

Human serum proteome profoundly overlaps with genetic signatures of disease

Valur Emilsson^{1,2, §,*}, Valborg Gudmundsdottir^{1,§}, Marjan Ilkov¹, James R. Staley³, Alexander Gudjonsson¹, Elias F. Gudmundsson¹, Lenore J. Launer⁴, Jan H. Lindeman⁵, Nicholas M. Morton⁶, Thor Aspelund¹, John R. Lamb⁷, Lori L. Jennings⁸ and Vilmundur Gudnason^{1,2,*}

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

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

³MRC/BHF Cardiovascular Epidemiology Unit, Department of Public Health and Primary Care, University of Cambridge, Cambridge, UK

⁴Laboratory of Epidemiology and Population Sciences, Intramural Research Program, National Institute on Aging, Bethesda, MD 20892-9205, USA.

⁵Department of General Surgery Leiden University Medical Center, Leiden. Holland

⁶Centre for Cardiovascular Sciences, Queen's Medical Research Institute, University of Edinburgh, Edinburgh EH16 4TJ, UK

⁷GNF Novartis, 10675 John Jay Hopkins Drive, San Diego, CA 92121, USA.

⁸Novartis Institutes for Biomedical Research, 22 Windsor Street, Cambridge, MA 02139, USA.

[§]These authors contributed equally

*Corresponding authors. Emails: valur@hjarta.is and v.gudnason@hjarta.is

Abstract

Circulating proteins are prognostic for human outcomes including cancer, heart failure, brain trauma and brain amyloid plaque burden. A deep serum proteome survey recently revealed close associations of serum protein networks and common diseases. The present study reveals unprecedented number of individual serum proteins that overlap genetic signatures of diseases emanating from different tissues of the body. Here, 55,932 low-frequency and common exome-array variants were compared with 4782 protein measurements in the serum of 5457 individuals of the deeply annotated AGES Reykjavik cohort. At a Bonferroni adjusted P-value threshold $< 2.16 \times 10^{-10}$, 5553 variants affecting levels of 1931 serum proteins were detected. These associated variants overlapped genetic loci for hundreds of complex disease traits, emphasizing the emerging role for serum proteins as biomarkers of and potential causative agents of multiple diseases.

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

like mRNA and/or protein levels with genetics and disease traits, the identification of the causal candidates can be enhanced³⁻⁶.

Proteins are arguably the ultimate players in all life processes in disease and health, however, high throughput detection and quantification of proteins has been hampered by the limitations of available proteomic technologies. Recently, a custom-designed Slow-Off rate Modified Aptamer (SOMAmer) protein profiling platform was developed to measure 4782 proteins encoded by 4137 human genes in the serum of 5457 individuals from the AGES Reykjavik study (AGES-RS)¹⁰, resulting in 26.1 million individual protein measurements. Various metrics related to the performance of the proteomic platform including aptamer specificity, assay variability and reproducibility have already been described¹⁰. We demonstrated that the human serum proteome is under strong genetic control¹⁰, in line with findings of others applying identical or different proteomics technologies^{11,12}. Moreover, serum proteins were found to exist in regulatory groups of network modules composed of members synthesized in all tissues of the body, suggesting that system level coordination or homeostasis is mediated to a significant degree by thousands of proteins in blood. Importantly, the deep serum and plasma proteome is associated with and prognostic for various diseases as well as human life span^{10,13-19}.

Here, we regressed levels of 4782 proteins on 55,932 low-frequency and common variants from the HumanExome BeadChip exome array, in sera from 5457 individuals of the deeply phenotyped AGES-RS cohort. Further cross-referencing of all significant genotype-to-protein associations to hundreds of genetic loci for various disease endpoints and clinical traits, demonstrated profound overlap between the genetics of circulating proteins and disease related phenotypes. We highlight how triangulation of data from different sources can link

genetics, protein levels and disease(s), in order to cross-validate one another and point to potentially causal relationship between proteins and complex disease(s).

Using genotype data from an exome array (HumanExome BeadChip) enriched for structural variants and tagged for many GWAS risk loci (**Methods**), the effect of low-frequency and common variants on the deep serum proteome was examined. Quality control filters²⁰, and exclusion of monomorphic variants reduced the available variants to 72,766. Additionally, we excluded variants at minor allele frequency (MAF) < 0.001 as they provide insufficient power for single-point association analysis²¹. This resulted in 55,932 low-frequency and common variants that were tested for association to each of the 4782 human serum protein measurements using linear regression analysis adjusted for the confounders age and sex (**Methods**). The current platform targets the serum proteome arising largely from active or passive secretion, ectodomain shedding, lysis and/or cell death^{10,22}. **Figure 1a** highlights the classification of the protein population targeted by the aptamer-based profiling platform, showing over 70% of the proteins are secreted or single pass transmembrane (SPTM) receptors.

Applying a Bonferroni corrected significance threshold of $P < 2.16 \times 10^{-10}$ ($0.05/55932/4137$) we detected 5553 exome array variants that were associated with variable levels of 1931 serum proteins (**Supplementary Table 1** and **Fig. 1b**). These protein quantitative trait loci (pQTLs) were *cis* and/or *trans* acting including several *trans* acting hotspots with pleiotropic effects on multiple co-regulated proteins (**Fig. 1b**). When compared to other protein classes, secreted proteins were enriched for pQTLs (53.2% vs. 43.6%, FET $P = 1.8 \times 10^{-7}$), underscoring that proteins destined for the systemic environment are under stronger genetic control than other proteins detected by the platform. Next, we cross-referenced all the 5553 pQTLs with a comprehensive collection of genetic loci associated with diseases and clinical

traits from the curated PhenoScanner database²³, revealing profound overlap of pQTLs with known GWAS loci, or 60% of all pQTLs linked to at least single trait (**Supplementary Tables 1 and 2**). Moreover, we observed a significant correlation between the number of serum proteins affected by a given variant and the number of associated phenotypes (Spearman $\rho = 0.570$, $P < 0.001$) (**Fig. 2a**). When we exclude the many associations driven by variants located at chromosome 6 (**Supplementary Table 1**), the correlation was weaker yet significant (Spearman $\rho = 0.22$, $P < 0.001$). In summary, this suggests that greater regulatory pleiotropy of pQTLs is associated with greater chance of disease trait pleiotropy, which agrees well with recent gene expression eQTL studies linked to common disease traits^{24,25}. **Figure 2b** highlights an example of a pleiotropic effect at the locus rs2251219 affecting several proteins and sharing genetics with different diseases and traits. **Table 1** highlights a selected set of pQTLs that share genetics with diseases of different etiologies including disorders of the brain, metabolism, immune and cardiovascular systems and cancer. In the sections that follow we specify examples of serum pQTLs overlapping disease risk loci and demonstrate how triangulation of data from different sources can cross-validate one another.

Variable levels of the anti-inflammatory protein TREM2 were associated with two distinct genomic regions at $P < 2.16 \times 10^{-10}$ (**Fig. 3a**). This included the missense variant rs75932628 (NP_061838.1: p.R47H) in TREM2 at chromosome 6 (**Fig. 3b**), known to confer a strong risk of late-onset Alzheimer's disease (LOAD)²⁶. The variant was also associated with IGFBL1 ($P = 3 \times 10^{-18}$) in serum (**Supplementary Table 1**), which has recently been implicated in axonal growth²⁷. Intriguingly, the region at chromosome 11 associated with soluble TREM2 levels harbors variants adjacent to the genes *MS4A4A* and *MS4A6A* including rs610932 known to influence genetic susceptibility for LOAD²⁸ (**Table 1, Fig. 3a, b**). The variant rs610932 was also associated with the proteins GLTPD2 and A4GALT (**Supplementary Table 1**). The alleles increasing risk of LOAD for both the common variant rs610932 and the

low-frequency variant rs75932628 were associated with low levels of soluble TREM2 (**Fig. 3b**). Consistently, we find that the high-risk allele for rs75932628 was associated with accelerated mortality post incident LOAD in the AGES-RS (**Fig. 3c**). It is of note that the levels of TREM2 in the cerebrospinal fluid (CSF) reflect the activity of brain TREM2-triggered microglia^{4,29}, while high levels of CSF TREM2 have been associated with improved cognitive functioning³⁰. **Figure 3d** highlights the correlation relationship (Spearman rank) between the different proteins affected by the LOAD risk loci at chromosomes 6 and 11. The accumulated data show a directionally consistent effect at independent risk loci for LOAD converging on the same causal candidate TREM2. In summary, these results demonstrate that the effect of genetic drivers on major brain-linked disease like LOAD can be readily detected in serum to both inform on the causal relationship and the directionality of the risk mediating effect. This would also suggest that serum may be an accessible proxy for microglia function and cognition.

Variable levels of the cell adhesion protein SVEP1 are associated with variants located at chromosomes 1 and 9 (**Supplementary Table 1** and **Fig. 4a**). Genetic associations to SVEP1 levels at chromosome 9 include the low-frequency missense variant rs111245230 in SVEP1 (NP_699197.3: pD2702G) (**Fig. 4b**), which was recently linked to coronary heart disease (CHD), blood pressure and type-2-diabetes (T2D)³¹. Overall, we found eight different missense mutations in *SVEP1* that were associated with SVEP1 serum levels (**Supplementary Table 1**). The risk allele C of rs111245230 was associated with elevated levels of SVEP1 in CHD and T2D patients compared to a group of controls free of either disease (**Fig. 4c**). Furthermore, high SVEP1 levels were positively correlated with systolic blood pressure ($\beta = 2.10$, $P = 4 \times 10^{-12}$) (**Fig. 4c**), but not with diastolic blood pressure ($\beta = 0.115$, $P = 0.413$). Consistently, higher serum levels of SVEP1 were associated with increased mortality post-incident CHD in the AGES-RS ($HR = 1.27$, $P = 9 \times 10^{-9}$) (**Fig. 4d**). The variants

at chromosome 1 linked to SVEP1 levels (**Fig. 4a**), have not previously been linked to any disease. Our data triangulation links genetics, protein levels and disease(s) and indicates that SVEP1 may be a point of intervention to therapeutically target CHD and T2D.

The ILMN exome array contains a number of tags related to previous GWAS findings³², including many risk loci for cancer. For example, 21 loci associated with melanoma³³ and 50 loci associated with colorectal cancer³⁴. The exome array variant rs910873 located in an intron of the GPI transamidase gene *PIGU* was previously linked to melanoma risk³⁵. The reported candidate gene *PIGU* is the gene most proximal to the lead SNP rs910873 and may be a novel candidate gene involved in melanoma. However, a more biologically relevant candidate is the agouti-signaling protein (*ASIP*) gene that is located 314kb downstream of the lead SNP rs910873. *ASIP* is a competitive inhibitor of MC1R³⁶, and is thus strongly biologically implicated in melanoma risk³⁷. We found that the melanoma risk allele for rs910873 was associated with elevated ASIP serum levels ($P = 3 \times 10^{-179}$) and the variant had no effect on other proteins measured with the current proteomic platform (**Fig. 5a**, **Supplementary Table 1** and **Table 1**). Interestingly, the pQTL rs910873 is also an eQTL for *ASIP* gene expression in skin²⁴, showing directionally consistent effect on the mRNA and protein. Our data point to the ASIP protein underlying the risk at rs910873, thus providing supportive evidence for the hypothesis that ASIP mediated inhibition of MC1R results in suppression of melanogenesis and increased risk of melanoma³⁸. An additional example is the colorectal cancer locus at rs1800469³⁹, which is a proxy to the pQTL rs2241714 ($r^2=0.978$) (**Table 1** and **Fig. 5b**). While, the TMEM91 gene was the reported candidate gene for the melanoma risk at the rs1800469 (**Table 1**), we find that the risk variant affected three proteins in either *cis* or *trans*, notably B3GNT2, B3GNT8 and TGFB1 (**Fig. 5b**). Intriguingly, all three proteins have previously been implicated in colorectal cancer⁴⁰⁻⁴². Although we cannot rule out *PIGU* or *TMEM91* as candidate genes for melanoma or colorectal cancer risk,

respectively, these results provide alternate, experimentally supported and perhaps more biologically relevant candidates.

We report here that 60% of the serum proteome that is under genetic control shares genetics with reported clinical traits including major diseases emanating from different tissues of the body. This is in line with a recent population-scale survey of human induced pluripotent stem cells, demonstrating that pQTLs are 1.93-fold enriched in disease risk variants compared to a 1.36-fold enrichment for eQTLs¹², underscoring the added value in pQTL mapping. We reaffirm widespread associations between genetic variants and their cognate proteins as well as distant *trans*-acting effects on serum proteins and demonstrate that many proteins are often involved in mediating the biological effect of a single causal variant affecting complex disease. It remains a possibility that the influence of some structural variants on protein levels reflect effects mediated by regulatory variants that are in linkage disequilibrium with the pQTL of interest. In addition, protein coding variants may cause technical artifacts in both affinity proteomics and mass spectrometry^{43,44}. However, systematic conditional and colocalization analyses in causality testing using the aptamer-based technology have shown that pQTLs driven by common missense variants being artefactual is an unlikely event^{11,45}.

We note that with the ever-increasing availability of large-scale omics data aligned with the human genome, cross-referencing different datasets can result in findings that occurred by sheer chance. Hence, a systematic colocalization analysis has been proposed in causality tests between intermediate traits and disease endpoints⁴⁶. This is, however, not feasible for application of the exome array given its sparse genomic coverage. Instead, multi-omics data triangulation to infer consistency in directionality, the approach used in the present study, can enhance confidence in the causal call and offer insights and guidelines for experimental follow-up studies. For the serum proteins TREM2, SVEP1, ASIP and other examples

highlighted in the presented study (see **Table 1**), further colocalization analyses and tests of causality are warranted. We previously asserted that serum proteins are intimately connected to and mediate global homeostasis¹⁰. The accumulated data show that serum proteins are under strong genetic control and closely associated with diseases of different aetiologies, which in turn suggests that serum proteins may be significant mediators of systemic homeostasis in human health and disease.

METHODS

Study population

Participants aged 66 through 96 are from the Age, Gene/Environment Susceptibility Reykjavik Study (AGES-RS) cohort⁴⁷. AGES-RS is a single-center prospective population-based study of deeply phenotyped subjects (5764, mean age 75±6 years) and survivors of the 40-year-long prospective Reykjavik study (n~18,000), an epidemiologic study aimed to understand aging in the context of gene/environment interaction by focusing on four biologic systems: vascular, neurocognitive (including sensory), musculoskeletal, and body composition/metabolism. Descriptive statistics of this cohort as well as detailed definition of the various disease end-points and relevant phenotypes measured have been published^{10,47}. The AGES-RS was approved by the NBC in Iceland (approval number VSN-00-063), and by the National Institute on Aging Intramural Institutional Review Board, and the Data Protection Authority in Iceland.

Genotyping platform

Genotyping was conducted using the exome-wide genotyping array Illumina HumanExome-24 v1.1 Beadchip for all the 5457 subjects with protein data. The exome array was enriched for exonic variants selected from over 12,000 individual exome and whole-genome sequences from different study populations³², and includes as well tags for previously described GWAS

hits, ancestry informative markers, mitochondrial SNPs and human leukocyte antigen tags³².

A total of 244,883 variants were included on the exome array. Genotype call and quality control filters including call rate, heterozygosity, sex discordance and PCA outliers were performed as previously described^{2,20}. Variants with call rate <90% or with Hardy–Weinberg P values <1×10⁻⁷ were removed from the study. 72,766 variants were detected in at least one individual of the AGES-RS cohort. Of these variants, 55,932 had a minor allele frequency > 0.001 and were examined for association against each of the 4782 human serum protein measurements (see below).

Protein measurements

Each protein has its own detection reagent selected from chemically modified DNA libraries, referred to as Slow Off-rate Modified Aptamers (SOMAmers)⁴⁸. We designed an expanded custom version of the SOMApanel platform to include proteins known or predicted to be found in the extracellular milieu, including the predicted extracellular domains of single- and certain multi-pass transmembrane proteins¹⁰. The new aptamer-based platform measures 5034 protein analytes in a single serum sample, of which 4782 SOMAmers bind specifically to 4137 human proteins (some proteins are detected by more than one aptamer) and 250 SOMAmers that recognize non-human targets (47 non-human vertebrate proteins and 203 targeting human pathogens). Only human protein targets were analyzed in the present study. The levels of the 4782 peripheral proteins in 5457 serum samples from the AGES-RS were determined at SomaLogic Inc. (Boulder, US). Direct validation of 779 SOMAmers (mass spectrometry) and inferential validation (proximal genetic *cis* effects, biomarker and pQTL replication studies) across different study populations and proteomic technologies indicated consistent target specificity across the platform¹⁰. To avoid batch or time of processing biases, both sample collection and sample processing for protein measurements were randomized and all samples run as a single set¹⁰. The 5034 SOMAmers that passed quality

control had median intra-assay and inter-assay coefficient of variation, CV < 5%.

Hybridization controls were used to correct for systematic variability in detection and calibrator samples of three dilution sets (40%, 1% and 0.005%) were included so that the degree of fluorescence was a quantitative reflection of protein concentration.

Statistical analysis

Prior to the analysis of the proteins measurements, we applied a Box-Cox transformation on all proteins to improve normality, symmetry and to maintain all protein variables on a similar scale⁴⁹. In the association analysis, we obtained residuals after controlling for sex and age and for all single-variant associations to serum proteins tested under an additive genetic model applying linear regression analysis. We applied Bonferroni correction for multiple comparisons by adjusting for the 55,932 variants and 4137 human proteins targeted: single variant associations with $P < 2.16 \times 10^{-10}$ were considered significant. For the associations of individual proteins to different phenotypic measures we used linear or logistic regression or Cox proportional hazards regression, depending on the outcome being continuous, binary or a time to an event. Given consistency in terms of sample handling including time from blood draw to processing (between 9-11 am), same personnel handling all specimens and the ethnic homogeneity of the population we adjusted only for age and sex in all our regression analyses.

Figure Legends

Figure. 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 genomic locations of all significant pQTLs ($P < 2.16 \times 10^{-10}$), 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 appear at the diagonal while *trans* acting pQTL effects including *trans* hot spots (highlighted on top of the graph) show up off-diagonally.

Fig. 2. Pleiotropy of variants affecting many proteins and disease traits. a. The 5553 pQTLs were cross-referenced against risk loci for hundreds of diseases and relevant traits from the PhenoScanner²³ database. The database consists of disease and clinical trait data including that from the UKBB, GRASP and GWAS catalogue databases. A significant correlation between number of proteins affected and number of traits at shared GWAS and pQTL loci (Spearman rank correlation $\rho = 0.57$, $P < 0.0001$). Outliers above 150 are not shown on the y-axis. **b.** Circos plot showing the pleiotropy of the locus rs2251219 (Supplementary Tables 1 and 2) affecting 14 proteins in *cis* and *trans* and sharing genetics with various diseases of different etiologies. Lines going from rs225121 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 at $P < 2.16 \times 10^{-10}$ (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 locus 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.** The figure shows the Spearman rank correlation among the four serum proteins affected by the two distinct LOAD risk loci.

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 at genome- $P < 2.16 \times 10^{-10}$ (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 high 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 high levels of SVEP1 were associated with higher rates of mortality post-incident CHD.

Figure 5. a. The melanoma risk allele A (highlighted in bold) for the variant rs910873 is associated with high serum levels of ASIP. **b.** 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). The variant rs2241714 (and rs1800469) regulate three serum proteins, B3GNT2 (in *trans*), B3GNT8 (in *cis*) and TGFB1 (in *cis*).

References

- 1 Buniello, A. *et al.* The NHGRI-EBI GWAS Catalog of published genome-wide association studies, targeted arrays and summary statistics 2019. *Nucleic Acids Research* **47**, D1005-D1012, doi:10.1093/nar/gky1120 (2019).
- 2 Liu, D. J. *et al.* Exome-wide association study of plasma lipids in >300,000 individuals. *Nature Genetics* **49**, 1758-1766, doi:10.1038/ng.3977 (2017).
- 3 Schadt, E. E. Molecular networks as sensors and drivers of common human diseases. *Nature* **461**, 218-223, doi:10.1038/nature08454 (2009).
- 4 Zhang, B. *et al.* Integrated Systems Approach Identifies Genetic Nodes and Networks in Late-Onset Alzheimer's Disease. *Cell* **153**, 707-720, doi:10.1016/j.cell.2013.03.030 (2013).
- 5 Emilsson, V. *et al.* Genetics of gene expression and its effect on disease. *Nature* **452**, 423-U422, doi:10.1038/nature06758 (2008).
- 6 Chen, Y. Q. *et al.* Variations in DNA elucidate molecular networks that cause disease. *Nature* **452**, 429-435, doi:10.1038/nature06757 (2008).
- 7 Pires, D. E., Chen, J., Blundell, T. L. & Ascher, D. B. In silico functional dissection of saturation mutagenesis: Interpreting the relationship between phenotypes and changes in protein stability, interactions and activity. *Sci Rep* **6**, 19848, doi:10.1038/srep19848 (2016).
- 8 Ho, J. E. *et al.* Common genetic variation at the IL1RL1 locus regulates IL-33/ST2 signaling. *Journal of Clinical Investigation* **123**, 4208-4218, doi:10.1172/JCI67119 (2013).
- 9 Interleukin-6 Receptor Mendelian Randomisation Analysis, C. *et al.* The interleukin-6 receptor as a target for prevention of coronary heart disease: a mendelian randomisation analysis. *Lancet* **379**, 1214-1224, doi:10.1016/s0140-6736(12)60110-x (2012).
- 10 Emilsson, V. *et al.* Co-regulatory networks of human serum proteins link genetics to disease. *Science* **361**, 769-773, doi:10.1126/science.aag1327 (2018).
- 11 Sun, B. B. *et al.* Genomic atlas of the human plasma proteome. *Nature* **558**, 73-79, doi:10.1038/s41586-018-0175-2 (2018).
- 12 Mirauta, B. A. *et al.* Population-scale proteome variation in human induced pluripotent stem cells. *bioRxiv*, 439216, doi:10.1101/439216 (2018).
- 13 Emilsson, V., Gudnason, V. & Jennings, L. L. Predicting health and life span with the deep plasma proteome. *Nat Med* **25**, 1815-1816, doi:10.1038/s41591-019-0677-y (2019).
- 14 Lehallier, B. *et al.* Undulating changes in human plasma proteome profiles across the lifespan. *Nat Med* **25**, 1843-1850, doi:10.1038/s41591-019-0673-2 (2019).
- 15 Williams, S. A. *et al.* Plasma protein patterns as comprehensive indicators of health. *Nat Med* **25**, 1851-1857, doi:10.1038/s41591-019-0665-2 (2019).
- 16 Nakamura, A. *et al.* High performance plasma amyloid-beta biomarkers for Alzheimer's disease. *Nature* **554**, 249-254, doi:10.1038/nature25456 (2018).
- 17 Dodgson, S. E. There Will Be Blood Tests. *Cell* **173**, 1-3, doi:10.1016/j.cell.2018.03.012 (2018).
- 18 Cohen, J. D. *et al.* Detection and localization of surgically resectable cancers with a multi-analyte blood test. *Science* **359**, 926-930, doi:10.1126/science.aar3247 (2018).
- 19 Kristensen, S. L. *et al.* Prognostic Value of N-Terminal Pro-B-Type Natriuretic Peptide Levels in Heart Failure Patients With and Without Atrial Fibrillation. *Circ Heart Fail* **10**, doi:10.1161/circheartfailure.117.004409 (2017).

- 20 Peloso, G. M. *et al.* Association of low-frequency and rare coding-sequence variants with blood lipids and coronary heart disease in 56,000 whites and blacks. *American Journal of Human Genetics* **94**, 223-232, doi:10.1016/j.ajhg.2014.01.009 (2014).
- 21 Richards, A. L. *et al.* Exome arrays capture polygenic rare variant contributions to schizophrenia. *Human Molecular Genetics* **25**, 1001-1007, doi:10.1093/hmg/ddv620 (2016).
- 22 Armengaud, J., Christie-Oleza, J. A., Clair, G., Malard, V. & Duport, C. Exoproteomics: exploring the world around biological systems. *Expert Rev Proteomics* **9**, 561-575, doi:10.1586/epr.12.52 (2012).
- 23 Staley, J. R. *et al.* PhenoScanner: a database of human genotype-phenotype associations. *Bioinformatics* **32**, 3207-3209, doi:10.1093/bioinformatics/btw373 (2016).
- 24 Aguet, F. *et al.* The GTEx Consortium atlas of genetic regulatory effects across human tissues. *bioRxiv*, 787903, doi:10.1101/787903 (2019).
- 25 Jordan, D. M., Verbanck, M. & Do, R. The landscape of pervasive horizontal pleiotropy in human genetic variation is driven by extreme polygenicity of human traits and diseases. *bioRxiv*, 311332, doi:10.1101/311332 (2018).
- 26 Jonsson, T. *et al.* Variant of TREM2 associated with the risk of Alzheimer's disease. *N Engl J Med* **368**, 107-116, doi:10.1056/NEJMoa1211103 (2013).
- 27 Guo, C. *et al.* IGFBPL1 Regulates Axon Growth through IGF-1-mediated Signaling Cascades. *Sci Rep* **8**, 2054, doi:10.1038/s41598-018-20463-5 (2018).
- 28 Hollingworth, P. *et al.* Common variants at ABCA7, MS4A6A/MS4A4E, EPHA1, CD33 and CD2AP are associated with Alzheimer's disease. *Nature Genetics* **43**, 429-435, doi:10.1038/ng.803 (2011).
- 29 Suarez-Calvet, M. *et al.* sTREM2 cerebrospinal fluid levels are a potential biomarker for microglia activity in early-stage Alzheimer's disease and associate with neuronal injury markers. *EMBO Mol Med* **8**, 466-476, doi:10.15252/emmm.201506123 (2016).
- 30 Ewers, M. *et al.* Increased soluble TREM2 in cerebrospinal fluid is associated with reduced cognitive and clinical decline in Alzheimer's disease. *Sci Transl Med* **11** (2019).
- 31 Myocardial Infarction, G. *et al.* Coding Variation in ANGPTL4, LPL, and SVEP1 and the Risk of Coronary Disease. *N Engl J Med* **374**, 1134-1144, doi:10.1056/NEJMoa1507652 (2016).
- 32 Huyghe, J. R. *et al.* Exome array analysis identifies new loci and low-frequency variants influencing insulin processing and secretion. *Nature Genetics* **45**, 197-201, doi:10.1038/ng.2507 (2013).
- 33 Ransohoff, K. J. *et al.* Two-stage genome-wide association study identifies a novel susceptibility locus associated with melanoma. *Oncotarget* **8**, 17586-17592, doi:10.18632/oncotarget.15230 (2017).
- 34 Lu, Y. *et al.* Large-Scale Genome-Wide Association Study of East Asians Identifies Loci Associated With Risk for Colorectal Cancer. *Gastroenterology*, doi:10.1053/j.gastro.2018.11.066 (2018).
- 35 Brown, K. M. *et al.* Common sequence variants on 20q11.22 confer melanoma susceptibility. *Nature Genetics* **40**, 838-840, doi:10.1038/ng.163 (2008).
- 36 Blanchard, S. G. *et al.* Agouti antagonism of melanocortin binding and action in the B16F10 murine melanoma cell line. *Biochemistry* **34**, 10406-10411 (1995).
- 37 Taylor, N. J. *et al.* Inherited variation at MC1R and ASIP and association with melanoma-specific survival. *Int J Cancer* **136**, 2659-2667, doi:10.1002/ijc.29317 (2015).

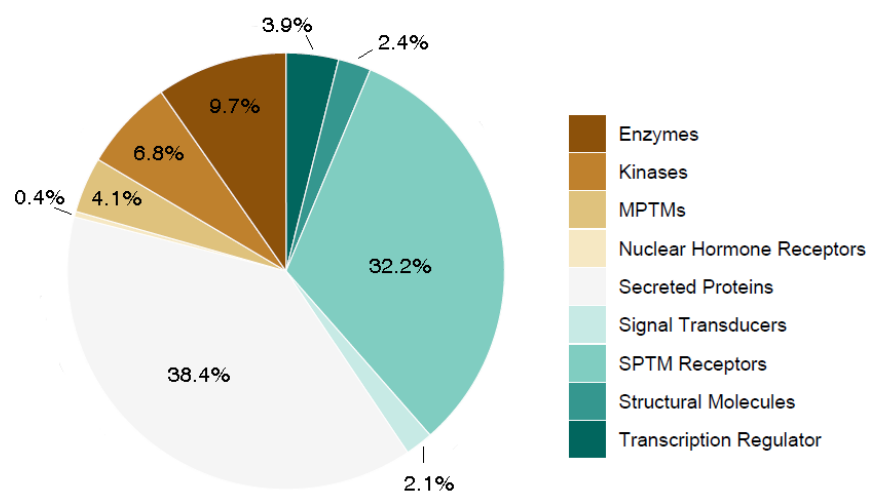
- 38 Wolf Horrell, E. M., Boulanger, M. C. & D'Orazio, J. A. Melanocortin 1 Receptor: Structure, Function, and Regulation. *Front Genet* **7**, 95, doi:10.3389/fgene.2016.00095 (2016).
- 39 Zhang, B. *et al.* Large-scale genetic study in East Asians identifies six new loci associated with colorectal cancer risk. *Nature Genetics* **46**, 533-542, doi:10.1038/ng.2985 (2014).
- 40 Calon, A. *et al.* Dependency of colorectal cancer on a TGF-beta-driven program in stromal cells for metastasis initiation. *Cancer Cell* **22**, 571-584, doi:10.1016/j.ccr.2012.08.013 (2012).
- 41 Venkitachalam, S. *et al.* Biochemical and functional characterization of glycosylation-associated mutational landscapes in colon cancer. *Sci Rep* **6**, 23642, doi:10.1038/srep23642 (2016).
- 42 Ishida, H. *et al.* A novel beta1,3-N-acetylglucosaminyltransferase (beta3Gn-T8), which synthesizes poly-N-acetyllactosamine, is dramatically upregulated in colon cancer. *Febs Letters* **579**, 71-78, doi:10.1016/j.febslet.2004.11.037 (2005).
- 43 Solomon, T. *et al.* Identification of Common and Rare Genetic Variation Associated With Plasma Protein Levels Using Whole-Exome Sequencing and Mass Spectrometry. *Circ Genom Precis Med* **11**, e002170, doi:10.1161/circgen.118.002170 (2018).
- 44 Smith, J. G. & Gerszten, R. E. Emerging Affinity-Based Proteomic Technologies for Large-Scale Plasma Profiling in Cardiovascular Disease. *Circulation* **135**, 1651-1664, doi:10.1161/circulationaha.116.025446 (2017).
- 45 Zheng, J. *et al.* Phenome-wide Mendelian randomization mapping the influence of the plasma proteome on complex diseases. *bioRxiv*, 627398, doi:10.1101/627398 (2019).
- 46 Liu, B., Gloudemans, M. J., Rao, A. S., Ingelsson, E. & Montgomery, S. B. Abundant associations with gene expression complicate GWAS follow-up. *Nature Genetics* **51**, 768-769, doi:10.1038/s41588-019-0404-0 (2019).
- 47 Harris, T. B. *et al.* Age, Gene/Environment Susceptibility-Reykjavik Study: multidisciplinary applied phenomics. *Am J Epidemiol* **165**, 1076-1087, doi:10.1093/aje/kwk115 (2007).
- 48 Candia, J. *et al.* Assessment of Variability in the SOMAscan Assay. *Sci Rep* **7**, 14248, doi:10.1038/s41598-017-14755-5 (2017).
- 49 Max Kuhn, K. J. *Applied Predictive Modeling*. (Springer, 2013).
- 50 MacArthur, J. *et al.* The new NHGRI-EBI Catalog of published genome-wide association studies (GWAS Catalog). *Nucleic Acids Research* **45**, D896-D901, doi:10.1093/nar/gkw1133 (2017).

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 $P < 2.16 \times 10^{-10}$.

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	4	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 Catalog⁵⁰. ^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 2.

a



b

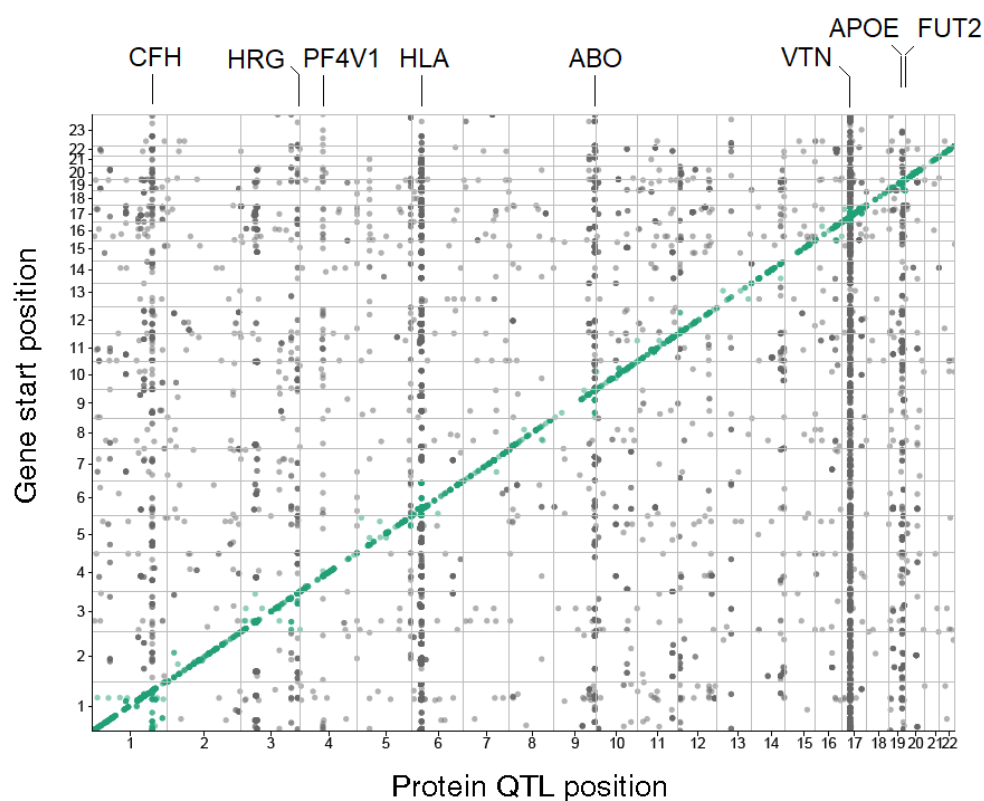


Figure 1

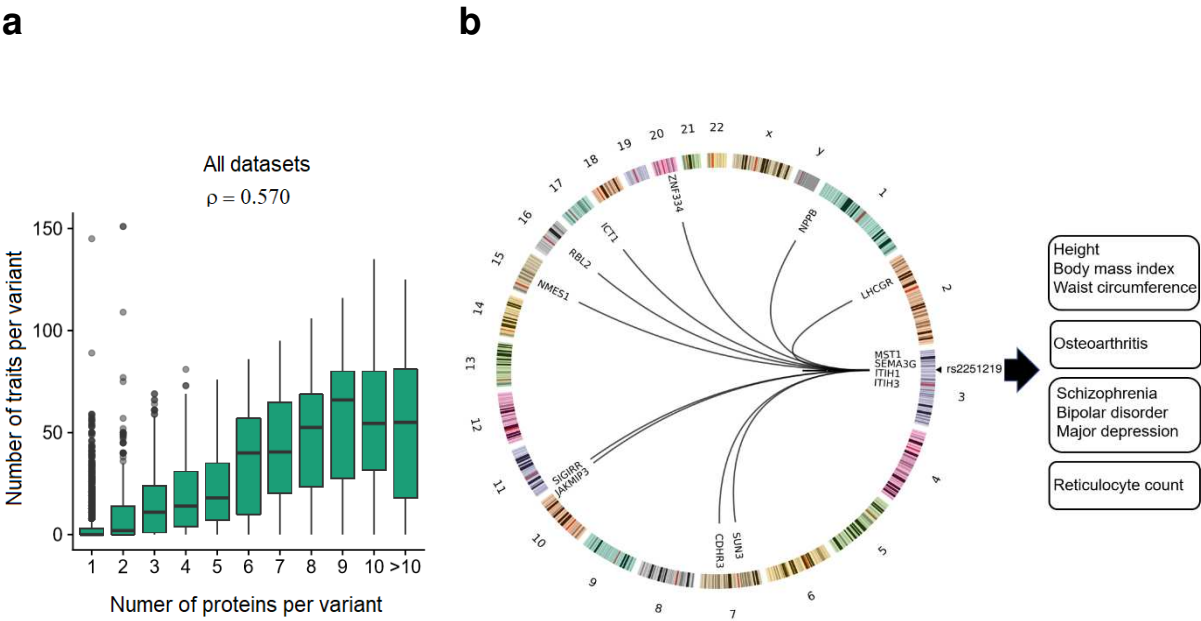


Figure 2

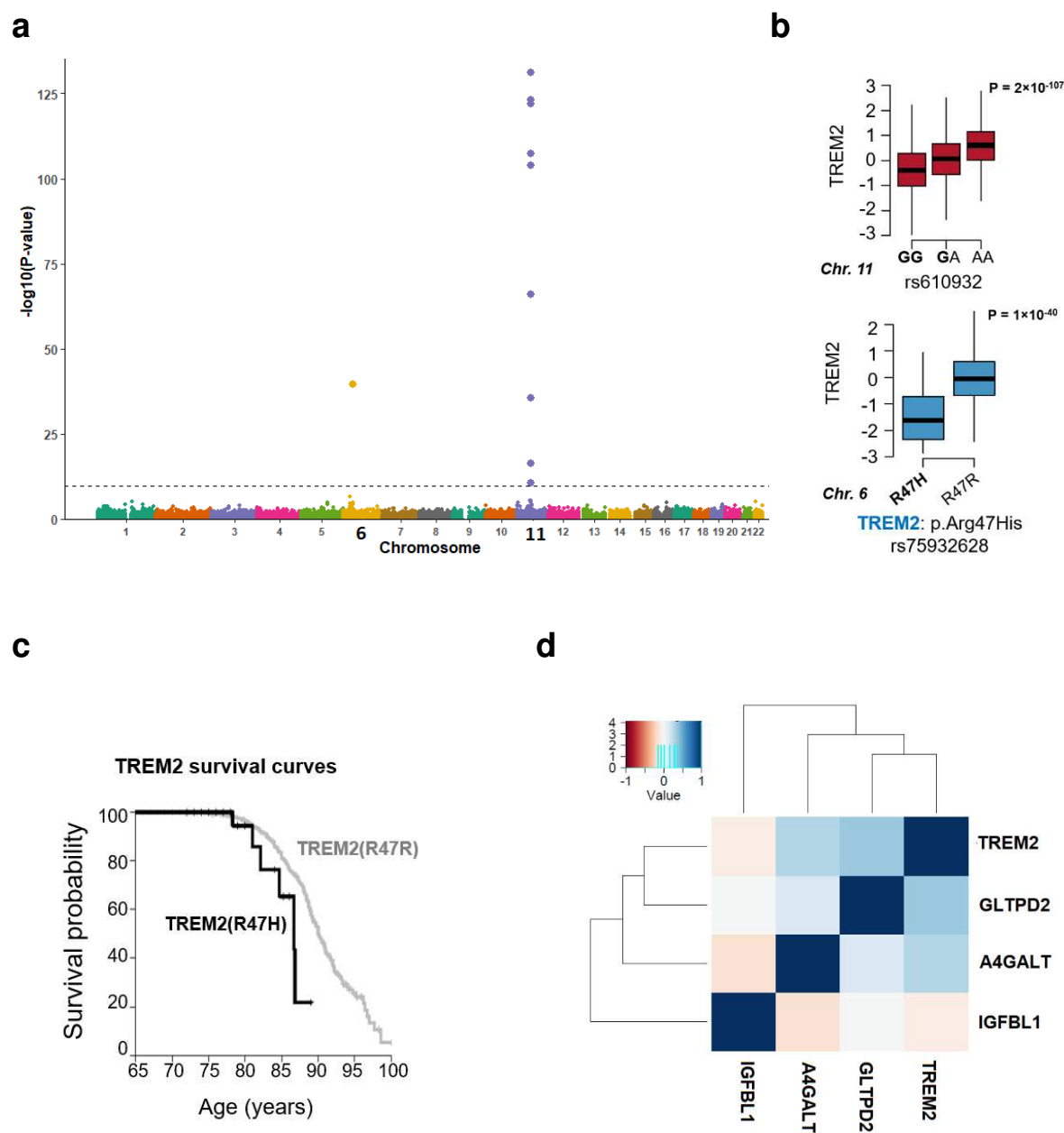


Figure 3

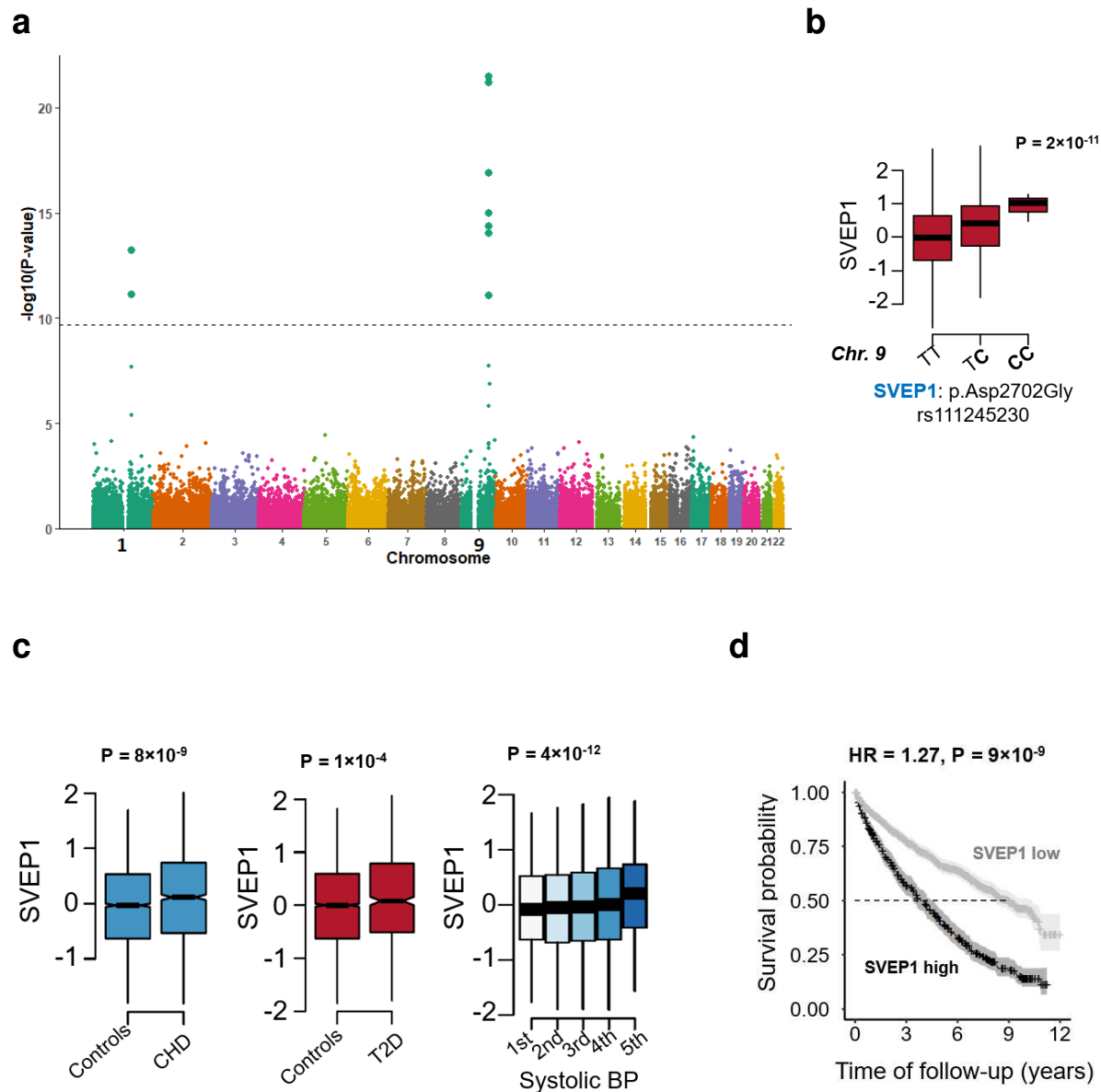


Figure 4

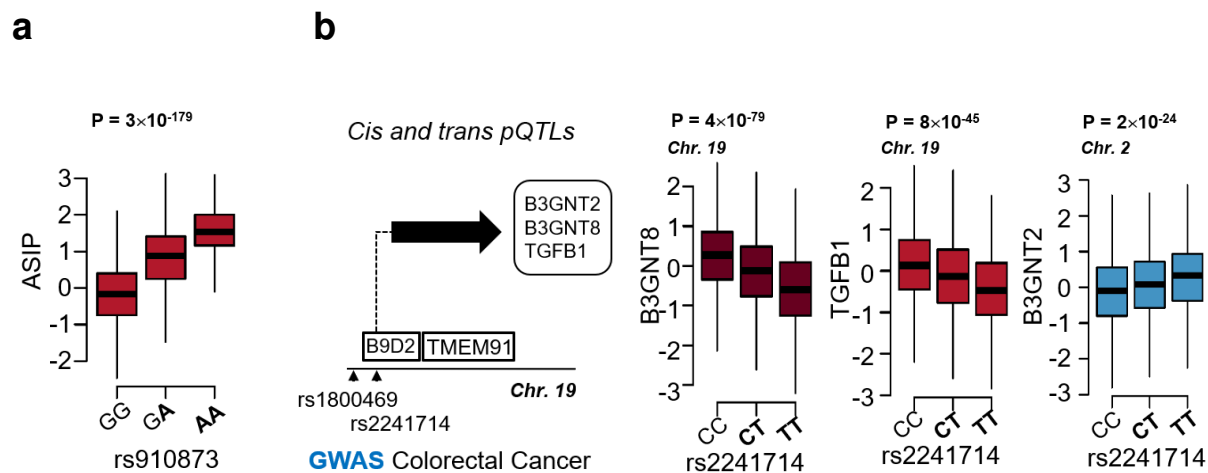


Figure 5