1 Multi-ancestry proteome-phenome-wide Mendelian randomization offers a comprehensive 2 protein-disease atlas and potential therapeutic targets

- 3
- 4 Chen-Yang Su^{1–3}, Adriaan van der Graaf⁴, Wenmin Zhang⁵, Dong-Keun Jang⁶, Susannah Selber-Hnatiw^{1,2,7}, Ta-Yu Yang⁸, Guillaume Butler-Laporte^{3,9}, Kevin Y. H. Liang^{3,9}, Yiheng
- Chen^{9,10}, Fumihiko Matsuda⁸, Maria C. Costanzo⁶, J. Brent Richards^{3,7,9–12}, Noel P. Burtt⁶, Jason Flannick^{6,13}, Sirui Zhou^{1–3,7}, Vincent Mooser^{1,2,7}, Tianyuan Lu^{14–17*}, Satoshi Yoshiji^{1–3,6,7,9*}
- 5 6 7 8 9 10
- ¹Canada Excellence Research Chair in Genomic Medicine, McGill University, Montréal, Québec, Canada
- ²McGill Genome Centre, McGill University, Montréal, Québec, Canada
- 11 12 ³Quantitative Life Sciences, McGill University, Montreal, Québec, Canada
- ⁴Department of Computational Biology, University of Lausanne, Lausanne, Switzerland
- 13 ⁵Montreal Heart Institute, Université de Montréal, Montreal, Québec, Canada
- 14 15 ⁶Programs in Metabolism and Medical & Population Genetics, The Broad Institute of MIT and Harvard,
- Cambridge, MA, USA
- 16 ⁷Department of Human Genetics, McGill University, Montreal, Québec, Canada
- 17 18 ⁸Center for Genomic Medicine, Graduate School of Medicine, Kyoto University, Kyoto, Japan
- ⁹Lady Davis Institute, Jewish General Hospital, McGill University, Montréal, Québec, Canada
- 19 ¹⁰5 Prime Sciences, Montréal, Québec, Canada
- 20 ¹¹Department of Epidemiology, Biostatistics and Occupational Health, McGill University, Montréal, Québec, Canada
- 21 22 23 24 25 26 ¹²Department of Twin Research and Genetic Epidemiology, King's College London, London, United Kingdom
- ¹³Harvard Medical School, Boston, MA, USA
- ¹⁴Department of Population Health Sciences, University of Wisconsin-Madison, Madison, WI, USA
- ¹⁵Department of Biostatistics and Medical Informatics, University of Wisconsin-Madison, Madison, WI, 27 USA
- 28 ¹⁶Center for Genomic Science Innovation, University of Wisconsin-Madison, Madison, WI, USA
- 29 ¹⁷Center for Demography of Health and Aging, University of Wisconsin-Madison, Madison, WI, USA 30

31 *Correspondence and equal contribution:

32 Tianyuan Lu (tianyuan.lu@wisc.edu) and Satoshi Yoshiji (satoshi.yoshiji@mcgill.ca)

33 Abstract

34 Circulating proteins influence disease risk and are valuable drug targets. To enhance the 35 discovery of protein-phenotype associations and identify potential therapeutic targets across 36 diverse populations, we conducted proteome-phenome-wide Mendelian randomization in three 37 ancestries, followed by comprehensive sensitivity analyses. We tested the potential causal effects 38 of up to 2,265 unique proteins on a curated list of 355 distinct phenotypes, identifying 726,035 39 protein-phenotype pairs in European, 33,078 in African, and 115,352 in East Asian ancestries. 40 Notably, 119 proteins were instrumentable only in African ancestry and 17 proteins only in East 41 Asian ancestry due to allele frequency differences that are common in these ancestries but rare 42 in European ancestry. We identified 3,949, 56, and 325 unique protein-phenotype pairs in 43 European, African, and East Asian ancestries, respectively, and assessed their druggability using 44 multiple databases. We highlighted the causal role of IL1RL1 in inflammatory bowel diseases, 45 supported by multiple orthogonal lines of evidence. Taken together, this study underscores the 46 importance of multi-ancestry inclusion and offers a comprehensive atlas of protein-phenotype 47 associations across three ancestries, enhancing our understanding of proteins involved in disease 48 etiology and potential therapeutic targets. Results are available at the Common Metabolic 49 Diseases Knowledge Portal (https://broad.io/protein mr atlas).

51 Introduction

Circulating proteins play a major role in a multitude of biological pathways¹⁻³, are important 52 biomarkers for disease diagnosis, prognosis, and prevention⁴⁻⁶, and serve as valuable drug 53 targets⁷⁻¹⁰. Current high-throughput proteomics platforms measure thousands of circulating 54 55 plasma proteins. With measurements in large cohorts, recent studies have conducted genome-56 wide association studies (GWAS) to evaluate genetic variants associated with the abundances of 57 thousands of proteins^{11–14}. We can leverage these proteomic GWAS to find causal proteins for 58 diseases and prioritize drug targets through Mendelian randomization (MR)^{15,16}. In this case, MR 59 uses genetic variants associated with variation in plasma protein levels (known as protein 60 quantitative trait loci, pQTLs) as instruments to measure the causal effect of an exposure (protein 61 abundance) on an outcome (a complex trait). This is shown to reduce confounding and reverse 62 causation biases affecting many epidemiological studies, provided that three key assumptions 63 are met: (1) the genetic variant is associated with the exposure, (2) there is no confounding of the 64 instrument-outcome association, and (3) the genetic variant influences the outcome solely 65 through the exposure.

66

Despite applications of proteogenomics in elucidating disease mechanisms and identifying 67 potential therapeutic targets^{17–23}, previous pQTL studies^{18,19} have been based on smaller sample 68 sizes and measured fewer proteins, assessed limited outcomes, and most importantly, have 69 70 predominantly focused on individuals of European ancestry. African ancestry proteomics cohorts 71 have emerged in recent years^{13,14,24}, yet existing studies have assessed a limited number of 72 outcomes and comparisons with other ancestries have also been limited¹⁹. Similarly, in East Asian ancestries, few large-scale proteome-phenome wide MR studies exist^{25,26}, largely due to the lack 73 74 of publicly available pQTLs, despite the presence of existing East Asian ancestry biobanks providing hundreds of publicly available GWAS outcomes^{27,28}. 75 76

77 The inclusion of multiple ancestries in a proteome-phenome wide atlas of associations can 78 potentially offer significant benefits. Diverse ancestry inclusion in MR leverages the natural 79 variations in genetic architecture present across populations, which allows analyses on otherwise 80 unseen genetic variation, increases in statistical power due to allele frequency increases, and differentiation of causal effect magnitude across populations^{29,30}. Combined, this approach may 81 82 be able to identify a greater number of instrumentable proteins (i.e., proteins that can be tested 83 in MR), which may lead to an increase in discoveries. By leveraging insights obtained through 84 instrumentable proteins in one ancestry, the identified protein-phenotype associations could 85 contribute to our understanding of disease mechanisms, which may be generalizable across 86 different ancestries and potentially benefit all populations. For instance, PCSK9 loss-of-function 87 variants Y142X and C679X predispose individuals to naturally lower LDL cholesterol levels. 88 These mutations were found to be common in African Americans but rare in European 89 Americans³¹ and inspired the development of PCSK9 inhibitors mimicking these variants to 90 effectively reduce LDL cholesterol levels and risk of cardiovascular events^{32,33}. In addition to 91 providing insights into novel therapeutic targets, enhancing diversity in genomic studies is crucial 92 to ensure equitable health outcomes and address imbalances in health disparities across 93 populations³⁴.

94

Here, we combined four of the largest European ancestry proteomics cohorts (n = up to 35,559), two of the largest African proteomics cohorts (n = up to 1,871), and a new East Asian proteomics cohort (n = 1,823). We then performed MR and colocalization analyses on a curated list of the most recent and largest ancestry-specific outcome GWAS to date for 179 European and 26 African ancestry outcomes, as well as 206 East Asian ancestry outcomes from Biobank Japan to construct an atlas of protein-phenotype associations. We integrated our findings with multiple drug databases to assess the druggability of the associations and highlighted novel targets.

Overall, our study supports the prioritization of thousands of protein-phenotype associations and
 provides a comprehensive, updated resource for the community, significantly expanding our
 understanding of these associations. Results are publicly available at the Common Metabolic
 Diseases Knowledge Portal (https://broad.io/protein mr atlas).

107 Results

108 The overall study design is shown in **Fig. 1**.

109



111 Figure 1. Study design.

112 We assessed the causal role of 2,265 circulating plasma proteins across three ancestries (four 113 European, two African, and a new East Asian ancestry cohort) on up to 355 phenotypes/outcomes 114 with extensively curated GWAS in each ancestry. For European and African ancestry outcomes, 115 we collected the largest and most recent GWAS available as of February 2024 for 179 and 26 116 outcomes, respectively, while 206 GWAS from BioBank Japan were used for East Asian ancestry 117 outcomes. We implemented a unique approach for defining *cis*-pQTLs using a multi-step process 118 combining a strict *cis*-pQTL definition with Open Targets Genetics variant-to-gene score³⁵ filtering 119 to minimize risk of horizontal pleiotropy, which we term "strict variant-to-gene (V2G) cis-pQTLs". 120 Then, we performed multi-ancestry proteome-wide MR and colocalization analyses using 121 PWCoCo^{18,36} and SharePro³⁷, a new colocalization method we developed, on the curated 122 phenotypes to identify causal protein-phenotype associations. Next, we overlapped prioritized 123 protein-phenotype pairs with drug databases such as the druggable genome³⁸, DrugBank³⁹, and the Open Targets platform⁴⁰. Finally, as an illustrative example of the value of our atlas, we 124 125 pinpoint IL1RL1 and explore its role in inflammatory bowel disease using multiple lines of 126 evidence. EUR: European, AFR: African, EAS: East Asian; MR: Mendelian randomization; IBD: 127 Inflammatory bowel disease; CD: Crohn's disease; UC: Ulcerative colitis.

- 128
- 129

130

131 1. Genetic instrument selection for proteins

132 We summarize the selected genetic instruments in **Supplementary Table 1**. Briefly, we used 133 proteomic GWAS from four European ancestry cohorts: ARIC (4,657 proteins measured in up to 134 7,213 individuals)¹³, deCODE (4,719 proteins measured in up to 35,559 individuals)¹², Fenland 135 (4,775 proteins measured in up to 10,708 individuals)¹¹, and UKB-PPP (2,923 proteins measured 136 in up to 34,557 individuals)¹⁴. For African ancestry, we analyzed proteomic GWAS from two 137 cohorts: ARIC (4657 proteins measured in up to 1,871 individuals)¹³ and UKB-PPP (2923 proteins 138 measured in up to 931 individuals)¹⁴. Additionally, we included a new East Asian ancestry cohort, 139 the Kyoto University Nagahama cohort which measured 4,196 proteins in up to 1,823 individuals.

140

141 1.1. Unifying cis-pQTL across cohorts

We re-defined *cis*-pQTLs across proteomics cohorts as pQTLs within 500 kb of the transcription start site (TSS) of the protein-coding gene, with independence and significance defined with linkage disequilibrium (LD) $r^2 < 0.001$ and $P < 5 \times 10^{-8}$, respectively. Independent pQTLs outside the *cis*-region were labeled as *trans*-pQTLs. We analyzed *cis*-pQTLs because they are more likely to have direct biological effects on the proteins of interest^{41,42}. We verified that newly defined *cis*pQTLs had high concordance within ancestries (**Supplementary Note 1A**).

148

149 1.2. Identifying strict V2G cis-pQTLs

150 We mitigated the risk of horizontal pleiotropy by using a unique approach to select genetic 151 instruments which we term strict V2G cis-pQTLs (Supplementary Note 1B). Strict V2G cis-152 pQTLs are associated with a single protein-coding gene ("strict") and have the strongest link to 153 the corresponding protein-coding gene based on the Open Targets Genetics V2G score, which 154 uses multiple sources of evidence to map variants to genes ("V2G")³⁵ (Extended Data Fig. 1 and 155 Supplementary Tables 2–8). We also assessed if these strict V2G cis-pQTLs were protein altering variants (PAVs), as PAVs may affect protein structure and lead to bias in effect size 156 157 estimation¹. We found that only 70 of 7,399 (0.95%), 30 of 1,684 (1.8%), and 9 of 663 (1.4%) strict V2G cis-pQTLs in European, African, and East Asian ancestries, respectively, were PAVs 158 159 or in high LD with PAVs of high impact. Notably, strict V2G cis-pQTLs were more enriched for cis-160 eQTLs, indicating their role in local gene regulation impacting protein levels (Extended Data Fig. 161 2a). They had significantly larger effects on protein levels (Extended Data Fig. 2b) and were

significantly closer to the TSS of the corresponding protein-coding gene (Extended Data Fig. 3).
 Hence, strict V2G selection of pQTLs limits the risks of violating MR assumptions (specifically no
 horizontal pleiotropy) while still retaining adequate statistical power. Specifically, all proteins in
 each cohort had F-statistics above 10, suggesting that the risk of weak instrument bias is limited⁴³
 (Supplementary Table 9).

168 2. Multi-ancestry inclusion is important in instrumenting proteins and reveals population-specific169 variants

170 Across the three ancestries, we were able to instrument 2,265 unique proteins with 449 of these 171 being shared across all three ancestries (Fig. 2a). Specifically, we instrumented 2,110 proteins in 172 European, 1,144 in African, and 581 in East Asian ancestries (see Supplementary Note 1C for 173 cohort level description and Extended Data Fig. 4). We identified 1,018 proteins that were unique 174 to individuals of European ancestry. Including African ancestries allowed for an additional 138 175 proteins, with 119 unique to African ancestry and 19 shared with East Asian ancestries. Including 176 an East Asian ancestry cohort allowed an additional 17 unique proteins (Fig. 2a). Altogether. 177 these findings underscore the value of including multiple ancestries in proteomic analyses.

178

167

179 Next, to better understand the unique proteins in non-European ancestries, we compared the 180 allele frequencies of their genetic instruments using gnomAD⁴⁴. Of the 130 genetic instruments 181 for the 119 proteins unique to African ancestries, 89 (68.5%) had a minor allele frequency (MAF) 182 below 0.01 in European ancestries (Fig. 2b). Similarly, for the 18 genetic instruments for the 17 183 proteins unique to East Asian ancestry, 13 (72.2%) had a MAF below 0.01 in European ancestry 184 (Fig. 2c). The majority (29, 63.0%) of the 46 unique genetic instruments for the 19 proteins 185 instrumentable by both African and East Asian but not by European ancestries, were rare (MAF 186 < 0.01) in European ancestry (Fig. 2d). These results suggest that populational allele frequency 187 differences allow for more proteins to be included in MR analyses. We refer to them as uniquely 188 instrumentable proteins.



190
 191
 192
 193
 193
 194
 195
 195
 196
 190
 190
 190
 191
 192
 193
 194
 195
 195
 196
 190
 190
 190
 191
 191
 192
 193
 194
 195
 195
 196
 196
 190
 190
 191
 192
 193
 194
 195
 195
 196
 196
 196
 196
 197
 198
 198
 199
 199
 190
 190
 191
 192
 193
 194
 195
 195
 195
 196
 196
 197
 198
 198
 199
 199
 199
 190
 190
 191
 192
 193
 194
 195
 195
 195
 196
 196
 197
 198
 198
 199
 199
 190
 190
 190
 191
 192
 193
 194
 195
 195
 195
 196
 196
 197
 198
 198
 199
 199
 199
 190
 190
 191
 191
 192
 193
 194
 194
 195
 195
 195
 196
 198
 199
 199
 190
 190
 190
 190
 191
 192

- (b) gnomAD minor allele frequencies (MAF) of *cis*-pQTLs instrumenting the 119 proteins uniquely instrumentable in African ancestry when plotted against corresponding gnomAD non-Finnish European allele frequency (AF). The black box in the top plot is zoomed in and shown in the bottom plot. Blue: ARIC African *cis*-pQTLs. Red: UKB-PPP African ancestry cohort *cis*-pQTLs.
 (c) gnomAD minor allele frequencies (MAF) of *cis*-pQTLs instrumenting the 17 proteins
- (c) gnomAD minor allele frequencies (MAF) of *cis*-pQTLs instrumenting the 17 proteins uniquely instrumentable in East Asian ancestry when plotted against corresponding gnomAD non-Finnish European AF. The black box in the top plot is zoomed in and shown in the bottom plot. Purple: Kyoto University Nagahama East Asian ancestry cohort *cis*-pQTLs.
 (d) gnomAD MAF minor allele frequencies (MAF) of *cis*-pQTLs instrumenting the 19 proteins
 - (d) gnomAD MAF minor allele frequencies (MAF) of *cis*-pQTLs instrumenting the 19 proteins common to African and East Asian but not European ancestries when plotted against corresponding gnomAD non-Finnish AF. The black box in the top plot is zoomed in and shown in the bottom plot. Blue: ARIC African *cis*-pQTLs. Red: UKB-PPP African *cis*pQTLs. Purple: Kyoto University Nagahama East Asian ancestry cohort *cis*-pQTLs.
- 212 213 214

208

209

210

211

215 3. Two-sample Mendelian randomization

We performed two-sample MR using strict V2G *cis*-pQTLs as instrumental variables to determine causal proteins implicated in human complex traits and diseases. Across three ancestries, we considered 355 unique outcomes (**Fig. 3a**) pertaining to 24 phenotype categories (**Fig. 3b**). These phenotypes/outcomes (as of February 2024) included 179 of the most up to date and largest phenotypic GWAS available for European ancestries (**Supplementary Table 10**), 26 of the largest available GWAS for African ancestries (**Supplementary Table 11**), and 206 from BioBank Japan for East Asian ancestries (**Supplementary Table 12**).



224 225 226

227

onique protein-prienotype pairs

5 Figure 3. MR analyses to determine the effects of proteins on phenotypes.

(a) Curated GWAS phenotypes common between different ancestries. European (red), African (green), and East Asian (blue) ancestries.

- (b) Phenotype categories for outcomes used in all three ancestries. 229
 - (c) Flowchart summary of MR and colocalization analyses to identify protein-phenotype associations for each cohort.

230 231 232

228

233 234 We tested a total of 726,035 protein-phenotype pairs in European ancestries (173,748 for ARIC, 235 182,433 for deCODE, 180,018 for Fenland, and 189,836 for UKB-PPP), 33,078 in African ancestries (20,442 for ARIC and 12,636 for UKB-PPP), and 115,352 in East Asian ancestries 236 237 (Kyoto University Nagahama cohort) (see Data availability). To control for false positives, we 238 applied a Benjamini-Hochberg-corrected P value (false discovery rate, FDR)⁴⁵ threshold of 0.05 239 (5%) per cohort in each ancestry¹⁹. We note that Bonferroni correction is overly stringent given 240 that (i) proteins are correlated with one another, and (ii) we tested the same protein-phenotype 241 associations across cohorts, making these tests not independent. Nevertheless, we also provide 242 the most stringent Bonferroni corrected associations in section 5. Results were further filtered 243 based on multiple sensitivity analyses robust to MR assumption violations and retained 244 associations are now termed "MR-passing". Assessment of sample overlap between European 245 ancestry GWAS outcomes that were based on the UK Biobank and proteomics from the European 246 ancestry UKB-PPP cohort was performed (Supplementary Note 2). Of the tested associations, 247 a total of 12,247 associations were considered MR-passing in European, 83 in African, and 387 248 in East Asian ancestries (Supplementary Table 13).

250 4. Colocalization analyses

251 MR results may be confounded by independent causal variants in LD^{46,47}. To guard against such bias, for each MR-passing protein-phenotype pair, we performed colocalization analysis using 252 two methods, PWCoCo^{18,36} and SharePro³⁷, to verify that the protein abundance and the tested 253 254 outcome share the same genetic signals (Fig. 3c). An MR-passing protein-phenotype pair was 255 considered putatively causal if it was supported by at least one method with a posterior 256 colocalization probability (PP_{max}) \geq 0.8. Across all cohorts and all three ancestries, 56.5% (7,182) 257 out of 12,717) of MR-passing associations were supported by colocalization evidence 258 (Supplementary Table 13).

259

249

260 5. Putatively causal protein-phenotype associations

261 Upon MR, sensitivity analyses, and colocalization, we identified 3,949 unique putatively causal 262 protein-phenotype pairs in European (Extended Data Fig. 5a and Supplementary Table 14), 56 263 in African (Extended Data Fig. 5b and Supplementary Table 15), and 325 in East Asian 264 ancestries (Extended Data Fig. 5c and Supplementary Table 16). Here, we use protein to refer 265 to protein-coding genes to harmonize across SomaScan and Olink platforms. Results are also 266 hosted at https://broad.io/protein mr atlas. We described cohort level associations and proteins 267 implicated in a multitude of phenotypes in Supplementary Notes 3 and 4. Particularly, 1,617, 30, 268 and 135 unique associations in European, African, and East Asian ancestries further withstood 269 the Bonferroni correction accounting for the total number of tests across all cohorts and all ancestries ($P < 0.05 / 874,465 = 5.7 \times 10^{-8}$), however, this threshold is likely overly conservative 270 271 due to many proteins being correlated with one another as well as non-independent tests from 272 the same protein-phenotypes associations being tested across cohorts (Supplementary Tables 273 **14 – 16**). 274

275 In European ancestries, the 3.949 significant unique protein-phenotype pairs identified in 276 European ancestries involved putatively causal effects between 995 proteins and 146 phenotypes, 277 of which only 56 (1.4%) showed discordant MR estimates in one of the tested cohorts 278 (Supplementary Table 14). These discrepancies could be attributed to population differences or

variations in proteomic assays, such as differential effects from SomaScan aptamers targeting
different domains compared to Olink assays. Of the 3,893 remaining putatively causal European
ancestry pairs between 991 proteins and 146 outcomes, 1,692 (43.5%) of the identified putatively
causal associations were from 452 proteins uniquely instrumentable by European ancestries.
Further, 3,853 (99.0%) associations have not been previously reported by earlier proteomephenome wide MR studies from Zheng et al.¹⁸ and Zhao et al.¹⁹.

285

286 5.1. Cardiovascular and autoimmune diseases

287 We demonstrate the highly interconnected nature between proteins and outcomes by highlighting 288 cardiovascular (Fig. 4a and Supplementary Note 5A) and autoimmune phenotypes (Fig. 4b) in 289 European ancestries given their significant impact on health. Cardiovascular outcomes were 290 influenced by up to 103 proteins (median: 7) while some proteins influenced up to 8 cardiovascular 291 outcomes (median: 1) (Fig. 4a). For instance, we found that an s.d. increase in genetically 292 predicted ULK3 levels increases systolic and diastolic blood pressure, pulse pressure, and 293 hypertension. ULK3 is a nuclear kinase which may contribute to vascular disease by mediating 294 autophagy dysregulation⁴⁸. In concordance, a recent study showed that functional splicing effects of ULK3 can contribute to coronary artery disease (CAD)⁴⁹ suggesting effective modulation of 295 296 ULK3 may be beneficial for reducing cardiovascular risk.

297

298 Similarly, we found high interconnectedness between autoimmune phenotypes and proteins (Fig. 299 **4b**). Autoimmune phenotypes were influenced by a median of 7 proteins: the highest number of 300 associations were for Crohn's disease (CD) (39 proteins), inflammatory bowel disease (IBD) (35 301 proteins), and ulcerative colitis (UC) (28 proteins). Proteins, on the other hand, influenced a 302 median of 2 autoimmune phenotypes. IL1RL1 and IL12B influenced the largest number of 303 outcomes (6 and 4 outcomes, respectively) (Fig. 4b) with both increasing risk of IBD, CD, and 304 UC. IL12B is targeted by a commercially available drug ustekinumab for CD⁵⁰; however, there are 305 currently no approved drugs targeting IL1RL1, although there has been increasing interest in modulating IL1RL1 for treating various conditions^{51,52}. For example, tozorakimab targets IL-33, 306 307 the interleukin that binds IL1RL1 (ST2)⁵³ and anakinra targets IL-1⁵⁴, a closely related interleukin. 308

309 Other potentially novel findings involving other disease categories in European ancestries are 310 highlighted in **Supplementary Note 5B**.



312 313

325

326

Figure 4. Protein-phenotype network plots.

Red arrows indicate a positive causal estimate of the protein on the phenotype while blue arrowsindicate a negative causal estimate of the protein on the phenotype.

- 316 (a) Significant estimates between proteins (orange circles) and cardiovascular traits (green 317 rectangles) both derived from European ancestry individuals. All cardiovascular traits are 318 named, as well as the protein traits that have a significant effect on 4 or more 319 cardiovascular traits. Arrow thickness representing number of significant estimates 320 indicates how often a protein measurement has a causal effect on the outcome trait. The 321 maximum number of significant estimates is 7 which is larger than the total number of 322 European ancestry proteomics studies, 4, due to the presence of more than one 323 SomaScan aptamer for that protein. For simplicity, we only depict protein-phenotype pairs 324 in which all European ancestry cohorts showed concordant direction of effect estimates.
 - (b) Significant estimates between proteins (orange circles) and autoimmune traits (green rectangles) both derived from European ancestry individuals. All autoimmune traits are

named, as well as the protein traits that have a significant effect on 4 or more autoimmune traits. Arrow thickness representing number of significant estimates indicates how often a protein measurement has a causal effect on the outcome trait. For simplicity, we only depict protein-phenotype pairs in which all European ancestry cohorts showed concordant direction of effect estimates.

- (c) Significant estimates between proteins (orange circles) and all traits (green rectangles) both derived from African ancestry individuals. All arrows are the same thickness to indicate support of the association in a single study.
- (d) Significant estimates between proteins (orange circles) and binary traits (green rectangles) both derived from East Asian ancestry individuals. All arrows are the same thickness and indicate support of the association in a single study.
- 337 338

327

328

329

330

331

332

333

334

335

336

339 340

341 5.2. Protein-phenotype associations in African and East Asian ancestries

342 In African ancestries, we identified 56 unique protein-phenotype associations involving 28 343 proteins and 11 phenotypes (Fig. 4c and Supplementary Table 15). Of these 56 pairs, 55 344 (98.2%) have not been previously reported by earlier proteome-phenome wide MR studies in 345 African ancestries¹⁹. Notably, 11 (19.6%) protein-phenotype pairs involving four proteins, APOE, 346 C7orf50, CD300LG, and PON3, were uniquely instrumentable in African ancestries (Extended 347 Data Fig. 6a). For instance, increased PON3 levels was associated with increased total cholesterol ($\beta_{UKB-PPP} = 0.03$, 95% CI = 0.02–0.05, $P = 1.1 \times 10^{-5}$, $PP_{max} = 0.95$). PON3 is highly 348 expressed in the liver and previously implicated in cholesterol metabolism^{55,56} and atherosclerosis 349 350 progression. This is notable as PON3 was not instrumentable in European ancestries, thus, 351 African ancestries are uniquely suitable to identify biologically plausible protein-phenotype 352 associations.

353

Furthermore, the Million Veteran Program⁵⁷ recently released summary statistics for additional 354 355 traits in up to 635,969 individuals. Using this resource, we tested the causal effect of the 119 356 proteins uniquely instrumentable in African ancestries on cardiovascular and autoimmune-related 357 binary outcomes. Using the African ARIC proteomics cohort, we identified 7 associations with 358 outcomes and no associations with autoimmune-related cardiovascular outcomes 359 (Supplementary Note 5C). Notably, increased PCYOX1 levels were associated with a reduced 360 risk of coronary atherosclerosis, atrial fibrillation, and flutter, indicating its protective effect. 361 PCYOX1 plays a role in oxidative stress and lipid metabolism and has been implicated in atherosclerosis in rodent studies⁵⁸, supporting our findings. This suggests that as more outcomes 362 363 with larger sample size, particularly binary disease outcomes, become available in African 364 ancestries, we can better leverage uniquely instrumentable proteins to uncover additional protein-365 disease associations.

366

In East Asian ancestries, 339 putatively causal associations were identified with 325 unique protein-phenotype pairs involving 110 proteins and 86 phenotypes (**Fig. 4d** and **Supplementary Table 16**). Among them, increased SMOC2 level was associated with decreased risk of peripheral artery disease (PAD) (OR = 0.82, 95% CI = 0.74-0.90, $P = 5.5 \times 10^{-5}$, PP_{max} = 0.97). *SMOC2* is highly expressed in arteries and plays roles in endothelial cell proliferation, angiogenesis, and matrix assembly⁵⁹. Notably, in European ancestries, increased SMOC2 was associated with decreased pulse pressure, implicating favorable cardiovascular effects across ancestries.

We also found 8 proteins which were only instrumentable in East Asian ancestry due to lack of genome-wide significant *cis*-pQTLs in other ancestries or stringent strict V2G instrument selection (**Supplementary Note 5D**). These 8 proteins were ALDH2, ANXA7, APOA1, DDOST, GSS,

PLA2G7, PRSS2, and UGT1A1 which together accounted for 67 (20.6%) protein-phenotype associations (**Extended Data Fig. 6b** and **Supplementary Note 5D**). For instance, increased PRSS2 levels in East Asian ancestries was associated with an OR of 2.05 for acute pancreatitis (95% CI: 1.42–2.98, $P = 1.43 \times 10^{-4}$, $PP_{max} = 0.95$) which has been previously validated in European only proteome-wide analyses⁶⁰ suggesting concordant effects across ancestries.

383

405 406

384 6. Concordant effects across ancestries

385 We evaluated the direction of effect for protein-phenotype associations that passed MR and 386 colocalization analyses across ancestries since concordance across ancestries may strengthen 387 the evidence of broader applicability of therapeutic targets. We found 186 total protein-388 phenotype/outcome pairs with 63 unique pairs between 43 proteins and 13 outcomes 389 (Supplementary Table 17 and Fig. 5). Among these, 51 pairs (81.0%) had concordant effects 390 across ancestries, while 12 pairs (19.0%) had inconsistent effects (Supplementary Figure 3 and 391 Supplementary Note 6). Notable concordant associations included PCSK9 and LDL cholesterol 392 in European and East Asian ancestries, and haptoglobin (HP), which binds free hemoglobin to 393 prevent oxidative damage, with LDL cholesterol and total cholesterol across all three ancestries. 394 Angiopoietin-related protein 4 (ANGPTL4) was negatively associated with HDL cholesterol and 395 positively with triglycerides in European and African ancestries, while lipoprotein-lipase (LPL) 396 showed opposite associations in these ancestries. ANGPTL4 has been shown to act as a local 397 inhibitor of LPL⁶¹ which serves as the rate-limiting enzyme in the degradation of triglycerides⁶², 398 concordant with our findings. While no approved drugs exist for ANGPTL4, inhibition of a closely 399 related protein ANGPTL3 through an RNA interference therapy zodasiran is currently undergoing 400 clinical trials for cholesterol lowering⁶³. Thus, many of these biological effects were validated in 401 our multi-ancestry analyses. We note that while concordant effects across ancestries provide 402 strong evidence of support, discordant effects in MR do not automatically indicate biologically 403 discordant effects across ancestries, which could be due to epitope-binding effects or other technical variations⁶⁴. 404



407 Figure 5. Multi-ancestry network plot for protein-phenotype pairs present in two or more

408 ancestries.

All estimates shown had evidence in European ancestry and at least one other ancestry (African or East Asian ancestry). Significant estimates between proteins (orange circles) and traits (green rectangles). Arrow thickness indicates the number of ancestries in which a protein measurement has a causal effect on the phenotype. For simplicity, we only depict estimates when cohorts within the same ancestry and across different ancestries showed concordant direction of effect estimates. Red arrows indicate a positive causal estimate of the protein on the phenotype while blue arrows indicate a negative causal estimate of the protein on the phenotype.

- 416
- 417 418
- 419 7. Druggability
- 420 7.1. Druggability of the instrumentable protein-coding genes

Across all ancestries, between 60% to 67% of instrumentable protein-coding genes overlapped
with the druggable genome³⁸, which classifies genes into Tier 1, 2, or 3 according to druggability,
and between 20.5% and 23.6% overlapped with Tiers 1 and 2 (Supplementary Table 18 and
Supplementary Note 7A). Druggability tiers for instrumentable proteins are in Supplementary
Tables 19–25.

426

Next, we incorporated the druggable genome³⁸, DrugBank³⁹, and Open Targets Platform⁴⁰ to 427 428 determine which instrumentable proteins overlapped at least one database. Proteins overlapping 429 with DrugBank³⁹ had approved or investigational drugs available while those overlapping Open 430 Targets have information on their clinical development phase and status of protein-drug-disease 431 combinations. Cross-ancestry comparison stratified by proteomics platform for instrumentable 432 proteins overlapping at least one database shows that in SomaScan v4, African ancestry adds 68 433 additional targets beyond European ancestry cohorts (Extended Data Fig. 7a), while East Asian 434 ancestry contributes 34 more targets (**Extended Data Fig. 7b**). In Olink Explore 3072, including 435 African ancestry presents 62 additional targets (**Extended Data Fig. 7c**). These findings suggest 436 that data from African and East Asian ancestries could enhance drug development by offering 437 more potential therapeutic targets.

- 438
- 439 7.2. Druggability of protein-phenotype pairs by integrating the druggable genome, DrugBank,
 440 and Open Targets Platform
- 441 Across three ancestries, among the 1,037 number of protein-coding genes that have at least one 442 protein-phenotype association, 669 (64.5%) were present in at least one database. Specifically, 443 579 (55.8%) protein-coding genes overlapped with the druggable genome³⁸, 350 (33.8%) had 444 approved or investigational drugs in DrugBank³⁹, and 191 (18.4%) overlapped with the Open 445 Targets Platform⁴⁰ (**Supplementary Table 26**). Notably 32 (3.1%) were unique to non-Europeans. 446 This highlights that multi-ancestry inclusion can expand the list of actionable druggable 447 associations. Overlap with the druggable genome and DrugBank stratified by ancestry is 448 presented in Supplementary Note 7B.
- 449

We found that higher levels of MANBA, a Tier 2 target, increased risk of atrial fibrillation in European ancestry (**Fig. 6a**) and myocardial infarction in East Asian ancestry (**Fig. 6b**). Currently no drugs exist for MANBA but its role in lysosomal metabolism suggests that modulating its activity could have therapeutic potential. Further, increased ANGPTL4, a Tier 3 target, leads to increased triglycerides and decreased HDL cholesterol levels in African ancestries (**Fig. 6c**) and increased CAD risk in European ancestries (**Supplementary Note 7C**) supporting its potential as a

456 therapeutic target. Notably, we found that increased STAT3, a Tier 1 target, increased risk of IBD

and its subtypes, CD, and UC (Fig. 6d). Danvatirsen, a STAT-3 mRNA 3'UTR antisense inhibitor
has been undergoing phase 1 and 2 clinical trials for multiple cancer types and may be potentially
repurposed for IBD. Currently, astegolimab, an inhibitor of IL1RL1, a Tier 3 target, has completed
phase 2 trials for eczema and asthma; in our study, we find genetic support where increased
IL1RL1 increased risk of eczema in East Asian ancestry (Fig. 6e) and was concordant in
European ancestries (Supplementary Table 17). Druggability visualization for remaining
diseases for European and East Asian ancestries are provided in Supplementary Note 7C.





468 Legend information:

469 Each cell displays a putatively causal protein-phenotype association. Cell color displays the MR 470 effect estimate based on Z score averaged across cohorts capped at -10 to +10 with red showing 471 a positive Z score indicating a positive MR effect of the protein (displayed as the gene name) on 472 the phenotype and blue showing a negative Z score indicating a negative MR effect of the protein 473 on the phenotype. For simplicity, in European ancestries, we only display protein-phenotypes with 474 consistent effect across European cohorts.

475 The y-axis shows the three drug databases. DrugBank (yellow square): DrugBank³⁹ shows 476 whether the protein has an available drug in the database. Open Targets (pink square): Open 477 Targets Platform⁴⁰ shows whether the protein has available clinical trial information. Druggability: The druggable genome from Finan et al.³⁸ is shown for Tiers 1 (dark green, representing direct 478 479 targets of approved small molecules and biotherapeutic drugs), Tier 2 (dark purple, representing 480 proteins closely related to approved drug targets or which have associated drug-like compounds), 481 Tier 3 (light purple, representing secreted or extracellular proteins, those distantly related to 482 approved drug targets, and members of important druggable gene families not covered in Tier 1 483 or Tier 2), and Unclassified (gray, all other proteins not in Tiers 1 to 3). Proteins on the y-axis 484 within each tier are sorted based on the number of supported databases.

- 485 (a) European ancestry putatively causal protein-phenotype pairs stratified to Tier 1 and Tier
 486 2 druggable proteins and cardiovascular phenotypes.
 487 (b) East Asian ancestry putatively causal protein-phenotype pairs stratified to Tier 1 and Tier
 - (b) East Asian ancestry putatively causal protein-phenotype pairs stratified to Tier 1 and Tier 2 druggable proteins and cardiovascular phenotypes.
 - (c) All of the African ancestry putatively causal protein-phenotype pairs with no druggability stratification.
 - (d) European ancestry putatively causal protein-phenotype pairs stratified to Tier 1 and Tier 2 druggable proteins and autoimmune phenotypes.
 - (e) East Asian ancestry putatively causal protein-phenotype pairs stratified to autoimmune phenotypes.
- 494 495

488

489

490

491

492

493

496

497498 8. Converging evidence of the causal effect of IL1RL1 on IBD

499 Upon further stringent filtering to find protein-phenotype pairs with strong evidence (see **Methods**). 500 we found that increased IL1RL1 levels was associated with increased risk of IBD, CD, and UC in 501 European ancestries across three cohorts for IBD and CD and two cohorts for UC with consistent 502 directions, which we validated using the largest available East Asian ancestry GWAS for IBD 503 (14,393 cases and 15,456 controls), CD (7,372 cases and 15,456 controls), and UC (6,862 cases and 15,456 controls) from Liu et al.⁶⁵ (Fig. 7a). In African ancestries, the number of cases were 504 505 limited in the largest publicly available GWAS for IBD (1,285 cases and 119,314 controls)⁵⁷ and 506 UC (857 cases and 119,909 controls)⁵⁷, thus estimates were not significant likely due to 507 insufficient power.



509 510

- a) MR for the effect of IL1RL1 and IBD (top), CD (middle), and UC (bottom) in European and East Asian ancestries. PP_{max} is the maximum colocalization posterior probability between PWCoCo and SharePro.
 b) Kaplan Meier estimates for cumulative incident of IBD (top), CD (middle), and UC (bottom)
 - b) Kaplan Meier estimates for cumulative incident of IBD (top), CD (middle), and UC (bottom) by baseline IL1RL1 level quantiles in the UK Biobank. P values were computed using the log-rank test.
 - c) Bulk RNA sequencing of IL1RL1 in the ileum (left), colon (middle), and rectum (right).
 - d) Single-cell RNA sequencing analyses of *IL1RL1*. *IL1RL1* expression patterns showing 720,633 cells collected from the terminal ileum and colon of 71 donors with different levels of inflammation. Single-cell transcriptomic data was obtained from Kong et al.⁶⁶ (SCP1884 https://singlecell.broadinstitute.org/).
- 521 522

515 516

517

518

519

520

523 524

525 To triangulate the evidence, we performed supplementary observational analysis in the UK 526 Biobank using Cox proportional hazards models. We adjusted for age, sex, recruitment center, 527 Olink measurement batch, Olink processing time, and the first 10 genetic principal components. 528 Notably, none of the genetic principal components were significant, suggesting that the 529 association is not specific to ancestry. Over 10 years of follow-up, a one s.d. increase in IL1RL1 530 was associated with elevated risk of IBD (hazard ratio, HR = 1.18; 95% CI: 1.05–1.33; P = 6.6 × 531 10⁻³). CD (HR = 1.21; 95% CI: 1.00–1.47; P = 0.047), and UC (HR = 1.15; 95% CI: 1.00–1.33; P 532 = 0.048), consistent with MR findings (Supplementary Table 27). Kaplan-Meier estimates for 533 cumulative incidence of disease stratified by baseline IL1RL1 level (lowest 25% versus highest 534 25% in the UK Biobank population) also showed differences in IBD (log-rank test $P = 2.0 \times 10^{-4}$). 535 CD (log-rank test $P = 7.0 \times 10^{-3}$), and UC (log-rank test P = 0.02) (Fig. 7b). We also performed 536 an alternative, less stringent filter and prioritized proteins involved in CAD and type 2 diabetes 537 which we provide in Supplementary Note 8.

538

539 8.1. IL1RL1 expression analyses

540 To further assess the role of IL1RL1 in IBD, we used the IBD Transcriptome and 541 Metatranscriptome Meta-Analysis (IBD TaMMA) platform⁶⁷ to compare expression of *IL1RL1* 542 transcripts between IBD patients and healthy controls in the ileum, colon, and rectum. We found 543 significantly higher *IL1RL1* gene expression in all three tissues (**Fig. 7c**), suggesting that 544 increased *IL1RL1* expression is a consistent feature of IBD regardless of the specific location 545 within the gastrointestinal tract.

546

547 To gain further insights into the role of IL1RL in IBD, CD and UC, we analyzed single-cell IL1RL1 548 expression in 720,633 cells from the terminal ileum and colon of 71 participants with different 549 levels of inflammation status from Kong et al.⁶⁶ (SCP1884 https://singlecell.broadinstitute.org/). In 550 single-cell RNA sequencing, IL1RL1 showed significant enrichment in mast cells compared to 24 551 other cell types (permutation $P < 2.0 \times 10^{-4}$) (Fig. 7d). Mast cells, key players in allergic reactions and inflammation and a key cell type involved in the pathogenesis of IBD^{68,69}, may contribute to 552 553 chronic IBD by releasing inflammatory mediators like histamine and cytokines when activated by 554 IL1RL1 in inflamed tissues. These findings align with our MR analyses showing that increased 555 IL1RL1 leads to increased risk of IBD, CD and UC.

557 **Discussion**

558 In this study, we conducted comprehensive multi-ancestry proteome-phenome analyses across 559 three ancestries. Using seven large proteomics cohorts including European, African, and East 560 Asian ancestries, we analyzed 355 complex traits or diseases, and identified 3,949, 56, and 325 561 putative causal effects of protein abundance on diseases and traits, respectively. By integrating 562 data from druggable genomes and drug databases, we prioritized potential protein targets for 563 drug development. Our findings offer a comprehensive atlas of protein-phenotype associations 564 and an evidence-based resource to support drug discovery and development, expand insight into 565 disease, and highlight potential targets for therapeutic intervention.

566

567 Our study provides an updated map to earlier phenome-wide MR studies of the human plasma 568 proteome on complex diseases which were either limited to European ancestries¹⁸ or considered only a few diseases in European and African ancestries¹⁹. The significance of incorporating 569 570 multiple ancestries is underscored by our identification of several proteins that are uniquely 571 instrumentable by each ancestry due to allele frequency differences. Specifically, we 572 instrumented an additional 119 proteins exclusively in African ancestry, 17 in East Asian ancestry, 573 and 19 shared between African and East Asian ancestries. Moreover, the finding that a significant 574 proportion of population-specific genetic variants-68.5% in African and 72.2% in East Asian-575 have a MAF below 0.01 in European ancestries highlights the potential for missed genetic 576 discoveries when studies focus on a single ancestry. Research on the genetic architecture across 577 different ancestries reveals both commonalities and differences, influenced by evolutionary history, genetic diversity, and population-specific factors^{26,30}. Thus, by including diverse African 578 579 and East Asian ancestries in proteomic analyses, we were able to instrument more proteins by 580 leveraging common genetic variants in these underrepresented populations which were mostly 581 rare in European ancestries, enhancing the potential for novel discoveries and exemplifying the 582 value of including non-European individuals for comprehensive and inclusive proteomic and 583 genetic analyses.

584

585 The inclusion of African and East Asian ancestries allowed discovery of protein-phenotype 586 associations. We emphasize that these findings were attributable to uniquely instrumentable 587 proteins in each ancestry and do not necessarily indicate ancestry-specific biological mechanisms. 588

589 As an illustrative example of the value of our atlas, we found that increased circulating 590 abundances of IL1RL1 was causal for IBD, CD, and UC in European and East Asian ancestries, 591 which was supported by observational analyses and gene expression analyses. IL1RL1 could be 592 a promising therapeutic approach for IBD, potentially reducing mast cell-driven inflammation. 593 Potential drugs targeting IL1RL1 include astegolimab which has completed phase 2 trials for 594 eczema and asthma, tozorakimab which neutralizes IL-33, the interleukin that binds IL1RL1 (ST2), 595 and anakinra which targets IL-1, a closely related interleukin. However, further research is 596 required to assess the safety and efficacy of potential IL1RL1 inhibition. Many such findings may 597 exist, and this atlas may be used as a tool to facilitate the selection of targets during primary or 598 pre-clinical drug development, exploring drug repurposing opportunities, and improve 599 understanding of proteins implicated in complex traits and diseases.

600

Our study has several key strengths. We curated GWAS for a wide range of complex traits and
diseases, decreasing the overlap and redundancy and increasing the power for discovery in MR.
Second, our study included diverse proteomics cohorts, including African ancestry and a new
East Asian ancestry cohort, the Kyoto University Nagahama cohort. This latter cohort, which used
the SomaScan v4 platform, is the largest to date aside from the China Kadoorie Biobank^{25,26}.
Further, our study combined both ARIC SomaScan and UKB-PPP Olink African proteomics
cohorts at a phenome-wide scale and to include three ancestries. Notably, we identified uniquely

instrumentable proteins in African and East Asian ancestries that were not found in European
ancestries, highlighting the value of including cohorts from diverse ancestral backgrounds. Third,
we performed extensive stringent filtering on genetic instruments with strict V2G criteria. Fourth,
we harnessed novel state-of-the-art colocalization methods to reduce the risk of confounding from
LD while increasing the statistical power to support more protein-phenotype associations with
colocalization evidence³⁷.

614

615 This study has several limitations. First, while two separate proteomics cohorts were included for 616 African ancestries, the number of outcomes considered was still limited. Additionally, African 617 ancestry phenotypes were curated from various cohorts with potentially finer genetic architecture 618 differences than controlled for by using continental ancestries, thereby potentially biasing our 619 analyses. Second, differences in measurement units make direct comparison of MR effect size 620 estimates for continuous outcomes difficult. Nevertheless, direction of effect should be robust to 621 this limitation. Third, while we reduced the risk of horizontal pleiotropy by using strict V2G cis-622 pQTLs, this resulted in many proteins being instrumented by a single *cis*-pQTL, limiting the ability 623 to perform MR sensitivity analyses. Nonetheless, we used robust colocalization methods to 624 mitigate risk of reporting false positives. Fourth, although we used assays measuring nearly 5,000 625 proteins from SomaScan and 3,000 proteins from Olink, coverage is still limited with regard to the 626 entire proteome. Lastly, while we analyzed three diverse ancestries, the sample sizes for both 627 proteomic GWAS and outcome GWAS were much larger for European ancestry, leading to 628 differences in the number of associations. Greater coverage of proteins and larger sample sizes 629 in non-European ancestries are needed.

630

631 In conclusion, through integrative multi-ancestry plasma proteome-phenome MR and extensive 632 sensitivity analyses, we provided a comprehensive atlas of protein-phenotype associations 633 across three ancestries and highlighted the value of multi-ancestry inclusion, as illustrated by 634 uniquely instrumentable proteins in non-European ancestries. This study serves as a valuable 635 resource for understanding disease mechanisms and prioritizing potential new targets.

637 Methods

638 1. Proteomics cohorts

We analyzed proteomics cohorts from three ancestries consisting of European (four cohorts:
ARIC, deCODE, Fenland, and UKB-PPP), African (two cohorts: ARIC and UKB-PPP), and East
Asian (one cohort: Kyoto University Nagahama East Asian cohort). All cohorts had proteomics
measured on the aptamer-based SomaScan assay v4 except for the UKB-PPP study, which used
the antibody-based Olink Explore 3072 platform.

644 645 *1.1. European ancestry cohorts*

646 We analyzed the GWAS of protein levels in individuals of European ancestry using four different 647 studies. Three of these four studies (ARIC, deCODE, and Fenland described below) had 648 proteomics measurements from the aptamer-based SomaScan assay v4 from SomaLogic 649 (Boulder, Colorado, USA). In brief, SomaScan assay v4 uses aptamers, which are single-650 stranded oligonucleotides that have specific binding affinities to protein targets and can measure 651 up to 5,000 unique proteins. The UKB-PPP study had proteomics measurements from the 652 antibody-based Olink Explore 3072 platform, which measures up to 3,000 proteins. Briefly, Olink 653 (Uppsala, Sweden) uses proximity extension assay (PEA) technology, which detects proteins 654 through the binding of two separate antibodies carrying complementary oligonucleotide tags. 655 which hybridize to the protein target. We restricted our analyses to proteins encoded by autosomal 656 genes and analyzed a list of 4,687 proteins from SomaLogic and 2,823 proteins from Olink. 657

658 1.1.1. ARIC

The Atherosclerosis Risk in Communities (ARIC)¹³ measured protein levels from 9,084 American
 participants of European and African ancestries using the SomaScan assay v4. Of these
 participants, 4,657 plasma proteins were measured for 7,213 European American individuals.

663 1.1.2. deCODE

The deCODE study¹² provided 4,907 aptamers that measure 4,719 proteins in 35,559 Icelandic
 individuals of European ancestry using SomaScan assay v4.

667 1.1.3. Fenland

The Fenland study¹¹ measured 4,775 proteins in 10,708 individuals of European ancestry using
SomaScan assay v4.

671 *1.1.4. UKB-PPP*

The UK Biobank Pharma Proteomics Project (UKB-PPP)¹⁴ conducted proteomic profiling on 54,219 individuals of multiple genetic ancestries in the UK Biobank using the Olink Explore 3072 platform. From this cohort, 34,557 European individuals, each with 2,923 unique proteins measured, were utilized in the UKB-PPP as the discovery cohort, and we used this discovery cohort in our study.

677

678 1.2. African ancestry cohorts

- 679 For African ancestries, we used GWAS of protein levels from two different studies (ARIC and 680 UKB-PPP) measured on SomaScan assay v4 and Olink Explore 3072, respectively.
- 681

682 *1.2.1. ARIC*

The ARIC cohort was previously described in the European cohort section and consists of 9,084
 European American and African American individuals. Of these participants, 4,657 proteins from
 the SomaScan v4 assay were measured for 1,871 African American individuals.

- 686
- 687 1.2.2. UKB-PPP

688 We used the African ancestry individuals from the UKB-PPP study, which consists of 931 689 individuals with 2,923 unique proteins measured with the Olink Explore 3072 platform.

- 690
- 691 1.3. East Asian ancestry cohort
- 692 Kyoto University Nagahama cohort

693 The Nagahama Primary Prevention Cohort Project (Kyoto-Nagahama cohort) is a joint project 694 between the Kyoto University Graduate School of Medicine and Nagahama City, Shiga Prefecture 695 involved 10,000 residents Nagahama (https://w3.genome.med.kyotothat of 696 u.ac.jp/en/nagahama-project/). Data generation was performed at the Kyoto University Center for 697 Genome Medicine, where 1.823 Japanese individuals of East Asian ancestry were whole 698 genome-sequenced and had 4,196 proteins measured using the SomaScan assay v4. Further 699 details can be found in Supplementary Note 9.

700

2. Identification of strict variant-to-gene (v2g) *cis*-protein quantitative trait loci (pQTLs)

For example, the ARIC¹³ study defined *cis*pQTL as those within ±500 kb of the TSS of the protein-coding gene with FDR < 5%. The deCODE study¹² defined *cis*-pQTL as those within ±1 Mb of the TSS of the protein-coding gene with $P < 1.8 \times 10^{-9}$. The Fenland Study¹¹ defined *cis*-pQTL as those within ±500 kb of the protein-coding gene with $P < 1.004 \times 10^{-11}$. The UKB-PPP¹⁴ noted that of their identified pQTLs, 66.9% of proteins tested (1,954 of 2,922 proteins) had a *cis*-pQTL within ±1 Mb of the proteincoding gene with $P < 1.7 \times 10^{-11}$. Thus, we created a common *cis* definition as follows.

710 2.1. Linkage disequilibrium (LD) clumping

711 We performed LD clumping (clumping window of 1 Mb, significance level of 5×10^{-8} , and clumping 712 r^2 threshold of 0.001) on each proteomic GWAS in each ancestry cohort. For European proteomic 713 GWAS, we used a reference panel composed of 50,000 randomly sampled unrelated UK 714 Biobank⁷⁰ individuals of European ancestry (UKB 50k). For African proteomic GWAS, we curated 715 an LD reference panel from the Human Genome Diversity Project and 1000 Genomes Project 716 (HGDP + 1kGP) reference panel⁷¹ for 994 African ancestry individuals. In East Asian ancestries, 717 we used the 1000 Genomes East Asian (1kGP EAS) reference panel. We retained variants with 718 a MAF > 0.01 in all reference panels.

719

720 2.2. Identification of cis- and trans-pQTLs

To determine *cis*-pQTLs, we used the Ensembl BioMart⁷³ package version Ensembl 105: Dec 721 722 2021 (Genome Reference Consortium Human Build 38, GRCh38.p13) to generate a human 723 protein-coding genes file. We considered relevant attributes such as the canonical TSS, gene 724 start, and gene end. Since Fenland proteomic GWAS were in GRCh37 coordinates, we 725 separately curated a protein-coding genes file using BioMart on the GRCh37 assembly in 726 Ensembl using version 110. For analysis, we excluded protein-coding genes located on sex 727 chromosomes and those located within the major histocompatibility complex (MHC) (GRCh38: 728 chr6 28,510,020– 33,480,577) due to the complex LD structure, high allelic diversity, and strong 729 pleiotropy in this region⁷⁴. We defined a *cis*-pQTL as a pQTL within ± 500 kb of the TSS of the 730 protein-coding gene. All other pQTLs were trans. We used cis-pQTLs since they are more likely 731 to directly impact the transcription and translation of the protein of interest.

732

733 2.3. Strict variant-to-gene cis-pQTL definition

734 To minimize potential horizontal pleiotropic effects, we defined a unique strict V2G definition for

735 *cis*-pQTLs whereby a *cis*-pQTL is a *strict V2G cis-pQTL* if it is a *cis*-pQTL for only one protein-

- coding gene (*strict*) and it has the strongest link to the corresponding protein-coding gene based
- 737 on multiple sources of evidence and concomitantly has the highest Open Targets Genetics V2G

739

740 2.3.1. Strict cis-pQTL

For genome-wide significant independent *cis*-pQTLs in each cohort, we retained those associated with a single protein-coding gene. We defined this as a "strict" *cis*-pQTL definition. Here, we used *protein-coding genes* instead of *proteins* (*aptamers*) for the strict *cis* definition due to SomaScan assay v4 having multiple instances where two or more aptamers target a single protein, which would result in the unwarranted scenario where pQTLs that were *cis* for more than one aptamer of the same protein were removed.

747

748 2.3.2. Open Targets Genetics Variant-to-Gene (V2G)

749 We used the Open Targets Genetics³⁵ database (https://genetics.opentargets.org/) to determine 750 whether each strict *cis*-pQTL held the highest V2G score for its corresponding associated 751 protein-coding gene. This ensures that the variant is a suitable proxy for the plasma levels of the 752 protein-coding gene. For instance, the variant may directly impact the protein-coding gene, 753 potentially by altering its transcription, thereby influencing its plasma abundance. Briefly, Open 754 Targets Genetics V2G scores are generated through a model trained on molecular QTLs 755 (eQTL, sQTL, pQTL), chromatin interaction experiments such as promoter capture Hi-C (PCHi-756 C), in silico functional predictions such as Ensembl Variant Effect Predictor (VEP), and the 757 distance between variants and genes' canonical transcription start sites. This composition of 758 evidence enables accurate assignment of variants to genes. We note that the pQTL datasets 759 utilized by Open Targets Genetics for model training are from earlier studies and do not overlap 760 with the proteomics cohorts we analyzed in this study, mitigating the risk of overfitting.

761

Together, we combined the strict *cis* definition and V2G score from Open Targets Genetics to
 curate strict V2G *cis*-pQTLs, which decreases the chance of horizontal pleiotropy in MR.

764

765 3. Protein altering variant (PAV) and expression quantitative trait loci (eQTL)

766 Since pQTLs altering the binding epitope of a protein may reflect assay specificity instead of the 767 true biological function of the protein^{1,64,75}, we annotated strict V2G *cis*-pQTLs as protein-altering 768 variants (PAVs) of moderate or high impact if they or any variants in high LD ($r^2 > 0.8$) were 769 identified as PAVs using VEP⁷⁶ (Supplementary Tables 2-8). VEP annotates the impact of 770 variants into four categories: Modifier. Moderate. Low, or High 771 (https://useast.ensembl.org/info/genome/variation/prediction/predicted data.html). "Modifier" 772 impact refers to variants in non-coding regions or affecting non-coding genes in which evidence 773 of impact is hard to predict or limited. "Low" impact refers to variants that may not change protein 774 behavior. "Moderate" impact variants are non-disruptive and may change protein effectiveness 775 and include missense variants. "High" impact variants are disruptive and can cause truncation of 776 proteins, loss of function, or trigger nonsense-mediated decay. If a strict V2G cis-pQTL and all of 777 its LD proxies were labelled as Modifier or Low impact, we considered this strict V2G cis-pQTL to 778 have no PAV. If a strict V2G cis-pQTL or any of its LD proxies were labeled Moderate or High, we 779 considered this strict V2G cis-pQTL to be a PAV of -Moderate or -High impact, respectively.

780

781 We also conducted a *cis*-eQTL enrichment analysis using 49 tissues from GTEx v8 in European 782 ancestries since if a *cis*-pQTL overlaps with the *cis*-eQTL of the same gene, it strengthens the 783 evidence that the *cis*-pQTL acts directly on the gene products, reducing the risk of horizontal 784 pleiotropy. To do so, we determined whether the strict V2G cis-pQTL was a cis-eQTL for the protein-coding gene of interest with $P < 1 \times 10^{-5}$ by querying 49 tissues in GTEx v8 European 785 eQTLs. Querving was performed based on CHR:POS:EA:NEA and CHR:POS:NEA:EA (CHR: 786 787 chromosome; POS: position; EA: effect allele; NEA: non-effect or other allele). In order to 788 compare effect sizes across different cohorts, we aligned the strict V2G cis-pQTL effect allele in 789 each cohort to the corresponding ancestry-specific reference panel's alternative allele (UKB 50k

for European, HGDP + 1kGP for African, and 1kGP EAS for East Asian ancestries). We note that pQTLs in the Fenland cohort were based on GRCh37 coordinates and were lifted to GRCh38 prior to querying for *cis*-eQTLs across GTEx v8 tissues. After synchronizing the genome assembly, 7 out of 3,045 strict V2G *cis*-pQTLs in Fenland were automatically labeled as having no eQTL since no corresponding GRCh38 coordinate existed for these variants, and 3,038 *cis*-pQTLs were analyzed.

796

797 4. GWAS outcome curation and selection

We manually curated the latest and largest GWAS (as of February 2024) for European and African
 ancestry outcomes. East Asian ancestry outcomes were selected from BioBank Japan¹⁵. All
 GWAS in GRCh37 were lifted to GRCh38 with the liftOver tool. We outline the curation steps for
 each ancestry's GWAS outcomes in detail:

802

803 *4.1. European ancestry outcomes*

804 During the curation, we considered 510 outcomes downloaded from 50 studies and one database. 805 We removed outcomes that were duplicated and had a larger GWAS available, had ambiguous 806 or broad definitions, were likely heterogeneous, were sex-specific, had missing relevant columns, 807 had no download link available, and were not relevant to our outcomes of interest. We retained 808 179 outcomes for analysis (Supplementary Table 10). We labeled rsids, chromosome, and 809 position if they were missing. Cases and sample sizes for each outcome were manually extracted 810 from the original manuscript or the supplementary tables of each corresponding study, and we 811 further categorized each outcome into one of 23 "Type" categories pertaining to the human 812 system the outcome was based on or most likely to fall under. 813

814 4.2. African ancestry outcomes

815 We manually curated 26 of the most up-to-date and publicly available African ancestry GWAS 816 summary statistics (Supplementary Table 11). Restricted access GWAS from dbGap were not 817 considered due to data access difficulties. Missing sample sizes were manually annotated with 818 the sample size by inspecting the original manuscript and Supplementary Tables, while missing 819 rsids were labeled using the HGDP + 1kGP reference panel. We annotated missing chromosomes 820 and positions and used VEP to annotate variants missing effect allele frequency with 821 gnomADg AFR AF (gnomAD genomes for African/American populations). We categorized each 822 outcome into one of 8 "Type" categories, including Respiratory, Musculoskeletal, Cardiovascular, 823 Eye, Anthropometry, Biomarker, Psychiatric, and Metabolic/endocrine.

824

As exploratory analyses, we included 114 binary cardiovascular and 9 binary autoimmune-related outcomes from the Million Veteran Program (n = 635,969) for African individuals.

828 4.3. East Asian ancestry outcomes

We curated 220 outcomes from Biobank Japan (<u>https://pheweb.jp/</u>)²⁷. We excluded 14 sexspecific outcomes, including Abortion, Breast cancer, Cervical cancer, Cesarian section, Ectopic pregnancy, Endometriosis, Endometrial cancer, Mastopathy, Ovarian cancer, Ovarian cyst, Preeclampsia, Prostate cancer, Uterine fibroid, and Uterine prolapse. We analyzed 206 outcomes (**Supplementary Table 12**). We also provide "Type" labels for each outcome, denoting the human system the outcome was based on or was most likely related to.

835

836 5. Two-sample Mendelian randomization

To assess the putative causal effect of protein abundance on outcomes in European, African, and
 East Asian ancestries, we performed two-sample MR using TwoSampleMR v.0.5.7⁷⁷. To mitigate
 horizontal pleiotropy, we used strict V2G *cis*-pQTLs as instrumental variables to proxy protein-

840 level exposures, as defined in the earlier sections of the methods. We harmonized exposure and

841outcome GWAS using the harmonise_data() function and performed a proxy search if an842instrument was absent in the outcome GWAS. For European, African, and East Asian ancestries,843we used the UKB 50k, HGDP + 1kGP, and 1kGP East Asian reference panels that were844previously used for LD clumping for the proxy search, respectively. We searched for proxies using845PLINK v.1.9⁷² parameters --Id-window=5000, --Id-window-kb=5000, --Id-window-r2=0.8 and846retained proxies with minor allele frequencies ≤ 0.42.

847

848 MR analyses were performed using the mr() function. For proteins with a single genetic instrument, 849 the association between the protein and outcome was evaluated using a Wald ratio estimate. For 850 proteins with ≥ 2 genetic instruments, we used an inverse variance weighted random effects 851 estimate. We determined whether genetic instrumental variables had F-statistics > 10^{43,78}, 852 indicating strong associations with the exposure and thus less chance of weak instrument bias, 853 which may bias the causal effect estimates towards the null in two-sample MR. F-statistics are 854 shown in Supplementary Table 9. We corrected for multiple testing per cohort in each ancestry 855 by applying a Benjamini-Hochberg-corrected P threshold (FDR)⁴⁵ of 0.05 (5%) as done 856 previously¹⁹.

857

We note that beta estimates from MR for continuous traits are not directly comparable across
different outcomes because the units used in GWAS vary. For instance, some GWAS use clinical
units, whereas others use standardized and/or residualized values.

862 6. MR sensitivity analyses

To increase the robustness of MR findings, we further filtered MR results based on multiple sensitivity analyses, including heterogeneity tests, MR using alternative approaches (including weighted median, weighted mode, MR-Egger), and Steiger directionality test⁷⁹ to assess reverse causality. Following these filtering steps, we refer to retained protein-phenotype associations as "MR-passing". The sensitivity analyses are described for proteins with \geq 2 instruments and proteins with \geq 3 instruments below:

869

870 6.1. Sensitivity analyses for proteins with two or more instruments

The heterogeneity test was performed for proteins with ≥ 2 instruments and describes whether strict V2G *cis*-pQTLs of the same protein are likely to show comparable effects on the tested outcomes. For heterogeneity testing, we used the mr_heterogeneity() function to compute a heterogeneity *P* value (Q_pval), and we calculated l^2 statistics using the "Isq()" function. If an association had an l^2 threshold ≥ 0.5 and a heterogeneity *P* value (Q_pval) < 0.05, this indicated considerable heterogeneity.

877

878 6.2. Sensitivity analyses for proteins with three or more instruments

879 For proteins with \geq 3 genetic instruments, we performed additional sensitivity analyses with 880 alternative MR methods such as MR weighted median, MR weighted mode, and MR-Egger 881 methods, as well as Steiger directionality testing. To check the consistency of MR estimates, we 882 required the MR estimate, as well as the sensitivity analyses estimates from the MR weighted 883 median, MR weighted mode, and MR-Egger approaches, to all have the same sign. For 884 directional pleiotropy, we used the mr pleiotropy test() function to perform the MR-Egger 885 intercept test and considered a P value < 0.05 as a statistically significant deviation from the null 886 and an indication of directional pleiotropy. We also performed Steiger filtering on all proteins 887 (directionality test() function). Any pQTLs that explain more variance in the outcome than in the 888 exposure potentially indicate reverse causation and were removed from further analysis.

- 889
- 890 7. Colocalization of proteomic GWAS with outcome GWAS

891 To assess whether plasma protein levels share the same causal variant with GWAS outcomes, 892 we employed two colocalization methods to ensure the robustness of our findings. We performed PWCoCo^{18,36} and SharePro⁸⁰ for MR associations that passed all MR sensitivity analyses 893 894 described in the previous section. PWCoCo and SharePro are recent methods allowing multiple 895 independent associations to be assessed. Both improve on the original coloc method⁸¹, which 896 was limited by the assumption that a single variant exists per GWAS, wherein the method only 897 considers the strongest of these distinct association signals when multiple independent 898 associations exist. A detailed description of both methods is provided in Supplementary Note 899 **10**. Colocalization analyses were performed around a 1-Mb region centered on the lead (lowest 900 *P* value) *cis*-pQTL. We set a colocalization posterior probability (PP) of a shared causal variant \geq 901 0.8 in any of PWCoCo or SharePro as evidence of colocalization. For simplicity, we report the maximum PP between PWCoCo and SharePro (PP_{max}) in the main text. We reported putatively 902 903 causal associations (Supplementary Tables 14-16) as associations which pass all MR 904 sensitivity analyses, Steiger filtering, and also colocalized with a PP ≥ 0.8 in any one of PWCoCo 905 or SharePro. Due to the difficulty in verifying the corresponding Olink assay target for each 906 SomaScan aptamer, we counted unique protein-phenotype associations using protein-coding 907 genes to harmonize between proteomics platforms. Supplementary Tables 14-16 contain MR 908 and colocalization summaries, summaries of effect direction consistency across cohorts, flags of 909 proteins instrumented by PAVs of high impact, and whether the protein-phenotype association 910 came from a protein uniquely instrumentable in that ancestry. We also annotate whether protein-911 phenotype associations passed the most stringent Bonferroni correction for the total number of 912 MR tests across three ancestries (P < 0.05 / 874,465).

913

8. Distinguishing between previously reported and unreported protein-phenotype pairs

- 8.1. Comparing against earlier studies identifying putatively causal protein-phenotype pairs with
 MR and colocalization evidence
- 917 To identify the status of protein-phenotype associations as reported or unreported (not found from 918 pre-existing proteome-phenome-wide MR studies) in European ancestries, we overlapped our 919 associations with recent proteome-phenome-wide MR analyses from Zheng et al. 2020¹⁸ and 920 multi-ancestry proteome-wide MR analyses from Zhao et al. 2022¹⁹. We used the 111 identified 921 putatively causal associations (65 proteins on 52 phenotypes) from "Table S7" of Zheng et al.¹⁸ 922 and the 45 associations from "ST7A" of Zhao et al.¹⁹.
- 923

We note that Zhao et al. used three colocalization methods and a more relaxed threshold of PP
 > 0.7 as evidence of colocalization. To do so, we harmonized the outcomes from Zheng and Zhao
 to our outcomes and matched protein-phenotype pairs using ensembl ID and outcome name for
 the Zheng study and UniProt ID and outcome name for the Zhao study. Any of our identified
 putatively causal European ancestry association pairs not from these two studies were identified
 as unreported.

930

For African ancestries, we overlapped our protein-phenotype pairs with the single protein phenotype pair passing FDR correction identified in "ST8A" of Zhao et al¹⁹. Any pairs that did not
 overlap were considered unreported.

- 934
- 935 9. Protein-phenotype network plots
- We used Python igraph (v.0.10.8) to generate networks. Placement of nodes was generated
 in Cytoscape (v.3.10.2) and manual editing was performed in Adobe Illustrator (v.28.2).
- 939 10. Effect concordance within-ancestry
- 940 To determine the concordance of MR effect estimates within European and African ancestries, 941 which both included more than one proteomics cohort, we compared protein-coding genes to

harmonize assay names across SomaScan assay v4 and Olink 3072 Explore platforms. However,
we could not verify whether Olink assays for a particular protein targeted the same domain as its
corresponding SomaScan assay, which may lead to differences in MR estimates. Moreover, since
some proteins were instrumented by *cis*-pQTLs that may be PAVs of high impact, which could
also lead to discordance in effects, we annotated these associations with a flag and advise caution
in the interpretation of these flagged results (Supplementary Tables 14–16).

948

956

949 11. Druggability assessment

We performed a druggability assessment on instrumentable protein-coding genes and protein phenotype associations that showed putatively causal relationships.

953 11.1. Druggability assessment of instrumentable proteins

954 For instrumentable protein-coding genes, we determined druggability based on Finan et al.³⁸, as 955 described below.

957 11.1.1. Finan et al.

Finan et al.³⁸ considered 20,300 protein-coding genes annotated using Ensembl version 73 and 958 classified 4,479 (22%) into three tiers (Tier 1, 2, and 3) as drugged or druggable. Tier 1 (1,427 959 960 genes) encompasses the primary targets of approved small molecules and biotherapeutic drugs, 961 along with those influenced by clinical-phase drug candidates. Tier 2 (682 gene) involves proteins 962 closely associated with drug targets or linked to drug-like compounds. Meanwhile, Tier 3 (2,370 963 genes) comprises secreted or extracellular proteins that are distantly related to approved drug 964 targets, and those in important druggable gene families not covered in Tiers 1 or 2. We denoted 965 all other protein-coding genes that did not fall in Tier 1, 2, or 3 categories as "Unclassified". To 966 check the overlap of protein-coding genes across all cohorts when stratifying into Finan et al. tiers 967 UpSetR⁸² (Supplementary Figure 5 and 6), we used the package v.1.4.0 968 (https://github.com/hms-dbmi/UpSetR).

969

970 11.2. Druggability assessment and enrichment of protein-phenotype pairs

We first assessed how many proteins from the identified pairs of putatively causal proteinphenotype associations had existing drugs by querying DrugBank³⁹. We used Ensembl ID to match protein targets in DrugBank. Next, we created heatmaps using the pheatmap package v.1.0.12 in R incorporating the druggable genome, Drugbank, and Open Targets Platform (described below). We overlapped protein-phenotype pairs with the DrugBank database to determine whether any drugs existed for these disease-implicated proteins while Open Targets was used to determine whether any clinical trial information existed for these proteins.

979 11.2.1. DrugBank

980 We used DrugBank database v.5.1.12 (<u>https://go.drugbank.com/releases/latest</u>) and R package 981 dbparser v.2.0.2 to parse the DrugBank database xml file. We aggregated drugs and target 982 information and overlapped this with putatively causal protein-phenotype associations to 983 determine which proteins had an available drug. 984

985 11.2.2. Open Targets

986 We used Open Targets v.24.03 (<u>https://platform.opentargets.org/downloads</u>). We used the 987 knownDrugsAggregated dataset, which provides information on known drugs for a given disease 988 and contains protein target information. Open Targets also includes information on clinical trials 989 and its phases, which we used to determine the status of a protein target and its corresponding 990 drug. We overlapped this dataset by matching Ensembl ID with our identified protein-phenotype 991 associations.

993 12. Follow-up analyses showing evidence for IL1RL1

We filtered protein-phenotype associations in European ancestries by ensuring that the protein was instrumented in two or more ancestries, was targeted by both SomaScan and Olink assays, had an MR effect that was concordant across all cohorts, had MR and colocalization evidence in three or more cohorts for the protein-phenotype pair, and was implicated in at least one binary phenotype. Following these filtering steps, we identified 53 candidate protein-phenotype pairs and highlighted IL1RL1 in IBD, CD, and UC as an illustrative example.

1000

1001 12.1. MR and colocalization of IL1RL1 in East Asian ancestry with IBD, CD, and UC

1002 To validate that IL1RL1 was also putatively causal for IBD, CD, and UC in non-European 1003 ancestries, we performed two-sample MR and colocalization with PWCoCo and SharePro using 1004 the IL1RL1 strict V2G cis-pQTL (rs12712135) identified in the Kyoto University Nagahama East 1005 Asian ancestry proteomics cohort. We used the largest East Asian ancestry GWAS for IBD 1006 (14,393 cases and 15,456 controls), CD (7,372 cases and 15,456 controls), and UC (6,862 cases and 15,456 controls) from Liu et al.⁶⁵. Harmonization was performed similarly to the primary 1007 1008 analyses, and we used the Wald ratio to obtain MR effect estimates. Colocalization was 1009 performed as described in earlier sections of the methods, and we used PP \geq 0.8 as the threshold 1010 for evidence of colocalization in any of PWCoCo or SharePro.

1011

1012 12.2. MR of IL1RL1 in African ancestry with IBD and UC

We estimated the causal effect of IL1RL1 on IBD (1,285 cases and 119,314 controls) and UC (857 cases and 119,909 controls) in African ancestry using outcome GWAS from the Million Veteran Program⁵⁷. We performed MR using the IL1RL1 strict V2G cis-pQTL (rs1420101) in the ARIC and UKB-PPP African ancestry cohorts. The Wald ratio was used to obtain MR effect estimates.

1019 12.3. Cox regression analysis for 10-year cumulative events of IBD, CD, and UC in the UK 1020 Biobank

1021 We used multivariable Cox proportional hazards regression to determine whether baseline 1022 plasma IL1RL1 protein level was associated with cumulative events of IBD, CD, or UC. We 1023 adjusted for age, sex, recruitment center, Olink measurement batch, Olink processing time, and 1024 the first 10 genetic principal components (UKB field: 22009) to adjust for genetic ancestry while 1025 protein levels were rank-based inverse normal transformed. We used the coxph() function from 1026 the survival R package v.3.2.13 and considered P < 0.05 as nominal significance of association. 1027 None of the genetic principal components were significantly associated with the IBD, CD, and UC 1028 outcome in each analysis. 1029

1030 We checked the proportional hazards assumption using the cox.zph() function for the IL1RL1 1031 association analysis for IBD, CD, and UC for the rank-based inverse normal transformed IL1RL1 1032 covariate. The proportional hazards assumption tests the null hypothesis that each covariate's 1033 effect estimate does not vary with time. We used the "GLOBAL" variable from cox.zph() which 1034 tests the null hypothesis of whether all inputted covariates meet the proportional hazards 1035 assumption. We considered P < 0.05 as evidence that the proportional hazards assumption was 1036 not fulfilled.

1037

1038 We defined IBD using ICD10 codes K50-K51, CD using K50, and UC using K51. We calculated 1039 the time to event by subtracting the date of event registration from the date of enrollment (data 1040 field: 53), focusing on events occurring within 10 years of enrollment. We excluded cases of 1041 prevalent IBD that met these criteria before enrollment, and those without a recorded event date 1042 for IBD, and performed the same steps for CD and UC. Controls for the IBD analysis were defined 1043 as individuals without an IBD, UC, or CD record based on self-reported medical history. Controls

for the CD analysis were defined as individuals without a CD record, and controls for the UC
analysis were defined as individuals without a UC record, based on self-reported medical history.
We analyzed 333 cases and 40,001 controls for IBD, 130 cases and 40,388 controls for CD, and
240 cases and 40,178 controls for UC.

- We plotted Kaplan-Meier curves by stratifying individuals into the bottom 25% and the top 25% based on baseline plasma IL1RL1 levels. We performed a log-rank test to assess whether there is a statistically significant difference in survival between these two groups, with a nominal P <0.05.
- 1053

1054 12.4. Bulk RNA-sequencing

1055 We used the IBD Transcriptome and Metatranscriptome Meta-Analysis (IBD TaMMA) platform⁶⁷ 1056 (https://ibd-meta-analysis.herokuapp.com) to evaluate changes in IL1RL1 gene expression. IBD 1057 TaMMA encompasses 3,853 RNA-Seg datasets from 26 studies on IBD and control samples 1058 across various tissues. All datasets were processed using a uniform computational pipeline and 1059 underwent batch correction for harmonizing data, enabling consistent comparison across studies. 1060 Differential expression results for ileum, colon, and rectum biopsies from CD and UC patients 1061 versus healthy controls were downloaded from IBD TaMMA. We assessed the log₂ fold change 1062 to create forest plots.

1063

1064 12.5. Single-cell RNA-sequencing

1065 To gain a better understanding of the enrichment of *IL1RL1* in specific cell types, we obtained 1066 single-cell **RNA** sequencing data from Kong et al.66 (SCP1884 from https://singlecell.broadinstitute.org/), which profiled 720,633 cells from the terminal ileum and 1067 1068 colon of 71 CD individuals with different levels of inflammation. Specific details of sample 1069 collection, data processing, and single-cell profiling have been described previously⁶⁶. We 1070 evaluated the normalized gene expression levels of IL1RL1 in 25 different cell types and replotted 1071 the first two dimensions of Uniform Manifold Approximation and Projection (UMAP) coordinates 1072 to visualize the cell clusters. To determine if *IL1RL1* was more significantly expressed in certain 1073 cell types, we conducted 5,000 permutations of the cell type labels. We assessed how often a 1074 specific cell type had the same or a higher proportion of cells expressing *IL1RL1* compared to all 1075 the cells in the overall population (permutation P value). 1076

1077 13. STROBE-MR statement

1078 Our study closely adheres to the STROBE-MR guidelines and the STROBE-MR checklist is 1079 attached in **Supplementary Note 11**. 1080

1081 14. Ethics declarations

All contributing cohorts obtained ethical approval from their institutional ethics review boards. The
contributing proteomics cohorts include the Atherosclerosis Risk in Communities (ARIC) Study,
deCODE study, Fenland study, UK Biobank, and Kyoto University Nagahama study. The UK
Biobank has approval from the North West Multi-centre Research Ethics Committee as a
Research Tissue Bank.

1087

1088 15. Data availability

1089 We will provide unfiltered proteome-phenome wide MR results for European (ARIC, deCODE, 1090 Fenland, UKB-PPP), African (ARIC, UKB-PPP), and East Asian (Kyoto University Nagahama 1091 cohort) ancestries on FigShare upon publication. We caution against directly comparing MR effect 1092 estimates across continuous outcomes, as the outcomes were collected from various sources 1093 and may not be scaled to the same units.

1095 15.1. Proteomic GWAS

- 1096 ARIC summary statistics (EUR and AFR): <u>http://nilanjanchatterjeelab.org/pwas/</u>
- 1097 deCODE summary summary statistics (EUR): <u>https://www.deCODE.com/summarydata/</u>
- 1098 Fenland summary statistics (EUR): <u>https://omicscience.org/apps/pgwas/</u>
- 1099 UKB-PPP summary statistics (EUR, AFR): <u>http://ukb-ppp.gwas.eu/</u>
- 1100 Kyoto University Nagahama cohort summary statistics (EAS): Available through contacting the 1101 authors of this study.
- 1102
- 1103 15.2. Outcome GWAS summary statistics
- 1104 Information on the 179 European outcomes used in this study and the link to the original summary 1105 statistics is available in **Supplementary Table 10**.
- 1106 The 26 African outcomes are available in **Supplementary Table 11**.
- 1107 The 206 East Asian outcomes are available in **Supplementary Table 12**.
- 1108 Million Veteran Program outcomes can be found in the original study⁵⁷.
- 1109 The largest East Asian ancestry IBD, CD, and UC GWAS are publicly available from Liu et al.⁶⁵
- 1110 and were downloaded from https://www.ibdgenetics.org/
- 1111
- 1112 15.3. Variant-to-gene score
- 1113 Open Targets Genetics³⁵ (<u>https://genetics-docs.opentargets.org/data-access/data-download</u>),
- 1114 1115 *15.4. Reference panels*
- 1116 European ancestries: UKB 50k (<u>https://www.ukbiobank.ac.uk/</u>)
- 1117 East Asian ancestries: 1000 Genomes Project (<u>https://www.internationalgenome.org/data</u>)
- 1118 African ancestries: HGDP+1KG (https://gnomad.broadinstitute.org/news/2020-10-gnomad-v3-1-
- 1119 <u>new-content-methods-annotations-and-data-availability/#the-gnomad-hgdp-and-1000-genomes-</u> 1120 callset)
- 1120

1122 *15.5. Druggability*

- 1123 We used Finan et al. 2017³⁸ for the list of 4,479 protein-coding genes in each druggability tier, 1124 DrugBank database v.5.1.12 (https://go.drugbank.com/releases/latest) for information on drugs 1125 specific proteins, and Open Targets Platform database v.24.03 targeting (https://platform.opentargets.org/downloads) for clinical trial phase and status information for 1126 1127 protein-drug-disease triplets.
- 1128
- 1129 15.6. Expression analyses
- 1130 For gene expression data, we used data from Kong et al.⁶⁶ (SCP1884 at Single Cell Portal 1131 <u>https://singlecell.broadinstitute.org/</u>). 1132
- 1133 16. Code availability
- 1134 We used R v.4.1.2 (https://www.r-project.org/),
- 1135 Python 3.10 (<u>https://www.python.org/downloads/release/python-3100/</u>)
- 1136 PLINK v.1.9⁷² (<u>http://pngu.mgh.harvard.edu/purcell/plink/</u>),
- 1137 TwoSampleMR v.0.5.6 (https://mrcieu.github.io/TwoSampleMR/),
- 1138 coloc v.5.2.3⁸¹ (<u>https://chr1swallace.github.io/coloc/</u>),
- 1139 PWCoCo^{18,36} (<u>https://github.com/jwr-git/pwcoco</u>),
- 1140 SharePro v.5.0.0⁸⁰ (<u>https://github.com/zhwm/SharePro_coloc/</u>),
- 1141 Cytoscape v.3.10.2 (https://cytoscape.org/),
- 1142 LocusZoom⁸³ (<u>https://my.locuszoom.org/</u>)
- 1143
- 1144 Code used in this study will be made available at <u>https://github.com/chenyangsu/pQTL-MR</u> upon
- 1145 publication.

1146

1147 17. Acknowledgments

This research has been conducted using the UK Biobank Resource under Application Number 27449. C.-Y.S. is supported by a CIHR Canada Graduate Scholarship Doctoral Award (Funding Reference Number: 187673), an FRQS doctoral training scholarship, and a Lady Davis Institute/TD-Bank Scholarship. T.L. has been supported by start-up funding from the Office of the Vice Chancellor for Research and Graduate Education, School of Medicine and Public Health, and Department of Population Health Sciences at the University of Wisconsin-Madison. S.Y. is supported by the Japan Society for the Promotion of Science.

1155

1156 The funders had no role in the study design, data collection and analysis, decision to publish, or 1157 preparation of the manuscript. We acknowledge Servier Medical Art (https://smart.servier.com/) 1158 for providing images that were used to create diagrams in this study.

- 1159
- 1160 18. Author contributions
- 1161 Conception and design: C.-Y.S., T.L., S.Y.
- 1162 Methodology: C.-Y.S., W. Z., T.L., S.Y.
- 1163 Data curation: C.-Y.S., S.S.-H., T.-Y.Y., K.Y.H.L., Y.C., F. M., T.L., S.Y.
- 1164 Data Analysis: C.-Y.S., T.L., S.Y.
- 1165 Visualization: C.-Y.S., A.v.d.G., T.L., S.Y.
- 1166 Knowledge portal: C.-Y.S., D.-K.J., M.C., S.Y.
- 1167 Writing—Original Draft: C.-Y.S.
- 1168 Writing—Review and Editing: all authors
- 1169 Supervision: T.L., S.Y.
- 1170 Project administration: T.L., S.Y.
- 1171 Funding acquisition: V.M., S.Z., T.L., S.Y.
- 1172
- 1173 19. Competing Interests
- 1174 Y.C. is an employee, J.B.R. is the CEO, and W.Z., G.B.-L., and T.L. have been consulting for 5
- 1175 Prime Sciences. However, this study was performed separately with no relationship to 5 Prime
- 1176 Sciences. J.B.R.'s institution has received investigator-initiated grant funding from Eli Lilly,
- 1177 GlaxoSmithKline and Biogen for projects unrelated to this research. The other authors declare no
- 1178 conflict of interest.
- 1179

1180 Extended Data Figures



- 1182 1183
- 1183 Extended Data Fig. 1. Flow diagram showing the definition of strict variant-to-gene *cis*-
- 1184 pQTLs.
- 1185 Flow diagram showing the selection of strict variant-to-gene (V2G) *cis*-pQTLs used as instruments
- 1186 for MR starting from the proteomic GWAS.



- 1187 1188
- 1189

1194

1195

Extended Data Fig. 2. eQTL enrichment analysis comparing strict V2G cis-pQTLs against 1190 all other cis-pQTLs.

- 1191 (a) Pie charts showing proportion of eQTL enrichment in the ARIC, deCODE, Fenland, and 1192 UKB-PPP European ancestry cohorts for strict V2G cis-pQTLs compared to all other cis-1193 pQTLs. Red: presence of a *cis*-eQTL; Yellow: absence of a *cis*-eQTL.
 - (b) Absolute value of effect of the minor allele on protein level broken down by the presence or absence of cis-eQTLs in the ARIC, deCODE, Fenland, and UKB-PPP European

1196ancestry cohorts. Strict V2G, strict variant-to-gene *cis*-pQTLs; All other, all other *cis*-
pQTLs that were removed due to strict V2G filtering. Boxplots show the median, lower,
and upper quartiles; whiskers end at 1.5 times the interquartile range from the top and
bottom of the box; points outside the whisker boundaries are plotted individually; smaller
black dots represent individual points and are used to show the number of samples
included in each boxplot; significance level *P* value is based on the Mann-Whitney U test.



1203 1204

1204 Extended Data Fig. 3. Absolute distance from transcription start site versus effect size for 1205 strict V2G *cis*-pQTLs and all other *cis*-pQTLs.

1206 The x-axis shows the distance of the pQTL from the canonical transcription start site of the 1207 associated protein-coding gene while the y-axis shows the absolute value of the effect size 1208 estimate of the effect allele aligned to the minor allele of each ancestry's respective reference 1209 panel. (Note, in European, African, and East Asian ancestry proteomics cohorts, the effect allele

of *cis*-pQTLs in each cohort was aligned to the minor allele of the corresponding variant in their respective reference panels—UKB 50k for European, HGDP+1kGP for African, and 1kGP for East Asian ancestry—to harmonize alleles across each ancestral cohort for plotting). Strict V2G *cis*-pQTLs are highlighted in blue while all other *cis*-pQTLs are highlighted in orange. Note that Fenland European *cis*-pQTLs are presented in GRCh37 coordinates while all other cohorts are presented in GRCh38 coordinates. *P* values show a one-sided t-test testing whether strict V2G *cis*-pQTLs have smaller absolute distance to the TSS compared to all other *cis*-pQTLs.

1216 *cis*-pQTLs have smaller absolute distance to the TSS compared to all other *cis*-p 1217

1219



n = 2,110



1223 (a) European cohorts involving ARIC, deCODE, Fenland, and UKB-PPP (4 cohorts).

(b) African cohorts involving ARIC and UKB-PPP (2 cohorts).

(c) East Asian ancestry cohort from Kyoto University Nagahama (single cohort).

1225 1226

1224

1220 1221



1227 1228 1229

Extended Data Fig. 5. Putatively causal protein-phenotype associations across European,
 African, and East Asian ancestries.

1230 Miami plots displaying chromosomal position (x-axis) of significant putatively causal protein-1231 outcome associations (MR-passing and colocalized with $PP_{max} \ge 0.8$) in (a) European, (b) African, 1232 (c) East Asian ancestry. The y-axis shows P values from the MR causal estimates where the 1233 exposure is protein level and outcome is the complex trait or disease. Colors indicate the type of 1234 complex trait or disease. Cancer types were harmonized under a single "Cancer" group. Ancestry 1235 is denoted by filled circle (European), filled diamond (African), and filled square (East Asian). Each 1236 data point is plotted based on the chromosome and transcription start-site of the protein-coding 1237 gene. For simplicity, Z scores and P values are averaged across cohorts in the European ancestry 1238 plot and the African ancestry plot and shown as a single data point. In European ancestry, only 1239 associations that were consistent across all cohorts are shown. 1240



1241 1242

1243 Extended Data Fig. 6. Uniquely instrumentable protein-phenotype pairs in African and East 1244 Asian ancestries.

- Red arrows indicate a positive causal estimate of the protein on the outcome while blue arrowsindicate a negative causal estimate of the protein on the outcome.
 - (a) Protein-phenotype pairs from 4 proteins uniquely instrumentable in African ancestry. Significant estimates between proteins (orange circles) and traits (green rectangles).
 - (b) Protein-phenotype pairs from 8 proteins uniquely instrumentable in East Asian ancestry. Significant estimates between proteins (orange circles) and traits (green rectangles).
- 1250 1251

1247

1248



C Olink Explore 3072 European vs. African ancestry



1252 Extended Data Fig. 7. Cross-ancestry comparison stratified by proteomics platform of 1254 instrumentable proteins overlapping at least one drug database (druggable genome, 1255 DrugBank, or Open Targets Platform).

(a) Comparison of SomaScan v4 platform instrumentable proteins overlapping at least one drug
 database between three European cohorts (ARIC, deCODE, and Fenland) and one African cohort
 (ARIC).

(b) Comparison of SomaScan v4 platform instrumentable proteins overlapping at least one drug
 database between three European cohorts (ARIC, deCODE, and Fenland) and one East Asian
 cohort (Kyoto University Nagahama).

- 1262 (c) Comparison of Olink Explore 3072 platform instrumentable proteins overlapping at least one
- 1263 drug database between one European cohort (UKB-PPP) and one African cohort (UKB-PPP).
- 1264

1265 **References**

- 1266 1. Sun, B. B. *et al.* Genomic atlas of the human plasma proteome. *Nature* **558**, 73– 1267 79 (2018).
- 1268 2. Emilsson, V. *et al.* Co-regulatory networks of human serum proteins link genetics
 1269 to disease. *Science* 361, 769–773 (2018).
- 1270 3. Suhre, K. *et al.* Connecting genetic risk to disease end points through the human 1271 blood plasma proteome. *Nat Commun* **8**, 14357 (2017).
- 1272 4. Williams, S. A. *et al.* Plasma protein patterns as comprehensive indicators of 1273 health. *Nat Med* **25**, 1851–1857 (2019).
- 1274 5. Su, C.-Y. *et al.* Circulating proteins to predict COVID-19 severity. *Sci Rep* **13**, 6236 (2023).
- 1276 6. Carrasco-Zanini, J. *et al.* Proteomic signatures improve risk prediction for common and rare diseases. *Nat Med* **30**, 2489–2498 (2024).
- 1278 7. Hopkins, A. L. & Groom, C. R. The druggable genome. *Nat Rev Drug Discov* 1, 1279 727–730 (2002).
- 1280 8. Overington, J. P., Al-Lazikani, B. & Hopkins, A. L. How many drug targets are 1281 there? *Nat Rev Drug Discov* **5**, 993–996 (2006).
- 1282 9. Bakheet, T. M. & Doig, A. J. Properties and identification of human protein drug targets. *Bioinformatics* 25, 451–457 (2009).
- 1284 10. Santos, R. *et al.* A comprehensive map of molecular drug targets. *Nat Rev Drug* 1285 *Discov* **16**, 19–34 (2017).
- 1286 11. Pietzner, M. *et al.* Mapping the proteo-genomic convergence of human diseases.
 1287 Science **374**, eabj1541 (2021).
- 1288 12. Ferkingstad, E. *et al.* Large-scale integration of the plasma proteome with 1289 genetics and disease. *Nat Genet* **53**, 1712–1721 (2021).
- 1290 13. Zhang, J. *et al.* Plasma proteome analyses in individuals of European and
 1291 African ancestry identify cis-pQTLs and models for proteome-wide association
 1292 studies. *Nat Genet* 54, 593–602 (2022).
- 1293 14. Sun, B. B. *et al.* Plasma proteomic associations with genetics and health in the UK Biobank. *Nature* **622**, 329–338 (2023).
- 1295 15. Skrivankova, V. W. *et al.* Strengthening the reporting of observational studies in
 epidemiology using mendelian randomisation (STROBE-MR): explanation and
 elaboration. *BMJ* **375**, n2233 (2021).
- Skrivankova, V. W. *et al.* Strengthening the Reporting of Observational Studies in
 Epidemiology Using Mendelian Randomization: The STROBE-MR Statement. *JAMA* **326**, 1614–1621 (2021).
- 1301 17. Chong, M. *et al.* Novel Drug Targets for Ischemic Stroke Identified Through
 1302 Mendelian Randomization Analysis of the Blood Proteome. *Circulation* **140**, 819–830
 1303 (2019).
- 1304 18. Zheng, J. *et al.* Phenome-wide Mendelian randomization mapping the influence 1305 of the plasma proteome on complex diseases. *Nat Genet* **52**, 1122–1131 (2020).
- 1306
 19. Zhao, H. *et al.* Proteome-wide Mendelian randomization in global biobank meta1307 analysis reveals multi-ancestry drug targets for common diseases. *Cell Genomics* 2,
 1308 100195 (2022).
- 1309 20. Zhou, S. *et al.* A Neanderthal OAS1 isoform protects individuals of European
 1310 ancestry against COVID-19 susceptibility and severity. *Nat Med* 27, 659–667 (2021).

- 1311 21. Yoshiji, S. *et al.* Proteome-wide Mendelian randomization implicates
 1312 nephronectin as an actionable mediator of the effect of obesity on COVID-19 severity.
 1313 Nat Metab 5, 248–264 (2023).
 1214 22 Yoshiji S. *et al.* COL 6A3 derived endetrophin mediates the effect of obesity on
- 1314 22. Yoshiji, S. *et al.* COL6A3-derived endotrophin mediates the effect of obesity on
 1315 coronary artery disease: an integrative proteogenomics analysis.
 1216 2022 04 10 22288706 Droprint at https://doi.org/10.1101/2022.04.10.22288706
- 13162023.04.19.23288706 Preprint at https://doi.org/10.1101/2023.04.19.232887061317(2023).
- 1318 23. Lu, T., Forgetta, V., Greenwood, C. M. T., Zhou, S. & Richards, J. B. Circulating
 1319 Proteins Influencing Psychiatric Disease: A Mendelian Randomization Study.
 1320 Biological Psychiatry 93, 82–91 (2023).
- 1321 24. Katz, D. H. *et al.* Whole Genome Sequence Analysis of the Plasma Proteome in
 1322 Black Adults Provides Novel Insights Into Cardiovascular Disease. *Circulation* 145,
 1323 357–370 (2022).
- 1324
 1325. Wang, B. *et al.* Comparative studies of genetic and phenotypic associations for
 2,168 plasma proteins measured by two affinity-based platforms in 4,000 Chinese
 1326 adults. 2023.12.01.23299236 Preprint at
- 1327 https://doi.org/10.1101/2023.12.01.23299236 (2023).
- Said, S. *et al.* Ancestry diversity in the genetic determinants of the human plasma
 proteome and associated new drug targets. 2023.11.13.23298365 Preprint at
 https://doi.org/10.1101/2023.11.13.23298365 (2023).
- 1331 27. Sakaue, S. *et al.* A cross-population atlas of genetic associations for 220 human 1332 phenotypes. *Nat Genet* **53**, 1415–1424 (2021).
- 1333 28. Feng, Y.-C. A. *et al.* Taiwan Biobank: A rich biomedical research database of the 1334 Taiwanese population. *Cell Genomics* **2**, 100197 (2022).
- 1335 29. Carlson, C. S. *et al.* Generalization and Dilution of Association Results from
 1336 European GWAS in Populations of Non-European Ancestry: The PAGE Study. *PLOS*1337 *Biology* **11**, e1001661 (2013).
- 1338 30. Martin, A. R. *et al.* Human Demographic History Impacts Genetic Risk Prediction
 1339 across Diverse Populations. *The American Journal of Human Genetics* **100**, 635–649
 1340 (2017).
- 1341 31. Cohen, J. *et al.* Low LDL cholesterol in individuals of African descent resulting 1342 from frequent nonsense mutations in PCSK9. *Nat Genet* **37**, 161–165 (2005).
- 1343 32. Robinson, J. G. *et al.* Efficacy and Safety of Alirocumab in Reducing Lipids and
 1344 Cardiovascular Events. *New England Journal of Medicine* **372**, 1489–1499 (2015).
- 1345 33. Sabatine, M. S. *et al.* Evolocumab and Clinical Outcomes in Patients with
 1346 Cardiovascular Disease. *New England Journal of Medicine* **376**, 1713–1722 (2017).
- 1347 34. Fatumo, S. *et al.* A roadmap to increase diversity in genomic studies. *Nat Med* 1348 **28**, 243–250 (2022).
- 1349 35. Ghoussaini, M. *et al.* Open Targets Genetics: systematic identification of trait1350 associated genes using large-scale genetics and functional genomics. *Nucleic Acids*1351 *Research* 49, D1311–D1320 (2021).
- 1352
 1353
 1353
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354
 1354</l
- 1355 37. Zhang, W. *et al.* SharePro: an accurate and efficient genetic colocalization 1356 method accounting for multiple causal signals. *Bioinformatics* **40**, btae295 (2024).

- 1357 38. Finan, C. *et al.* The druggable genome and support for target identification and
 1358 validation in drug development. *Science Translational Medicine* 9, eaag1166 (2017).
- 1359 39. Wishart, D. S. *et al.* DrugBank: a comprehensive resource for in silico drug
 1360 discovery and exploration. *Nucleic Acids Research* 34, D668–D672 (2006).
- 40. Koscielny, G. *et al.* Open Targets: a platform for therapeutic target identification
 and validation. *Nucleic Acids Research* 45, D985–D994 (2017).
- 1363 41. Schmidt, A. F. *et al.* Genetic drug target validation using Mendelian
 1364 randomisation. *Nat Commun* **11**, 3255 (2020).
- 1365 42. Gkatzionis, A., Burgess, S. & Newcombe, P. J. Statistical methods for cis1366 Mendelian randomization with two-sample summary-level data. *Genetic Epidemiology*1367 47, 3–25 (2023).
- 1368
 43. Lawlor, D. A., Harbord, R. M., Sterne, J. A. C., Timpson, N. & Davey Smith, G.
 1369
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 <li
- 1371 44. Karczewski, K. J. *et al.* The mutational constraint spectrum quantified from variation in 141,456 humans. *Nature* **581**, 434–443 (2020).
- 1373 45. Benjamini, Y. & Hochberg, Y. Controlling the False Discovery Rate: A Practical and Powerful Approach to Multiple Testing. *Journal of the Royal Statistical Society.* 1375 *Series B (Methodological)* 57, 289–300 (1995).
- 1376
 1376
 1377
 1377
 1378
 1378
 1378
 1378
 1379
 1379
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370
 1370</l
- 1379 47. van der Graaf, A. *et al.* Mendelian randomization while jointly modeling cis
 1380 genetics identifies causal relationships between gene expression and lipids. *Nat*1381 *Commun* **11**, 4930 (2020).
- 48. Goruppi, S. *et al.* The ULK3 Kinase Is Critical for Convergent Control of CancerAssociated Fibroblast Activation by CSL and GLI. *Cell Reports* 20, 2468–2479
 (2017).
- Hodonsky, C. J. *et al.* Multi-ancestry genetic analysis of gene regulation in
 coronary arteries prioritizes disease risk loci. *Cell Genomics* 4, 100465 (2024).
- 1387 50. Verstockt, B. *et al.* IL-12 and IL-23 pathway inhibition in inflammatory bowel 1388 disease. *Nat Rev Gastroenterol Hepatol* **20**, 433–446 (2023).
- 1389 51. O'Meara, E. *et al.* Independent Prognostic Value of Serum Soluble ST2
 1390 Measurements in Patients With Heart Failure and a Reduced Ejection Fraction in the
 1391 PARADIGM-HF Trial (Prospective Comparison of ARNI With ACEI to Determine
 1392 Impact on Global Mortality and Morbidity in Heart Failure). *Circulation: Heart Failure*1393 **11**, e004446 (2018).
- Jiang, Y. *et al.* An IL1RL1 genetic variant lowers soluble ST2 levels and the risk
 effects of APOE-ε4 in female patients with Alzheimer's disease. *Nat Aging* 2, 616–634 (2022).
- 1397 53. England, E. *et al.* Tozorakimab (MEDI3506): an anti-IL-33 antibody that inhibits
 1398 IL-33 signalling via ST2 and RAGE/EGFR to reduce inflammation and epithelial
 1399 dysfunction. *Sci Rep* 13, 9825 (2023).
- 1400 54. Cohen, S. B. The use of anakinra, an interleukin-1 receptor antagonist, in the
 1401 treatment of rheumatoid arthritis. *Rheumatic Disease Clinics of North America* **30**,
 1402 365–380 (2004).

- Aviram, M. *et al.* Paraoxonase inhibits high-density lipoprotein oxidation and
 preserves its functions. A possible peroxidative role for paraoxonase. *J Clin Invest* **101**, 1581–1590 (1998).
- 1406 56. Reddy, S. T. *et al.* Human Paraoxonase-3 Is an HDL-Associated Enzyme With
 1407 Biological Activity Similar to Paraoxonase-1 Protein but Is Not Regulated by Oxidized
 1408 Lipids. *Arteriosclerosis, Thrombosis, and Vascular Biology* 21, 542–547 (2001).
- 1409 57. Verma, A. *et al.* Diversity and scale: Genetic architecture of 2068 traits in the VA 1410 Million Veteran Program. *Science* **385**, eadj1182 (2024).
- 1411 58. Banfi, C. *et al.* Prenylcysteine oxidase 1, an emerging player in atherosclerosis.
 1412 *Commun Biol* 4, 1–17 (2021).
- 1413 59. Rocnik, E. F., Liu, P., Sato, K., Walsh, K. & Vaziri, C. The Novel SPARC Family
 1414 Member SMOC-2 Potentiates Angiogenic Growth Factor Activity*. *Journal of*1415 *Biological Chemistry* 281, 22855–22864 (2006).
- 1416 60. Bourgault, J. *et al.* Proteome-Wide Mendelian Randomization Identifies Causal
 1417 Links Between Blood Proteins and Acute Pancreatitis. *Gastroenterology* 164, 9531418 965.e3 (2023).
- 1419
 61. Lafferty, M. J., Bradford, K. C., Erie, D. A. & Neher, S. B. Angiopoietin-like
 1420
 1421
 1421
 1421
 FORMATION* *This work was supported by a grant from the Pew Foundation (to S.
- 1422 B. N.). Journal of Biological Chemistry **288**, 28524–28534 (2013).
- 1423 62. Wang, H. & Eckel, R. H. Lipoprotein lipase: from gene to obesity. *American*1424 *Journal of Physiology-Endocrinology and Metabolism* 297, E271–E288 (2009).
- 1425 63. Rosenson Robert S. *et al.* Zodasiran, an RNAi Therapeutic Targeting ANGPTL3,
 1426 for Mixed Hyperlipidemia. *New England Journal of Medicine* **0**,.
- Suhre, K., McCarthy, M. I. & Schwenk, J. M. Genetics meets proteomics:
 perspectives for large population-based studies. *Nat Rev Genet* 22, 19–37 (2021).
- 1429 65. Liu, Z. *et al.* Genetic architecture of the inflammatory bowel diseases across East
 1430 Asian and European ancestries. *Nat Genet* 55, 796–806 (2023).
- 1431 66. Kong, L. *et al.* The landscape of immune dysregulation in Crohn's disease
 1432 revealed through single-cell transcriptomic profiling in the ileum and colon. *Immunity*1433 56, 444-458.e5 (2023).
- 1434 67. Massimino, L. *et al.* The Inflammatory Bowel Disease Transcriptome and
 1435 Metatranscriptome Meta-Analysis (IBD TaMMA) framework. *Nat Comput Sci* 1, 511–
 1436 515 (2021).
- 1437 68. Hamilton, M. J., Frei, S. M. & Stevens, R. L. The Multifaceted Mast Cell in 1438 Inflammatory Bowel Disease. *Inflammatory Bowel Diseases* **20**, 2364–2378 (2014).
- 1439 69. Boeckxstaens, G. Mast cells and inflammatory bowel disease. *Current Opinion in* 1440 *Pharmacology* **25**, 45–49 (2015).
- 1441 70. Bycroft, C. *et al.* The UK Biobank resource with deep phenotyping and genomic data. *Nature* **562**, 203–209 (2018).
- 1443 71. Koenig, Z. *et al.* A harmonized public resource of deeply sequenced diverse
 1444 human genomes. 2023.01.23.525248 Preprint at
- 1445 https://doi.org/10.1101/2023.01.23.525248 (2024).
- 1446 72. Purcell, S. *et al.* PLINK: A Tool Set for Whole-Genome Association and
 1447 Population-Based Linkage Analyses. *The American Journal of Human Genetics* 81,
 1448 550 575 (2007)
- 1448 559–575 (2007).

- 1449 73. Durinck, S. *et al.* BioMart and Bioconductor: a powerful link between biological
 1450 databases and microarray data analysis. *Bioinformatics* 21, 3439–3440 (2005).
- 1451 74. Butler-Laporte, G. *et al.* HLA allele-calling using multi-ancestry whole-exome
 1452 sequencing from the UK Biobank identifies 129 novel associations in 11 autoimmune
 1453 diseases. *Commun Biol* 6, 1–17 (2023).
- Holmes, M. V., Richardson, T. G., Ference, B. A., Davies, N. M. & Davey Smith,
 G. Integrating genomics with biomarkers and therapeutic targets to invigorate
 cardiovascular drug development. *Nat Rev Cardiol* 18, 435–453 (2021).
- 1457 76. McLaren, W. *et al.* The Ensembl Variant Effect Predictor. *Genome Biology* **17**, 1458 122 (2016).
- 1459 77. Hemani, G. *et al.* The MR-Base platform supports systematic causal inference 1460 across the human phenome. *eLife* **7**, e34408 (2018).
- 1461 78. Pierce, B. L., Ahsan, H. & VanderWeele, T. J. Power and instrument strength
 1462 requirements for Mendelian randomization studies using multiple genetic variants.
 1463 *International Journal of Epidemiology* **40**, 740–752 (2011).
- Hemani, G., Tilling, K. & Smith, G. D. Orienting the causal relationship between
 imprecisely measured traits using GWAS summary data. *PLOS Genetics* 13,
 e1007081 (2017).
- 1467
 1467
 80. Zhang, W. *et al.* SharePro: an accurate and efficient genetic colocalization method accounting for multiple causal signals. 2023.07.24.550431 Preprint at https://doi.org/10.1101/2023.07.24.550431 (2023).
- 1470 81. Giambartolomei, C. *et al.* Bayesian Test for Colocalisation between Pairs of
 1471 Genetic Association Studies Using Summary Statistics. *PLOS Genetics* 10,
 1472 e1004383 (2014).
- 1473 82. Conway, J. R., Lex, A. & Gehlenborg, N. UpSetR: an R package for the
 1474 visualization of intersecting sets and their properties. *Bioinformatics* 33, 2938–2940
 1475 (2017).
- 1476 83. Pruim, R. J. *et al.* LocusZoom: regional visualization of genome-wide association 1477 scan results. *Bioinformatics* **26**, 2336–2337 (2010).
- 1478