

1 **Editor summary:**

2 Here the authors identify genetic effectors of the level of inflammation-related plasma proteins and
3 use Mendelian randomization to identify proteins that contribute to immune-mediated disease risk.

4 **Peer Review Information:**

5 Nature Immunology thanks Tom Richardson and the other, anonymous, reviewer(s) for their contri-
6 bution to the peer review of this work.

7 **Editor recognition statement:**

8 N. Bernard was the primary editor on this article and managed its editorial process and peer review
9 in collaboration with the rest of the editorial team.

10

11

1. Extended Data

Figure or Table # Please group Extended Data items by type, in sequential order. Total number of items (Figs. + Tables) must not exceed 10.	Figure/Table title One sentence only	Filename Whole original file name including extension. i.e.: Smith_ED_Fig1.jpg	Figure/Table Legend If you are citing a reference for the first time in these legends, please include all new references in the main text Methods References section, and carry on the numbering from the main References section of the paper. If your paper does not have a Methods section, include all new references at the end of the main Reference list.
Extended Data Fig. 1	Overview of the pQTL analysis	EDF 1.eps	Schematic of the analysis pipeline.
Extended Data Fig. 2	Plasma protein abundance and pQTL detection	EDF 2.eps	a) Proteins with low abundance are more likely to have no detectable pQTL. Y-axis: percentage of samples above lower limit of detection for each protein, calculated using the INTERVAL data (n=4,896) for which we had individual-level protein data available. Blue and red points indicate presence or absence of at least 1 significant pQTL in the GWAS meta-analysis, respectively. b) Manhattan plot for genetic associations with plasma IL17C, where the red horizontal line indicates the statistical significance threshold (5×10^{-10}). P-values from linear regression.
Extended Data Fig. 3	pQTL replication in the ARISTOTLE cohort	EDF 3.eps	Comparison of effect sizes between pQTLs from the discovery pQTL meta-analysis (n=14,824) and the ARISTOTLE cohort (n=1,585). Each point represents a genetic variant that was a significant pQTL in the discovery meta-analysis. Effect size = standard deviation (sd) increase in protein per allele.

			174 of 180 genetic variants were available for testing in the ARISTOTLE data. Red= cis, Blue= trans.
Extended Data Fig. 4	Genetic architecture of circulating inflammation-related proteins	EDF 4.eps	Genetic architecture of circulating inflammation-related proteins. a) Relationship between minor allele frequency (MAF), pQTL effect size and proportion of variance explained ($2MAF(1-MAF)Effect^2$), for 227 conditionally independent pQTLs (red=cis, blue=trans). b) Proportion of variance explained (PVE) by the conditionally independent variants associated with each protein. Proteins are annotated using the gene symbol of their encoding genes. Protein names are coloured in red if over 80% of samples have levels below the lower limit of detection in the INTERVAL dataset.
Extended Data Fig. 5	Chemokine <i>trans</i>-pQTL hotspot	EDF 5.eps	Forest plot showing the associations for the pleiotropic <i>trans</i> -pQTL at rs12075 (GRCh37, 1:158175353-160525679) with plasma levels of chemokines and blood cell counts. Centre of bar = effect size estimate, whiskers = 95% confidence interval (ci). WBC = white blood cell count. P = p-value, b= beta (effect size). SE = standard error. Blood cell association data from Chen <i>et al.</i> ¹ . P-values from linear regression.
Extended Data Fig. 6	Colocalisation of pleiotropic chemokine <i>trans</i>-pQTL and blood cell count trait signals	EDF 6.eps	Regional association plots in the region around rs12075 (GRCh37, 1:158175353-160525679). Left side: association with plasma chemokine levels. Right: associations with basophil, monocyte and white blood cell (WBC) counts using data from Chen <i>et al.</i> ¹ . P-values from linear regression.
Extended Data Fig. 7	Interactions between the candidate mediators for multi-locus-regulated proteins	EDF 7.jpg	a) TNFSF10 (also known as TRAIL), b) KITLG (also known as stem cell factor), and c) IL12B. The graphs were generated using the STRINGdb (v11.5) webtool. The colouring of the edges indicates the type of evidence supporting an interaction, as shown in the legend above.
Extended Data Fig. 8	Protein-disease connections from overlap of pQTLs and disease GWAS	EDF 8.eps	The protein and the corresponding pQTL sentinel variant are indicated in the format of protein-rsid. The nearest gene to the pQTL sentinel variant is shown in brackets. Red lettering= cis-pQTL, blue lettering= <i>trans</i> -pQTL. Asterix indicates the genetic variant lies

			in the <i>HLA</i> region. Red squares: genetic susceptibility to increased plasma levels of the protein is associated with increased disease risk. Blue squares: decreased disease risk.
Extended Data Fig. 9	Protein and immune-mediated disease (IMD) connections from overlap of pQTLs and disease GWAS	EDF 9.eps	The protein and the corresponding pQTL sentinel variant are indicated in the format of protein-rsid. The nearest gene to the pQTL sentinel variant is shown in brackets. Red lettering= cis-pQTL, blue lettering= trans-pQTL. Asterix indicates the genetic variant lies in the <i>HLA</i> region. Red squares: genetic susceptibility to increased plasma levels of the protein is associated with increased disease risk. Blue squares: decreased disease risk.
Extended Data Fig. 10	Mendelian randomisation analysis for CXCL5 and ulcerative colitis	EDF 10.eps	a) Scatterplot showing the 13 variants used in the GSMR analysis assessing the effect of CXCL5 on ulcerative colitis (UC) risk from the GWAS by de Lange <i>et al</i> (ref 49) Each point represents a genetic variant, and indicates the effect size of the variant on CXCL5 levels versus UC risk (log odds ratio). Vertical and horizontal lines represent 95% confidence intervals. b) Directional concordance between CXCL5 pQTL and blood and colon tissue eQTLs. Forest plots showing effect size estimates for rs450373 pQTL in plasma (from our discovery meta-analysis) and eQTLs in whole blood and transverse colon tissue (GTEx v8 data). OR= odds ratio, calculated from beta estimate (representing the change in inverse-rank normalised plasma protein level in standard deviations associated with each copy of the effect allele). CI = confidence interval. P = p-value. Centre of bar = OR estimate, whiskers = 95% CI.

12

13 **2. Supplementary Information:**

14 **A. PDF Files**

15

Item	Present?	Filename	A brief, numerical description of file contents.
		Whole original file name including extension. i.e.: Smith_SI.pdf. The	i.e.: <i>Supplementary Figures 1-4, Supplementary Discussion, and Supplementary Tables 1-4.</i>

		extension must be .pdf	
Supplementary Information	Yes	Supplementary-merged.pdf	Supplementary Note Supplementary Figures 1-4
Reporting Summary	Yes	NI-A35592B_Peters_RSf.pdf	
Peer Review Information	Yes	NI-A35592B_Peters_TPR.pdf	

16
17
18

B. Additional Supplementary Files

Type	Number Each type of file (Table, Video, etc.) should be numbered from 1 onwards. Multiple files of the same type should be listed in sequence, i.e.: Supplementary Video 1, Supplementary Video 2, etc.	Filename Whole original file name including extension. i.e.: <i>Smith_Supplementary_Video_1.mov</i>	Legend or Descriptive Caption Describe the contents of the file
Supplementary Video	Supplementary Item 1	SI.html	3-dimensional interactive genomic map of pQTLs. The html file (supplied separately) shows pQTL sentinel variant position in relation to the gene encoding the target protein and the strength of the statistical association (two-sided <i>P</i> -values are from meta-analysis of linear regression estimates). Hover over a point to see detailed information. The image can be rotated by holding at left clicking the mouse.
Supplementary Table	Supplementary Table 1	Supplementary-Tables.xlsx	Excel file containing multiple tabs (1 summary tab, then 1 tab per Supplementary Table).

19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86

Genetics of circulating inflammatory proteins identifies drivers of immune-mediated disease risk and therapeutic targets

Jing Hua Zhao^{1,2,3,7}, David Stacey^{1,2,3,4,37}, Niclas Eriksson⁵, Erin Macdonald-Dunlop⁶, Åsa K Hedman^{7,8}, Anette Kalnapenkis^{9,1}, Stefan Enroth¹⁰, Domenico Cozzetto¹¹, Jonathan Digby-Bell¹², Jonathan Marten¹, Lasse Folkersen¹³, Christian Herder^{14,15,16}, Lina Jonsson¹⁷, Sarah E Bergen⁷, Christian Geiger^{18,19}, Elise J Needham^{1,2}, Praveen Surendran^{1,20,21}, Estonian Biobank Research Team⁹, Dirk S Paul^{1,20,2,22}, Ozren Polasek²³, Barbara Thorand^{16,18}, Harald Grallert^{16,18,19}, Michael Roden^{14,24,16}, Urmo Vösa⁹, Tonu Esko⁹, Caroline Hayward²⁵, Åsa Johansson¹⁰, Ulf Gyllensten¹⁰, Nicholas Powell¹¹, Oskar Hansson^{26,27}, Niklas Mattsson-Carlgren^{28,29,30}, Peter K Joshi⁶, John Danesh^{1,2,20,21,31,32}, Leonid Padyukov^{33,34}, Lars Klareskog^{33,34}, Mikael Landén^{17,7}, James F Wilson^{6,25}, Agneta Siegbahn³⁵, Lars Wallentin³⁵, Anders Mälarstig^{7,8}, Adam S Butterworth^{1,2,20,21,31,38,*}, James E Peters^{1,21,36,38,*}

¹British Heart Foundation Cardiovascular Epidemiology Unit, Department of Public Health and Primary Care, University of Cambridge, Cambridge, United Kingdom

²Victor Phillip Dahdaleh Heart and Lung Research Institute, University of Cambridge, Papworth Road, Cambridge, United Kingdom

³Australian Centre for Precision Health, Unit of Clinical and Health Sciences, University of South Australia, Adelaide, South Australia, Australia

⁴South Australian Health and Medical Research Institute, Adelaide, South Australia, Australia

⁵Uppsala Clinical Research Center, Uppsala University, Uppsala, Sweden

⁶Centre for Global Health Research, Usher Institute, University of Edinburgh, Teviot Place, Edinburgh, United Kingdom

⁷Department of Medical Epidemiology and Biostatistics, Karolinska Institutet, Stockholm, Sweden

⁸Pfizer Worldwide Research, Development and Medical, Stockholm, Sweden

⁹Estonian Genome Center, Institute of Genomics, University of Tartu, Tartu, Estonia

¹⁰Department of Immunology, Genetics, and Pathology, Biomedical Center, SciLifeLab Uppsala, Uppsala University, Uppsala, Sweden

¹¹Department of Metabolism, Digestion and Reproduction, Faculty of Medicine, Imperial College London, London, United Kingdom

¹²School of Immunology and Microbial Sciences, King's College London, London, United Kingdom

¹³Nucleus Genomics Ltd, New York, NY, USA

¹⁴Institute for Clinical Diabetology, German Diabetes Center, Leibniz Center for Diabetes Research at Heinrich Heine

University Düsseldorf, Düsseldorf, Germany

¹⁵Department of Endocrinology and Diabetology, Medical Faculty and University Hospital Düsseldorf, Heinrich Heine University Düsseldorf, Düsseldorf, Germany

¹⁶German Center for Diabetes Research (DZD), Munich-Neuherberg, Germany

¹⁷Institute of Neuroscience and Physiology, University of Gothenburg, Gothenburg, Sweden

¹⁸Institute of Epidemiology, Helmholtz Zentrum München, German Research Center for Environmental Health, Neuherberg, Germany

¹⁹Research Unit of Molecular Epidemiology, Helmholtz Zentrum München, German Research Center for Environmental Health, Neuherberg, Germany

²⁰British Heart Foundation Centre of Research Excellence, School of Clinical Medicine, Addenbrooke's Hospital, University of Cambridge, Cambridge, United Kingdom

²¹Health Data Research UK, Wellcome Genome Campus and University of Cambridge, Hinxton, United Kingdom

²²Centre for Genomics Research, Discovery Sciences, BioPharmaceuticals R&D, AstraZeneca, Cambridge, United Kingdom

²³University of Split Medical School, University of Split, Split, Croatia

²⁴Division of Endocrinology and Diabetology, Medical Faculty, Heinrich Heine University Düsseldorf, Düsseldorf, Germany

²⁵MRC Human Genetics Unit, Institute of Genetics and Cancer, University of Edinburgh, Western General Hospital, Crewe Road, Edinburgh, United Kingdom

²⁶Clinical Memory Research Unit, Department of Clinical Sciences Malmö, Lund University, Lund, Sweden

²⁷Skåne University Hospital, Malmö, Sweden

²⁸Wallenberg Centre for Molecular Medicine, Lund University, Lund, Sweden

²⁹Clinical Memory Research Unit, Faculty of Medicine, Lund University, Lund, Sweden

³⁰Department of Neurology, Skåne University Hospital, Lund University, Lund, Sweden

³¹NIHR Blood and Transplant Research Unit in Donor Health and Behaviour, University of Cambridge, Cambridge, United Kingdom

³²Department of Human Genetics, Wellcome Sanger Institute, Wellcome Genome Campus, Hinxton, United Kingdom

³³Division of Rheumatology, Department of Medicine (Solna), Karolinska Institutet and Karolinska University Hospital, Stockholm, Sweden

³⁴Center for Molecular Medicine, Karolinska Institutet, Stockholm, Sweden

³⁵Department of Medical Sciences and Uppsala Clinical Research Center, Uppsala University, Uppsala, Sweden

³⁶Department of Immunology and Inflammation, Imperial College London, London, United Kingdom

³⁷These authors contributed equally

³⁸These authors jointly supervised the study

*e-mail: Adam S Butterworth asb38@medschl.cam.ac.uk or James E Peters j.peters@imperial.ac.uk

87

88 **ABSTRACT**

89 Circulating proteins have important functions in inflammation and a broad range of diseases.
90 To identify genetic influences on inflammation-related proteins, we conducted a genome-wide
91 protein quantitative trait locus (pQTL) study of 91 plasma proteins measured using the Olink
92 Target platform in 14,824 participants. We identify 180 pQTLs (59 *cis*, 121 *trans*). Integration
93 of pQTL data with eQTL and disease GWAS provided insight into pathogenesis, implicating
94 lymphotoxin- α in multiple sclerosis. Using Mendelian randomisation (MR) to assess causality
95 in aetiology, we identify both shared and distinct effects of specific proteins across immune-
96 mediated diseases, including directionally discordant functions for CD40 in rheumatoid
97 arthritis versus multiple sclerosis and inflammatory bowel disease. MR implicated CXCL5 in
98 the aetiology of ulcerative colitis (UC) and we show elevated gut *CXCL5* transcript expression
99 in patients with UC. These results identify targets for existing drugs and provide a powerful
100 resource to facilitate future drug target prioritisation.

101

102 **INTRODUCTION**

103

104 Inflammation is a physiological host response to infection or injury. However, aberrant
105 inflammatory responses result in tissue damage and are central to the pathogenesis of
106 multiple diseases including sepsis, autoimmunity and atherothrombosis. Inflammatory
107 responses are orchestrated by a complex network of cells and mediators, including circulating
108 proteins such as cytokines and soluble receptors. Therefore, discovery of the genetic
109 determinants of abundance of inflammation-related circulating proteins should yield valuable
110 insights into both physiology and the aetiology of a broad range of diseases.

111

112 Proteomic studies are informative as proteins are the effector molecules of most biological
113 processes, and from a translational perspective, proteins are the targets of most drugs. The
114 development of high-throughput proteomic technologies now allows for profiling of the plasma
115 proteome at epidemiological scale. Coupling genomic and proteomic data enables
116 identification of genetic variants associated with protein abundance, protein quantitative trait
117 loci (pQTLs). pQTLs provide valuable insights into the molecular basis of complex traits and
118 diseases, by identifying proteins that lie between genotype and phenotype. Recent years have
119 seen a rapid increase in both the number and size of pQTL studies, transforming our
120 understanding of the genetic architecture of the circulating proteome¹⁻¹¹.

121

122 Here we extend previous work by performing pQTL mapping for 91 inflammation-related

123 proteins in 14,824 participants. We integrate these data with disease genome-wide
124 association studies (GWAS) to characterise the functional effects of disease-associated
125 variants. Using Mendelian randomisation and colocalisation analyses, we identify proteins that
126 play a causal role in immune-mediated disease aetiology. Our results reveal both pathways
127 that are known to be therapeutically important and new putative drug targets, including CD40
128 in rheumatoid arthritis, lymphotoxin- α (LTA) in multiple sclerosis, and the chemokine CXCL5
129 in ulcerative colitis.

130

131 RESULTS

132

133 **Genetic architecture of circulating inflammatory proteins.** We performed genome-wide
134 pQTL mapping for 91 plasma proteins measured using the Olink Target Inflammation panel in
135 11 cohorts totalling 14,824 European-ancestry participants (**Supplementary Table 1,**
136 **Supplementary Note 1**), and meta-analysed the results (**Extended Data Figure 1**). In order
137 to provide a succinct and standardised nomenclature, we report proteins by the non-italicised
138 symbols of the genes encoding them (see **Supplementary Table 2** for a mapping of symbols
139 to full protein names and UniProt identifiers). We identified a total of 180 significant ($P \leq 5 \times 10^{-10}$,
140 fixed-effects meta-analysis) associations between 108 genomic regions (see **Methods** for
141 locus definition) and 70 proteins (**Figure 1, Supplementary Table 3, Supplementary Item,**
142 **Supplementary Figures 1-2**). Of the 180 significant locus-protein associations, 59 (33%)
143 were local-acting ('*cis*' pQTLs; defined here as a genetic variant lying within +/- 1 megabase
144 of the gene encoding the associated protein) and 121 (67%) were distant-acting ('*trans*'). We
145 found evidence of *trans*-pQTL hotspots associated with multiple proteins (e.g. rs3184504 at
146 the *SH2B3* locus was associated with CXCL9, CXCL10, CXCL11, CD5, CD244, and IL12B)
147 (**Figure 2a**).

148

149 For 70 (77%) of the 91 proteins studied, we identified at least 1 significant pQTL, including 59
150 (65%) proteins that had a *cis*-pQTL. Of these 70 proteins, 19 had only *cis*-pQTL(s), 11 had
151 only *trans*-pQTL(s), and 40 had both *cis*- and *trans*-pQTLs. For 18 of the 21 proteins for which
152 no pQTL was detected, >50% of samples had levels below the lower limit of detection (LLOD),
153 suggesting that the lack of genetic signal is due to low protein abundance in plasma
154 (**Extended Data Figure 2a**). The number of genomic loci associated with each protein ranged
155 between 1 and 8 (**Figure 2b**), but was less than 4 for the majority of proteins. Examples of
156 multi-locus-regulated proteins include TNFSF10 and IL12B, both of which had 1 *cis*- and 7
157 *trans*-pQTLs (**Figure 2c,d**). Conditional analyses revealed the presence of an additional 47
158 independent signals, raising the total number of pQTL signals from 180 (59 *cis*, 108 *trans*) to
159 227 (99 *cis*, 128 *trans*) (**Supplementary Table 4**).

160

161 To validate our pQTL results, we tested significant associations from our discovery meta-
162 analysis for replication in an independent cohort (ARISTOTLE) comprising 1,585 participants
163 with Olink plasma proteomic data¹². Of the 180 pQTL signals, we were able to test 174 in the
164 ARISTOTLE data, of which 168 had a directionally consistent effect estimate. There was a
165 strong correlation (Pearson $r=0.97$) between the pQTL effect estimates in ARISTOTLE and in
166 the discovery meta-analysis; this correlation was consistent for both *cis*- and *trans*-pQTL effect
167 sizes ($r=0.99$ and $r=0.94$, respectively) (**Extended Data Figure 3**). Out of the 174 pQTL
168 signals, 32 replicated at $P\leq 5\times 10^{-10}$ (linear regression) and 72 at $P\leq 2.8\times 10^{-4}$ (a Bonferroni-
169 corrected threshold), respectively (**Supplementary Table 5**). We also tested our significant
170 pQTLs for replication in 35,556 Icelanders from the deCODE study⁹, which assayed plasma
171 proteins using the aptamer-based SomaScan platform (**Supplementary Note 2**). 72 of the 91
172 proteins in our study were measured in the deCODE study. Of the 158 locus-protein
173 associations that could be tested, 75 were significant at $P\leq 5\times 10^{-10}$ (linear regression), and 96
174 were significant at $P\leq 2.8\times 10^{-4}$. Overall, we replicated 126 (71%) of the 178 testable pQTLs in
175 either ARISTOTLE or deCODE at $P\leq 2.8\times 10^{-4}$ (linear regression) (**Supplementary Note Table**
176 **1**).

177

178 In line with other GWAS, we observed an inverse relationship between effect size and minor
179 allele frequency (MAF), with rarer pQTL variants generally showing larger effect sizes
180 (**Extended Data Figure 4a**). The proportion of variance explained by the significant sentinel
181 variants from our discovery meta-analysis varied from 0.003 for NTF3 to 0.285 for CCL8
182 (**Extended Data Figure 4b**).

183

184 **Annotation and characterisation of *cis*-pQTLs.** Of the 59 *cis*-pQTLs identified, 11 sentinel
185 variants were protein-altering variants (PAVs) (10 missense and 1 splice acceptor). A further
186 10 sentinel variants were in high linkage disequilibrium ($r^2>0.8$) with a protein-altering variant
187 (all missense); of these, 7 were variants in the gene encoding the target protein itself and 3 in
188 another nearby gene (**Supplementary Note 3**). PAVs can result in false positive *cis*-pQTL
189 signals by altering protein epitopes recognised by antibodies used in proteomic assays¹³.
190 However, they can also impact the abundance of plasma proteins through several
191 mechanisms, including protein translation, secretion into the circulation, enzymatic cleavage
192 of pre-proteins, and protein clearance and degradation. Alternatively, plasma protein
193 abundance can also be affected by altered transcriptional regulation in blood cells or other
194 tissues.

195

196 We next examined the degree to which our 59 *cis*-pQTLs were explained by corresponding

197 *cis*-eQTLs, by comparing our findings with publicly available *cis*-eQTL data. In a meta-analysis
198 of whole blood eQTL data from the eQTLGen Consortium¹⁴, we found a genome-wide
199 significant ($P \leq 5 \times 10^{-8}$; meta-analysis) *cis*-eQTL for 40 of the 59 *cis*-pQTLs, where the *cis*-eQTL
200 target gene encodes the *cis*-pQTL target protein. However, statistical colocalisation analyses
201 showed that only six (rs34790908-*TNFSF12*, rs72912115-*TGFA*, rs471994-*MMP1*, rs674379-
202 *CD5*, rs450373-*CXCL5*, rs5744249-*IL18*) of these *cis*-eQTLs colocalised (posterior probability
203 [PP] ≥ 0.8) with their cognate *cis*-pQTLs (**Supplementary Table 6**), indicating that the
204 remaining 34 eQTL-pQTL pairs may not share the same underlying causal genetic variant.
205 Examination of regional association plots confirmed that the majority of blood eQTL and pQTL
206 signals were distinct (**Supplementary Figure 3**). Of the 6 colocalising eQTL-pQTL pairs, 5
207 were directionally consistent. However, the eQTL and pQTL for *IL18* at rs5744249 were
208 oppositely associated with the mRNA and protein levels. rs5744249 resides in intron 2 of *IL18*
209 and is in high LD ($r^2 > 0.8$) with a 3' UTR variant (rs5744292, $r^2 = 0.98$ | 1000G EUR), but no
210 PAVs. Therefore, the directional discordance is not easily explained either by an artefactual
211 pQTL signal due to altered antibody binding or by a difference in the release of *IL18* into the
212 circulation due to differences in protein structure, but may instead relate to differential post-
213 transcriptional regulatory mechanisms or contributions of different cell types to the plasma
214 pQTL versus whole blood eQTL. Indeed, directional uncoupling of eQTL-pQTL pairs has been
215 previously reported⁸, and eQTL directional discordance has been observed between different
216 tissues¹⁵ or even within different leucocytes¹⁶.

217

218 Since tissues other than blood are the primary source of many plasma proteins, we explored
219 eQTL data across a range of tissues and cell types from the Genotype-Tissue Expression
220 (GTEx) (v8) project¹⁵ and the eQTL Catalogue¹⁷. Systematic COLOC analyses revealed
221 colocalising (PP ≥ 0.8) *cis*-eQTLs in at least one tissue or cell type for 32 of the 59 *cis*-pQTLs
222 (**Supplementary Tables 7-8**); 16 were highlighted by both eQTL resources, 12 by GTEx only,
223 and the remaining 4 by the eQTL Catalogue only. These included all 6 colocalising *cis*-eQTLs
224 from eQTLGen. These findings suggest that at least 50% of our *cis*-pQTLs may be driven by
225 underlying cognate *cis*-eQTLs. In most cases, colocalisation (PP $H_4 \geq 0.8$) between *cis*-eQTL-
226 pQTL pairs was observed across two or more distinct tissues or cell types, up to a maximum
227 of 41 (for rs1883832-*CD40*). In other cases, colocalisation was observed in just a single tissue
228 or cell type (e.g. the colocalising *cis*-eQTL signal (rs62360376) for *GDNF* was found only in
229 skeletal muscle). Of the 27 *cis*-pQTLs without a corresponding colocalising *cis*-eQTL, for 12
230 the sentinel variant or a proxy in high LD was a PAV (**Supplementary Note Table 3**).

231

232 **Identifying the mediators of *trans*-pQTLs.** We sought to identify the most likely gene
233 mediators of the *trans*-pQTLs using the ProGeM bioinformatics tool¹⁸, which utilises genomic

234 (e.g., *cis*-eQTL) and biological (e.g., gene ontology and pathways) annotation data from
235 multiple sources. For some *trans*-pQTLs, we identified strong evidence to implicate a gene
236 encoded near the pQTL as mediating the distant association with the target protein. Examples
237 included receptor-ligand pairs such as IL-6-IL6-R, IL-10-IL-10RA, CCL2-CCR2, CCL4-CCR5,
238 and CCL11-CCR3. We also identified genes mediating pQTLs through intracellular signalling
239 pathways rather than direct ligand-receptor interactions. An example is rs385076, an intronic
240 variant in *NLRC4*, which is a *trans*-pQTL for IL-18. IL-18 is synthesised as an inactive
241 precursor (pro-IL-18), which is cleaved by caspase-1 in the NLRC4 inflammasome to produce
242 the active form of IL-18 (**Figure 3a**). Since rs385076 is also a *cis*-eQTL for the inflammasome
243 gene *NLRC4* (**Figure 3b**), together, these QTL data suggest that genetic variation in *NLRC4*
244 alters its expression and thereby inflammasome activity, with consequent effects on circulating
245 IL-18 levels.

246

247 Following a manual literature review to refine the ProGeM output, we narrowed down the most
248 likely mediating gene(s) to either one or two candidates for 100 of the 121 *trans*-pQTLs
249 (**Supplementary Table 9**). For 94, one of the three nearest genes to the sentinel variant was
250 the primary candidate. In several instances where either one or two candidate genes were
251 prioritised, ProGeM revealed functional links between both (i) the sentinel variant and the
252 nearby candidate mediating gene (e.g., *cis*-eQTL) and (ii) the same candidate mediating gene
253 and the *trans*-affected protein(s) (e.g., through protein-protein interaction). We have previously
254 shown that such convergence on the same gene is indicative of a strong candidate¹⁸. An
255 example of this is the *trans*-pQTL at rs12075, which is associated with multiple chemokines
256 (CCL2, CCL7, CCL8, CCL11, CCL13, CXCL6) that attract and activate leucocytes. rs12075
257 is a missense variant and a *cis*-eQTL in whole-blood for the *DARC* gene, which encodes the
258 atypical chemokine receptor 1 (ACKR1) protein. STRINGdb analysis revealed that ACKR1 is
259 an interacting partner for 3 (CCL2, CCL7, CCL8) of the 6 *trans*-affected chemokines. Previous
260 studies have shown that ACKR1 acts as a negative regulator of inflammation by non-
261 specifically binding both the CCL and CXCL family of chemokines¹⁹, suggesting an
262 explanation for the multiple chemokine associations at this variant. Potentially downstream of
263 its effects on chemokines, rs12075 is also associated with white blood cell count, as well as
264 monocyte and basophil count²⁰ (**Extended Data Figures 5-6**).

265

266 We found that plasma levels of some proteins were associated with numerous genetic loci,
267 with IL12B, KITLG, and TNFSF10 regulated by seven genetic loci each. We hypothesised that
268 the mediating genes at each of the regulatory loci for a given protein might be functionally
269 related, enabling identification of shared pathways and/or the most likely mediating gene(s).

270 We therefore generated protein-protein interaction networks for each of these multi-locus-
271 regulated proteins and their respective candidate mediating genes (**Extended Data Figure**
272 **7**). For TNFSF10, the network analysis linked genetic regulators of TNFSF10 to the
273 plasminogen-activating system (**Extended Data Figure 7a, Supplementary Note 4**). For
274 KITLG, a driver of hematopoiesis²¹, we found a cluster of interacting proteins, including PON1,
275 ABCA1, PLTP, (**Extended Data Figure 7b**) converging on cholesterol metabolism.
276 Supporting this, we found that 5 of the 7 *trans*-pQTLs for KITLG were significantly ($P \leq 5 \times 10^{-8}$,
277 linear regression) associated with levels of either HDL or LDL cholesterol, and some with other
278 lipids such as triglycerides and blood cell traits (**Supplementary Table 10**). Our findings
279 therefore suggest a link between plasma KITLG levels, cholesterol metabolism and altered
280 hematopoiesis.

281

282 **Overlap with GWAS of traits and diseases.** Genome-wide association studies (GWAS)
283 have identified thousands of genomic regions associated with common diseases²², including
284 immune-mediated diseases (IMDs). Many of these disease-associated loci lie outside protein-
285 coding regions, leaving the effector molecules and pathways by which these genetic variants
286 confer disease risk unclear. Integration of pQTL and GWAS data can help bridge this
287 knowledge gap by linking disease risk loci to specific proteins. To this end, we looked for
288 overlap between pQTLs, or proxy variants in high LD ($r^2 \geq 0.8$) with our sentinel variants, and
289 disease-associated variants from GWAS. This revealed overlap between our pQTLs and
290 disease-associated variants for 73 diseases (**Extended Data Figure 8, Supplementary**
291 **Table 11**). Examples of genetically anchored protein-disease connections included: TNFSF11
292 (RANKL) with osteoporosis and hypothyroidism; NGF (nerve growth factor) with migraine;
293 TNFSF12 (TWEAK) with hypertension; and FGF5 with hypertension and cardiovascular
294 diseases.

295

296 We next focussed on IMDs in more detail, intersecting our pQTLs with IMD GWAS data to
297 identify proteins linking genotype and disease phenotypes. We found that 31 of our pQTLs
298 overlap GWAS hits for at least one common IMD, with 76 unique pQTL-protein-disease
299 associations (**Supplementary Table 12, Extended Data Figure 9**). For example, we
300 observed that a *cis*-pQTL for IL-10 was also associated with risk of inflammatory bowel
301 disease (IBD), with the allele associated with higher plasma IL-10 correlating with reduced
302 IBD risk, consistent with the anti-inflammatory effects of IL-10. Some pQTLs showed diverging
303 directions of effect on different diseases (e.g. the *trans*-pQTL at *IL6R* for plasma IL-6 levels
304 described earlier had opposing directions of effect on risk of rheumatoid arthritis and allergic
305 diseases (**Extended Data Figure 9**), as previously described^{23,24}).

306

307 **Trans-pQTL implicates the *LTBR-LTA* axis in multiple sclerosis.** We identified a *trans*-
308 pQTL for lymphotoxin α (LTA, also known as TNF- β) at rs2364485 on chromosome 12 (**Table**
309 **1**), an intergenic variant previously found to be associated with multiple sclerosis²⁵. We found
310 that the multiple sclerosis risk allele, rs2364485:A, was associated with higher plasma levels
311 of LTA. We next applied the ProGeM algorithm which revealed two candidate genes in the
312 region near the pQTL that might mediate the *trans*-pQTL: *TNFRSF1A* (encoding tumour
313 necrosis factor receptor 1, TNFR1) and *LTBR* (encoding lymphotoxin beta receptor, LTBR).
314 LTA is a ligand for TNFR1, but also can bind the membrane-bound receptor LTBR when in
315 complex with LTB. Functional studies have shown that *TNFRSF1A* is the causal gene
316 underlying a neighbouring independent multiple sclerosis association in the region, about 70kb
317 upstream from rs2364485. The sentinel variant at this neighbouring signal, rs1800693, results
318 in an alternative *TNFRSF1A* isoform due to skipping of exon 6²⁶. We therefore sought to
319 determine if *TNFRSF1A* is also the likely mediating gene for the LTA *trans*-pQTL at
320 rs2364485, or whether *LTBR* is the more likely candidate. Through mining of eQTL databases,
321 we found that rs2364485 is a *cis*-eQTL for *LTBR* (but not *TNFRSF1A*) in multiple tissues,
322 including in the eQTLGen consortium meta-analysis of whole-blood¹⁴, with the multiple
323 sclerosis risk allele (rs2364485:A) associated with reduced *LTBR* mRNA. Pairwise statistical
324 colocalisation analyses using conditioned *LTBR* eQTL data (from eQTLGen) and multiple
325 sclerosis GWAS data²⁵ (**Methods**) showed that the rs2364485 *trans*-pQTL signal for LTA
326 colocalises with both *LTBR* mRNA expression in whole blood (PP=0.79) and multiple sclerosis
327 (PP=0.86) (**Figure 4**). Taken together, these data are consistent with a pathogenic model
328 whereby the multiple sclerosis risk allele results in lower abundance of LTBR (the receptor)
329 and consequently higher circulating levels of the ligand LTA.

330

331 **Mendelian randomisation to identify protein drivers of IMDs.** Observational studies
332 comparing patients with IMDs to healthy controls have identified many proteins that are
333 dysregulated. However, it is often unclear whether such proteins play causal roles in the
334 disease process or are merely downstream markers. Distinguishing these possibilities is
335 important therapeutically, as pharmacological targeting of the latter is unlikely to be beneficial.
336 We therefore applied Mendelian randomisation (MR), an approach that tests the causal role
337 of a risk factor ('exposure') in a disease in observational data using genetic variants as
338 instrumental variables²⁷. We used the 58 proteins with *cis*-pQTLs outside the HLA region in
339 our study as exposures and 14 IMDs as outcomes (**Methods**). By restricting our use of genetic
340 instruments to *cis*-pQTLs, we reduce the likelihood of violating MR assumptions through
341 horizontal pleiotropy. Using Generalised Summary-data-based MR (GSMR)²⁸, we found 22
342 significant (FDR<0.01) putative causal associations (**Figure 5, Supplementary Table 13**). To
343 evaluate the robustness of these associations, we performed additional checks including

344 evaluating the strength of the disease association in the GWAS summary statistics and
345 whether there might be confounding due to LD (**Methods, Supplementary Table 14**). After
346 applying these filters, 10 disease-protein pairs with robust evidence remained (**Table 1**).
347 These results highlighted a number of established links between proteins and inflammatory
348 diseases that are supported by other lines of evidence. For example, our finding that genetic
349 predisposition to higher plasma levels of IL-12B (a subunit of IL-12) was associated with
350 increased risk of IBD is consistent with the therapeutic benefit of ustekinumab (a monoclonal
351 antibody targeting the p40 subunit common to IL-12 and IL-23) in IBD (**Supplementary Table**
352 **15**).

353

354 Our MR analysis implicated CXCL5, a chemokine that acts on neutrophils, in the aetiology of
355 ulcerative colitis (UC). The plasma *cis*-pQTL for CXCL5, colocalised with *cis*-eQTLs for CXCL5
356 in both blood and gut tissue, and with the UC GWAS signal (**Figure 6a**). To further explore
357 the role of CXCL5 in UC, we compared expression of CXCL5 transcripts in gut samples from
358 patients with IBD and healthy controls using IBD TAMMA, an open-access resource for
359 interrogating transcriptomic data across multiple datasets²⁹. We observed that CXCL5 gene
360 expression was significantly increased in mucosal biopsies from patients with UC in
361 comparison with biopsies from healthy control participants (\log_2 fold change 7.07, P 1.98x10⁻
362 ¹⁷⁴, Wald test) (**Figure 6b**). Indeed, CXCL5 was the third most highly upregulated transcript
363 across the transcriptome (**Figure 6c**). We replicated these findings in three independent
364 datasets (**Figure 6d**). Of note, our MR analysis revealed the association of CXCL5 was
365 restricted to UC (unadjusted $P=2.3 \times 10^{-6}$, GSMR), with no significant association in CD
366 (unadjusted $P=0.4$) (**Figures 6a,e**). Supporting this specific pathogenic effect, CXCL5 gene
367 expression in gut samples from IBD patients was higher in UC than in CD (**Figure 6b**).
368 Counter-intuitively (given the upregulation of CXCL5 in UC patient tissue samples), evaluation
369 of the direction of MR association effect revealed that genetic susceptibility to higher plasma
370 CXCL5 reduces risk of UC (**Figure 6e**). This effect was consistent across 12 of the 13 the
371 individual genetic variants used in the MR score (**Extended Data Figure 10a**). We found
372 consistent directions of effect for the CXCL5 plasma pQTL and the blood and gut eQTLs
373 (**Extended Data Figure 10b**), indicating that our results are generalizable at both the mRNA
374 and protein level and across local and systemic sites. Together these data indicate that genetic
375 tendency to lower CXCL5 is a causal risk factor for development of UC, despite the strong
376 upregulation of CXCL5 once disease develops.

377

378 We observed that genetic predisposition to higher plasma CD40 levels was associated with
379 increased rheumatoid arthritis risk, consistent with evidence from both animal models and

380 humans implicating the CD40 pathway in rheumatoid pathogenesis³⁰. In addition, our MR
381 analysis identified a potential causal role for the CD40 pathway in IBD (including both Crohn's
382 and UC) and multiple sclerosis. Interestingly, however, the MR associations for these diseases
383 had the opposite direction of effect compared to rheumatoid arthritis i.e. genetic predisposition
384 to *lower* plasma CD40 levels was associated with *higher* risk of IBD and multiple sclerosis.
385 These findings highlight how the same pathway can have pleiotropic effects on disease
386 susceptibility, but also point to the complexity of immune-mediated disease pathogenesis, with
387 opposing effects on different diseases.

388

389

390 **DISCUSSION**

391

392 Here, we performed a large-scale pQTL GWAS of 91 circulating inflammation-related proteins
393 measured using Olink immunoassays, identifying 180 significant primary pQTL signals (59 *cis*,
394 121 *trans*). Colocalisation analysis suggested that only a small proportion of the plasma *cis*-
395 pQTLs reported here are underpinned by the same causal genetic variant as the whole blood
396 *cis*-eQTL for the corresponding gene. Of note, the plasma proteome is not the direct corollary
397 of the whole blood transcriptome: plasma pQTL studies examine genetic effects on
398 extracellular protein levels whereas blood eQTL studies examine the effects on intracellular
399 RNA levels (predominantly in leucocytes). This has several implications. First, plasma protein
400 levels can be affected by non-transcriptional mechanisms including cleavage, secretion and
401 clearance. Second, a wide range of tissues other than blood cells (e.g. the liver) contribute to
402 the plasma proteome. This is evident when considering circulating proteins that are measured
403 as biomarkers in clinical practice (e.g. albumin produced by the liver, troponin by the heart,
404 PSA by the prostate). Indeed, by extending our comparison across multi-tissue eQTL
405 databases, we showed that at least 50% of the *cis*-pQTLs we observed are likely driven by
406 cognate *cis*-eQTLs in a diverse range of tissues and cell-types. Blood eQTL studies have been
407 carried out using sample sizes similar to the sample size in our pQTL study. eQTL studies in
408 other tissues tend to be smaller, and so it is likely that some of the plasma *cis*-pQTLs observed
409 here are underpinned by tissue-specific eQTLs that have not yet been detected due to lack of
410 statistical power. Finally, other mechanisms such as alternative splicing might account for
411 some *cis*-pQTLs without corresponding eQTLs.

412

413 Our pQTL study identified twice as many *trans* associations compared to *cis* (121 versus 59,
414 respectively), in keeping with other well-powered pQTL studies (e.g., ⁷⁻⁹). The integration of
415 *cis*-pQTLs (and *cis*-eQTLs) with GWAS data provides useful, if sometimes obvious, insights
416 into the upstream mechanisms of disease, since the mediating gene has usually already been

417 suspected by virtue of the location of the GWAS signal. In contrast, *trans*-pQTLs represent a
418 double-edged sword for interpreting genetic associations with disease. On the one hand, they
419 often represent a less direct link from genotype to disease than *cis*-pQTLs, and from the
420 perspective of causal inference analysis, are more vulnerable to violating the assumptions of
421 MR through horizontal pleiotropy. On the other hand, they can reveal important molecular
422 mediators of disease encoded by genes distant from the disease GWAS signal. For example,
423 we identified a *trans*-pQTL (rs2364485) for LTA at a multiple sclerosis risk locus. This multiple
424 sclerosis risk locus contains two plausible causal genes (*TNFRSF1A*, *LTBR*) and two
425 independent signals for multiple sclerosis risk (rs1800693, rs2364485). By integrating whole
426 blood eQTL and multiple sclerosis GWAS data, we showed that *LTBR* is the most likely gene
427 mediating our LTA *trans*-pQTL at rs2364485, and one of the multiple sclerosis signals at the
428 locus. LTA is a member of the TNF superfamily of proteins and is the only member of this
429 superfamily that is generated as a secreted protein rather than through cleavage of a
430 membrane-bound protein. The multiple sclerosis risk allele is associated with lower expression
431 of *LTBR* and higher circulating protein levels of LTA, a component of its ligand. This raises the
432 question as to whether elevated LTA is secondary to lower *LTBR*, or vice versa (e.g. through
433 compensatory receptor downregulation). The distinction between *cis*- and *trans*-QTLs enables
434 us to address this. Given that the eQTL for *LTBR* is *cis*, and the pQTL for LTA is *trans*, it is
435 highly likely that former is the upstream effect, with the higher levels of soluble LTA occurring
436 as a result of reduced binding to its receptor. This demonstrates the value of pairing QTLs for
437 ligands and their receptors for deconvoluting the ordering of biological pathways.

438

439 Integration of pQTLs with GWAS disease signals revealed disease-protein connections
440 reflecting both established and plausible putative mechanisms of pathophysiology. For
441 example, a *cis*-pQTL for *TNFSF11* (*RANKL*) overlapped with GWAS signals for osteoporosis
442 and hypothyroidism. The former is consistent with *RANKL*'s well-established role in bone
443 biology and *RANKL* is the target of the anti-osteoporosis drug denosumab³¹. However, *RANKL*
444 also plays a role in the immune system³², and these effects may be relevant to risk of
445 autoimmune hypothyroidism. A *cis*-pQTL for *TNFSF12* (*TWEAK*) was associated with risk of
446 hypertension. *TWEAK* is a cytokine predominantly produced by leucocytes and has pleiotropic
447 actions, including on the endothelium^{33,34}, potentially explaining the association with blood
448 pressure. A *cis*-pQTL for *FGF5* was also associated with susceptibility to hypertension and
449 cardiovascular diseases, with the allele associated with higher plasma *FGF5* levels
450 associating with lower risk of cardiovascular diseases. Consistent with this, there are reports
451 that *FGF5* has cardioprotective effects in pig models³⁵.

452

453 31 of our pQTLs overlap GWAS hits for at least one common IMD. Disease-protein links

454 identified from this analysis highlighted commonalities in pathogenesis between specific IMDs,
455 mirroring the overlap in clinical manifestations. However, the contributions of proteins to IMD
456 risk were sometimes complex, with the same protein conferring risk of one IMD but protecting
457 from another. For example, genetic predisposition to higher levels of soluble IL6 had opposing
458 effects on risk of rheumatoid arthritis and allergic disease. We observed a similar phenomenon
459 for CD40, with genetic predisposition to higher CD40 increasing risk of RA but protecting
460 against IBD and multiple sclerosis.

461

462 The development of biologic therapies targeting specific inflammatory proteins has
463 transformed the clinical management of immune-mediated diseases³⁶. Understanding which
464 proteins are drivers of disease and distinguishing these from proteins that are simply markers
465 of inflammation is therefore important for development of new treatments. To this end, we
466 used Mendelian randomisation to evaluate the causal contributions of proteins to different
467 IMDs. Our results identify pathways that are already the target of existing drugs (e.g. IL-12B
468 in IBD), providing confirmation of the utility of this approach, and also highlight new potential
469 therapeutic targets.

470

471 One such example was the CD40 pathway in rheumatoid arthritis. CD40 is a stimulatory
472 receptor constitutively or inducibly expressed on both immune and non-immune cells³⁷. Its
473 ligand, CD40L, is expressed primarily on activated T cells but also on a range of other immune
474 and non-immune cells. CD40L-CD40 binding triggers immune cell activation and proliferation
475 and inflammatory cytokine production, and the differentiation of B cells into IgG-secreting
476 plasma cells, making it central to antibody responses. In a murine model of inflammatory
477 arthritis, knock out or inhibition of the CD40 pathway resulted in reduced inflammation³⁸.
478 Observational studies have demonstrated upregulation of CD40L in the blood and tissues of
479 patients with RA and other autoimmune rheumatic diseases^{30,39}. These findings previously
480 motivated development of drugs targeting the CD40 pathway in RA and other IMDs, but anti-
481 CD40L therapy was complicated by thrombosis due to cross-linking CD40L on platelets.
482 Therapeutic targeting of CD40 rather than CD40L may avoid this. Our MR results suggest
483 rheumatoid arthritis as a candidate for this approach. However, the directionally discordant
484 effects we observed of CD40 on RA versus multiple sclerosis and IBD raises the possibility of
485 triggering other forms of immune-mediated diseases as a side-effect of anti-CD40 therapy.
486 This has some parallels with therapies targeting TNF, which are effective in rheumatoid
487 arthritis but not in multiple sclerosis, and indeed can worsen multiple sclerosis or provoke *de*
488 *novo* central nervous system demyelination^{40,41}.

489

490 Our MR findings implicate CXCL5 in the aetiology of UC, where genetic susceptibility to higher

491 levels of plasma CXCL5 was associated with lower UC risk. Examination of eQTL data
492 revealed this observation was consistent at the RNA level in both the blood and gut tissue. By
493 contrast, in our case-control analysis comparing gut tissue from UC patients versus controls,
494 CXCL5 is one of the most upregulated transcripts. A previous study reported that serum levels
495 of CXCL5 are higher in IBD patients than controls⁴². Recent studies using UC gut tissue have
496 implicated upregulation of genes encoding neutrophil-targeting chemokines, including CXCL5,
497 by non-immune cells as correlating with important histopathological features, such as
498 ulceration, and differentiating patient trajectories, including their responsiveness to different
499 treatments^{43,44}. Targeting CXCR2, the receptor for CXCL5, significantly attenuates animal
500 models of UC⁴⁴. One possible explanation that may reconcile these apparently contradictory
501 findings is that genetic tendency to lower CXCL5 expression increases UC risk through
502 impaired mucosal immune homeostasis, but that elevated CXCL5 is an important driver of
503 tissue injury once disease is initiated. By analogy, a non-coding genetic variant associated
504 with lower gene and protein expression of TNFSF15 (encoding the inflammatory cytokine
505 TL1A) in monocytes and macrophages increases IBD susceptibility⁴⁵, but TL1A is upregulated
506 both systemically and in the gut in patients with active IBD^{46,47}, and anti-TL1A therapies have
507 recently shown efficacy in IBD in phase 2 randomised trials (NCT05013905 and
508 NCT04996797⁴⁸).

509

510 Our study has several limitations. Our pQTL analysis was restricted to 91 proteins, limiting the
511 generalisability of our findings, particularly with regards to genetic architecture. Since this was
512 a pQTL meta-analysis, study-level technical variation resulted in heterogeneity, which
513 necessitated filtering out of potentially spurious associations that were inconsistent across
514 cohorts. There is a risk that some true biological signals were also removed in this process.
515 Very large single cohorts with standardised sample processing such as UK Biobank will avoid
516 this issue. Our meta-analysis consisted predominantly of general population cohorts without
517 inflammatory disease. There may be context-specific pQTLs that are only present during
518 infection or inflammation, which our study may not have detected. By analogy, eQTL studies
519 using human immune cells stimulated *in vitro* (e.g. with lipopolysaccharide or interferon) have
520 demonstrated eQTLs that are not present in resting cells but become apparent in the context
521 of cellular activation^{49,50}. Conducting well-powered pQTL studies in patients with inflammation
522 will be an important future research endeavour. Where proteins exist in both membrane-bound
523 and cleaved states, it is not always clear whether plasma proteomic assays are exclusively
524 capturing the soluble form or also protein from cell membranes (e.g. arising from *in vivo*
525 sources such as exo-/ectosomes, or *ex vivo* processes such as venepuncture or sample
526 processing). This complicates the interpretation of the direction of effect from MR analysis.
527 Future well-powered studies examining genetic determinants of cell surface protein

528 expression measured through flow cytometry would provide valuable complementary
529 information to aid the interpretation of plasma pQTL studies. Finally, as with all
530 epidemiological-scale pQTL studies, proteins were measured in plasma (i.e. the extracellular
531 component of blood), which may not always be the disease-relevant biological compartment,
532 and where the direction of genotype-expression association may even be opposite to the site
533 of inflammation. Thus, future tissue- and cell-specific pQTL studies will be valuable to
534 understand differences in genetic signals across tissues.

535

536 In summary, we have used a large international consortium to identify the genetic
537 determinants of a set of inflammation-related proteins, providing insight into the aetiology of
538 immune-mediated diseases. The pQTL summary statistics generated in this study will be a
539 valuable resource for interrogating future disease GWAS and guiding drug target identification
540 and prioritisation.

541

542 **Acknowledgements**

543

544 This work was performed under the auspices of the SCALLOP Consortium.

545 We thank:

- 546 -study participants from the contributing cohorts;
- 547 -the International Multiple Sclerosis Genetics Consortium, who provided multiple sclerosis
548 GWAS summary statistics used in our analyses;
- 549 -Aikaterina Siopi and Dianna McLeod for support with SCALLOP Consortium administration;
- 550 -the authors of the GCTA software for advice;
- 551 -Bram Prins for help with INTERVAL study genotype data quality control;
- 552 -Arianne Richard for comments on the manuscript.

553

554 J.E.P was supported by a grant and a fellowship from the Medical Research Foundation grant
555 (MRF-042-0001-RG-PETE-C0839, MRF-057-0003-RG-PETE-C0799). E.J.N. was supported
556 by the Schmidt Science Fellows, in partnership with the Rhodes Trust. P.S. was supported by
557 a Rutherford Fund Fellowship from the Medical Research Council (grant no. MR/S003746/1).
558 J.D. holds a British Heart Foundation Professorship and a NIHR Senior Investigator Award*.
559 C.H. is supported by an MRC University Unit Programme Grant "QTL in Health and Disease"
560 (U.MC_UU_00007/10).

561

562 Funding of the GWAS and proteomics studies of STABILITY and ARISTOTLE were supported
563 by GlaxoSmithKline, BristolMyersSquibb, and the Swedish Foundation for Strategic Research
564 (grant number RB13-0197).

565

566 The Orkney Complex Disease Study (ORCADES) was supported by the Chief Scientist Office
567 of the Scottish Government (CZB/4/276, CZB/4/710), a Royal Society URF to J.F.W., the MRC
568 Human Genetics Unit quinquennial programme "QTL in Health and Disease", Arthritis
569 Research UK and the European Union framework program 6 EUROSPAN project (contract
570 no. LSHG-CT-2006-018947). DNA extractions were performed at the Edinburgh Clinical
571 Research Facility, University of Edinburgh. We would like to acknowledge the invaluable
572 contributions of the research nurses in Orkney, the administrative team in Edinburgh and the
573 people of Orkney. For the purpose of open access, the author has applied a Creative

574 Commons Attribution (CC BY) licence to any Author Accepted Manuscript version arising from
575 this submission.

576 Participants in the INTERVAL trial were recruited with the active collaboration of NHS Blood
577 and Transplant England (www.nhsbt.nhs.uk), which has supported field work and other
578 elements of the trial. DNA extraction and genotyping were co-funded by the National Institute
579 for Health and Care Research (NIHR), the NIHR BioResource (<http://bioresource.nihr.ac.uk>)
580 and the NIHR Cambridge Biomedical Research Centre (BRC-1215-20014). The academic
581 coordinating centre for INTERVAL was supported by core funding from the: NIHR Blood and
582 Transplant Research Unit (BTRU) in Donor Health and Genomics (NIHR BTRU-2014-10024),
583 NIHR BTRU in Donor Health and Behaviour (NIHR203337), UK Medical Research Council
584 (MR/L003120/1), British Heart Foundation (SP/09/002; RG/13/13/30194; RG/18/13/33946)
585 and NIHR Cambridge BRC (BRC-1215-20014; NIHR203312) [*] and has received funding
586 from an EC-Innovative Medicines Initiative (BigData@Heart). The academic coordinating
587 centre thank blood donor centre staff and blood donors for participating in the INTERVAL trial.
588 This work was supported by Health Data Research UK, which is funded by the UK Medical
589 Research Council, Engineering and Physical Sciences Research Council, Economic and
590 Social Research Council, Department of Health and Social Care (England), Chief Scientist
591 Office of the Scottish Government Health and Social Care Directorates, Health and Social
592 Care Research and Development Division (Welsh Government), Public Health Agency
593 (Northern Ireland), British Heart Foundation and Wellcome. For the purpose of open access,
594 the author(s) has applied a Creative Commons Attribution (CC BY) licence to any Author
595 Accepted Manuscript version arising from this submission.

596 Estonian Biobank work was supported by the European Regional Development Fund and the
597 programme Mobilitas Pluss (MOBTP108), No. 2014-2020.4.01.15-0012 GENTRANSMED
598 and 2014-2020.4.01.16-0125 This study was also funded by the EU H2020 grant 692145, by
599 the Estonian Research Council Grant PUT1660 and by the Estonian Research Council grant
600 PUT (PRG1291). Data analyzes with Estonian datasets were carried out in part in the High-
601 Performance Computing Center of University of Tartu.

602 The SWEBIC biobank was supported by the Stanley Medical Research Institute. The
603 proteomic analyses in SWEBIC was funded by the Swedish foundation for Strategic Research
604 (KF10-0039). For RECOMBINE and SWEBIC, the data handling and analysis were enabled
605 by resources provided by the Swedish National Infrastructure for Computing (SNIC), partially
606 funded by the Swedish Research Council through grant agreement no. 2018-05973.

607
608 The CROATIA-Vis study was funded by grants from the Medical Research Council (UK), from
609 the Republic of Croatia Ministry of Science, Education and Sports (108-1080315-0302; 216-
610 1080315-0302) and the Croatian Science Foundation (8875). We thank the staff of several
611 institutions in Croatia that supported the field work, including Zagreb Medical Schools, the
612 Institute for Anthropological Research in Zagreb, the recruitment team from the Croatian
613 Centre for Global Health, University of Split and all the study participants.

614 The KORA study was initiated and financed by the Helmholtz Zentrum München – German
615 Research Center for Environmental Health, which is funded by the German Federal Ministry
616 of Education and Research and by the State of Bavaria. Furthermore, KORA research was
617 supported within the Munich Center of Health Sciences (MC-Health), Ludwig-Maximilians-
618 Universität, as part of LMUinnovativ.

619 The measurement of inflammatory biomarkers was funded by a grant from the German Center
620 for Diabetes Research (DZD; to C. Herder and B. Thorand). This work was also supported by
621 the Ministry of Culture and Science of the State of North Rhine-Westphalia and the German
622 Federal Ministry of Health. This study was supported in part by a grant from the German
623 Federal Ministry of Education and Research to the German Center for Diabetes Research

624 (DZD).

625 NP is supported by a Wellcome Trust Discovery award (225875/Z/22/Z). DC is supported by
626 the NIHR Imperial Biomedical Research Centre (BRC). Infrastructure support for this research
627 was provided by the NIHR Imperial Biomedical Research Centre (BRC).

628 Support for title page creation and format was provided by AuthorArranger, a tool developed
629 at the National Cancer Institute.

630 **Author Contributions**

631 J.E.P., N.E., E.M-D., A.K.H., A.K., S.E., L.F., C.H., L.J., S.E.B., P.S. conducted study-level
632 analyses.

633 C.G., D.S.P., O.P., B.T., H.G., M.R., U.V., T.O., C.H., A.J., U.G., N.P., O.H., N.M-C., P.K.J.,
634 J.D., L.P., L.K., M.L., J.F.W., A.S., L.W., A.M., A.S.B., and J.E.P. provided data and study
635 supervision.

636 D.C., J.D-B., N.P. collected IBD samples and generated and analysed the IBD RNA-seq data.

637 J.H.Z., D.S., A.K., J.M., P.S, conducted the meta-analysis and downstream analyses.

638 J.H.Z, D.S., E.N., A.S.B, and J.E.P. drafted the manuscript.

639 J.F.W., A.M., A.S.B., and J.E.P. conceived the project.

640 All authors critically reviewed the manuscript and gave final approval to publish.

641

642 **Competing interests**

643 J.D. serves on scientific advisory boards for AstraZeneca, Novartis, and UK Biobank, and has
644 received multiple grants from academic, charitable and industry sources outside of the
645 submitted work. A.S.B. has received grants unrelated to this work from AstraZeneca, Bayer,
646 Biogen, BioMarin, Bioverativ, Novartis and Sanofi. J.E.P. has received hospitality and travel
647 expenses to speak at Olink-sponsored academic meetings (none within the past 5 years).
648 During the drafting of the manuscript, D.S.P. became a full-time employee of AstraZeneca,
649 and P.S. became a full-time employee of GlaxoSmithKline. M.L. has received lecture
650 honoraria from Lundbeck pharmaceutical. The other authors declare no competing interests.

651

652

653

Protein	Disease	OR (95% CI)	P
CD40	Rheumatoid arthritis	1.28 (1.21-1.37)	1.4x10 ⁻¹⁵
CD40	Multiple sclerosis	0.75 (0.70-0.82)	1.2x10 ⁻¹²
CD40	Crohn's disease	0.81 (0.75-0.87)	2.2x10 ⁻⁸
CD40	IBD	0.87 (0.82-0.92)	1.9x10 ⁻⁶
CD5	Primary sclerosing cholangitis	0.50 (0.35-0.70)	8.1x10 ⁻⁵
CD6	IBD	1.10 (1.06-1.14)	2.1x10 ⁻⁷
CXCL5	Ulcerative colitis	0.79 (0.72-0.87)	2.3x10 ⁻⁶
IL-12B	IBD	1.38 (1.31-1.46)	1.5x10 ⁻³⁰
IL-12B	Ulcerative colitis	1.38 (1.29-1.48)	4.7x10 ⁻²⁰
IL-18R1	Eczema	1.15 (1.10-1.20)	2.1x10 ⁻¹⁰

655 **Table 1. Putative causal protein-disease associations from MR analysis.** IBD =
656 inflammatory bowel disease (i.e. based on GWAS where Crohn's disease and ulcerative colitis
657 cases are grouped together). P = two-sided P-value for the causal estimate of protein on
658 disease from the GSMR package. FDR = false discovery rate. OR = odds ratio associated
659 with a 1 standard deviation increase in the protein level. OR>1 indicates genetic propensity to
660 higher levels of the plasma protein is associated with higher disease risk, and OR<1 with
661 reduced risk. CI = confidence interval.

662

663 **Figure 1. Genomic map of genetic determinants of inflammation-related proteins.** Circos
664 plot linking the location of pQTLs to the gene encoding their associated proteins. Labels for
665 the *cis*-pQTLs (red) indicate the gene encoding the target protein. For the *trans*-pQTLs (blue),
666 the gene symbols of the target proteins are indicated along with the putative mediating gene(s)
667 at the *trans*-pQTLs in square brackets. $-\log_{10}(P\text{-values})$ are capped at 150 for readability. Two-
668 sided P-values are from meta-analysis of linear regression estimates.

669

670 **Figure 2. Genetic architecture of 91 inflammation-related proteins.** a) Distribution of the
671 number of identified pQTLs per protein. The HLA region was treated as a single region. b)
672 Circos plot showing the six *trans*-pQTLs for the *SH2B3* 'hotspot' locus on chromosome 12. c)
673 Manhattan plots showing genetic associations with plasma abundance of IL-12B and d)
674 TNFSF10 (TRAIL). The horizontal red line indicates statistical significance ($P=5\times 10^{-10}$). Two-
675 sided P-values are from meta-analysis of linear regression estimates. Nearest genes in the
676 region of pQTL signals are annotated.

677

678 **Figure 3. Genetic regulation of the inflammasome affects plasma IL-18 levels.** a)
679 Schematic illustrating the cleavage of pro-IL-18 by caspase 1 and subsequent secretion of
680 mature IL18 from the cell into the extracellular space. b) Regional association plots around
681 *NLRC4* showing: the *trans*-pQTL signal for plasma IL18 protein (top) from this study
682 (n=14,824), and the *cis*-eQTL signal for *NLRC4* (bottom) in whole blood from the eQTLGen
683 study (n=31,684)¹⁴. Purple diamond = sentinel pQTL variant. Other variants coloured by LD to

684 sentinel pQTL. Two-sided *P*-values are from meta-analysis of linear regression estimates.

685

686 **Figure 4. The *LTBR-LTA* axis in the aetiology of multiple sclerosis.** **a)** Unconditioned and
687 **b)** conditioned regional association plots at the *TNFRSF1A-LTBR* locus (rs2364485 +/- 100kb)
688 for MS (top), plasma LTA protein levels (middle), and *LTBR* mRNA expression in whole blood
689 from eQTLGen¹⁴ (bottom). MS associations were conditioned on rs1800693 (the strongest
690 disease signal in the region). *LTBR* mRNA expression levels were conditioned on the following
691 independent eQTLs: rs3759322, rs1800692, rs2228576, rs10849448, rs2364480, and
692 rs12319859. Purple diamond = sentinel pQTL variant. Other variants coloured by LD to
693 sentinel pQTL. Two-sided *P*-values are from meta-analysis of linear regression estimates.

694

695 **Figure 5. MR analysis of circulating proteins in immune-mediated disease aetiology.**
696 GSMR analysis²⁸ using *cis*-pQTLs as genetic instruments to test the causal role of plasma
697 proteins across IMDs. Cells are coloured according to effect size and direction: red = higher
698 genetically-predicted plasma protein levels are associated with higher risk of disease; blue =
699 higher genetically-predicted plasma protein levels are associated with lower risk of disease;
700 grey = fewer than 3 variants available to run the GSMR analysis. Associations with FDR≤0.01
701 are denoted with dots, with filled circles indicating those that are robust to confounding by LD
702 and open circles indicating those that were not.

703

704 **Figure 6. *CXCL5* in UC pathogenesis.** **a)** Genetic associations in the *CXCL5* gene region.
705 From top to bottom: plasma *CXCL5* pQTL (n=14,824 participants), whole blood eQTL (from
706 eQTLGen data, n=31,684 participants), colon tissue eQTL (GTEx, n=368 individuals),
707 ulcerative colitis (cases=12,366, controls=33,609) and Crohn's disease (cases=12,194,
708 controls=28,072)(from the IBD Genetics Consortium⁵¹). Purple diamond = sentinel pQTL
709 variant. Other variants coloured by LD to sentinel pQTL. *P*-values from linear regression. **b)**
710 Violin plots showing *CXCL5* expression in gut mucosal samples from patients with UC or CD
711 and healthy controls (HC) in IBD TaMMA. **c)** Volcano plot showing differential expression
712 analysis comparing colonic tissue from UC versus healthy controls (IBD TaMMA). Red and
713 blue points represent significantly (5% FDR) up- and down-regulated transcripts, respectively.
714 Grey = non-significant. P_{BH} = Benjamini-Hochberg adjusted *P*-values. *P*-values in c-d from
715 Wald test (two-sided). **d)** Replication. Left hand panel: *CXCL5* differential expression in colon
716 biopsies in UC versus HC from transcriptome-wide analysis across 3 cohorts. GSE numbers
717 are Gene Expression Omnibus accession numbers. Imperial = Imperial UC cohort. Each
718 lollipop represents a separate cohort. GSE16879 n=24 UC patients versus n=6 HC.
719 GSE206285 n=550 UC patients versus n=18 HC. Imperial n=16 UC vs 6 HC. Circle colour
720 indicates the log₂ fold change (FC) in *CXCL5* expression between UC and HC. Right hand
721 panel: *CXCL5* expression in colon biopsies sampled at baseline during the UNIFI clinical trial.
722 Each point represents an individual. **e)** Forest plot showing Mendelian randomisation analysis
723 for UC and CD. OR = odds ratio for the risk associated with a 1 standard deviation increase
724 in the level of the protein. Centre of bar = point estimate for OR, whiskers = 95% confidence
725 intervals (CI).

726 REFERENCES

- 727 1. Sun, B.B. *et al.* Genomic atlas of the human plasma proteome. *Nature* **558**, 73-79
728 (2018).
729

- 730 2. Enroth, S., Johansson, A., Enroth, S.B. & Gyllensten, U. Strong effects of genetic
731 and lifestyle factors on biomarker variation and use of personalized cutoffs. *Nat*
732 *Commun* **5**, 4684 (2014).
733
- 734 3. Suhre, K. *et al.* Connecting genetic risk to disease end points through the human
735 blood plasma proteome. *Nat Commun* **8**, 14357 (2017).
736
- 737 4. Emilsson, V. *et al.* Co-regulatory networks of human serum proteins link genetics to
738 disease. *Science* **361**, 769-773 (2018).
739
- 740 5. Melzer, D. *et al.* A genome-wide association study identifies protein quantitative trait
741 loci (pQTLs). *PLoS Genet* **4**, e1000072 (2008).
742
- 743 6. Lourdasamy, A. *et al.* Identification of cis-regulatory variation influencing protein
744 abundance levels in human plasma. *Hum Mol Genet* **21**, 3719-3726 (2012).
745
- 746 7. Folkersen, L. *et al.* Genomic and drug target evaluation of 90 cardiovascular proteins
747 in 30,931 individuals. *Nat Metab* **2**, 1135-1148 (2020).
748
- 749 8. Pietzner, M. *et al.* Mapping the proteo-genomic convergence of human diseases.
750 *Science* **374**, eabj1541 (2021).
751
- 752 9. Ferkingstad, E. *et al.* Large-scale integration of the plasma proteome with genetics
753 and disease. *Nature Genetics* **53**, 1712-1721 (2021).
754
- 755 10. Zhang, J. *et al.* Plasma proteome analyses in individuals of European and African
756 ancestry identify cis-pQTLs and models for proteome-wide association studies.
757 *Nature Genetics* **54**, 593-602 (2022).
758
- 759 11. Gudjonsson, A. *et al.* A genome-wide association study of serum proteins reveals
760 shared loci with common diseases. *Nature Communications* **13**, 480 (2022).
761
- 762 12. Siegbahn, A. *et al.* Multiplex protein screening of biomarkers associated with major
763 bleeding in patients with atrial fibrillation treated with oral anticoagulation. *J Thromb*
764 *Haemost* **19**, 2726-2737 (2021).
765
- 766 13. Pietzner, M. *et al.* Synergistic insights into human health from aptamer- and antibody-
767 based proteomic profiling. *Nature Communications* **12**, 6822 (2021).
768
- 769 14. Vösa, U. *et al.* Large-scale cis- and trans-eQTL analyses identify thousands of
770 genetic loci and polygenic scores that regulate blood gene expression. *Nat Genet* **53**,
771 1300-1310 (2021).
772
- 773 15. The GTEx Consortium *et al.* The GTEx Consortium atlas of genetic regulatory effects
774 across human tissues. *Science* **369**, 1318-1330 (2020).
775
- 776 16. Peters, J.E. *et al.* Insight into Genotype-Phenotype Associations through eQTL

- 777 Mapping in Multiple Cell Types in Health and Immune-Mediated Disease. *PLoS*
778 *Genet* **12**, e1005908 (2016).
779
- 780 17. Kerimov, N. *et al.* A compendium of uniformly processed human gene expression
781 and splicing quantitative trait loci. *Nat Genet* **53**, 1290-1299 (2021).
782
- 783 18. Stacey, D. *et al.* ProGeM: a framework for the prioritization of candidate causal
784 genes at molecular quantitative trait loci. *Nucleic Acids Res* **47**, e3 (2019).
785
- 786 19. Rappoport, N., Simon, A.J., Amariglio, N. & Rechavi, G. The Duffy antigen receptor
787 for chemokines, ACKR1, - 'Jeanne DARC' of benign neutropenia. *Br J Haematol* **184**,
788 497-507 (2019).
789
- 790 20. Chen, M.-H. *et al.* Trans-ethnic and Ancestry-Specific Blood-Cell Genetics in 746,667
791 Individuals from 5 Global Populations. *Cell* **182**, 1198-1213.e1114 (2020).
792
- 793 21. Hassan, H.T. & Zander, A. Stem cell factor as a survival and growth factor in human
794 normal and malignant hematopoiesis. *Acta Haematol* **95**, 257-262 (1996).
795
- 796 22. Claussnitzer, M. *et al.* A brief history of human disease genetics. *Nature* **577**, 179-
797 189 (2020).
798
- 799 23. Ferreira, R.C. *et al.* Functional IL6R 358Ala allele impairs classical IL-6 receptor
800 signaling and influences risk of diverse inflammatory diseases. *PLoS Genet* **9**,
801 e1003444 (2013).
802
- 803 24. Rosa, M. *et al.* A Mendelian randomization study of IL6 signaling in cardiovascular
804 diseases, immune-related disorders and longevity. *NPJ Genom Med* **4**, 23 (2019).
805
- 806 25. Patsopoulos, A. Multiple sclerosis genomic map implicates peripheral immune cells
807 and microglia in susceptibility. *Science* **365**, eaav7188 (2019).
808
- 809 26. Gregory, A.P. *et al.* TNF receptor 1 genetic risk mirrors outcome of anti-TNF therapy
810 in multiple sclerosis. *Nature* **488**, 508-511 (2012).
811
- 812 27. Smith, G.D. & Ebrahim, S. 'Mendelian randomization': can genetic epidemiology
813 contribute to understanding environmental determinants of disease? *Int J Epidemiol*
814 **32**, 1-22 (2003).
815
- 816 28. Zhu, Z. *et al.* Causal associations between risk factors and common diseases
817 inferred from GWAS summary data. *Nature Communications* **9**, 224 (2018).
818
- 819 29. Massimino, L. *et al.* The Inflammatory Bowel Disease Transcriptome and
820 Metatranscriptome Meta-Analysis (IBD TaMMA) framework. *Nature Computational*
821 *Science* **1**, 511-515 (2021).
822
- 823 30. Croft, M. & Siegel, R.M. Beyond TNF: TNF superfamily cytokines as targets for the

- 824 treatment of rheumatic diseases. *Nat Rev Rheumatol* **13**, 217-233 (2017).
- 825
- 826 31. Yasuda, H. Discovery of the RANKL/RANK/OPG system. *J Bone Miner Metab* **39**, 2-
827 11 (2021).
- 828
- 829 32. Walsh, M.C. & Choi, Y. Biology of the RANKL-RANK-OPG System in Immunity,
830 Bone, and Beyond. *Front Immunol* **5**, 511 (2014).
- 831
- 832 33. Jakubowski, A. *et al.* Dual role for TWEAK in angiogenic regulation. *J Cell Sci* **115**,
833 267-274 (2002).
- 834
- 835 34. Donohue, P.J. *et al.* TWEAK is an endothelial cell growth and chemotactic factor that
836 also potentiates FGF-2 and VEGF-A mitogenic activity. *Arterioscler Thromb Vasc*
837 *Biol* **23**, 594-600 (2003).
- 838
- 839 35. Domouzoglou, E.M. *et al.* Fibroblast growth factors in cardiovascular disease: The
840 emerging role of FGF21. *Am J Physiol Heart Circ Physiol* **309**, H1029-1038 (2015).
- 841
- 842 36. Schett, G., McInnes, I.B. & Neurath, M.F. Reframing Immune-Mediated Inflammatory
843 Diseases through Signature Cytokine Hubs. *N Engl J Med* **385**, 628-639 (2021).
- 844
- 845 37. Peters, A.L., Stunz, L.L. & Bishop, G.A. CD40 and autoimmunity: the dark side of a
846 great activator. *Semin Immunol* **21**, 293-300 (2009).
- 847
- 848 38. Durie, F.H. *et al.* Prevention of collagen-induced arthritis with an antibody to gp39,
849 the ligand for CD40. *Science* **261**, 1328-1330 (1993).
- 850
- 851 39. Guo, Y. *et al.* CD40L-Dependent Pathway Is Active at Various Stages of Rheumatoid
852 Arthritis Disease Progression. *J Immunol* **198**, 4490-4501 (2017).
- 853
- 854 40. The Lenercept Multiple Sclerosis Study Group and The University of British Columbia
855 MS/MRI Analysis Group. TNF neutralization in MS: results of a randomized, placebo-
856 controlled multicenter study. The Lenercept Multiple Sclerosis Study Group and The
857 University of British Columbia MS/MRI Analysis Group. *Neurology* **53**, 457-465
858 (1999).
- 859
- 860 41. Bosch, X., Saiz, A., Ramos-Casals, M. & Group, B.S. Monoclonal antibody therapy-
861 associated neurological disorders. *Nat Rev Neurol* **7**, 165-172 (2011).
- 862
- 863 42. Singh, U.P. *et al.* Chemokine and cytokine levels in inflammatory bowel disease
864 patients. *Cytokine* **77**, 44-49 (2016).
- 865
- 866 43. Friedrich, M. *et al.* IL-1-driven stromal-neutrophil interactions define a subset of
867 patients with inflammatory bowel disease that does not respond to therapies. *Nat*
868 *Med* **27**, 1970-1981 (2021).
- 869
- 870 44. Pavlidis, P. *et al.* Interleukin-22 regulates neutrophil recruitment in ulcerative colitis

- 871 and is associated with resistance to ustekinumab therapy. *Nat Commun* **13**, 5820
872 (2022).
873
- 874 45. Richard, A.C. *et al.* Reduced monocyte and macrophage TNFSF15/TL1A expression
875 is associated with susceptibility to inflammatory bowel disease. *PLoS Genet* **14**,
876 e1007458 (2018).
877
- 878 46. Bamias, G. *et al.* Differential expression of the TL1A/DcR3 system of TNF/TNFR-like
879 proteins in large vs. small intestinal Crohn's disease. *Dig Liver Dis* **44**, 30-36 (2012).
880
- 881 47. Bamias, G. *et al.* High intestinal and systemic levels of decoy receptor 3 (DcR3) and
882 its ligand TL1A in active ulcerative colitis. *Clin Immunol* **137**, 242-249 (2010).
883
- 884 48. Sands, B. *et al.* OP40 PRA023 Demonstrated Efficacy and Favorable Safety as
885 Induction Therapy for Moderately to Severely Active UC: Phase 2 ARTEMIS-UC
886 Study Results 2023 [cited 2023 2/7] Available from: [https://www.ecco-](https://www.ecco-ibd.eu/publications/congress-abstracts/item/op40-pra023-demonstrated-efficacy-and-favorable-safety-as-induction-therapy-for-moderately-to-severely-active-uc-phase-2-artemis-uc-study-results.html)
887 [ibd.eu/publications/congress-abstracts/item/op40-pra023-demonstrated-efficacy-and-](https://www.ecco-ibd.eu/publications/congress-abstracts/item/op40-pra023-demonstrated-efficacy-and-favorable-safety-as-induction-therapy-for-moderately-to-severely-active-uc-phase-2-artemis-uc-study-results.html)
888 [favorable-safety-as-induction-therapy-for-moderately-to-severely-active-uc-phase-2-](https://www.ecco-ibd.eu/publications/congress-abstracts/item/op40-pra023-demonstrated-efficacy-and-favorable-safety-as-induction-therapy-for-moderately-to-severely-active-uc-phase-2-artemis-uc-study-results.html)
889 [artemis-uc-study-results.html](https://www.ecco-ibd.eu/publications/congress-abstracts/item/op40-pra023-demonstrated-efficacy-and-favorable-safety-as-induction-therapy-for-moderately-to-severely-active-uc-phase-2-artemis-uc-study-results.html)
890
- 891 49. Fairfax, B.P. *et al.* Innate immune activity conditions the effect of regulatory variants
892 upon monocyte gene expression. *Science* **343**, 1246949 (2014).
893
- 894 50. Lee, M.N. *et al.* Common genetic variants modulate pathogen-sensing responses in
895 human dendritic cells. *Science* **343**, 1246980 (2014).
896
- 897 51. de Lange, K.M. *et al.* Genome-wide association study implicates immune activation
898 of multiple integrin genes in inflammatory bowel disease. *Nat Genet* **49**, 256-261
899 (2017).
900

901 **METHODS**

902 **Cohorts.** We recruited 11 cohorts, totalling 14,824 participants, with genome-wide genetic
903 data and plasma proteomic data measured using the Olink Target Inflammation panel. All
904 participants provided written, informed consent. No statistical methods were used to pre-
905 determine sample sizes but our sample sizes are similar to or larger than those reported in
906 previous publications^{1-4,7-9}. Cohort details are provided in the **Supplementary Note 1**.

907 **Protein assays.** Plasma proteins were measured using the Olink Target-96 Inflammation
908 multiplexed immunoassay panel, which measures 92 inflammation-related proteins.
909 Proteomic data for each cohort were generated at Olink laboratories in Uppsala. During the
910 course of the project, BDNF was removed from the inflammation panel by Olink due to assay
911 problems therefore 91 proteins were included in our study (**Supplementary Table 2**).
912 Normalised Protein eXpression (NPX), is Olink's normalised relative units in log₂ scale. Olink

913 defines the lower limit of detection (LLOD) for quantification of each protein as 3 standard
914 deviations above background (determined using blank control samples), but provides NPX as
915 continuous data which can include values below the calculated LLOD. We had access to
916 individual-level data for INTERVAL, the largest contributing cohort (n= 4,896), and used this
917 to calculate the proportion of samples <LLOD for each protein (**Extended Data Figure 2a**).

918 **Genotyping.** Each cohort was genotyped on a SNP (single nucleotide polymorphism) array
919 and imputed using either a 1000Genomes or Haplotype Reference Consortium (HRC) panel
920 (**Supplementary Table 1**).

921 **Cohort-level pQTL mapping.** In each cohort, a GWAS analysis was run for each protein
922 using linear regression (additive genetic association model) with protein level as the
923 dependent variable. Proteins were inverse-rank normalised prior to linear regression, and thus
924 met the assumptions of the statistical test. Population substructure was adjusted for by
925 including genetic principal components as covariates. We also included age, sex and other
926 study-specific covariates in the model (see **Supplementary Table 1**). To avoid proteins with
927 truncated distributions due to LLOD with multiple tied values that would violate linear
928 regression assumptions, pQTL analysis was performed using continuous protein values
929 (including those below the LLOD where relevant). We illustrate the value of this approach in
930 recovering biological signals in **Extended Data Figure 2b**.

931 **pQTL meta-analysis.** We meta-analysed pQTL summary statistics from each cohort
932 (**Supplementary Table 1**), representing a total of 14,824 participants. A schematic of our
933 analysis pipeline is shown in **Extended Data Figure 1**. Prior to the meta-analysis, we applied
934 cohort-level filters to pQTL GWAS summary statistics with respect to MAF (≥ 0.001), HWE
935 ($P > 10^{-6}$), and imputation score ($r^2 \geq 0.3$ or SNPTEST proper_info ≥ 0.4). For each cohort, we
936 generated QQ plots and Manhattan plots for visual examination using the R packages qqman
937 v0.1.4 and QCGWAS v1.0-8. We performed the fixed-effects meta-analysis with the METAL
938 software (version 28.8.2018), using inverse-variance weighted analysis of regression betas
939 and standard errors from the cohort-level summary statistics. From the meta-analysis
940 summary statistics, we calculated the genomic inflation factor for each protein GWAS, and
941 generated QQ and Manhattan plots (**Supplementary Figure 1**). We generated Forest plots
942 to examine inter-cohort heterogeneity using the gap package v1.2.3-6. Regional association
943 plots were generated using LocusZoom 1.4 (**Supplementary Figure 2**). We defined statistical
944 significance as $P \leq 5 \times 10^{-10}$ (based on Bonferroni correction of the conventional 'genome-wide'
945 significance threshold $P \leq 5 \times 10^{-8}$ for approximately 100 proteins).

946 To remove potentially erroneous meta-analysis signals arising due to a strong association in

947 a single cohort, we examined the meta-analysis results at each sentinel variant by visual
948 inspection of the forest plot and imposed the following criteria: 1) to be included in the meta-
949 analysis, a variant was required to be available in at least three studies and in at least 3,500
950 participants; 2) in order to be declared significant, we required a meta-analysis $P \leq 5 \times 10^{-10}$,
951 and, if there was evidence of heterogeneity with $I^2 > 30\%$, then we required the P -value in at
952 least three studies to be < 0.05 and the direction of effect in those studies to be consistent with
953 the overall meta-analysis results. These were implemented through modification of METAL
954 source code.

955 **Replication cohort.** We compared the results from our primary meta-analysis to pQTL results
956 generated in an independent set of 1,585 participants from the ARISTOTLE study^{12,52}.

957 **Definition of pQTL sentinel variants and regions.** We defined a pQTL as a genetic locus
958 significantly ($P \leq 5 \times 10^{-10}$) associated with protein abundance. We defined the sentinel variant
959 at a locus as the variant with the lowest P -value in the region for a given protein. We used the
960 following approach for each protein to define genomic regions and the sentinel variant in each:
961 1) we first obtained a list of significant ($P \leq 5 \times 10^{-10}$) variants and the flanking region (± 1 Mb)
962 for each variant; 2) overlapping regions were then iteratively merged until no overlapping
963 regions remained; 3) the most significant variant in each resulting region was then defined as
964 the sentinel variant. This approach has the flexibility to cope with long stretches of LD whilst
965 avoiding the drawback of setting a longer than necessary region for all variants. The algorithm
966 was implemented using bedtools v2.27.0. Signals within 1Mb of the transcription start site
967 (TSS) of the gene encoding the target protein were defined as *cis* and those beyond 1 Mb as
968 *trans*.

969 **Protein variance explained by pQTLs.** We used the following equation to estimate the
970 proportion of variance explained (PVE) by (T) pQTLs from the meta-analysis summary
971 statistics for each protein:

$$972 \text{ PVE} = \sum_{i=1}^T \frac{\chi_i^2}{\chi_i^2 + N_i - 2} \quad (1)$$

973 where χ_i^2 is the chi-squared statistic based on each associated variant's effect size and its
974 standard error and N_i is the associated sample size.

975 **Conditional analysis.** To identify conditionally independent signals within a genomic region,
976 we performed approximate stepwise conditional analyses using GCTA v1.93.0beta with the '-
977 -cojo-slc1' option, using effect sizes and standard errors from the meta-analysis. We estimated
978 the correlation between variants using individual-level data from the INTERVAL study. As
979 GCTA imputes LD from mean genotypes when they are missing, to avoid bias we excluded

980 variants with $MAF < 1\%$ (unless they were sentinel variants). For stepwise selection, we
981 considered only those variants passing the genome-wide threshold ($P \leq 5 \times 10^{-10}$), rather than
982 all variants in the region. As in certain cases GCTA conditional analysis yielded results
983 involving pairs of variants in relatively high LD ($r^2 \geq 0.7$), we restricted to independent genetic
984 variants ($r^2 \leq 0.1$)⁵³ in the INTERVAL imputed genotype data whilst forcing the inclusion of the
985 sentinel variants in the pruned set⁵⁴ (**Supplementary Table 4**).

986 **Identification of known pQTLs.** To identify previously reported pQTLs, we manually curated
987 published results from literature obtained from the NCBI web interface
988 (<https://pubmed.ncbi.nlm.nih.gov/>) through its Entrez programming utility R/rentrez⁵⁵,
989 PhenoScanner v2⁵⁶, and the NHGRI-EBI GWAS Catalog with phenotypes mapped to the
990 Experimental Factor Ontology (EFO) EFO_0004747 (protein measurement), restricting to
991 associations reported in European-ancestry populations. We selected variants that reached
992 the conventional genome-wide significance threshold $P \leq 5 \times 10^{-8}$ and that were in high LD
993 ($r^2 \geq 0.8$) with the meta-analysis pQTL sentinel variant.

994 **Variant annotation.** We obtained the absolute distance of sentinel variants to the TSS of the
995 gene encoding the target protein using the rGREAT (Genomic Regions Enrichment of
996 Annotations Tool)⁵⁷ R package. We annotated sentinel variants and LD proxies (defined as
997 $r^2 \geq 0.8$, using the INTERVAL dataset as the LD reference panel) using Ensembl's Variant Effect
998 Predictor (VEP, v98.3) including the LOFTEE plugin.

999 **eQTL-pQTL colocalisation analysis.** We performed pairwise statistical colocalisation
1000 analyses of *cis*-pQTLs identified in the meta-analysis with cognate *cis*-eQTL data from
1001 eQTLGen^{14,58}, the eQTL Catalogue¹⁷, and GTEx v8¹⁵. We extracted the meta-analysis
1002 summary statistics for each *cis*-pQTL sentinel and their +/-1Mb flanking regions, then
1003 extracted the same genomic windows from their cognate *cis*-eQTL data. eQTLGen comprises
1004 eQTL data from 31,684 participants on 19,250 genes that are robustly expressed in blood
1005 (<https://www.eqtlgen.org/cis-eqtls.html>). Of our 59 *cis*-pQTLs, there was genome-wide
1006 significant ($P \leq 5 \times 10^{-8}$) *cis*-eQTL for 40 genes in the eQTLGen data. 1 gene (*TGFB1*) had a *cis*-
1007 eQTL at $FDR < 0.05$ but that was not genome-wide significant ($p = 1.8 \times 10^{-7}$), and 2 had no eQTL
1008 association (*IL17C*, *TNFSF11*). 16 genes had no eQTL data in the eQTLGen summary
1009 statistics, presumably due to lack of robust expression in blood. These were: *CCL11*, *CCL13*,
1010 *CCL19*, *CCL20*, *CCL7*, *CST5*, *CX3CL1*, *CXCL11*, *DNER*, *FGF21*, *FGF5*, *GDNF*, *IL12B*,
1011 *MMP10*, *NGF*, *TNFRSF11B*.

1012 For GTEx v8 and eQTL Catalogue, all 59 *cis*-pQTLs had corresponding eQTL summary
1013 statistics available for colocalisation testing across one or more tissues. We performed
1014 colocalisation analyses using the coloc R package as implemented in v.5.2.2 of the eQTL

1015 Catalogue/colocalisation workflow¹⁷ ([https://github.com/kauralasoo/eQTL-Catalogue-](https://github.com/kauralasoo/eQTL-Catalogue-resources)
1016 [resources](https://github.com/kauralasoo/eQTL-Catalogue-resources)). Coloc returns posterior probabilities indicating the likelihood that the following
1017 scenarios are true: (i) there is no association at the locus with either protein or mRNA (PP0),
1018 (ii) there is an association with either the protein (PP1) or the mRNA (PP2), but not both, (iii)
1019 there is an association with both the protein and the mRNA but with distinct causal variants
1020 (PP3) or with a shared causal variant (PP4). We considered a $PP4 \geq 0.8$ to be robust evidence
1021 of colocalisation between a *cis*-pQTL and its cognate *cis*-eQTL. Since eQTLGen data only
1022 provides allele frequency (*f*) and z-score statistic (*z*) for a particular variant, we obtained the
1023 effect size (*b*) and its standard error (*se*) as follows⁵⁹:

$$1024 \quad b = z/d \quad (2)$$

$$1025 \quad se = 1/d \quad (3)$$

1026 where

$$1027 \quad d = \sqrt{2f(1-f)(z^2 + N)} \quad (4)$$

1028 and *N* is the sample size.

1029 **Prioritizing likely mediating genes at *trans*-pQTLs.** To prioritize likely mediating genes at
1030 *trans*-pQTLs, we used the ProGeM tool¹⁸. To identify *cis*-eQTLs that could mediate our *trans*-
1031 pQTLs, we queried the *trans*-pQTL sentinel variants in eQTLGen¹⁴, the eQTL Catalogue¹⁷,
1032 and the Genotype-Tissue Expression project (GTEx, v8) data. To determine whether the *trans*-
1033 pQTL sentinels are likely to be causal *cis*-eQTL variants in the eQTL Catalogue and GTEx
1034 data, we used the fine-mapped eQTL credible sets available at the eQTL Catalogue
1035 (https://www.ebi.ac.uk/eql/Data_access/). For the eQTLGen data, where credible sets were
1036 not available, we used a manual approach whereby we: (i) first defined a region around each
1037 *trans*-pQTL sentinel variant of +/-500kb; (ii) identified the variant with the lowest *cis*-eQTL *P*-
1038 value in this region for the *cis*-affected gene(s); and (iii) checked to see whether this sentinel
1039 *cis*-eQTL variant is the same sentinel variant for the *trans*-pQTL, or if the two are in high LD
1040 ($r^2 \geq 0.8$).

1041 For the “top-down” component of ProGeM, we first identified locally-encoded genes using a
1042 window around each *trans*-pQTL sentinel variant of +/-500kb. We then probed the proteins
1043 encoded by these local genes using: (1) protein:protein interaction (PPI) data; and (2) data
1044 from functional annotation databases. With the PPI data we aimed to determine whether there
1045 is evidence to indicate that genes residing close to each sentinel variant might physically
1046 interact with the corresponding *trans*-affected protein. We used the Bioconductor package
1047 STRINGdb (v2.8.4) to identify any pairwise interactions. We used data from functional
1048 annotation databases to determine whether any local genes encode proteins that might be
1049 functionally related to the corresponding *trans*-affected protein(s). For both the *trans*-affected
1050 proteins and locally encoded proteins, all assigned Gene Ontology terms, Reactome
1051 pathways, and KEGG pathways were extracted using the Bioconductor biomaRt (v2.52) and

1052 KEGGREST (v1.36) packages. To assess whether there is significant overlap between the
1053 functional annotation terms/pathways assigned to locally encoded proteins and the
1054 corresponding *trans*-affected proteins, we determined the number of shared and non-shared
1055 terms for each local gene and the corresponding *trans*-affected protein. Fisher's exact test
1056 was then applied for each local gene/*trans*-protein combination, and *P*-values were
1057 Bonferroni-corrected for the number of local genes at each given *trans*-pQTL. The background
1058 set of terms for each *trans*-pQTL was composed of all terms assigned to all local genes at the
1059 locus (i.e., all protein-coding genes within 500kb from the sentinel variant).

1060 To determine the most likely mediating genes for the multi-locus regulated proteins IL12B,
1061 KITLG, and TNFSF10 (TRAIL), we used the STRINGdb webtool to identify interactions or
1062 functional relationships between genes residing at distinct loci. This is based on the
1063 assumption that if the mediating genes at distinct loci are all associated with plasma levels of
1064 the same protein, then they may share some other functional relationship. As input to
1065 STRINGdb, we used all proteins encoded by candidate mediating genes identified by ProGeM
1066 (**Supplementary Table 9**) at each of the loci for a given protein, as well as the relevant *trans*-
1067 affected protein. We deemed clusters of proteins residing at distinct loci with multiple functional
1068 interactions to be the most likely mediating genes at their respective loci. We performed a
1069 phenome-scan of the *trans*-pQTLs for *KITLG* using the Open Targets Genetics webtool⁶⁰.

1070 **Overlap of pQTL and disease traits.** We used a PhenoScanner v2-based R code to look up
1071 associations of our pQTL sentinels and their LD proxies ($r^2 \geq 0.8$) in disease GWAS summary
1072 statistics.

1073 To investigate potential colocalisation between a *trans*-pQTL (rs2364485) for LTA identified in
1074 our meta-analysis, a multiple sclerosis GWAS signal²⁵, and a *cis*-eQTL for *LTBR* from
1075 eQTLGen¹⁴, we used HyPrColoc for multi-trait colocalisation⁶¹. We obtained multiple sclerosis
1076 summary statistics (MSchip, "discovery_metav3.0.meta.gz") from Patsopoulos *et al* by request
1077 to the International Multiple Sclerosis Genetics Consortium (IMSGC). Due to a lack of
1078 genotype coverage at the *LTBR/TNFRSF1A* locus in the extended and replication samples
1079 from Patsopoulos *et al*, we selected the summary statistics from the "discovery" sample
1080 ($n=41,505$) for colocalisation analyses, not the full meta-analysis. As a result, the *P*-value for
1081 association between the variant of interest (rs2364458) and multiple sclerosis in the discovery
1082 subset ($P=5.78 \times 10^{-6}$, logistic regression) was higher than reported in Patsopoulos *et al*²⁵
1083 ($P=2.0 \times 10^{-20}$, fixed-effects meta-analysis). We then extracted summary statistics for
1084 rs2364458 (+/-1Mb) (chr12:5514963-7514963) from each of the 3 datasets, and performed
1085 conditional analyses to adjust for any independent signals at the locus using GCTA-COJO.
1086 We ran this using a two-step approach: we first used the COJO-slc function to identify

1087 independent signals at the locus, and then for datasets with independent signals (i.e., in
1088 addition to rs2364485), we used COJO-cond to generate conditioned summary statistics for
1089 use in HyPrColoc. HyPrColoc returns the posterior probability that 2 or more traits colocalise,
1090 akin to PP4 from coloc. We considered a $PP \geq 0.8$ as robust evidence of colocalisation between
1091 traits.

1092 **Mendelian randomisation analyses.** We performed MR analyses using the proteins with *cis*-
1093 pQTLs identified in this meta-analysis as exposures, and immune-mediated diseases (IMDs)
1094 as outcomes. All MR analyses were run using the Generalized Summary-data-based
1095 Mendelian Randomisation (GSMR) method²⁸, which implements two-sample MR accounting
1096 for correlation between variants. For each protein analysed, we defined a +/- 1 Mb window
1097 around the gene encoding it and extracted pQTL summary statistics for this region. For
1098 outcome data, we downloaded GWAS summary statistics for IMDs from OpenGWAS
1099 (<https://gwas.mrcieu.ac.uk/datasets/>) or from the GWAS Catalog
1100 (<https://www.ebi.ac.uk/gwas/downloads>) where studies with larger sample sizes or more
1101 variants were available. For IMDs with several alternative datasets available, we selected the
1102 dataset with the largest number of cases, provided it (i) had genotype data with sufficient
1103 coverage at the loci of interest, (ii) was generated in European-ancestry samples so that it
1104 matched the ancestry of the participants in our pQTL meta-analysis, and (ii) had effect
1105 estimates and standard errors either available or calculable. Proteins encoded by genes in the
1106 *HLA* region were excluded because MR analysis would be confounded by complex LD. The
1107 analysis involved 57 proteins and 14 diseases. We used the GSMR implementation in GCTA
1108 with the following parameters: 1) at least 3 (--gsmr-snp-min 3) genome-wide significant (--
1109 gwas-thresh 5e-8), quasi-independent variants (--clump-r2 0.1); 2) difference in the allele
1110 frequency of each effect allele between the GWAS summary datasets and the LD reference
1111 sample is at most 0.4 (--diff-freq 0.4); 3) a *P*-value threshold of 0.05 for the HEIDI-outlier
1112 filtering analysis (--heidi-thresh 0.05), which is used to identify potential confounding by LD
1113 (<https://yanglab.westlake.edu.cn/software/gcta/#Mendelianrandomisation>). *P*-values were
1114 corrected for the number of models tested using Benjamini-Hochberg correction, with
1115 $FDR < 0.01$ used to define statistical significance.

1116 To evaluate the robustness of significant associations, we performed additional checks. First,
1117 we checked the strength of the disease association in the GWAS summary statistics. Of the
1118 22 significant protein-disease MR associations, we eliminated 5 due to the lack of convincing
1119 disease association (smallest *P*-value at the locus $P > 1 \times 10^{-4}$). For the remaining 17 MR
1120 associations, we then evaluated whether there might be confounding due to LD. We first
1121 evaluated the r^2 between the sentinel pQTL and the disease-associated variant. For 12 of 17
1122 disease-protein pairs, r^2 was > 0.8 (**Supplementary Table 14**). We next performed visual

1123 inspection of regional association plots of these 12 pQTL–disease pairs (**Supplementary**
1124 **Figure 4**) and colocalisation testing using PWCoCo^{62,63} which accounts for the presence of
1125 multiple independent signals within a locus (see below).

1126 **Pairwise conditional colocalisation analyses.** PairWise COnditional and COlocalisation
1127 analysis (PWCoCo)^{62,63} integrates conditional analyses (GCTA-COJO) to identify independent
1128 signals for each of two tested traits associated with a genomic region, followed by pairwise
1129 colocalisation analyses (COLOC) to test all possible pairs of conditionally independent signals
1130 across the traits. We ran PWCoCo for the 12 significant protein-disease pairs that resulted
1131 from our MR filtering steps using the default parameters, detailed as follows: 1) *P*-value cutoff
1132 for variants to be selected by the stepwise selection process, --p_cutoff 5e-8 for disease and
1133 protein summary statistics; 2) a large number of variants subject to the stepwise selection
1134 process, --top_snp 1e10; 3) distance beyond which variants are treated as in linkage
1135 equilibrium, --ld_window 1e7 (kb); 4) collinearity threshold for variants, --collinear 0.9; 5)
1136 variant frequency filter for the reference dataset according to this threshold, --maf 0.1; 6)
1137 exclusion threshold for variants with allele frequency difference between the phenotype and
1138 the reference datasets, --freq_threshold 0.2; 7) stop criteria, --init_h4 80 (i.e. 80%); 8) the three
1139 prior probabilities, --coloc_pp 1e-4, 1e-4 and 1e-5.

1140 **CXCL5 differential expression analysis in UC cohorts.** Changes in *CXCL5* gene
1141 expression levels were evaluated in four independent cohorts, including the IBD
1142 Transcriptome and Metatranscriptome Meta-Analysis (TaMMA) platform²⁹, the Gene
1143 Expression Omnibus (GEO) series GSE16879, GSE206285, as well as the Imperial UC
1144 cohort. IBD TaMMA (<https://ibd-meta-analysis.herokuapp.com/>) gives access to 3,853
1145 transcriptomic profiles from 26 independent studies including IBD and control samples across
1146 different tissues, all processed with the same pipeline and batch-corrected²⁹. Pre-computed
1147 differential expression results between colon biopsies from UC patients versus healthy donors
1148 were downloaded and plotted.

1149
1150 Data from GEO Series GSE16879 used in this study consist of colonic mucosa microarray
1151 expression profiles from healthy donors (n=6) and UC patients (n=24) sampled before first
1152 infliximab treatment⁶⁴. CEL file import into R, and background correction, normalization and
1153 RMA calibration of the raw intensity data were carried out using the oligo package. Only probe
1154 sets with median expression greater than 4 and associated to only one ENTREZ gene
1155 identifier were kept for analysis. Intensity data for different probe sets mapped to the same
1156 ENTREZ gene identifier were combined by taking the geometric mean sample wise. Tests of
1157 differential gene expression of ulcerative colitis samples compared to healthy control samples

1158 were performed with the limma package. *P*-values were adjusted for multiple testing with the
1159 Benjamini and Hochberg procedure.

1160

1161 GEO Series GSE206285 contains array-based transcriptomic data collected at baseline as
1162 part of UNIFI, a randomised placebo-controlled phase 3 clinical trial evaluating the efficacy
1163 and safety of ustekinumab⁶⁵. RMA signal intensity profiles and associated donor information
1164 were downloaded from NCBI GEO. Only probe sets associated to only one ENTREZ gene
1165 identifier were kept for analysis. Intensity data for different probe sets mapped to the same
1166 ENTREZ gene identifier were combined by taking the geometric mean sample wise. Genes
1167 with median expression greater than 3 across all samples were tested for differential
1168 expression between ulcerative colitis samples (n=550) versus healthy control samples (n=18)
1169 using the R limma package. *P*-values were adjusted for multiple testing with the Benjamini
1170 and Hochberg procedure.

1171

1172 The Imperial UC cohort includes whole tissue biopsies from ulcerative colitis patients (n=16)
1173 and healthy volunteers (n=6). RNA was extracted (Qiagen RNeasy mini kit) and sequencing
1174 libraries were generated using NEBNext® Ultra™ RNA Library Prep Kit for Illumina® (NEB,
1175 USA) following manufacturer's recommendations. Briefly, mRNA was purified from total RNA
1176 using poly-T oligo-attached magnetic beads. Fragmentation was carried out using divalent
1177 cations under elevated temperature in NEBNext First Strand Synthesis Reaction Buffer (5X).
1178 First strand cDNA was synthesised using random hexamer primer and M-MuLV Reverse
1179 Transcriptase (RNase H-). Second strand cDNA synthesis was subsequently performed
1180 using DNA Polymerase I and RNase H. Remaining overhangs were converted into blunt ends
1181 via exonuclease/polymerase activities. After adenylation of 3' ends of DNA fragments,
1182 NEBNext Adaptor with hairpin loop structure were ligated to prepare for hybridization. Library
1183 fragments were purified with AMPure XP system (Beckman Coulter, Beverly, USA) and
1184 treated with 3 µl USER Enzyme (NEB, USA) at 37°C for 15 min, followed by 5 min at 95 °C.
1185 Then PCR was performed with Phusion High-Fidelity DNA polymerase, Universal PCR
1186 primers and Index (X) Primer. Library quality was assessed on Agilent Bioanalyzer 2100 and
1187 on Nanodrop ND-1000 Spectrophotometer. The library preparations were sequenced on an
1188 Illumina HiSeq platform, generating 150 bp paired end reads. The resulting fastq files were
1189 processed with trimmomatic⁶⁶ (v. 0.39) to remove adaptor contamination and poor-quality
1190 bases. The output read files were mapped to the GRCh38 assembly of the human genome
1191 using Hisat2⁶⁷ (v. 2.2.1) with default parameters. The number of reads mapping to the genomic
1192 features annotated in Ensembl with a MAPQ score higher than or equal to 10 was calculated
1193 for all samples using htseq-count⁶⁸(v. 0.11.3) with default parameters. Data for Ensembl genes

1194 with no associated ENTREZ gene identifier were discarded; the read counts for Ensembl
1195 genes mapped to the same ENTREZ gene identifier were summed up sample wise.
1196 Differential expression analysis between ulcerative colitis versus healthy biopsies was
1197 performed in R (v. 3.6.1) using the Wald test as implemented in DESeq2. Only genes with an
1198 average read count across samples greater than or equal to 10 were tested for differential
1199 expression. *P*-values were adjusted for multiple testing with the Benjamini-Hochberg
1200 procedure.

1201

1202 **Data availability.** Full per-protein GWAS summary statistics are available for download at
1203 <https://www.phpc.cam.ac.uk/ceu/proteins/> and the EBI GWAS Catalog
1204 <https://www.ebi.ac.uk/gwas/studies/GCST90270765> (accession numbers GCST90270765-
1205 GCST90270855). Individual-level genetic and proteomic data available for the INTERVAL
1206 cohort are deposited in the European-Genome Phenome Archive under the accession code
1207 <https://ega-archive.org/studies/EGAS00001002555>. Gene expression data are in Gene
1208 Expression Omnibus (GEO) under the accession code GSE16879 for mucosal expression in
1209 IBD patients (<https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE16879>), GSE206285
1210 for the UNIFI trial (<https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE206285>), and
1211 the IBD Transcriptome and Metatranscriptome Meta-Analysis (TaMMA - [https://ibd-meta-
1212 analysis.herokuapp.com/](https://ibd-meta-analysis.herokuapp.com/)). Whole blood *cis*-eQTL summary statistics from the eQTLGen
1213 Consortium were downloaded from <https://www.eqtlgen.org/cis-eqtls.html>. Fine-mapped
1214 eQTL credible sets were downloaded from the eQTL Catalogue
1215 (https://www.ebi.ac.uk/eqtl/Data_access/). MR GWAS summary statistics for IMDs were
1216 downloaded from OpenGWAS (<https://gwas.mrcieu.ac.uk/datasets/>) or from the GWAS
1217 Catalog (<https://www.ebi.ac.uk/gwas/downloads>).

1218

1219 **Code availability.** GitHub: <https://jinghuazhao.github.io/INF/>, cambridge-ceu:
1220 <https://cambridge-ceu.github.io/public> (modified METAL, pQTLtools)

1221 **Methods only REFERENCES**

1222

1223 52. Hijazi, Z. *et al.* Screening of Multiple Biomarkers Associated With Ischemic Stroke in
1224 Atrial Fibrillation. *J Am Heart Assoc* **9**, e018984 (2020).

1225

1226 53. Sanna, S. *et al.* Causal relationships among the gut microbiome, short-chain fatty
1227 acids and metabolic diseases. *Nature Genetics* **51**, 600-605 (2019).

1228

1229 54. Yang, J. *et al.* Conditional and joint multiple-SNP analysis of GWAS summary
1230 statistics identifies additional variants influencing complex traits. *Nature Genetics* **44**,
1231 369-375 (2012).

1232

1233 55. Winter, D.J. rentrez: an R package for the NCBI eUtils API. *The R Journal* **9**, 520-526
1234 (2017).

1235

1236 56. Kamat, M.A. *et al.* PhenoScanner V2: an expanded tool for searching human
1237 genotype-phenotype associations. *Bioinformatics* **35**, 4851-4853 (2019).

1238

1239 57. McLean, C.Y. *et al.* GREAT improves functional interpretation of cis-regulatory
1240 regions. *Nature Biotechnology* **28**, 495-501 (2010).

1241

1242 58. Vösa, U. *et al.* Unraveling the polygenic architecture of complex traits using blood
1243 eQTL metaanalysis. *bioRxiv*, 447367 (2018).

1244

1245 59. Zhu, Z. *et al.* Integration of summary data from GWAS and eQTL studies predicts
1246 complex trait gene targets. *Nat Genet* **48**, 481-487 (2016).

1247

1248 60. Ochoa, D. *et al.* Open Targets Platform: supporting systematic drug-target
1249 identification and prioritisation. *Nucleic Acids Res* **49**, D1302-D1310 (2021).

1250

1251 61. Foley, C.N. *et al.* A fast and efficient colocalization algorithm for identifying shared
1252 genetic risk factors across multiple traits. *Nature Communications* **12**, 764 (2021).

1253

1254 62. Zheng, J. *et al.* Phenome-wide Mendelian randomization mapping the influence of
1255 the plasma proteome on complex diseases. *Nature Genetics* **52**, 1122-1131 (2020).

1256

1257 63. Robinson, J., W. *et al.* An efficient and robust tool for colocalisation: Pair-wise
1258 Conditional and Colocalisation (PWCoCo). *bioRxiv*, 2022.2008.2008.503158 (2022).

1259

1260 64. Arijs, I. *et al.* Mucosal gene expression of antimicrobial peptides in inflammatory
1261 bowel disease before and after first infliximab treatment. *PLoS One* **4**, e7984 (2009).

1262

1263 65. Sands, B.E. *et al.* Ustekinumab as Induction and Maintenance Therapy for Ulcerative
1264 Colitis. *N Engl J Med* **381**, 1201-1214 (2019).

1265

1266 66. Bolger, A.M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for Illumina
1267 sequence data. *Bioinformatics* **30**, 2114-2120 (2014).

1268

1269 67. Kim, D., Langmead, B. & Salzberg, S.L. HISAT: a fast spliced aligner with low
1270 memory requirements. *Nat Methods* **12**, 357-360 (2015).

1271

1272 68. Anders, S., Pyl, P.T. & Huber, W. HTSeq--a Python framework to work with high-
1273 throughput sequencing data. *Bioinformatics* **31**, 166-169 (2015).

1274

1275

1276

1277 Estonian Biobank Research Team:

1278

1279 Andres Metspalu⁹

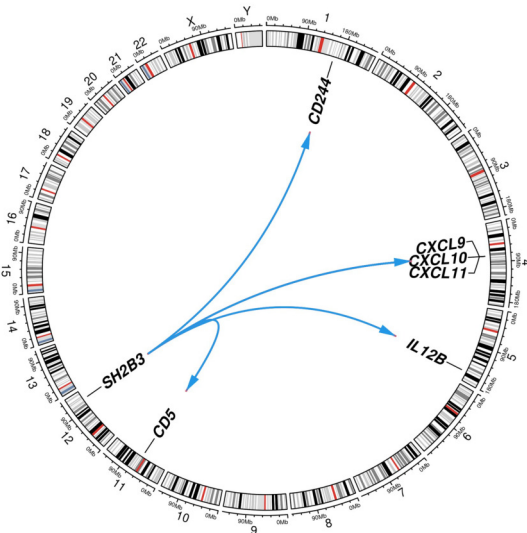
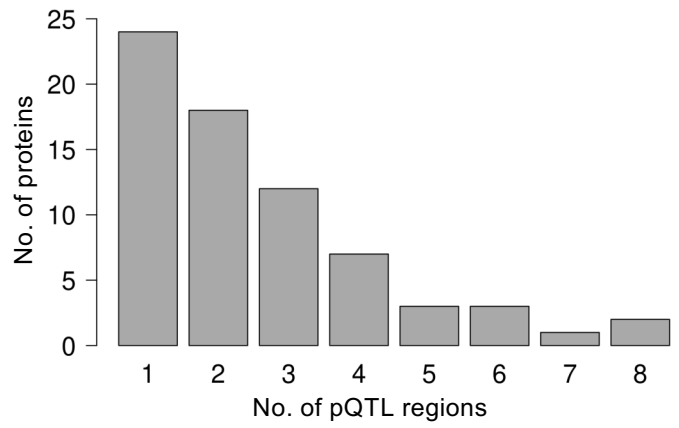
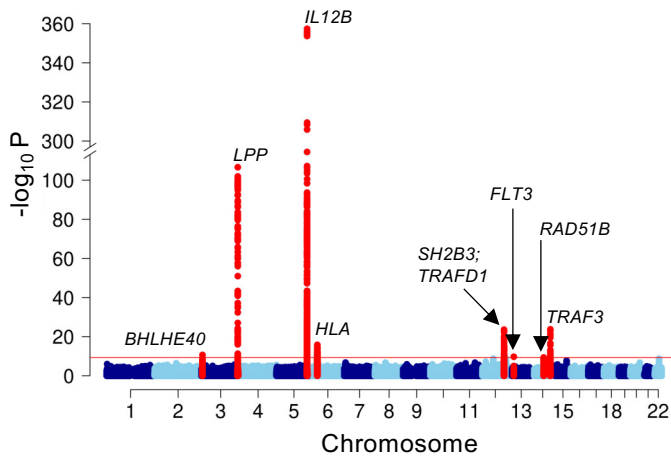
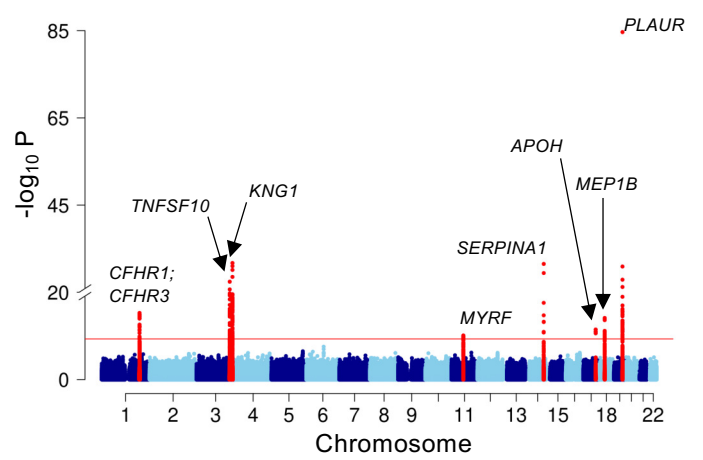
1280 Lili Milani⁹

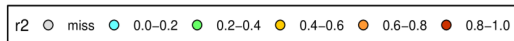
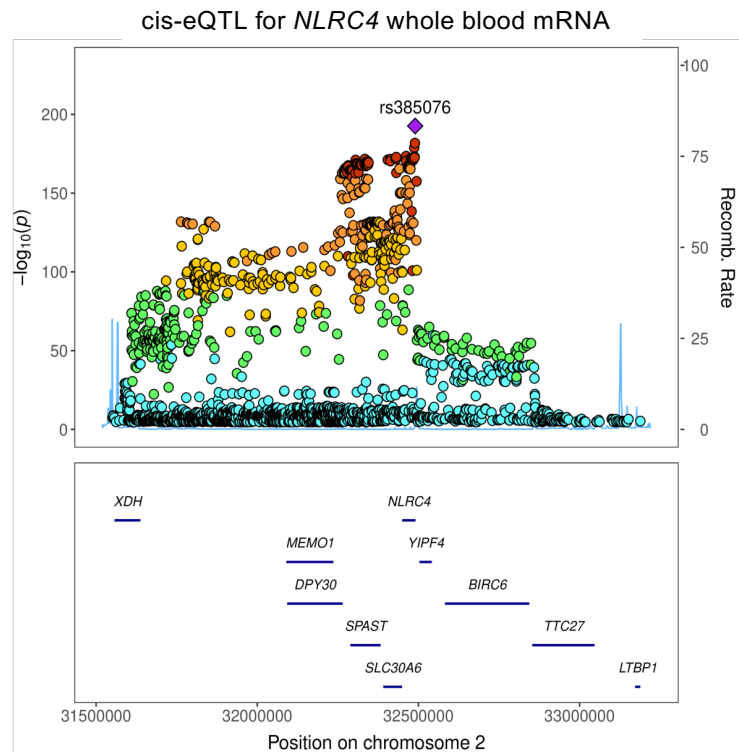
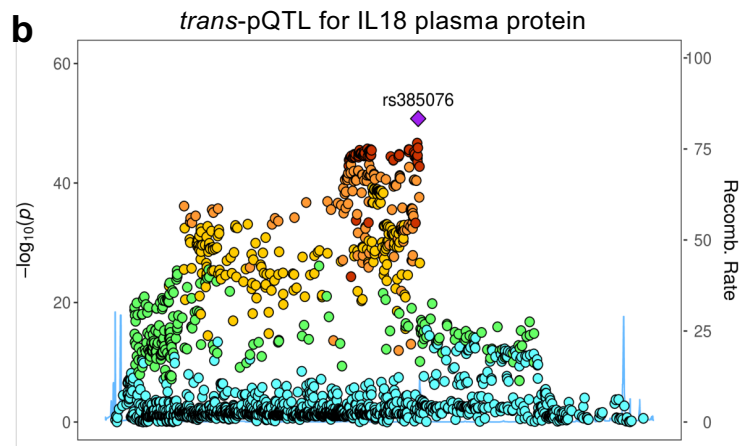
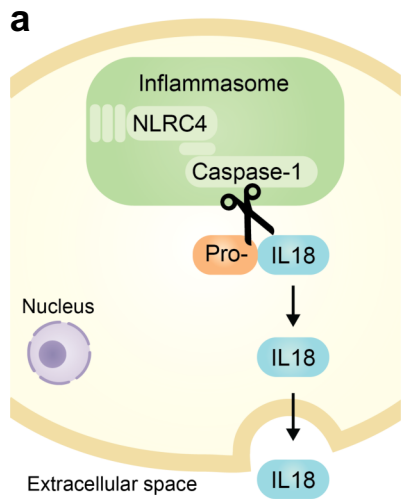
1281 Reedik Mägi⁹

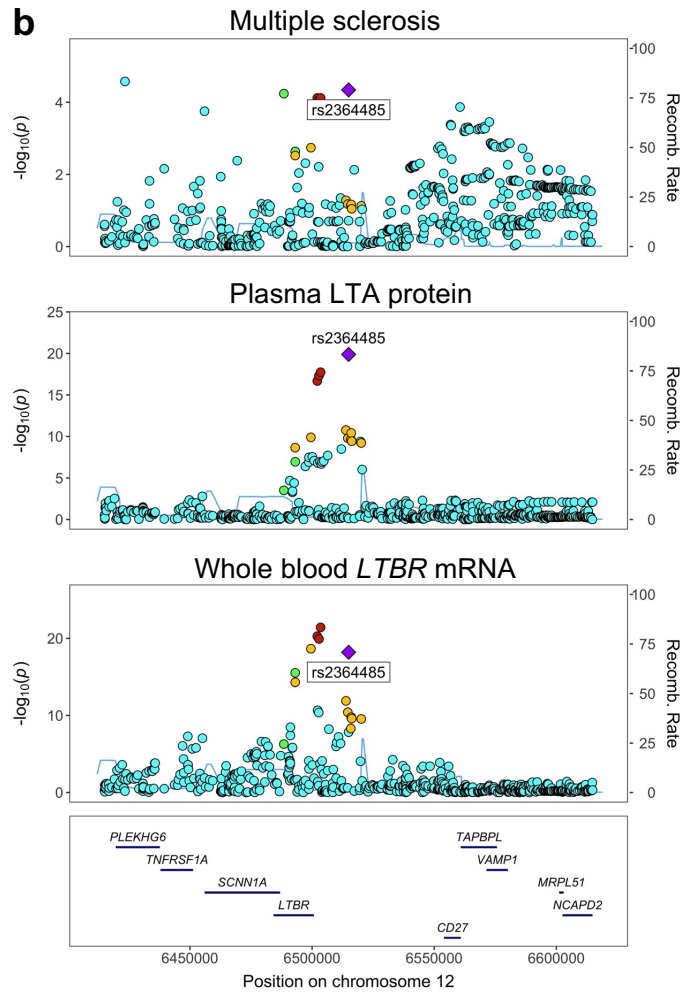
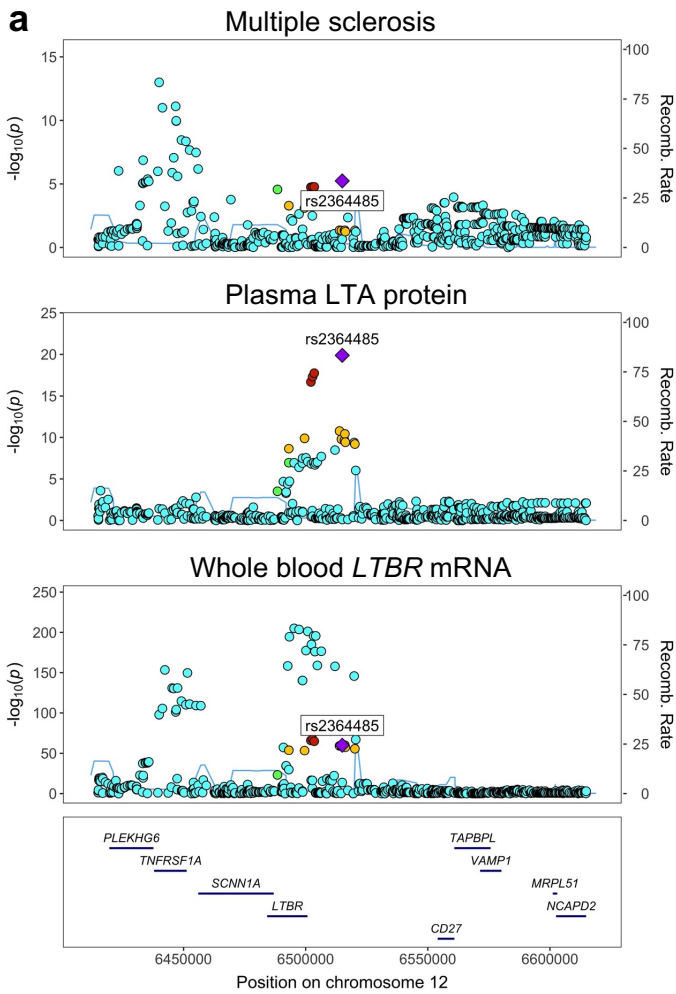
1282 Mari Nelis⁹

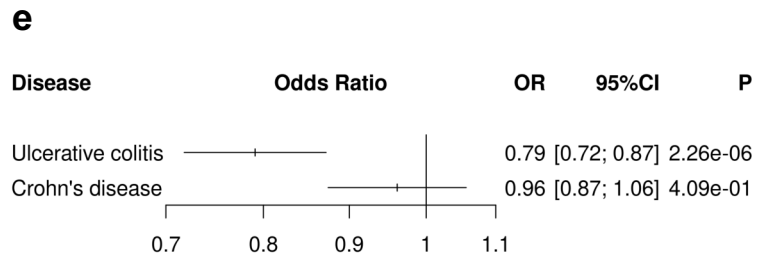
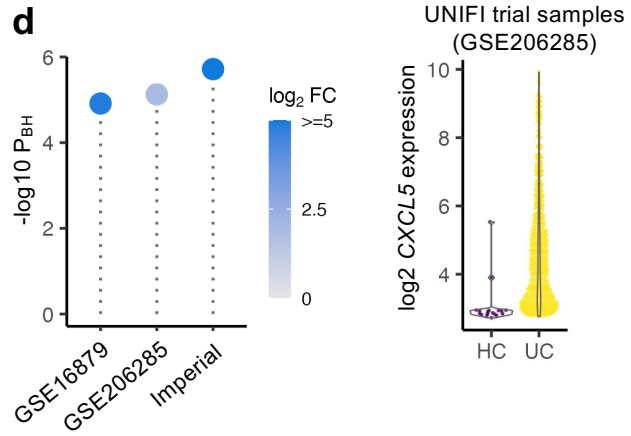
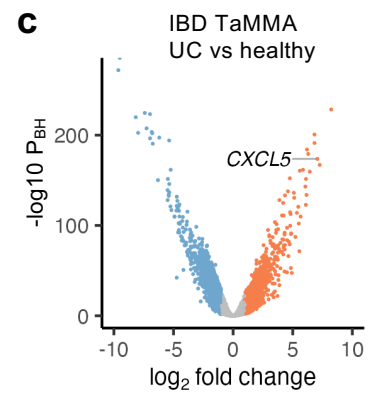
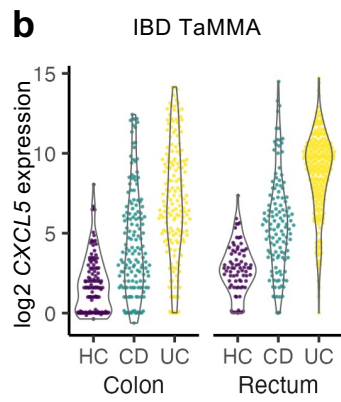
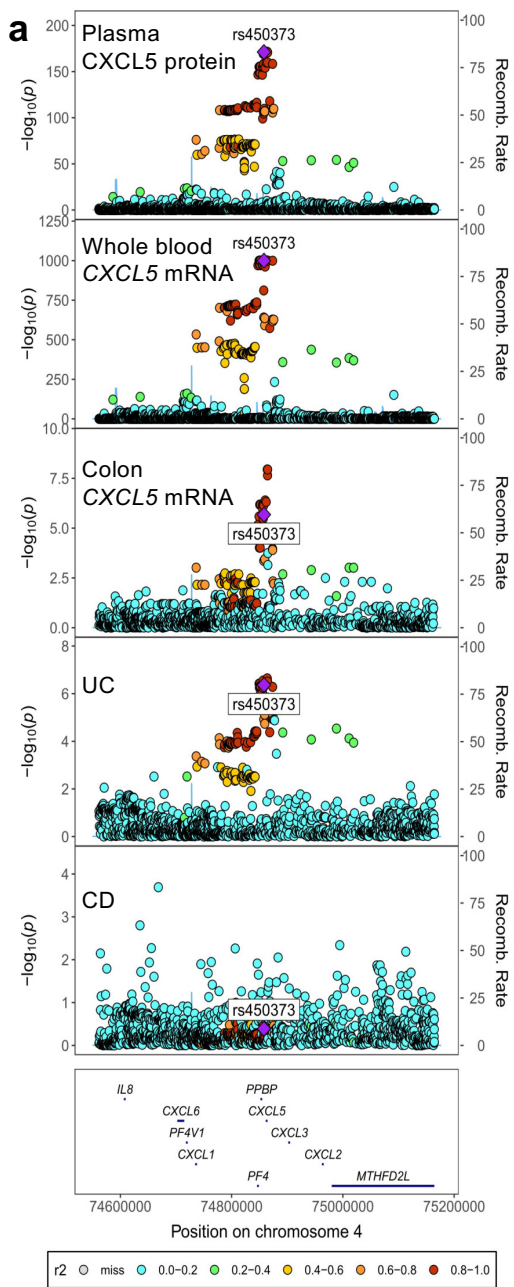
1283 Georgi Hudjašov⁹

1284

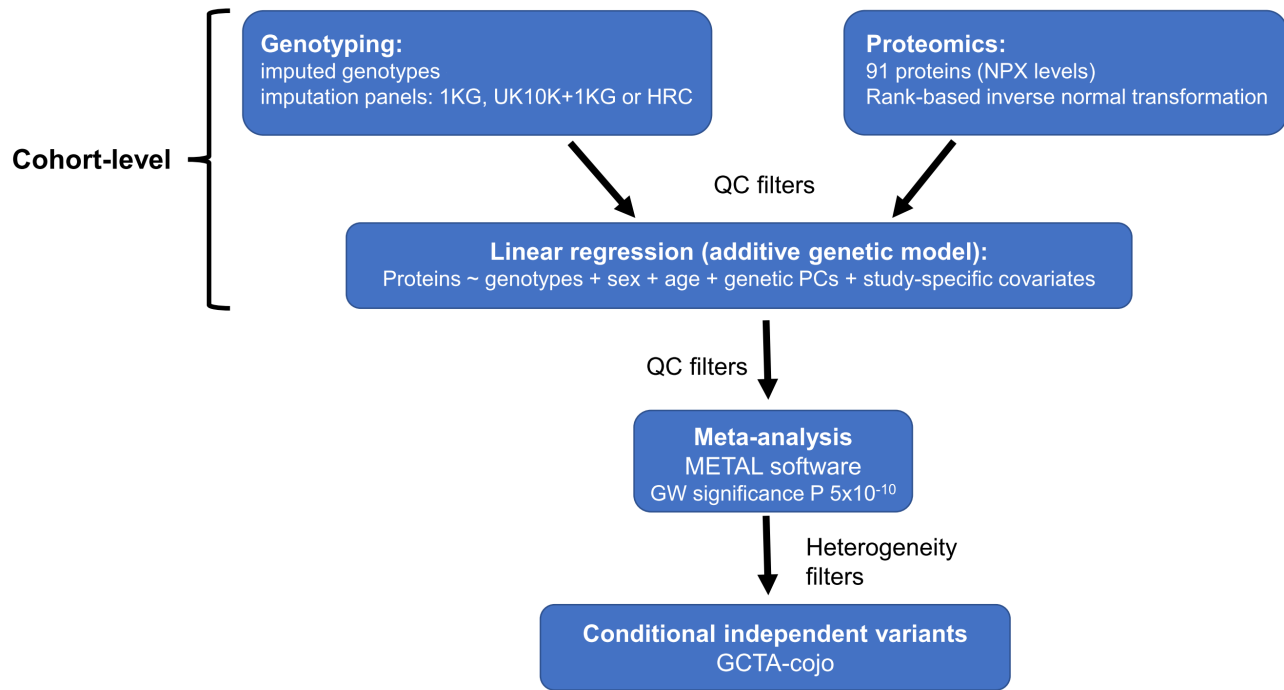
a**b****c****d**

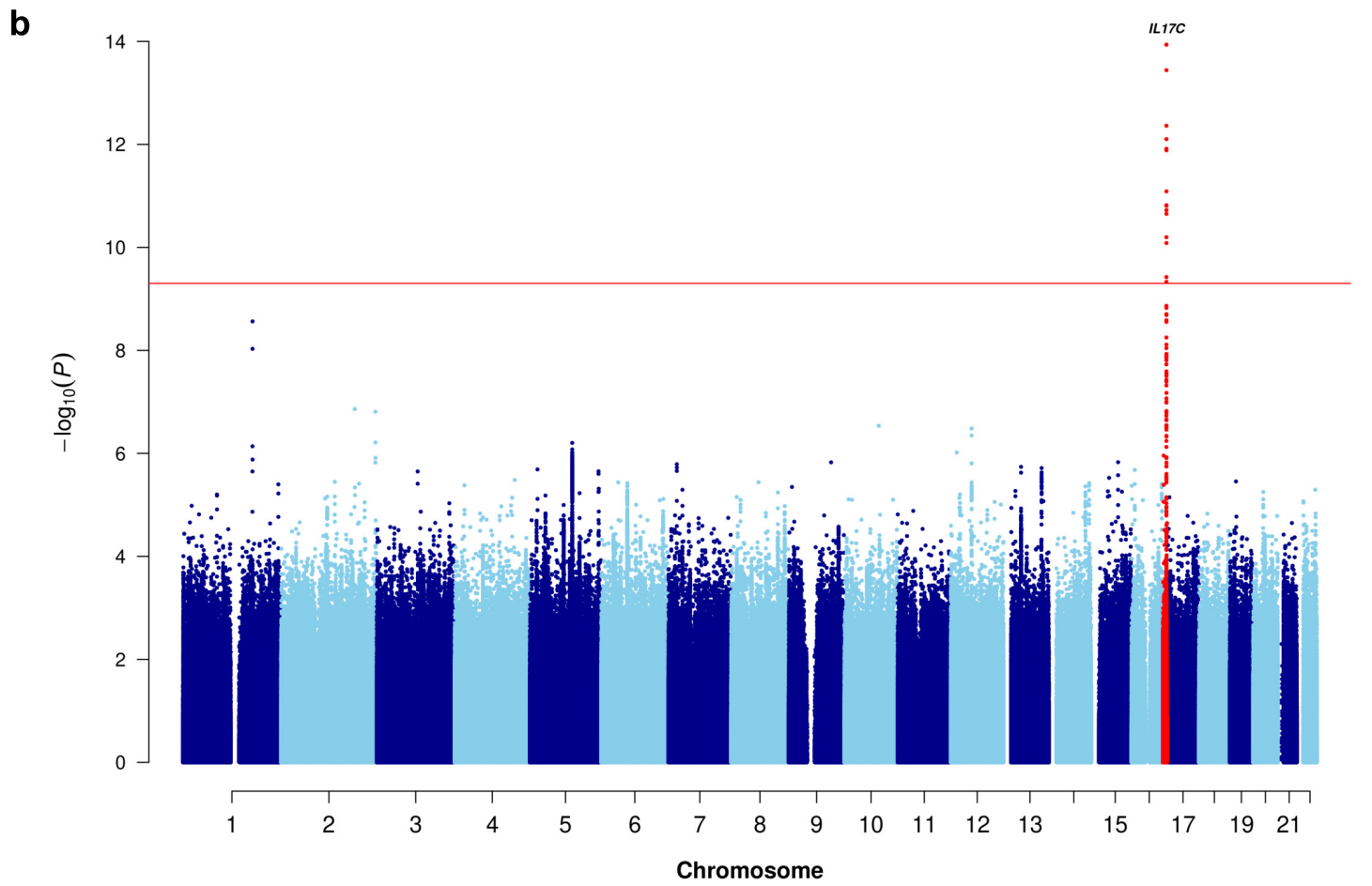
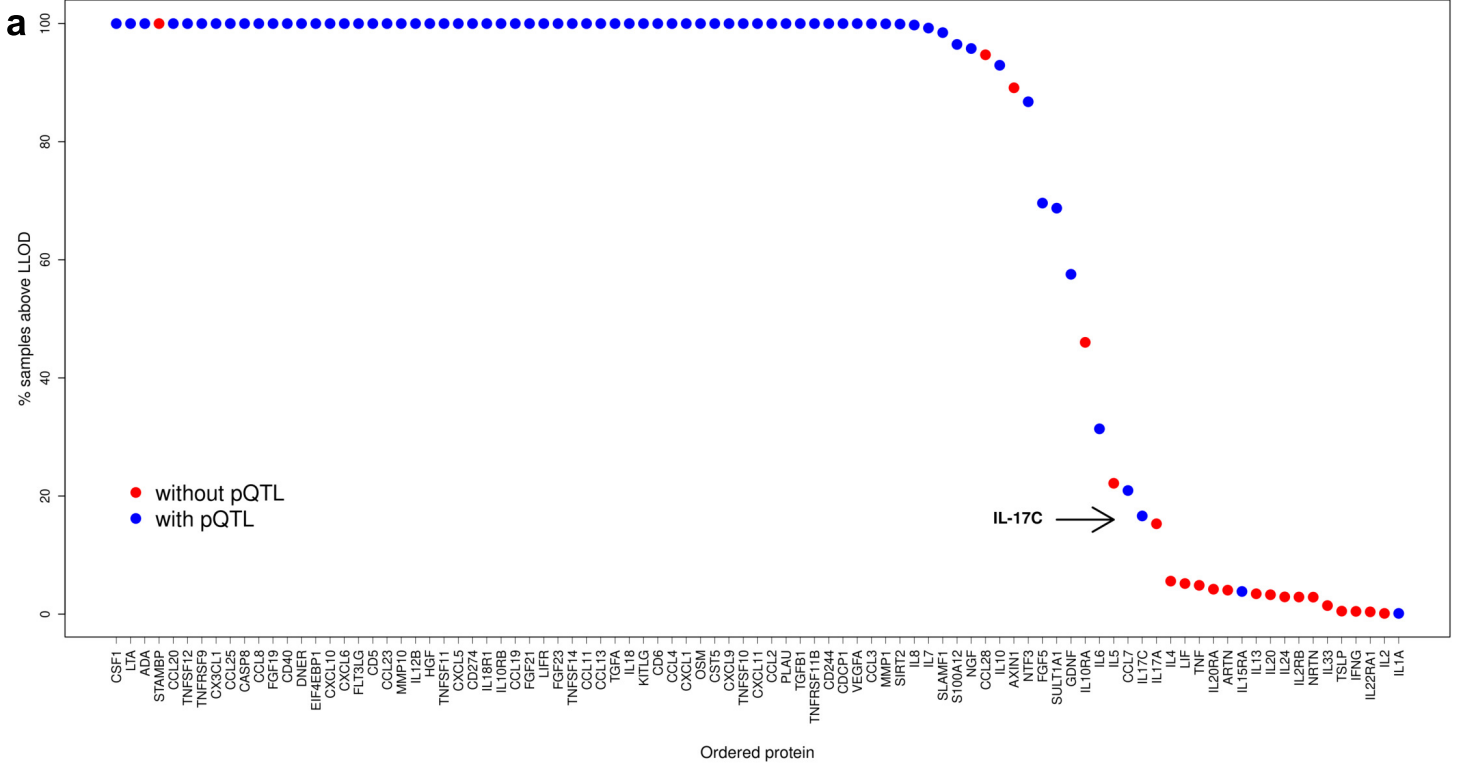


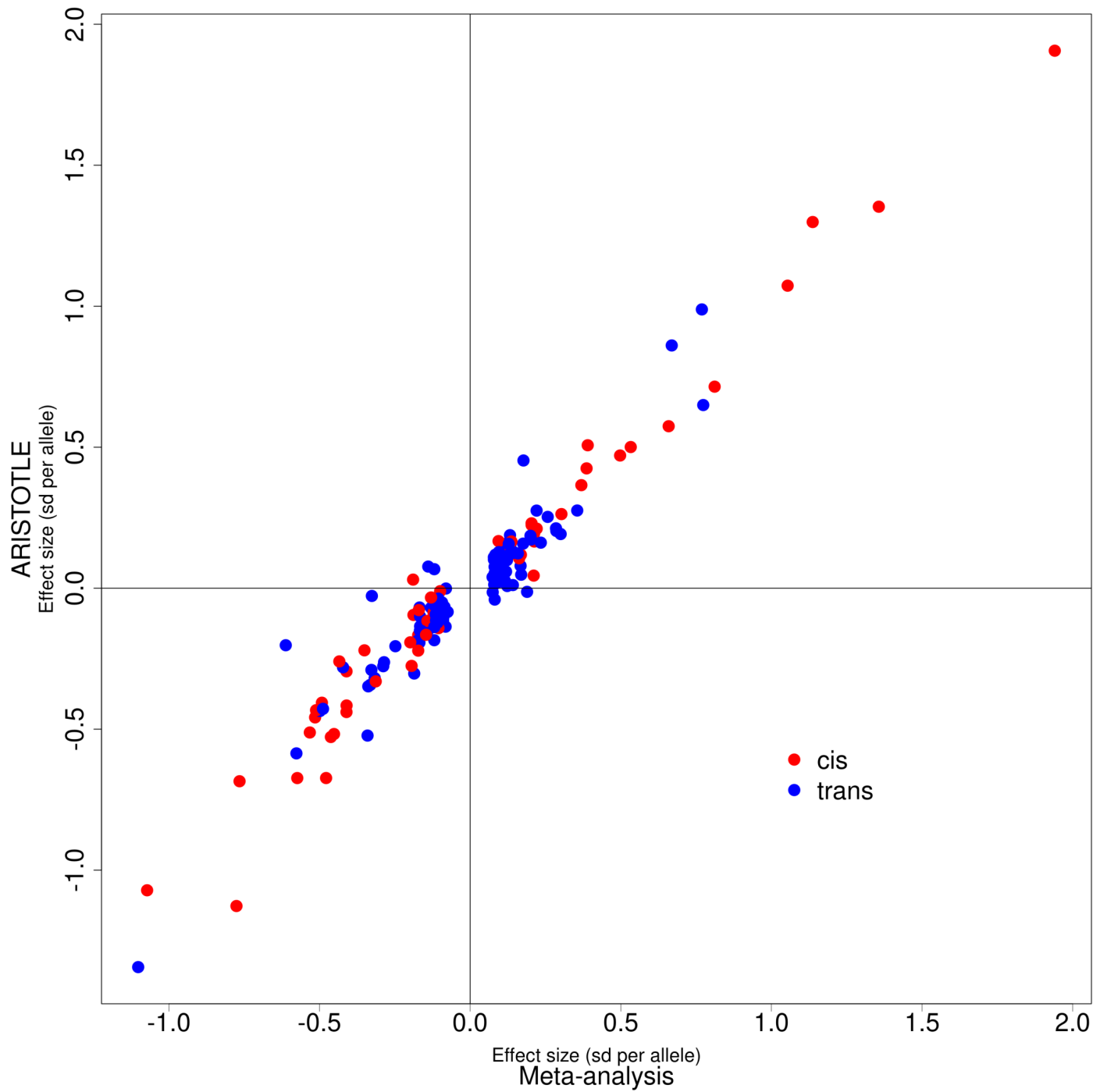


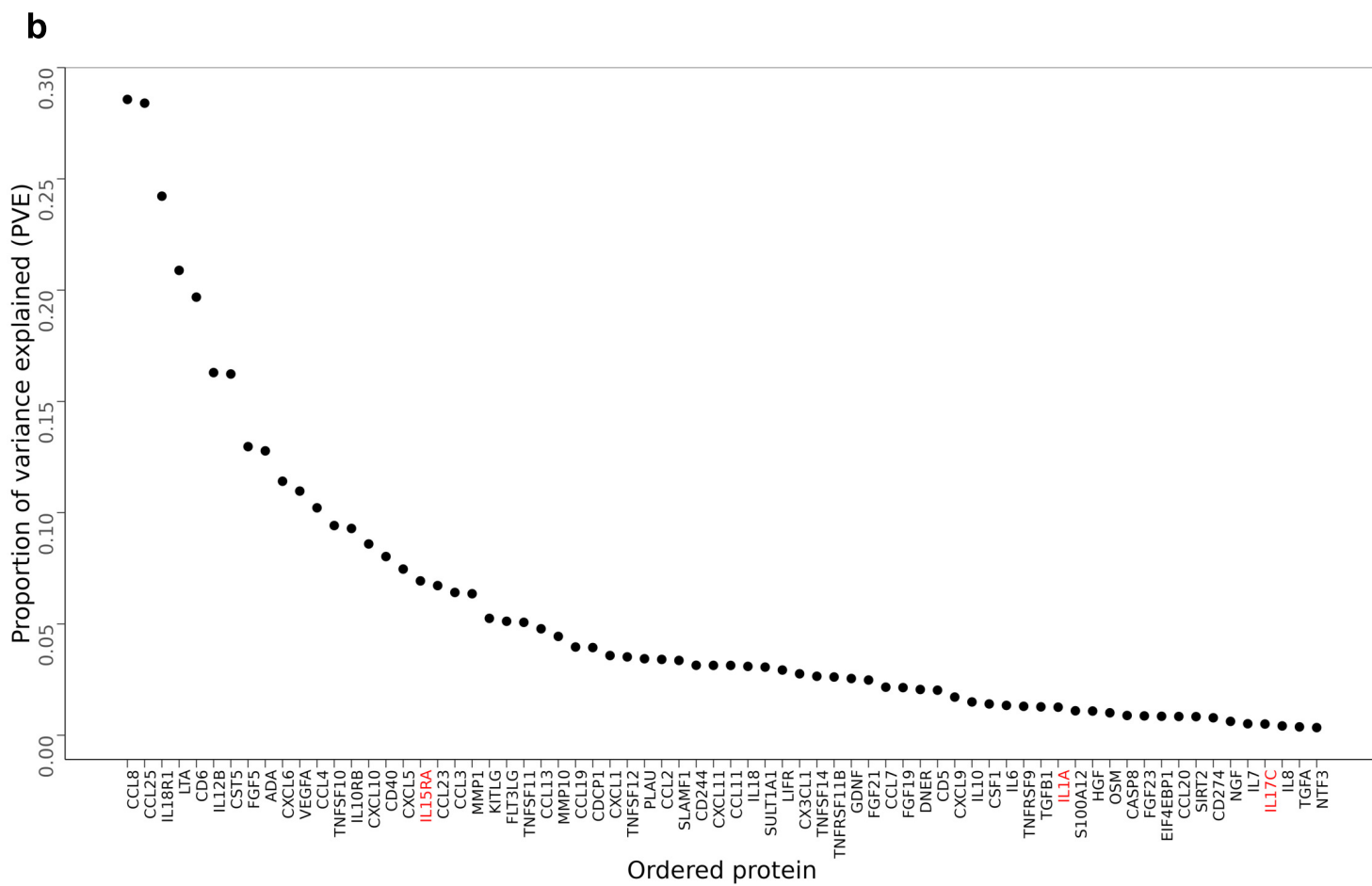
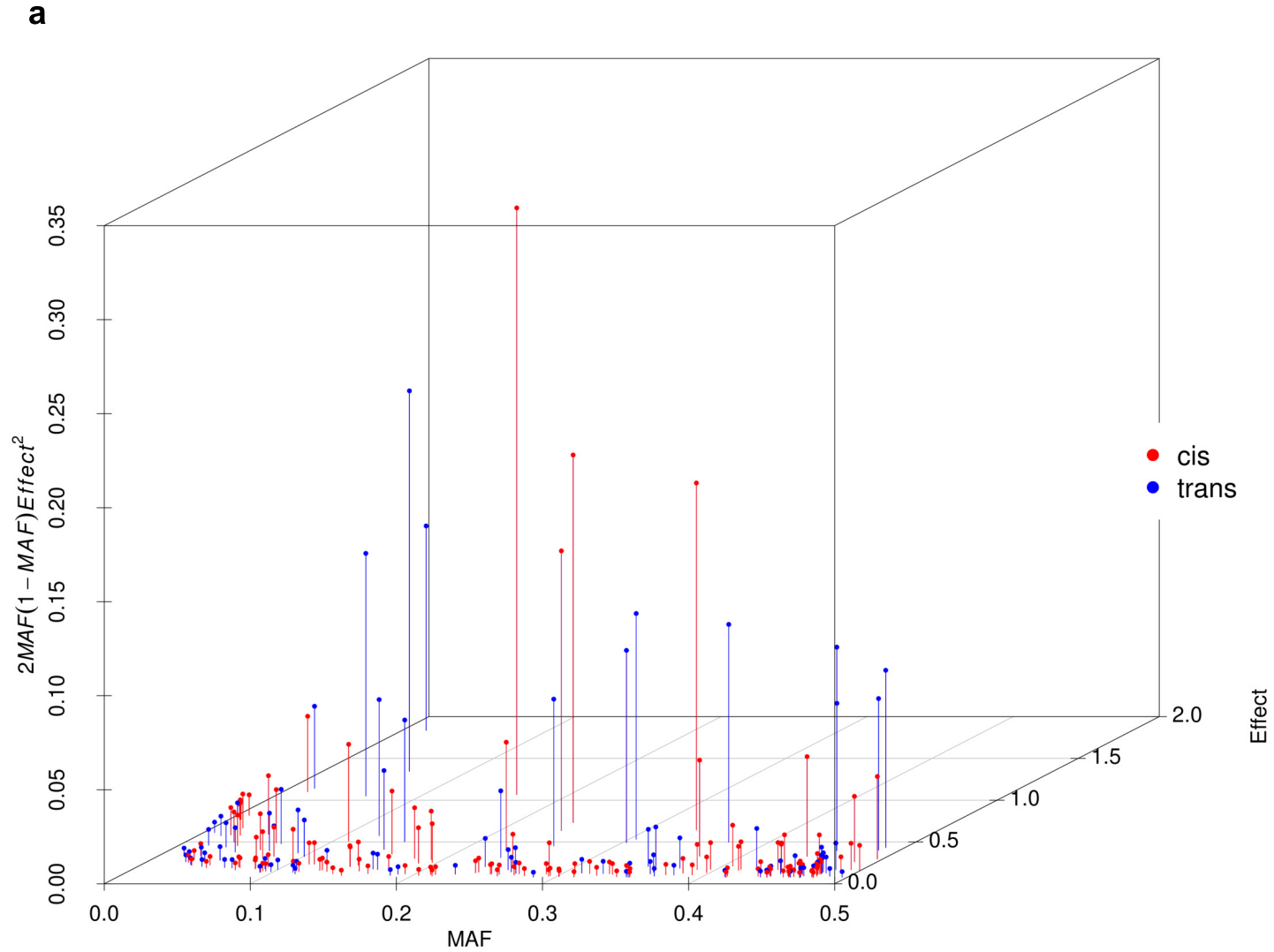


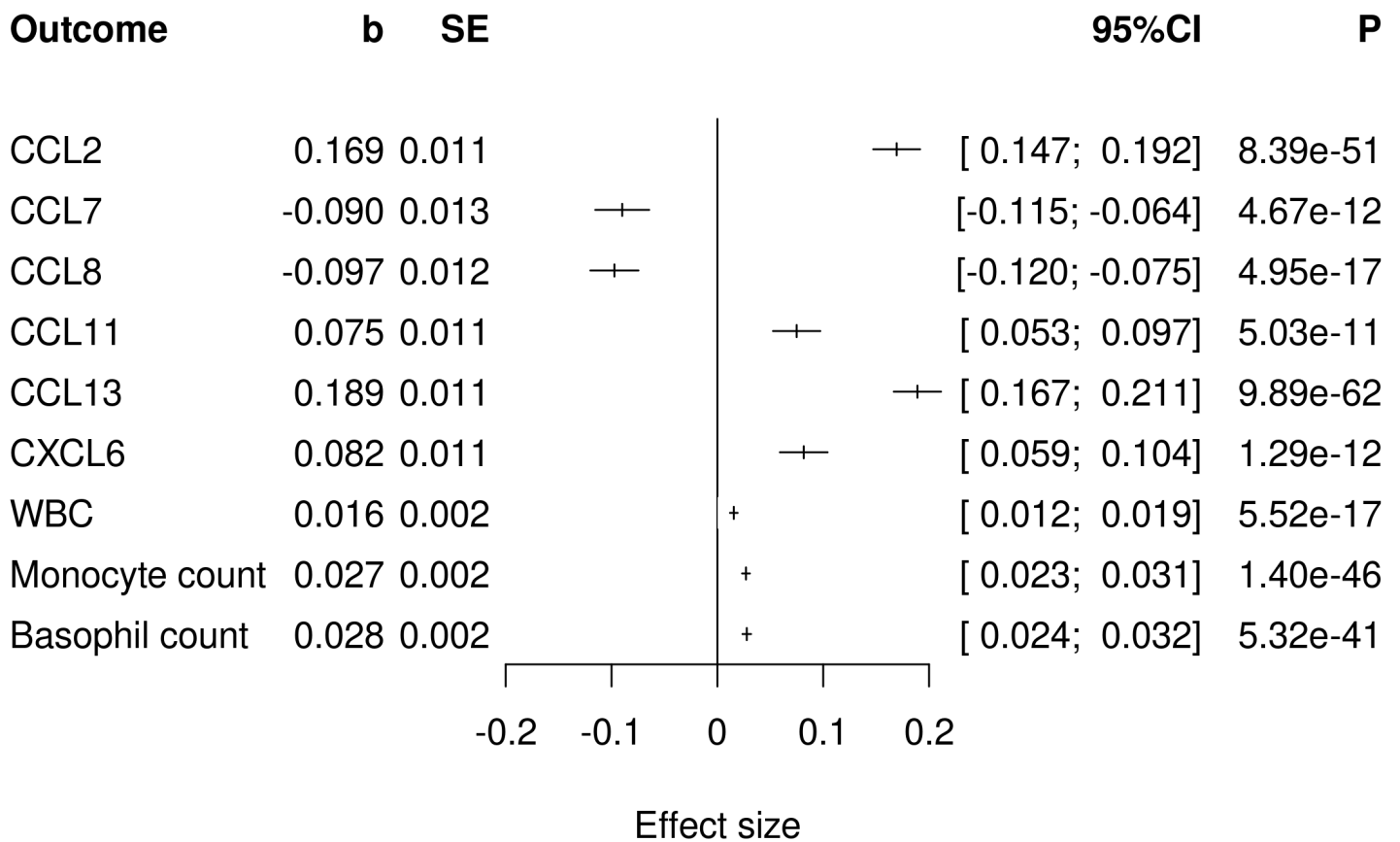
Analysis outline

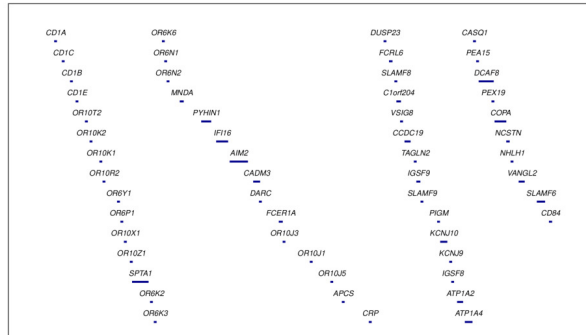
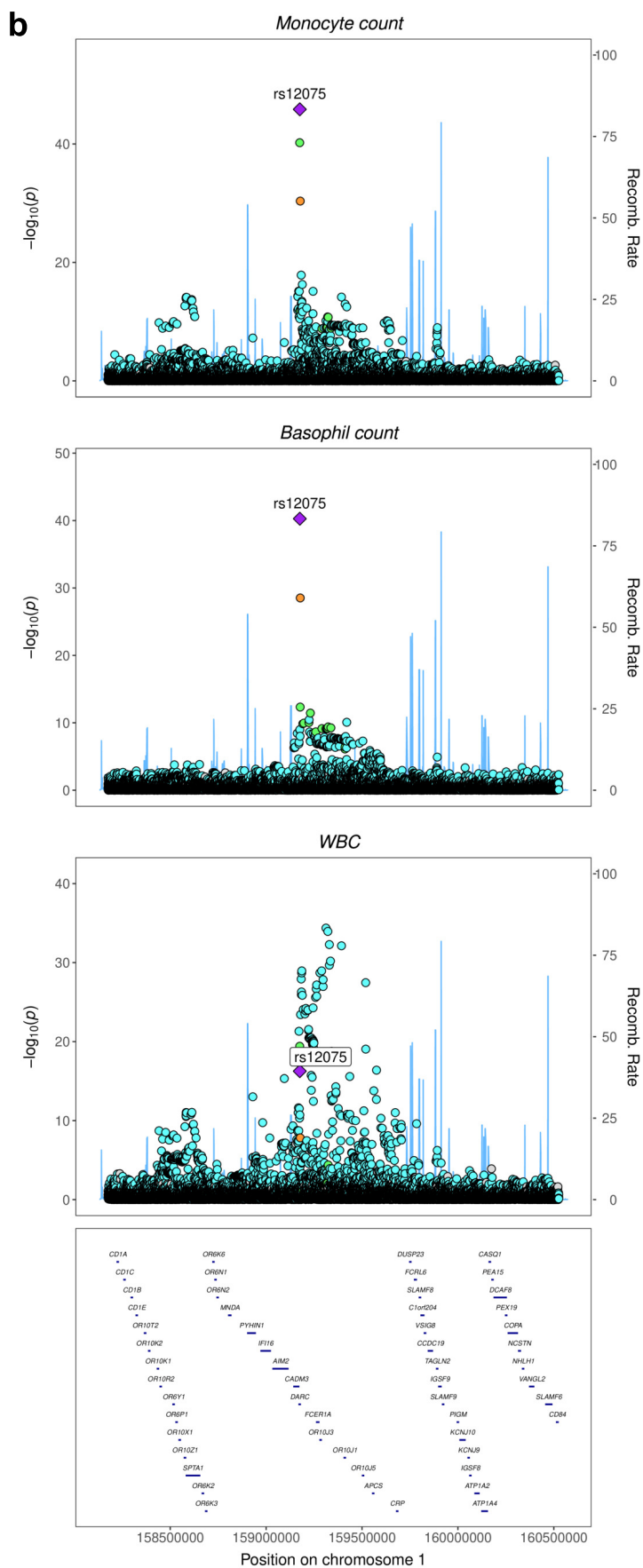
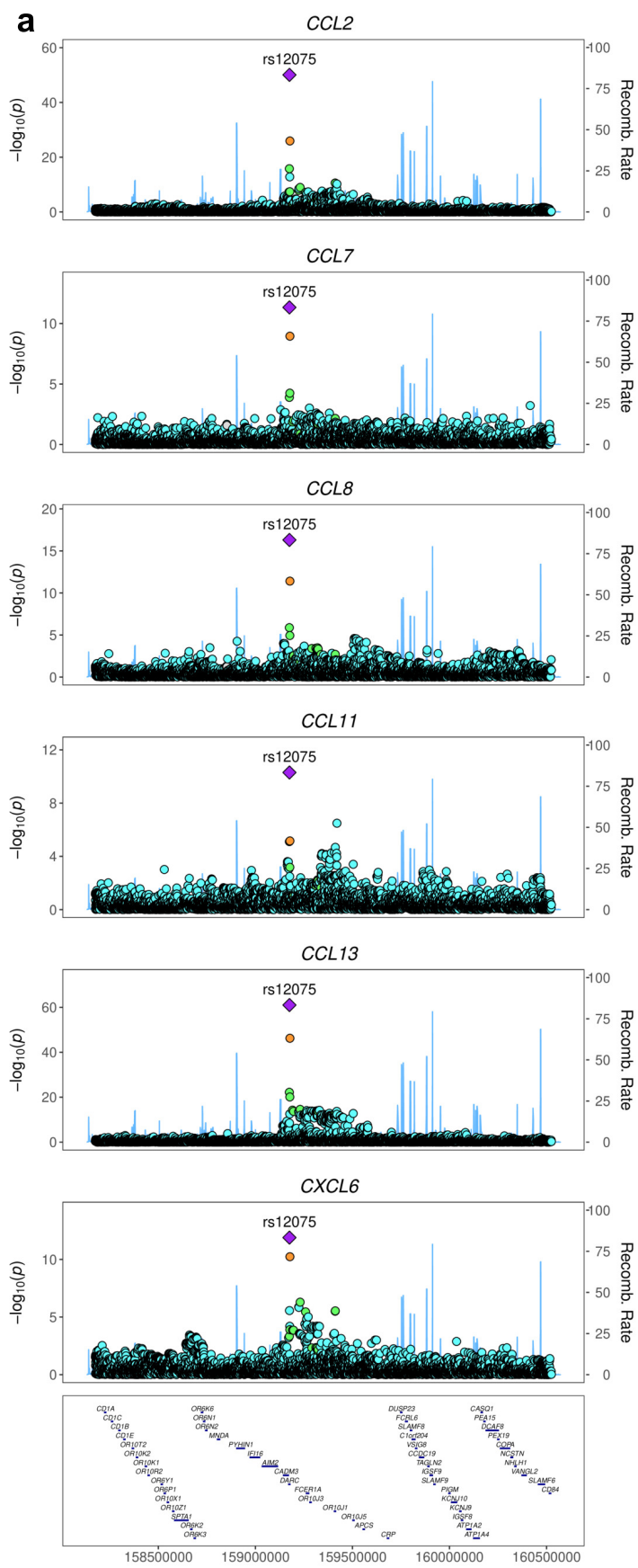


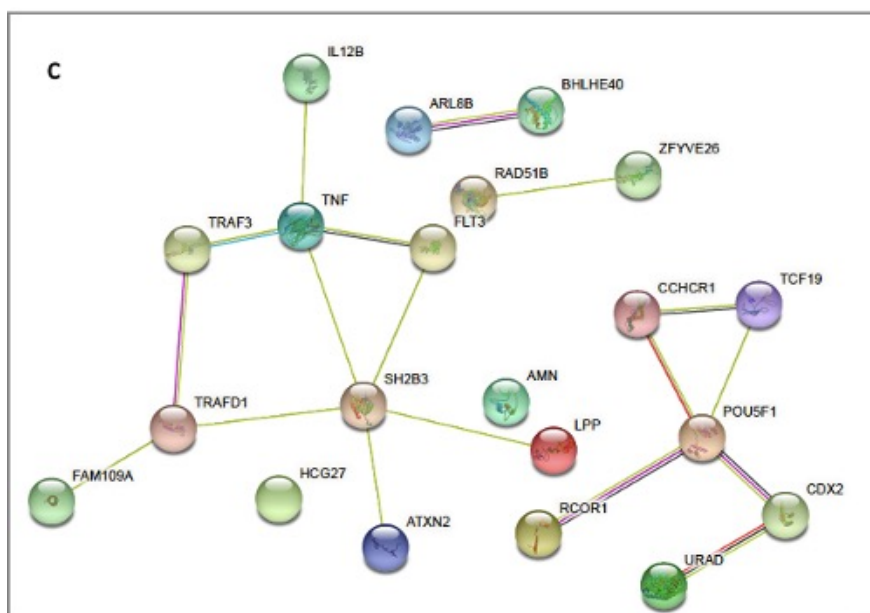
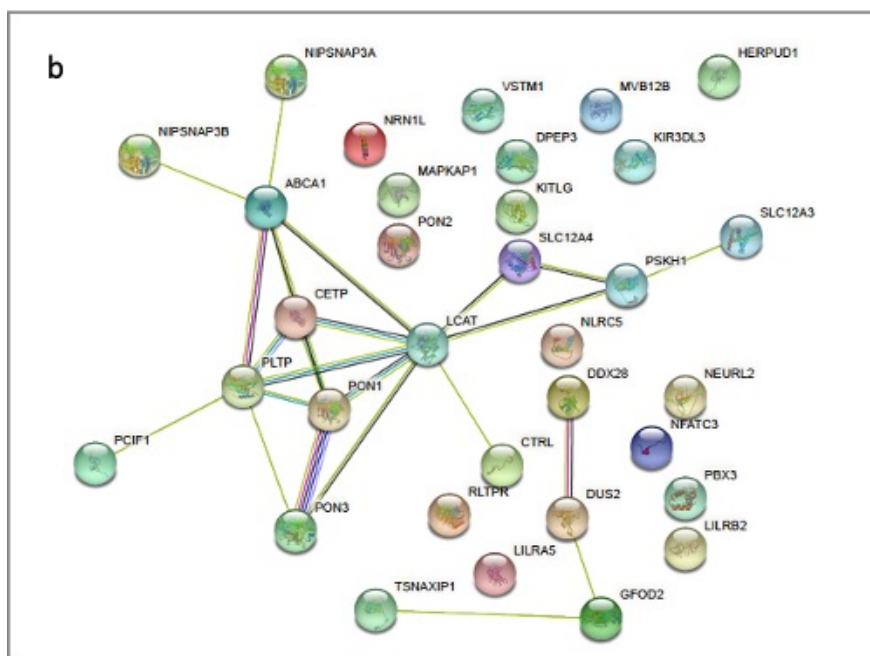
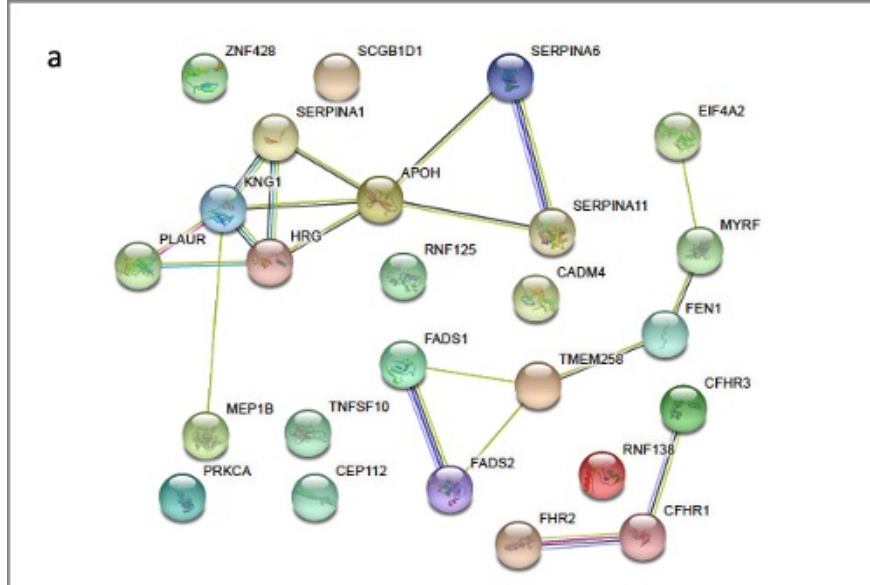












Known Interactions

- from curated databases
- experimentally determined

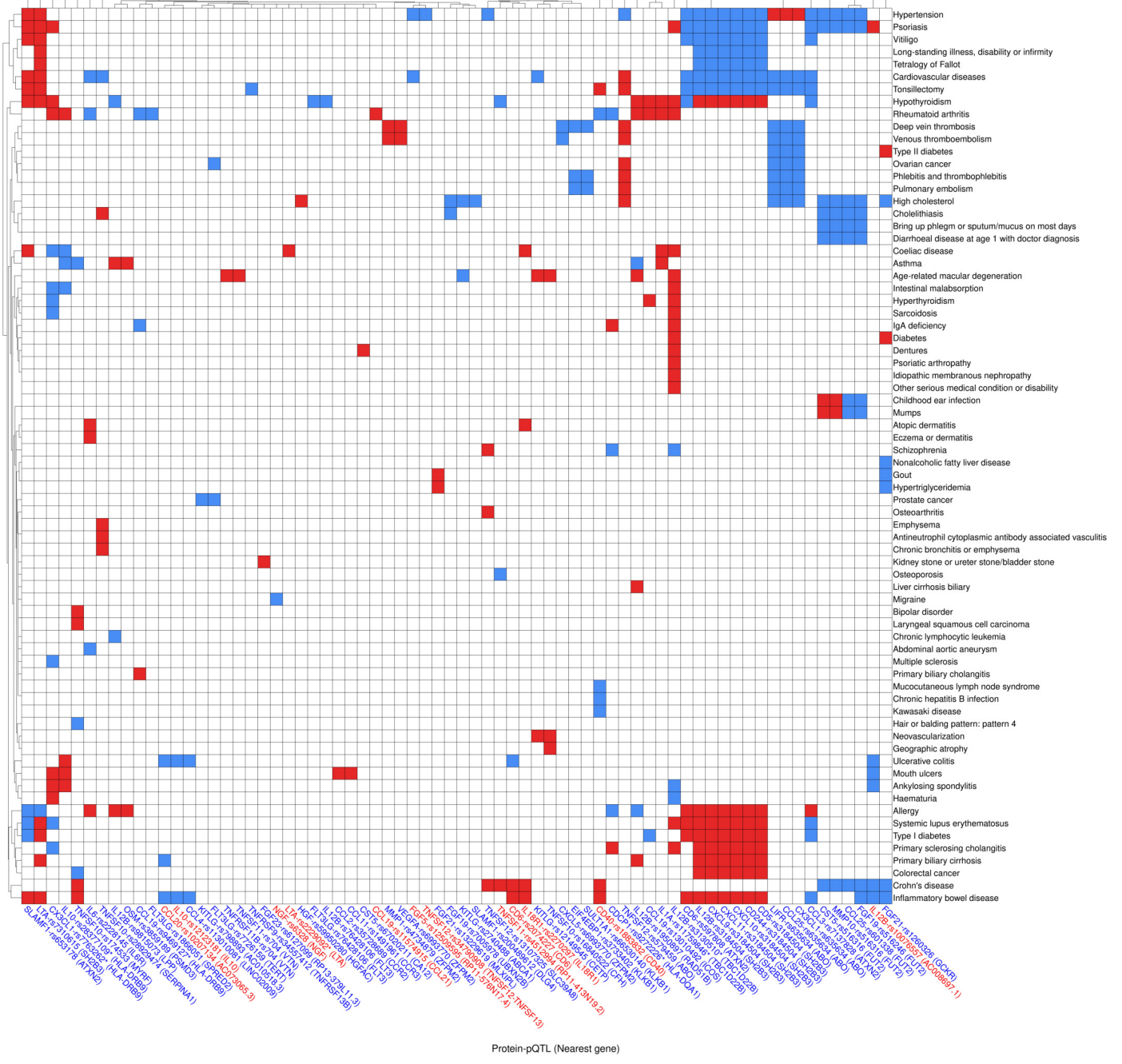
Predicted Interactions

- gene neighborhood
- gene fusions
- gene co-occurrence

Others

- textmining
- co-expression
- protein homology

GWAS diseases



Protein-pQTL (Nearest gene)

