes encoding corresponding proteins are shown in parentheses.



**Figure 3.** Overview of (a) 10 *cis*-pQTL and (b) *trans*-pQTL proteins with the most colocalising eQTLs from the GTEx database (version 8; GTEx Consortium, 2020). Colours indicate eQTL tissues of origin. Brain tissues are pooled; a complete list is provided in the Supplementary Table S9.

From the epigenetic QTL datasets, we identified 6236 meQTLs, 267 histone modification QTLs and 129 exon-inclusion PSI associations for 193 primary pQTLs (from 142 *cis* and 60 *trans* associations for 130 proteins; Supplementary Table S11). Most ( $n = 256$ ) meQTLs were associated with the ADAM8 *cis*-pQTL rs2995310. The variant with the most ( $n = 10$ ) histone modifications was rs10415777, a *cis*-pQTL for OSCAR. Methylation data originates from five tissues: cord blood, monocytes, neutrophils, T cells and whole blood; due to tissue availability, 78.7% (4906/6236) of the identified meQTLs were from whole blood studies.

### Common SNP pQTLs and complex traits

#### *PheWAS*

The queries for the 268 unique variants corresponding to 278 significant pQTL lead variants and their high-LD proxies led to the identification of 135 (50.4%) variants with 5046 significant associations for 432 complex traits (Supplementary Table S12). Of these associations, 1538 (30.5%) were with various blood cell traits from the study conducted by Astle et al.<sup>66</sup>. As expected, given the targeted nature of our protein panels, coronary artery disease (CAD) and rheumatoid arthritis were most often linked to pQTLs with 118 and 99 associations, respectively. For example, 5 of 145 significant independent signals for CAD from mixed-ancestry samples<sup>67</sup> and 2 of 7 significant loci for rheumatoid arthritis from the study conducted by Stahl et al.<sup>68</sup> were pQTLs in our dataset. In terms of

the most associations per pQTL lead variant, *ABO* intronic variant rs507666 had the most associations per lead pQTL variant [ $n = 332$ , 85 (25.6%) with blood cell traits]. No associated traits were found for 62 proteins.

For 61 proteins (64 lead pQTL variants, 36 *cis*- and 28 *trans*-pQTLs), significant associations were detected in both the eQTL colocalisation analysis and PheWAS. We restricted this set to 27 proteins (28 variants) which were not coded from the *HLA* region but showed associations with diagnosis, treatment, or other phenotypes linked directly to health status (excluding haematological and biochemical measurements). Six of these proteins (CD6, PRSS27, CEACAM5, CD40, TNFRSF6B and IL1RL1) had significant colocalisations with eQTLs from brain tissue, but no evidence of shared conditions with direct effects on the brain tissue in the PheWAS.

For example, based on pQTL-eQTL colocalisation analysis, IL6R pQTL signal was also an eQTL of the *IL6R* gene in macrophages, monocytes, T cells, whole blood and pancreatic islets. A previous study has shown a link between IL6R and CAD<sup>69</sup>. We also identified associations between IL6R pQTLs and CAD, rheumatoid arthritis and 7 other disease traits (Supplementary Table S12), thereby supporting the findings of the study<sup>69</sup>. As another example, *IL1RL1* pQTLs colocalised with *IL1RL1*, *IL18R1* and *IL18RAP* eQTLs detected in multiple cell types with direct effects on the immune system (e.g. T-cells; Supplementary Table S9); these variants were associated with asthma and allergic reactions in the PheWAS.

Eleven out of 27 proteins had *trans*-associations. *Trans*-pQTLs for the *CTRC* and *TEK* proteins were in the *ABO* locus and colocalised with *ABO* eQTLs; in the PheWAS, they were linked to multiple self-reported diagnoses (e.g. 'blood clot in the leg') from the UK Biobank sample, and to haematological traits.

Most [ $n = 140/157$  (89.2%)] of the proteins with significant pQTLs belonged to the druggable genome category. These proteins were associated with 1365 drug-gene interactions.

### Colocalisation analysis

Based on the pQTL associations with genetic regions ( $P < 5 \times 10^{-8}$ ), PheWAS (PhenoScanner  $P < 1 \times 10^{-5}$ ) and eQTL colocalisation results (colocalisation with a  $PP_4 \geq 0.8$ ), we chose five *cis*-pQTL effects (affecting FGF5, IL1RL2, TNFRSF6B, IL2RA, and IL6R) that were associated with clinical traits and had significant pQTL-eQTL colocalisations. Furthermore, *SULT1A1* was chosen due to additional CNV-pQTL associations in its region which enabled to analyse colocalisation with respective complex traits. All selected proteins except IL6R had synonymous lead pQTL variants. Therefore, the input data for colocalisation analyses comprised of region-based summary statistics for 6 protein traits and 61 clinical complex traits (83 pQTL-complex trait pairs).

We identified 46 significant colocalisation events (Supplementary Table S13). FGF5 had 25 colocalisations with cardiovascular phenotypes and medications, such as CAD and perindopril use. IL6R had a total of 11 significant colocalisations, which included colocalisations with CAD as well as immunological conditions such as asthma. TNFRSF6B and *SULT1A1* colocalised with inflammatory bowel disease, and TNFRSF6B also separately colocalised with its two main forms: Crohn's disease and ulcerative colitis. IL2RA colocalised with tonsillectomy +/- adenoid operation. The PheWAS revealed associations of IL1RL2 with immune diseases which were not supported by the colocalisation results.

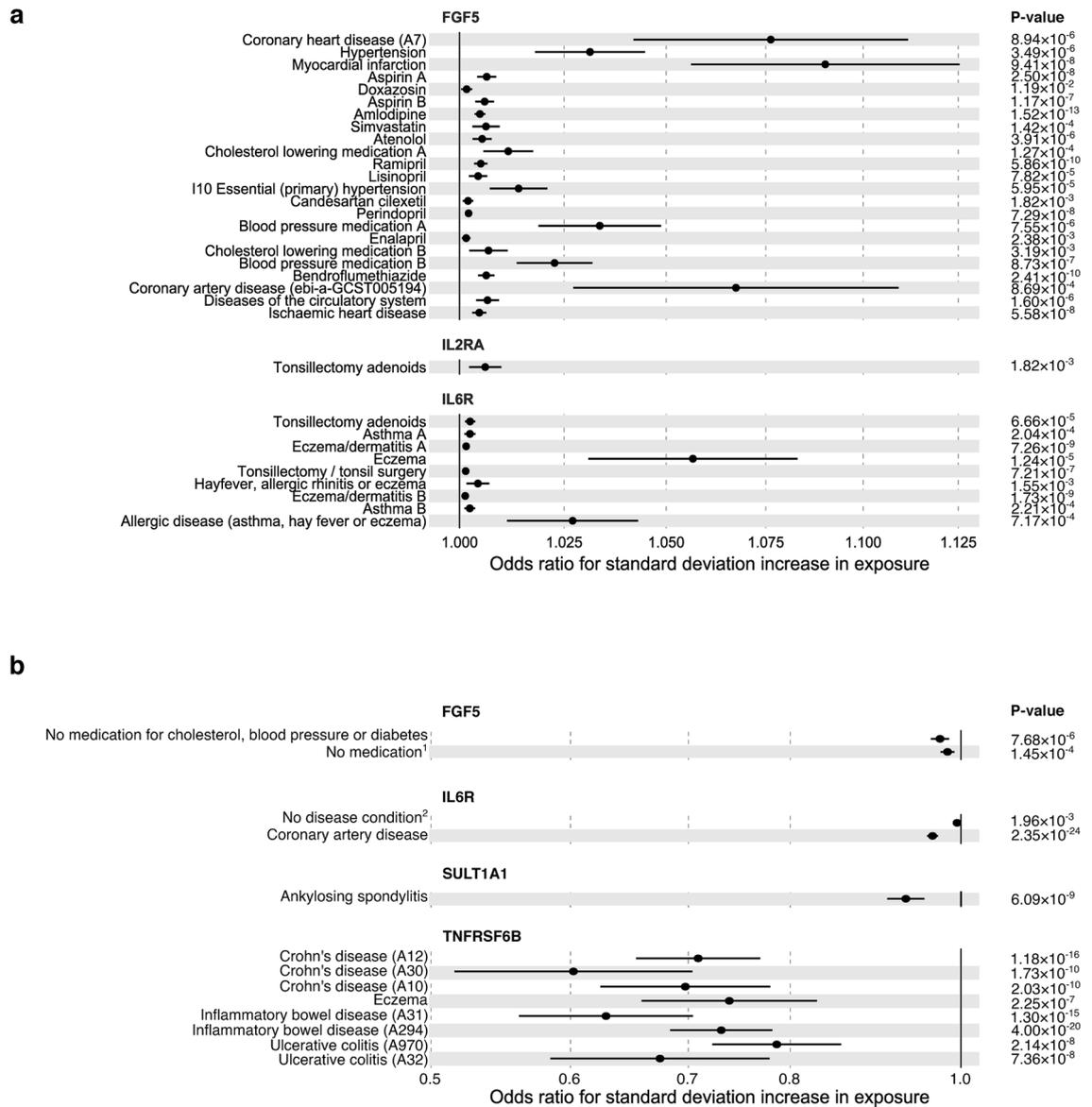
### MR findings

We conducted MR analyses using 46 significant (FDR-corrected) pQTL-complex trait pairs from the colocalisation analysis (Fig. 4, Supplementary Table S14). We found a causal relationship between the elevated level of soluble IL6R and a lower risk of cardiovascular disease ( $P = 2.35 \times 10^{-24}$ , Benjamini-Hochberg FDR =  $1.08 \times 10^{-22}$ ). Higher IL6R levels were also associated with an increased risk of inflammatory conditions such as asthma and eczema ( $P = 2.04 \times 10^{-4}$ , Benjamini-Hochberg FDR =  $2.60 \times 10^{-4}$ ;  $P = 1.24 \times 10^{-5}$ , Benjamini-Hochberg FDR =  $1.96 \times 10^{-5}$ , respectively). The TNFRSF6B level was causally linked to a reduced risk of inflammatory bowel disease and its subtypes (inflammatory bowel disease (A294),  $P = 4.00 \times 10^{-20}$ , Benjamini-Hochberg FDR =  $9.19 \times 10^{-19}$ ; Crohn's disease (A12),  $P = 1.18 \times 10^{-16}$ , Benjamini-Hochberg FDR =  $1.82 \times 10^{-15}$ ; ulcerative colitis (A970),  $P = 2.14 \times 10^{-8}$ , Benjamini-Hochberg FDR =  $7.56 \times 10^{-8}$ ). Elevated levels of FGF5 were associated with a significantly increased risk of coronary disease ( $P = 8.94 \times 10^{-6}$  and Benjamini-Hochberg FDR =  $1.47 \times 10^{-5}$ ).

### Rare variant pQTLs

The gene-based association analysis revealed 19 significant associations [5 (26.3%) *cis* and 14 (73.7%) *trans*] emanating from 19 genes containing rare nonsynonymous SNPs and affecting the levels of 7 proteins (Supplementary Table S15). The majority of identified rare variant effects (13, 68.4%) were with the level of GDF-15. The most significant rare variant association was a *trans* signal between *JAKMIP1* on chromosome 4 and the level of GDF-15 ( $P = 5.41 \times 10^{-12}$ ). We also assessed if rare nonsynonymous SNPs affect the expression of same genes encoding the corresponding pQTL proteins, however we did not detect any nominally significant (Benjamini-Hochberg FDR < 0.05) associations (Supplementary Table S15).

We next conducted GeneMANIA network analysis<sup>39,40</sup> to identify functional connections between genes harbouring rare SNPs and proteins affected by *trans* associations. First, we studied the potential connection between rare variant genes associated with the GDF-15 level. Ten of the identified genes harbouring rare SNPs (*CKAP5*, *GDF15*, *JAKMIP1*, *KRT19*, *STAT5B*, *SLC35E1*, *RNF112*, *TUBGCP4*, *ZNF766* and *PPAPDC1B*), including gene encoding identified pQTL protein, formed shared network with GDF-15, based on co-expression (57.85%), pathway (19.97%), physical (18.45%) and genetic (3.73%) interactions, according to GeneMANIA. However, no functional connection to GDF-15 was found for *LY6G6E*, *RPL7L1* and *EFR3B*. We conducted additional network analysis using the STRING database<sup>41</sup>; at the medium confidence interaction score of 0.4, there was no shared network with GDF-15. On the other hand, at the low confidence interaction score of 0.15, connections were established between *JAKMIP1*, *RNF112*, *CKAP5* and *TUBGCP4* based on co-expression and co-mentioning in PubMed abstracts but no links were detected with GDF-15. GDF-15 formed shared links



**Figure 4.** Forest plots of Mendelian randomisation results for proteins with (a) positive and (b) negative effects on complex traits. Protein (exposure) names are indicated on top of the section, complex traits (outcomes) are on the left side. Multiple instances of traits with the same name for one protein, indicating MR signal replication across multiple studies of the same trait, have been marked 'A' and 'B'. Error bars denote standard errors and all presented results are significant at a Benjamini–Hochberg FDR < 0.05. Details of causal associations are provided in Supplementary Table S14. <sup>1</sup>“Medication for cholesterol, blood pressure, diabetes, or take exogenous hormones: None of the above” (MRC IEU UK Biobank); <sup>2</sup>“Blood clot, DVT, bronchitis, emphysema, asthma, rhinitis, eczema, allergy diagnosed by doctor: None of the above” (MRC IEU UK Biobank).

with KRT19 and STAT5B based on experimental or biochemical data, co-mentioning in PubMed abstracts and co-expression. In addition, the functional link between SLC35E1 and RPL7L1 was due to co-expression. *KRT19* and *STAT5B* associations with GDF-15 were supported by GeneMANIA as well as STRING database. In the case of BioGRID<sup>42,43</sup>, none of the *trans* associations had support for a functional connection with GDF-15. *Trans* associations between rare variants and SELPLG and MUC-16 levels were supported by the GeneMANIA-based identification of two shared networks: between *TMEM119* and SELPLG, as well as *GAL3ST2* and MUC-16, respectively. Those connections were based mainly on physical interactions (77.64%) and co-expression (8.01%). Based on the STRING database<sup>41</sup>, at the medium confidence interaction score of 0.4, *TMEM119* and SELPLG do not have a shared link. However, at the low confidence interaction score of 0.15, *TMEM119* and SELPLG have a shared connection due to co-expression and co-mentioning in PubMed abstracts, functional link was also supported by GeneMANIA. According to the STRING database<sup>41</sup>, even at the low confidence interaction score of 0.15, *GAL3ST2* and MUC-16 do not share a functional connection. Based on BioGRID<sup>42,43</sup>, *trans* associations between *TMEM119* and SELPLG, as well as *GAL3ST2* and MUC-16 had no supporting interactions.

Four proteins (CTSZ, GDF-15, PON3 and SELPLG), had significant associations from both, common variant and rare variant pQTL analyses. For CTSZ and GDF-15, the genetic regions detected from the rare variant analysis were not the same as identified by SNP pQTL analysis. However, PON3 had direct *cis* associations emanating from from 7q21.3 locus in both analyses: nonsynonymous variants of *PON3* in the rare variant pQTL analysis and rs10953142 in the common variant pQTL analysis. Similarly, SELPLG had *cis* associations emanating from 12q24.11 locus: nonsynonymous variants of the *TMEM119* for rare variant analysis and an intergenic rs11114010 for common variant analysis.

### CNV pQTLs

We detected 12 significant (Bonferroni-corrected  $P$ -value threshold  $1.12 \times 10^{-7}$ ) pQTL associations between CNVs and plasma protein levels (7 *cis* and 5 *trans*, 11 proteins; Supplementary Table S16), with two *cis* associations detected for the MICA-MICB heterodimer. The CNV eQTL analysis in the overlapping set of samples identified 673 significant (Bonferroni-corrected  $P$ -value threshold  $1.61 \times 10^{-9}$ ) CNV eQTLs for 244 unique genes (Supplementary Table S17). 16.67% (2/12) of significant CNV pQTLs had significant CNV eQTL associations with a corresponding gene.

For example, the deletion in the 3q12.1 intergenic region (chromosome 3: 98,410,653–98,414,807 bp; frequency = 0.651) acted as a hub, having multiple *trans* associations with protein levels: ICAM2 ( $P = 1.31 \times 10^{-29}$ ), FLT4 ( $P = 2.34 \times 10^{-24}$ ), PDCD1LG2 ( $P = 2.88 \times 10^{-15}$ ) and IL1R1 ( $P = 8.19 \times 10^{-8}$ ). Three of these associations (with ICAM2, FLT4 and PDCD1LG2) were also detected by the SNP pQTL analysis but did not remain significant after conditioning of the model on the CNVs, suggesting that CNV may underlie the observed associations. However, eQTL analysis indicated that none of the genes encoding those proteins is regulated by this locus, and a follow-up GeneMANIA network analysis<sup>39,40</sup> revealed a shared network based on physical interactions (77.64%), co-expression (8.01%), predicted functional relationship between genes (5.37%), co-localisation (3.63%), genetic interactions (2.87%), pathway (1.88%) and shared protein domains (0.60%).

Another *trans* association example was between a 5q13.2 CNV (chromosome 5: 70,305,253–70,312,310 bp; deletion frequency = 0.074, duplication frequency = 0.195) overlapping the *NAIP* gene but affecting IL-18 level ( $P = 7.9 \times 10^{-10}$ ). This locus was also an eQTL for *NAIP* ( $P = 6.4 \times 10^{-48}$ ), but not for IL18 expression ( $P > 0.001$ ). We also detected moderate correlation between IL-18 protein expression and *NAIP* gene expression (Spearman's  $R = 0.17$ ); Spearman correlation coefficient between IL-18 protein and gene expressions was 0.05. MR analysis using *NAIP* gene expression as exposure and IL-18 level as an outcome confirmed the causal effect of the CNV on the IL-18 protein level (Wald test;  $Z = 6.26$ ,  $P = 3.8 \times 10^{-10}$ ). This association was not observed in the SNP-based analyses, highlighting the case where the pQTL signal would not be detected.

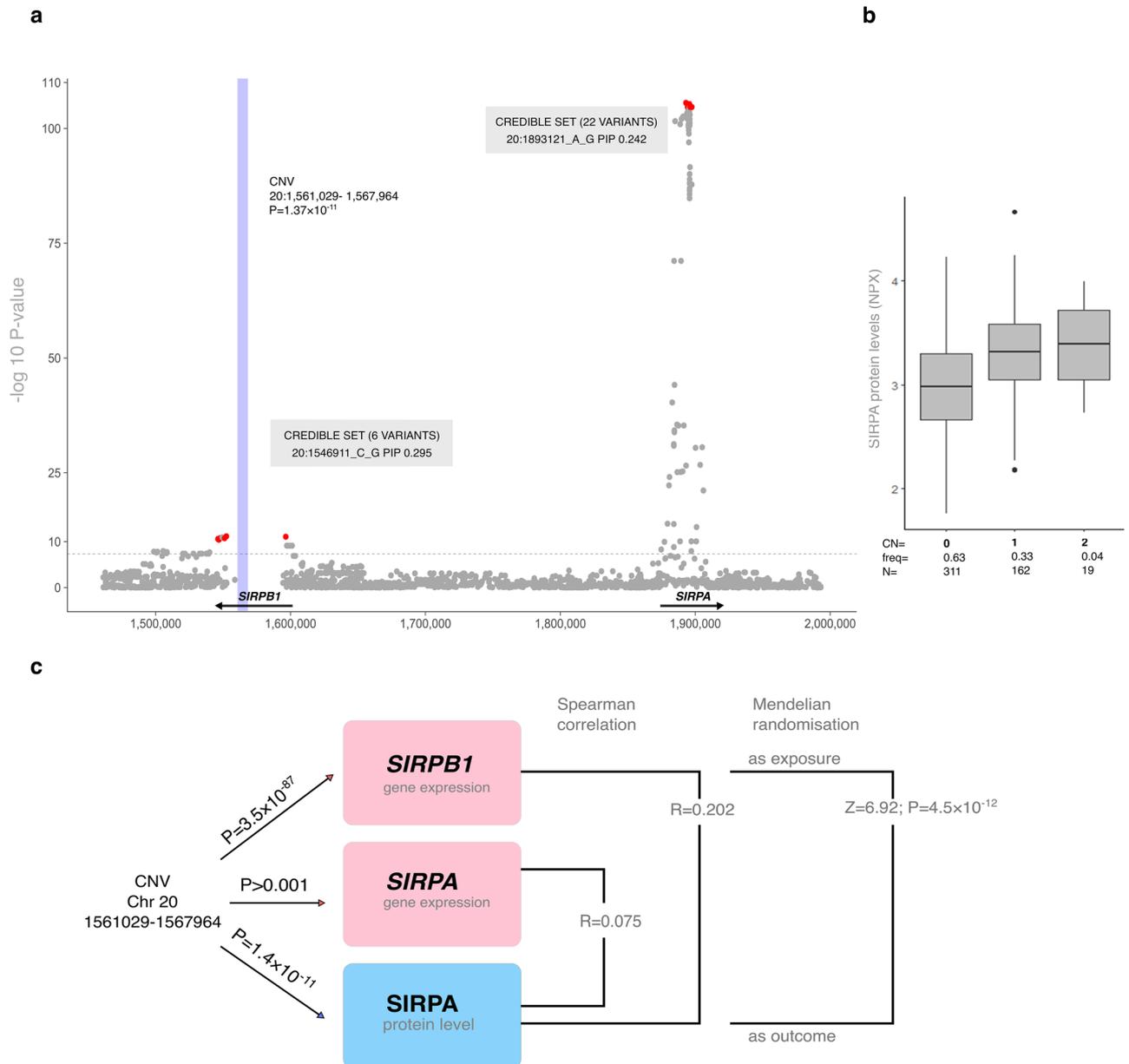
From *cis* effects, we detected an association between CNV in the 16p11.2 region (deletion frequency = 0.022, duplication frequency = 0.382; partially overlapping *SULT1A1*; pQTL,  $P = 3.46 \times 10^{-21}$ ; eQTL,  $P = 4.74 \times 10^{-119}$ ) and *SULT1A1* protein and gene expression. Similarly, we determined that a 19q13.42 deletion (frequency = 0.291) overlapping the *VSTM1* gene was an eQTL and a pQTL for nearby gene *OSCAR* ( $P = 1.77 \times 10^{-14}$  and  $P = 5.64 \times 10^{-9}$ , respectively). However, the CNV was also associated with the expression of *VSTM1* itself ( $P = 1.81 \times 10^{-39}$ ) and both gene–protein expression pairs showed moderate correlation (*OSCAR*–*OSCAR*, Spearman's  $R = 0.32$ ; *VSTM1*–*OSCAR*, Spearman's  $R = 0.34$ ). The effect of the CNV through gene expression is supported by the MR analysis, when using a CNV as an instrument, *OSCAR* expression as an exposure and *OSCAR* level as an outcome ( $Z = 5.94$ ;  $P = 2.81 \times 10^{-9}$ ) and secondly, *VSTM1* as an exposure and *OSCAR* level as an outcome ( $Z = 5.92$ ;  $P = 3.27 \times 10^{-9}$ ). Those results suggest that CNV works through gene expression, although it remains unclear whether the effect on the *OSCAR* level is through *OSCAR* or *VSTM1* gene expression.

Additionally, we identified an association between the *SIRPA* level and a high-frequency (frequency = 0.955) 20p13 deletion overlapping *SIRPB1*, a paralog of *SIRPA* ( $P = 1.4 \times 10^{-11}$ ; Fig. 5a and b). eQTL analysis indicated that the deletion was also associated with *SIRPB1*, but not *SIRPA*, expression ( $P = 3.5 \times 10^{-87}$ ). The correlation between *SIRPA* protein and gene expression was weaker than that between *SIRPA* protein and *SIRPB1* expression (Spearman's  $R = 0.075$  and 0.202, respectively). Colocalisation was confirmed by the Wald test ( $Z = 6.92$ ,  $P = 4.5 \times 10^{-12}$ ; Fig. 5c). In SNP pQTL fine mapping, we detected two independent CSs, overlapping *SIRPB1* [variant with the largest posterior inclusion probability (PIP) = 0.295] and at *SIRPA* (variant with the largest PIP = 0.242; Fig. 5a). When conditioned on the deletion, the significance of pQTLs from only the *SIRPB1* CNV region was reduced dramatically (chr 20 position 1,546,911 variant pQTL mapping,  $P_{\text{primary}} = 3.75 \times 10^{-11}$ ,  $P_{\text{conditional}} = 0.41$ , regional pQTL mapping with EMMAX linear mixed-model<sup>30</sup> and the occurrence of the CNV and the number of its copies used as an additional covariate). This example highlights that the second signal from the primary pQTL analysis *SIRPA* locus was due to CNV-tagging variants rather than an independent signal.

Associations for nine proteins significant in both, CNV and pQTL mapping, were emanating from the same loci in both analyses. For example, ICAM2 and FLT4 had *trans* associations with rs12493830 on chromosome 3 and a CNV (chromosome 3: 98,410,653–98,414,807 bp) in the same intergenic region, separated by 3859 bp.

### PheWAS for CNV pQTLs

Significant PheWAS associations were detected for three CNVs. For the MICA-MICB dimer pQTL, associations were detected between CNV on chromosome 6 (31,292,078–31,293,977 bp; deletion frequency = 0.876) and medium HDL triglycerides ( $P = 8.82 \times 10^{-5}$ ), and between a CNV on chromosome 6 (31,337,848–31,341,642 bp; deletion frequency = 0.074) and lower-limb oedema (ICD-10 code R60;  $P = 9.06 \times 10^{-5}$ ). Additionally, we detected nominally significant associations for a CNV on chromosome 19 (41,381,588–41,387,347 bp, deletion frequency = 0.054 and duplication frequency = 0.022) with the pQTL of the MIA protein level ( $P = 2.38 \times 10^{-6}$ ) and migraine (ICD-10 code G43;  $P = 3.14 \times 10^{-5}$ ).



**Figure 5.** (a) Regional plot combining SNP- and CNV-based results for the SIRPA level with additional single-variant fine-mapping information. The blue rectangle indicates the genetic location of the CNV. The horizontal dashed line indicates the genome-wide significance threshold of  $P = 5 \times 10^{-8}$ . Genetic variants identified by fine mapping as belonging to 95% credible sets are coloured red. The number of variants and the variant with the highest PIP in the credible set are indicated in grey boxes. (b) Box plot of SIRPA levels based on the CNV number of copies and frequencies. Error bars indicate 95% confidence intervals; the bottoms and tops of the boxes are the 25th and 75th percentiles, respectively; the lines inside the boxes indicate medians. Outliers are depicted as circles. (c) Overview of SIRPA level analyses.  $P$ -values are from the CNV-based pQTL analysis for SIRPA and eQTL analyses for *SIRPB1* and *SIRPA*.

### CNV-tagging SNPs

To further interpret the of CNV-pQTL results, we examined additional pQTLs for proteins that were not measured in our study. For that, we leveraged LD between the EstBB CNVs and previously reported pQTL SNPs and prioritised CNVs which could underlie the previously reported pQTL associations ( $R^2$  between SNP and CNV  $> 0.8$ ). We identified eight CNVs with possible effects on protein levels (Table 2) from the Sun et al. 2018 study<sup>4</sup>. Only one of those associations [proxy SNP rs10935473 with the CNV on chromosome 3 (98,410,653–98,414,807; deletion frequency = 0.651)] affecting FLT4/VEGF-sR3 levels, was identified in our study because the other proteins were not measured in our cohort.

We also detected 76 tagging SNP–CNV pairs for 33 unique CNVs and 72 proteins (Supplementary Table S18) from a more recent Sun et al. 2022 study<sup>9</sup>. Twenty-nine (40.3%) of the proteins were also measured in the EstBB cohort, of which six proteins had significant CNV pQTLs ( $P < 1.12 \times 10^{-7}$ ). However, CNV-based pQTLs of the

chr	Position	marker	CNV	CNV Frequency	R <sup>2</sup>	Type	gene	protein
1	55,097,068	rs11206397	1:55,092,289-55,095,991	deletion 0.538	0.90	<i>cis</i>	<i>FAM151A</i>	F151A
1	159,004,851	rs72709516	1:159,016,577-159,019,397	duplication 0.001, deletion 0.122	0.97	<i>cis</i>	<i>IFI16</i>	IP16
1	196,821,380	rs115094736	1:196,728,841-196,730,702	deletion 0.265	0.97	<i>trans</i>	<i>CANX</i>	Calnexin
1	196,825,287	rs7519758	1:196,728,841-196,730,702	deletion 0.265	0.96	<i>trans</i>	<i>LRRC19</i>	LRC19
3	98,416,900	rs10935473	3:98,410,653-98,414,807	deletion 0.651	1.00	<i>trans</i>	<i>FLT4</i>	VEGF sR3
6	32,587,859	rs9271421	6:32,461,274-32,468,482	deletion 0.973	0.86	<i>trans</i>	<i>H6PD</i>	G6PE
8	57,876,576	rs112433249	8:57,918,258-57,925,230	deletion 0.031	0.90	<i>cis</i>	<i>IMPAD1</i>	IMPA3
16	89,781,756	rs34714188	16:89,896,104-89,898,445	duplication 0.001, deletion 0.108	0.91	<i>trans</i>	<i>PMEL</i>	GP100

**Table 2.** Overview of SNPs tagging CNVs for proteins reported by Sun et al. (2018). CNV frequencies are derived from the EstBB data.

MICA-MICB heterodimer and SIRPA were not associated with the same CNVs in the EstBB cohort as tagged by SNPs in Sun et al.'s<sup>9</sup> study. Twenty-five (32.9%) of the tagging SNP–CNV pairs were associated with a deletion in the 3q12.1 intergenic region (chromosome 3: 98,410,653–98,414,807 bp, frequency = 0.651; the closest gene is *ST3GAL6*), a *trans* association hub (Fig. 6), and the same deletion was associated with four proteins (ICAM2, FLT4, PDCD1LG2 and IL1R1) in the EstBB dataset.

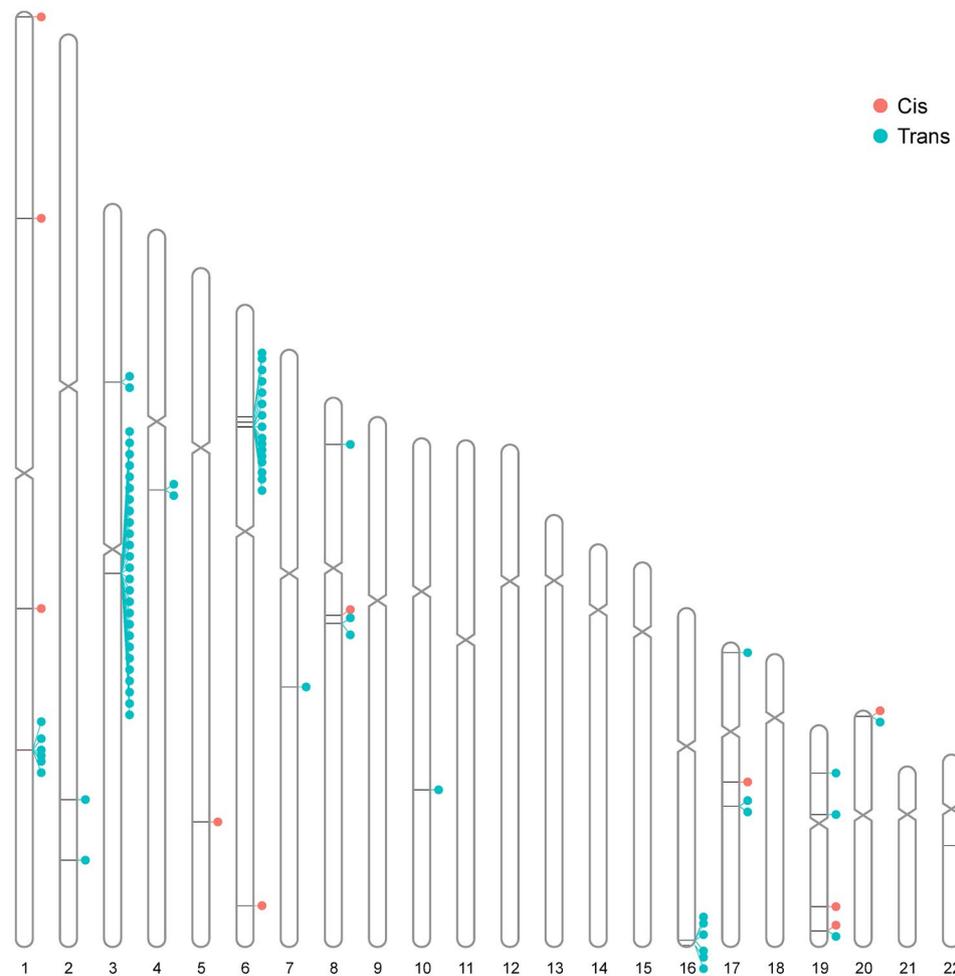
None of the pQTLs tagging the CNV has known associations with complex traits which are not cell type or metabolite related, according to the GWAS Catalog<sup>90</sup>. In addition, 19.7% (15/76) of the CNVs paired with tagging SNP were located in the *HLA* region on chromosome 6. The proteins TACSTD2, CLEC5A, IL15 and SIGLEC9 were affected by multiple *trans*-pQTL SNPs tagging CNVs. Whereas we detected a CNV associated with the SIRPB1 level on chromosome 20 (1,556,917–1,561,028 bp, deletion frequency = 0.336) and a deletion in the same locus overlapping *SIRPB1* and affecting the SIRPA level and (more strongly) *SIRPB1* gene expression, based on Sun et al. tagging–CNV analysis, the SIRPB1 protein level was associated with a different CNV than was its gene expression.

## Discussion

The SNV-pQTL analyses conducted in this study revealed 278 genetic variants (184 *cis* and 94 *trans*, including indels), that were associated with the levels of 157 unique proteins. Consistent with previous findings<sup>4,6,8</sup>, the largest proportion of *cis*-pQTLs was located in intronic and intergenic regions. The analysis of individual-level WGS data together with in-sample LD information, enabled us to pinpoint the likely causal variants with a good resolution through statistical fine mapping. This mapping led to the identification of at least one 95% CS for each of 98 (53%) *cis* and 87 (47%) *trans* signals. For 16 *cis* and 28 *trans* associations, we identified 95% CSs consisting of the single most likely causal variants, which are good candidates for further functional studies. Notably, the prioritised variants for nine (56%) of the single-variant CSs for *cis*-pQTLs had protein-altering effects. This observation outlines that it is important to consider technical epitope effects in the *cis*-pQTL analyses<sup>70</sup>. However, the identification of PAVs demonstrates that fine mapping is also helpful for prioritising biologically causal variants, because PAVs are likely to have a direct, albeit technical, effect on protein levels. Only a limited number of pQTL studies have conducted fine mapping<sup>9,71</sup> as one of the post-GWAS analyses. We and Zhang et al.<sup>71</sup> detected CSs for 58 (59.2%) protein *cis* regions using data from cohorts of European ancestry, and Sun et al.<sup>9</sup> fine mapped CSs in 127 (67.6%) genetic regions for 117 proteins, matching our findings. The 95% CSs contained an average of 15.7 variants in our study and 22.7 variants (9.6 *cis* and 29.4 *trans*) in that of Sun et al.<sup>9</sup>. Our CSs for *cis* associations contained an average of 15.76 variants, whereas Zhang et al. used imputed genotyping data and reported an average of 21.29 variants<sup>71</sup>. Generally smaller credible sets might outline the added value of WGS data on fine mapping performance.

To support our findings with orthogonal data, we used the most comprehensive publicly available eQTL resource, the eQTL Catalogue<sup>18</sup>, to conduct eQTL–pQTL colocalisation analyses. Detected colocalisations were 56.4% for *cis*- and 43.6% for *trans*-pQTLs. Of the *cis* associations, 25.5% (2021/7936) colocalised with the eQTLs for the corresponding protein-encoding gene from the full eQTL Catalogue, while for the GTEx dataset alone it was 54.3% (273/503). Given the use of eQTL data from different tissues, this analysis reflects how pQTLs may originate through active secretion or/and passive leakage, as 42.68% of all significant SNV-pQTL proteins identified are actively secreted into the blood at least in one isoform (Supplementary Table S19)<sup>72</sup>, meaning that more than half of these proteins do not originate from the blood. Similar to our findings, Pietzner et al.<sup>7</sup> recently detected a significant colocalisation of 50.1% of the *cis*-pQTLs with corresponding gene eQTLs using GTEx.

We sought to systematically identify links between proteins and phenotypes by conducting a PheWAS followed by a colocalisation analysis, in order to find signals likely driven by the same causal variant. We then applied MR to significant colocalisation events to assess causality, a strategy recommended by Zuber et al.<sup>17</sup>. As



**Figure 6.** Overview of SNP-tagged CNV and protein *cis* and *trans* associations. Each line depicts the CNV which is in LD ( $R^2 > 0.8$ ) with pQTL SNP previously reported by Sun et al. (2018) or Sun et al. (2022) study. Each dot indicates corresponding pQTL protein and colour depicts the type of association.

they have highlighted, a positive colocalisation finding typically implies a non-zero MR estimate, the reverse is not generally true<sup>17</sup>. For example, *FGF5* plays essential roles in the regulation of cell proliferation, including in cardiac myocytes, and cell differentiation<sup>73</sup>; it has also been associated with cardiac angiogenesis<sup>74</sup>. The *FGF5* locus has been linked to cardiovascular conditions in previous GWASs<sup>67,75</sup>. We detected a *cis* signal for the *FGF5* level and associated variants in the region, which overlapped with previous GWAS findings for cardiovascular diseases and medications used to treat them. Our colocalisation and MR results suggest that the *FGF5* level shares common causal SNPs with various heart-related conditions and treatments, prioritising it as an interesting target for future follow-up studies. However, the translation of PheWAS results to a molecular level is complicated by the nature of associated disease phenotype. Plasma proteins are potentially more relevant for circulatory diseases where the blood is in contact with the affected tissue, such as in the *FGF5* example, rather than for conditions with a limited number of affected tissues.

The availability of the high-quality WGS data also gave us a unique opportunity to investigate the effect of CNVs on protein expression. To the best of our knowledge, one study has previously studied CNVs in this context, focusing only on deletions<sup>15</sup>. We conducted the first comprehensive CNV-based pQTL mapping and identified 12 associations (7 *cis* and 5 *trans*) between plasma proteins and CNVs, including those with a *trans*-association hub CNV in the 3q12.1 region. We further interpreted the CNV-pQTLs using a CNV-tagging SNP approach with external data on a broader range of proteins. This strategy yielded additional CNV-based pQTLs for 79 proteins and determined that the 3q12.1-region hub CNV was associated with 25 proteins. Signals from the SNV and CNV analyses overlapped for nine proteins, which constitute interesting loci where QTL associations were likely driven by CNVs, rather than SNVs. This emphasises the value of the CNV data, especially if the purpose is to prioritise causal genetic variation underlying the pQTL signal. None of the associations reported by Png et al.<sup>15</sup> were replicated in this study, possibly because there was only a partial overlap between the assayed protein sets, differences between cohorts (European ancestry vs a Greek population isolate with population-specific CNVs)<sup>76</sup>, and differences in the approach used for CNV detection.

As an example, we outline IL-18, a pro-inflammatory cytokine that plays important roles in natural killer cell activation and the T-helper 1 response<sup>77</sup>. We found that a CNV on chromosome 5 overlapping with *NAIP* has *trans* effects on the IL18 protein level and a *cis* effect on the *NAIP* gene expression level, but there is no significant effect on the *IL18* gene expression. The *NAIP* eQTL signal was stronger than the IL18 pQTL signal, suggesting that the CNV affects the protein level through gene expression. As the *NAIP* level was not measured in our cohort, it remains unclear whether the main effect of the CNV is on *NAIP*. To our knowledge, there are no previous studies analysing the effect of genetic variants on *NAIP* level. *NAIP* is an anti-apoptotic protein and sensor component of the NLRC4 inflammasome that protects against bacterial pathogens, and *NAIP*-NLRC4 inflammasome activation has been reported to lead to elevated IL-18 expression in enterocytes and monocyte-derived macrophages<sup>78</sup>. This example highlights the importance of including structural variants in addition to SNVs in studies of the genetic basis of molecular traits, as also exemplified by the CNV-tagging SNP approach.

We identified 19 significant rare variant effects on the levels of seven proteins that would not have been detected by the SNV pQTL analysis alone. Gene-based pQTL analyses of rare variants constitute an emerging approach<sup>10–13</sup>, and no golden standard for their performance has been established, making the replication of findings difficult. Previous studies indicate that few proteins are driven by rare variants<sup>11–13</sup>. Kierczak et al.<sup>13</sup> detected *cis*-region rare variant associations for four proteins (CTS2, CYR61, GDF-15 and PON3) and *trans* associations of rare *GAL3ST2* variants affecting the MUC16 level, the effect also detected in our study; they used a maximum MAF threshold 0.0239, whereas we used a standard conservative threshold of 0.01. The significant rare variant associations detected in our study were not reported in the largest gene-based rare variant pQTL study conducted to date which included three isolated European cohorts with a total sample size of  $n = 4422$ <sup>12</sup>. As an example, we found a rare-variant effect on GDF-15, which regulates food intake, energy expenditure and body weight in response to metabolic and toxin-induced stress<sup>79–81</sup>. The most significant association with the GDF-15 level was a *trans* association with rare variants in *JAKMIP1*, associated with type 2 diabetes and medications used to treat it<sup>82–84</sup>. Additionally, GDF-15 has been reported to be involved in inflammation, metabolism and cancer<sup>85</sup>, and recent findings support its role as a biomarker of metabolic stress<sup>86</sup>. Whereas we detected rare variant *trans* associations emanating from GDF-15 for nine proteins, only SNP-based *cis* associations with GDF-15 itself have been identified in previous pQTL studies<sup>9,86</sup>. This demonstrates that gene-based rare variant pQTLs complement single variant analyses and help to unravel novel biologically interpretable associations.

Our study has several limitations. First, the sample size was small relative to those of recent pQTL studies, which made the detection of *trans* effects with greater multiple-testing burden and weak effects of common and rare variants more difficult. Rare genetic variants tend to have greater population specificity<sup>87</sup>, making replication of findings from rare variant analyses more difficult. Same applies to common CNVs we reported in our pQTL analyses; structural variants are currently understudied in terms of pQTL detection, limiting replicability. Second, most pQTL studies have been conducted using serum or plasma measurements from blood samples<sup>4,6,8,51</sup> and only a limited number of studies has involved the examination of liver and brain tissue-specific pQTLs<sup>88,89</sup>. Therefore, it is often challenging to understand whether observed pQTL effects manifest in the blood cells or reflect the regulation happening in some distal tissue. Finally, although we showed that CNVs affect plasma protein levels, to our knowledge no large-scale CNV-based association database is currently available to overlap the identified CNV-pQTL associations with CNV-phenotype associations. However, CNV-tagging SNPs could be used as a proxy method to assess the effect of CNVs on complex traits.

In conclusion, we have generated a comprehensive pQTL resource and interpreted it by using eQTL, as well as publicly available GWAS data. We have demonstrated the importance of including structural variants in addition to SNVs, to fully characterise the genetic background of plasma proteins and their links to health-related phenotypes.

## Data availability

All relevant data are within the paper and its Supporting Information files. All single variant protein analysis full results will be available for download at the EBI GWAS Catalog upon publication (accession numbers GCST90277554-GCST90277894). Due to the sensitive nature of participant data, the access to the EstBB individual-level data is restricted and regulated by the Estonian Human Genes Research Act, GDPR, and procedures established by the EstBB (<https://genomics.ut.ee/en/content/estonian-biobank>). Custom code is available on GitHub (<https://github.com/kalnapekis/pQTLs>).

Received: 22 June 2023; Accepted: 23 March 2024

Published online: 02 April 2024

## References

1. MacArthur, J. *et al.* The new NHGRI-EBI Catalog of published genome-wide association studies (GWAS Catalog). *Nucleic Acids Res.* **45**, D896–D901 (2017).
2. Maurano, M. T. *et al.* Systematic localization of common disease-associated variation in regulatory DNA. *Science* **337**, 1190–1195 (2012).
3. Geyer, P. E., Holdt, L. M., Teupser, D. & Mann, M. Revisiting biomarker discovery by plasma proteomics. *Mol. Syst. Biol.* **13**, 942 (2017).
4. Sun, B. B. *et al.* Genomic atlas of the human plasma proteome. *Nature* **558**, 73–79 (2018).
5. Emilsson, V. *et al.* Co-regulatory networks of human serum proteins link genetics to disease. *Science* **361**, 769–773 (2018).
6. Folkersen, L. *et al.* Genomic and drug target evaluation of 90 cardiovascular proteins in 30,931 individuals. *Nat. Metab.* **2**, 1135–1148 (2020).
7. Pietzner, M. *et al.* Mapping the proteo-genomic convergence of human diseases. *Science* **374**, eabj1541 (2021).
8. Ferkingstad, E. *et al.* Large-scale integration of the plasma proteome with genetics and disease. *Nat. Genet.* **53**, 1712–1721 (2021).

9. Sun, B. B. *et al.* Genetic regulation of the human plasma proteome in 54,306 UK Biobank participants. 2022.06.17.496443 Preprint at <https://doi.org/10.1101/2022.06.17.496443> (2022).
10. Solomon, T. *et al.* Identification of common and rare genetic variation associated with plasma protein levels using whole-exome sequencing and mass spectrometry. *Circ. Genom. Precis. Med.* **11**, e002170 (2018).
11. Gilly, A. *et al.* Whole-genome sequencing analysis of the cardiometabolic proteome. *Nat. Commun.* **11**, 6336 (2020).
12. Gilly, A. *et al.* Gene-based whole genome sequencing meta-analysis of 250 circulating proteins in three isolated European populations. *Mol. Metab.* **61**, 101509 (2022).
13. Kierczak, M. *et al.* Contribution of rare whole-genome sequencing variants to plasma protein levels and the missing heritability. *Nat. Commun.* **13**, 2532 (2022).
14. Dhindsa, R. S. *et al.* Influences of rare protein-coding genetic variants on the human plasma proteome in 50,829 UK Biobank participants. 2022.10.09.511476 Preprint at <https://doi.org/10.1101/2022.10.09.511476> (2022).
15. Png, G. *et al.* Population-wide copy number variation calling using variant call format files from 6,898 individuals. *Genet. Epidemiol.* **44**, 79–89 (2020).
16. Zheng, J. *et al.* Phenome-wide Mendelian randomization mapping the influence of the plasma proteome on complex diseases. *Nat. Genet.* **52**, 1122–1131 (2020).
17. Zuber, V. *et al.* Combining evidence from Mendelian randomization and colocalization: Review and comparison of approaches. *Am. J. Hum. Genet.* **109**, 767–782 (2022).
18. Kerimov, N. *et al.* A compendium of uniformly processed human gene expression and splicing quantitative trait loci. *Nat. Genet.* **53**, 1290–1299 (2021).
19. Leitsalu, L. *et al.* Cohort Profile: Estonian Biobank of the Estonian Genome Center, University of Tartu. *Int. J. Epidemiol.* **44**, 1137–1147 (2015).
20. Mitt, M. *et al.* Improved imputation accuracy of rare and low-frequency variants using population-specific high-coverage WGS-based imputation reference panel. *Eur. J. Hum. Genet. EJHG* **25**, 869–876 (2017).
21. Handsaker, R. E. *et al.* Large multiallelic copy number variations in humans. *Nat. Genet.* **47**, 296–303 (2015).
22. Lepamets, M. *et al.* Omics-informed CNV calls reduce false-positive rates and improve power for CNV-trait associations. *Hum. Genet. Genom. Adv.* **3**, 100133 (2022).
23. Assarsson, E. *et al.* Homogenous 96-Plex PEA immunoassay exhibiting high sensitivity, specificity, and excellent scalability. *PLoS One* **9**, e95192 (2014).
24. Bolger, A. M., Lohse, M. & Usadel, B. Trimmomatic: A flexible trimmer for Illumina sequence data. *Bioinform. Oxf. Engl.* **30**, 2114–2120 (2014).
25. Andrews, S. FastQC: A quality control tool for high throughput sequence data 2010. Available from: <https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>.
26. Dobin, A. *et al.* STAR: ultrafast universal RNA-seq aligner. *Bioinform. Oxf. Engl.* **29**, 15–21 (2013).
27. Robinson, M. D. & Oshlack, A. A scaling normalization method for differential expression analysis of RNA-seq data. *Genome Biol.* **11**, R25 (2010).
28. Robinson, M. D., McCarthy, D. J. & Smyth, G. K. edgeR: A Bioconductor package for differential expression analysis of digital gene expression data. *Bioinform. Oxf. Engl.* **26**, 139–140 (2010).
29. Lepik, K. *et al.* C-reactive protein upregulates the whole blood expression of CD59—An integrative analysis. *PLoS Comput. Biol.* **13**, e1005766 (2017).
30. Kang, H. M. *et al.* Variance component model to account for sample structure in genome-wide association studies. *Nat. Genet.* **42**, 348–354 (2010).
31. Chang, C. C. *et al.* Second-generation PLINK: Rising to the challenge of larger and richer datasets. *GigaScience* **4**, 7 (2015).
32. Boyle, A. P. *et al.* Annotation of functional variation in personal genomes using RegulomeDB. *Genome Res.* **22**, 1790–1797 (2012).
33. Moore, C. M., Jacobson, S. A. & Fingerlin, T. E. Power and sample size calculations for genetic association studies in the presence of genetic model misspecification. *Hum. Hered.* **84**, 256–271 (2020).
34. Yao, C. *et al.* Genome-wide mapping of plasma protein QTLs identifies putatively causal genes and pathways for cardiovascular disease. *Nat. Commun.* **9**, 3268 (2018).
35. Gao, X., Starmer, J. & Martin, E. R. A multiple testing correction method for genetic association studies using correlated single nucleotide polymorphisms. *Genet. Epidemiol.* **32**, 361–369 (2008).
36. 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).
37. Lê, S., Josse, J. & Husson, F. FactoMineR: An R package for multivariate analysis. *J. Stat. Softw.* **25**, 1–18 (2008).
38. Harrow, J. *et al.* GENCODE: The reference human genome annotation for The ENCODE Project. *Genome Res.* **22**, 1760–1774 (2012).
39. Franz, M. *et al.* GeneMANIA update 2018. *Nucleic Acids Res.* **46**, W60–W64 (2018).
40. Warde-Farley, D. *et al.* The GeneMANIA prediction server: Biological network integration for gene prioritization and predicting gene function. *Nucleic Acids Res.* **38**, W214–220 (2010).
41. Szklarczyk, D. *et al.* The STRING database in 2023: Protein–protein association networks and functional enrichment analyses for any sequenced genome of interest. *Nucleic Acids Res.* **51**, D638–D646 (2023).
42. Stark, C. *et al.* BioGRID: A general repository for interaction datasets. *Nucleic Acids Res.* **34**, D535–D539 (2006).
43. Oughtred, R. *et al.* The BioGRID database: A comprehensive biomedical resource of curated protein, genetic, and chemical interactions. *Protein Sci.* **30**, 187–200 (2021).
44. Wang, G., Sarkar, A., Carbonetto, P. & Stephens, M. A simple new approach to variable selection in regression, with application to genetic fine mapping. *J. R. Stat. Soc. Ser. B Stat. Methodol.* **82**, 1273–1300 (2020).
45. Zou, Y., Carbonetto, P., Wang, G. & Stephens, M. Fine-mapping from summary data with the “Sum of Single Effects” model. *PLoS Genet.* **18**, e1010299 (2022).
46. Di Tommaso, P. *et al.* Nextflow enables reproducible computational workflows. *Nat. Biotechnol.* **35**, 316–319 (2017).
47. Benner, C. *et al.* Prospects of fine-mapping trait-associated genomic regions by using summary statistics from genome-wide association studies. *Am. J. Hum. Genet.* **101**, 539–551 (2017).
48. Kamat, M. A. *et al.* PhenoScanner V2: An expanded tool for searching human genotype-phenotype associations. *Bioinformatics Oxf. Engl.* **35**, 4851–4853 (2019).
49. Staley, J. R. *et al.* PhenoScanner: A database of human genotype-phenotype associations. *Bioinformatics Oxf. Engl.* **32**, 3207–3209 (2016).
50. Folkersen, L. *et al.* Mapping of 79 loci for 83 plasma protein biomarkers in cardiovascular disease. *PLoS Genet.* **13**, e1006706 (2017).
51. Suhre, K. *et al.* Connecting genetic risk to disease end points through the human blood plasma proteome. *Nat. Commun.* **8**, 14357 (2017).
52. Chris, F. *et al.* The druggable genome and support for target identification and validation in drug development. *Sci. Transl. Med.* **9**, eaag1166 (2017).
53. Freshour, S. L. *et al.* Integration of the Drug-Gene Interaction Database (DGIdb 4.0) with open crowdsourcing efforts. *Nucleic Acids Res.* **49**, D1144–D1151 (2021).

54. Elsworth, B. *et al.* The MRC IEU OpenGWAS data infrastructure. 2020.08.10.244293 <https://www.biorxiv.org/content/https://doi.org/10.1101/2020.08.10.244293v1> (2020). <https://doi.org/10.1101/2020.08.10.244293>.
55. Giambartolomei, C. *et al.* Bayesian test for colocalisation between pairs of genetic association studies using summary statistics. *PLoS Genet.* **10**, e1004383 (2014).
56. Wallace, C. Eliciting priors and relaxing the single causal variant assumption in colocalisation analyses. *PLoS Genet.* **16**, e1008720 (2020).
57. Kasela, S. *et al.* Pathogenic implications for autoimmune mechanisms derived by comparative eQTL analysis of CD4+ versus CD8+ T cells. *PLoS Genet.* **13**, e1006643 (2017).
58. GTEx Consortium. The GTEx Consortium atlas of genetic regulatory effects across human tissues. *Science.* **369**, 1318–1330 (2020).
59. Bretherick, A. D. *et al.* Linking protein to phenotype with Mendelian Randomization detects 38 proteins with causal roles in human diseases and traits. *PLoS Genet.* **16**, e1008785 (2020).
60. Hemani, G. *et al.* The MR-Base platform supports systematic causal inference across the human phenome. *eLife* **7**, e34408 (2018).
61. Hemani, G., Tilling, K. & Smith, G. D. Orienting the causal relationship between imprecisely measured traits using GWAS summary data. *PLoS Genet.* **13**, e1007081 (2017).
62. Shabalin, A. A. Matrix eQTL: Ultra fast eQTL analysis via large matrix operations. *Bioinformatics* **28**, 1353–1358 (2012).
63. Zhu, Z. *et al.* Integration of summary data from GWAS and eQTL studies predicts complex trait gene targets. *Nat. Genet.* **48**, 481–487 (2016).
64. Hao, Z. *et al.* Rldeogram: Drawing SVG graphics to visualize and map genome-wide data on the idiograms. *PeerJ Comput. Sci.* **6**, e251 (2020).
65. Macdonald-Dunlop, E. *et al.* Mapping genetic determinants of 184 circulating proteins in 26,494 individuals to connect proteins and diseases. *medRxiv* 2021.08.03.21261494 (2021). <https://doi.org/10.1101/2021.08.03.21261494>.
66. Astle, W. J. *et al.* The allelic landscape of human blood cell trait variation and links to common complex disease. *Cell* **167**, 1415–1429.e19 (2016).
67. van der Harst, P. & Verweij, N. Identification of 64 novel genetic loci provides an expanded view on the genetic architecture of coronary artery disease. *Circ. Res.* **122**, 433–443 (2018).
68. Stahl, E. A. *et al.* Genome-wide association study meta-analysis identifies seven new rheumatoid arthritis risk loci. *Nat. Genet.* **42**, 508–514 (2010).
69. Interleukin-6 Receptor Mendelian Randomisation Analysis (IL6R MR) Consortium *et al.* The interleukin-6 receptor as a target for prevention of coronary heart disease: A mendelian randomisation analysis. *Lancet Lond. Engl.* **379**, 1214–1224 (2012).
70. Suhre, K., McCarthy, M. I. & Schwenk, J. M. Genetics meets proteomics: Perspectives for large population-based studies. *Nat. Rev. Genet.* **22**, 19–37 (2021).
71. Zhang, J. *et al.* Plasma proteome analyses in individuals of European and African ancestry identify cis-pQTLs and models for proteome-wide association studies. *Nat. Genet.* **54**, 593–602 (2022).
72. Uhlén, M. *et al.* The human secretome. *Sci. Signal.* **12**, eaaz0274 (2019).
73. Ornitz, D. M. *et al.* Receptor specificity of the fibroblast growth factor family. *J. Biol. Chem.* **271**, 15292–15297 (1996).
74. Vatner, S. F. FGF induces hypertrophy and angiogenesis in hibernating myocardium. *Circ. Res.* **96**, 705–707 (2005).
75. Nikpay, M. *et al.* A comprehensive 1,000 Genomes-based genome-wide association meta-analysis of coronary artery disease. *Nat. Genet.* **47**, 1121–1130 (2015).
76. Panoutsopoulou, K. *et al.* Genetic characterization of Greek population isolates reveals strong genetic drift at missense and trait-associated variants. *Nat. Commun.* **5**, 5345 (2014).
77. Tominaga, K. *et al.* IL-12 synergizes with IL-18 or IL-1beta for IFN-gamma production from human T cells. *Int. Immunol.* **12**, 151–160 (2000).
78. Kay, C., Wang, R., Kirkby, M. & Man, S. M. Molecular mechanisms activating the NAIP-NLRC4 inflammasome: Implications in infectious disease, autoinflammation, and cancer. *Immunol. Rev.* **297**, 67–82 (2020).
79. Emmerson, P. J. *et al.* The metabolic effects of GDF15 are mediated by the orphan receptor GFRAL. *Nat. Med.* **23**, 1215–1219 (2017).
80. Hsu, J.-Y. *et al.* Non-homeostatic body weight regulation through a brainstem-restricted receptor for GDF15. *Nature* **550**, 255–259 (2017).
81. Yang, L. *et al.* GFRAL is the receptor for GDF15 and is required for the anti-obesity effects of the ligand. *Nat. Med.* **23**, 1158–1166 (2017).
82. Mahajan, A. *et al.* Refining the accuracy of validated target identification through coding variant fine-mapping in type 2 diabetes. *Nat. Genet.* **50**, 559–571 (2018).
83. Wu, Y. *et al.* Genome-wide association study of medication-use and associated disease in the UK Biobank. *Nat. Commun.* **10**, 1891 (2019).
84. Vujkovic, M. *et al.* Discovery of 318 new risk loci for type 2 diabetes and related vascular outcomes among 1.4 million participants in a multi-ancestry meta-analysis. *Nat. Genet.* **52**, 680–691 (2020).
85. Breit, S. N. *et al.* The TGF- $\beta$  superfamily cytokine, MIC-1/GDF15: A pleiotropic cytokine with roles in inflammation, cancer and metabolism. *Growth Factors Chur Switz.* **29**, 187–195 (2011).
86. Lemmälä, S. *et al.* Integrated analyses of growth differentiation factor-15 concentration and cardiometabolic diseases in humans. *eLife* **11**, e76272 (2022).
87. Momozawa, Y. & Mizukami, K. Unique roles of rare variants in the genetics of complex diseases in humans. *J. Hum. Genet.* **66**, 11–23 (2021).
88. He, B., Shi, J., Wang, X., Jiang, H. & Zhu, H.-J. Genome-wide pQTL analysis of protein expression regulatory networks in the human liver. *BMC Biol.* **18**, 97 (2020).
89. Robins, C. *et al.* Genetic control of the human brain proteome. *Am. J. Hum. Genet.* **108**, 400–410 (2021).
90. Sollis, E. *et al.* The NHGRI-EBI GWAS Catalog: knowledgebase and deposition resource. *Nucleic Acids Res.* **51**, D977–D985 (2023).

## Acknowledgements

The authors would like to thank Hanna Maria Kariis for helpful comments on the manuscript. Data analyses with Estonian datasets were carried out in part in the High-Performance Computing Center of University of Tartu.

## Author contributions

A.K.: conceptualisation, data curation, formal analysis, investigation, methodology, software, visualisation, writing—original draft, writing—review and editing; M.J.: conceptualisation, data curation, formal analysis, investigation, methodology, software, visualisation, writing—review and editing; K.L.: formal analysis, resources, writing—review and editing; V.K.: resources; M.K.: resources; K.A.: resources, writing—review and editing; E.B.R.T.: resources; R.M.: methodology, supervision; T.E.: conceptualisation, data curation, funding acquisition, methodology, project administration, resources, supervision, writing—review and editing; U.V.:

conceptualisation, data curation, methodology, project administration, software, supervision, writing—review and editing. All authors read and approved the final manuscript.

### Funding

AK, TE and UV were supported by the Estonian Research Council grant PUT (PRG1291). MJ and RM were supported by the Estonian Research Council grant PUT (PRG1911). MJ and MK were supported by the European Union through the European Regional Development Fund (Project No. 2014-2020.4.01.16-0125). MJ and RM received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 101016775. KA was supported by a grant from the Estonian Research Council (PSG415). UV was supported by the European Regional Development Fund and the programme Mobilitas Pluss (MOBTP108). This study was funded by the European Union through the European Regional Development Fund Project No. 2014-2020.4.01.15-0012 GENTRANSMED.

### Competing interests

The authors declare no competing interests.

### Additional information

**Supplementary Information** The online version contains supplementary material available at <https://doi.org/10.1038/s41598-024-57966-3>.

**Correspondence** and requests for materials should be addressed to A.K., T.E. or U.V.

**Reprints and permissions information** is available at [www.nature.com/reprints](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 licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence 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 licence, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2024

---

### Estonian Biobank Research Team

Andres Metspalu<sup>1</sup>, Lili Milani<sup>1</sup>, Tõnu Esko<sup>1,7</sup>✉ & Reedik Mägi<sup>1</sup>, Mari Nelis<sup>1</sup> & Georgi Hudjashov<sup>1</sup>