

1 **Coding and regulatory variants affect serum protein levels and common disease**

2

3

4

5 Valur Emilsson^{1,2,ξ*}, Valborg Gudmundsdottir^{1,ξ}, Alexander Gudjonsson^{1,ξ}, Mohd A
6 Karim^{3,4}, Marjan Ilkov¹, James R. Staley⁵, Elias F. Gudmundsson¹, Brynjolfur G. Jonsson¹,
7 Lenore J. Launer⁶, Jan H. Lindeman⁷, Nicholas M. Morton⁸, Thor Aspelund¹, John R. Lamb⁹,
8 Lori L. Jennings¹⁰ and Vilmundur Gudnason^{1,2,*}

9

10

11

12 ¹Icelandic Heart Association, Holtasmari 1, IS-201 Kopavogur, Iceland.

13 ²Faculty of Medicine, University of Iceland, 101 Reykjavik, Iceland

14 ³Wellcome Trust Sanger Institute, Wellcome Genome Campus, Hinxton, Cambridgeshire
15 CB10 1SA, UK.

16 ⁴Open Targets, Wellcome Genome Campus, Hinxton, Cambridgeshire CB10 1SD, UK

17 ⁵MRC/BHF Cardiovascular Epidemiology Unit, Department of Public Health and Primary
18 Care, University of Cambridge, Cambridge, UK

19 ⁶Laboratory of Epidemiology and Population Sciences, Intramural Research Program,
20 National Institute on Aging, Bethesda, MD 20892-9205, USA.

21 ⁷Department of General Surgery Leiden University Medical Center, Leiden. Holland

22 ⁸Centre for Cardiovascular Sciences, Queen's Medical Research Institute, University of
23 Edinburgh, Edinburgh EH16 4TJ, UK

24 ⁹GNF Novartis, 10675 John Jay Hopkins Drive, San Diego, CA 92121, USA.

25 ¹⁰Novartis Institutes for Biomedical Research, 22 Windsor Street, Cambridge, MA 02139,
26 USA.

27

28

29

30

31

32

33

34 ^ξThese authors contributed equally.

35 ^{*}Corresponding authors. Emails: valur@hjarta.is and v.gudnason@hjarta.is

36

37 **Abstract**

38 **Circulating proteins are prognostic for human outcomes including cancer, heart failure,**
39 **brain trauma and brain amyloid plaque burden. A deep serum proteome survey**
40 **recently revealed close associations of serum protein networks and common diseases.**
41 **The present study reveals unprecedented number of individual serum proteins that**
42 **overlap genetic signatures of diseases emanating from different tissues of the body.**
43 **Here, 54,469 low-frequency and common exome-array variants were compared with**
44 **4782 protein measurements in the serum of 5343 individuals of the deeply annotated**
45 **AGES Reykjavik cohort. Using a study-wide significant threshold, 2019 independent**
46 **exome array variants affecting levels of 2135 serum proteins were identified. These**
47 **variants overlapped genetic loci for hundreds of complex disease traits, emphasizing the**
48 **emerging role for serum proteins as biomarkers of and potential causative agents of**
49 **multiple diseases.**

50 Large-scale genome-wide association studies (GWASs) have expanded our knowledge of the
51 genetic basis of complex disease. As of 2018, approximately 5687 GWASs have been
52 published revealing 71,673 DNA variants to phenotype associations¹. Furthermore, exome-
53 wide genotyping arrays have linked rare and common variants to many complex traits. For
54 example, 444 independent risk variants were recently identified for lipoprotein fractions
55 across 250 genes². Despite the overall success of GWAS, the common lead SNPs rarely point
56 directly to a clear causative polymorphism, making determination of the underlying disease
57 mechanism difficult³⁻⁶. Regulatory variants affecting mRNA and/or protein levels and
58 structural variants like missense mutations can point directly to the causal candidate.
59 Alteration of the amino acid sequence may affect protein activity and/or influence
60 transcription, translation, stability, processing, and secretion of the protein in question⁷⁻⁹.

61 Thus, by integrating intermediate traits like mRNA and/or protein levels with genetics and
62 disease traits, the identification of the causal candidates can be enhanced³⁻⁶.

63 Proteins are arguably the ultimate players in all life processes in disease and health, however,
64 high throughput detection and quantification of proteins has been hampered by the limitations
65 of available proteomic technologies. Recently, a custom-designed Slow-Off rate Modified
66 Aptamer (SOMAmer) protein profiling platform was developed to measure 4782 proteins
67 encoded by 4137 human genes in the serum of 5457 individuals from the AGES Reykjavik
68 study (AGES-RS)¹⁰, resulting in 26.1 million individual protein measurements. Various
69 metrics related to the performance of the proteomic platform including aptamer specificity,
70 assay variability and reproducibility have already been described¹⁰. We demonstrated that the
71 human serum proteome is under strong genetic control¹⁰, in line with findings of others
72 applying identical or different proteomics technologies^{11,12}. Moreover, serum proteins were
73 found to exist in regulatory groups of network modules composed of members synthesized in
74 all tissues of the body, suggesting that system level coordination or homeostasis is mediated
75 to a significant degree by thousands of proteins in blood¹³. Importantly, the deep serum and
76 plasma proteome is associated with and prognostic for various diseases as well as human life
77 span^{10,14-20}.

78 Here, we regressed levels of 4782 proteins on 54,469 low-frequency and common variants
79 from the HumanExome BeadChip exome array, in sera from 5343 individuals of the deeply
80 phenotyped AGES-RS cohort. Further cross-referencing of all significant genotype-to-protein
81 associations to hundreds of genetic loci for various disease endpoints and clinical traits,
82 demonstrated profound overlap between the genetics of circulating proteins and disease
83 related phenotypes. We highlight how triangulation of data from different sources can link

84 genetics, protein levels and disease(s), with the intention of cross-validating one another and
85 point to potentially causal relationship between proteins and complex disease(s).

86 Using genotype data from an exome array (HumanExome BeadChip) enriched for structural
87 variants and tagged for many GWAS risk loci (**Methods**), the effect of low-frequency and
88 common variants on the deep serum proteome was examined. Quality control filters²¹, and
89 exclusion of monomorphic variants reduced the available variants to 76,891. Additionally, we
90 excluded variants at minor allele frequency (MAF) < 0.001 as they provide insufficient power
91 for single-point association analysis²². This resulted in 54,469 low-frequency (54%,
92 MAF<0.05) and common variants (46%, MAF ≥ 0.05) that were tested for association to each
93 of the 4782 human serum protein measurements using linear regression analysis adjusted for
94 the confounders age and sex (**Methods**). The current platform targets the serum proteome
95 arising largely from active or passive secretion, ectodomain shedding, lysis and/or cell
96 death^{10,23}. **Figure 1a** highlights the classification of the protein population targeted by the
97 aptamer-based profiling platform, showing over 68.7% of the proteins are secreted or single
98 pass transmembrane (SPTM) receptors.

99 Applying a Bonferroni corrected significance threshold of $P < 1.92 \times 10^{-10}$ ($0.05/54469/4782$)
100 we detected 5472 exome array variants that were associated with variable levels of 2135
101 serum proteins (**Supplementary Table 1** and **Fig. 1b**), of which 2019 variants are
102 independent (**Supplementary Table S2**). **Supplementary Table 1** lists all associations at P-
103 value < 1×10^{-6} , or 10,200 exome array variants affecting 3096 human proteins. These protein
104 quantitative trait loci (pQTLs) were *cis* and/or *trans* acting including several *trans* acting
105 hotspots with pleiotropic effects on multiple co-regulated proteins (**Fig. 1b**). Secreted proteins
106 were enriched for pQTLs (P-value < 0.0001) as compared to non-secreted proteins using
107 10,000 permutations to obtain the empirical distribution of the χ^2 test of equality of

108 proportions (**Supplementary Fig. S1**). This suggests that proteins bound for the systemic
109 environment are subject to more genetic regulation than other proteins identified by the
110 current platform. **Supplementary Table S3** summarizes various pathogenicity prediction
111 scores for all independent study-wide significant pQTLs in **Supplementary Table S2**, using
112 the Ensembl Variant Effect Predictor (VEP)^{24,25}. Next, we cross-referenced all the 5472
113 study-wide significant pQTLs with a comprehensive collection of genetic loci associated with
114 diseases and clinical traits from the curated PhenoScanner database²⁶, revealing that 60% of
115 all pQTLs were linked to at least one disease-related trait (**Supplementary Table S4**). We
116 have shown in our previous studies that genetic loci affecting several serum proteins exhibit
117 pleiotropy in relation to complex diseases¹⁰. An example of a possible pleiotropic effect
118 mediated by the variant rs2251219 within the gene *PBRM1* affecting multiple proteins and
119 sharing genetics with various diseases and clinical features is illustrated in **Fig. 2**.
120 **Supplementary Fig. S2** depicts the relationship between all proteins and some quantitative
121 traits associated with rs2251219. **Table 1** highlights a selected set of pQTLs that share
122 genetics with diseases of different etiologies including disorders of the brain, metabolism,
123 immune and cardiovascular system and cancer. In the sections that follow, we give examples
124 of serum pQTLs that overlap disease risk loci and demonstrate how different data sources can
125 cross-validate one another. Although data triangulation can be used to infer directional
126 consistency, it cannot tell whether the relationship is causal or reactive to a given outcome. As
127 a result, we used two-sample Mendelian randomization analysis (MR) on highlighted
128 examples to test support for a protein's causality to an outcome.

129 Variable levels of the anti-inflammatory protein TREM2 were associated with two distinct
130 genomic regions (**Fig. 3a** and **Supplementary Fig. S3**). This included the missense variant
131 rs75932628 (NP_061838.1: p.R47H) in *TREM2* at chromosome 6 (**Fig. 3b**), known to confer
132 a strong risk of late-onset Alzheimer's disease (LOAD)²⁷. The variant was also associated

133 with IGF1 (P = 3×10^{-18}) in serum (**Supplementary Table 1**), a protein recently
134 implicated in axonal growth²⁸. Intriguingly, the region at chromosome 11 associated with
135 soluble TREM2 levels harbors variants adjacent to the genes *MS4A4A* and *MS4A6A* including
136 rs610932 known to influence genetic susceptibility for LOAD²⁹ (**Table 1** and **Fig. 3a, b**).
137 The variant rs610932 was also associated with the proteins GLTPD2 and A4GALT
138 (**Supplementary Table 1**). The alleles increasing risk of LOAD for both the common variant
139 rs610932 and the low-frequency variant rs75932628 were associated with low levels of
140 soluble TREM2 (**Fig. 3b**). Consistently, we find that the high-risk allele for rs75932628 was
141 associated with accelerated mortality post incident LOAD in the AGES-RS (**Fig. 3c**). It is of
142 note that the levels of TREM2 in the cerebrospinal fluid (CSF) reflect the activity of brain
143 TREM2-triggered microglia^{4,30}, while high levels of CSF TREM2 have been associated with
144 improved cognitive functioning³¹. **Supplementary Fig. S4** highlights the correlation
145 (Spearman rank) between the different proteins affected by the LOAD risk loci at
146 chromosomes 6 and 11. The accumulated data show a directionally consistent effect at
147 independent risk loci for LOAD converging on the same causal candidate TREM2.
148 Furthermore, a two-sample MR analysis using genetic instruments across the *TREM2* and
149 *MS4A4A/MS4A6A* loci and GWAS associations for LOAD in Europeans as outcome³²,
150 provided evidence that variable TREM2 protein levels are causally related to LOAD (P =
151 7.6×10^{-5}) (**Fig 3d**). In summary, these results demonstrate that the effect of genetic drivers on
152 major brain-linked disease like LOAD can be readily detected in serum to both inform on the
153 causal relationship and the directionality of the risk mediating effect. This would also suggest
154 that serum may be an accessible proxy for microglia function and cognition.

155 Variable levels of the cell adhesion protein SVEP1 are associated with variants located at
156 chromosomes 1 and 9 (**Supplementary Table 1, Fig. 4a** and **Supplementary Fig. S5**).
157 Genetic associations to SVEP1 levels at chromosome 9 include the low-frequency missense

158 variant rs111245230 in *SVEP1* (NP_699197.3: pD2702G) (**Fig. 4b**), which was recently
159 linked to coronary heart disease (CHD), blood pressure and type-2-diabetes (T2D)³³. Overall,
160 we found eight different missense mutations in *SVEP1* that were associated with *SVEP1*
161 serum levels (**Supplementary Table 1**). The CHD and T2D risk allele (C) of rs111245230
162 was associated with elevated levels of *SVEP1*, and *SVEP1* levels were consistently elevated
163 in CHD and T2D patients (**Fig. 4c**). Furthermore, high *SVEP1* levels were positively
164 correlated with systolic blood pressure ($\beta = 2.10$, $P = 4 \times 10^{-12}$) (**Fig. 4c**), but not with diastolic
165 blood pressure ($\beta = 0.115$, $P = 0.413$). Consistently, higher serum levels of *SVEP1* were
166 associated with increased mortality post-incident CHD in the AGES-RS (HR = 1.27, $P =$
167 9×10^{-9}) (**Fig. 4d**). The variants at chromosome 1 linked to *SVEP1* levels (**Fig. 4a**), have not
168 previously been linked to any disease. Given the currently available GWAS summary
169 statistics, a two-sample MR analysis using *cis*-variants on chromosome 9 for *SVEP1* as
170 instruments and a GWAS associations for T2D³⁴ support a causal relationship of *SVEP1* with
171 T2D ($P = 1.2 \times 10^{-5}$) (**Fig. 4e**), but not with CHD³⁵ or systolic blood pressure³⁶ ($P > 0.05$). Our
172 data triangulation and causal tests integrating genetics, serum protein levels and disease(s),
173 indicate that *SVEP1* may be a therapeutic target for T2D.

174 The ILMN exome array contains several tags related to previous GWAS findings³⁷, including
175 many risk loci for cancer. For example, 21 loci associated with melanoma³⁸ and 50 loci
176 associated with colorectal cancer³⁹. The exome array variant rs910873 located in an intron of
177 the GPI transamidase gene *PIGU* was previously linked to melanoma risk⁴⁰. The reported
178 candidate gene *PIGU* is the gene most proximal to the lead SNP rs910873 and may be a novel
179 candidate gene involved in melanoma. However, a more biologically relevant candidate is the
180 agouti-signaling protein (*ASIP*) gene that is located 314kb downstream of the lead SNP
181 rs910873. *ASIP* is a competitive inhibitor of MC1R⁴¹, and is thus strongly biologically
182 implicated in melanoma risk⁴². We found that the melanoma risk allele for rs910873 was

183 associated with elevated ASIP serum levels ($P = 3 \times 10^{-179}$) and the variant had no effect on
184 other proteins measured with the current proteomic platform (**Fig. 5a, Supplementary Table**
185 **1** and **Table 1**). Interestingly, the pQTL rs910873 is also an eQTL for *ASIP* gene expression
186 in skin⁴³, showing directionally consistent effect on the mRNA and protein. Importantly, we
187 found that serum ASIP levels were supported as causally related to malignant melanoma ($P =$
188 4.8×10^{-26}) using a two-sample MR analysis on the protein-to-outcome causal sequence of
189 events (**Fig. 5b**). Our data point to the *ASIP* protein underlying the risk at rs910873, thus
190 providing supportive evidence for the hypothesis that *ASIP* mediated inhibition of *MC1R*
191 results in suppression of melanogenesis and increased risk of melanoma⁴⁴. An additional
192 example is the susceptibility variant rs1800469 for colorectal cancer⁴⁵, which is a proxy to the
193 pQTL rs2241714 ($r^2=0.978$) (**Table 1** and **Fig. 5b**). While the *TMEM91* gene was the
194 reported candidate gene for the colorectal cancer risk at the rs1800469 (**Table 1**), we find that
195 the risk variant affected three proteins in either *cis* (*B3GNT8* and *TGFB1*) or *trans* (*B3GNT2*)
196 (**Fig. 5b**). Intriguingly, all three proteins have previously been implicated in colorectal
197 cancer⁴⁶⁻⁴⁸. Due to a lack of available and powered GWAS summary statistics data, we were
198 unable to formally test the causality of these proteins to colorectal cancer. In conclusion,
199 while we cannot rule out *PIGU* as a candidate gene for malignant melanoma, these findings
200 point to an alternate, and possibly more biologically relevant, candidate, *ASIP*.

201 We outlined the construction of the serum protein network in our previous report and
202 identified common genetic variants underlying the network structure¹⁰. This included a
203 targeted study of the effects of common *cis* and *cis-to-trans* acting variants on levels of serum
204 proteins. The comparison between that study and the current one using all independent study-
205 wide significant associations (**Supplementary Table S2**) and linkage disequilibrium (LD)
206 threshold of $r^2 > 0.50$ for known associations, shows that 77.2% of the current study's variant-to-
207 protein associations are novel. Importantly, while 70% of the variants detected with the

208 exome array are exonic and 59% of mapped pQTLs in the current study are exonic, only 7%
209 of the identified pQTLs were exonic in our earlier report¹⁰. Previously, we discovered that
210 80% of *cis* pQTL effects and 74% of *trans* pQTL effects were replicated across populations
211 and proteomics platforms measuring common variants¹⁰. Given that the exome array platform
212 is enriched for rare and low-frequency variants, a comparable test of replication is not
213 straightforward. Examining the proteins and variants measured across studies, we find that
214 76.0% of SNP-to-protein associations are novel in the present study when compared to, say,
215 Sun et al.¹¹, and 60.1% are novel when compared to the majority of studies published to date
216 (**Supplementary Table S5**), for all independent associations in the current study and LD of
217 $r^2 < 0.5$ between study specific markers.

218 We report here that many of the measured serum proteins under genetic control share genetics
219 with a variety of clinical features, including major diseases arising from various body tissues.
220 This is in line with a recent population-scale survey of human induced pluripotent stem cells,
221 demonstrating that pQTLs are 1.93-fold enriched in disease risk variants compared to a 1.36-
222 fold enrichment for eQTLs¹², underscoring the added value in pQTL mapping. We reaffirm
223 widespread associations between genetic variants and their cognate proteins as well as distant
224 *trans*-acting effects on serum proteins and demonstrate that many proteins are often involved
225 in mediating the biological effect of a single causal variant affecting complex disease. Protein
226 coding variants may cause technical artifacts in both affinity proteomics and mass
227 spectrometry^{49,50}. Systematic conditional and colocalization studies have shown, however,
228 that pQTLs powered by common missense variants being artifactual are not a common event
229 using the aptamer-based technology^{11,51}, however, given the enrichment of missense variants
230 in the present study, it may occur in some cases.

231 We note that with the ever-increasing availability of large-scale omics data aligned with the
232 human genome, cross-referencing different datasets can result in findings that occurred by
233 sheer chance. Hence, a systematic colocalization analysis has been proposed for identifying
234 shared causal variants between intermediate traits and disease endpoints⁵². This is, however,
235 not feasible for application of the exome array given its sparse genomic coverage. Instead,
236 multi-omics data triangulation to infer consistency in directionality, the approach used in the
237 present study, can enhance confidence in the causal call and offer insights and guidelines for
238 experimental follow-up studies. In fact, the causal calls for TREM2 (LOAD), SVEP1 (T2D)
239 and ASIP (melanoma) were validated, using a two-sample MR analysis. We previously
240 asserted that serum proteins are intimately connected to and may mediate global
241 homeostasis¹⁰. The accumulated data show that serum proteins are under strong genetic
242 control and closely associated with diseases of different aetiologies, which in turn suggests
243 that serum proteins may be significant mediators of systemic homeostasis in human health
244 and disease.

245 **METHODS**

246 **Study population**

247 Participants aged 66 through 96 are from the Age, Gene/Environment Susceptibility
248 Reykjavik Study (AGES-RS) cohort⁵³. AGES-RS is a single-center prospective population-
249 based study of deeply phenotyped subjects (5764, mean age 75±6 years) and survivors of the
250 40-year-long prospective Reykjavik study (n~18,000), an epidemiologic study aimed to
251 understand aging in the context of gene/environment interaction by focusing on four biologic
252 systems: vascular, neurocognitive (including sensory), musculoskeletal, and body
253 composition/metabolism. Descriptive statistics of this cohort as well as detailed definition of
254 the various disease endpoints and relevant phenotypes measured have been published^{10,53}. The

255 AGES-RS was approved by the NBC in Iceland (approval number VSN-00-063), and by the
256 National Institute on Aging Intramural Institutional Review Board, and the Data Protection
257 Authority in Iceland.

258 **Genotyping platform**

259 Genotyping was conducted using the exome-wide genotyping array Illumina HumanExome-
260 24 v1.1 Beadchip from Illumina (San Diego, CA, USA) for all AGES-RS participants as
261 previously described⁵⁴. The exome array was enriched for exonic variants selected from over
262 12,000 individual exome and whole-genome sequences from different study populations³⁷,
263 and includes as well tags for previously described GWAS hits, ancestry informative markers,
264 mitochondrial SNPs and human leukocyte antigen tags³⁷. A total of 244,883 variants were
265 included on the exome array. Genotype call and quality control filters including call rate,
266 heterozygosity, sex discordance and PCA outliers were performed as previously described^{2,21}.
267 Variants with call rate <90% or with Hardy–Weinberg P values $<1 \times 10^{-7}$ were removed from
268 the study. 72,766 variants were detected in at least one individual of the AGES-RS cohort. Of
269 these variants, 54,469 had a minor allele frequency > 0.001 and were examined for
270 association against each of the 4782 human serum protein measurements (see below).

271 **Protein measurements**

272 Each protein has its own detection reagent selected from chemically modified DNA libraries,
273 referred to as Slow Off-rate Modified Aptamers (SOMAmers)⁵⁵. The design and quality
274 control of the SOMApanel platform's custom version to include proteins known or predicted
275 to be present in the extracellular milieu have been described in detail elsewhere¹⁰. Briefly
276 though, the aptamer-based platform measures 5034 protein analytes in a single serum sample,
277 of which 4782 SOMAmers bind specifically to 4137 human proteins (some proteins are
278 identified by more than one aptamer) and 250 SOMAmers that recognize non-human targets
279 (47 non-human vertebrate proteins and 203 targeting human pathogens)¹⁰. Consistent target

280 specificity across the platform was indicated by direct (through mass spectrometry) and/or
281 indirect validation of the SOMAmers¹⁰. Both sample selection and sample processing for
282 protein measurements were randomized, and all samples were run as a single set to prevent
283 batch or time of processing biases.

284 **Statistical analysis**

285 Prior to the analysis of the proteins measurements, we applied a Box-Cox transformation on
286 all proteins to improve normality, symmetry and to maintain all protein variables on a similar
287 scale⁵⁶. In the association analysis, we obtained residuals after controlling for sex, age,
288 potential population stratification using principal component (PCs) analysis⁵⁷, and for all
289 single-variant associations to serum proteins tested under an additive genetic model applying
290 linear regression analysis (protein ~ SNP + age + sex + PC1 + PC2 + ...PC5). We report both
291 variants to protein associations at $P < 1 \times 10^{-6}$ for suggestive evidence and Bonferroni
292 correction for multiple comparisons by adjusting for the 54,469 variants and 4782 human
293 protein analytes where single variant associations with $P < 1.9 \times 10^{-10}$ were considered study-
294 wide significant (**Supplementary Table S1**). P-values corresponding to the estimated effect
295 size and standard errors of the genotypes, were recalculated to increase accuracy. Independent
296 genetic signals were found through a stepwise conditional and joint association analysis for
297 each protein analyte separately with the GCTA-COJO software^{58,59}. We conditioned on the
298 current lead variant listed in **Supplementary Table S1**, defined as the variant with the lowest
299 P-value, and then kept track of any new variants that were not in LD (the default GCTA-
300 COJO option $r^2 < 0.9$ for colinearity) with previously chosen lead variants and reported
301 findings at P-value $< 1 \times 10^{-6}$ (**Supplementary Table S2**). In the joint model all conditionally
302 significant SNPs for each protein analyte were combined in the regression model.

303 **Supplementary Table S3** summarizes, through use of VEP^{24,25}, various pathogenicity
304 prediction scores for all independent study-wide significant pQTLs in **Supplementary Table**
305 **S2**, including the Likelihood Ratio Test (LRT)⁶⁰, Variant Effect Scoring Tool (VEST)⁶¹,
306 MutationAssessor⁶² and MutationTaster⁶³. To test whether the percentage of secreted proteins
307 among pQTLs is equal to the percentage of secreted proteins among non-pQTLs, 10,000
308 permutations were performed to obtain the empirical distribution of the χ^2 test of equality of
309 proportions. Our null and alternate hypotheses were:

310 $H_0: P(\text{pQTL} | \text{Secreted}) = P(\text{pQTL} | \text{Not Secreted})$ and $H_1: P(\text{pQTL} | \text{Secreted}) > P(\text{pQTL} | \text{Not Secreted})$

311 The test statistics calculated from our data was compared to the quantiles of this distribution
312 to obtain $P(\text{Data}|H_0)$ (**Supplementary Fig. S1**).

313 We applied the “TwoSampleMR” R package⁶⁴ to perform a two-sample MR analysis to test
314 for causal associations between protein and outcome (protein-to-outcome). For different
315 outcomes we used GWAS associations for LOAD in Europeans³², malignant melanoma in
316 European individuals from the UK biobank data (UKB-b-12915)⁶⁵ and T2D in Europeans³⁴.
317 Genetic variants (SNPs) associated with serum protein levels at a genome-wide significant
318 threshold ($P < 5 \times 10^{-8}$) identified in the AGES dataset and filtered to only include uncorrelated
319 variants ($r^2 < 0.2$) were used as instruments. The inverse variance weighted (IVW) method⁶⁶
320 was used for the MR analysis, with P-values < 0.05 considered significant.

321 For the associations of individual proteins to different phenotypic measures we used linear or
322 logistic regression or Cox proportional hazards regression, depending on the outcome being
323 continuous, binary or a time to an event. Given consistency in terms of sample handling
324 including time from blood draw to processing (between 9-11 am), same personnel handling
325 all specimens and the ethnic homogeneity of the population we adjusted only for age and sex
326 in all our regression analyses. All statistical analysis was performed using R version 3.6.0 (R
327 Foundation for Statistical Computing, Vienna, Austria).

328 We compared our pQTL results to 19 previously published proteogenomic studies
329 (**Supplementary Table 5**), including the protein GWAS in the INTERVAL study¹¹, and our
330 previously reported genetic analysis of 3,219 AGES cohort participants¹⁰. In the previous
331 proteogenomic analysis of AGES participants, one *cis* variant was reported per protein using a
332 locus-wide significance threshold, as well as *cis-to-trans* variants at a Bonferroni corrected
333 significance threshold. Due to these differences in reporting criteria, we only considered the
334 associations in previous AGES results that met the current study-wide P-value threshold. For
335 all other studies we retained the pQTLs at the reported significance threshold. In addition, we
336 performed a lookup of all independent pQTLs from the current study available in summary
337 statistics from the INTERVAL study, considering them known if they reached a study-wide
338 significance in their data. We calculated the LD structure between the reported significant
339 variants for all studies, using 1000 Genomes v3 EUR samples, but using AGES data when
340 comparing to previously reported AGES results. We considered variants in LD at $r^2 > 0.5$ to
341 represent the same signal across studies. Comparison was performed on protein level, by
342 matching the reported Entrez gene symbol from each study.

343 **Acknowledgements**

344 V.E. and Va.G. are supported by the Icelandic Research Fund (IRF grants 195761-051 and
345 184845-053). The Age, Gene/Environment Susceptibility-Reykjavik Study (AGES-RS) was
346 supported by NIH contracts N01-AG-1-2100 and HHSN27120120022C, the NIA Intramural
347 Research Program, Hjartavernd (the Icelandic Heart Association), and the Althingi (the
348 Icelandic Parliament). M.A.K. was funded by Open Targets and by the Wellcome Trust Grant
349 206194.

350

351

352 **References**

- 353 1 Buniello, A. *et al.* The NHGRI-EBI GWAS Catalog of published genome-wide
354 association studies, targeted arrays and summary statistics 2019. *Nucleic Acids*
355 *Research* **47**, D1005-D1012, doi:10.1093/nar/gky1120 (2019).
- 356 2 Liu, D. J. *et al.* Exome-wide association study of plasma lipids in >300,000
357 individuals. *Nature Genetics* **49**, 1758-1766, doi:10.1038/ng.3977 (2017).
- 358 3 Schadt, E. E. Molecular networks as sensors and drivers of common human diseases.
359 *Nature* **461**, 218-223, doi:10.1038/nature08454 (2009).
- 360 4 Zhang, B. *et al.* Integrated Systems Approach Identifies Genetic Nodes and Networks
361 in Late-Onset Alzheimer's Disease. *Cell* **153**, 707-720, doi:10.1016/j.cell.2013.03.030
362 (2013).
- 363 5 Emilsson, V. *et al.* Genetics of gene expression and its effect on disease. *Nature* **452**,
364 423-U422, doi:10.1038/nature06758 (2008).
- 365 6 Chen, Y. Q. *et al.* Variations in DNA elucidate molecular networks that cause disease.
366 *Nature* **452**, 429-435, doi:10.1038/nature06757 (2008).
- 367 7 Pires, D. E., Chen, J., Blundell, T. L. & Ascher, D. B. In silico functional dissection of
368 saturation mutagenesis: Interpreting the relationship between phenotypes and changes
369 in protein stability, interactions and activity. *Sci Rep* **6**, 19848, doi:10.1038/srep19848
370 (2016).
- 371 8 Ho, J. E. *et al.* Common genetic variation at the IL1RL1 locus regulates IL-33/ST2
372 signaling. *Journal of Clinical Investigation* **123**, 4208-4218, doi:10.1172/JCI67119
373 (2013).
- 374 9 Interleukin-6 Receptor Mendelian Randomisation Analysis, C. *et al.* The interleukin-6
375 receptor as a target for prevention of coronary heart disease: a mendelian
376 randomisation analysis. *Lancet* **379**, 1214-1224, doi:10.1016/s0140-6736(12)60110-x
377 (2012).
- 378 10 Emilsson, V. *et al.* Co-regulatory networks of human serum proteins link genetics to
379 disease. *Science* **361**, 769-773, doi:10.1126/science.aaq1327 (2018).
- 380 11 Sun, B. B. *et al.* Genomic atlas of the human plasma proteome. *Nature* **558**, 73-79,
381 doi:10.1038/s41586-018-0175-2 (2018).
- 382 12 Mirauta, B. A. *et al.* Population-scale proteome variation in human induced
383 pluripotent stem cells. *Elife* **9**, doi:10.7554/eLife.57390 (2020).
- 384 13 Lamb, J. R., Jennings, L. L., Gudmundsdottir, V., Gudnason, V. & Emilsson, V. It's in
385 Our Blood: A Glimpse of Personalized Medicine. *Trends Mol Med*,
386 doi:10.1016/j.molmed.2020.09.003 (2020).
- 387 14 Emilsson, V., Gudnason, V. & Jennings, L. L. Predicting health and life span with the
388 deep plasma proteome. *Nat Med* **25**, 1815-1816, doi:10.1038/s41591-019-0677-y
389 (2019).
- 390 15 Lehallier, B. *et al.* Undulating changes in human plasma proteome profiles across the
391 lifespan. *Nat Med* **25**, 1843-1850, doi:10.1038/s41591-019-0673-2 (2019).
- 392 16 Williams, S. A. *et al.* Plasma protein patterns as comprehensive indicators of health.
393 *Nat Med* **25**, 1851-1857, doi:10.1038/s41591-019-0665-2 (2019).
- 394 17 Nakamura, A. *et al.* High performance plasma amyloid-beta biomarkers for
395 Alzheimer's disease. *Nature* **554**, 249-254, doi:10.1038/nature25456 (2018).
- 396 18 Dodgson, S. E. There Will Be Blood Tests. *Cell* **173**, 1-3,
397 doi:10.1016/j.cell.2018.03.012 (2018).
- 398 19 Cohen, J. D. *et al.* Detection and localization of surgically resectable cancers with a
399 multi-analyte blood test. *Science* **359**, 926-930, doi:10.1126/science.aar3247 (2018).

- 400 20 Kristensen, S. L. *et al.* Prognostic Value of N-Terminal Pro-B-Type Natriuretic
401 Peptide Levels in Heart Failure Patients With and Without Atrial Fibrillation. *Circ*
402 *Heart Fail* **10**, doi:10.1161/circheartfailure.117.004409 (2017).
- 403 21 Peloso, G. M. *et al.* Association of low-frequency and rare coding-sequence variants
404 with blood lipids and coronary heart disease in 56,000 whites and blacks. *American*
405 *Journal of Human Genetics* **94**, 223-232, doi:10.1016/j.ajhg.2014.01.009 (2014).
- 406 22 Richards, A. L. *et al.* Exome arrays capture polygenic rare variant contributions to
407 schizophrenia. *Human Molecular Genetics* **25**, 1001-1007, doi:10.1093/hmg/ddv620
408 (2016).
- 409 23 Armengaud, J., Christie-Oleza, J. A., Clair, G., Malard, V. & Duport, C.
410 Exoproteomics: exploring the world around biological systems. *Expert Rev*
411 *Proteomics* **9**, 561-575, doi:10.1586/epr.12.52 (2012).
- 412 24 McLaren, W. *et al.* The Ensembl Variant Effect Predictor. *Genome Biol* **17**, 122,
413 doi:10.1186/s13059-016-0974-4 (2016).
- 414 25 McLaren, W. *et al.* Deriving the consequences of genomic variants with the Ensembl
415 API and SNP Effect Predictor. *Bioinformatics* **26**, 2069-2070,
416 doi:10.1093/bioinformatics/btq330 (2010).
- 417 26 Staley, J. R. *et al.* PhenoScanner: a database of human genotype-phenotype
418 associations. *Bioinformatics* **32**, 3207-3209, doi:10.1093/bioinformatics/btw373
419 (2016).
- 420 27 Jonsson, T. *et al.* Variant of TREM2 associated with the risk of Alzheimer's disease. *N*
421 *Engl J Med* **368**, 107-116, doi:10.1056/NEJMoa1211103 (2013).
- 422 28 Guo, C. *et al.* IGFBPL1 Regulates Axon Growth through IGF-1-mediated Signaling
423 Cascades. *Sci Rep* **8**, 2054, doi:10.1038/s41598-018-20463-5 (2018).
- 424 29 Hollingworth, P. *et al.* Common variants at ABCA7, MS4A6A/MS4A4E, EPHA1,
425 CD33 and CD2AP are associated with Alzheimer's disease. *Nature Genetics* **43**, 429-
426 435, doi:10.1038/ng.803 (2011).
- 427 30 Suarez-Calvet, M. *et al.* sTREM2 cerebrospinal fluid levels are a potential biomarker
428 for microglia activity in early-stage Alzheimer's disease and associate with neuronal
429 injury markers. *EMBO Mol Med* **8**, 466-476, doi:10.15252/emmm.201506123 (2016).
- 430 31 Ewers, M. *et al.* Increased soluble TREM2 in cerebrospinal fluid is associated with
431 reduced cognitive and clinical decline in Alzheimer's disease. *Sci Transl Med* **11**
432 (2019).
- 433 32 Kunkle, B. W. *et al.* Genetic meta-analysis of diagnosed Alzheimer's disease identifies
434 new risk loci and implicates A β , tau, immunity and lipid processing. *Nat Genet* **51**,
435 414-430, doi:10.1038/s41588-019-0358-2 (2019).
- 436 33 Myocardial Infarction, G. *et al.* Coding Variation in ANGPTL4, LPL, and SVEP1 and
437 the Risk of Coronary Disease. *N Engl J Med* **374**, 1134-1144,
438 doi:10.1056/NEJMoa1507652 (2016).
- 439 34 Mahajan, A. *et al.* Fine-mapping type 2 diabetes loci to single-variant resolution using
440 high-density imputation and islet-specific epigenome maps. *Nat Genet* **50**, 1505-1513,
441 doi:10.1038/s41588-018-0241-6 (2018).
- 442 35 Nikpay, M. *et al.* A comprehensive 1,000 Genomes-based genome-wide association
443 meta-analysis of coronary artery disease. *Nat Genet* **47**, 1121-1130,
444 doi:10.1038/ng.3396 (2015).
- 445 36 Evangelou, E. *et al.* Genetic analysis of over 1 million people identifies 535 new loci
446 associated with blood pressure traits. *Nat Genet* **50**, 1412-1425, doi:10.1038/s41588-
447 018-0205-x (2018).

- 448 37 Huyghe, J. R. *et al.* Exome array analysis identifies new loci and low-frequency
449 variants influencing insulin processing and secretion. *Nature Genetics* **45**, 197-201,
450 doi:10.1038/ng.2507 (2013).
- 451 38 Ransohoff, K. J. *et al.* Two-stage genome-wide association study identifies a novel
452 susceptibility locus associated with melanoma. *Oncotarget* **8**, 17586-17592,
453 doi:10.18632/oncotarget.15230 (2017).
- 454 39 Lu, Y. *et al.* Large-Scale Genome-Wide Association Study of East Asians Identifies
455 Loci Associated With Risk for Colorectal Cancer. *Gastroenterology*,
456 doi:10.1053/j.gastro.2018.11.066 (2018).
- 457 40 Brown, K. M. *et al.* Common sequence variants on 20q11.22 confer melanoma
458 susceptibility. *Nature Genetics* **40**, 838-840, doi:10.1038/ng.163 (2008).
- 459 41 Blanchard, S. G. *et al.* Agouti antagonism of melanocortin binding and action in the
460 B16F10 murine melanoma cell line. *Biochemistry* **34**, 10406-10411 (1995).
- 461 42 Taylor, N. J. *et al.* Inherited variation at MC1R and ASIP and association with
462 melanoma-specific survival. *Int J Cancer* **136**, 2659-2667, doi:10.1002/ijc.29317
463 (2015).
- 464 43 Aguet, F. *et al.* The GTEx Consortium atlas of genetic regulatory effects across human
465 tissues. *bioRxiv*, 787903, doi:10.1101/787903 (2019).
- 466 44 Wolf Horrell, E. M., Boulanger, M. C. & D'Orazio, J. A. Melanocortin 1 Receptor:
467 Structure, Function, and Regulation. *Front Genet* **7**, 95, doi:10.3389/fgene.2016.00095
468 (2016).
- 469 45 Zhang, B. *et al.* Large-scale genetic study in East Asians identifies six new loci
470 associated with colorectal cancer risk. *Nature Genetics* **46**, 533-542,
471 doi:10.1038/ng.2985 (2014).
- 472 46 Calon, A. *et al.* Dependency of colorectal cancer on a TGF-beta-driven program in
473 stromal cells for metastasis initiation. *Cancer Cell* **22**, 571-584,
474 doi:10.1016/j.ccr.2012.08.013 (2012).
- 475 47 Venkitachalam, S. *et al.* Biochemical and functional characterization of glycosylation-
476 associated mutational landscapes in colon cancer. *Sci Rep* **6**, 23642,
477 doi:10.1038/srep23642 (2016).
- 478 48 Ishida, H. *et al.* A novel beta1,3-N-acetylglucosaminyltransferase (beta3Gn-T8),
479 which synthesizes poly-N-acetyllactosamine, is dramatically upregulated in colon
480 cancer. *Febs Letters* **579**, 71-78, doi:10.1016/j.febslet.2004.11.037 (2005).
- 481 49 Solomon, T. *et al.* Identification of Common and Rare Genetic Variation Associated
482 With Plasma Protein Levels Using Whole-Exome Sequencing and Mass Spectrometry.
483 *Circ Genom Precis Med* **11**, e002170, doi:10.1161/circgen.118.002170 (2018).
- 484 50 Smith, J. G. & Gerszten, R. E. Emerging Affinity-Based Proteomic Technologies for
485 Large-Scale Plasma Profiling in Cardiovascular Disease. *Circulation* **135**, 1651-1664,
486 doi:10.1161/circulationaha.116.025446 (2017).
- 487 51 Zheng, J. *et al.* Phenome-wide Mendelian randomization mapping the influence of the
488 plasma proteome on complex diseases. *bioRxiv*, 627398, doi:10.1101/627398 (2019).
- 489 52 Liu, B., Gloudemans, M. J., Rao, A. S., Ingelsson, E. & Montgomery, S. B. Abundant
490 associations with gene expression complicate GWAS follow-up. *Nature Genetics* **51**,
491 768-769, doi:10.1038/s41588-019-0404-0 (2019).
- 492 53 Harris, T. B. *et al.* Age, Gene/Environment Susceptibility-Reykjavik Study:
493 multidisciplinary applied phenomics. *Am J Epidemiol* **165**, 1076-1087,
494 doi:10.1093/aje/kwk115 (2007).
- 495 54 Grove, M. L. *et al.* Best practices and joint calling of the HumanExome BeadChip: the
496 CHARGE Consortium. *PLoS One* **8**, e68095, doi:10.1371/journal.pone.0068095
497 (2013).

- 498 55 Candia, J. *et al.* Assessment of Variability in the SOMAscan Assay. *Sci Rep* **7**, 14248,
499 doi:10.1038/s41598-017-14755-5 (2017).
- 500 56 Max Kuhn, K. J. *Applied Predictive Modeling*. (Springer, 2013).
- 501 57 Price, A. L. *et al.* Principal components analysis corrects for stratification in genome-
502 wide association studies. *Nat Genet* **38**, 904-909, doi:10.1038/ng1847 (2006).
- 503 58 Yang, J. *et al.* Conditional and joint multiple-SNP analysis of GWAS summary
504 statistics identifies additional variants influencing complex traits. *Nat Genet* **44**, 369-
505 375, s361-363, doi:10.1038/ng.2213 (2012).
- 506 59 Yang, J., Lee, S. H., Goddard, M. E. & Visscher, P. M. GCTA: a tool for genome-
507 wide complex trait analysis. *Am J Hum Genet* **88**, 76-82,
508 doi:10.1016/j.ajhg.2010.11.011 (2011).
- 509 60 Chun, S. & Fay, J. C. Identification of deleterious mutations within three human
510 genomes. *Genome Res* **19**, 1553-1561, doi:10.1101/gr.092619.109 (2009).
- 511 61 Carter, H., Douville, C., Stenson, P. D., Cooper, D. N. & Karchin, R. Identifying
512 Mendelian disease genes with the variant effect scoring tool. *BMC Genomics* **14**
513 **Suppl 3**, S3, doi:10.1186/1471-2164-14-s3-s3 (2013).
- 514 62 Reva, B., Antipin, Y. & Sander, C. Predicting the functional impact of protein
515 mutations: application to cancer genomics. *Nucleic Acids Res* **39**, e118,
516 doi:10.1093/nar/gkr407 (2011).
- 517 63 Schwarz, J. M., Rödelberger, C., Schuelke, M. & Seelow, D. MutationTaster
518 evaluates disease-causing potential of sequence alterations. *Nat Methods* **7**, 575-576,
519 doi:10.1038/nmeth0810-575 (2010).
- 520 64 Hemani, G. *et al.* The MR-Base platform supports systematic causal inference across
521 the human phenome. *Elife* **7**, doi:10.7554/eLife.34408 (2018).
- 522 65 Sudlow, C. *et al.* UK biobank: an open access resource for identifying the causes of a
523 wide range of complex diseases of middle and old age. *PLoS Med* **12**, e1001779,
524 doi:10.1371/journal.pmed.1001779 (2015).
- 525 66 Burgess, S., Butterworth, A. & Thompson, S. G. Mendelian randomization analysis
526 with multiple genetic variants using summarized data. *Genet Epidemiol* **37**, 658-665,
527 doi:10.1002/gepi.21758 (2013).
- 528 67 MacArthur, J. *et al.* The new NHGRI-EBI Catalog of published genome-wide
529 association studies (GWAS Catalog). *Nucleic Acids Research* **45**, D896-D901,
530 doi:10.1093/nar/gkw1133 (2017).

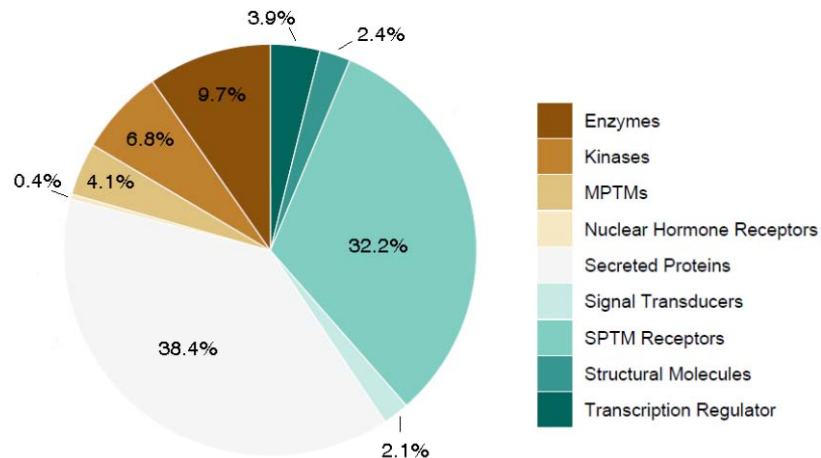
531

Table 1 | Selected examples of exome array variants affecting serum protein levels and complex disease. CHD, coronary heart disease; VTE, venous thromboembolism; CKD, chronic kidney disease; T2D, type 2 diabetes; VAT, visceral adipose tissue; LOAD, late-onset Alzheimer’s disease; SLE, systemic lupus erythematosus; IBD, inflammatory bowel disease; AMD, age-related macular degeneration; N/A, not applicable. All reported effects are genome-wide significant at $P < 1.92 \times 10^{-10}$.

Disease class	Disease trait	PMID or database	pQTL	GWAS lead SNP(s) ^a	Function pSNP ^b	Mapped GWAS locus ^c	#Proteins affected	Example of <i>cis</i> and/or <i>trans</i> affected proteins ^d
<i>Cardiovascular</i>								
	CHD	28714975	rs12740374	rs12740374	3'-UTR	CELSR2	8	C1QTNF1, IGFBP1
	VTE	UKBB, 28373160	rs2343596	rs16873402, rs4602861	Intron	ZFPM2	7	VEGFA, DKK1
	Stroke	26708676	rs653178	rs653178	Intron	ATXN2	2	THPO, CXCL11
<i>Metabolic</i>								
	T2D	22885922	rs7202877	rs7202877	Intergenic	CTRB1	5	CTRB1, PRSS2, CPB1
	VAT	20935629	rs9491696	rs9491696	Intron	RSPO3	1	RSPO3
	Triglyceride	21386085	rs2266788	rs2266788	3'-UTR	APOA5	5	APOA5, PCSK7, ANGPTL3
<i>CNS</i>								
	LOAD	21460840	rs610932	rs610932	3'-UTR	MS4A6A	3	TREM2, GLTPD2
	Parkinson	21738487	rs6599389	rs6599389	Intron	GAK	1	IDUA
	Schizophrenia	25056061	rs3617	rs3617	Q315K	ITIH3	8	ITIH3, JAKMIP3
<i>Inflammatory</i>								
	SLE, T1D	26502338	rs2304256	rs2304256	V362F	TYK2	2	ICAM1, ICAM5
	Crohn’s, IBD	21102463	rs11209026	rs11209026	R381Q	IL23R	1	IL23R
	AMD	2355636	rs10737680	rs10737680	Intron	CFH	22	CFH, CFHR1, CFB
<i>Cancer</i>								
	Colorectal	24836286	rs2241714	rs1800469	I11M	TMEM91	3	B3GNT2, TGFB1
	Lung	18978787	rs3117582	rs3117582	Intron	APOM	10	MICB, ISG15
	Melanoma	18488026	rs910873	rs910873	Intron	PIGU	1	ASIP

^aProtein QTLs overlapping GWAS lead SNPs using the PhenoScanner database²³. No SNP proxies were applied except when the lead pSNP was not in the query then we used the best proxy ($r^2 \geq 0.8$ between markers). ^bThe functional annotation of pQTLs was obtained from the PhenoScanner database²³. ^cReported causal candidates are from the GWAS Catalog⁶⁷. ^dThe definition of *cis* vs. *trans* effects is somewhat arbitrary depending on the window size chosen across the protein gene in question. However, all affected proteins located at other chromosomes than the pQTL location, were considered *trans* acting and are highlighted in bold letters. All significant pQTLs are listed in Supplementary Table 1 and the overlap with GWAS risk loci summarized in Supplementary Table 4.

a



b

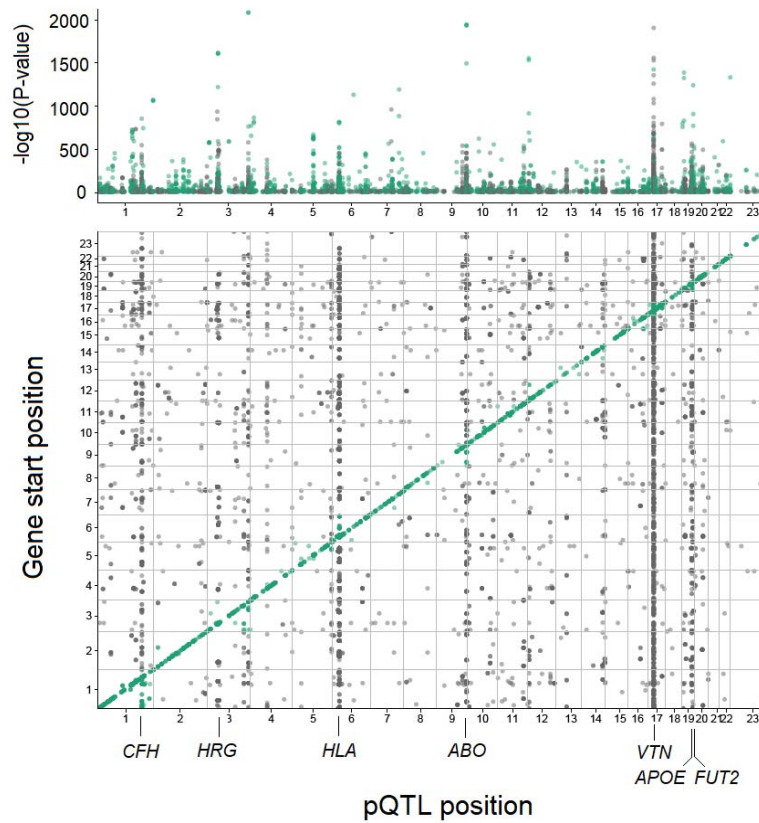


Fig. 1. Classification of the target protein population and genomic locations of observed pQTLs. a. Pie chart showing the relative distribution (percentage) of the different protein classes targeted by the present proteomics platform, with secreted proteins (38.4%) and single

pass transmembrane (SPTM) receptors (32.2%) dominating the target protein population.

Protein classes were manually curated based on information from the SecTrans, Gene Ontology (GO) and Swiss-Prot databases, and were composed of secreted proteins (e.g. cytokines, adipokines, hormones, chemokines and growth factors), SPTM receptors (e.g. tyrosine and serine/threonine kinase receptors), multi-pass transmembrane (MPTM) receptors (e.g. GPCR, ion channels, transporters), enzymes (intracellular), kinases, nuclear hormone receptors, structural molecules, transcriptional regulators and signal transducers. **b.** The Manhattan plot in the top panel uses precise P-values to highlight all study-wide significant associations in Supplementary Table S1. The bottom panel shows the genomic locations of all study-wide significant pQTLs ($P < 1.92 \times 10^{-10}$), where the start position of the protein encoding gene is shown on the y-axis and the location of the pSNP at the x-axis. *Cis* acting effects, using a 300kb window, appear at the diagonal while *trans* acting pQTL effects including *trans* hot spots show up off-diagonally. The genetic loci highlighted across the x-axis are *trans*-acting hotspots.

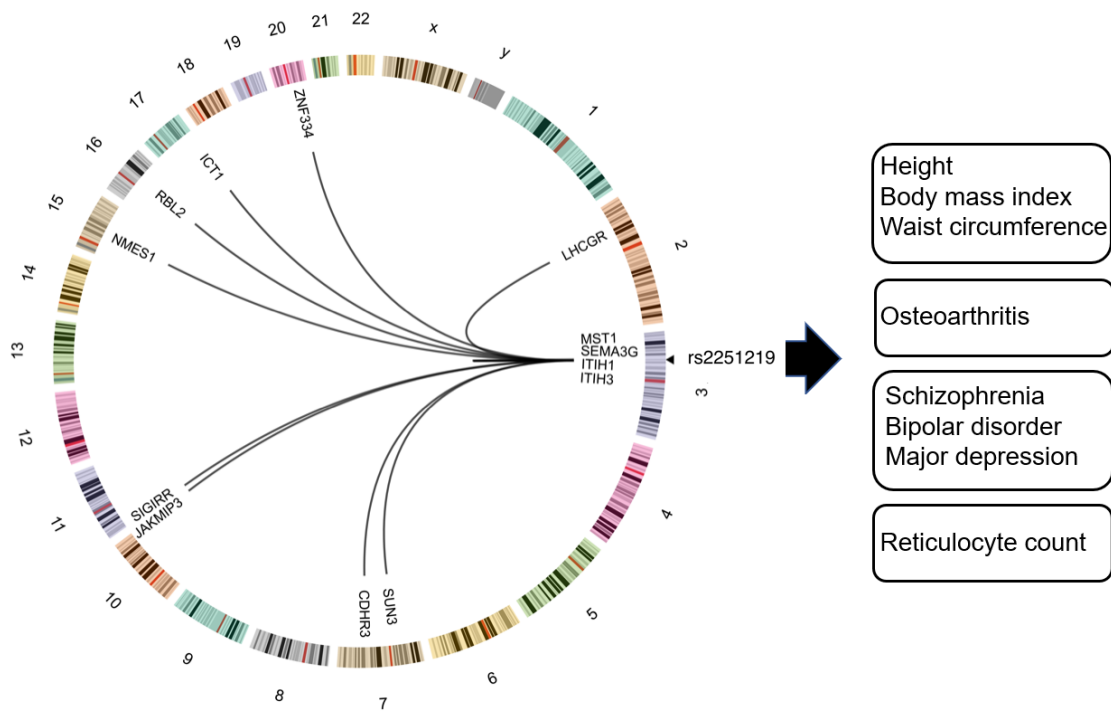


Fig. 2. Pleiotropy of rs2251219 affecting many proteins and disease traits. a. Circos plot showing the effect of the variant rs2251219 (Supplementary Tables 1 and 2) on 13 proteins acting in *cis* or *trans* and sharing genetics with various diseases of different etiologies. Only study-wide significant ($P < 1.92 \times 10^{-10}$) genotype-to-protein associations are shown. Lines going from rs2251219 show links to genomic locations of the protein encoding genes affected while numbers refer to chromosomes.

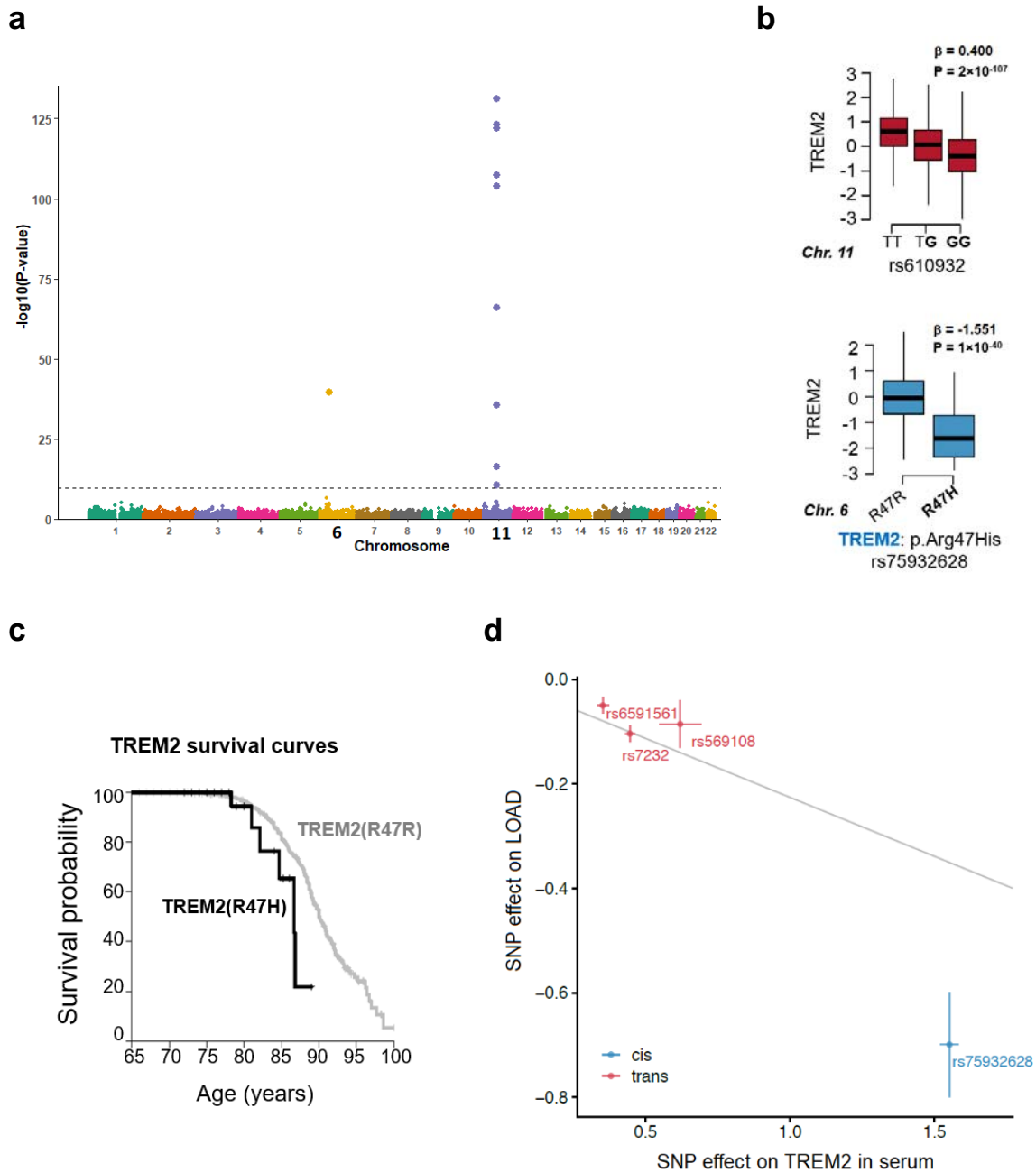


Fig. 3. Effects of distinct risk loci for LOAD converge on the protein TREM2. **a.** The Manhattan plot highlights variants at two distinct chromosomes associated with serum TREM2 levels. Study-wide significant associations at $P < 1.92 \times 10^{-10}$ are indicated by the horizontal line. The y-axis shows the $-(\log_{10})$ of the P-values for the association of each genetic variant on the exome array present along the x-axis. Variants at both chromosomes 6

and 11 associated with TREM2 have been independently linked to risk of LOAD including the rs75932628 (NP_061838.1: p.R47H) in TREM2 at chromosome 6 and the variant rs610932 at chromosome 11. **b.** Boxplot of the *trans* effect of the well-established GWAS risk variant rs610932 for LOAD on TREM2 serum levels (upper panel), where the LOAD risk allele G (highlighted in bold) is associated with lower levels of TREM2. Similarly, the LOAD causing p.R47H mutation was associated with low levels of TREM2 (lower panel). **c.** TREM2p.R47H carriers demonstrated lower survival probability post-incident LOAD compared to TREM2p.R47R carriers ($P = 0.04$). **d.** Scatterplot for the TREM2 protein supported as having a causal effect on LOAD in a two sample MR analysis. The figure demonstrates the estimated effects (with 95% confidence intervals) of their respective *cis*- and *trans*-acting genetic instruments on the serum TREM2 levels in AGES-RS (x-axis) and risk of LOAD through a GWAS by Kunkle et al.³² (y-axis), using 21,982 LOAD cases and 41,944 controls. The line indicates the inverse variance weighted causal estimate ($\beta = -0.226$, $SE = 0.057$, $P = 7.6 \times 10^{-5}$).

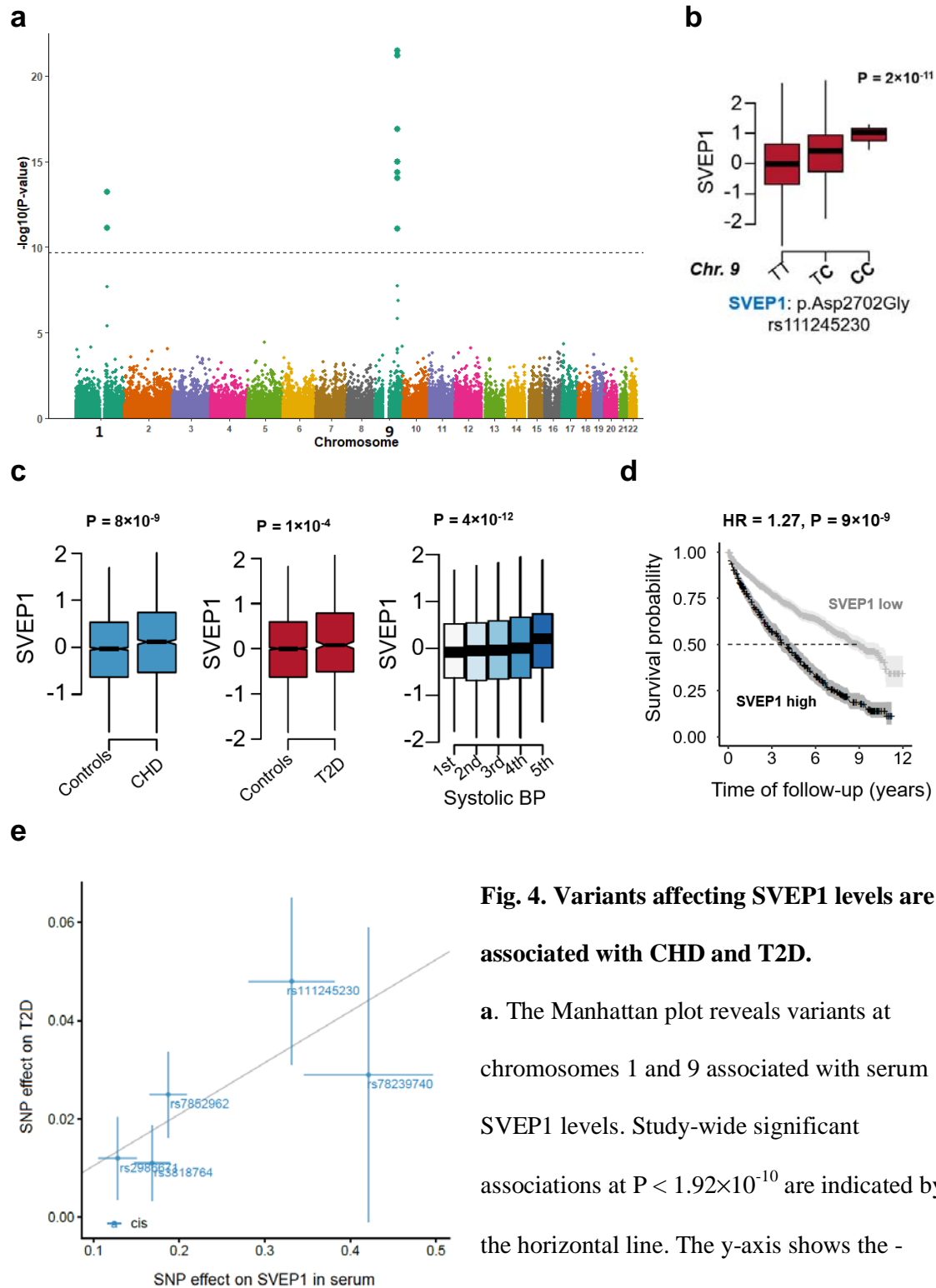


Fig. 4. Variants affecting SVEP1 levels are associated with CHD and T2D.

a. The Manhattan plot reveals variants at chromosomes 1 and 9 associated with serum SVEP1 levels. Study-wide significant associations at $P < 1.92 \times 10^{-10}$ are indicated by the horizontal line. The y-axis shows the $-\log_{10}$ of the P-values for the association of each genetic variant on the exome array present along the x-axis. **b.** One of the variants associated with SVEP1 levels and underlying the peak at chromosome 9 is the low-frequency

CHD risk variant rs111245230 (NP_699197.3: pAsp2702Gly). The CHD risk allele **C** (highlighted in bold) is associated with increased serum SVEP1 levels. **c.** Serum levels of SVEP1 were associated with CHD ($P = 8 \times 10^{-9}$), T2D ($P = 1 \times 10^{-4}$) and systolic blood pressure ($P = 4 \times 10^{-12}$) in the AGES-RS, all in a directionally consistent manner. **d.** Consistent with the directionality of the effects described above, we find that elevated levels of SVEP1 were associated with higher rates of mortality post-incident CHD. **e.** Scatterplot for the SVEP1 protein supported as having a causal effect on T2D in a two-sample MR analysis. The figure demonstrates the estimated effects (with 95% confidence intervals) of the SNP effect on serum SVEP1 levels and T2D from a GWAS in Europeans³⁴ (y-axis), with 74,124 T2D patients and 824,006 controls. The line indicates the inverse variance weighted causal estimate ($\beta = 0.105$, $SE = 0.024$, $P = 1.2 \times 10^{-5}$).

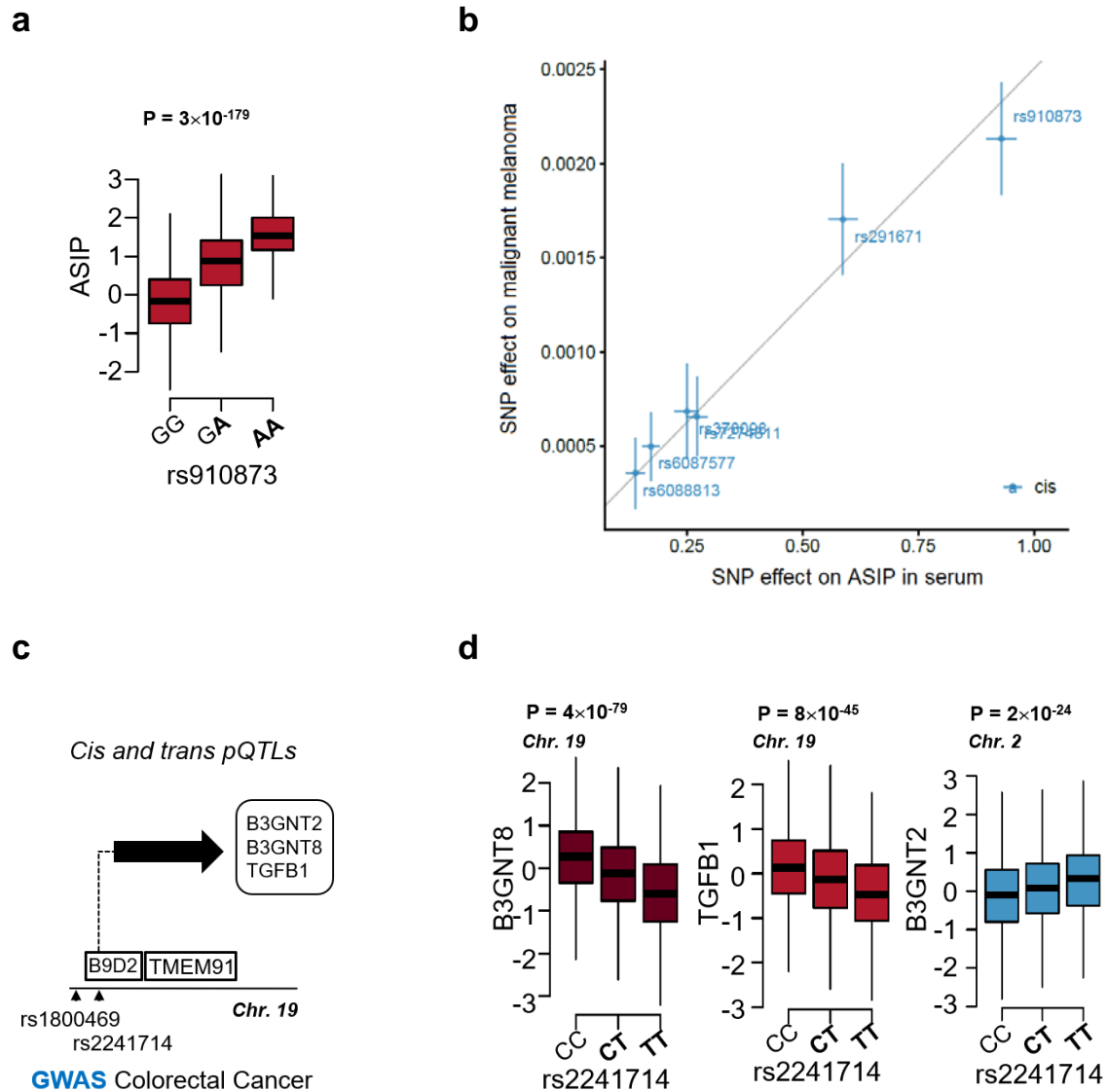


Figure 5. Proteins associated with malignant melanoma and colorectal cancer. **a.** The melanoma risk allele A (highlighted in bold) for the variant rs910873 is associated with high serum levels of ASIP. **b.** Scatterplot for the ASIP protein supported as having a causal effect on malignant melanoma in a two sample MR analysis. The figure demonstrates the estimated effects (with 95% confidence intervals) of their respective genetic instruments on the serum ASIP levels in AGES (x-axis) and risk of melanoma in GWAS by UK biobank data (UKB-b-12915)⁶⁵ (y-axis), that included 3598 melanoma cases and 459,335 controls. The line indicates the inverse variance weighted causal estimate ($\beta = 0.0025$, $SE = 0.0002$, $P = 4.8 \times 10^{-$

²⁶). **c.** The pQTL rs2241714 is a proxy for the colorectal cancer associated variant rs1800469 ($r^2 = 0.978$) (Supplementary Table 2), located within the gene *B9D2* and proximal to *TMEM91* which is the reported candidate gene at this locus (see Table 1). **d.** The variant rs2241714 (and rs1800469) regulate three serum proteins, B3GNT2 (in *trans*), B3GNT8 (in *cis*) and TGFB1 (in *cis*).

SUPPLEMENTARY MATERIAL

Coding and regulatory variants affect serum protein levels and common disease

Valur Emilsson^{1,2, ξ,*}, Valborg Gudmundsdottir^{1,ξ}, Alexander Gudjonsson^{1,ξ}, Mohd A Karim^{3,4}, Marjan Ilkov¹, James R. Staley⁵, Elias F. Gudmundsson¹, Brynjolfur G. Jonsson¹, Lenore J. Launer⁶, Jan H. Lindeman⁷, Nicholas M. Morton⁸, Thor Aspelund¹, John R. Lamb⁹, Lori L. Jennings¹⁰ and Vilmundur Gudnason^{1,2,*}

¹Icelandic Heart Association, Holtasmari 1, IS-201 Kopavogur, Iceland.

²Faculty of Medicine, University of Iceland, 101 Reykjavik, Iceland

³Wellcome Trust Sanger Institute, Wellcome Genome Campus, Hinxton, Cambridgeshire CB10 1SA, UK.

⁴Open Targets, Wellcome Genome Campus, Hinxton, Cambridgeshire CB10 1SD, UK

⁵MRC/BHF Cardiovascular Epidemiology Unit, Department of Public Health and Primary Care, University of Cambridge, Cambridge, UK

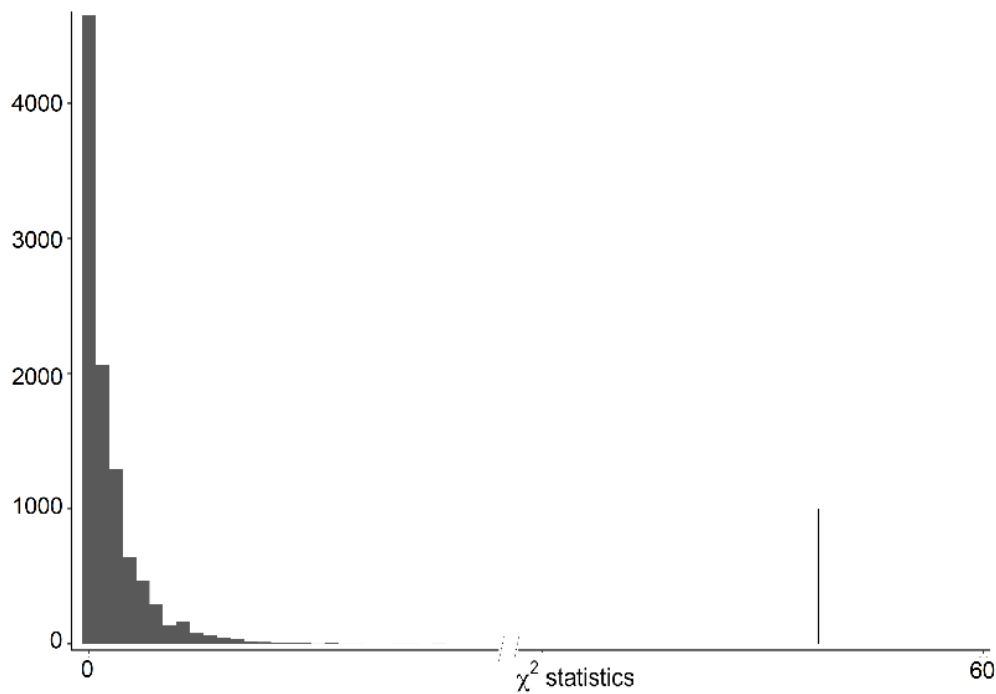
⁶Laboratory of Epidemiology and Population Sciences, Intramural Research Program, National Institute on Aging, Bethesda, MD 20892-9205, USA.

⁷Department of General Surgery Leiden University Medical Center, Leiden. Holland

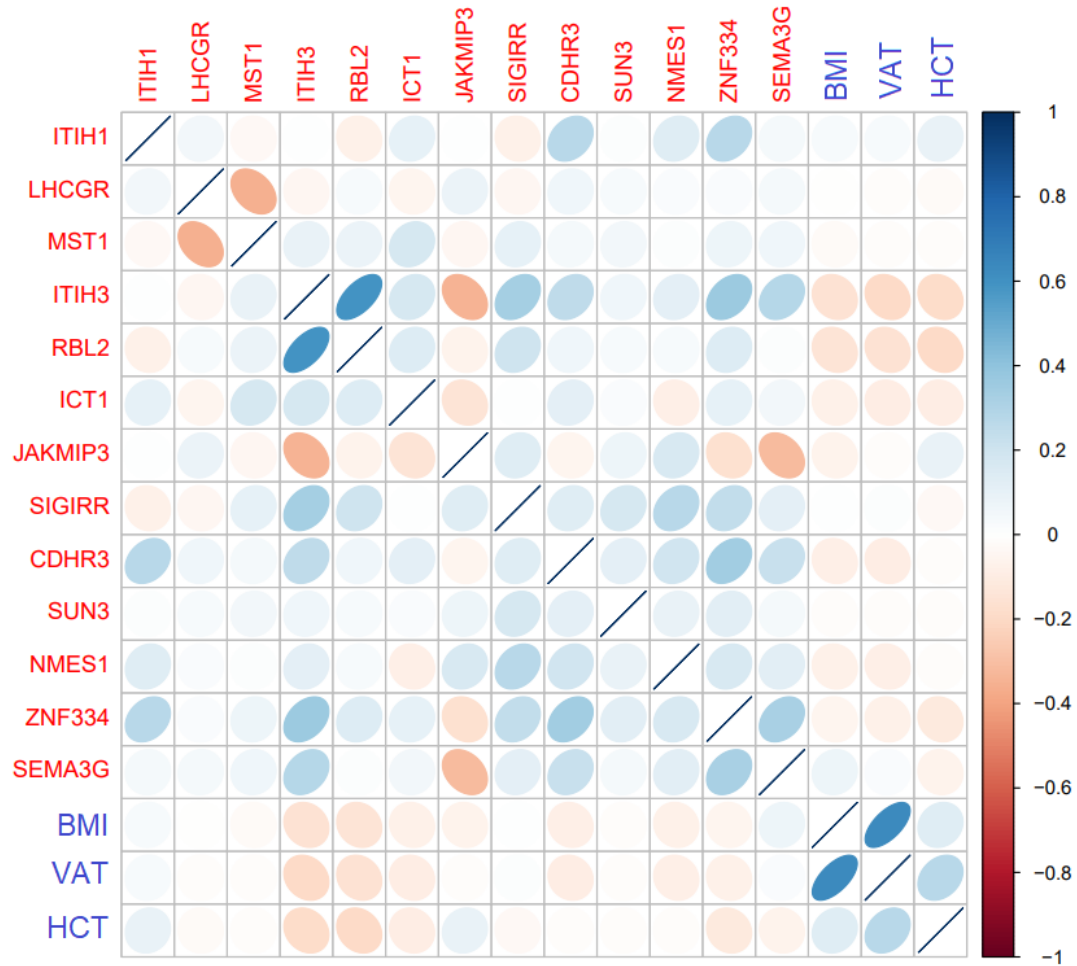
⁸Centre for Cardiovascular Sciences, Queen's Medical Research Institute, University of Edinburgh, Edinburgh EH16 4TJ, UK

⁹GNF Novartis, 10675 John Jay Hopkins Drive, San Diego, CA 92121, USA.

¹⁰Novartis Institutes for Biomedical Research, 22 Windsor Street, Cambridge, MA 02139, USA.

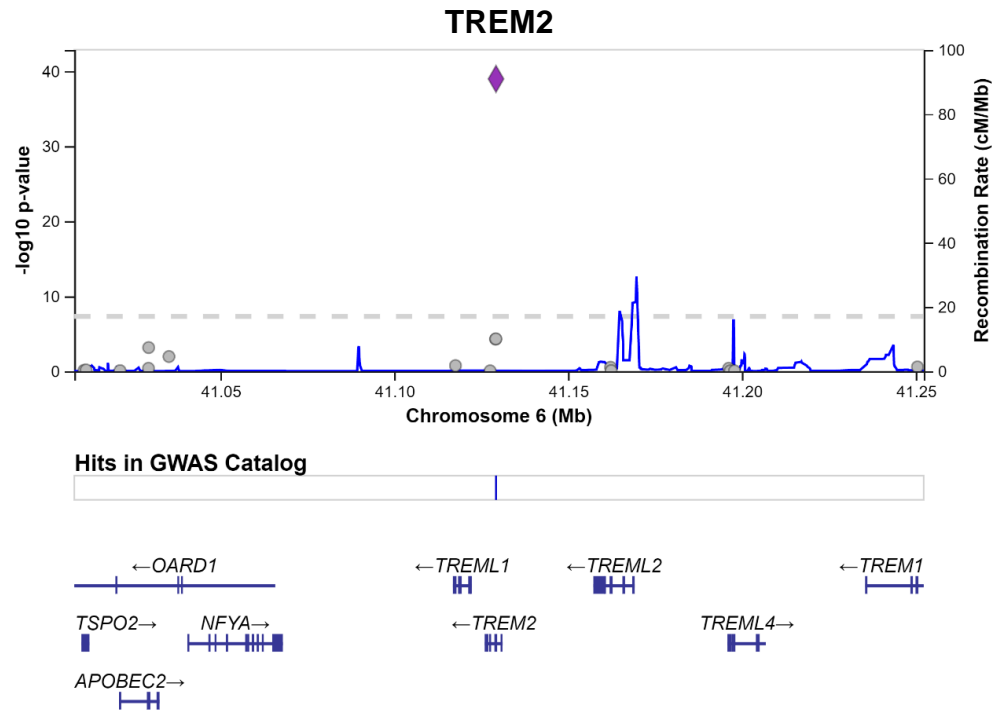


Supplementary Fig. S1. Empirical distribution of the test statistic as a histogram and the observed statistics calculated from our data as a vertical line. 10,000 permutations were performed to obtain the empirical distribution of the χ^2 test of equality of proportions of pQTLs among secreted versus non-secreted proteins. Here, the test statistics calculated from our data to the quantiles of this distribution to obtain $P(\text{Data}|\text{H}_0)$ were compared. Of 10,000 permutations none gave a value greater than the observed statistic leading us to $P\text{-value} = P(\text{Data}|\text{H}_0) < 0.0001$.

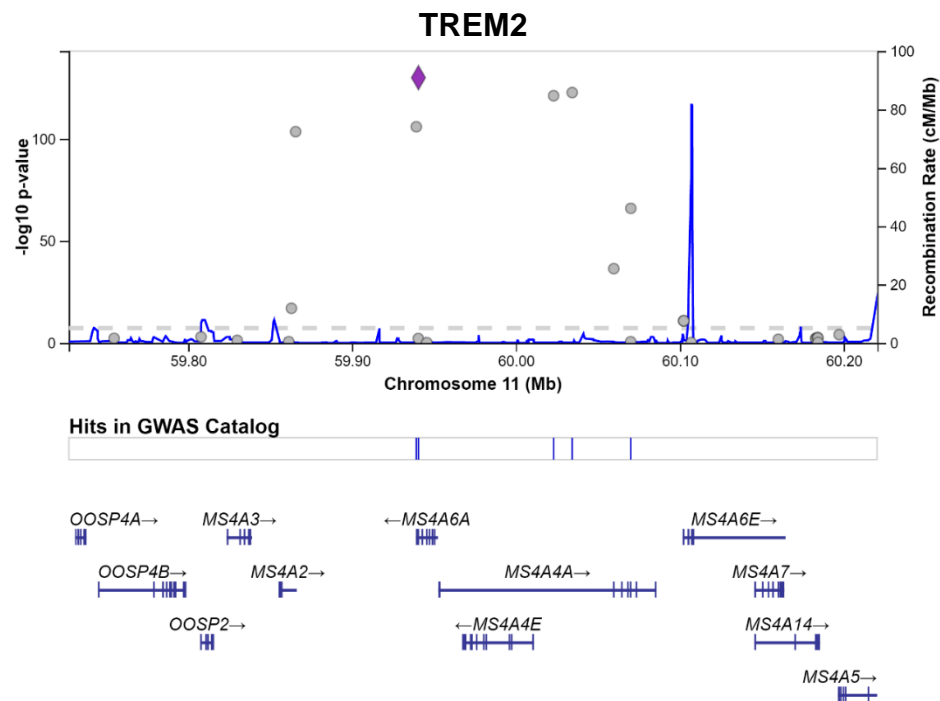


Supplementary Fig. S2. A Spearman rank correlation between all proteins as well as some quantitative traits including body mass index (BMI, kg/m²), visceral adipose tissue (VAT, measured *via* computed tomography) and hematocrit (HCT), that were associated with rs2251219.

a

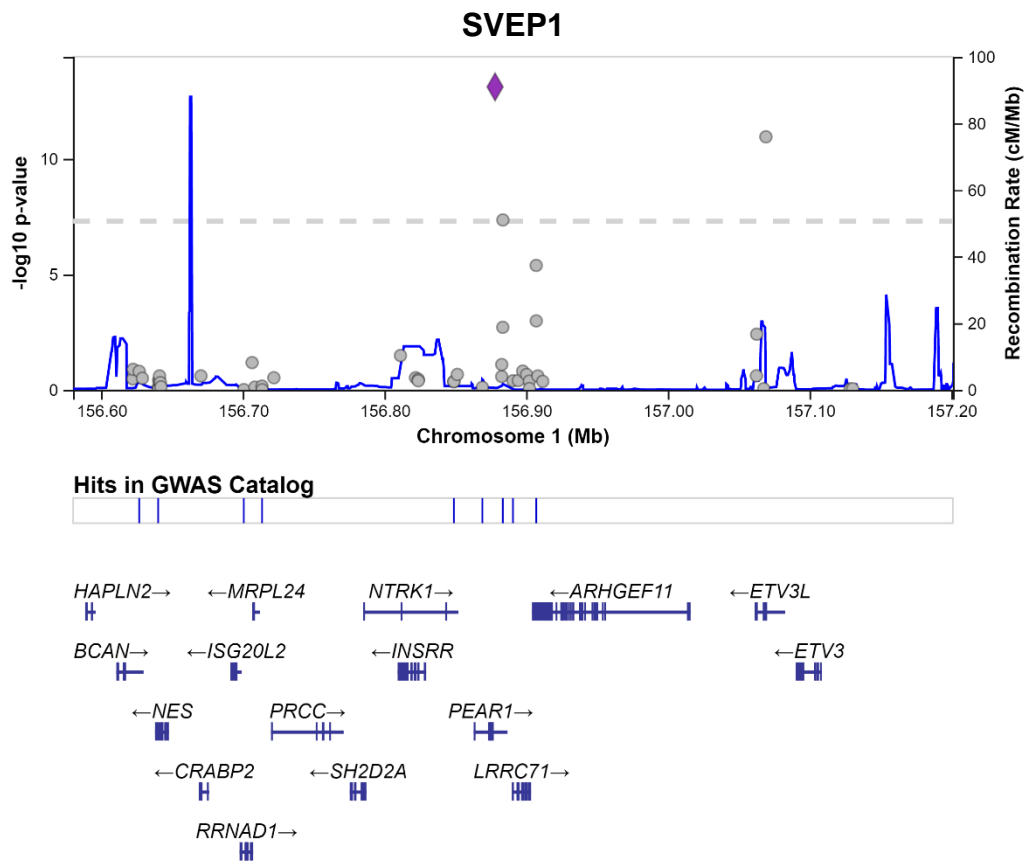


b

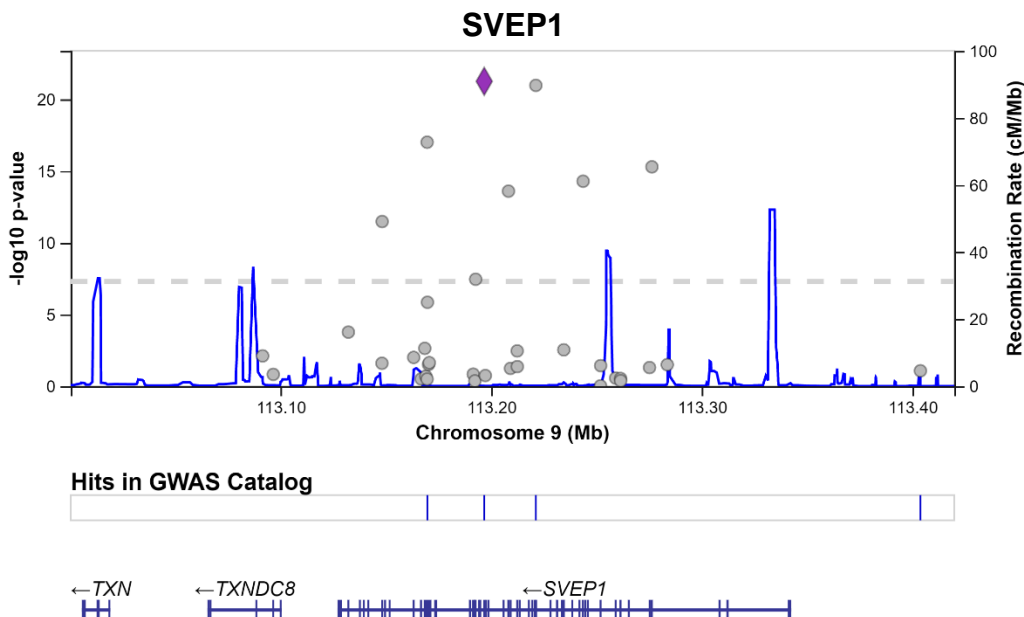


Supplementary Fig. S3. TREM2 regional plots (LocusZoom) based on exome array variants at chromosomes 6 and 11.

a



b



Supplementary Fig. S5. SVEP1 regional plots (LocusZoom) based on exome array variants at chromosomes 1 and 9.