

1 **From genome to phenome via the proteome: broad capture, antibody-based proteomics to**  
2 **explore disease mechanisms**

3 Mine Koprulu<sup>1</sup>, Julia Carrasco-Zanini<sup>1</sup>, Eleanor Wheeler<sup>1</sup>, Sam Lockhart<sup>1,2</sup>, Nicola D. Kerrison<sup>1</sup>,  
4 Nicholas J. Wareham<sup>1</sup>, Maik Pietzner<sup>1,3</sup>, Claudia Langenberg<sup>1,3</sup>

5 1. MRC Epidemiology Unit, University of Cambridge School of Clinical Medicine, Institute of  
6 Metabolic Science, Cambridge, CB2 0QQ, UK.

7 2. MRC Metabolic Diseases Unit, Wellcome-MRC Institute of Metabolic Science, University of  
8 Cambridge, Cambridge, UK

9 3. Computational Medicine, Berlin Institute of Health at Charité-Universitätsmedizin Berlin,  
10 10117 Berlin, Germany.

11

12 Correspondence to:

13 Claudia Langenberg ([claudia.langenberg@mrc-epid.cam.ac.uk](mailto:claudia.langenberg@mrc-epid.cam.ac.uk)), MRC Epidemiology Unit,  
14 University of Cambridge School of Clinical Medicine, Institute of Metabolic Science, Cambridge,  
15 CB2 0QQ, UK.

16 **Abstract**

17 Studying the plasma proteome as the intermediate layer between the genome and the phenome  
18 has the potential to identify disease causing genes and proteins and to improve our  
19 understanding of the underlying mechanisms. Here, we conducted a *cis*-focused proteogenomic  
20 analysis of 2,923 plasma proteins measured in 1,180 individuals using novel antibody-based  
21 assays (Olink® Explore 1536 and Explore Expansion) to identify disease causing genes and  
22 proteins across the human phenome. We describe 1,553 distinct credible sets of protein  
23 quantitative trait loci (pQTL), of which 256 contained *cis*-pQTLs not previously reported. We  
24 identify 224 *cis*-pQTLs shared with 578 unique health outcomes using statistical colocalization,  
25 including, gastrin releasing peptide (GRP) as a potential therapeutic target for type 2 diabetes.  
26 We observed convergence of phenotypic consequences of *cis*-pQTLs and rare loss-of-function  
27 gene burden for twelve protein coding genes (e.g., *TIMD4* and low-density lipoprotein  
28 metabolism), highlighting the complementary nature of both approaches for drug target  
29 prioritization. Proteogenomic evidence also improved causal gene assignment at 40% (n=192) of  
30 overlapping GWAS loci, including *DKKL1* as the candidate causal gene for multiple sclerosis.

31 Our findings demonstrate the ability of broad capture, high-throughput proteomic technologies  
32 to robustly identify new gene-protein-disease links, provide mechanistic insight, and add value  
33 to existing GWASs by enabling and refining causal gene assignment.

34

## 35 Introduction

36 Rare and common sequence variation across the genome contributes to the risk of most human  
37 diseases investigated to date (1). However, the translation of the many established and emerging  
38 genome-to-phenome links is limited by the uncertainty around the underlying causal genes. This  
39 presents a major limitation for experimental follow-up, mechanistic understanding, and use of  
40 the emerging genomic evidence in drug development. Different approaches, such as integration  
41 of tissue-specific gene expression data (2), experimentally derived functional genomic data such  
42 as ChIP-seq or ATAC-seq (3), or functional characterization of candidate variants using CRISPR  
43 screens in cellular models (4) have been used to address this gap and to identify likely causal  
44 genes at risk loci. However, complex regulatory processes take place at each stage of  
45 transcription and translation, which often leads to low correlation between transcripts and  
46 proteins, and cellular models can only approximate complex human biology. Compared to these  
47 methods, the proteogenomic approach has the advantage of focusing on the biologically active  
48 entity - the protein.

49 The development of broad-capture proteomic assays, targeting thousands of proteins in parallel,  
50 now enables proteogenomic approaches which can efficiently identify causal genes by  
51 systematically testing for shared genetic regulation of protein levels or function and disease  
52 susceptibility. This has catalyzed substantial advances in the identification of a) causal genes and  
53 proteins underlying established disease 'loci', and b) molecular 'hubs' that connect the genome  
54 not to one but many diseases through the encoded protein (5-18). Previous large-scale  
55 proteogenomic studies covering thousands of proteins have almost exclusively used aptamer-  
56 based assays (10, 11, 15, 16). Correlations of protein measures from aptamer versus antibody-  
57 based technologies have been shown to vary widely, and proteogenomic results are concordant  
58 for around only 65% based on around 900 overlapping proteins targets (16). To date, antibody-  
59 based proteomic assays have only been available for selected protein panels at scale (9, 14, 18),  
60 but this is changing with the availability of the Olink® Explore 1536 and Olink® Explore Expansion  
61 assays measuring ~1,400 proteins each.

62 The UK Biobank Pharma Proteomics Project (UKB-PPP) project which measured ~1,400 proteins  
63 using Olink® Explore 1536 assay in over 50,000 participants successfully demonstrated the power  
64 of scaling up by cataloguing over 10,000 mainly novel pQTLs (17). However, this study provided  
65 few insights about the translational potential of pQTLs to systematically inform candidate gene  
66 annotation at known risk loci and more importantly, to reveal novel biological roles of proteins  
67 for human health at scale. UKB-PPP and others did demonstrate that genuine and biologically  
68 relevant protein quantitative trait loci (pQTL) can be discovered in as few as hundreds of  
69 individuals (14, 17, 19), suggesting that broader proteomic coverage in even small-scale  
70 proteogenomic studies can make substantial advances to the understanding of diseases if  
71 integrated with large-scale phenomic data.

72 Here we generate antibody-based proteomic data using the Olink® Explore 1536 and Explore  
73 Expansion assays to capture 2,923 proteins in 1,180 individuals. We perform genetic fine-  
74 mapping at protein coding genes ( $\pm 500\text{kb}$ ) and enhance the understanding of disease  
75 mechanisms by systematically integrating cis-pQTLs with thousands of diseases and health  
76 measures to (a) refine the candidate causal gene assignment at existing disease susceptibility loci  
77 at scale and (b) identify novel disease mechanisms in phenome-wide colocalization analyses.

78

## 79 Results

### 80 Identification and fine-mapping of cis-proteogenomic signals for 2,923 protein targets

81 We adopted a Bayesian fine-mapping strategy (20) to identify proximal acting genetic variants  
82 (cis-pQTLs,  $\pm 500$ kb around the protein coding gene) that were associated with plasma abundance  
83 of 2,923 proteins measured in 1,180 participants of the EPIC-Norfolk cohort (21) (**Supplementary**  
84 **Table 1**). We identified a total of 1,553 independent credible sets for 914 unique protein targets  
85 for which sentinel variants reached genome-wide significance ( $p < 5 \times 10^{-8}$ ) when modelled jointly  
86 at each protein coding locus (**Fig. 1A, Supplementary Table 2**). The number of independent  
87 credible sets for each protein target ranged between one and eight (mean=1.64, IQR=1-2),  
88 illustrating wide-spread allelic heterogeneity at protein coding loci. We observed a high  
89 replication rate (89.9%, 910 out of 1,013) for credible sets of 590 protein targets overlapping with  
90 the UKBB-PPP effort. Conversely, we identified 4.5% of the 20,540 reported signals in as few as  
91 1,200 participants that were reported based on more than 35,000 participants of the UKB-PPP.  
92 A total of 256 (16.5%) credible sets contained cis-pQTLs not previously reported, including 131  
93 proteins that have not been measured by previous platforms (**Fig. 1B**) (5-18). Notably 125 signals  
94 were for 101 previously targeted proteins, the majority of which ( $n=92$  proteins) have been  
95 targeted using non-antibody-based technologies in samples sizes up to 30 times larger than ours  
96 (10, 11, 15, 16) (**Fig. 1B**).

97 Effect size and minor allele frequency distributions of unreported cis-pQTLs were comparable to  
98 the 1,297 (83.5%) successfully replicated cis-pQTLs (5-18) (**Supplementary Table 2**), illustrating  
99 that complementary proteomic technologies can still identify genetic variants that would have  
100 been anticipated to be seen in previous studies (**Fig. 1A, Fig. 1D**).

101 We observed a strong inverse relationship between the absolute effect sizes of cis-pQTLs and the  
102  $\log_{10}$ -transformed frequencies of their minor alleles ( $r = -0.78$ ;  $p < 1 \times 10^{-300}$ ), likely due to the more  
103 severe predicted consequences of rarer alleles, such as stop-gain mutations (**Fig. 1B-C**). We also  
104 report 482 cis-pQTLs with a minor allele frequency (MAF) above 5% with large absolute effect  
105 sizes (range 0.5-1.72 s.d. per allele), suggesting strong genetic control of the associated proteins.  
106 Of these, less than half (35.9%) were protein altering variants themselves or were in strong

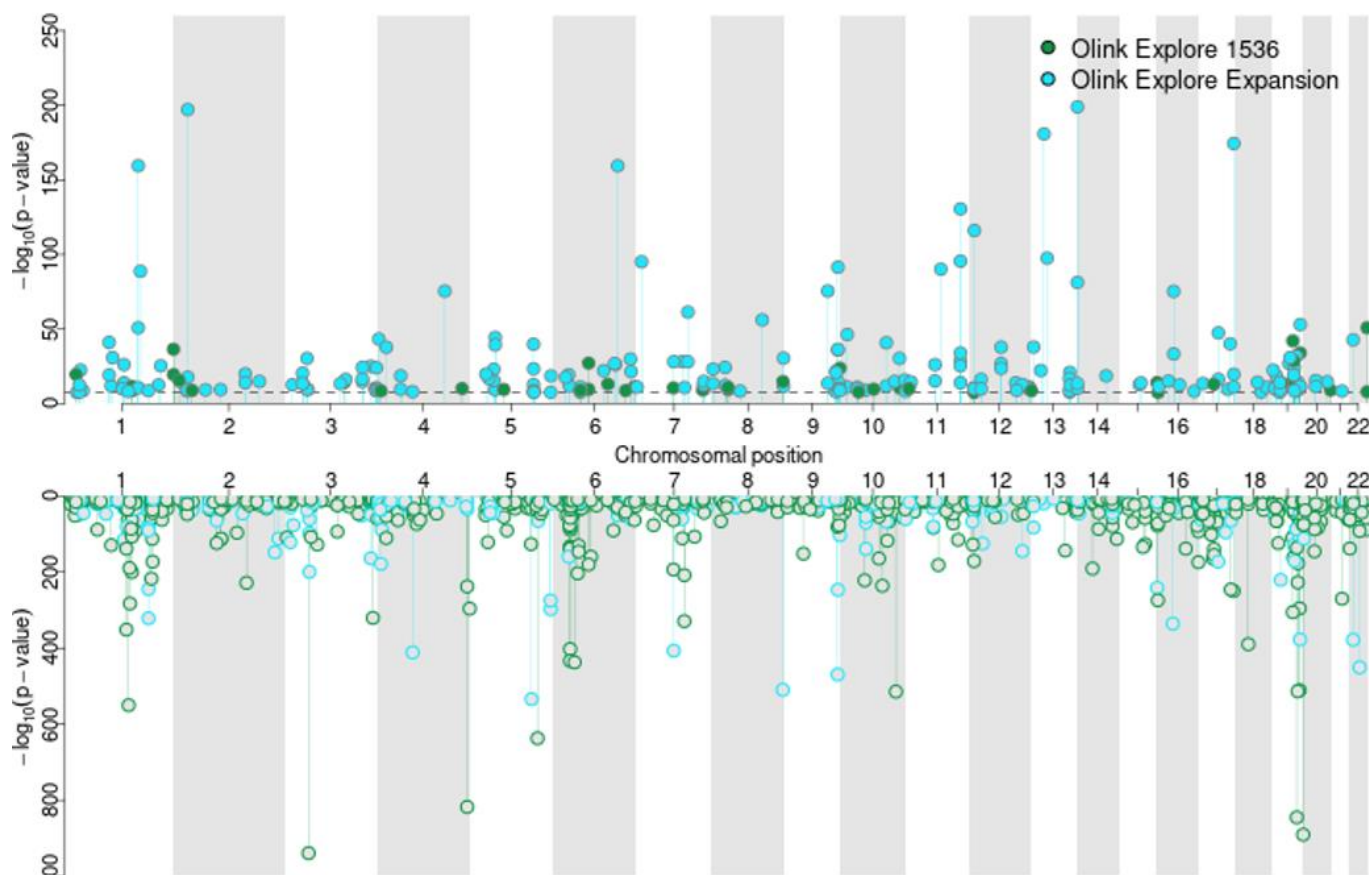
107 linkage disequilibrium (LD;  $r^2 > 0.6$ ) with one, potentially affecting the binding affinity of  
108 antibodies. Proteins with at least one significant cis-pQTL were enriched for characteristics of  
109 secreted proteins, like the presence of disulfide-bonds (odds ratio: 4.47; p-value= $3.0 \times 10^{-74}$ ) or  
110 glycosylation sites (odds ratio: 2.22; p-value= $1.1 \times 10^{-13}$ ), but depleted of sites for posttranslational  
111 modifications that are important for intracellular signaling, like phosphorylation (odds ratio: 0.42;  
112 p-value= $5.5 \times 10^{-14}$ ) or ubiquitination (odds ratio: 0.30; p-value= $4.2 \times 10^{-11}$ ).

113 Finally, for more than half of the protein targets (n=532) with at least one pQTL, we observed  
114 strong evidence of colocalization (PP>80%) between a cis-pQTL and the corresponding gene  
115 expression QTL (eQTL) signal in at least one out of 49 tissues of the GTEx resource  
116 (**Supplementary Table 3**). These results suggest altered expression of protein coding genes in one  
117 or multiple tissues as the major source for cis associations observed with plasma protein levels.

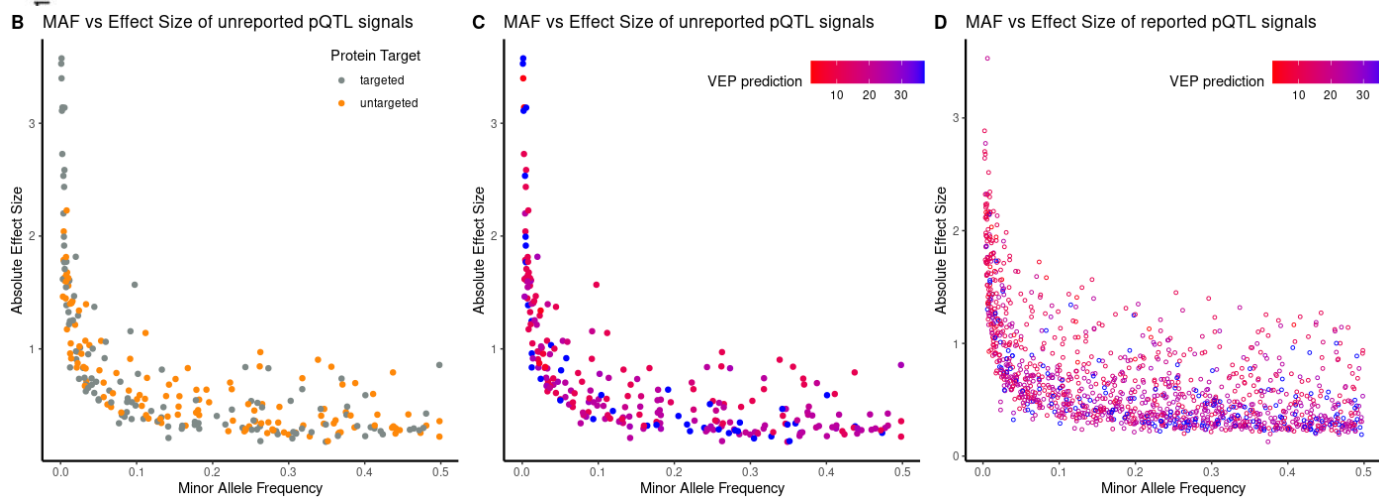
118

119

120



121



122

123 **Figure 1: Genetic regulation of 2,923 proteins measured by the Olink Explore 1536 and Olink Explore**  
 124 **Expansion platforms in 1,180 individuals.** Previously unreported and reported pQTLs are represented  
 125 with a filled and hollow circle, respectively. Only the variants which are genome-wide significant ( $p$ -  
 126 value  $< 5 \times 10^{-8}$ ) in the joint model (see Methods) are presented. **A. Miami plot representing the**  
 127 **independent lead cis-pQTLs identified through Bayesian fine-mapping for 914 unique proteins.** Shown  
 128 are p-values from a linear regression model modelling all identified credible set variants for a given protein  
 129 target jointly. *Top*: Lead cis-pQTL signals unreported to date. *Bottom*: Lead cis-pQTL signals which were in

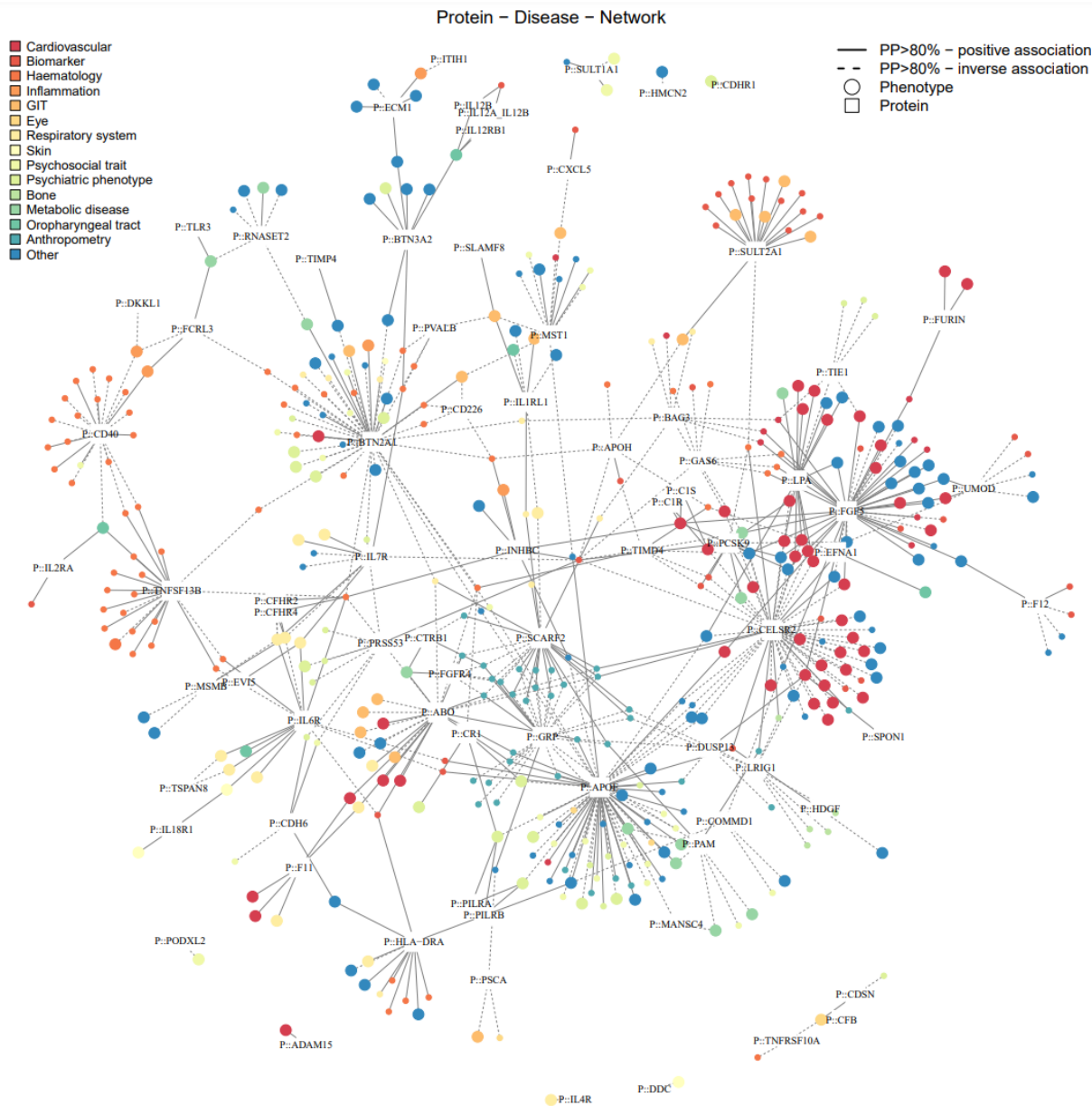
130 linkage disequilibrium (LD;  $r^2 > 0.5$ ) with a previously reported pQTL. **B. Minor allele frequency vs effect**  
131 **size of unreported pQTL signals, coloured by whether the protein has previously been targeted.**  
132 Unreported pQTL signals for a previously targeted protein are coloured grey and those for a previously  
133 untargeted protein are coloured orange. **C. Minor allele frequency vs effect size of unreported pQTL**  
134 **signals, coloured by most severe variant consequence prediction.** The colour coding represents the most  
135 severe Variant Effect Predictor (22) consequence of the lead cis-pQTL, or variants in LD ( $r^2 > 0.6$ ) within the  
136 protein encoding gene. The most severe consequence is coloured red (Ensembl consequence rank = 1)  
137 and the least severe consequence is coloured blue (Ensembl consequence rank = 37). **D. Minor allele**  
138 **frequency vs effect size of reported pQTL signals, coloured by most severe variant prediction.** The colour  
139 coding represents the most severe Variant Effect Predictor (22) consequence of the lead cis-pQTL, or  
140 variants in LD ( $r^2 > 0.6$ ) with the lead cis-pQTL within the protein encoding gene. The most severe  
141 consequence is coloured red (Ensembl consequence rank = 1) and the least severe consequence is  
142 coloured blue (Ensembl consequence rank = 37).

143

#### 144 **From genome to phenome via the proteome**

145 The genome is linked to the phenome via the proteome and the translational potential of pQTLs  
146 is due to their ability to link insights about the genetic regulation of protein levels and function  
147 to diseases (15). We identified 1,110 robust protein – phenotype pairs (**Fig. 2**; posterior  
148 probability [PP] > 80% of a shared genetic signal) comprising 224 protein targets for 575 unique  
149 traits by systematically testing for a shared genetic architecture at protein coding loci ( $\pm 500$ kb)  
150 across the phenome (see **Methods**; **Supplementary Table 4**). This included well-described  
151 examples, such as UMOD and kidney disease or established drug targets like PCSK9 and LDL-  
152 cholesterol, but importantly 93 protein targets connected with at least one phenotype that have  
153 been missed by previous aptamer-based efforts.





154

155

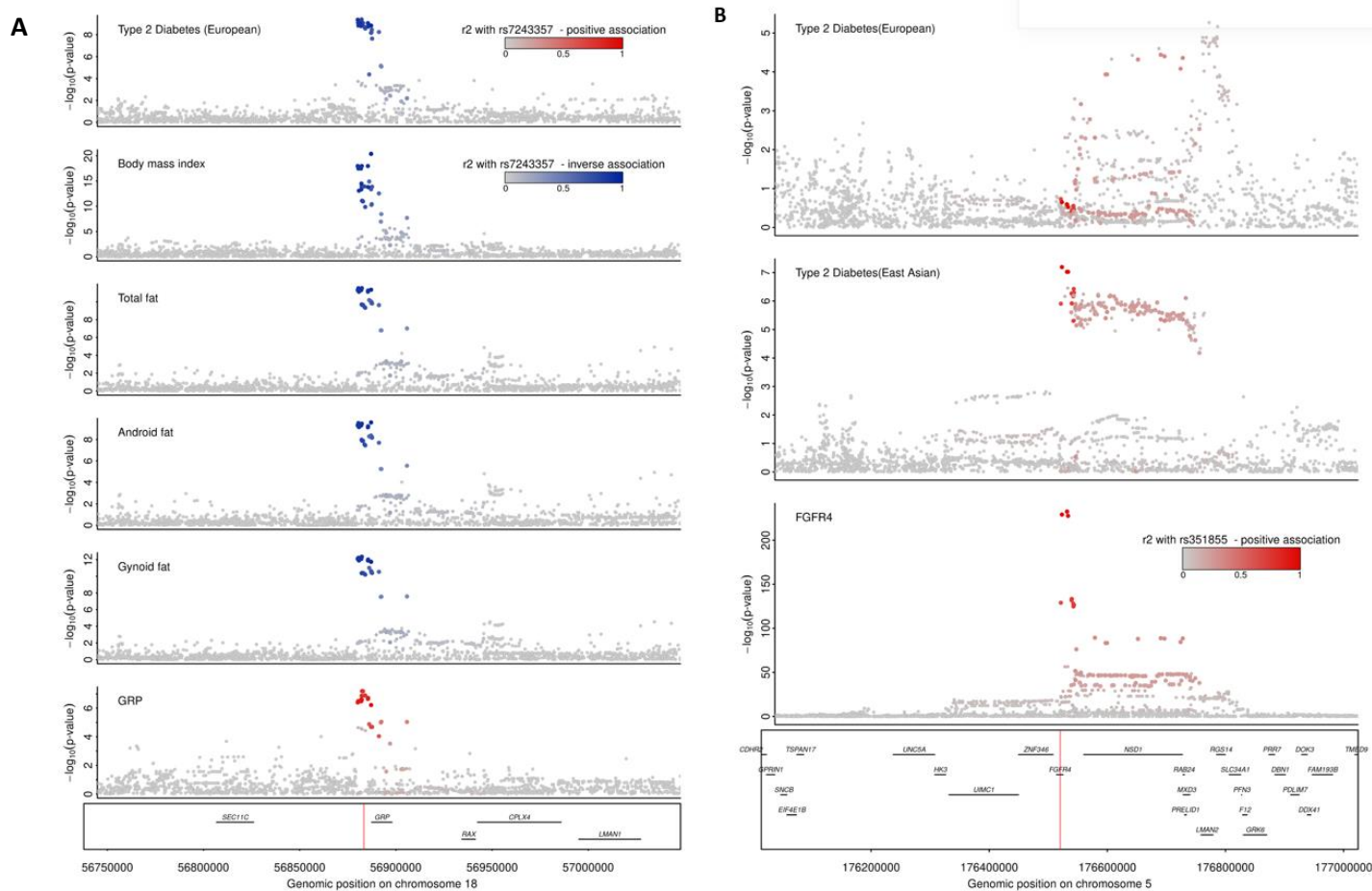
156 **Figure 2: Protein – disease network.** Results from phenome-wide colocalization at protein coding  
 157 loci ( $\pm 500\text{kb}$ ) are shown. For simplicity, only proteins with at least one binary outcome (i.e.,  
 158 mainly diseases) association are included. Proteins are presented with a square, binary outcomes  
 159 are presented with large circles, and continuous outcomes are presented with small circles. The  
 160 colour for the circles present the trait category. Edges between proteins and phenotypes  
 161 represent strong evidence for a shared genetic signal ( $PP > 80\%$  and  $LD$  between regional sentinel  
 162 variants  $> 0.8$ ). Effect directions are indicated by the line type (solid = higher protein abundance,  
 163 increased risk, dashed = higher protein abundance, reduced risk) and derived based on the lead  
 164 cis-pQTL at the corresponding locus. The full list of colocalization results can be found in  
 165 Supplementary Table 4. Abbreviations: GIT, gastrointestinal tract.

166

167 One of the examples is gastrin releasing peptide (GRP, encoded by *GRP*), for which we observed  
168 strong evidence of colocalization (posterior probability [PP] =82.5%) between plasma levels and  
169 type 2 diabetes (T2D) risk at an established GWAS locus (18q21) for which different genes had  
170 been prioritized, including *SEC11C*, *GRP*, and *MC4R* (23-25). The GRP-increasing G-allele of the  
171 lead cis-pQTL (rs1517035; MAF=0.18) was associated with a reduced risk for T2D (odds  
172 ratio=0.96, p-value= $7.8 \times 10^{-10}$ ). GRP is a neuropeptide named for its ability to stimulate secretion  
173 of the gastric acid secretagogue, gastrin, in the stomach (26, 27), but it is likely involved in other  
174 metabolic pathways. We obtained strong evidence that GRP likely mediates T2D risk via an effect  
175 on overall obesity, based on the convergence of evidence from mice studies, human trials, and  
176 human genetic data. We established a shared genetic signal between plasma GRP, body mass  
177 index and fat, and T2D risk using multi-trait colocalization with coherent effect directions (**Fig. 3**).  
178 GRP induces satiety in mice via its cognate GRP receptor (*Grpr*) (28, 29). Further, mice lacking  
179 *Grpr* show impaired glucose tolerance after gastric glucose administration (30) and gain excess  
180 body weight under *ad libidum* conditions (29). These observations have been corroborated by  
181 human trials, in which treatment with human recombinant GRP (hrGRP) led to weight loss  
182 through reduced food intake (31). In summary, our results motivate investigations into hrGRP for  
183 appetite control and body weight lowering to possibly assist in T2D management and remission,  
184 an approach similar to recently implemented treatment strategies targeting incretins, like GLP-  
185 1, and associated receptors, with preliminary evidence of an additive effect in rats (32).

186

187



188

189 **Figure 3: Stacked regional association plots for the multi-trait colocalization. A. Stacked regional**  
 190 **association plots for the multi-trait colocalization of the GRP cis-pQTL with gynoid fat, android fat, total**  
 191 **body fat, body mass index and type 2 diabetes.** The top candidate SNP highlighted by multi-trait  
 192 colocalization (rs7243357) and lead cis-pQTL for GRP (rs1517035) are in strong LD ( $r^2=0.8$ ). Gynoid fat,  
 193 android fat and total body fat phenotypes are based on UK Biobank and were analysed in-house using  
 194 BOLT-LMM (33). **B. Stacked regional association plot the multi-trait colocalization of the FGFR4 cis-pQTL**  
 195 **with type 2 diabetes in East Asian populations.** Red colouring represents a positive effect direction in  
 196 reference to the protein increasing allele for GRP whereas blue represent an inverse association. The hue  
 197 of the colour represents the strength of  $r^2$  representing the LD structure, as indicated on the legend.  
 198 European Type 2 diabetes summary statistics were obtained from dbGAP Million Veteran Program (MVP)  
 199 European subset ( $n_{\text{cases}}=148,726$ ,  $n_{\text{controls}}=965,732$ ) (34). East Asian Type 2 diabetes summary statistics  
 200 were obtained from Mahajan et al (2022) ( $n_{\text{cases}}=56,268$ ,  $n_{\text{controls}}=227,155$ ) (25). The body mass index  
 201 summary statistics were obtained from Pulit et al. (2019) ( $n=806,834$ ) (35).

202

203

204 Several T2D loci have been reported to be specific to certain ancestries (23-25). In the absence  
205 of strong differences in allele frequency, such ancestry specific effects could be caused by a  
206 variety of different factors, including environmental factors such as dietary intake. We obtained  
207 robust evidence that *FGFR4* is the candidate causal gene at the East Asian-specific *FGFR4-NSD1*  
208 locus supported by a high posterior probability (PP=97%) for a shared genetic signal with plasma  
209 levels of the gene product fibroblast growth factor receptor 4 (FGFR4) and trans-ancestral  
210 conserved LD between regional sentinel variants ( $r^2 > 0.96$ ; **Fig. 3**). The protein-increasing A-allele  
211 of the lead cis-pQTL (rs351855, beta= 1.01, p-value= $9.8 \times 10^{-234}$ ,  $EAF_{\text{European}}=0.30$ ,  $EAF_{\text{EastAsian}}=0.46$ )  
212 was associated with an increased risk for T2D (beta=0.05, p-value= $1.1 \times 10^{-7}$ ). Candidate gene  
213 studies have implicated rs351855 (p.G388R) in cancer susceptibility (36-38), and subsequent  
214 mechanistic studies showed a gain of function of the mutant FGFR4 by binding transducer and  
215 activator of transcription 3 (STAT3) (39). While we found no evidence for an association to cancer,  
216 there are different studies that support our observation of FGFR4 in T2D-related pathways  
217 including hepatic glucose, bile, and lipid metabolism, and possibly insulin signaling in a diet-  
218 dependent manner (40-43). Briefly, *Fgfr4*<sup>-/-</sup> mice fed a normal chow diet exhibit insulin resistance  
219 and impaired glucose tolerance compared to wild-type controls, however, this difference is not  
220 observed in high-fat diet fed mice. A similar masked genetic effect is seen with the mutant protein  
221 in mice and small observational studies in humans (44). The ability of diet to obscure genetic  
222 effects may explain the ancestral-specific effect in the absence of strong differences in allele  
223 frequencies, with high-fat diet conditions being substantially more common in Western-style  
224 countries of predominantly European ancestry compared to East Asia (45), in particular Japan, in  
225 line with Biobank Japan (p-value<sub>T2D</sub> <  $7.6 \times 10^{-11}$ ) being the largest contributing population to the  
226 East Asian T2D meta-analysis (25).

227

## 228 **Proteogenomic guided annotation of genes at loci reported for diseases and traits related to** 229 **human health**

230 Annotation of the candidate causal genes at disease susceptibility loci is the major bottleneck in  
231 the translation of GWAS into biological and possibly clinical insights (46). We exploited the

232 genomic proximity between cis-pQTLs and the protein coding gene for gene annotation by  
233 systematically overlapping identified credible sets in this study with reported risk loci ( $p < 5 \times 10^{-8}$ )  
234 from the GWAS catalog (downloaded on 23/03/2022; (1)). We identified 480 credible sets  
235 targeting 395 unique proteins (43.2% of all, 914 unique protein targets) for which the lead cis-  
236 pQTL or a proxy ( $r^2 > 0.8$ ) had been reported for one or more of 5,391 collated traits in the GWAS  
237 catalog (**Fig. 4 and Supplemental Tab. 5, see Methods**). For 40% ( $n=192$ ) of those, we prioritized  
238 a gene that was different from the one originally reported, of which 50% ( $n=96$ ) were not the  
239 gene nearest to the GWAS sentinel variant. We further refined a longer list of putative causal  
240 genes to a single one for an additional 31 cis-regions (6.5%). These results exemplify the unique  
241 potential of cis-pQTLs for gene annotation of loci reported across diseases and traits related to  
242 human health (**Fig. 4 and Supplemental Tab. 5**), with one example outlined in more detail below.

243 Multiple Sclerosis (MS) is an autoimmune, inflammatory, and neurodegenerative disease of the  
244 central nervous system that is caused by both genetic and strong environmental factors (47). A  
245 strong signal at 19q13.33 is one of 233 reported GWAS loci (48). Several variants in high LD  
246 ( $r^2 > 0.9$ ) reported for this locus have been linked to different candidate causal genes, including  
247 *DKKL1*, *CD37*, and *SLC6A16* which was the most recently annotated gene based on a sophisticated  
248 ensemble of methods (48). We identify a shared genetic signal (PP=96.1%, **Supplemental Fig. 1**)  
249 between dickkopf like acrosomal protein 1 (DKKL1), encoded by *DKKL1*, and MS at this locus, led  
250 by a cis-pQTL (rs2288480; MAF=0.25) in high LD ( $r^2=0.97$ ) with the lead MS variant (rs1465697;  
251 MAF=0.33, OR=1.09,  $p$ -value= $3 \times 10^{-18}$ ). We note that the lead cis-pQTL was also in LD ( $r^2=0.97$ )  
252 with a recently identified variant at the same locus for systemic lupus erythematosus (SLE) among  
253 East Asians (49).

254 The lead cis-pQTL is in strong LD ( $r^2 > 0.8$ ) with a cluster of three common missense variants  
255 (rs2288481, rs2303759, and rs1054770) that might impair protein function or processing. Little  
256 is known about the biological role of DKKL1 in general, but a non-essential role in  
257 spermatogenesis has been described (50). However, a link towards MS and/or SLE might be  
258 conceivable via a possible role of DKKL1 in adaptive immunity and hence the inflammatory  
259 component of MS. Briefly, *DKKL1* expression is enriched among memory B-cells (51) and an  
260 independent secondary cis-pQTL (rs66532151, MAF=22.4%) for DKKL1 tagged ( $r^2 > 0.96$ ) a cluster

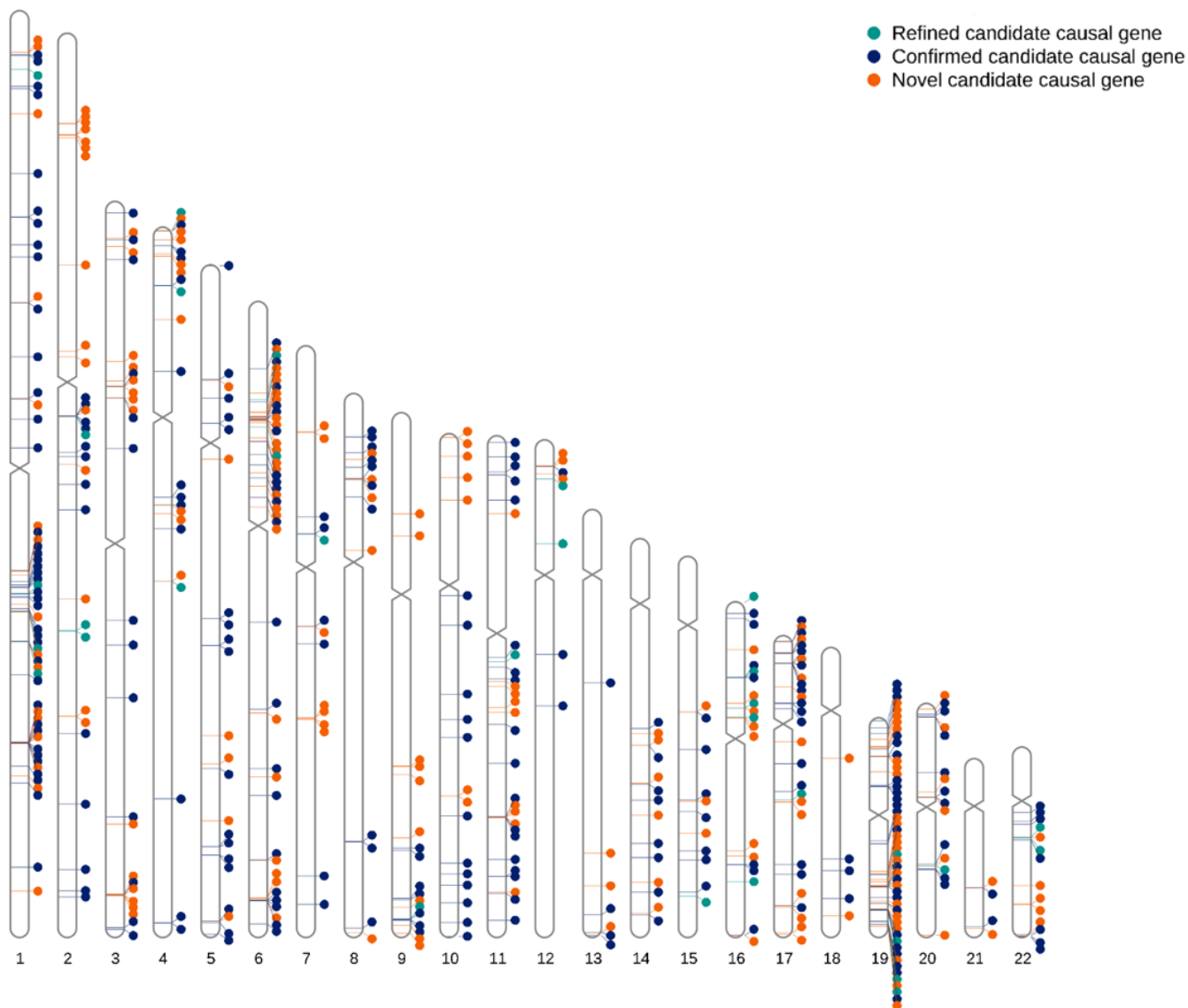
261 of variants associated with different characteristics of CD20<sup>+</sup> memory B-cells (52). This cis-pQTL  
262 (rs66532151) was associated with MS at p-value= $3.4 \times 10^{-7}$ , providing late genetic evidence for  
263 depletion of B-cells being one of the most effective treatments for MS, a therapeutic strategy  
264 that originally emerged from clinical and neuropathological studies (53, 54). Further follow-up  
265 studies are needed to clarify a possible role of DKKL1 in immune cells and whether DKKL1 may  
266 play a role in B-cell hyperactivity observed in MS (53).

267

268

269

270



271

272 **Figure 4: Candidate causal gene assignment at reported GWAS loci using pQTLs. The overlap between**  
273 **existing GWAS risk loci and pQTL loci (n=480) are marked on the human chromosome karyotypes**  
274 **(chromosomes 1-22). The locus is coloured orange if the pQTL provides a novel candidate causal gene**  
275 **assignment for one or more traits, light blue if it refines a candidate causal gene from a longer list of**  
276 **reported or closest genes, and dark blue if it confirms the candidate causal gene assignment provided by**  
277 **the GWAS.**

278

279 Multiple independent genetic variants associated with the same protein target at the same locus,  
280 so-called allelic heterogeneity, provides the highest confidence in gene assignment but can also  
281 highlight differential biological roles for the same protein. We observed 73 such protein targets

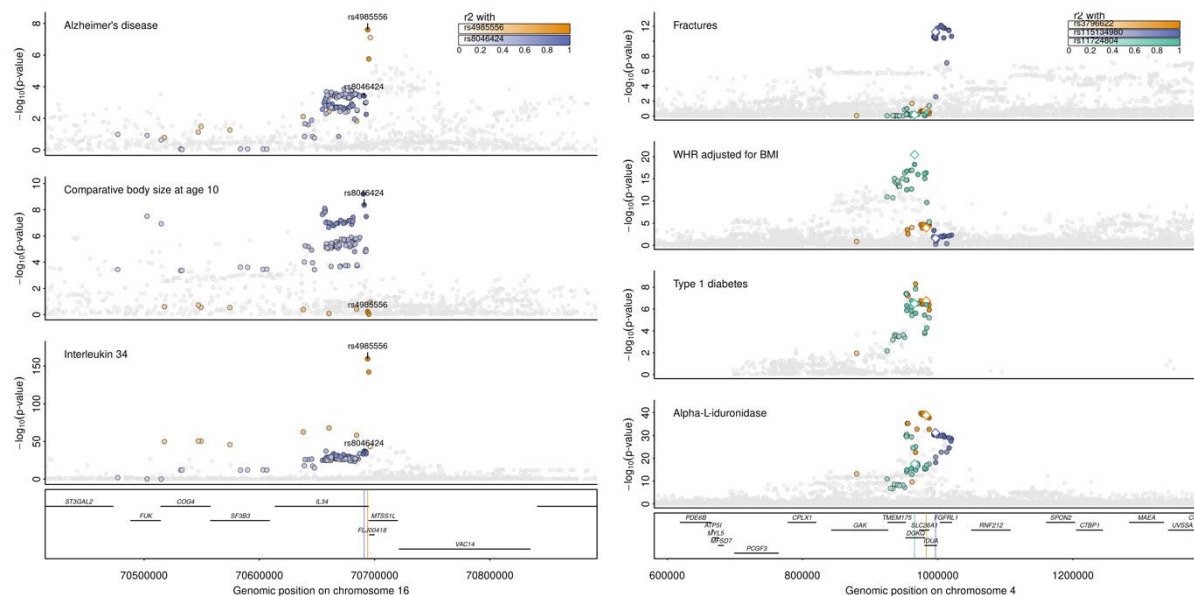


282 with two or more credible sets including distinct GWAS variants for related and unrelated traits.  
283 For example, we discovered three distinct credible sets for plasma levels of interleukin 34 (IL-34)  
284 at 16q22.1. Two contained independent ( $r^2=0.07$ ) lead cis-pQTL variants with distinct structural  
285 consequences on the protein, that were also associated with two distinct outcomes – Alzheimer’s  
286 disease (55) and childhood obesity (56) (**Fig. 5**). rs4985556 is associated with increased risk for  
287 Alzheimer’s disease (MAF=12.2%; beta=0.07, p-value< $2.3 \times 10^{-8}$ ) and introduces a premature stop  
288 (p.Tyr213Ter), truncating the protein and likely affecting dimerization and possibly secretion.  
289 rs8046424 (alternate allele: C;  $r^2=0.96$  with lead sentinel childhood obesity variant rs4985555) is  
290 associated with reduced childhood obesity (MAF = 48.2%; beta=-0.008, p-value< $4.4 \times 10^{-9}$ ). It is a  
291 missense variant (p.Glu123Gln) of moderate consequence (CADD score 11.6) that maps to a  
292 binding domain of the cognate IL-34 receptor CSF-1R (57). Therefore, both variants will likely  
293 strongly (rs4985556) or moderately (rs8046424) attenuate signaling via CSF-1R, which has been  
294 shown to drive cerebrovascular pathologies that are common in Alzheimer’s disease (58). While  
295 the gradient of structural consequences translates into a graded effect on Alzheimer’s disease  
296 (rs8046424, beta=0.03, p< $3.6 \times 10^{-4}$ ), the absence of any effect of the more detrimental variant  
297 (rs4985556) on childhood obesity (beta=-0.001, p=0.59) might point to a different, yet to be  
298 defined, pathway.

299 We observed a similar segregation of phenotypes across distinct cis-pQTLs for alpha-L-  
300 iduronidase encoded at *IDUA*. Briefly, three out of four detected credible sets contained GWAS  
301 risk loci or strong proxies ( $r^2>0.8$ ) for fractures (59) (rs115134980; MAF=16.1%; beta=-0.06, p-  
302 value= $7.4 \times 10^{-12}$ ), waist-to-hip ratio adjusted for BMI (60) and inflammatory diseases (61)  
303 (rs11724804; MAF=44.7%; beta=-0.017, p-value< $7.6 \times 10^{-21}$ ), as well as type 1 diabetes (62)  
304 (rs3796622; MAF=35.2%; beta=-0.07, p-value< $1.7 \times 10^{-7}$ ) (**Fig. 5**). Alpha-L-iduronidase is essential  
305 for the breakdown of glycosaminoglycans within lysosomes and numerous rare pathogenic  
306 variants within *IDUA* are known to cause accumulation of glycosaminoglycans in lysosomes  
307 (mucopolysaccharidosis type I [MPS-1]). Patients present with a wide spectrum of complications,  
308 such as skeletal deformities or organomegaly, that has been attributed to the variable impact of  
309 mutations on enzyme activity, with nonsense mutations causing most severe diseases (Hurler  
310 syndrome) (63). While skeletal abnormalities in rare disease patients may relate to bone



311 phenotypes seen for the common cis-pQTL, there are no reports for an elevated risk for  
 312 inflammatory or autoimmune disease among MPS-1 patients or other evidence from rare variant  
 313 analysis. Tissue-dependent effects of common variants might be one explanation for the  
 314 different phenotypes linked to distinct cis-pQTLs for alpha-L-iduronidase.  
 315



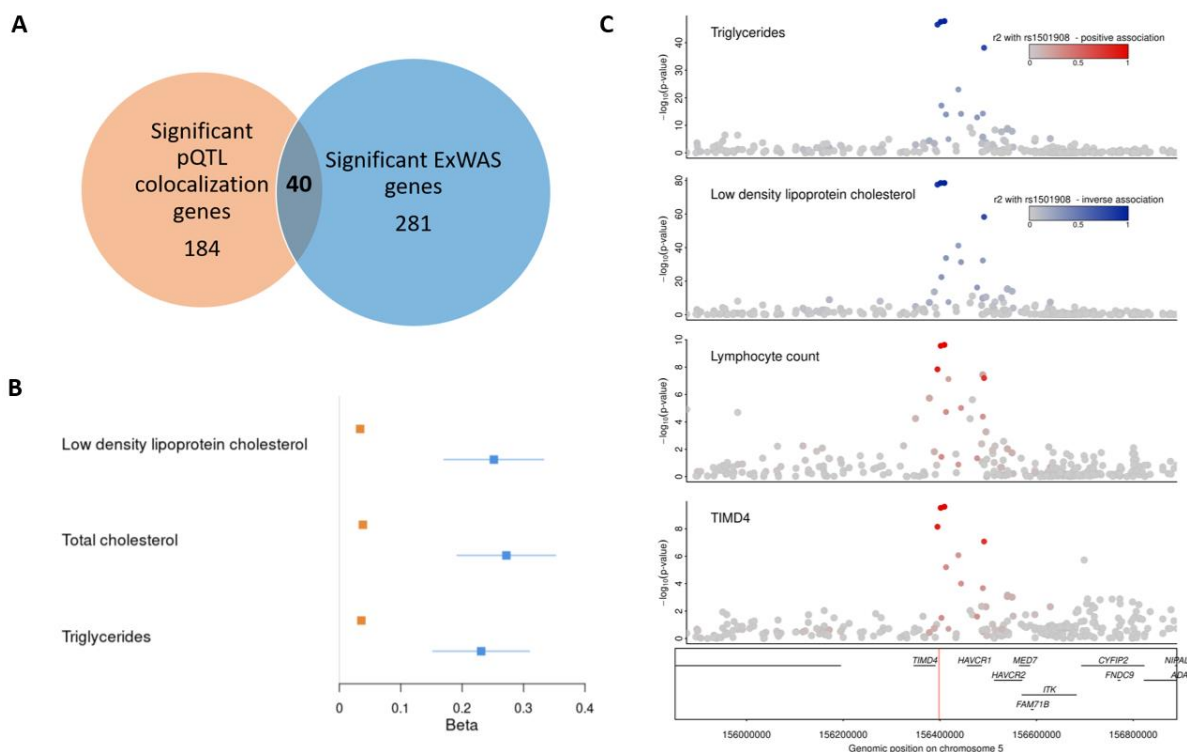
316 **Figure 5: Allelic heterogeneity at protein coding loci translates into distinct phenotypic consequences**  
 317 **at *IL34* and *IDUA*.** *Left* Regional association plots centered around *IL34* ( $\pm 200$ kb) for plasma interleukin 34  
 318 levels, comparative body size at age 10 (56), and Alzheimer’s disease (55). Shown are association statistics  
 319 ( $p$ -values) from genome-wide association analysis. Single genetic variants were coloured based on LD with  
 320 two distinct cis-pQTLs (rs4985556 – orange; rs8046424 – purple). *Right* Regional associations plots  
 321 centered around *IDUA* ( $\pm 400$ kb) for plasma alpha-L-iduronidase levels, type 1 diabetes (62), waist-to-hip  
 322 ratio (WHR) adjusted for body mass index (BMI) (60), and risk of fractures (59). Shown are association  
 323 statistics ( $p$ -values) from genome-wide association analysis. Single genetic variants were coloured based  
 324 on LD with three distinct cis-pQTLs (rs3796522 – orange; rs115134980 – purple; rs11724804 – green).  
 325 Lead cis-pQTLs are highlighted by hollow diamonds.  
 326

327  
 328 **Phenotypic convergence of rare variant burden and common cis-pQTLs for protein coding**  
 329 **genes**

330 Much effort and funding has been invested into biobank-scale whole-exome sequencing studies  
 331 (ExWAS) to identify rare deleterious genetic variants and novel disease candidate genes for the  
 332 development of treatment strategies (64, 65). However, it is unknown how efforts focusing on

333 the rare deleterious end of gene (and protein) dysfunction relate to less extreme alterations of  
 334 protein levels or function. To explore whether evidence from ExWAS and our cis-based phenome-  
 335 wide colocalization analyses converge for disease-linked genes, we systematically integrated our  
 336 results with those from a recent ExWAS among ~450,000 UK Biobank participants across almost  
 337 4,000 phenotypes (64).

338 Among 2,939 protein coding genes covered by the Olink Explore 1536 and Explore Expansion  
 339 platforms, 40 (1.3%) showed evidence for phenotypic associations with a rare variant gene-  
 340 burden and statistical colocalization with a cis-pQTL, whereas 281 and 184 protein coding genes  
 341 were linked to phenotypes through ExWAS or cis-pQTLs only, respectively (**Fig. 6**). Out of the 40  
 342 overlapping genes, we observed phenotypic convergence for only 12 genes across 21 phenotypes  
 343 following manual review to harmonize phenotype definitions (**Supplementary Tab. 6**). These  
 344 results clearly exemplify the complementary nature of both approaches and the unique ability of  
 345 bespoke proteogenomic experiments to prioritize disease mediators and hence putative  
 346 therapeutic targets.



347

348 **Figure 6: Phenotypic convergence of rare variant burden and common cis-pQTLs for protein coding**  
349 **genes and TIMD4 as an example.** A. Venn diagram showing the number of genes with a significant rare  
350 variant gene burden association ( $p < 1E-06$ ) with at least one trait (64) in blue and the number of genes  
351 with a significant pQTL colocalization ( $PP > 80\%$ ) with at least one trait in orange. All 2,939 unique genes  
352 covered by Olink Explore 1536 and Explore Expansion assays were investigated. **B. Forest plot comparing**  
353 **the effect size estimates between TIMD4 cis-pQTL (rs58198139) and rare TIMD4 loss of function (LoF)**  
354 **gene-burden results (variant group: missense and loss of function variants with a minor allele frequency**  
355 **< 1%) for low density lipoprotein cholesterol, total cholesterol and triglyceride levels.** Rare TIMD4 loss  
356 of function (LoF) gene-burden results are shown in blue and TIMD4 cis-pQTL associations are shown in  
357 orange. C. **Stacked regional plot of the multi-trait colocalization of TIMD4 cis-pQTL with lymphocyte**  
358 **count, low density lipoprotein cholesterol, and triglycerides.** Red colouring represents a positive effect  
359 direction with protein increasing allele with TIMD4 whereas blue represent an inverse association. The  
360 hue of the colour represents the strength of  $r^2$  representing the LD structure, as indicated on the legend.

361

362 Convergence of phenotypic consequences from rare gene burden and common cis-pQTLs not  
363 only provides compelling evidence for causal gene assignment but can establish dose-response  
364 relationships that are an essential prerequisite for genetically informed drug discovery (66). We  
365 observed such a dose-response relationship between putative functional consequences for T-cell  
366 immunoglobulin and mucin domain containing 4 (TIMD4) and LDL-cholesterol as well as total  
367 triglyceride, but not HDL-cholesterol levels in blood (**Fig. 6C**). The protein-decreasing T-allele of  
368 the lead cis-pQTL (rs58198139) was associated with moderate effects on LDL-cholesterol in UK  
369 Biobank ( $MAF = 0.26$ ;  $\beta_{LDL} = 0.03$ ,  $p\text{-value}_{LDL} = 7 \times 10^{-44}$ ), likely mediated by altered protein  
370 expression, while the cumulative burden of rare loss-of-function variants was associated with  
371 substantially higher LDL levels ( $\beta_{LDL} = 0.25$ ,  $p\text{-value}_{LDL} = 1.51 \times 10^{-9}$ , variant mask: predicted loss  
372 of function and deleterious missense variants with  $MAF < 1\%$ ), in line with this locus being one of  
373 the earliest discovered loci for polygenic dyslipidemia but with few functional insights gained  
374 since (67). TIMD4 is best known for its role in tissue-dependent macrophage efferocytosis of  
375 apoptotic cells (68, 69) but does also participate in T-cell activation and recruitment (70).  
376 Accordingly, *Timd4*<sup>-/-</sup> mice show impaired macrophage phagocytosis and increased lymphocyte  
377 cell counts (71), an observation recapitulated by our phenome-wide colocalization analyses  
378 identifying an inverse association for the protein-decreasing T-allele for lymphocyte counts  
379 ( $\beta = 1.02$ ,  $p\text{-value} = 1.4 \times 10^{-12}$ ) with high certainty ( $PP = 97.5\%$ ). Circulating leucocytes and  
380 resident M2 macrophages can take up cholesterol from circulating LDL particles and sequestered  
381 lipoproteins in the vasculature, but classical pathways, like the LDL-receptor mediated uptake,

382 were shown, at least in mice, to have no substantial effect on plasma LDL-cholesterol levels (72).  
383 In contrast, more recent work demonstrated the ability of TIMD4<sup>+</sup> adipose tissue macrophages  
384 to significantly contribute to the regulation of post-prandial HDL-cholesterol levels in mice (73).  
385 While there was no difference in triglycerides or non-HDL cholesterol following TIMD4  
386 blockade, TIMD4 blockade inhibited LDL-induced lysosomal activity *in vitro*, suggesting a role for  
387 TIMD4 in peripheral LDL cholesterol processing. These findings provide evidence of a role for  
388 TIMD4 in the regulation of systemic lipoprotein metabolism, and taken together with our  
389 proteogenomic findings, provide a compelling rationale to explore the role of TIMD4<sup>+</sup>  
390 macrophages in systemic LDL cholesterol metabolism. In general, increased uptake of modified  
391 LDL-cholesterol particles by resident macrophages contributes to atherosclerotic foam cell  
392 formation, a major cardiovascular risk factor (74). We observed no conclusive evidence that  
393 either the cis-pQTL (protein decreasing T-allele; odds ratio [95% CI] = 1.04 [1.02-1.07], p-  
394 value=0.002) or the cumulative burden of the loss of function variation in the gene (odds ratio  
395 [95% CI] = 1.3 [0.90-1.88], p-value=0.16) were associated with coronary artery disease (CAD; the  
396 currently most powered GWAS for atherosclerotic consequences). However, our findings urge  
397 further investigations into the functional role of TIMD4 in immune cell-mediated LDL-cholesterol  
398 turnover and foam cell formation. Although the extent to which this mechanism can contribute  
399 to addressing CAD burden is currently unclear, blocking of TIMD4 using monoclonal antibodies  
400 increased atherosclerotic lesion size in *Ldlr*<sup>-/-</sup> mice although in a possibly cholesterol-independent  
401 manner (75).

402

## 403 Discussion

404 Proteogenomic approaches have the potential to establish a direct link from rare and common  
405 variation in or closeby protein-encoding genes to human health via the protein product (5-18).  
406 Despite recent advances and early successes, the field is still in its infancy with respect to the  
407 scale and protein capture, with existing broad-capture technologies currently targeting less than  
408 a third of all proteins encoded in the human genome (5-18), not capturing posttranslational  
409 modifications, or providing absolute protein quantification.

410 Here we identified more than 200 novel cis-pQTLs that have not been reported so far, even in  
411 studies 30-times larger than ours using non-antibody-based technologies, by capitalizing on  
412 recent assay developments. The fact that we identified hundreds of cis-pQTLs for proteins that  
413 have been investigated in studies much larger than ours might be best explained by the need to  
414 develop further orthogonal methods to measure protein targets as we have outlined previously  
415 (14).

416 We demonstrate that systematic application of cis-pQTLs to large-scale genetic studies of human  
417 diseases can 1) guide causal gene annotation at GWAS loci (e.g., *DKKL1* for multiple sclerosis), 2)  
418 identify pathways that link genes to diseases guided by a protein-phenotype network, and 3)  
419 complement gene-burden testing of rare variants to discover novel biology. We highlight specific  
420 examples in more detail and share a large number of high-confidence protein-phenotype  
421 associations that provide a direct guide for functional follow-up and future investigations of  
422 variant protein with disease relevance about which little is known to date.

423 The vast majority (~90%) of genetic variants identified in GWASs reside in non-coding regions of  
424 the genome (76), creating a challenge for variant-to-function annotation. In line with previous  
425 studies, we demonstrate the efficiency and ability of cis-pQTLs to prioritize causal candidate  
426 genes including reassignments at 40% of overlapping loci. In contrast to other annotation  
427 approaches, the particular value of the integration of proteogenomic studies lies in the  
428 instrumentalization of the likely biological effector molecules. Studies using proteomic profiling  
429 in disease relevant tissues or single cell-types are needed to further elucidate the mechanisms  
430 underlying the many thousands of unassigned GWAS loci.

431 This study is a powerful demonstrating that even moderately sized proteomic studies can result  
432 in the identification of novel biology when combined with bespoke analysis pipelines designed  
433 for the identification of cis-pQTL and systematic integration and follow-up of disease GWAS  
434 summary statistics. Eventually multiple different technologies will be needed at scale to capture  
435 not only proteins of interest but also the vast spectrum of proteoforms with possible distinct  
436 phenotypic consequences (14). This prediction is supported by power calculations of the UKB-  
437 PPP (17), but also based on the observation that our study identified cis-pQTLs for genes that are  
438 under less evolutionary constraint as indicated by higher observed/expected scores for missense  
439 (+0.15; p-value= $5.0 \times 10^{-50}$ ) and loss-of-function (+0.22; p-value= $7.5 \times 10^{-51}$ ) variation in gnomAD  
440 (77). This observation is in line with recent findings among eQTL studies (78).

441 We observed convergence of gene–phenotype associations between ExWAS and our  
442 proteogenomic approach only at a small number of genes. Gene identification with overlapping  
443 or converging evidence, as shown for *TIMD4*, provides high confidence about the underlying  
444 causal gene, while the incomplete overlap clearly indicates the complementary nature of both  
445 approaches for drug target prioritization. An important distinction between both approaches,  
446 beyond the different genetic variants covered, is the ability of proteogenomics to emulate  
447 protein variation across the whole spectrum of abundance and in some cases function, and not  
448 only putative loss-of-function (rarely gain of function), which might explain differences seen in  
449 phenotypic consequences between both approaches. In addition, in terms of practicality,  
450 integration of pQTLs into colocalization and GWAS loci annotation enabled us to uncover  
451 unreported disease biology with a small sample size of 1,180 individuals, whereas substantially  
452 larger sample sizes, even millions of individuals, are needed to reach enough power to detect  
453 rare variant associations in ExWAS studies for disease endpoints (79).

454 Our study has some limitations that need to be considered. Affinity-based reagents allow for the  
455 quantification of protein abundance but are inherently limited to quantify the level of activity,  
456 although a general correspondence between the two can be assumed. This limits insights about  
457 the role of protein targets using a proteogenomic approach. Further, numerous posttranslational  
458 modifications can change the function and abundance of proteins but are currently not  
459 distinguishable using affinity reagents at scale. We deliberately decided to restrict genetic

460 analysis of protein targets to the corresponding protein coding regions ( $\pm 500\text{kb}$ ) for two reasons:  
461 1) the high biological prior to identify genetic variants directly linked to protein  
462 function/abundance, and 2) to increase power for statistical analysis by limiting the multiple  
463 testing burden. Larger studies are needed to explore the spectrum of trans-pQTLs that have  
464 generally smaller effect sizes but can identify protein interaction partners and facilitate  
465 systematic dose-response analysis in the two-sample Mendelian randomization frameworks.

466 In summary, we demonstrate the clear potential of broad-capture proteogenomic studies to  
467 identify novel biological pathways that link protein-encoding genes to human diseases.  
468 Systematic integration of human genomic with proteomic and phenomic data enables such  
469 investigation even in relatively moderately sized studies and can help to prioritize targets and  
470 indications for the development of safe and effective therapeutic interventions.

471

472



## 473 **Materials and Methods**

474

### 475 *Study participants*

476 We measured protein levels among 1,180 participants of the European Prospective Investigation  
477 into Cancer (EPIC)-Norfolk study, a cohort of 25,639 middle-aged, individuals from the general  
478 population of Norfolk, a county in Eastern England which is a component of EPIC (21). The study  
479 was approved by the Norfolk Research Ethics Committee (ref. 05/Q0101/191) and all participants  
480 gave their informed written consent before entering the study. Information on lifestyle factors  
481 and medical history was obtained from questionnaires as reported previously (21). We selected  
482 a random sub-cohort of 771 participants and a set of 429 participants with selected incident  
483 events during follow-up for cost-efficient proteomic profiling.

484

### 485 *Proteomic profiling*

486 We used serum samples from the baseline assessment (1993 - 1997) that had been stored in  
487 liquid nitrogen for proteomic profiling using the Olink Explore 1536 and Explore Expansion  
488 platforms targeting 2925 unique proteins by 2943 assays, of which 2923 unique proteins mapped  
489 to a protein encoding locus in genome assembly GRCh37. Details regarding the assay have been  
490 have been described in detail (80). Briefly, proteins are targeted by two separate unique  
491 antibodies, each of which are labelled with complementary single stranded oligonucleotides  
492 (proximity extension assays (81)). These proximity extension assays hybridization occurs  
493 subsequent to the binding of antibody pairs with complementary oligonucleotides which can be  
494 quantified using next generation sequencing (NGS). NGS read-outs undergo quality control  
495 procedures where internal (incubation, extension and amplification controls) and external  
496 (negative, plate and sample controls) controls are included. Normalized protein expression (NPX)  
497 units are generated by normalization to the extension control and further normalization to the  
498 plate control and reported on a log<sub>2</sub> scale. We excluded 3 samples as they were shown to be  
499 extreme outliers using principal component analysis from their entire proteomic profiles. For  
500 downstream genetic analysis (fine-mapping and region-based association analysis) we first rank-  
501 inverse normal transformed NPX-values and corrected for age, sex, and the first ten genetic



502 principal components using linear regression models. The residuals of this analysis were used  
503 throughout the study.

504

#### 505 *Genotyping*

506 EPIC-Norfolk samples (n=21,448) were genotyped on the Affymetrix UK Biobank Axiom array chip  
507 by Cambridge Genomic Services, University of Cambridge, UK. Sample and variant QC followed  
508 the Affymetrix Best Practices guidelines. Samples were excluded based on DishQC < 0.82  
509 (fluorescence signal contrast), call-rate <97%, heterozygosity outliers and sex discordance  
510 checks. Variants were excluded if call-rate <95% or  $HWE \leq 1e-6$ . Monomorphic variants and those  
511 with cluster problems detected using Affymetrix SNPolisher were excluded. Genotype imputation  
512 was performed using two different reference panels, the Haplotype Reference Consortium (HRC)  
513 (release 1) reference panel and the combined UK10K+1000 Genomes Phase 3 reference panel.  
514 After pre-imputation QC, 21,044 samples remained for imputation. All SNPs imputed using the  
515 HRC reference panel were included, and additional variants imputed using only the UK10K+1000  
516 Genomes reference panel were added to create a combined imputed set. Variants with  
517 imputation quality INFO < 0.4 or MAF of < < 0.0001 were excluded. All positions are on genome  
518 assembly GRCh37. After excluding ancestry outliers, individuals without a high-quality proteomic  
519 profile for each panel and pruning the sample set for related individuals, 1,180 and 1,178  
520 individuals were included in proteogenomic analyses for Olink Explore 1536 and Explore  
521 Expansion platforms, respectively.

522

#### 523 *Fine mapping*

524 We used statistical fine-mapping as implemented in the 'sum of single effects' model (SuSiE) (82)  
525 using individual level genotype and protein data to identify credible sets at protein encoding loci  
526 ( $\pm 500$ kb). Briefly, SuSiE employs a Bayesian framework for variable selection in a multiple  
527 regression problem with the aim to identify sets of independent variants each of which likely  
528 contain the true causally underlying genetic variant (82). We implemented the workflow using  
529 the R package *susieR* (v.0.11.92) and default prior and parameter settings. However, we noticed  
530 that SuSiE sometimes reports overlapping credible sets or credible sets that contained variants

531 in high LD with already selected ones. Therefore, we adopted a grid search by first iterating the  
532 maximum number of credible sets from 2 to 10 ( $L$  in SuSiE terminology) and subsequently  
533 selecting the output for the maximum  $L$  so that none of the credible sets reported variants in LD  
534 ( $r^2 > 0.1$ ). We further tested for independent effects of all lead credible set variants (selecting using  
535 highest posterior inclusion probability) by including them in a joint regression model. We only  
536 report credible sets which showed genome-wide significance ( $p < 5 \times 10^{-8}$ ) in those joint models.  
537 We used R v.4.1 to compute regression models.

538

### 539 *Region-based association testing*

540 To complement fine-mapping analysis, we computed regional association statistics at protein  
541 coding loci ( $\pm 500\text{kb}$ ) using fastGWA software provided by GCTA (v. 1.93.2beta) (83). To account  
542 for the different selection designs of the sub- and the case-cohort, we performed these analyses  
543 within each cohort separately and combined in an inverse-variance fixed-effects meta-analysis in  
544 METAL (84).

545

### 546 *Gene, Variant, and Protein annotation*

547 We obtained conservation scores for all protein coding genes from gnomAD. We used the Variant  
548 Effect Predictor software (22) (version 98.3) with the --pick option to annotate all independent  
549 lead variants and proxies ( $r^2 > 0.6$ ) of identified pQTLs in our data set and report possible  
550 functional consequences. We collapsed pQTLs mapping to the same functional variant to reduce  
551 redundancy. We further obtained protein characteristics, e.g., glycosylation sites, from UniProt  
552 (85).

553

### 554 *Annotation of GWAS catalog loci*

555 We downloaded genome-wide significant summary statistics from the GWAS catalog (date  
556 23/03/2022; (1)) and tested whether any of the lead credible set variants (cis-pQTLs) or proxies  
557 ( $r^2 > 0.8$ ) have been reported to be associated with any non-proteomic trait, that is omitting any  
558 results that related to multiplex proteomic assays. Out of 347,165 entries ( $n=9,997$  unique traits),  
559 212,628 entries ( $n=5,391$  unique traits) passed this and additional filtering steps (missing effect

560 estimates, missing risk allele, and not passing genome-wide significance). For each cis-pQTL –  
561 GWAS variant mapping, we compared the reported or mapped gene (closest gene assigned by  
562 the GWAS catalog) to the protein-encoding gene at the locus.

563

#### 564 *Phenome-wide analyses at protein-encoding loci*

565 We performed phenome-wide analyses using statistical colocalization for 914 protein targets  
566 where we had evidence for at least one cis-pQTL. To this end, we queried the Open GWAS  
567 database (86, 87) using a defined region ( $\pm 500$  kb) around the protein-encoding gene body and  
568 tested whether any of the traits in the databases showed a high posterior probability (PP) of  
569 shared genetic signal with plasma concentrations of the encoded protein target using statistical  
570 colocalization (88). We chose a cut-off of  $PP > 80\%$  to declare that a protein target and a  
571 phenotypic trait are highly likely to share a genetic signal at a locus. We used a conservative prior  
572 setting with  $p_{12} = 1 \times 10^{-6}$  and further ensured that regional sentinel variants were in strong LD  
573 ( $r^2 > 0.8$ ). To avoid spurious colocalization results due to imperfect overlap of SNPs, we filter all  
574 results for which the strongest cis-pQTL or sufficient proxy ( $r^2 > 0.8$ ) in the overlapping set was not  
575 included in the overlapping set of SNPs or if less than 500 SNPs were overlapping. We used the  
576 *igraph* package in R to visualize protein – disease colocalization results as a network to account  
577 for cross-disease dependencies established by proteins.

578

#### 579 *Incorporation of gene expression data*

580 We systematically tested for a shared genetic signal between plasma abundances of a protein  
581 and gene expression levels (eQTL) of the protein coding gene in 49 tissues from the GTEx project  
582 (v8) (89). We used a similar colocalization framework as described above but adopting a less  
583 stringent  $p_{12}$  prior ( $p_{12} = 1 \times 10^{-5}$ ) to account for the higher biological prior of genetic signals in the  
584 protein encoding region. All GTEx variant-gene cis-eQTL and cis-sQTL associations from each  
585 tissue were downloaded in January 2020 from  
586 <https://console.cloud.google.com/storage/browser/gtex-resources>.

587

588 *Phenotypic convergence between pQTL colocalization and rare loss of function gene-burden*  
589 *associations*

590 To compare the phenotypic convergence of rare loss of function gene-burden and cis-pQTLs  
591 colocalization results, we downloaded single variant and gene-burden results for 3,986  
592 phenotypic outcomes from UK Biobank respectively which were analysed by Backman et al.  
593 (2021) (downloaded on: 07/12/2021); (64). We filtered the results for 2,939 protein coding genes  
594 covered by the Olink Explore 1536 and Explore Expansion platforms. We compared the  
595 phenotypic convergence of genes that were significant for at least one phenotypic outcome in  
596 the exome-wide association analysis at exome-wide significance ( $p < 1 \times 10^{-6}$ ) with the pQTLs that  
597 showed significant statistical colocalization for at least one trait (PP > 80%). If ExWAS results were  
598 significant for more than one variant group for the same gene – trait association, we have filtered  
599 the results to only take forward the most significant finding.

600 *Multitrait colocalisation*

601 We used hypothesis prioritisation in multi-trait colocalisation (HyPrColoc) (90) at selected protein  
602 loci to identify a shared genetic signal across various traits, including gene expression, plasma  
603 protein levels, and prioritized phenotypes from the disease-wise colocalization framework.  
604 HyPrColoc provides for each cluster three different types of output: 1) a posterior probability (PP)  
605 that all phenotypes in the cluster share a common genetic signal, 2) a regional association  
606 probability, that it, that all the phenotypes share an association with one or more variants in the  
607 region, and 3) the proportion of the PP explained by the candidate variant. We considered a  
608 highly likely alignment of a genetic signal across various phenotypes if the PP > 80% and report  
609 obtained PPs otherwise.

## 610 **References**

- 611 1. Buniello A, MacArthur JAL, Cerezo M, Harris LW, Hayhurst J, Malangone C, et al. The  
612 NHGRI-EBI GWAS Catalog of published genome-wide association studies, targeted arrays and  
613 summary statistics 2019. *Nucleic Acids Res.* 2019;47(D1):D1005-D12.
- 614 2. Barbeira AN, Bonazzola R, Gamazon ER, Liang Y, Park Y, Kim-Hellmuth S, et al. Exploiting  
615 the GTEx resources to decipher the mechanisms at GWAS loci. *Genome Biol.* 2021;22(1):49.
- 616 3. Moore JE, Purcaro MJ, Pratt HE, Epstein CB, Shores N, Adrian J, et al. Expanded  
617 encyclopaedias of DNA elements in the human and mouse genomes. *Nature.*  
618 2020;583(7818):699-710.
- 619 4. Abell NS, DeGorter MK, Gloudemans MJ, Greenwald E, Smith KS, He Z, et al. Multiple  
620 causal variants underlie genetic associations in humans. *Science.* 2022;375(6586):1247-54.
- 621 5. Suhre K, Arnold M, Bhagwat AM, Cotton RJ, Engelke R, Raffler J, et al. Connecting genetic  
622 risk to disease end points through the human blood plasma proteome. *Nat Commun.*  
623 2017;8:14357.
- 624 6. Sun BB, Maranville JC, Peters JE, Stacey D, Staley JR, Blackshaw J, et al. Genomic atlas of  
625 the human plasma proteome. *Nature.* 2018;558(7708):73-9.
- 626 7. Folkersen L, Fauman E, Sabater-Lleal M, Strawbridge RJ, Frånberg M, Sennblad B, et al.  
627 Mapping of 79 loci for 83 plasma protein biomarkers in cardiovascular disease. *PLoS Genet.*  
628 2017;13(4):e1006706.
- 629 8. Yao C, Chen G, Song C, Keefe J, Mendelson M, Huan T, et al. Author Correction: Genome-  
630 wide mapping of plasma protein QTLs identifies putatively causal genes and pathways for  
631 cardiovascular disease. *Nat Commun.* 2018;9(1):3853.
- 632 9. Gilly A, Park YC, Png G, Barysenka A, Fischer I, Bjørnland T, et al. Whole-genome  
633 sequencing analysis of the cardiometabolic proteome. *Nat Commun.* 2020;11(1):6336.
- 634 10. Ferkingstad E, Sulem P, Atlason BA, Sveinbjornsson G, Magnusson MI, Styrismiddottir EL, et  
635 al. Large-scale integration of the plasma proteome with genetics and disease. *Nat Genet.*  
636 2021;53(12):1712-21.

- 637 11. Gudjonsson A, Gudmundsdottir V, Axelsson GT, Gudmundsson EF, Jonsson BG, Launer LJ,  
638 et al. A genome-wide association study of serum proteins reveals shared loci with common  
639 diseases. *Nat Commun.* 2022;13(1):480.
- 640 12. Katz DH, Tahir UA, Bick AG, Pampana A, Ngo D, Benson MD, et al. Whole Genome  
641 Sequence Analysis of the Plasma Proteome in Black Adults Provides Novel Insights Into  
642 Cardiovascular Disease. *Circulation.* 2022;145(5):357-70.
- 643 13. Png G, Barysenka A, Repetto L, Navarro P, Shen X, Pietzner M, et al. Mapping the serum  
644 proteome to neurological diseases using whole genome sequencing. *Nat Commun.*  
645 2021;12(1):7042.
- 646 14. Pietzner M, Wheeler E, Carrasco-Zanini J, Kerrison ND, Oerton E, Koprulu M, et al.  
647 Synergistic insights into human health from aptamer- and antibody-based proteomic profiling.  
648 *Nat Commun.* 2021;12(1):6822.
- 649 15. Pietzner M, Wheeler E, Carrasco-Zanini J, Cortes A, Koprulu M, Wörheide MA, et al.  
650 Mapping the proteo-genomic convergence of human diseases. *Science.*  
651 2021;374(6569):eabj1541.
- 652 16. Zhang J, Dutta D, Köttgen A, Tin A, Schlosser P, Grams ME, et al. Plasma proteome  
653 analyses in individuals of European and African ancestry identify cis-pQTLs and models for  
654 proteome-wide association studies. *Nat Genet.* 2022;54(5):593-602.
- 655 17. Sun BB, Chiou J, Traylor M, Benner C, Hsu Y-H, Richardson TG, et al. Genetic regulation of  
656 the human plasma proteome in 54,306 UK Biobank participants. *bioRxiv.*  
657 2022:2022.06.17.496443.
- 658 18. Folkersen L, Gustafsson S, Wang Q, Hansen DH, Hedman Å, Schork A, et al. Genomic and  
659 drug target evaluation of 90 cardiovascular proteins in 30,931 individuals. *Nat Metab.*  
660 2020;2(10):1135-48.
- 661 19. Enroth S, Johansson A, Enroth SB, Gyllensten U. Strong effects of genetic and lifestyle  
662 factors on biomarker variation and use of personalized cutoffs. *Nat Commun.* 2014;5:4684.
- 663 20. Wang G, Sarkar A, Carbonetto P, Stephens M. A simple new approach to variable selection  
664 in regression, with application to genetic fine mapping. *Journal of the Royal Statistical Society:*  
665 *Series B (Statistical Methodology).* 2020;82(5):1273-300.

- 666 21. Day N, Oakes S, Luben R, Khaw KT, Bingham S, Welch A, et al. EPIC-Norfolk: study design  
667 and characteristics of the cohort. *European Prospective Investigation of Cancer. Br J Cancer.*  
668 1999;80 Suppl 1:95-103.
- 669 22. McLaren W, Gil L, Hunt SE, Riat HS, Ritchie GR, Thormann A, et al. The Ensembl Variant  
670 Effect Predictor. *Genome Biol.* 2016;17(1):122.
- 671 23. Vujkovic M, Keaton JM, Lynch JA, Miller DR, Zhou J, Tcheandjieu C, et al. Discovery of 318  
672 new risk loci for type 2 diabetes and related vascular outcomes among 1.4 million participants in  
673 a multi-ancestry meta-analysis. *Nat Genet.* 2020;52(7):680-91.
- 674 24. Spracklen CN, Horikoshi M, Kim YJ, Lin K, Bragg F, Moon S, et al. Identification of type 2  
675 diabetes loci in 433,540 East Asian individuals. *Nature.* 2020;582(7811):240-5.
- 676 25. Mahajan A, Spracklen CN, Zhang W, Ng MCY, Petty LE, Kitajima H, et al. Multi-ancestry  
677 genetic study of type 2 diabetes highlights the power of diverse populations for discovery and  
678 translation. *Nat Genet.* 2022;54(5):560-72.
- 679 26. McDonald TJ, Nilsson G, Vagne M, Ghatei M, Bloom SR, Mutt V. A gastrin releasing peptide  
680 from the porcine nonantral gastric tissue. *Gut.* 1978;19(9):767-74.
- 681 27. McDonald TJ, Jörnvall H, Nilsson G, Vagne M, Ghatei M, Bloom SR, et al. Characterization  
682 of a gastrin releasing peptide from porcine non-antral gastric tissue. *Biochem Biophys Res*  
683 *Commun.* 1979;90(1):227-33.
- 684 28. Ladenheim EE, Taylor JE, Coy DH, Moore KA, Moran TH. Hindbrain GRP receptor blockade  
685 antagonizes feeding suppression by peripherally administered GRP. *Am J Physiol.* 1996;271(1 Pt  
686 2):R180-4.
- 687 29. Ladenheim EE, Hampton LL, Whitney AC, White WO, Battey JF, Moran TH. Disruptions in  
688 feeding and body weight control in gastrin-releasing peptide receptor deficient mice. *J*  
689 *Endocrinol.* 2002;174(2):273-81.
- 690 30. Persson K, Gingerich RL, Nayak S, Wada K, Wada E, Ahrén B. Reduced GLP-1 and insulin  
691 responses and glucose intolerance after gastric glucose in GRP receptor-deleted mice. *Am J*  
692 *Physiol Endocrinol Metab.* 2000;279(5):E956-62.

- 693 31. Gutzwiller JP, Drewe J, Hildebrand P, Rossi L, Lauper JZ, Beglinger C. Effect of intravenous  
694 human gastrin-releasing peptide on food intake in humans. *Gastroenterology*. 1994;106(5):1168-  
695 73.
- 696 32. Mhalhal TR, Washington MC, Newman KD, Heath JC, Sayegh AI. Combined gastrin  
697 releasing peptide-29 and glucagon like peptide-1 reduce body weight more than each individual  
698 peptide in diet-induced obese male rats. *Neuropeptides*. 2018;67:71-8.
- 699 33. Loh PR, Tucker G, Bulik-Sullivan BK, Vilhjálmsson BJ, Finucane HK, Salem RM, et al.  
700 Efficient Bayesian mixed-model analysis increases association power in large cohorts. *Nat Genet*.  
701 2015;47(3):284-90.
- 702 34. Gaziano JM, Concato J, Brophy M, Fiore L, Pyarajan S, Breeling J, et al. Million Veteran  
703 Program: A mega-biobank to study genetic influences on health and disease. *J Clin Epidemiol*.  
704 2016;70:214-23.
- 705 35. Pulit SL, Stoneman C, Morris AP, Wood AR, Glastonbury CA, Tyrrell J, et al. Meta-analysis  
706 of genome-wide association studies for body fat distribution in 694 649 individuals of European  
707 ancestry. *Hum Mol Genet*. 2019;28(1):166-74.
- 708 36. Frullanti E, Berking C, Harbeck N, Jézéquel P, Haugen A, Mawrin C, et al. Meta and pooled  
709 analyses of FGFR4 Gly388Arg polymorphism as a cancer prognostic factor. *Eur J Cancer Prev*.  
710 2011;20(4):340-7.
- 711 37. Chou CH, Hsieh MJ, Chuang CY, Lin JT, Yeh CM, Tseng PY, et al. Functional FGFR4  
712 Gly388Arg polymorphism contributes to oral squamous cell carcinoma susceptibility. *Oncotarget*.  
713 2017;8(56):96225-38.
- 714 38. Xiong SW, Ma J, Feng F, Fu W, Shu SR, Ma T, et al. Functional FGFR4 Gly388Arg  
715 polymorphism contributes to cancer susceptibility: Evidence from meta-analysis. *Oncotarget*.  
716 2017;8(15):25300-9.
- 717 39. Ulaganathan VK, Sperl B, Rapp UR, Ullrich A. Germline variant FGFR4 p.G388R exposes a  
718 membrane-proximal STAT3 binding site. *Nature*. 2015;528(7583):570-4.
- 719 40. Shin DJ, Osborne TF. FGF15/FGFR4 integrates growth factor signaling with hepatic bile  
720 acid metabolism and insulin action. *J Biol Chem*. 2009;284(17):11110-20.



- 721 41. Ge H, Zhang J, Gong Y, Gupte J, Ye J, Weizmann J, et al. Fibroblast growth factor receptor  
722 4 (FGFR4) deficiency improves insulin resistance and glucose metabolism under diet-induced  
723 obesity conditions. *J Biol Chem*. 2014;289(44):30470-80.
- 724 42. Wu X, Ge H, Lemon B, Weizmann J, Gupte J, Hawkins N, et al. Selective activation of  
725 FGFR4 by an FGF19 variant does not improve glucose metabolism in ob/ob mice. *Proc Natl Acad  
726 Sci U S A*. 2009;106(34):14379-84.
- 727 43. Huang X, Yang C, Luo Y, Jin C, Wang F, McKeenan WL. FGFR4 prevents hyperlipidemia and  
728 insulin resistance but underlies high-fat diet induced fatty liver. *Diabetes*. 2007;56(10):2501-10.
- 729 44. Lutz SZ, Hennige AM, Peter A, Kovarova M, Totsikas C, Machann J, et al. The  
730 Gly385(388)Arg Polymorphism of the FGFR4 Receptor Regulates Hepatic Lipogenesis Under  
731 Healthy Diet. *J Clin Endocrinol Metab*. 2019;104(6):2041-53.
- 732 45. Micha R, Khatibzadeh S, Shi P, Fahimi S, Lim S, Andrews KG, et al. Global, regional, and  
733 national consumption levels of dietary fats and oils in 1990 and 2010: a systematic analysis  
734 including 266 country-specific nutrition surveys. *BMJ*. 2014;348:g2272.
- 735 46. Lappalainen T, MacArthur DG. From variant to function in human disease genetics.  
736 *Science*. 2021;373(6562):1464-8.
- 737 47. Filippi M, Bar-Or A, Piehl F, Preziosa P, Solari A, Vukusic S, et al. Multiple sclerosis. *Nat Rev  
738 Dis Primers*. 2018;4(1):43.
- 739 48. Consortium IMSG. Multiple sclerosis genomic map implicates peripheral immune cells  
740 and microglia in susceptibility. *Science*. 2019;365(6460).
- 741 49. Yin X, Kim K, Suetsugu H, Bang SY, Wen L, Koido M, et al. Meta-analysis of 208370 East  
742 Asians identifies 113 susceptibility loci for systemic lupus erythematosus. *Ann Rheum Dis*.  
743 2021;80(5):632-40.
- 744 50. Kaneko KJ, Kohn MJ, Liu C, DePamphilis ML. The acrosomal protein Dickkopf-like 1 (DKKL1)  
745 is not essential for fertility. *Fertil Steril*. 2010;93(5):1526-32.
- 746 51. Uhlen M, Karlsson MJ, Zhong W, Tebani A, Pou C, Mikes J, et al. A genome-wide  
747 transcriptomic analysis of protein-coding genes in human blood cells. *Science*. 2019;366(6472).
- 748 52. Orrù V, Steri M, Sidore C, Marongiu M, Serra V, Olla S, et al. Complex genetic signatures  
749 in immune cells underlie autoimmunity and inform therapy. *Nat Genet*. 2020;52(10):1036-45.

- 750 53. Cencioni MT, Mattoscio M, Magliozzi R, Bar-Or A, Muraro PA. B cells in multiple sclerosis  
751 - from targeted depletion to immune reconstitution therapies. *Nat Rev Neurol*. 2021;17(7):399-  
752 414.
- 753 54. Granqvist M, Boremalm M, Poorghobad A, Svenningsson A, Salzer J, Frisell T, et al.  
754 Comparative Effectiveness of Rituximab and Other Initial Treatment Choices for Multiple  
755 Sclerosis. *JAMA Neurol*. 2018;75(3):320-7.
- 756 55. de Rojas I, Moreno-Grau S, Tesi N, Grenier-Boley B, Andrade V, Jansen IE, et al. Common  
757 variants in Alzheimer's disease and risk stratification by polygenic risk scores. *Nat Commun*.  
758 2021;12(1):3417.
- 759 56. Richardson TG, Sanderson E, Elsworth B, Tilling K, Davey Smith G. Use of genetic variation  
760 to separate the effects of early and later life adiposity on disease risk: mendelian randomisation  
761 study. *BMJ*. 2020;369:m1203.
- 762 57. Liu H, Leo C, Chen X, Wong BR, Williams LT, Lin H, et al. The mechanism of shared but  
763 distinct CSF-1R signaling by the non-homologous cytokines IL-34 and CSF-1. *Biochim Biophys*  
764 *Acta*. 2012;1824(7):938-45.
- 765 58. Delaney C, Farrell M, Doherty CP, Brennan K, O'Keeffe E, Greene C, et al. Attenuated CSF-  
766 1R signalling drives cerebrovascular pathology. *EMBO Mol Med*. 2021;13(2):e12889.
- 767 59. Morris JA, Kemp JP, Youlten SE, Laurent L, Logan JG, Chai RC, et al. An atlas of genetic  
768 influences on osteoporosis in humans and mice. *Nat Genet*. 2019;51(2):258-66.
- 769 60. Lotta LA, Wittemans LBL, Zuber V, Stewart ID, Sharp SJ, Luan J, et al. Association of Genetic  
770 Variants Related to Gluteofemoral vs Abdominal Fat Distribution With Type 2 Diabetes, Coronary  
771 Disease, and Cardiovascular Risk Factors. *JAMA*. 2018;320(24):2553-63.
- 772 61. Acosta-Herrera M, Kerick M, González-Serna D, Wijmenga C, Franke A, Gregersen PK, et  
773 al. Genome-wide meta-analysis reveals shared new *loci* in systemic seropositive rheumatic  
774 diseases. *Ann Rheum Dis*. 2019;78(3):311-9.
- 775 62. Robertson CC, Inshaw JRJ, Onengut-Gumuscu S, Chen WM, Santa Cruz DF, Yang H, et al.  
776 Fine-mapping, trans-ancestral and genomic analyses identify causal variants, cells, genes and  
777 drug targets for type 1 diabetes. *Nat Genet*. 2021;53(7):962-71.

- 778 63. Clarke LA, Giugliani R, Guffon N, Jones SA, Keenan HA, Munoz-Rojas MV, et al. Genotype-  
779 phenotype relationships in mucopolysaccharidosis type I (MPS I): Insights from the International  
780 MPS I Registry. *Clin Genet*. 2019;96(4):281-9.
- 781 64. Backman JD, Li AH, Marcketta A, Sun D, Mbatchou J, Kessler MD, et al. Exome sequencing  
782 and analysis of 454,787 UK Biobank participants. *Nature*. 2021;599(7886):628-34.
- 783 65. Szustakowski JD, Balasubramanian S, Kvikstad E, Khalid S, Bronson PG, Sasson A, et al.  
784 Advancing human genetics research and drug discovery through exome sequencing of the UK  
785 Biobank. *Nat Genet*. 2021;53(7):942-8.
- 786 66. Plenge RM, Scolnick EM, Altshuler D. Validating therapeutic targets through human  
787 genetics. *Nat Rev Drug Discov*. 2013;12(8):581-94.
- 788 67. Kathiresan S, Willer CJ, Peloso GM, Demissie S, Musunuru K, Schadt EE, et al. Common  
789 variants at 30 loci contribute to polygenic dyslipidemia. *Nat Genet*. 2009;41(1):56-65.
- 790 68. Lemke G. How macrophages deal with death. *Nat Rev Immunol*. 2019;19(9):539-49.
- 791 69. Miyanishi M, Tada K, Koike M, Uchiyama Y, Kitamura T, Nagata S. Identification of Tim4  
792 as a phosphatidylserine receptor. *Nature*. 2007;450(7168):435-9.
- 793 70. Kuchroo VK, Dardalhon V, Xiao S, Anderson AC. New roles for TIM family members in  
794 immune regulation. *Nat Rev Immunol*. 2008;8(8):577-80.
- 795 71. Mouse Genome Database (MGD) at the Mouse Genome Informatics website, The  
796 Jackson Laboratory 2022 [Available from: <http://www.informatics.jax.org/>].
- 797 72. Fazio S, Hasty AH, Carter KJ, Murray AB, Price JO, Linton MF. Leukocyte low density  
798 lipoprotein receptor (LDL-R) does not contribute to LDL clearance in vivo: bone marrow  
799 transplantation studies in the mouse. *J Lipid Res*. 1997;38(2):391-400.
- 800 73. Magalhaes MS, Smith P, Portman JR, Jackson-Jones LH, Bain CC, Ramachandran P, et al.  
801 Role of Tim4 in the regulation of ABCA1. *Nat Commun*. 2021;12(1):4434.
- 802 74. Koelwyn GJ, Corr EM, Erbay E, Moore KJ. Regulation of macrophage immunometabolism  
803 in atherosclerosis. *Nat Immunol*. 2018;19(6):526-37.
- 804 75. Foks AC, Engelbertsen D, Kuperwaser F, Alberts-Grill N, Gonen A, Witztum JL, et al.  
805 Blockade of Tim-1 and Tim-4 Enhances Atherosclerosis in Low-Density Lipoprotein Receptor-  
806 Deficient Mice. *Arterioscler Thromb Vasc Biol*. 2016;36(3):456-65.

- 807 76. Alsheikh AJ, Wollenhaupt S, King EA, Reeb J, Ghosh S, Stolzenburg LR, et al. The landscape  
808 of GWAS validation; systematic review identifying 309 validated non-coding variants across 130  
809 human diseases. *BMC Med Genomics*. 2022;15(1):74.
- 810 77. Gudmundsson S, Karczewski KJ, Francioli LC, Tiao G, Cummings BB, Alfoldi J, et al.  
811 Addendum: The mutational constraint spectrum quantified from variation in 141,456 humans.  
812 *Nature*. 2021;597(7874):E3-E4.
- 813 78. Mostafavi H, Spence JP, Naqvi S, Pritchard JK. Limited overlap of eQTLs and GWAS hits  
814 due to systematic differences in discovery. *bioRxiv*. 2022:2022.05.07.491045.
- 815 79. Akbari P, Gilani A, Sosina O, Kosmicki JA, Khrimian L, Fang YY, et al. Sequencing of 640,000  
816 exomes identifies *GPR75* variants associated with protection from obesity. *Science*.  
817 2021;373(6550).
- 818 80. Zhong W, Edfors F, Gummesson A, Bergström G, Fagerberg L, Uhlén M. Next generation  
819 plasma proteome profiling to monitor health and disease. *Nat Commun*. 2021;12(1):2493.
- 820 81. Assarsson E, Lundberg M, Holmquist G, Björkesten J, Thorsen SB, Ekman D, et al.  
821 Homogenous 96-plex PEA immunoassay exhibiting high sensitivity, specificity, and excellent  
822 scalability. *PLoS One*. 2014;9(4):e95192.
- 823 82. Wang G, Sarkar A, Carbonetto P, Stephens M. A simple new approach to variable  
824 selection in regression, with application to genetic fine mapping. *Royal Statistical Society*; 2020.  
825 p. 1273-300.
- 826 83. Jiang L, Zheng Z, Qi T, Kemper KE, Wray NR, Visscher PM, et al. A resource-efficient tool  
827 for mixed model association analysis of large-scale data. *Nat Genet*. 2019;51(12):1749-55.
- 828 84. Willer CJ, Li Y, Abecasis GR. METAL: fast and efficient meta-analysis of genomewide  
829 association scans. *Bioinformatics*. 2010;26(17):2190-1.
- 830 85. Consortium U. UniProt: the universal protein knowledgebase in 2021. *Nucleic Acids Res*.  
831 2021;49(D1):D480-D9.
- 832 86. Elsworth B, Lyon M, Alexander T, Liu Y, Matthews P, Hallett J, et al. The MRC IEU  
833 OpenGWAS data infrastructure. *bioRxiv*. 2020:2020.08.10.244293.
- 834 87. Hemani G, Zheng J, Elsworth B, Wade KH, Haberland V, Baird D, et al. The MR-Base  
835 platform supports systematic causal inference across the human phenome. *Elife*. 2018;7.

836 88. Giambartolomei C, Vukcevic D, Schadt EE, Franke L, Hingorani AD, Wallace C, et al.  
837 Bayesian test for colocalisation between pairs of genetic association studies using summary  
838 statistics. PLoS Genet. 2014;10(5):e1004383.

839 89. Consortium G. The GTEx Consortium atlas of genetic regulatory effects across human  
840 tissues. Science. 2020;369(6509):1318-30.

841 90. Foley CN, Staley JR, Breen PG, Sun BB, Kirk PDW, Burgess S, et al. A fast and efficient  
842 colocalization algorithm for identifying shared genetic risk factors across multiple traits. Nat  
843 Commun. 2021;12(1):764.

844

845

## 846 **Acknowledgements**

847 The EPIC-Norfolk study (DOI 10.22025/2019.10.105.00004) has received funding from the  
848 Medical Research Council (MR/N003284/1 MC-UU\_12015/1 and MC\_UU\_00006/1) and Cancer  
849 Research UK (C864/A14136). The genetics work in the EPIC-Norfolk study was funded by the  
850 Medical Research Council (MC\_PC\_13048). We are grateful to all the participants who have been  
851 part of the project and to the many members of the study teams at the University of Cambridge  
852 including the EPIC-Norfolk investigators, the Study Co-ordination team, the Epidemiology Field,  
853 Data and Laboratory teams who have enabled this research. Proteomics measurements were  
854 supported by a collaboration agreement between the University of Cambridge and Olink. We  
855 thank Philippa Pettingill, Ida Grundberg and Janet Kenyon for their support with quality control  
856 on the proteomic data. M.K. is supported by Gates Cambridge Trust. JCZS is supported by a 4-  
857 year Wellcome Trust PhD Studentship and the Cambridge Trust. C.L., E.W., M.P., N.K. and N.J.W.  
858 are funded by the Medical Research Council (MC\_UU\_00006/1). The authors thank Million  
859 Veteran Program (MVP) staff, researchers, and volunteers, who have contributed to MVP, and  
860 especially participants who previously served their country in the military and now generously  
861 agreed to enroll in the study (see <https://www.research.va.gov/mvp/> for more details). We thank  
862 Friedemann Paul, Aroon Hingorani and Siamon Gordon for sharing their expertise on disease-  
863 specific examples.

864

## 865 **Competing interests**

866 E.W. is now an employee of AstraZeneca.

867

## 868 **Author contributions**

869 MK, MP and CL designed the analysis and drafted the manuscript. MK, MP and EW have  
870 performed the bioinformatics analyses. JCZS and NK have performed the quality control and data  
871 preparation of the proteomic data. SL contributed to the interpretation and curation of disease

872 examples. NJW is PI of the EPIC-Norfolk study. All authors contributed to the interpretation of  
873 the results and critically reviewed the manuscript.

874

#### 875 **Data availability**

876 The EPIC-Norfolk data can be requested by bona fide researchers for specified scientific purposes  
877 via the study website (<https://www.mrc-epid.cam.ac.uk/research/studies/epic-norfolk/>). Data  
878 will either be shared through an institutional data sharing agreement or arrangements will be  
879 made for analyses to be conducted remotely without the need for data transfer.

880 Fine-mapped summary statistics for protein coding regions will be shared publicly following  
881 publication.

882 Genome-wide association studies for anthropometric phenotypes have been conducted using  
883 the UK Biobank resource (application no. 44448). Access to the UK Biobank genotype and  
884 phenotype data is open to all approved health researchers (<http://www.ukbiobank.ac.uk/>).

885