

## Supplementary Data Legends

### Supplementary Data 1. a: Study-wide significant cis-pQTL. b: Study-wide significant trans pQTL.

column	description
<b>panel</b>	OLINK protein panel on which the protein was assayed
<b>protein</b>	protein name as provided by OLINK
<b>lead SNP</b>	index SNP defining a 2Mb locus
<b>independent SNP</b>	independently associated SNP at locus, as defined by COJO
<b>novelty_status</b>	does the p-value of the independent SNP conditioned on all previous signals pass the study-wide significance threshold of $7.45e-11$
<b>reference MANOLIS</b>	reference allele GRCh38 in discovery
<b>alternate MANOLIS</b>	alternate allele GRCh38 in discovery
<b>AF MANOLIS</b>	alternate allele frequency in discovery
<b>effect allele MANOLIS</b>	effect allele
<b>MAF MANOLIS</b>	minor allele frequency in discovery
<b>beta MANOLIS</b>	discovery effect size
<b>se MANOLIS</b>	discovery standard error of the effect size
<b>p MANOLIS</b>	score test discovery p-value
<b>reference Pomak</b>	reference allele GRCh38 in replication
<b>alternate Pomak</b>	alternate allele GRCh38 in replication
<b>AF Pomak</b>	allele frequency in replication
<b>effect allele Pomak</b>	effect allele in replication
<b>MAF Pomak</b>	minor allele frequency in replication
<b>beta Pomak</b>	effect size in replication
<b>se Pomak</b>	standard error in replication
<b>p Pomak</b>	score test replication p-value
<b>proxy if not present in Pomak</b>	tagging variant ( $r^2 > 0.8$ ) used for meta-analysis if independent SNP not present in replication cohort
<b>beta</b>	meta-analysis effect size
<b>se</b>	meta-analysis standard error
<b>p</b>	inverse variance based meta-analysis p-value
<b>rsids</b>	Rsids of the discovery variant if present
<b>replicates</b>	does the meta-analysis p-value pass the Bonferroni-corrected threshold of nominal significance
<b>cis consequence</b>	consequence on the <i>cis</i> gene as per Ensembl VEP

<b>cis gene ensembl id</b>	Ensembl ID of the <i>cis</i> gene
<b>distance to cis (bp)</b>	distance to the above, as per Ensembl REST API
<b>cis gene name</b>	Gene symbol, as per REST API
<b>uniprot id</b>	UniProt protein name
<b>uniprot short name</b>	UniProt short name
<b>uniprot long name</b>	UniProt long name
<b>P conditional</b>	maximum conditional score test p-value across all previously associated variants in the region for the same protein. "same variant" indicates the variant itself is known.
<b>not studied before</b>	literature search and GWAS catalog did not produce evidence of previous genetic associations for this protein
<b>novel locus</b>	have there never been any reported associations for this protein at this locus
<b>mapped gene consequence</b>	for <i>trans</i> loci, consequence of the discovery variant on the mapped gene
<b>mapped gene ID</b>	for <i>trans</i> loci, Ensembl ID of the mapped gene
<b>mapped gene distance (bp)</b>	for <i>trans</i> loci, distance to the mapped gene as per REST API
<b>mapped gene</b>	for <i>trans</i> loci, gene that the signal has been mapped to. Details on how this is done are provided in the Methods section of the paper.
<b>receptor/ligand</b>	does the association reflect a documented receptor/ligand pair

**Supplementary Data 2: Allele frequencies in MANOLIS compared to reference European populations.** Alternate allele counts according to GRCh38 are compared to an absolute population alternate frequency. If gnomAD non-Finnish European counts are available for exomes, those are used for comparison. If they are not, gnomAD non-Finnish European genomes are used. If the variant is not present in gnomAD, TOPMed frequencies are used as reference.

**Supplementary Data 3. a: Phenotypes used for 2-sample MR in MRBase.** All the information displayed in this table is a filtered output from the MRBase `get_available_outcomes` function. "id" is the MRBase ID. **b: Phenotypes used for 2-sample MR using manually downloaded statistics.** For every study, we display the file name used, sample sizes broken down by cases and controls, a description of the trait and a PMID for the supporting paper.

**Supplementary Data 4: 2-Sample MR results.** Exposures are reported according to their OLINK panel name and OLINK protein name separated by a dot. The method is Wald ratio for single SNP instruments and Inverse-variance weighted for multi-SNV instruments. The FDR-corrected p-value is displayed in the pBH column, which is filtered for significant results (pBH<0.05). Outcome names are as given by MR-Base for that trait, or equal to the filenames used in case of manually downloaded summary statistics, as described in Supplementary Data 3b.

**Supplementary Data 5: Directions of effect for outcomes and instruments in Mendelian randomization analysis. In case a pQTL locus included in MR involved variants with both positive and negative effects, their direction is denoted as "mixed".**

**Supplementary Data 6. a: Effect of protein level genetic scores in logistic disease prediction models in UK Biobank, for most significantly associated protein scores. Effect: log-odds change per score unit increase  $\beta$  : effect size,  $\sigma$ : standard error. Wald P: P-value for the Wald test of predictor significance. Threshold: P threshold used for constructing the score. SR: Self-reported. ICD: International Classification of Diseases, Tenth Revision. b: Logistic regression models of disease states on protein PRS and clinical factors. When several scores for the same protein pass the Bonferroni-corrected p-value threshold, only the model with the strongest evidence of score effects across all p-value thresholds is reported. Only models where a protein score passes the significance threshold of  $0.05/(47*86)$  are reported here.**

**Supplementary Data 7: Proteins measured in this study, and reasons for exclusion if applicable. The symbol column denotes the internal name provided by the vendor (Olink) used throughout this study.**

**Supplementary Data 8: Two-sample MR results for the PDL2 gene association, when only the cis locus is included as an instrument. Rows are sorted by ascending p-value. P-values are reported for the test specified in the "method" column (one-sided). Highlighted rows pass the FDR-adjusted significance threshold.**

**Supplementary Data 9: Variance explained by MANOLIS PRS in the Pomak cohort as a function of MAF and P-value thresholds for inclusion**

**Supplementary Data 10. a: Optimal score selection and R-squared for protein level prediction in Pomak, using summary statistics in MANOLIS using the full dataset. This analysis is a repeat of that presented in ST12 using a more recent version of PRSice (2.3.1.e) and is presented for comparison purposes only. b: Optimal score selection and R-squared for protein level prediction in Pomak, using summary statistics in MANOLIS using repeated cross-validation (repeats:3, folds:5). c: Optimal score selection and R-squared for protein level prediction in Pomak, using summary statistics in MANOLIS, using the full dataset, with a 2Mb region around ABO excluded from the analysis. Only the four proteins found to be correlated with each other and with evidence of a SNV in ABO driving the top score are considered here.**

**Supplementary Data 11. a: Non-cancer illness codes, with description and number of cases, used for polygenic prediction in UK Biobank. b: ICD10 codes (primary and secondary) with more than 999 cases used for polygenic prediction in UK Biobank.**

**Supplementary Data 12: Association between protein scores and PheCodes. The table, which reports all coefficients, is filtered for significant protein scores: the variable column should be a protein and the P-value column should be smaller than  $0.05/(43*47)$ , correcting for the number of proteins and PheCodes tested. Suffixes in the variable column correspond to P-value thresholds for building scores, 1 corresponds to the most stringent, at  $7.45 \times 10^{-11}$ , with a step of  $1 \times 10^{-10}$ . The null model is the model without any protein scores.**

**Supplementary Data 13: List of GWAS individual accession numbers used in this study.**

**Supplementary Data 14: Variants included in the burdens of the selected associations. Alleles: reference/non-reference alleles; Consequence: the most severe consequence predicted by Ensembl VEP; MAF: minor allele frequency; AC: non-reference allele count; Weight: the weight assigned to the variant in the most significant test as per Supplementary Data 15; NA indicates that no weight was applied.**

**Supplementary Data 15: Non-proteomic phenotypes used for comparison in the MANOLIS cohort. "rbin" indicates inverse normal transformation, ln indicates natural logarithm. The last three columns indicate whether the corresponding covariate was adjusted for.**

**Supplementary Data 16: PheWAS results for all independent variants reported in this study. P4 is the posterior probability of a single colocalising signal. Alpha12 and alpha21 are the products of the effect sizes for the protein and the tested traits, n is the number of variants input in the colocalisation test.**

**Supplementary Data 17: Extract of the DrugBank database where target genes correspond to cis genes for proteins measured in the current study**

**Supplementary Data 18: Mouse orthologs and their phenotype associations as provided by MGI for genes coding for the proteins included in our study. P-values are hard-filtered using Bonferroni correction for the number of orthologs tested.**