



Sun, B. B., Maranville, J. C., Peters, J. E., Stacey, D., Staley, J. R., Blackshaw, J., ... Butterworth, A. S. (2018). Genomic atlas of the human plasma proteome. *Nature*, 558(7708), 73-79. <https://doi.org/10.1038/s41586-018-0175-2>

Peer reviewed version

Link to published version (if available):  
[10.1038/s41586-018-0175-2](https://doi.org/10.1038/s41586-018-0175-2)

[Link to publication record in Explore Bristol Research](#)  
PDF-document

This is the author accepted manuscript (AAM). The final published version (version of record) is available online via Nature Publishing Group at <https://www.nature.com/articles/s41586-018-0175-2> . Please refer to any applicable terms of use of the publisher.

## **University of Bristol - Explore Bristol Research**

### **General rights**

This document is made available in accordance with publisher policies. Please cite only the published version using the reference above. Full terms of use are available:  
<http://www.bristol.ac.uk/pure/about/ebr-terms>

# Genomic atlas of the human plasma proteome

Benjamin B. Sun<sup>1\*</sup>, Joseph C. Maranville<sup>2\*</sup>, James E. Peters<sup>1,3\*</sup>, David Stacey<sup>1</sup>, James R. Staley<sup>1</sup>, James Blackshaw<sup>1</sup>, Stephen Burgess<sup>1,4</sup>, Tao Jiang<sup>1</sup>, Ellie Paige<sup>1,5</sup>, Praveen Surendran<sup>1</sup>, Clare Oliver-Williams<sup>1,6</sup>, Mihir A. Kamat<sup>1</sup>, Bram P. Prins<sup>1</sup>, Sheri K. Wilcox<sup>7</sup>, Erik S. Zimmerman<sup>7</sup>, An Chi<sup>2</sup>, Narinder Bansal<sup>1,8</sup>, Sarah L. Spain<sup>9</sup>, Angela M. Wood<sup>1</sup>, Nicholas W. Morrell<sup>10</sup>, John R. Bradley<sup>11</sup>, Nebojsa Janjic<sup>7</sup>, David J. Roberts<sup>12,13</sup>, Willem H. Ouwehand<sup>3,14,15,16,17</sup>, John A. Todd<sup>18</sup>, Nicole Soranzo<sup>3,14,16,17</sup>, Karsten Suhre<sup>19</sup>, Dirk S. Paul<sup>1</sup>, Caroline S. Fox<sup>2</sup>, Robert M. Plenge<sup>2</sup>, John Danesh<sup>1,3,16,17</sup>, Heiko Runz<sup>2\*</sup>, Adam S. Butterworth<sup>1,17\*</sup>

1. MRC/BHF Cardiovascular Epidemiology Unit, Department of Public Health and Primary Care, University of Cambridge, Cambridge CB1 8RN, UK.
2. MRL, Merck & Co., Inc., Kenilworth, New Jersey, USA.
3. British Heart Foundation Cambridge Centre of Excellence, Division of Cardiovascular Medicine, Addenbrooke's Hospital, Cambridge CB2 0QQ, UK.
4. MRC Biostatistics Unit, University of Cambridge, Cambridge CB2 0SR, UK.
5. National Centre for Epidemiology and Population Health, The Australian National University, Canberra, ACT, Australia.
6. Homerton College, Cambridge, CB2 8PH, UK.
7. SomaLogic Inc., Boulder, Colorado 80301, USA.
8. Perinatal Institute, Birmingham B15 3BU, UK.
9. Wellcome Trust Sanger Institute, Wellcome Trust Genome Campus, Hinxton, Cambridge CB10 1RQ, UK.
10. Division of Respiratory Medicine, Department of Medicine, University of Cambridge, Cambridge CB2 0QQ, UK.
11. NIHR Cambridge Biomedical Research Centre / BioResource, Cambridge University Hospitals, Cambridge CB2 0QQ, UK.
12. National Health Service (NHS) Blood and Transplant and Radcliffe Department of Medicine, NIHR Oxford Biomedical Research Centre, University of Oxford, John Radcliffe Hospital, Oxford OX3 9DU, UK.
13. BRC Haematology Theme and Department of Haematology, Churchill Hospital, Oxford OX3 7LE, UK.
14. Department of Haematology, University of Cambridge, Cambridge Biomedical Campus, Long Road, Cambridge CB2 0PT, UK.
15. National Health Service (NHS) Blood and Transplant, Cambridge Biomedical Campus, Cambridge CB2 0PT, UK.
16. Department of Human Genetics, Wellcome Trust Sanger Institute, Wellcome Trust Genome Campus, Hinxton, Cambridge CB10 1RQ, UK.
17. NIHR Blood and Transplant Research Unit in Donor Health and Genomics, Department of Public Health and Primary Care, University of Cambridge, Cambridge CB1 8RN, UK.
18. JDRF/Wellcome Trust Diabetes and Inflammation Laboratory, Wellcome Trust Centre for Human Genetics, Nuffield Department of Medicine, NIHR Oxford Biomedical Research Centre, University of Oxford, Oxford OX3 7BN, UK.
19. Department of Physiology and Biophysics, Weill Cornell Medicine - Qatar, PO 24144 Doha, Qatar.

\* These authors contributed equally to this work.

Corresponding authors: [asb38@medschl.cam.ac.uk](mailto:asb38@medschl.cam.ac.uk) (A.S.B.), [jd292@medschl.cam.ac.uk](mailto:jd292@medschl.cam.ac.uk) (J.D.)

## Summary

Although plasma proteins play important roles in biological processes and are the direct targets of many drugs, there is limited knowledge of the genetic factors determining inter-individual variation in plasma protein levels. Here we characterize the genetic architecture of the human plasma proteome in healthy blood donors from the INTERVAL study. We identify 1,927 genetic associations with 1,478 proteins, a 4-fold increase on existing knowledge, including *trans* associations for 1,104 proteins. To understand consequences of perturbations in plasma protein levels, we apply an integrated approach that links genetic variation with biological pathway, disease, and drug databases. We show overlap of pQTL with eQTL, as well as with disease-associated loci, and show support for causal roles for protein biomarkers in disease using Mendelian randomisation analysis. By linking genetic factors to disease via specific proteins, our analyses suggests potential therapeutic targets, opportunities for matching existing drugs with new disease indications, and potential safety concerns for drugs under development.

Plasma proteins play key roles in various biological processes including signalling, transport, growth, repair, and defence against infection. They are frequently dysregulated in disease and are important drug targets. Identifying factors that determine inter-individual protein variability should, therefore, furnish biological and medical insights<sup>1</sup>. Despite evidence of the heritability of plasma protein abundance<sup>2</sup>, however, systematic assessment of how genetic variation influences plasma protein levels has been limited<sup>3-5</sup>. Studies have examined intracellular ‘protein quantitative trait loci’ (pQTLs)<sup>6,7</sup>, but they have tended to be small and involved cell lines rather than primary human tissues.

Here we create and interrogate a genetic atlas of the human plasma proteome, using an expanded version of an aptamer-based multiplex protein assay (SOMAscan)<sup>8</sup> to quantify 3,622 plasma proteins in 3,301 healthy participants from the INTERVAL study, a genomic bioresource of 50,000 whole blood donors from 25 centres across England recruited into a randomised trial of blood donation frequency<sup>9,10</sup>. We identify 1,927 genotype-protein associations, including *trans*-associated loci for 1,104 proteins, providing new understanding of the genetic control of protein regulation. 88 pQTLs overlap with disease susceptibility loci, suggesting the molecular effects of disease-associated variants. Using the principle of Mendelian randomisation<sup>11</sup>, we find evidence to support causal roles in disease for several protein pathways, and cross-reference our data with disease and drug databases to highlight potential therapeutic targets.

## RESULTS

### Genetic architecture of the plasma proteome

We performed genome-wide testing of 10.6 million imputed autosomal variants against levels of 2,994 plasma proteins in 3,301 European-ancestry individuals ([Methods](#), [Extended Data Figure 1](#)). We demonstrated robustness of protein measurements in several ways ([Supplementary Note](#)), including: highly consistent measurements in replicate samples; temporal consistency of protein levels within individuals over two years ([Extended Data Figure 2b](#)); replication of known associations with non-genetic factors ([Supplementary Tables 1-2](#)). To assess potential off-target cross-reactivity, we tested 920 aptamers (“SOMAmers”) for detection of proteins with  $\geq 40\%$  sequence homology to the target protein ([Methods](#)).

Although 126 (14%) SOMAmers showed comparable binding with a homologous protein ([Supplementary Table 3](#)), nearly half of these were binding to alternative forms of the same protein.

We found 1,927 significant ( $p < 1.5 \times 10^{-11}$ ) associations between 1,478 proteins and 764 genomic regions ([Figure 1a](#), [Supplementary Table 4](#), [Supplementary Video 1](#)), with 89% of pQTLs previously unreported. Of the 764 associated regions, 502 (66%) had local-acting ('*cis*') associations only, 228 (30%) *trans* only, and 34 (4%) both *cis* and *trans* ([Supplementary Note Table 1](#)). 95% and 87% of *cis* pQTL variants were located within 200Kb and 100Kb, respectively, of the relevant gene's canonical transcription start site (TSS) ([Figure 1b](#)), and 44% were within the gene itself. The *p*-values for *cis* associations increased with distance from the TSS, mirroring findings for expression QTLs (eQTLs)<sup>12</sup>. Of proteins with a significant pQTL, 88% had either *cis* (n=374) or *trans* (n=925) associations only, while 12% (n=179) had both ([Supplementary Note Table 1](#)). The majority of significantly associated proteins (75%; n=1,113) had a single pQTL, while 20% had two and 5% had >2 ([Figure 1c](#)). To detect multiple independent signals at the same locus we used stepwise conditional analysis, identifying 2,658 conditionally significant associations ([Supplementary Table 5](#)). Of the 1,927 locus-protein associations, 414 (21%) had multiple conditionally significant signals ([Figure 1d](#)), of which 255 were *cis*.

We tested replication of 163 pQTLs in 4,998 individuals using an alternative protein assay (Olink, [Methods](#))<sup>13</sup>. Effect-size estimates were strongly correlated between the SOMAscan and Olink platforms ( $r=0.83$ ; [Extended Data Figure 2c](#)). 106/163 (65% overall; 81% *cis*, 52% *trans*) pQTLs replicated after Bonferroni correction ([Supplementary Tables 4,6](#)). The lower replication rate of *trans* signals may reflect various factors, including differences between protein assays (e.g., detection of free versus complexed proteins) and the higher 'biological prior' for *cis* associations.

Of 1,927 pQTLs, 549 (28.5%) were *cis*-acting ([Supplementary Table 4](#)). Genetic variants that change protein structure may result in apparent *cis* pQTLs due to altered aptamer-binding rather than true quantitative differences in protein levels. We found evidence against such artefactual associations for 371

(67.6%) *cis* pQTLs ([Methods](#), [Supplementary Tables 4, 7-8](#)). Results were materially unchanged when we repeated downstream analyses excluding pQTLs without evidence against binding effects.

The median variation in protein levels explained by pQTLs was 5.8% (interquartile range: 2.6-12.4%, [Figure 1e](#)). For 193 proteins, genetic variants explained >20% of the variation. There was a strong inverse relationship between effect-size and minor allele frequency (MAF) ([Figure 1f](#)), consistent with previous genome-wide association studies (GWAS) of quantitative traits<sup>7,10,14</sup>. We found 23 and 208 associations with rare (MAF <1%) variants and low-frequency (MAF 1-5%) variants, respectively ([Supplementary Table 4](#)). Of the 36 strongest associations (per-allele effect-size >1.5 standard deviations), 29 were with rare or low-frequency variants.

Both *cis* and *trans* pQTLs were strongly enriched for missense variants ( $p < 0.0001$ ) and for location in 3' untranslated ( $p = 0.0025$ ) or splice sites ( $p = 0.0004$ ) ([Figure 1g](#), [Extended Data Figure 3a](#)). We found  $\geq 3$ -fold enrichment ( $p < 5 \times 10^{-5}$ ) of pQTLs at features indicative of transcriptional activation in blood cells and at hepatocyte regulatory elements, consistent with the liver's role in protein synthesis and secretion ([Methods](#), [Extended Data Figure 4](#), [Supplementary Table 9](#)).

## Overlap of eQTLs and pQTLs

To help evaluate the extent to which genetic associations with plasma protein levels are driven by effects at the transcription level rather than other mechanisms (e.g., altered protein clearance or secretion), we cross-referenced our *cis* pQTLs with previous eQTL studies ([Supplementary Table 10](#)), initially defining overlap between an eQTL and pQTL as high linkage disequilibrium (LD) ( $r^2 \geq 0.8$ ) between the lead pQTL and eQTL variants. 40% (n=224) of *cis* pQTLs were eQTLs for the same gene in  $\geq 1$  tissue or cell-type ([Supplementary Table 8](#)). The greatest overlaps were in whole blood (n=117), liver (n=70) and lymphoblastoid cell-lines (LCLs) (n=52), consistent with biological expectation, but also likely driven by the larger eQTL study sample sizes for these cell-types. To examine whether the same causal variant was likely to underlie overlapping eQTLs and pQTLs, we performed colocalisation testing ([Methods](#)). Of 228

non-*HLA* pQTLs for which testing was possible, colocalisation in  $\geq 1$  tissue or cell-type was highly likely (posterior probability[PP] $>0.8$ ) in 179 (78.5%) and the most likely explanation (PP $>0.5$ ) in 197 (86.4%) (Supplementary Table 8). *Cis* pQTLs were significantly enriched for eQTLs for the corresponding gene ( $p < 0.0001$ ) (Methods, Supplementary Table 11). To address the converse (i.e., to what extent are eQTLs also pQTLs), we selected well-powered eQTL studies in relevant tissues (whole blood, LCLs, liver and monocytes<sup>15-18</sup>). Of the strongest *cis* eQTLs ( $p < 1.5 \times 10^{-11}$ ) in whole blood, LCLs, liver and monocytes, 12.2%, 21.3%, 14.8% and 14.7%, respectively, were plasma *cis* pQTLs .

Comparisons between eQTL and pQTL studies have inherent limitations, including differences in the tissues, sample sizes and technological platforms used. Moreover, plasma protein levels may not reflect levels within tissues or cells. Nevertheless, our data suggest that genetic effects on plasma protein abundance are often, but not exclusively, driven by regulation of mRNA. *Cis* pQTLs without corresponding *cis* eQTLs may reflect genetic effects on processes other than transcription, including protein degradation, binding, secretion, or clearance from circulation.

## ***Trans* pQTLs identify pathways to disease**

Of the 764 protein-associated regions, 262 had *trans* associations with 1,104 proteins (Supplementary Table 4, 12). There was no enrichment of cross-reactivity in SOMAmers with a *trans* pQTL versus those without (Supplementary Note). We replicated known *trans* associations including *TMPRSS6* with transferrin receptor protein 1<sup>19</sup> and *SORT1* with granulins<sup>20</sup> and identified several novel biologically plausible *trans* associations (Supplementary Table 13), including known or presumed ligand:receptor pairs (e.g., the *CD320* locus, encoding the transcobalamin receptor, was associated with transcobalamin-2 levels).

Most (82%) *trans* loci were associated with  $<4$  proteins, but 12 ‘hotspot’ regions were associated with  $>20$  (Figure 1a, Extended Data Figure 3b), including well-known pleiotropic loci (e.g., *ABO*, *CFH*, *APOE*, *KLKB1*) and loci associated with many correlated proteins (e.g., the *ZFPM2* locus encoding the transcription factor FOG2). Similar pleiotropy at these loci has been seen in other plasma pQTL studies<sup>3-5</sup>, albeit with

fewer proteins due to limited assay breadth. rs28929474:T in *SERPINA1* was associated with 13 proteins at  $p < 1.5 \times 10^{-11}$  and a further six at  $p < 5 \times 10^{-8}$  (Figure 2). This missense variant (the ‘Z-allele’, p.Glu366Lys) results in defective secretion and intracellular accumulation of alpha1-antitrypsin (A1AT), an anti-protease. ZZ homozygotes have deficiency of circulating A1AT and increased risk of emphysema, liver cirrhosis and vasculitis. The ‘protease-antiprotease’ hypothesis posits that these clinical manifestations result from unchecked protease activity. However, our discovery of multiple *trans*-associated proteins at this locus highlights additional pathways potentially relevant to pathogenesis, a hypothesis supported by accumulating data<sup>21</sup>.

GWAS have identified thousands of loci associated with common diseases, but the mechanisms by which most variants influence disease susceptibility await discovery. To identify intermediate links between genotype and disease, we overlapped pQTLs with disease-associated variants from GWAS. 88 of our sentinel pQTL variants were in high LD ( $r^2 \geq 0.8$ ) with sentinel disease-associated variants (Supplementary Table 14), including 30 with *cis* associations, 54 with *trans*, and 4 with both. Since some genetic loci are associated with multiple diseases, these 88 genetic loci represent 253 distinct genotype-disease associations. Overlap of a pQTL and a disease association signal does not necessarily imply that the same genetic variant underlies both traits, since there may be distinct causal variants for each trait that are in LD. We therefore performed colocalisation testing (Methods). Of 108 non-MHC locus-disease associations for which testing was possible, colocalisation was highly likely ( $PP > 0.8$ ) for 96 (88.9%), and the most likely explanation ( $PP > 0.5$ ) for 106 (98.1%) (Supplementary Table 14).

*Trans* pQTLs that overlap with disease associations can highlight previously unsuspected candidate proteins through which genetic loci may influence disease risk. To help identify such candidates, we applied the ProGeM framework<sup>22</sup> (Methods, Supplementary Table 12, Extended Data Figure 5). We show that an inflammatory bowel disease (IBD) risk allele<sup>23</sup> (rs3197999:A, missense p.Arg703Cys) in *MST1* on chromosome 3, that decreases plasma MST1 levels<sup>24</sup>, is a *trans* pQTL for eight additional proteins (Supplementary Table 4, Figure 3). Notably, genes that encode three of these proteins (*PRDMI*, *FASLG*,



and *DOCK9*) each lie within 500kb of IBD GWAS loci where the causal gene is ambiguous<sup>25</sup>. For instance, the IBD-associated variant rs6911490 lies on chromosome 6 in the intergenic region between *PRDMI* (encoding BLIMP1, a master regulator of immune cell differentiation) and *ATG5* (involved in autophagy) (Figure 3c). Neither fine-mapping nor eQTL colocalisation analyses have unequivocally resolved the causal gene at this locus<sup>25</sup>; both *PRDMI* and *ATG5* are plausible candidates. Our data provide support for *PRDMI*.

Anti-neutrophil cytoplasmic antibody-associated vasculitis (AAV) is an autoimmune disease characterised by vascular inflammation and autoantibodies to the neutrophil proteases proteinase-3 (PR3) or myeloperoxidase. GWAS reveal distinct genetic signals according to antibody specificity<sup>26</sup>, with variants near *PRTN3* (encoding PR3) and at the Z-allele of *SERPINA1* (encoding alpha1-antitrypsin, an inhibitor of PR3) associated specifically with PR3-antibody positive AAV. The SOMAscan assay has two SOMAmers targeting PR3; we identified a *cis* pQTL signal immediately upstream of *PRTN3* for both, and replicated it with the Olink assay (Supplementary Table 4, Figures 4a-b). Conditional analysis revealed multiple independently associated variants (Supplementary Table 5), one of which (rs7254911) was in high LD with the PR3+ vasculitis tag SNPs (Supplementary Note). We show that the vasculitis risk allele at *PRTN3* is associated with higher plasma levels of PR3 (Supplementary Note Table 4).

For one PR3 SOMAmer, we also found a *trans* pQTL at *SERPINA1*, with the Z-allele associating with lower plasma PR3 (Figure 4a). To understand the SOMAmer-specific nature of this signal, we assayed the relative affinity of these SOMAmers for the free and complexed states of PR3 and A1AT (which binds and inhibits proteases including PR3). We found that the SOMAmer showing *cis* and *trans* associations predominantly measures the PR3:A1AT complex rather than free PR3, whereas the SOMAmer with only *cis* association measures both the free and complexed forms. Importantly, neither SOMAmer bound free A1AT, demonstrating that the *SERPINA1* pQTL did not reflect non-specific cross-reactivity (Supplementary Note).

These data show that the vasculitis risk allele at *PRTN3* increases total PR3 plasma levels, consistent with its effect on *PRTN3* mRNA abundance in whole blood in GTEx data<sup>27</sup>. The *SERPINA1* Z-allele results in a

reduced proportion of PR3 bound to A1AT. We thus demonstrate how altered availability of PR3, conferred by two independent genetic mechanisms, is a key susceptibility factor for breaking immune tolerance to PR3 and the development of PR3+ vasculitis ([Figure 4c](#)).

## Causal evaluation of candidate proteins in disease

Association of plasma protein levels with disease risk does not necessarily imply causation. To help establish causality, we used Mendelian randomisation (MR) analysis<sup>11</sup> ([Extended Data Figure 6](#)). The concept is that if a genetic variant that specifically influences levels of a protein is also associated with disease risk, then this provides evidence of the protein's causal role. For example, serum levels of PSP-94 (MSMB) are lower in patients with prostate cancer<sup>28</sup>, but it is debated whether this association is correlative or causal. We identified a *cis* pQTL associated with lower PSP-94 plasma levels that overlaps with the prostate cancer susceptibility variant rs10993994<sup>29</sup>, supporting a protective role for PSP-94 in prostate cancer ([Supplementary Table 14](#)).

Next, we leveraged multi-variant MR analysis methods to distinguish causal proteins among multiple plausible candidates, exemplified by the *IL1RL1-IL18R1* locus, which is associated with multiple immune-mediated diseases including atopic dermatitis<sup>30</sup>. We identified four proteins that each had *cis* pQTLs at this locus ([Supplementary Table 4](#)), and created a genetic score for each protein ([Methods](#)). Initial 'one-protein-at-a-time' analysis identified associations of the scores for IL18R1 ( $p=9.3 \times 10^{-72}$ ) and IL1RL1 ( $p=5.7 \times 10^{-27}$ ) with atopic dermatitis risk ([Figure 5a](#)), and a weak association for IL1RL2 ( $p=0.013$ ). We then mutually adjusted these associations for one another to account for the effects of the variants on multiple proteins. While the association of IL18R1 remained significant ( $p=1.5 \times 10^{-28}$ ), the association of IL1RL1 ( $p=0.01$ ) was attenuated. In contrast, the association of IL1RL2 ( $p=1.1 \times 10^{-69}$ ) became much stronger, suggesting that IL1RL2 and IL18R1 underlie atopic dermatitis risk at this locus.

MMP-12 plays a key role in lung tissue damage, and MMP-12 inhibitors are being tested for chronic obstructive pulmonary disease<sup>31</sup>. We created a multi-allelic genetic score that explains 14% of the variation

in plasma MMP-12 levels ([Methods](#)). Observational studies reveal an association of higher levels of plasma MMP-12 with recurrent cardiovascular events<sup>32</sup>, stimulating interest in MMP-12 inhibitors for cardiovascular disease. In contrast, we found that genetic predisposition to higher MMP-12 levels is associated with *decreased* coronary disease risk ( $p=2.8 \times 10^{-13}$ ) ([Figure 5b](#)) and *decreased* large artery atherosclerotic stroke risk<sup>33</sup>. Understanding the discordance between the observational epidemiology and the genetic risk score will be important given the therapeutic interest in this target.

## Drug target prioritisation

Drugs directed at therapeutic targets implicated by human genetic data have a greater likelihood of success<sup>34</sup>. Of the proteins for which we identified a pQTL, 244 (17%) are established drug targets in the Informa Pharmaprojects database ([Supplementary Table 15](#)). 31 pQTLs for drug target proteins were highly likely to colocalise (posterior probability > 0.8) with a GWAS disease locus, including some that are targets of approved drugs such as tocilizumab (anti-IL6R) and ustekinumab (anti-IL12/23) ([Supplementary Table 16a](#)).

To identify additional indications for existing drugs, we investigated disease associations of pQTLs for proteins already targeted by licensed drugs. Our results suggest potential drug ‘re-purposing’ opportunities. For example, we identified a *cis* pQTL for RANK (encoded by *TNFRSF11A*) at rs884205, a variant associated with Paget’s disease<sup>35</sup>, a condition characterised by excessive bone turnover, deformity and fracture ([Supplementary Table 16b](#)). Standard Paget’s disease treatment is osteoclast inhibition with bisphosphonates, originally developed as anti-osteoporotic drugs. Denosumab, another anti-osteoporosis drug, is a monoclonal antibody targeting RANKL, the ligand for RANK. Our data suggest denosumab may be an alternative for Paget’s disease when bisphosphonates are contra-indicated, a hypothesis supported by clinical case reports<sup>36</sup>.

Next we evaluated targets of drugs currently under development. Drugs targeting GPIBA, the receptor for von Willebrand factor, are in pre-clinical development as anti-thrombotic agents and in phase 2 trials for

thrombotic thrombocytopenic purpura. We found a *cis* pQTL associated with both higher GP1BA abundance and higher platelet count, suggesting a link between GP1BA and platelet count ([Supplementary Table 16](#)). Furthermore, we identified a *trans* pQTL for GP1BA at the *SH2B3/BRAP* locus, which colocalised with associations with platelet count<sup>10</sup>, myocardial infarction and stroke ([Supplementary Table 16b](#)). The risk allele for cardiovascular disease increases both plasma GP1BA and platelet count, suggesting GP1BA influences vascular risk via platelets. Collectively, these results support targeting GP1BA in conditions characterised by platelet aggregation such as arterial thrombosis. More generally, our data provide a substrate for generating hypotheses about potential therapeutic targets through linking genetic factors to disease via specific proteins.

## DISCUSSION

This study elucidates the genetic control of the human plasma proteome and uncovers intermediate molecular pathways connecting the genome to disease endpoints. We applied our discoveries to evaluate causal roles for proteins in human disease using the principle of Mendelian randomisation (MR). Proteins provide an ideal paradigm for MR analysis because they are under proximal genetic control. However, application of protein-based MR has been constrained by limited availability of suitable genetic instruments, a bottleneck remedied by our data. Our study provides a resource for studying complex traits and an example for application of novel bioassay technologies to population biobanks.

## REFERENCES

1. Albert, F. W. & Kruglyak, L. The role of regulatory variation in complex traits and disease. *Nat. Rev. Genet.* **16**, 197–212 (2015).
2. Liu, Y. *et al.* Quantitative variability of 342 plasma proteins in a human twin population. *Mol. Syst. Biol.* **11**, 786 (2015).
3. Suhre, K. *et al.* Connecting genetic risk to disease end points through the human blood plasma proteome. *Nat. Commun.* **8**, 14357 (2017).
4. Yao, C. *et al.* Genome-wide association study of plasma proteins identifies putatively causal genes, proteins, and pathways for cardiovascular disease. *bioRxiv* (2017). doi:10.1101/136523
5. de Vries, P. S. *et al.* Whole-genome sequencing study of serum peptide levels: the Atherosclerosis Risk in Communities study. *Hum. Mol. Genet.* **26**, 3442–3450 (2017).
6. Wu, L. *et al.* Variation and genetic control of protein abundance in humans. *Nature* **499**, 79–82 (2013).
7. Battle, A. *et al.* Impact of regulatory variation from RNA to protein. *Science*. **347**, 644–7 (2015).
8. Rohloff, J. C. *et al.* Nucleic acid ligands with protein-like side chains: modified aptamers and their use as diagnostic and therapeutic agents. *Mol. Ther. Nucleic Acids* **3**, e201 (2014).
9. Di Angelantonio, E. *et al.* Efficiency and safety of varying the frequency of whole blood donation (INTERVAL): a randomised trial of 45 000 donors. *Lancet* **390**, 2360–2371 (2017).
10. 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).
11. Burgess, S. *et al.* Using published data in Mendelian randomization: a blueprint for efficient identification of causal risk factors. *Eur. J. Epidemiol.* **30**, 543–52 (2015).
12. Stranger, B. E. *et al.* Patterns of cis regulatory variation in diverse human populations. *PLoS Genet.* **8**, e1002639 (2012).
13. Lundberg, M., Eriksson, A., Tran, B., Assarsson, E. & Fredriksson, S. Homogeneous antibody-based proximity extension assays provide sensitive and specific detection of low-abundant proteins in human blood. *Nucleic Acids Res.* **39**, e102 (2011).
14. Walter, K. *et al.* The UK10K project identifies rare variants in health and disease. *Nature* **526**, 82–90 (2015).
15. Westra, H.-J. *et al.* Systematic identification of trans eQTLs as putative drivers of known disease associations. *Nat Genet* **45**, 1238–1243 (2013).
16. Lappalainen, T. *et al.* Transcriptome and genome sequencing uncovers functional variation in humans. *Nature* **501**, 506–511 (2013).
17. Schadt, E. E. *et al.* Mapping the genetic architecture of gene expression in human liver. *PLoS Biol.* **6**, e107 (2008).

18. Zeller, T. *et al.* Genetics and beyond – the transcriptome of human monocytes and disease susceptibility. *PLoS One* **5**, e10693 (2010).
19. Nai, A. *et al.* TMPRSS6 rs855791 modulates hepcidin transcription in vitro and serum hepcidin levels in normal individuals. *Blood* **118**, 4459-62 (2011).
20. Carrasquillo, M. M. *et al.* Genome-wide screen identifies rs646776 near sortilin as a regulator of progranulin levels in human plasma. *Am. J. Hum. Genet.* **87**, 890–897 (2010).
21. Gooptu, B., Dickens, J. A. & Lomas, D. A. The molecular and cellular pathology of  $\alpha_1$ -antitrypsin deficiency. *Trends Mol. Med.* **20**, 116–27 (2014).
22. Stacey, D. *et al.* ProGeM: A framework for the prioritisation of candidate causal genes at molecular quantitative trait loci. *bioRxiv* 230094 (2017). doi:10.1101/230094
23. Liu, J. Z. *et al.* Association analyses identify 38 susceptibility loci for inflammatory bowel disease and highlight shared genetic risk across populations. *Nat. Genet.* **47**, 979–986 (2015).
24. Di Narzo, A. F. *et al.* High-throughput characterization of blood serum proteomics of IBD patients with respect to aging and genetic factors. *PLoS Genet.* **13**, e1006565 (2017).
25. Huang, H. *et al.* Fine-mapping inflammatory bowel disease loci to single-variant resolution. *Nature* **547**, 173–178 (2017).
26. Lyons, P. A. *et al.* Genetically distinct subsets within ANCA-associated vasculitis. *N. Engl. J. Med.* **367**, 214–223 (2012).
27. Aguet, F. *et al.* Genetic effects on gene expression across human tissues. *Nature* **550**, 204–213 (2017).
28. Grönberg, H. *et al.* Prostate cancer screening in men aged 50–69 years (STHLM3): a prospective population-based diagnostic study. *Lancet Oncol.* **16**, 1667–1676 (2015).
29. Eeles, R. A. *et al.* Multiple newly identified loci associated with prostate cancer susceptibility. *Nat. Genet.* **40**, 316–321 (2008).
30. Paternoster, L. *et al.* Multi-ancestry genome-wide association study of 21,000 cases and 95,000 controls identifies new risk loci for atopic dermatitis. *Nat. Genet.* **47**, 1449–56 (2015).
31. Dahl, R. *et al.* Effects of an oral MMP-9 and -12 inhibitor, AZD1236, on biomarkers in moderate/severe COPD: A randomised controlled trial. *Pulm. Pharmacol. Ther.* **25**, 169–177 (2012).
32. Ganz, P. *et al.* Development and validation of a protein-based risk score for cardiovascular outcomes among patients with stable coronary heart disease. *JAMA* **315**, 2532-41 (2016).
33. Traylor, M. *et al.* A novel MMP12 locus is associated with large artery atherosclerotic stroke using a genome-wide age-at-onset informed approach. *PLoS Genet.* **10**, e1004469 (2014).
34. Nelson, M. R. *et al.* The support of human genetic evidence for approved drug indications. *Nat. Genet.* **47**, 856–860 (2015).
35. Albagha, O. M. E. *et al.* Genome-wide association study identifies variants at CSF1, OPTN and TNFRSF11A as genetic risk factors for Paget’s disease of bone. *Nat. Genet.* **42**, 520–524 (2010).

36. Schwarz, P., Rasmussen, A. Q., Kvist, T. M., Andersen, U. B. & Jørgensen, N. R. Paget's disease of the bone after treatment with Denosumab: a case report. *Bone* **50**, 1023–5 (2012).

## FIGURE LEGENDS

### Figure 1. The genetic architecture of plasma protein levels.

n=3,301 participants. (a) Genomic location of pQTLs. Red=*cis*, blue=*trans*. X- and Y-axes indicate the positions of the sentinel variant and the gene encoding the associated protein, respectively. Highly pleiotropic genomic regions are annotated. (b) Significance of *cis* associations (linear regression) versus distance from TSS. (c) Number of significantly associated loci per protein. (d) Number of conditionally significant signals within each associated locus. (e) Histogram of variance explained by conditionally significant variants. (f) Effect-size versus MAF. (g) Distributions of the predicted functional annotation class of sentinel pQTL variants versus null sets of variants from permutation. Bar height represents the mean proportion of variants within each class and error bars reflect one standard deviation from the mean. \*=significant enrichment (permutation test, Bonferroni-corrected threshold,  $p < 0.005$ ).

### Figure 2. Missense variant rs28929474 in *SERPINA1* is a *trans* pQTL hotspot.

Numbers (outermost) indicate chromosomes. Lines link the genomic location of rs28929474 with genes encoding significantly associated proteins. Associations with and without asterixes indicate significance at  $p < 5 \times 10^{-8}$  and  $p < 1.5 \times 10^{-11}$ , respectively. Line thickness is proportional to effect-size (red=positive, blue=negative). n=3,301 participants.

### Figure 3. *Trans* pQTL for BLIMP1 at an inflammatory bowel disease (IBD) associated missense variant (rs3197999:A) in *MST1*.

(a) rs3197999:A is associated with multiple proteins. Lines link rs3197999 and the genes encoding significantly associated proteins. Line thickness is proportional to effect-size. Line thickness is proportional to effect-size of the IBD risk allele (red=positive, blue=negative). n=3,301 participants. \*=genes in IBD GWAS loci. (b) Regional association plots at *MST1*, showing IBD association (top) and *trans* pQTLs for BLIMP1, DOCK9 and FASLG. Colour key indicates  $r^2$  with rs3197999. (c) Regional association plot of the IBD susceptibility locus at *PRDM1*, which encodes BLIMP1. IBD association data are for European participants from Liu *et al.*, 2015.

### Figure 4. Proteinase-3, *SERPINA1*, and vasculitis.

(a) Manhattan plots for plasma PR3 measured with two SOMAmers and the Olink assay. (b) *PRTN3* regional association plots. Colour key indicates  $r^2$  with sentinel variant rs10425544. ‘Vasculitis GWAS’: previously reported vasculitis-associated variants (see [Supplementary Note](#)). EVGC=rs62132295 (from European Vasculitis Genetics Consortium<sup>39</sup>); VCRCi=rs138303849 and VCRCt=rs62132293, most significant imputed and genotyped variants, respectively, from Vasculitis Clinical Research Consortium<sup>69</sup>. ‘Independent pQTLs’: conditionally independent PR3 pQTL variants (black lettering=lead variant for both SOMAmers; purple=conditionally independent variant for SOMamer PRTN3.3514.49.2; green for PRTN3.13720.95.3). (c) Proposed mechanisms by which *PRTN3* and *SERPINA1* impact PR3 levels and thus vasculitis risk. Left: individuals without either the *PRTN3* or *SERPINA1* vasculitis risk alleles. Middle: *SERPINA1* Z-allele carriers have lower circulating A1AT, resulting in higher free plasma PR3. Right: *cis*-acting variant at the *PRTN3* locus results in higher total plasma PR3. Increases in either free or total PR3 predispose to loss of immune tolerance.

### Figure 5. Evaluation of causal role of proteins in disease.

n=3,301 participants. (a) MR estimates with 95% CIs (instrumental variable analysis) for proteins encoded in the *IL1RL1-IL18R1* locus and atopic dermatitis (AD) risk. Univariable MR not possible for IL1R1 and IL18RAP (no significant pQTLs to select as “genetic instruments”). (b) MMP-12 levels and risk of coronary heart disease (CHD). Above: MR estimates with 95% CIs. Below: estimated effect-sizes (with 95% CIs) on plasma MMP-12 (from linear regression) and CHD risk (from logistic regression) for each variant used in the genetic score.



## **ONLINE METHODS**

### **Study participants**

The INTERVAL study comprised about 50,000 participants nested within a randomised trial of varying blood donation intervals<sup>9</sup>. Between mid-2012 and mid-2014, whole-blood donors aged 18 years and older were recruited at 25 centres of England's National Health Service Blood and Transplant (NHSBT). All participants gave informed consent before joining the study and the National Research Ethics Service approved (11/EE/0538) this study. Participants completed an online questionnaire including questions about demographic characteristics (e.g., age, sex, ethnic group), anthropometry (height, weight), lifestyle (e.g., alcohol and tobacco consumption) and diet. Participants were generally in good health because blood donation criteria exclude people with a history of major diseases (such as myocardial infarction, stroke, cancer, HIV, and hepatitis B or C) and those who have had recent illness or infection. For protein assays, we randomly selected two non-overlapping subcohorts of 2,731 and 831 participants from INTERVAL. After genetic QC, 3,301 participants (2,481 and 820 in the two subcohorts) remained for analysis (Supplementary Table 17).

### **Plasma sample preparation**

Sample collection procedures for INTERVAL have been described previously<sup>37</sup>. In brief, blood samples for research purposes were collected in 6ml EDTA tubes using standard venepuncture protocols. The tubes were inverted three times and transferred at room temperature to UK Biocentre (Stockport, UK) for processing. Plasma was extracted into two 0.8ml plasma aliquots by centrifugation and subsequently stored at -80°C prior to use.

### **Protein measurements**

We used a multiplexed, aptamer-based approach (SOMAscan assay) to measure the relative concentrations of 3,622 plasma proteins/protein complexes assayed using 4,034 modified aptamers ("SOMAmer reagents", hereafter referred to as 'SOMAmers'; Supplementary Table 18). The assay extends the lower limit of

detectable protein abundance afforded by conventional approaches (e.g., immunoassays), measuring both extracellular and intracellular proteins (including soluble domains of membrane-associated proteins), with a bias towards proteins likely to be found in the human secretome (Extended Data Figure 7a)<sup>8,38</sup>. The proteins cover a wide range of molecular functions (Extended Data Figure 7b). The selection of proteins on the platform reflects both the availability of purified protein targets and a focus on proteins suspected to be involved in pathophysiology of human disease.

Aliquots of 150 µl of plasma were sent on dry ice to SomaLogic Inc. (Boulder, Colorado, US) for protein measurement. Assay details have been previously described<sup>38-40</sup> and a technical white paper with further information can be found at the manufacturer's website ([http://somallogic.com/wp-content/uploads/2017/06/SSM-002-Technical-White-Paper\\_010916\\_LSM1.pdf](http://somallogic.com/wp-content/uploads/2017/06/SSM-002-Technical-White-Paper_010916_LSM1.pdf)). In brief, modified single-stranded DNA SOMAmers are used to bind to specific protein targets that are then quantified using a DNA microarray. Protein concentrations are quantified as relative fluorescent units.

Quality control (QC) was performed at the sample and SOMAmer level using control aptamers, as well as calibrator samples. At the sample level, hybridisation controls on the microarray were used to correct for systematic variability in hybridisation, while the median signal over all features assigned to one of three dilution sets (40%, 1% and 0.005%) was used to correct for within-run technical variability. The resulting hybridisation scale factors and median scale factors were used to normalise data across samples within a run. The acceptance criteria for these values are between 0.4 and 2.5 based on historical runs. SOMAmer-level QC made use of replicate calibrator samples using the same study matrix (plasma) to correct for between-run variability. The acceptance criterion for each SOMAmer was that the calibration scale factor be less than 0.4 from the median for each of the plates run. In addition, at the plate level, the acceptance criteria were that the median of the calibration scale factors be between 0.8 and 1.2, and that 95% of individual SOMAmers be less than 0.4 from the median within the plate.

In addition to QC processes routinely conducted by SomaLogic, we measured protein levels of 30 and 10 pooled plasma samples randomly distributed across plates for subcohort 1 and subcohort 2, respectively. Laboratory technicians were blinded to the presence of pooled samples. This approach enabled estimation of the reproducibility of the protein assays. We calculated CVs for each SOMAmer within each subcohort by dividing the standard deviation by the mean of the pooled plasma sample protein read-outs. In addition to passing SomaLogic QC processes, we required SOMAmers to have a  $CV \leq 20\%$  in both subcohorts. Eight non-human protein targets were also excluded, leaving 3,283 SOMAmers (mapping to 2,994 unique proteins/protein complexes) for inclusion in the GWAS.

Protein mapping to UniProt identifiers and gene names was provided by SomaLogic. Mapping to Ensembl gene IDs and genomic positions was performed using Ensembl Variant Effect Predictor v83 (VEP)<sup>41</sup>. Protein subcellular locations were determined by exporting the subcellular location annotations from UniProt<sup>42</sup>. If the term ‘membrane’ was included in the descriptor, the protein was considered to be a membrane protein, whereas if the term ‘secreted’ (but not ‘membrane’) was included in the descriptor, the protein was considered to be a secreted protein. Proteins not annotated as either membrane or secreted proteins were classified (by inference) as intracellular proteins. Proteins were mapped to molecular functions using gene ontology annotations<sup>43</sup> from UniProt.

## **Non-genetic associations of proteins**

To provide confidence in the reproducibility of the protein assays, we attempted to replicate the associations with age or sex of 45 proteins previously reported by Ngo *et al* and 40 reported by Menni *et al*<sup>39,44</sup>. We used Bonferroni-corrected  $p$ -value thresholds of  $p=1.1 \times 10^{-3}$  (0.05/45) and  $p=1.2 \times 10^{-3}$  (0.05/40) respectively. Relative protein abundances were rank-inverse normalised within each subcohort and linear regression was performed using age, sex, BMI, natural log of estimated glomerular filtration rate (eGFR) and subcohort as independent variables.

## **Genotyping and imputation**

The genotyping protocol and QC for the INTERVAL samples (n~50,000) have been described previously in detail<sup>10</sup>. Briefly, DNA extracted from buffy coat was used to assay approximately 830,000 variants on the Affymetrix Axiom UK Biobank genotyping array at Affymetrix (Santa Clara, California, US). Genotyping was performed in multiple batches of approximately 4,800 samples each. Sample QC was performed including exclusions for sex mismatches, low call rates, duplicate samples, extreme heterozygosity and non-European descent. An additional exclusion made for this study was of one participant from each pair of close (first- or second-degree) relatives, defined as  $\hat{\pi} > 0.187$ . Identity-by-descent was estimated using a subset of variants with a call rate  $> 99\%$  and MAF  $> 5\%$  in the merged dataset of both subcohorts, pruned for linkage disequilibrium (LD) using PLINK v1.9<sup>45</sup>. Numbers of participants excluded at each stage of the genetic QC are summarised in Extended Data Figure 1. Multi-dimensional scaling was performed using PLINK v1.9 to create components to account for ancestry in genetic analyses.

Prior to imputation, additional variant filtering steps were performed to establish a high-quality imputation scaffold. In summary, 654,966 high quality variants (autosomal, non-monomorphic, bi-allelic variants with Hardy Weinberg Equilibrium (HWE)  $p > 5 \times 10^{-6}$ , with a call rate of  $> 99\%$  across the INTERVAL genotyping batches in which a variant passed QC, and a global call rate of  $> 75\%$  across all INTERVAL genotyping batches) were used for imputation. Variants were phased using SHAPEIT3 and imputed using a combined 1000 Genomes Phase 3-UK10K reference panel. Imputation was performed via the Sanger Imputation Server (<https://imputation.sanger.ac.uk>) resulting in 87,696,888 imputed variants.

Prior to genetic association testing, variants were filtered in each subcohort separately using the following exclusion criteria: (1) imputation quality (INFO) score  $< 0.7$ , (2) minor allele count  $< 8$ , (3) HWE  $p < 5 \times 10^{-6}$ . In the small number of cases where imputed variants had the same genomic position (GRCh37) and alleles, the variant with the lowest INFO score was removed. 10,572,788 variants passing all filters in both subcohorts were taken forward for analysis (Extended Data Figure 1).

## **Genome-wide association study**

Within each subcohort, relative protein abundances were first natural log-transformed. Log-transformed protein levels were then adjusted in a linear regression for age, sex, duration between blood draw and processing (binary,  $\leq 1$  day/ $>1$ day) and the first three principal components of ancestry from multi-dimensional scaling. The protein residuals from this linear regression were then rank-inverse normalised and used as phenotypes for association testing. Simple linear regression using an additive genetic model was used to test genetic associations. Association tests were carried out on allelic dosages to account for imputation uncertainty (“-method expected” option) using SNPTEST v2.5.2<sup>46</sup>.

## Meta-analysis and statistical significance

Association results from the two subcohorts were combined via fixed-effects inverse-variance meta-analysis combining the betas and standard errors using METAL<sup>47</sup>. Genetic associations were considered to be genome-wide significant based on a conservative strategy requiring associations to have (i) a meta-analysis  $p$ -value  $< 1.5 \times 10^{-11}$  (genome-wide threshold of  $p = 5 \times 10^{-8}$  Bonferroni-corrected for 3,283 aptamers tested), (ii) at least nominal significance ( $p < 0.05$ ) in both subcohorts, and (iii) consistent direction of effect across subcohorts. We did not observe significant genomic inflation (mean inflation factor was 1.0, standard deviation=0.01) ([Extended Data Figure 2d](#)).

## Refinement of significant regions

To identify distinct non-overlapping regions associated with a given SOMAmer, we first defined a 1Mb region around each significant variant for that SOMAmer. Starting with the region containing the variant with the smallest  $p$ -value, any overlapping regions were then merged and this process was repeated until no more overlapping 1Mb regions remained. The variant with the lowest  $p$ -value for each region was assigned as the “regional sentinel variant”. Due to the complexity of the Major Histocompatibility Region (MHC) region, we treated the extended MHC region (chr6:25.5-34.0Mb) as one region. To identify whether a region was associated with multiple SOMAmers, we used an LD-based clumping approach. Regional sentinel variants in high LD ( $r^2 \geq 0.8$ ) with each other were combined together into a single region.

## Conditional analyses

To identify conditionally significant signals, we performed approximate genome-wide step-wise conditional analysis using GCTA v1.25.2<sup>48</sup> using the “cojo-slc” option. We used the same conservative significance threshold of  $p=1.5 \times 10^{-11}$  as for the univariable analysis. As inputs for GCTA, we used the summary statistics (i.e. betas and standard errors) from the meta-analysis. Correlation between variants was estimated using the ‘hard-called’ genotypes (where a genotype was called if it had a posterior probability of  $>0.9$  following imputation or set to missing otherwise) in the merged genetic dataset, and only variants also passing the univariable genome-wide threshold ( $p < 1.5 \times 10^{-11}$ ) were considered for step-wise selection. As the conditional analyses use different data inputs to the univariable analysis (i.e. summarised rather than individual-level data), there were some instances where the conditional analysis failed to include in the step-wise selection sentinel variants that were only just statistically significant in the univariable analysis. In these instances ( $n=28$ ), we re-conducted the joint model estimation without step-wise selection in GCTA, using the variants identified by the conditional analysis in addition to the regional sentinel variant. We report and highlight these cases in [Supplementary Table 5](#).

## Replication of previous pQTLs

We attempted to identify all previously reported pQTLs from GWAS and to assess whether they replicated in our study. We used the NCBI Entrez programming utility in R (rentrez) to perform a literature search for pQTL studies published from 2008 onwards. We searched for the following terms: ‘pQTL’, ‘pQTLs’, and ‘protein quantitative trait locus’. We supplemented this search by filtering out GWAS associations from the NHGRI-EBI GWAS Catalog v.1.0.1<sup>49</sup> (<https://www.ebi.ac.uk/gwas/>, downloaded November 2017), which has all phenotypes mapped to the Experimental Factor Ontology (EFO)<sup>50</sup>, by restricting to those with EFO annotations relevant to protein biomarkers (e.g., ‘protein measurement’, EFO\_0004747). Studies identified through both approaches were manually filtered to include only studies that profiled plasma or serum samples and to exclude studies not assessing proteins. We recorded basic summary information for each study including the assay used, sample size and number of proteins with pQTLs ([Supplementary Table 19](#)). To reduce the impact of ethnic differences in allele frequencies on replication rate estimates, we filtered

studies to include only associations reported in European-ancestry populations. We then manually extracted summary data on all reported associations from the manuscript or the supplementary material. This included rsID, protein UniProt ID, *p*-values, and whether the association is *cis/trans* ([Supplementary Table 20](#)).

To assess replication we first identified the set of unique UniProt IDs that were also assayed on the SOMAscan panel. For previous studies that used SomaLogic technology, we refined this match to the specific aptamer used. We then clumped associations into distinct loci using the same method that we applied to our pQTLs (see **Refinement of significant regions**). For each locus, we asked if the sentinel SNP or a proxy ( $r^2 > 0.6$ ) was associated with the same protein/aptamer in our study at a defined significance threshold. For our primary assessment, we used a *p*-value threshold of  $10^{-4}$  ([Supplementary Table 21](#)). We also performed sensitivity analyses to explore factors that influence replication rate ([Supplementary Note](#)).

## Replication study using Olink assay

To test replication of 163 pQTLs for 116 proteins, we performed protein measurements using an alternative assay, i.e., a proximity extension assay method (Olink Bioscience, Uppsala, Sweden)<sup>51</sup> in an additional subcohort of 4,998 INTERVAL participants. Proteins were measured using three 92-protein ‘panels’ – ‘inflammatory’, ‘cvd2’ and ‘cvd3’ (10 proteins were assayed on more than 1 panel). 4,902, 4,947 and 4,987 samples passed quality control for the ‘inflammatory’, ‘cvd2’ and ‘cvd3’ panels, respectively, of which, 712, 715 and 721 samples were from individuals included in our primary pQTL analysis using the SOMAscan assay. Normalised protein levels (‘NPX’) were regressed on age, sex, plate, time from blood draw to processing (in days), and season (categorical – ‘Spring’, ‘Summer’, ‘Autumn’, ‘Winter’). The residuals were then rank-inverse normalized. Genotype data was processed as described earlier. Linear regression of the rank-inversed normalised residuals on genotype was carried out in SNPTTEST with the first three components of multi-dimensional scaling as covariates to adjust for ancestry. pQTLs were considered to have replicated if they met a *p*-value threshold Bonferroni-corrected for the number of tests ( $p < 3.1 \times 10^{-4}$ ;  $0.05/163$ ) and had a directionally concordant beta estimate with the SOMAscan estimate.

## Candidate gene annotation

We defined a pQTL as *cis* when the most significantly associated variant in the region was located within 1Mb of the transcription start site (TSS) of the gene(s) encoding the protein. pQTLs lying outside of the region were defined as *trans*. When considering the distance of the lead *cis*-associated variant from the relevant TSS, only proteins that map to single genes on the primary assembly in Ensembl v83 were considered.

For *trans* pQTLs, we sought to prioritise candidate genes in the region that might underpin the genotype-protein association. We applied the ProGeM framework<sup>22</sup> that leverages a combination of databases of molecular pathways, protein-protein interaction networks, and variant annotation, as well as functional genomic data including eQTL and chromosome conformation capture. In addition to reporting the nearest gene to the sentinel variant, ProGeM employs complementary ‘bottom up’ and ‘top down’ approaches, starting from the variant and protein respectively. For the ‘bottom up’ approach, the sentinel variant and corresponding proxies ( $r^2 > 0.8$ ) for each *trans* pQTL were first annotated using Ensembl VEP v83 (using the ‘pick’ option) to determine whether variants were (1) protein-altering coding variants; (2) synonymous coding or 5’/3’ untranslated region (UTR); (3) intronic or up/downstream; or (4) intergenic. Second, we queried all sentinel variants and proxies against significant *cis* eQTL variants (defined by beta distribution-adjusted empirical *p*-values using an FDR threshold of 0.05, see <http://www.gtexportal.org/home/documentationPage> for details) in any cell type or tissue from the Genotype-Tissue Expression (GTEx) project v6<sup>27</sup> (<http://www.gtexportal.org/home/datasets>). Third, we also queried promoter capture Hi-C data in 17 human primary hematopoietic cell types<sup>52</sup> to identify contacts (with a CHICAGO score  $> 5$  in at least one cell type) involving chromosomal regions containing a sentinel variant. We considered gene promoters annotated on either fragment (i.e., the fragment containing the sentinel variant or the other corresponding fragment) as potential candidate genes. Using these three sources of information, we generated a list of candidate genes for the *trans* pQTLs. A gene was considered a candidate if it fulfilled at least one of the following criteria: (1) it was proximal (intragenic or  $\pm 5$ Kb from the gene) or nearest to the sentinel variant; (2) it contained a sentinel or proxy variant ( $r^2 > 0.8$ ) that was protein-



altering; (3) it had a significant *cis* eQTL in at least one GTEx tissue overlapping with a sentinel pQTL variant (or proxy); or (4) it was regulated by a promoter annotated on either fragment of a chromosomal contact<sup>52</sup> involving a sentinel variant.

For the ‘top down’ approach, we first identified all genes with a TSS located within the corresponding pQTL region using the GenomicRanges Bioconductor package<sup>53</sup> with annotation from a GRCh37 GTF file from Ensembl ([ftp://ftp.ensembl.org/pub/grch37/update/gtf/homo\\_sapiens/](ftp://ftp.ensembl.org/pub/grch37/update/gtf/homo_sapiens/); file: ‘Homo\_sapiens.GRCh37.82.gtf.gz’, downloaded June 2016). We then identified any local genes that had previously been linked with the corresponding *trans*-associated protein(s) according to the following open source databases: (1) the Online Mendelian Inheritance in Man (OMIM) catalogue<sup>54</sup> (<http://www.omim.org/>); (2) the Kyoto Encyclopedia of Genes and Genomes (KEGG)<sup>55</sup> (<http://www.genome.jp/kegg/>); and (3) STRINGdb<sup>56</sup> (<http://string-db.org/>; v10.0). We accessed OMIM data via HumanMine web tool<sup>57</sup> (<http://www.humanmine.org/>; accessed June 2016), whereby we extracted all OMIM IDs for (i) our *trans*-affected proteins and (ii) genes local ( $\pm 500\text{Kb}$ ) to the corresponding *trans*-acting variant. We extracted all human KEGG pathway IDs using the KEGGREST Bioconductor package (<https://bioconductor.org/packages/release/bioc/html/KEGGREST.html>). In cases where a *trans*-associated protein shared either an OMIM ID or a KEGG pathway ID with a gene local to the corresponding *trans*-acting variant, we took this as evidence of a potential functional involvement of that gene. We interrogated protein-protein interaction data by accessing STRINGdb data using the STRINGdb Bioconductor package<sup>58</sup>, whereby we extracted all pairwise interaction scores for each *trans*-affected protein and all proteins with genes local to the corresponding *trans*-acting variants. We took the default interaction score of 400 as evidence of an interaction between the proteins, therefore indicating a possible functional involvement for the local gene. In addition to using data from open source databases in our top down approach we also adopted a ‘guilt-by-association’ (GbA) approach utilising the same plasma proteomic data used to identify our pQTLs. We first generated a matrix containing all possible pairwise Pearson’s correlation coefficients between our 3,283 SOMAmers. We then extracted the coefficients relating to our *trans*-associated proteins and any proteins encoded by genes local to their corresponding *trans*-acting variants (where available).

Where the correlation coefficient was  $\geq 0.5$  we prioritised the relevant local genes as being potential mediators of the *trans* signal(s) at that locus.

We report the potential candidate genes for our *trans* pQTLs from both the ‘bottom up’ and ‘top down’ approaches, highlighting cases where the same gene was highlighted by both approaches.

## **Functional annotation of pQTLs**

Functional annotation of variants was performed using Ensembl VEP v83 using the ‘pick’ option. We tested the enrichment of significant pQTL variants for certain functional classes by comparing to permuted sets of variants showing no significant association with any protein ( $p > 0.0001$  for all proteins tested). First, the regional sentinel variants were LD-pruned at  $r^2$  of 0.1. Each time the sentinel variants were LD-pruned, one of the pairs of correlated variants was removed at random and for each set of LD-pruned sentinel variants, 100 sets of equally sized null permuted variants were sampled matching for MAF (bins of 5%), distance to TSS (bins of 0-0.5Kb, 0.5-2Kb, 2-5Kb, 5-10Kb, 10-20Kb, 20-100Kb and >100Kb in each direction) and LD ( $\pm$  half the number of variants in LD with the sentinel variant at  $r^2$  of 0.8). This procedure was repeated 100 times resulting in 10,000 permuted sets of variants. An empirical  $p$ -value was calculated as the proportion of permuted variant sets where the proportion that is classified as a particular functional group exceeded that of the test set of sentinel pQTL variants, and we used a significance threshold of  $p = 0.005$  (0.05/10 functional classes tested).

## **Evidence against aptamer-binding effects at *cis* pQTLs**

All protein assays that rely on binding (e.g., of antibodies or SOMAmers) are susceptible to the possibility of binding-affinity effects, where protein-altering variants (PAVs) (or their proxies in LD) are associated with protein measurements due to differential binding rather than differences in protein abundance. To account for this potential effect, we performed conditional analysis at all *cis* pQTLs where the sentinel variant was in LD ( $r^2 \geq 0.1$  and  $r^2 \leq 0.9$ ) with a PAV in the gene(s) encoding the associated protein. First, variants were annotated with Ensembl VEP v83 using the “per-gene” option. Variant annotations were

considered protein-altering if they were annotated as coding sequence variant, frameshift variant, in-frame deletion, in-frame insertion, missense variant, protein altering variant, splice acceptor variant, splice donor variant, splice region variant, start lost, stop gained, or stop lost. To avoid multi-collinearity, PAVs were LD-pruned ( $r^2 > 0.9$ ) using PLINK v1.9 before including them as covariates in the conditional analysis on the meta-analysis summary statistics using GCTA v1.25.2. Coverage of known common (MAF > 5%) PAVs in our data was checked by comparison with exome sequences from ~60,000 individuals in the Exome Aggregation Consortium (ExAC [<http://exac.broadinstitute.org>], downloaded June 2016)<sup>59</sup>.

## Testing for regulatory and functional enrichment

We tested whether our pQTLs were enriched for functional and regulatory characteristics using GARFIELD v1.2.0<sup>60</sup>. GARFIELD is a non-parametric permutation-based enrichment method that compares input variants to permuted sets matched for number of proxies ( $r^2 \geq 0.8$ ), MAF and distance to the closest TSS. It first applies “greedy pruning” ( $r^2 < 0.1$ ) within a 1Mb region of the most significant variant. GARFIELD annotates variants with more than a thousand features, drawn predominantly from the GENCODE, ENCODE and ROADMAP projects, which includes genic annotations, histone modifications, chromatin states and other regulatory features across a wide range of tissues and cell types.

The enrichment analysis was run using all variants that passed our Bonferroni-adjusted significance threshold ( $p < 1.5 \times 10^{-11}$ ) for association with any protein. For each of the matching criteria (MAF, distance to TSS, number of LD proxies), we used five bins. In total we tested 25 combinations of features (classified as transcription factor binding sites, FAIRE-seq, chromatin states, histone modifications, footprints, hotspots, or peaks) with up to 190 cell types from 57 tissues, leading to 998 tests. Hence, we considered enrichment with a  $p < 5 \times 10^{-5}$  (0.05/998) to be statistically significant.

## Disease annotation

To identify diseases that our pQTLs have been associated with, we queried our sentinel variants and their strong proxies ( $r^2 \geq 0.8$ ) against publicly available disease GWAS data using PhenoScanner<sup>61</sup>. A list of

datasets queried is available at <http://www.phenoscaner.medschl.cam.ac.uk/information.html>. For disease GWAS, results were filtered to  $p < 5 \times 10^{-8}$  and then manually curated to retain only the entry with the strongest evidence for association (i.e. smallest  $p$ -value) per disease. Non-disease phenotypes such as anthropometric traits, intermediate biomarkers and lipids were excluded manually.

### ***Cis* eQTL overlap and enrichment of *cis* pQTLs for *cis* eQTLs**

For each regional sentinel *cis* pQTL variant, its strong proxies ( $r^2 \geq 0.8$ ) were queried against publicly available eQTL association data using PhenoScanner. *Cis* eQTL results were filtered to retain only variants with  $p < 1.5 \times 10^{-11}$ . Only *cis* eQTLs for the same gene as the *cis* pQTL protein were retained. We tested whether *cis* pQTLs were significantly enriched for eQTLs for the corresponding gene compared to null sets of variants appropriately matched for MAF and distance to nearest TSS. For this analysis, we restricted eQTL data to the GTEx project v6, since this project provided complete summary statistics across a wide range of tissues and cell-types, in contrast to many other studies which only report  $p$ -values below some significance level. GTEx results were filtered to contain only variants lying in *cis* (i.e., within 1Mb) of genes that encode proteins analysed in our study and only variants in both datasets were utilised.

For the enrichment analysis, the *cis* pQTL sentinel variants were first LD-pruned ( $r^2 < 0.1$ ) and the proportion of sentinel *cis* pQTL variants that are also eQTLs (at our pQTL significance threshold [ $p < 1.5 \times 10^{-11}$ ], conventional genomewide significance [ $p < 5 \times 10^{-8}$ ] or a nominal  $p$ -value threshold [ $p < 1 \times 10^{-5}$ ]) for the same protein/gene was compared to a permuted set of variants that were not pQTLs ( $p > 0.0001$  for all proteins). We generated 10,000 permuted sets of null variants for each significance threshold matched for MAF, distance to TSS and LD (as described for functional annotation enrichment in **Functional annotation of pQTLs**). An empirical  $p$ -value was calculated as the proportion of permuted variant sets where the proportion that are also *cis* eQTLs exceeded that of the test set of sentinel *cis* pQTL variants.

At a stringent eQTL significance threshold ( $p < 1.5 \times 10^{-11}$ ), we found significant enrichment of *cis* pQTLs for eQTLs ( $p < 0.0001$ ) ([Supplementary Table 11](#)) with 19.5% overlap observed compared to a mean overlap of

1.8% in the null sets. Results were similar in sensitivity analyses using the standard genome-wide or nominal significance thresholds as well as when using only the sentinel variants at *cis* pQTLs that were robust to adjusting for PAVs ([Supplementary Table 7](#)), suggesting our results are robust to the choice of threshold and potential differential binding effects.

## Colocalisation analysis

Colocalisation testing was performed using the coloc package<sup>62</sup>. For testing colocalisation of pQTLs and disease association signals, colocalisation testing was necessarily limited to disease traits where full GWAS summary statistics had been made available. We obtained GWAS summary statistics obtained through PhenoScanner. For testing colocalisation of pQTLs with eQTLs, we used publically available summary statistics for expression traits from GTEx<sup>27</sup>. We used the default priors. Regions for testing were determined by dividing the genome into 0.1cM chunks using recombination data. Evidence for colocalisation was assessed using the posterior probability (PP) for hypothesis 4 (that there is an association signal for both traits and they are driven by the same causal variant[s]). Signals with  $PP_4 > 0.5$  were deemed likely to colocalise as this gives hypothesis 4 the highest likelihood of being correct, while  $PP_4 > 0.8$  was deemed to be ‘highly likely to colocalise’.

## Selection of genetic instruments for Mendelian randomisation

In Mendelian randomisation (MR), genetic variants are used as ‘instrumental variables’ (IV) for assessing the causal effect of the exposure (here a plasma protein) on the outcome (here disease)<sup>11,63</sup> ([Extended Data Figure 6](#)).

### Proteins in the *IL1RL1-IL18R1* locus and atopic dermatitis

To identify the likely causal proteins that underpin the previous genetic association of the *IL1RL1-IL18R1* locus (chr11:102.5-103.5Mb) with atopic dermatitis (AD)<sup>30</sup>, we used the following approach. For each protein encoded by a gene in the *IL1RL1-IL18R1* locus, we took genetic variants that had a *cis* association at  $p < 1 \times 10^{-4}$  and ‘LD-pruned’ them at  $r^2 < 0.1$  to leave largely independent variants. We then used these genetic

variants to construct a genetic score for each protein. Formally, we used these variants as instrumental variables for their respective proteins in univariable MR. For multivariable MR, association estimates for all proteins in the locus were extracted for all instruments. We used PhenoScanner to obtain association statistics for the selected variants in the European-ancestry population of a recent large-scale GWAS meta-analysis<sup>30</sup>. Where the relevant variant was not available, the strongest proxy with  $r^2 \geq 0.8$  was used.

### **MMP-12 and coronary heart disease (CHD)**

To test whether plasma MMP-12 levels have a causal effect on risk of CHD, we selected genetic variants in the *MMP12* gene region to use as instrumental variables. We constructed a genetic score comprising 17 variants that had a *cis* association with MMP-12 levels at  $p < 5 \times 10^{-8}$  and that were not highly correlated with one another ( $r^2 < 0.2$ ). To perform multivariable MR, we used association estimates for these variants with other MMP proteins in the locus (MMP-1, MMP-7, MMP-8, MMP-10, MMP-13). Summary associations for variants in the score with CHD were obtained through PhenoScanner from a recent large-scale GWAS meta-analysis which consists mostly (77%) individuals of European ancestry<sup>64</sup>.

### **MR analysis**

Two-sample univariable MR was performed for each protein separately using summary statistics in the inverse-variance weighted method adapted to account for correlated variants<sup>65-66</sup>. For each of  $G$  genetic variants ( $g = 1, \dots, G$ ) having per-allele estimate of the association with the protein  $\beta_{Xg}$  and standard error  $\sigma_{Xg}$ , and per-allele estimate of the association with the outcome (here, AD or CHD)  $\beta_{Yg}$  and standard error  $\sigma_{Yg}$ , the IV estimate ( $\hat{\theta}_{XY}$ ) is obtained from generalised weighted linear regression of the genetic associations with the outcome ( $\beta_Y$ ) on the genetic associations with the protein ( $\beta_X$ ) weighting for the precisions of the genetic associations with the outcome and accounting for correlations between the variants according to the regression model:

$$\beta_Y = \theta_{XY} \beta_X + \varepsilon, \quad \varepsilon \sim N(0, \Omega)$$

where  $\beta_Y$  and  $\beta_X$  are vectors of the univariable (marginal) genetic associations, and the weighting matrix  $\Omega$  has terms  $\Omega_{g_1g_2} = \sigma_{Yg_1}\sigma_{Yg_2}\rho_{g_1g_2}$ , and  $\rho_{g_1g_2}$  is the correlation between the  $g_1$ th and  $g_2$ th variants.

The IV estimate from this method is:

$$\hat{\theta}_{XY} = (\beta_X^T \Omega^{-1} \beta_X)^{-1} \beta_X^T \Omega^{-1} \beta_Y$$

and the standard error is:

$$\text{se}(\hat{\theta}_{XY}) = \sqrt{(\beta_X^T \Omega^{-1} \beta_X)^{-1}}$$

where  $^T$  is a matrix transpose. This is the estimate and standard error from the regression model fixing the residual standard error to 1 (equivalent to a fixed-effects model in a meta-analysis).

Genetic variants in univariable MR need to satisfy three key assumptions to be valid instruments:

- (1) the variant is associated with the risk factor of interest (i.e., the protein level),
- (2) the variant is not associated with any confounder of the risk factor-outcome association,
- (3) the variant is conditionally independent of the outcome given the risk factor and confounders.

To account for potential effects of functional pleiotropy<sup>67</sup>, we performed multivariable MR using the weighted regression-based method proposed by Burgess *et al*<sup>68</sup>. For each of  $K$  risk factors in the model ( $k = 1, \dots, K$ ), the weighted regression-based method is performed by multivariable generalized weighted linear regression of the association estimates  $\beta_Y$  on each of the association estimates with each risk factor  $\beta_{Xk}$  in a single regression model:

$$\beta_Y = \theta_{XY1} \beta_{X1} + \theta_{XY2} \beta_{X2} + \dots + \theta_{XYK} \beta_{XK} + \varepsilon, \quad \varepsilon \sim N(0, \Omega)$$

where  $\beta_{X_1}$  is the vectors of the univariable genetic associations with risk factor 1, and so on. This regression model is implemented by first pre-multiplying the association vectors by the Cholesky decomposition of the weighting matrix, and then applying standard linear regression to the transformed vectors. Estimates and standard errors are obtained fixing the residual standard error to be 1 as above.

The multivariable MR analysis allows the estimation of the causal effect of a protein on disease outcome accounting for the fact that genetic variants may be associated with multiple proteins in the region. Causal estimates from multivariable MR represent direct causal effects, representing the effect of intervening on one risk factor in the model while keeping others constant.

### **MMP-12 genetic score sensitivity analyses**

We performed two sensitivity analyses to determine the robustness of the MR findings. First, we measured plasma MMP-12 levels using a different method (proximity extension assay; Olink Bioscience, Uppsala, Sweden<sup>51</sup>) in 4,998 individuals, and used this to derive genotype-MMP12 effect estimates for the 17 variants in our genetic score. Second, we obtained effect estimates from a pQTL study based on SOMAscan assay measurements in an independent sample of ~1,000 individuals<sup>3</sup>. In both cases the genetic score reflecting higher plasma MMP-12 was associated with lower risk of CHD.

### **Overlap of pQTLs with drug targets**

We used the Informa Pharmaprojects database from Citeline to obtain information on drugs that target proteins assayed on the SOMAscan platform. This is a manually curated database that maintains profiles for >60,000 drugs. For our analysis, we focused on the following information for each drug: protein target, indications, and development status. We included drugs across the development pipeline, including those in pre-clinical studies or with no development reported, drugs in clinical trials (all phases), and launched/registered drugs. For each protein assayed, we identified all drugs in the Informa Pharmaprojects



with a matching protein target based on UniProt ID. When multiple drugs targeted the same protein, we selected the drug with the latest stage of development.

For drug targets with significant pQTLs, we identified the subset where the sentinel variant or proxy variants in LD ( $r^2 > 0.8$ ) are also associated with disease risk through PhenoScanner. We used an internal Merck auto-encoding method to map GWAS traits and drug indications to a common set of terms from the Medical Dictionary for Regulatory Activities (MedDRA). MedDRA terms are organised into a hierarchy with five levels. We mapped each GWAS trait and indication onto the 'Lowest Level Terms' (i.e. the most specific terms available). All matching terms were recorded for each trait or indication. We matched GWAS traits to drug indications based on the highest level of the hierarchy, called 'System Organ Class' (SOC). We designated a protein as 'matching' if at least one GWAS trait term matched with at least one indication term for at least one drug.

## **Data availability**

Participant-level genotype and protein data, and full summary association results from the genetic analysis, are available through the European Genotype Archive (accession number EGAS00001002555). Summary association are also available via FTP and through PhenoScanner

(<http://www.phenoscanner.medschl.cam.ac.uk>).

## Online References

37. Moore, C. *et al.* The INTERVAL trial to determine whether intervals between blood donations can be safely and acceptably decreased to optimise blood supply: study protocol for a randomised controlled trial. *Trials* **15**, 363 (2014).
38. Gold, L. *et al.* Aptamer-based multiplexed proteomic technology for biomarker discovery. *PLoS One* **5**, e15004 (2010).
39. Menni, C. *et al.* Circulating proteomic signatures of chronological age. *J Gerontol A Biol Sci Med Sci* **70**, 809-16 (2014).
40. Sattlecker, M. *et al.* Alzheimer's disease biomarker discovery using SOMAscan multiplexed protein technology. *Alzheimer's Dement.* **10**, 724–734 (2014).
41. McLaren, W. *et al.* Deriving the consequences of genomic variants with the Ensembl API and SNP Effect Predictor. *Bioinformatics* **26**, 2069–70 (2010).
42. UniProt Consortium. UniProt: a hub for protein information. *Nucleic Acids Res.* **43**, D204–D212 (2015).
43. Ashburner, M. *et al.* Gene Ontology: tool for the unification of biology. *Nat. Genet.* **25**, 25–29 (2000).
44. Ngo, D. *et al.* Aptamer-based proteomic profiling reveals novel candidate biomarkers and pathways in cardiovascular disease. *Circulation* **134**, 270-285 (2016).
45. Chang, C. C. *et al.* Second-generation PLINK: rising to the challenge of larger and richer datasets. *Gigascience* **4**, 7 (2015).
46. Marchini, J., Howie, B., Myers, S., McVean, G. & Donnelly, P. A new multipoint method for genome-wide association studies by imputation of genotypes. *Nat. Genet.* **39**, 906–913 (2007).
47. Willer, C. J., Li, Y. & Abecasis, G. R. METAL: fast and efficient meta-analysis of genomewide association scans. *Bioinformatics* **26**, 2190–1 (2010).
48. Yang, J. *et al.* Conditional and joint multiple-SNP analysis of GWAS summary statistics identifies additional variants influencing complex traits. *Nat. Genet.* **44**, 369–75, S1-3 (2012).
49. Welter, D. *et al.* The NHGRI GWAS Catalog, a curated resource of SNP-trait associations. *Nucleic Acids Res.* **42**, D1001-6 (2014).
50. Malone, J. *et al.* Modeling sample variables with an Experimental Factor Ontology. *Bioinformatics* **26**, 1112–1118 (2010).
51. Enroth, S., Johansson, Å., Enroth, S. B. & Gyllensten, U. Strong effects of genetic and lifestyle factors on biomarker variation and use of personalized cutoffs. *Nat Commun* **5**, 4684 (2014).
52. Javierre, B. M. *et al.* Lineage-specific genome architecture links enhancers and non-coding disease variants to target gene promoters. *Cell* **167**, 1369–1384. e19 (2016).
53. Lawrence, M. *et al.* Software for computing and annotating genomic ranges. *PLoS Comput. Biol.* **9**, e1003118 (2013).

54. Amberger, J. S., Bocchini, C. A., Schiettecatte, F., Scott, A. F. & Hamosh, A. OMIM.org: Online Mendelian Inheritance in Man (OMIM(R)), an online catalog of human genes and genetic disorders. *Nucleic Acids Res.* **43**, D789–D798 (2015).
55. Kanehisa, M., Sato, Y., Kawashima, M., Furumichi, M. & Tanabe, M. KEGG as a reference resource for gene and protein annotation. *Nucleic Acids Res.* **44**, D457–D462 (2016).
56. Szklarczyk, D. *et al.* STRING v10: protein-protein interaction networks, integrated over the tree of life. *Nucleic Acids Res.* **43**, D447–D452 (2015).
57. Smith, R. N. *et al.* InterMine: a flexible data warehouse system for the integration and analysis of heterogeneous biological data. *Bioinformatics* **28**, 3163–3165 (2012).
58. Franceschini, A. *et al.* STRING v9.1: protein-protein interaction networks, with increased coverage and integration. *Nucleic Acids Res.* **41**, D808–D815 (2013).
59. Lek, M. *et al.* Analysis of protein-coding genetic variation in 60,706 humans. *Nature* **536**, 285–291 (2016).
60. Iotchkova, V. *et al.* GARFIELD - GWAS Analysis of Regulatory or Functional Information Enrichment with LD correction. *bioRxiv* (2016). doi:10.1101/085738
61. Staley, J. R. *et al.* PhenoScanner: a database of human genotype-phenotype associations. *Bioinformatics* **32**, 3207–3209 (2016).
62. Giambartolomei, C. *et al.* Bayesian test for colocalisation between pairs of genetic association studies using summary statistics. *PLoS Genet.* **10**, e1004383 (2014).
63. Hingorani, A. & Humphries, S. Nature's randomised trials. *Lancet* **366**, 1906–8 (2005).
64. Nikpay, M. *et al.* A comprehensive 1,000 Genomes-based genome-wide association meta-analysis of coronary artery disease. *Nat. Genet.* **47**, 1121–30 (2015).
65. Burgess, S., Butterworth, A. & Thompson, S. G. Mendelian randomization analysis with multiple genetic variants using summarized data. *Genet. Epidemiol.* **37**, 658–665 (2013).
66. Burgess, S., Dudbridge, F. & Thompson, S. G. Combining information on multiple instrumental variables in Mendelian randomization: comparison of allele score and summarized data methods. *Stat. Med.* **35**, 1880–906 (2016).
67. Burgess, S. & Thompson, S. G. Multivariable Mendelian randomization: the use of pleiotropic genetic variants to estimate causal effects. *Am. J. Epidemiol.* **181**, 251–60 (2015).
68. Burgess, S., Dudbridge, F. & Thompson, S. G. Re: 'Multivariable Mendelian randomization: the use of pleiotropic genetic variants to estimate causal effects'. *Am. J. Epidemiol.* **181**, 290–1 (2015).
69. Merkel, P. A. *et al.* Identification of functional and expression polymorphisms associated with risk for anti-neutrophil cytoplasmic autoantibody-associated vasculitis. *Arthritis Rheumatol.* **69**, 1054–1066 (2016).

## Extended Data Figure legends

**Extended Data Figure 1. Flowchart of sample processing and quality control stages for proteomic and genetic measurements prior to genetic analyses.**

**Extended Data Figure 2. Evidence for the reliability of protein measurements made using the SOMAscan assay.**

- (a) Distribution of coefficients of variation of all proteins on the SOMAscan assay in each subcohort.
- (b) Spearman's correlations for all proteins passing QC derived from contemporaneous assay of baseline and two-year samples from 60 participants.
- (c) Scatterplot of pQTL effect size estimates from SOMAscan versus Olink showing all 163 pQTLs tested (top) and the 106 that formally replicated (bottom).  $r$  = Pearson's correlation coefficient.
- (d) Distribution of inflation factors across proteins that underwent genome-wide association testing, stratified by subcohort and allele frequency (MAF $\geq$ 5%, MAF $<$ 5%)

**Extended Data Figure 3. Genetic architecture of the pQTLs.**

pQTL mapping in  $n=3,301$  individuals.

- (a) Distribution of the predicted consequences of the sentinel pQTL variants compared to matched permuted null sets of variants, stratified by *cis* and *trans*. Asterisks indicate empirical enrichment using a permutation test (10,000 permuted sets of non-associated variants) at a Bonferroni-corrected significance value ( $p < 0.005$ ). Bar height represents the mean proportion of variants within each class and error bars reflect one standard deviation from the mean.
- (b) Number of proteins associated ( $p < 1.5 \times 10^{-11}$ ) with each sentinel variant across the genome.

**Extended Data Figure 4. Enrichment of pQTLs at DNase I hypersensitive sites by tissue/cell-type.**

Circle shows enrichment for DNase I hypersensitive sites ("hotspots") for each of 55 tissues (183 cell-types) available from the ENCODE and Roadmap Epigenomics projects, with tissues/cell-types clustered and coloured by anatomical grouping. Some tissues have multiple values due to availability of multiple cell-type or multiple tests per cell-type. Radial lines show fold-enrichment, while dots around the inside edge of the circle denote statistically significant enrichment at a Bonferroni-corrected significant threshold  $p < 5 \times 10^{-5}$ . Enrichment testing performed using GARFIELD (which tests enrichment against permuted sets of variants matched for MAF, distance to TSS and LD). pQTL data from  $n=3,301$  individuals.

**Extended Data Figure 5. Scheme outlining the combined "bottom-up" and "top-down" process utilised for candidate gene annotation of *trans* pQTL regions** (see [Methods](#)).

GbA; guilt-by-association, KEGG; Kyoto Encyclopedia of Genes and Genomics, OMIM; Online Mendelian Inheritance in Man, STRINGdb; STRING database.

**Extended Data Figure 6. Comparison between a randomised controlled trial and Mendelian randomisation to assess the causal effect of changes in protein biomarker levels on disease risk.**

**Extended Data Figure 7. Characterisation of protein targets measured using the SOMAscan assay.**

- (a) Compartment distribution with annotations of all proteins in the Human Protein Atlas for comparison.
- (b) GO molecular functions.

**Extended Data Figure 8. Examples of protein targets for which the SOMAmer is highly specific.**

SDS-PAGE with Alexa-647-labeled proteins captured by the (a) IL1RL2 SOMAmer or (b) GP1BA SOMAmer. For each protein target, the protein captured by the SOMAmer is compared to the standard. The cognate targets are the only ones with protein visible in the capture lanes. These experiments were performed once.

**Extended Data Figure 9. Follow-up of PR3 SOMAmers.**

These experiments were repeated three times independently with similar results.

- (a) SOMAmer pulldowns with purified PR3, A1AT, and PR3:A1AT complex. SOMAmer PRTN3.3514.49.2 enriched PR3:A1AT complex to a much greater degree than free PR3. Conversely, SOMAmer PRTN3.13720.95.3 enriched free PR3 to a greater degree than the PR3:A1AT complex.
- (b) Solution Affinity of PRTN3.3514.49.2 and PRTN3.13720.95.3 for PR3, A1AT, and PR3:A1AT complex. SOMAmer PRTN3.3514.49.2 has a higher affinity for PR3:A1AT complex than for free PR3. SOMAmer PRTN3.13720.95.3, on the other hand, has a higher affinity for free PR3 than SOMAmer PRTN3.3514.49.2.
- (c) Competitive binding of SOMAmers PRTN3.13720.95.3 and PRTN3.3514.49.2 to PR3. Limiting amount of radiolabeled PRTN3.13720.95.3 was incubated with 1 nM Proteinase-3 and a titration of either cold PRTN3.13720.95.3 or cold PRTN3.3514.49.2.

**Extended Data Figure 10. The *WFIKKN2* region is a *trans* pQTL for GDF11/8 plasma levels.**

- (a) Regional association plots of the *trans* pQTL (sentinel variant rs11079936) for GDF11/8 before and after adjusting for levels of *WFIKKN2* (upper panels), and the *WFIKKN2 cis* pQTL after adjusting for GDF11/8 levels (bottom panel). A similar pattern of association for *WFIKKN2* was seen prior to GDF11/8 adjustment (not shown).
- (b) Attenuation of the GDF11/8 *trans* pQTL upon adjustment for plasma levels of the *cis* protein *WFIKKN2*.

**Supplementary Video 1. Three-dimensional interactive plot of sentinel variant-protein associations (red-*cis*, blue-*trans*).** X-axis (“pQTL position”) represents position of the sentinel variant along chromosomes 1-22. Y-axis (“Protein position”) represents the start position of the gene encoding the protein. Z-axis represents the  $-\log_{10}(p)$  of the association. Additional details can be viewed when hovering over the points. Clicking on *cis/trans* in the legend toggles display of points by *cis/trans*. Additional viewing controls are available at the top right of the window. For clarity, associations with  $p < 10^{-300}$  (diamonds) are plotted at  $-\log_{10}(p) = 300$ .

The plot is generated using “plotly” R package v4.5.6 (Plotly Technologies Inc., Montréal, Canada).

## Supplementary Information

Supplementary Information is available in the online version of the paper.

## Acknowledgements

Aaron Day-Williams, Joshua McElwee, Dorothee Diogo, William Astle, Emanuele Di Angelantonio, Ewan Birney, Arianne Richard, Justin Mason and Michael Inouye commented helpfully on the manuscript; Mark Sharp helped with mapping drug indications to GWAS traits. We thank: INTERVAL study participants; staff at recruiting NHSBT blood donation centres; the INTERVAL Study Co-ordination team, Operations Team (led by Richard Houghton and Carmel Moore) and Data Management Team (led by Matthew Walker).

Funding is listed in the Supplementary information.

## Author Contributions

Conceptualization and experimental design: J.D., A.S.B., B.B.S., H.R., R.M.P.;

Methodology: B.B.S., A.S.B., J.C.M., J.E.P., H.R., S.B.;

Analysis: B.B.S., J.C.M., J.E.P., D.S., J.B., J.R.S., T.J., E.P., P.S., C.O-W., M.A.K., S.K.W., A.C., N.B., S.L.S.;

Contributed reagents, materials, protocols or analysis tools: N.J., S.K.W., E.S.Z., J.B., M.A.K., J.R.S., B.P.P.;

Supervision: A.S.B., H.R., J.D., R.M.P., C.S.F., D.S.P., A.M.W.; Writing: A.S.B., J.E.P., B.B.S., J.C.M., H.R., J.D.;

Creation of the INTERVAL BioResource: J.R.B., D.J.R., W.H.O., N.W.M., J.D.;

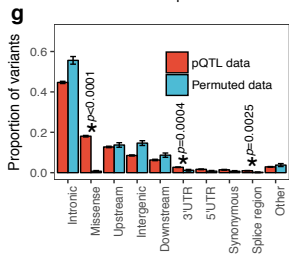
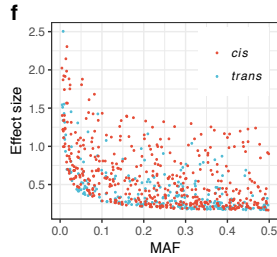
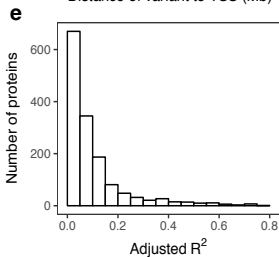
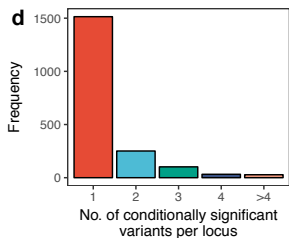
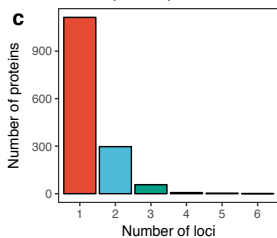
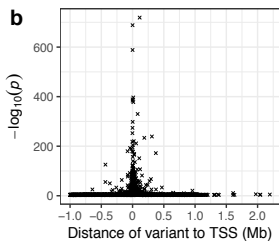
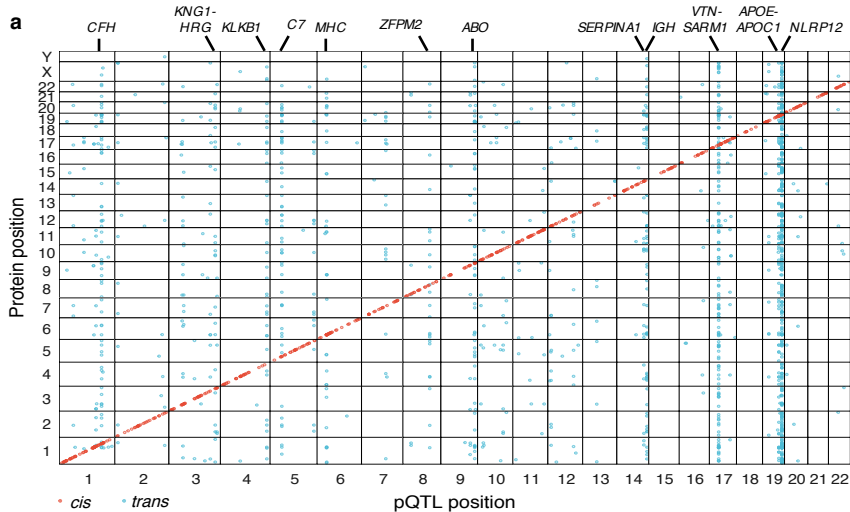
Funding: N.W.M., J.R.B., D.J.R., W.H.O., H.R., R.M.P., J.D.; all authors critically reviewed the manuscript.

## Author Information

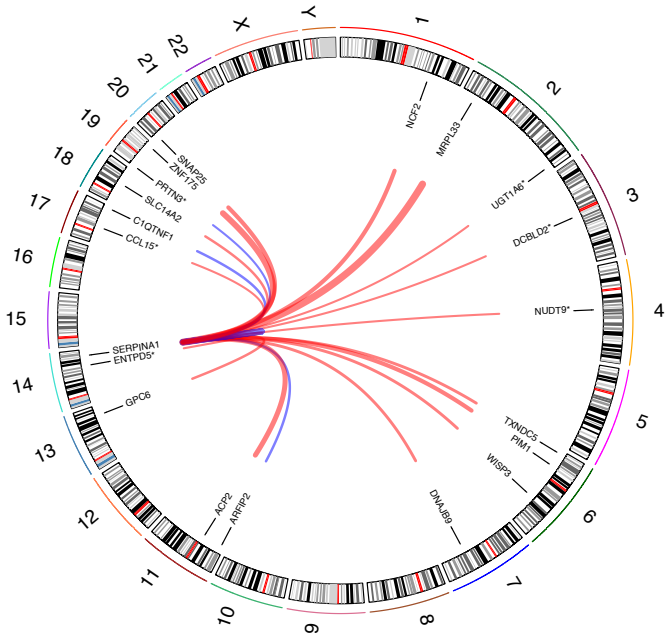
Reprints and permissions information is available at [www.nature.com/reprints](http://www.nature.com/reprints).

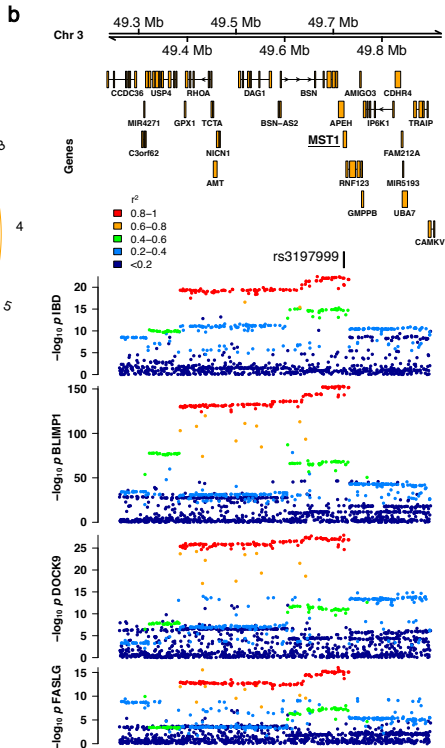
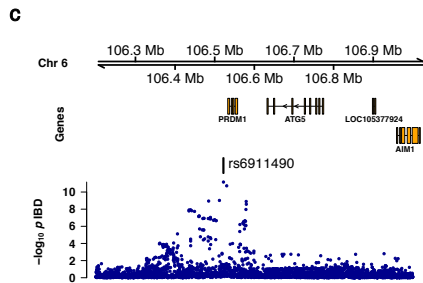
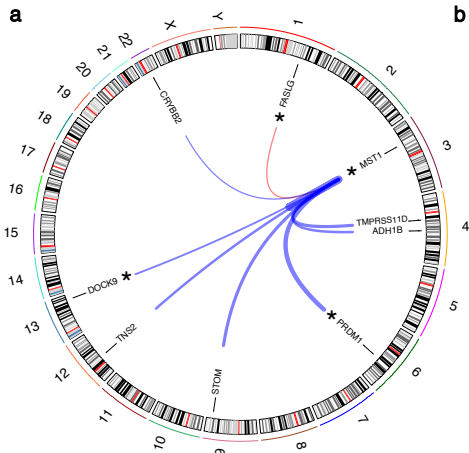
The authors declare the following competing interests: AC,CSF-Merck employees; NJ,SKW-SomaLogic Inc employees and stakeholders; ESZ-SomaLogic Inc employee. JCM,RMP-Merck employees during this study, now Celgene employees. HR-Merck employee during this study; JEP-travel and accommodation expenses and hospitality from Olink to speak at Olink-sponsored academic meetings; ASB-grants from Merck, Pfizer, Novartis, Biogen and Bioverativ and personal fees from Novartis; JD-sits on the Novartis Cardiovascular and Metabolic Advisory Board, had grant support from Novartis. Participant-level genotype and protein data, and full summary association results from the genetic analysis, are available through the European Genotype Archive (accession number EGAS00001002555). Summary association are also available via FTP and through PhenoScanner (<http://www.phenoscanter.medschl.cam.ac.uk>).

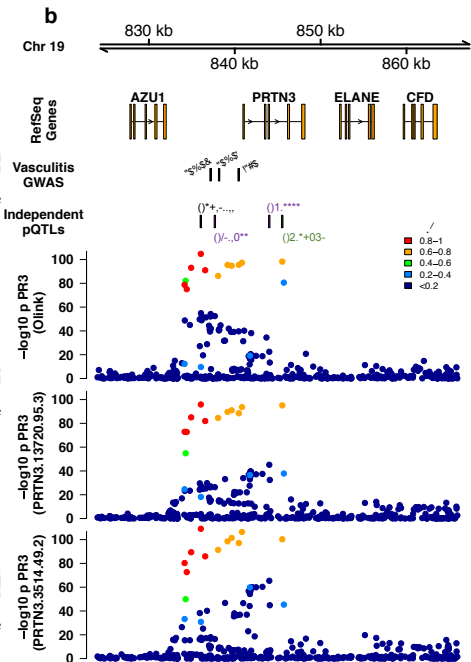
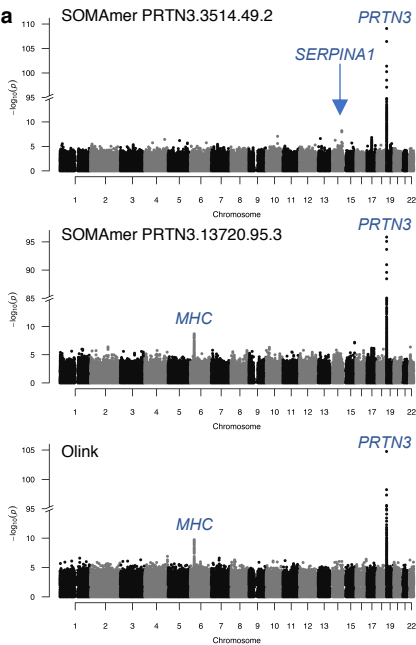
Correspondence to A.S.B. ([asb38@medschl.cam.ac.uk](mailto:asb38@medschl.cam.ac.uk)) and J.D. ([jd292@medschl.cam.ac.uk](mailto:jd292@medschl.cam.ac.uk)).











**c**

