

# Proteomic prediction of disease largely reflects environmental risk exposure

Kristin Tsuo<sup>1-3</sup>, M. Austin Argentieri<sup>1-3</sup>, Danni Gadd<sup>4,5</sup>, Mitja Kurki<sup>1-3,6</sup>, Zhili Zheng<sup>1-3,6</sup>, Denis Baird<sup>7</sup>, Riccardo E. Marioni<sup>4,5</sup>, Christopher Foley<sup>4,8</sup>, Hailiang Huang<sup>1-3</sup>, Benjamin B. Sun<sup>7,9</sup>, Chia-Yen Chen<sup>7</sup>, Mark J. Daly<sup>1-3,6,10</sup>, Alicia R. Martin<sup>1-3,10</sup>

1. Analytic and Translational Genetics Unit, Massachusetts General Hospital, Boston, MA, USA
2. Stanley Center for Psychiatric Research, Broad Institute of MIT and Harvard, Cambridge, MA, USA
3. Program in Medical and Population Genetics, Broad Institute of MIT and Harvard, Cambridge, MA, USA
4. Optima Partners, Edinburgh, EH2 4HQ, UK.
5. Centre for Genomic and Experimental Medicine, Institute of Genetics and Cancer, University of Edinburgh, Edinburgh, EH4 2XU, UK.
6. Institute for Molecular Medicine Finland, University of Helsinki, Helsinki, Finland
7. Biogen Inc., Cambridge, MA, USA.
8. Bayes Centre, The University of Edinburgh, Edinburgh, EH8 9BT, UK.
9. Cardiovascular Epidemiology Unit, Department of Public Health and Primary Care, University of Cambridge, Cambridge, CB1 8RN, UK.
10. These authors jointly supervised this work.

\*Corresponding authors: [ktsuo@broadinstitute.org](mailto:ktsuo@broadinstitute.org), [mjdaly@broadinstitute.org](mailto:mjdaly@broadinstitute.org), [armartin@broadinstitute.org](mailto:armartin@broadinstitute.org)

## 23 Abstract

24 Plasma proteomic signatures accurately predict disease risk, but our understanding of the mechanisms  
25 contributing to the predictive value of the proteome remains limited. Here, we characterized proteomic  
26 biomarkers of 19 age-related diseases, based on observational associations between 2,923 protein levels and  
27 incidence of these outcomes in the UK Biobank (N = 45,438). To identify the subset of these biomarkers that  
28 may represent causal drivers of disease, we first employed Mendelian Randomization (MR) and found that only  
29 8% of the protein-disease associations with genetic instruments showed suggestive evidence of causal  
30 relationships, and were more likely to pertain to only a single disease. We then tested the hypothesis that many  
31 proteomic biomarkers, particularly the non-causal proteins, are impacted by environmental factors that might  
32 independently affect disease risk and protein levels. We discovered that the vast majority (>90%) of proteins  
33 associated with diseases like lung cancer and COPD are also associated with smoking, and more than half of  
34 all disease-associated proteins tested in MR were associated with smoking. These proteins showed no evidence  
35 of causal effects on disease, suggesting their predictive value is as an environmental sensor. Given the sensitivity  
36 of the plasma proteome to smoking, we developed a proteomic score for smoking (SmokingPS) and  
37 demonstrated that the plasma proteome can serve as a quantitative index of smoking behavior and history.  
38 Extending this approach to alcohol intake phenotypes, our results generally suggest that many plasma proteins  
39 identified in observational associations are more likely to be readouts of environmental risk factors than disease-  
40 specific signals. We conclude that the plasma proteome may provide critical objective biomarkers for quantifying  
41 the impacts of environmental risk factors on human health and disease. Our results have significant implications  
42 for implementing predictive plasma protein biomarkers in disease prevention, and can help guide interpretation  
43 of putative protein-disease associations as actionable therapeutic targets or quantitative indications of upstream  
44 exposures that represent potential intervention points.

## 45 Introduction

46 Technological advances in high-throughput, broad-capture proteomic assays have highlighted the potential of  
47 leveraging plasma proteins for disease prediction, biomarker discovery, and therapeutic development. In the  
48 past few years, population-scale biobanks have used these assays to profile portions of the plasma proteome in  
49 thousands of individuals<sup>1-5</sup>. The UK Biobank Pharma Proteomics Project (UKB-PPP) generated one of the  
50 largest proteomic resources to date, with measurements of nearly 3,000 blood plasma analytes in more than  
51 54,000 UKB participants using the Olink Explore 3072 proximity extension assay<sup>6</sup>. Coupled with existing deep  
52 phenotyping data, UKB-PPP creates a unique opportunity to interrogate the links between the plasma proteome  
53 and a broad spectrum of diseases and disease-related traits. Indeed, studies have already leveraged this  
54 resource to identify thousands of associations between plasma protein levels and disease outcomes<sup>7-9</sup>.  
55 Proteomic predictors of incident disease, developed from these associations, outperform traditional clinical  
56 models and polygenic risk scores for a wide range of diseases<sup>7,8,10</sup> (e.g. AUC > 0.8 for 92 diseases<sup>9</sup>).

57

58 The ability of the plasma proteome to predict incident disease is promising for a range of clinical applications –  
59 for example building on ELISA (enzyme-linked immunosorbent assay) and other multiplex immunoassays for  
60 diagnostics and monitoring treatments<sup>11</sup>, as well as for patient stratification in clinical trials<sup>12</sup> – but our  
61 understanding of what contributes to the predictive value of the circulating proteome remains limited. A few  
62 overlapping mechanisms could explain the widespread associations between plasma proteins and incident  
63 disease outcomes: (1) plasma proteins have direct causal effects on disease onset; (2) plasma proteins are  
64 impacted by disease processes that begin before disease diagnosis; and (3) plasma proteins are impacted by  
65 external factors, like environmental exposures, that both contribute to disease onset and affect protein levels.  
66 Gaining a deeper understanding of the specific roles of proteomic biomarkers in disease prediction has important  
67 implications for the potential clinical utility of these biomarkers.

68

69 First, plasma proteins with causal effects on disease may be candidates for novel therapeutic targets. Several  
70 disease-specific studies have investigated the potential causal roles of plasma proteins, including for type II  
71 diabetes<sup>13</sup>, cardiometabolic diseases<sup>10,14,15</sup>, psychiatric disorders<sup>16,17</sup>, and gastrointestinal diseases<sup>18,19</sup>.  
72 Phenome-wide studies across many diseases and phenotypes have also investigated the causal basis of  
73 protein-disease associations<sup>20,21</sup>. These studies have identified putative drug targets such as SCARA5 for  
74 cardioembolic stroke<sup>15</sup> and IL1RL1 for inflammatory bowel diseases<sup>20</sup>.

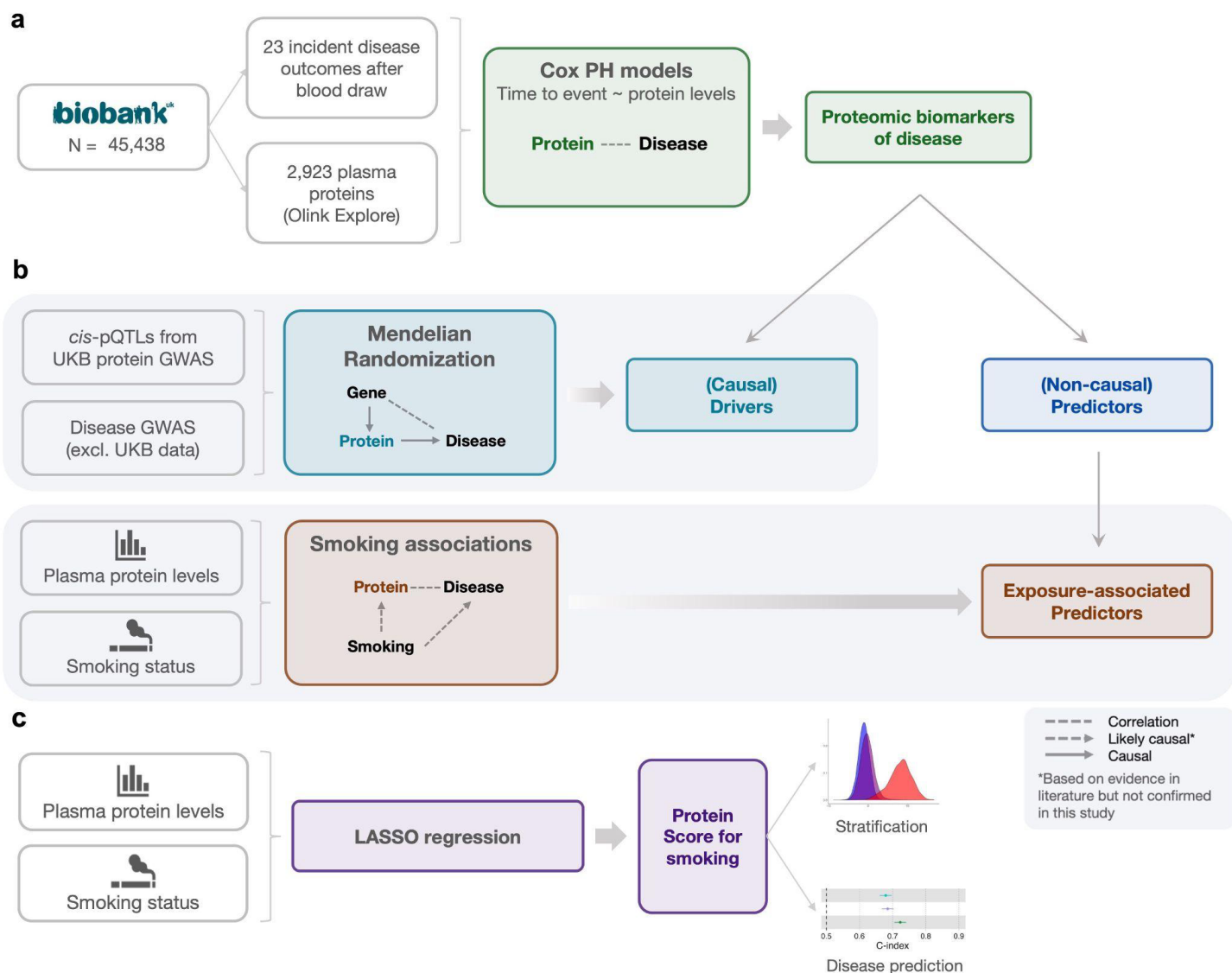
75

76 On the other hand, non-causal plasma proteins may serve as biomarkers for monitoring disease onset and  
77 progression. For example, a study on plasma proteomic biomarkers of dementia found that abundances of  
78 certain predictive proteins began changing years before diagnosis<sup>22</sup>. Non-causal plasma proteins may also

79 potentially help monitor impacts of changes in lifestyle and other environmental factors on disease risk. Previous  
80 work has speculated that the plasma proteome captures environmental effects<sup>7,23</sup>, and a handful of studies  
81 investigating factors contributing to the variance of individual proteins have found that environmental exposures  
82 can explain a substantial portion of variability in the levels of certain proteins<sup>24,25</sup>. However, major gaps remain  
83 in understanding how modifiable factors contribute to the utility of the overall proteome as disease predictors. In  
84 particular, few studies to date have simultaneously evaluated putatively causal proteins alongside other  
85 mechanisms that may explain the predictive value of non-causal associations. In this study, we aimed to conduct  
86 a comprehensive characterization of the relationships between plasma proteins and key disease outcomes in  
87 order to delineate proteins as potential therapeutic targets or non-causal biomarkers of environmental risk  
88 factors.

89  
90 Utilizing blood proteomic data available from a subset of the UKB-PPP cohort (N = 45,438, see **Methods**), we  
91 investigated an expanded set of associations between 2,923 unique proteins and 23 age-related incident disease  
92 outcomes (see study design in **Figure 1**). We first applied Mendelian Randomization (MR) to partition the  
93 disease-associated proteomic biomarkers into those with potential causal roles across diseases (i.e. drivers) vs.  
94 non-causal roles (i.e. predictors). Next, we examined the impacts of smoking on the plasma proteome because  
95 of its large environmental effect on disease risk. Integrating the MR and smoking analyses, we identified that a  
96 large proportion of protein predictors are exposure-associated and likely play no causal role in disease  
97 pathogenesis. To further evaluate the sensitivity of the plasma proteome to smoking, we developed a proteomic  
98 score for smoking, and tested this score alongside smoking status and other clinical biomarkers in disease  
99 prediction models. By replicating this proteomic score in an external dataset (FinnGen), we highlight the  
100 objective, quantitative, and generalizable nature of environmental risk quantification enabled by proteomics.  
101 Finally, we confirmed the generality of this paradigm by developing and testing an additional proteomic score for  
102 alcohol intake, which is nearly entirely independent from the smoking protein score.

103



104

105

106

107

108

109

110

111

112

113

114

115

116

117

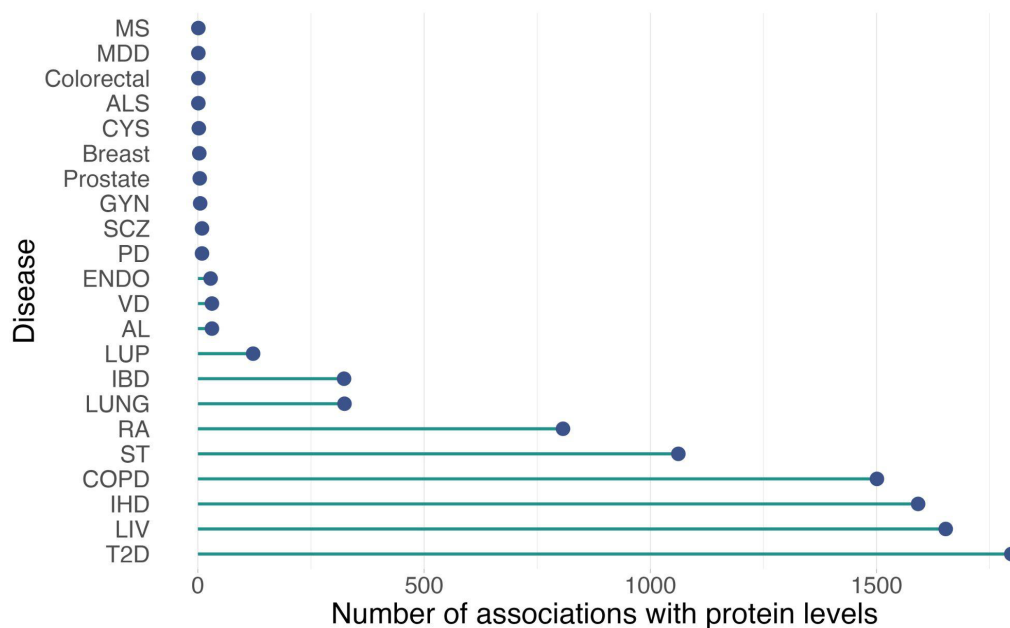
**Figure 1. Study design for characterizing proteomic biomarkers of disease.** **a**, We tested for associations between 2,923 plasma proteins measured in a subset of UKB-PPP participants and 23 incident disease outcomes using Cox proportional hazards (PH) models. We refer to the proteins in the significant protein-disease associations as biomarkers. **b**, To better understand the protein-disease associations identified in **(a)**, we took two approaches. First, we identified which of the protein-disease associations were likely causal by applying two-sample Mendelian Randomization (MR) using protein quantitative trait loci (pQTLs) as genetic instruments and disease GWAS that did not utilize UKB data; we refer to the putatively causal proteins as drivers. Second, we identified proteins that are likely not causal themselves but instead reflect the effects of disease-related exposures, by identifying proteins that were significantly associated with smoking but lacking evidence from MR; we refer to these proteins as exposure-associated predictors. **c**, We trained a LASSO regression model on the subset of UKB-PPP participants with smoking status data to develop a protein score for smoking (SmokingPS). We demonstrated that the SmokingPS accurately captures quantity and frequency of smoking and used SmokingPS to predict disease incidence.

## 118 Results

### 119 Plasma proteins are widely associated with incident disease

120 Previously, we tested associations between levels of 1,468 plasma proteins, measured in UKB-PPP participants,  
121 and 23 incident diseases over a 15-year follow-up period using Cox proportional hazards (PH) models. In this  
122 study, we expanded these analyses to include associations with an additional 1,455 proteins (total proteins =  
123 2,923), measured in the same individuals. Details on the plasma protein measurements, QC procedures, and  
124 UKB-PPP sample have been described previously<sup>6</sup> (see also **Methods**). We found 9,308 significant associations  
125 between 2,122 proteins and 22 of the 23 incident disease outcomes, adjusting for age and sex variables  
126 (Bonferroni-adjusted  $P$ -value  $< 0.05/(23 \times 2923) = 7.44 \times 10^{-7}$ ) (**Supplementary Table 1**). The number of  
127 associations ranged from 1 association each for ALS, multiple sclerosis, colorectal cancer, and major depressive  
128 disorder to 1,653 and 1,798 for liver disease and type 2 diabetes, respectively; brain/CNS cancer showed no  
129 significant protein associations (**Extended Data Fig. 1**).

130



131

132 **Extended Data Fig. 1. Number of associations between plasma proteins and 22 incident disease**  
133 **outcomes.** 9,308 significant associations involving 2,122 proteins and 22 incident disease outcomes are shown  
134 (Bonferroni-adjusted  $P$ -value  $< 7.44 \times 10^{-7}$ ). (MS = multiple sclerosis; MDD = major depressive disorder;  
135 Colorectal = colorectal cancer; ALS = amyotrophic lateral sclerosis; GYN = gynecological cancers; CYS = cystitis;  
136 Breast = breast cancer; Prostate = prostate cancer; PD = Parkinson's disease; SCZ = schizophrenia; ENDO =  
137 endometriosis; AL = Alzheimer's dementia; VD = vascular dementia; LUP = systemic lupus erythematosus;  
138 LUNG = lung cancer; IBD = inflammatory bowel disease; RA = rheumatoid arthritis; ST = ischemic stroke; COPD  
139 = chronic obstructive pulmonary disease; IHD = ischemic heart disease; LIV = Liver disease; T2D = Type 2  
140 diabetes)

## 141 A small subset of proteins have causal effects on disease

142 To identify which subset of the thousands of protein-disease associations discovered represent potentially causal  
143 relationships, we applied two-sample Mendelian Randomization (MR). 19 of the 23 disease outcomes that we  
144 tested in the Cox PH models had suitable GWAS summary statistics for MR, as they did not include UKB data  
145 which was used to identify genetic instruments (**Supplementary Table 2**). These diseases are the focal set for  
146 subsequent analyses. We leveraged the catalog of protein quantitative trait loci (pQTLs) that were previously  
147 mapped to these proteins<sup>6</sup> to first subset to proteins with robust genetic instruments. To minimize potential  
148 biases, we restricted genetic instruments to *cis*-pQTLs, and retained the proteins with *F*-statistics > 10 and that  
149 passed Steiger filtering. Ultimately, 2,907 of the significant protein-disease pairs (involving 782 unique proteins)  
150 had genetic instruments and were tested for causality via MR (**Supplementary Table 3**).

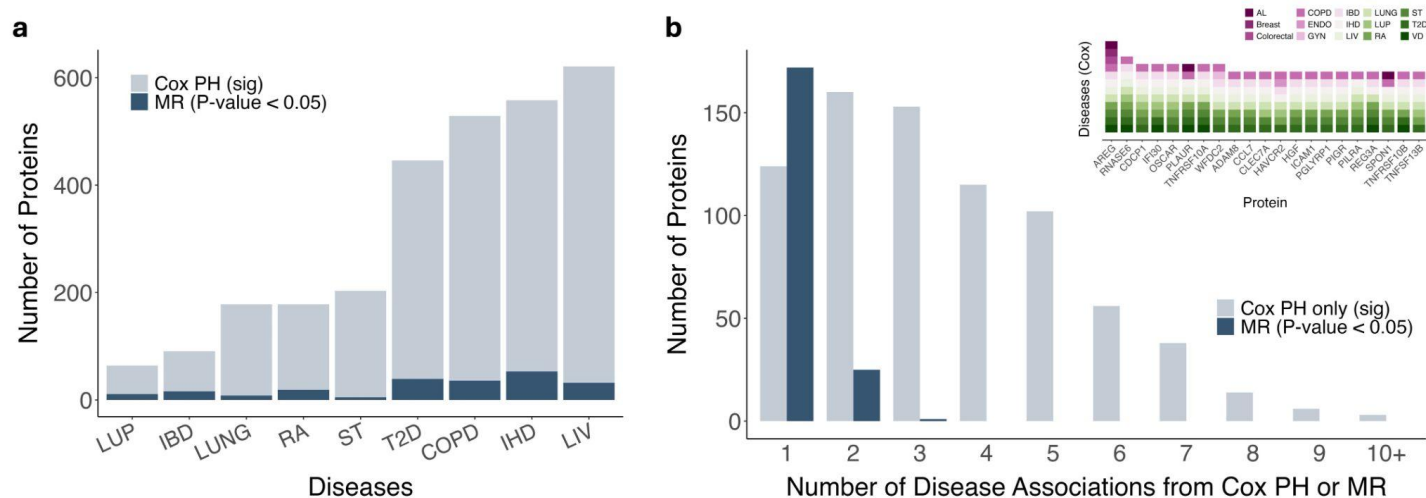
151  
152 Only 225/2,907 protein-disease pairs (8%) showed suggestive evidence for causality (*P*-value < 0.05) (**Figure**  
153 **2a**). In additional sensitivity tests for horizontal pleiotropy using MR-Egger<sup>26</sup>, we found that 9 of the 225 nominally  
154 significant protein-disease pairs may have pleiotropic genetic instruments (**Supplementary Table 4**), and thus  
155 further validation of these pairs may be needed. Overall, these *cis*-MR analyses did not reveal causal roles for  
156 the majority of the proteins identified in observational associations, indicating that the vast majority of protein-  
157 disease associations likely cannot be explained by causal effects of the protein on disease onset.

158  
159 Among the 225 protein-disease pairs, we replicated known causal relationships such as proprotein convertase  
160 subtilisin/kexin type 9 (PCSK9) and ischaemic heart disease<sup>27</sup>, and apolipoprotein E (APOE) and Alzheimer's  
161 disease<sup>28</sup>. Other protein-disease pairs with significant causal associations in this study and prior evidence for  
162 causality included tumor necrosis factor receptor superfamily member 6B (TNFRSF6B) and inflammatory bowel  
163 disease (IBD)<sup>29</sup>, peptidylglycine alpha-amidating monooxygenase (PAM) and type 2 diabetes<sup>30</sup>, and transforming  
164 growth factor beta 1 (TGFB1) and ischaemic heart disease<sup>31</sup>. Sortilin1 (SORT1) had robust and significant  
165 evidence for causal effects on type 2 diabetes in this study and its potential causal role has not been previously  
166 identified.

167  
168 We assessed the disease specificity of causal drivers identified in MR compared to the non-causal, disease-  
169 associated predictors. We found that hundreds of non-causal predictor proteins were associated with multiple  
170 diseases in the Cox PH models; several were associated with 9 or more diseases (e.g., amphiregulin (AREG),  
171 ribonuclease A family member k6 (RNASE6), and CUB domain-containing protein 1 (CDCP1)). The diseases  
172 associated with these highly non-specific proteins are involved in a broad range of mechanisms and pathways  
173 (**Figure 2b**). In contrast, most of the causal driver proteins (172/198, 87%) were associated with 1 disease. On

174 average, this group of proteins was associated with 1.1 diseases, while the disease-associated proteins lacking  
 175 MR support were associated with an average of 3.5 diseases.

176



177

178 **Figure 2. Comparison of proteins involved in observational vs. causal associations across diseases. a,**  
 179 Counts of proteins significantly associated with diseases in Cox PH models and that also had genetic instruments  
 180 (N = 782 proteins involved in 2,907 protein-disease pairs). Diseases with more than 20 significant protein  
 181 associations are shown. Counts of the subset of these proteins showing suggestive evidence of causality from  
 182 MR analyses (P-value < 0.05) are indicated by darker blue bars. **b,** Counts of proteins associated with 1 to more  
 183 than 10 diseases in the Cox PH models vs. MR tests. All proteins shown here are involved in the 2,907 protein-  
 184 disease pairs that were significant in Cox PH models and had genetic instruments. Light blue bars represent  
 185 non-causal predictors; dark blue bars represent causal drivers. Inset shows 27 proteins significantly associated  
 186 with 8 or more diseases in the Cox PH models (IBD = inflammatory bowel disease; LUP = systemic lupus  
 187 erythematosus; RA = rheumatoid arthritis; ST = ischemic stroke; LUNG = lung cancer; COPD = chronic  
 188 obstructive pulmonary disease; T2D = Type 2 diabetes; LIV = Liver disease; IHD = ischemic heart disease,  
 189 Breast = breast cancer; Colorectal = colorectal cancer; VD = vascular dementia).

190 **Many protein-disease associations are driven by smoking**

191 Given the widespread sharing of non-causal predictor proteins across diverse diseases, we next tested the  
 192 hypothesis that these proteins broadly associate with incident disease due to their associations with  
 193 environmental risk factors, and thus reflect environmental impacts without directly causing disease. Previously  
 194 developed proteomic scores for the incident disease outcomes showed better or similar prediction performance  
 195 compared to models comprised of age, sex, and other lifestyle factors, including BMI, alcohol intake, social  
 196 deprivation, educational attainment, physical activity, and smoking status<sup>7</sup>, pointing to a potential role of the  
 197 plasma proteome in partially capturing the effects of these risk factors. For lung cancer and COPD specifically,  
 198 we observed that incremental models with disease proteomic scores had minimal improvement beyond smoking  
 199 status (**Extended Data Fig. 2**). Previous literature, including studies using small, targeted microarrays, suggests  
 200 that smoking may alter plasma protein levels<sup>32,33</sup>. Thus, we further investigated the relationship between the

201 plasma proteome and smoking, a well-known environmental risk factor with large effects on many common  
202 complex diseases<sup>34</sup>.

203  
204 To first gauge the overall effects of smoking on the plasma proteome, we characterized the relationships between  
205 self-reported smoking status and the 2,923 proteins measured in the UKB-PPP cohort, adjusting for age and sex  
206 variables. 1,673 proteins (57%) were significantly associated with smoking status (Bonferroni-adjusted  $P$ -value  
207  $< 1.7 \times 10^{-5}$ ) (**Supplementary Table 5**). 30 proteins measured in the Olink panel mapped to 21 plasma proteins  
208 previously found to be associated with smoking<sup>32</sup> (**Supplementary Table 6**), and all were significantly associated  
209 with smoking in these analyses.

210  
211 We then evaluated the effects of smoking on the disease-associated proteins specifically. All 178 proteins  
212 associated with lung cancer in the Cox PH models were significantly associated with smoking, as were between  
213 82-97% of proteins associated with all other diseases examined (**Figure 3a**), a notable enrichment beyond the  
214 proportion of all measured plasma proteins associated with smoking. Additionally, we found that all proteins  
215 associated with 8 or more diseases in the Cox PH models were smoking-associated, compared to smaller  
216 proportions of the more disease-specific proteins (**Figure 3b**).

217  
218 After adjusting for self-reported smoking status in the Cox PH models, 8,326 out of 9,308 (89%) protein-disease  
219 pairs remained significant (**Supplementary Table 7**). Lung cancer showed the most substantial drop in the  
220 number of associated proteins after adjustment (from 324 to 33 proteins, representing a 90% decrease); COPD  
221 also showed a large decrease (from 1,501 to 1,250 proteins, representing a 17% decrease). Of note, most of  
222 the remaining lung cancer-associated proteins (28/33, 85%) had significantly attenuated hazard ratios after  
223 smoking adjustment (Wald test,  $p$ -value  $< 0.05$ ), with a mean attenuation of 48.3% of the  $|\log(HR)|$ . Taken  
224 together, these results indicate that although adjusting for smoking status yields large attenuations in the number  
225 of associated proteins for some diseases, basic measures of smoking may fail to capture the full effects of  
226 smoking in protein-disease associations.

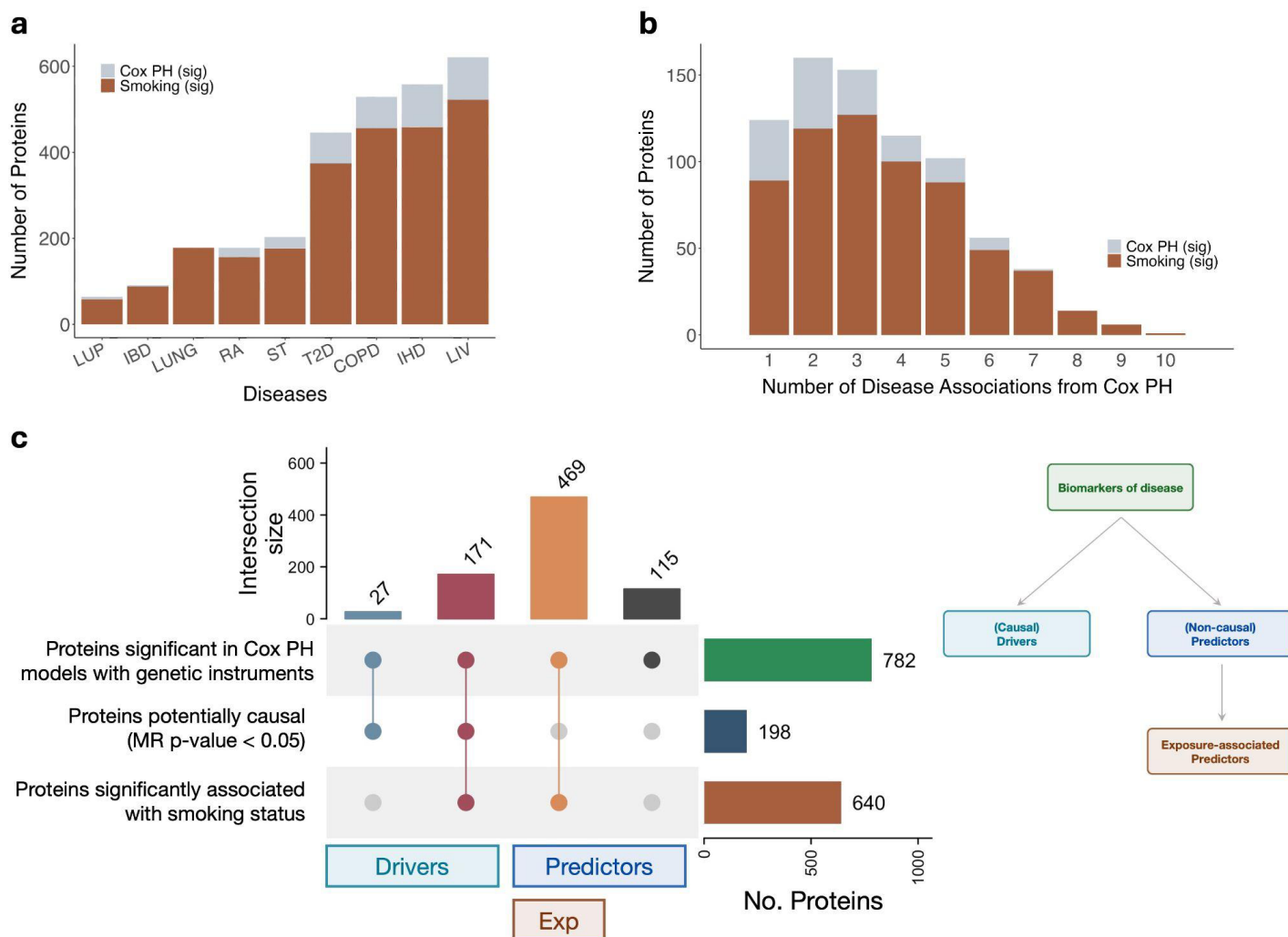
227  
228 Given the clear biological impacts of smoking, we then sought to systematically distinguish causal driver proteins  
229 from exposure-associated predictor proteins, focusing on the 782 disease-associated proteins with valid genetic  
230 instruments. We found that the vast majority of these proteins ( $N=640$ ) were significantly associated with smoking  
231 status. Moreover, most of these smoking-associated proteins ( $N=469$ ) showed no evidence for causal effects on  
232 disease from the *cis*-MR analyses (**Figure 3c**). Several proteins ( $N=171$ ) were both associated with smoking  
233 and nominated by MR, indicating that these proteins may represent (or be correlated with) causal mediators  
234 between smoking and incident disease.

235

236 Another group of disease-associated proteins were neither associated with smoking nor showed evidence for  
237 causality (**Figure 3c**). The majority of proteins in this group (86/115, 75%) were associated with 1 to 3 diseases  
238 in the Cox PH models, potentially representing a group of proteins enriched for non-causal, disease-specific  
239 predictors tagging incident disease processes. For example, 5 proteins in this group were associated with only  
240 COPD, including a transmembrane activin inhibitor (BAMBI) that was previously found to be involved in  
241 Th17/Treg pathway imbalances in COPD patients<sup>35</sup>.

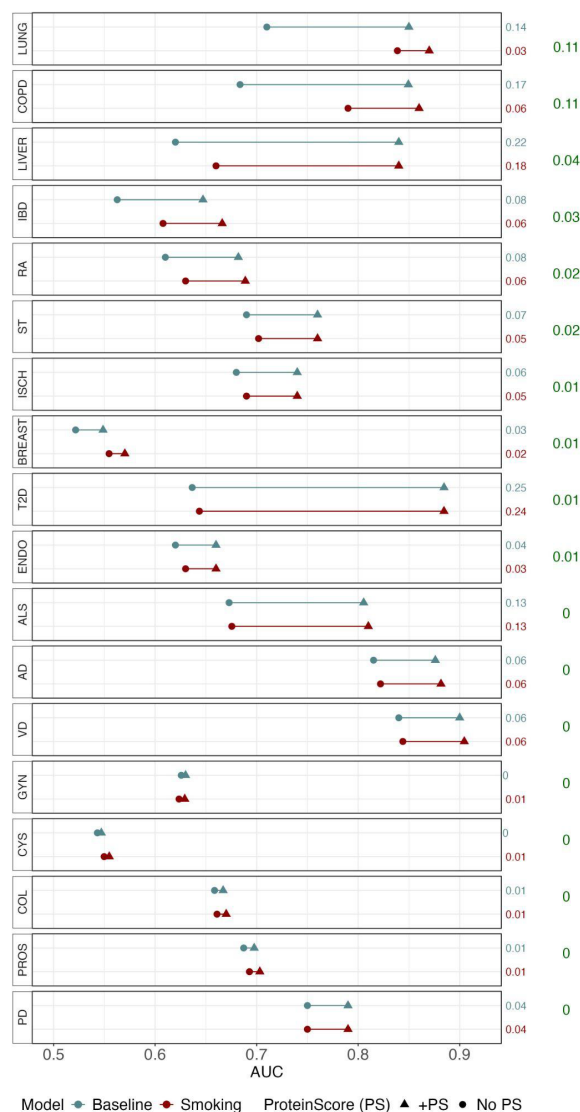
242

243 Based on these groupings, we characterized the proteins selected in our previously developed proteomic scores  
244 for COPD and lung cancer<sup>7</sup>. 47 proteins in the COPD score had at least one genetic instrument, and all except  
245 2 (CLSTN2 and LYPD8) were significantly associated with smoking; 10 proteins additionally showed evidence  
246 for causal effects. In the lung cancer proteomic score, 5 of the 6 proteins had at least one genetic instrument,  
247 and all 5 were significantly associated with smoking; carcinoembryonic antigen (CEACAM5) additionally showed  
248 suggestive evidence for causality. These disease proteomic scores appear to mostly capture exposure-  
249 associated predictor proteins and not causal drivers, which aligns with our previous finding that they have  
250 diminished predictive utility when smoking status is included in the prediction models (**Extended Data Fig. 2**).  
251 Additionally, the proteomic score for lung cancer had limited prediction performance over the baseline model of  
252 age and sex in individuals that reported never smoking compared to previous and current smokers ( $\Delta$ AUC  
253 between baseline plus lung cancer ProteinScore and baseline only = 0.02, 0.06, and 0.05 in never, previous,  
254 and current smokers, respectively).



255  
256

**Figure 3. Smoking associations categorize many biomarkers as exposure-associated predictors.** **a**, Counts of proteins significantly associated with diseases in Cox PH models, with counts of the subset of these proteins significantly associated with smoking status (across the entire UKB dataset) highlighted by orange bars. Diseases with more than 20 significant protein associations are shown (IBD = inflammatory bowel disease; LUP = systemic lupus erythematosus; RA = rheumatoid arthritis; ST = ischemic stroke; LUNG = lung cancer; COPD = chronic obstructive pulmonary disease; Diab = Type 2 diabetes; LIV = Liver disease; IHD = ischemic heart disease). **b**, Counts of proteins significantly associated with 1 to 10 diseases in the Cox PH models, with counts of the subset of proteins in each category significantly associated with smoking status highlighted by orange bars (Bonferroni-adjusted  $P$ -value <  $1.7 \times 10^{-5}$ ). **c**, UpSet plot showing the classification of 782 proteins significantly associated with a disease (Bonferroni-adjusted  $P$ -value <  $0.05/(19 \text{ diseases} \times 2,923 \text{ proteins}) = 9.00 \times 10^{-7}$  from Cox PH associations) for which a genetic instrument existed. Significant association with smoking status was determined based on Bonferroni-adjusted  $P$ -value <  $0.05/2,923 = 1.7 \times 10^{-5}$ . Groups delineated as causal drivers, non-causal predictors, and exposure-associated predictors are indicated by bars on the bottom, and illustrated in the schematic on the right.



271  
272  
273  
274  
275  
276  
277  
278

**Extended Data Fig. 2. Performance of disease-based proteomic scores from Gadd et al.<sup>7</sup>** Differences in AUC between models with standard covariates and models with the addition of the disease ProteinScores (PS). Disease ProteinScores were developed and described in Gadd et al.<sup>7</sup> In blue (baseline), models without PS consist of age and sex in which diseases were not sex-stratified. In red (smoking), models without PS consist of age, sex, and self-reported smoking status. To the right of the plot, the first column of numbers shows the differences in AUC with the addition of the PS per model; the second column shows  $\Delta$ +PS in baseline minus  $\Delta$ +PS in smoking.

279  
280  
281  
282  
283

### Proteomic score quantifies smoking behavior and history

The widespread effects of smoking on the plasma proteome indicate that plasma proteins themselves may serve as precise, quantitative measures of the cumulative biological effects of smoking and other exposures. To assess the utility of the plasma proteome as a quantitative readout of smoking, we developed a smoking protein score (SmokingPS). Since many of the measured proteins are likely correlated, we performed variable selection via

LASSO on the protein levels of current and never smokers in the UKB (N=14,585) (**Extended Data Fig. 3**), identifying 550 proteins that we combined into the SmokingPS. This score predicted smoking status with high accuracy (AUC = 0.955 [95% CI = 0.949-0.961]) in a hold-out set of UKB participants (N=14,586). We replicated the score in another biobank, FinnGen<sup>36</sup>, with plasma protein measurements from the Olink Explore panel in 1,990 participants, 1,863 of whom had smoking status information. The SmokingPS demonstrated good but lower ability to discriminate current vs. never smoking status in FinnGen participants (AUC = 0.844 [95% CI = 0.814-0.875]), and the SmokingPS distributions showed some separation between current/previous smokers and never smokers (**Extended Data Fig. 4a**). We note that because smoking status information was collected at recruitment several years before proteomic sampling in FinnGen, it was not possible to distinguish between current and previous smokers for this analysis.

We compared the distributions of the SmokingPS across individuals in UKB-PPP who reported current, previous, and never smoking status. As expected, the previous smokers had a SmokingPS distribution between the current and never smokers (**Figure 4a**). We then stratified the previous smokers by years since smoking cessation and pack years, which captures both quantity and frequency of smoking. As years since smoking cessation increased, the distributions of the SmokingPS in previous smokers shifted towards the never smokers (**Figure 4b**). As illustrated by the SmokingPS distributions across individuals who stopped smoking recently, the smoking-associated proteome starts to revert back to non-smoking levels within a couple years of smoking cessation; this pattern is also observed from the levels of the top-weighted proteins in the SmokingPS, averaged across bins of years since smoking cessation in former smokers (**Extended Data Fig. 5**). Similarly, those who smoked fewer pack years had SmokingPS distributions closer to the never smokers (**Extended Data Fig. 4b**). Stratifying current smokers by pack years reflected the same trend – individuals who reported fewer pack years had SmokingPS distributions closer to, or overlapping, that of previous smokers (**Figure 4c**). The average number of cigarettes smoked per day among current smokers ranged from 8.6 in the lowest quantile of SmokingPS distribution to 20.1 in the highest quantile, underlining that the amount of smoking has a dose-response impact on plasma protein levels (**Extended Data Fig. 4c**).

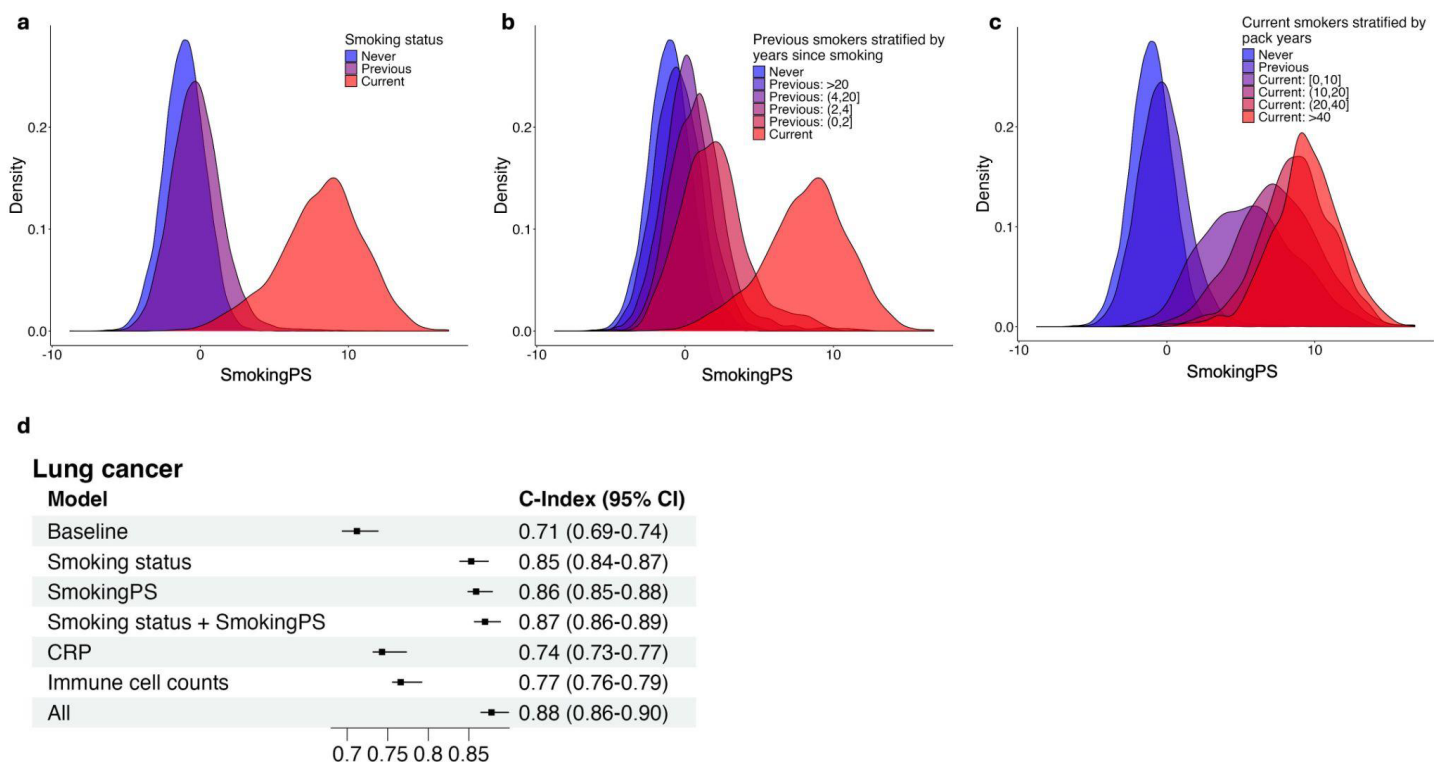
We further tested the value of this SmokingPS for disease prediction (**Figure 4d, Extended Data Fig. 6**). Among the 19 incident disease outcomes, adding smoking status to baseline variables of age, sex, and genetic principal components (PCs) resulted in the greatest improvement in performance for lung cancer, COPD, and liver disease (delta C-index between baseline model vs. baseline with smoking status: 0.141, 0.083, and 0.024, respectively). Adding the SmokingPS to the baseline model with smoking status further improved performance for these three diseases (delta C-index between baseline with smoking status vs. with smoking status and SmokingPS: 0.017, 0.011, and 0.008 for lung cancer, COPD, and liver disease), indicating that the plasma

318 proteome likely provides a more quantitative biological readout of smoking exposure than self-reported smoking  
319 status for some diseases.

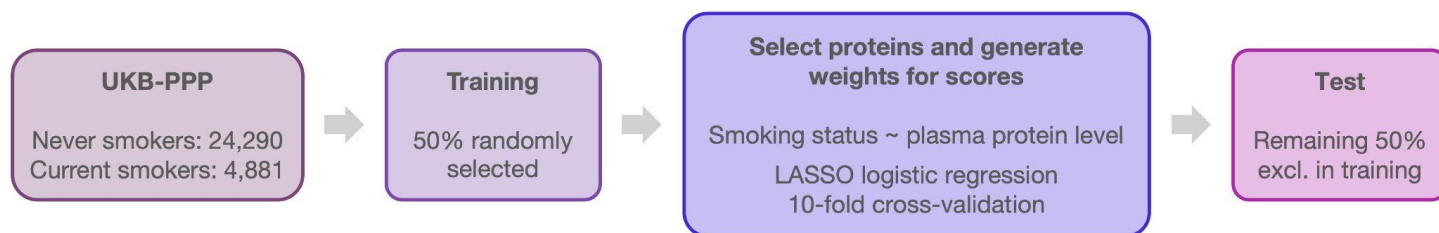
320  
321 Of the individuals with incident lung cancer diagnoses (N = 412), 58 were self-reported never smokers. The  
322 SmokingPS showed better prediction of lung cancer incidence over the baseline model in current smokers (delta  
323 C-index = 0.044, p-value = 0.069) compared to in the never smokers (delta C-index = 0.013, p-value = 0.71)  
324 (**Extended Data Fig. 7c**). Additionally, distributions of the SmokingPS between never smokers with and without  
325 lung cancer were not significantly different (two-sample Kolmogorov-Smirnov test  $D = 0.15$ , p-value = 0.11)  
326 (**Extended Data Fig. 7b**). We then considered the full set of non-smokers, and leveraged the breadth of  
327 exposure data collected in the UKB to investigate other factors that the SmokingPS may be associated with in  
328 never smokers. Across 71 environmental and lifestyle factors covering a range of variables (e.g. secondhand  
329 smoke exposure, physical and social activities, air pollution measures) (**Supplementary Table 8**), the  
330 SmokingPS was significantly associated with a few smoking-related factors, such as “attends pub/club”,  
331 “exposure to tobacco smoke at home”, and “number of smokers in household” (**Supplementary Table 9**).

332  
333 In additional sensitivity analyses, we compared the SmokingPS to commonly-used clinical biomarkers, C-  
334 reactive protein (CRP) and immune cell counts (**Figure 4d, Extended Data Fig. 6a**). For lung cancer and COPD,  
335 the SmokingPS outperformed both CRP and immune cell counts (delta C-index between baseline with  
336 SmokingPS vs. with CRP: 0.116 and 0.051 for lung cancer and COPD, respectively; delta C-index between  
337 baseline with SmokingPS vs. with immune cell counts: 0.093 and 0.015). Combining CRP and immune cell  
338 counts with both smoking status and the SmokingPS resulted in the best performing models for lung cancer,  
339 COPD, and liver disease, indicating that each measure may capture unique aspects of disease risk.

340

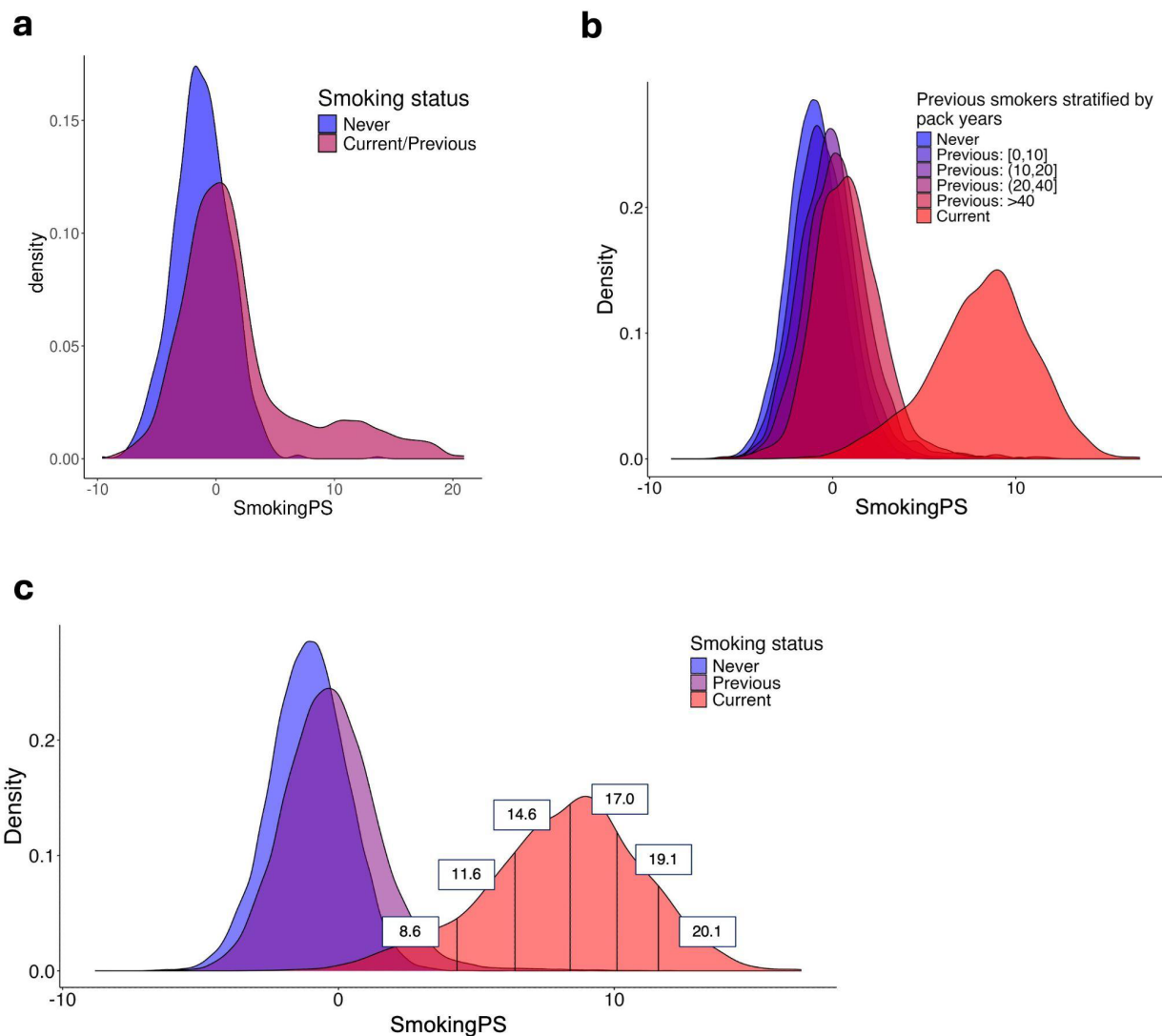


**Figure 4. SmokingPS captures different smoking measures and predicts incident disease.** **a**, Density plot of SmokingPS with individuals stratified by self-reported smoking status. **b**, Density plot of SmokingPS with self-reported previous smokers stratified by number of years since smoking cessation. **c**, Density plot of SmokingPS with self-reported current smokers stratified by pack years, calculated as number of cigarettes smoked per day, divided by twenty, multiplied by number of years smoking. **d**, Associations between incident lung cancer (cases = 405, controls = 46,968) and models with various predictors, shown on y-axis, using Cox PH. Baseline model includes age, sex, age<sup>2</sup>, age × sex, age<sup>2</sup> × sex, and first 10 genetic PCs. All models following baseline include these variables, as well as the predictor listed. “CRP” indicates C-reactive protein. “Immune cell counts” include neutrophil, eosinophil, basophil, monocyte, lymphocyte, and white blood cell counts. “All” indicates the baseline variables, smoking status, SmokingPS, CRP, and immune cell counts. C-index is shown on the x-axis and listed alongside 95% confidence intervals.



**Extended Data Fig. 3. Summary of SmokingPS development.** UKB-PPP participants who self-reported never and current smoking status were used for developing the SmokingPS. 50% were randomly assigned to the training group; the remaining 50% were assigned to the test group. Ratios of never to current smokers in each group were similar. LASSO regression with 10-fold cross validation was used to select proteins out of the 2,923 proteins measured and derive weighting coefficients for the selected proteins.

359



360

361

362

363

364

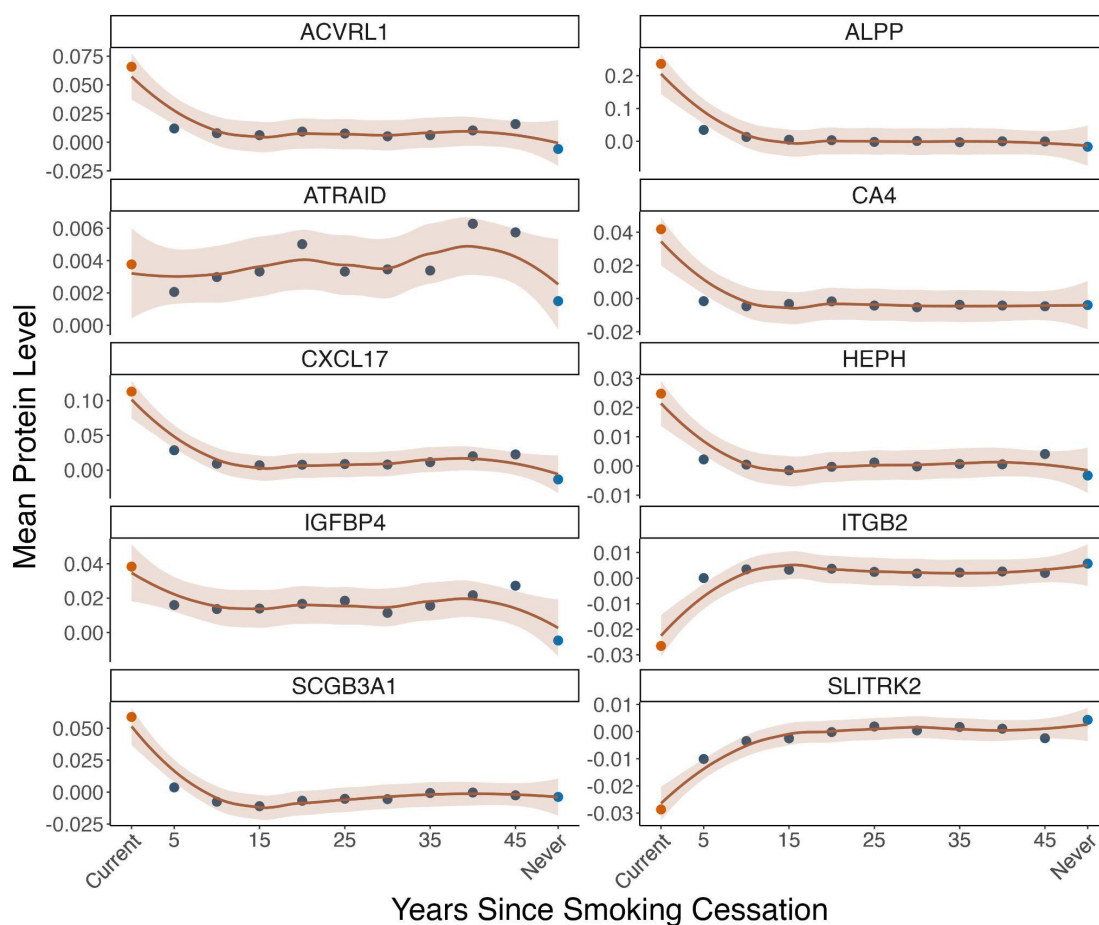
365

366

367

368

**Extended Data Fig. 4. SmokingPS distributions in FinnGen and across previous and current smokers in UKB.** **a**, Density plot of the SmokingPS in FinnGen participants (N = 1,862), stratified by self-reported never and current/previous smoking status. Current and previous smokers are grouped together due to gaps in timing between the collection of smoking information and proteomic sampling. **b**, Density plot of SmokingPS with self-reported previous smokers in UKB stratified by pack years, calculated as number of cigarettes smoked per day, divided by twenty, multiplied by number of years smoking. **c**, Density plot of SmokingPS with self-reported current smokers divided into bins of SmokingPS at 0.10, 0.25, 0.50, 0.75, and 0.90 quantiles. Cigarettes smoked per day were averaged across individuals in each bin; each average is reported on the density plot.



369

370

371

372

373

374

**Extended Data Fig. 5. Average protein levels of top-weighted proteins in SmokingPS across groups of previous smokers.** X-axis represents years since smoking cessation, grouped in 5-year intervals for former smokers: (0,5], (10,15], etc. Average protein levels in current smokers and never smokers are included for reference.

**a**

**COPD**

Model	C-Index (95% CI)
Baseline	0.70 (0.69-0.71)
Smoking status	0.78 (0.77-0.79)
SmokingPS	0.78 (0.77-0.79)
Smoking status + SmokingPS	0.79 (0.78-0.80)
CRP	0.73 (0.72-0.74)
Immune cell counts	0.76 (0.75-0.78)
All	0.82 (0.81-0.83)

**b**

**Liver disease**

Model	C-Index (95% CI)
Baseline	0.61 (0.59-0.63)
Smoking status	0.63 (0.62-0.66)
SmokingPS	0.64 (0.62-0.66)
Smoking status + SmokingPS	0.64 (0.63-0.66)
CRP	0.68 (0.66-0.70)
Immune cell counts	0.69 (0.67-0.70)
All	0.72 (0.71-0.74)

375

376

377

378

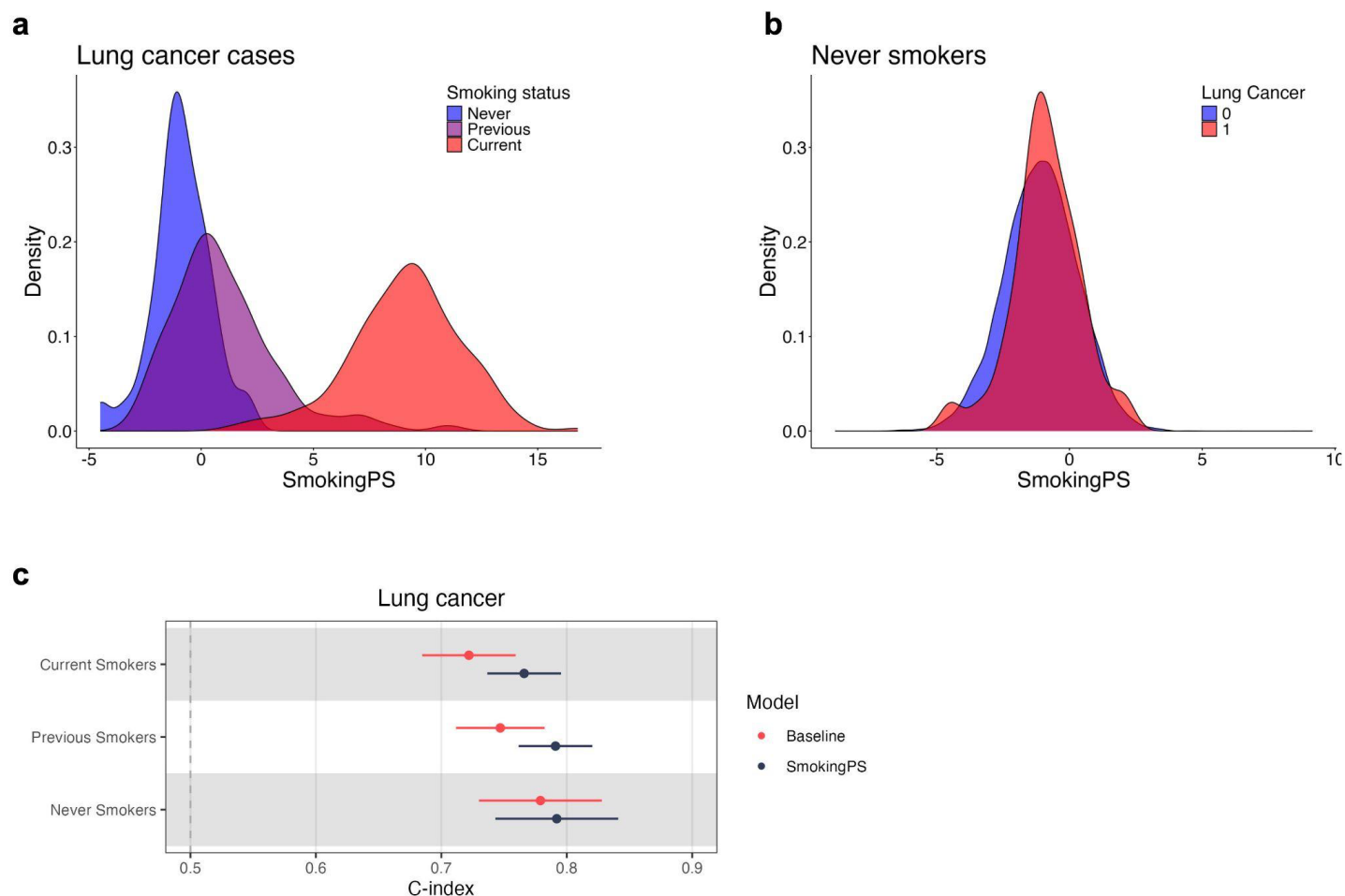
379

380

381

**Extended Data Fig. 6. Comparison of SmokingPS and other models for prediction of incident COPD and liver disease.** Associations between incident disease and models with various predictors, shown on y-axis, using Cox PH. Baseline model includes age, sex, age<sup>2</sup>, age × sex, age<sup>2</sup> × sex, and first 10 genetic PCs. All models following baseline include these variables, as well as the predictor listed. “CRP” indicates C-reactive protein. “Immune cell counts” include neutrophil, eosinophil, basophil, monocyte, lymphocyte, and white blood cell counts. “All” indicates the baseline variables, smoking status, SmokingPS, CRP, and immune cell counts.

382 C-index is shown on the x-axis and listed alongside 95% confidence intervals. **a**, Incident disease outcome is  
383 COPD (cases = 1,973, controls = 44,765). **b**, Incident disease outcome is liver disease (cases = 328, controls =  
384 46,913).  
385



386  
387 **Extended Data Fig. 7. SmokingPS in individuals with lung cancer stratified by smoking status.** **a**, Density  
388 plot of SmokingPS in individuals with lung cancer, stratified by self-reported smoking status. **b**, Density plot of  
389 SmokingPS in never smokers, stratified by incident lung cancer outcome. **c**, Cox PH associations between  
390 incident lung cancer and (1) baseline model of age, sex, age<sup>2</sup>, age × sex, age<sup>2</sup> × sex, and first 10 genetic PCs  
391 and (2) baseline model with SmokingPS. Cox PH associations were conducted in individuals stratified by self-  
392 reported smoking status.

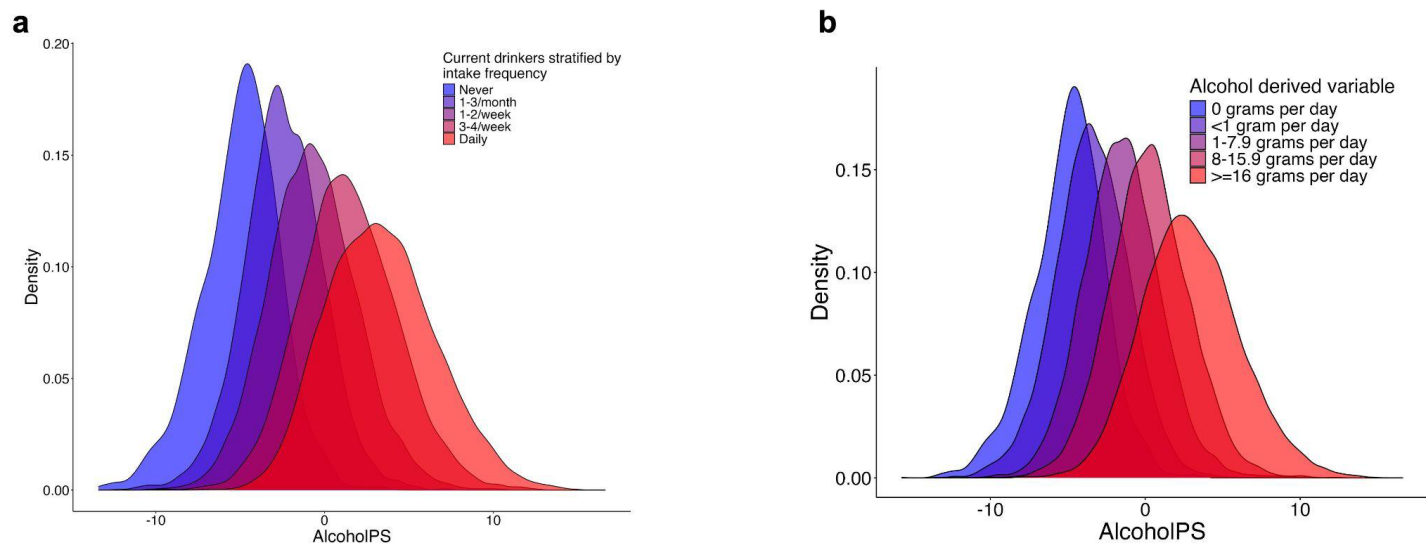
393 **Generalization of paradigm: proteomic score for alcohol use**

394 Having demonstrated that the proteome can effectively serve as a quantitative measure of disease-relevant  
395 environmental risk factors using smoking as an example, we sought to explore whether this approach was  
396 generalizable to other environmental exposures. We tested associations between the plasma proteome and  
397 alcohol intake status, and found that 1,228 (42%) of proteins were significantly associated with alcohol use  
398 (Bonferroni-adjusted P-value < 1.7 × 10<sup>-5</sup>) (**Supplementary Table 10**). We developed a proteomic score for

399 alcohol intake (AlcoholIPS), trained on the protein levels of UKB-PPP individuals who reported never drinking  
400 and daily drinking. The score consisted of 474 proteins selected by LASSO, and predicted current vs. never  
401 drinking status in a hold-out set of UKB participants with high accuracy (AUC = 0.975). Although 134 (28%)  
402 proteins overlapped with those in the SmokingPS, the AlcoholIPS distributions did not stratify according to  
403 smoking status and vice versa (two-sample Kolmogorov-Smirnov test for AlcoholIPS:  $D = 0.16$ , p-value <  $2.2e-$   
404  $16$  for never vs. previous smokers;  $D = 0.14$ , p-value <  $2.2e-16$  for never vs. current smokers;  $D = 0.05$ , p-value  
405 =  $3.96e-09$  for previous vs. current smokers; and SmokingPS:  $D = 0.17$ , p-value <  $2.2e-16$  for never vs. previous  
406 drinkers;  $D = 0.17$ , p-value <  $2.2e-16$  for never vs. current drinkers;  $D = 0.07$ , p-value =  $1.73e-08$  for previous  
407 vs. current drinkers) (**Extended Data Fig. 8a-b**), indicating that the two scores are largely independent.

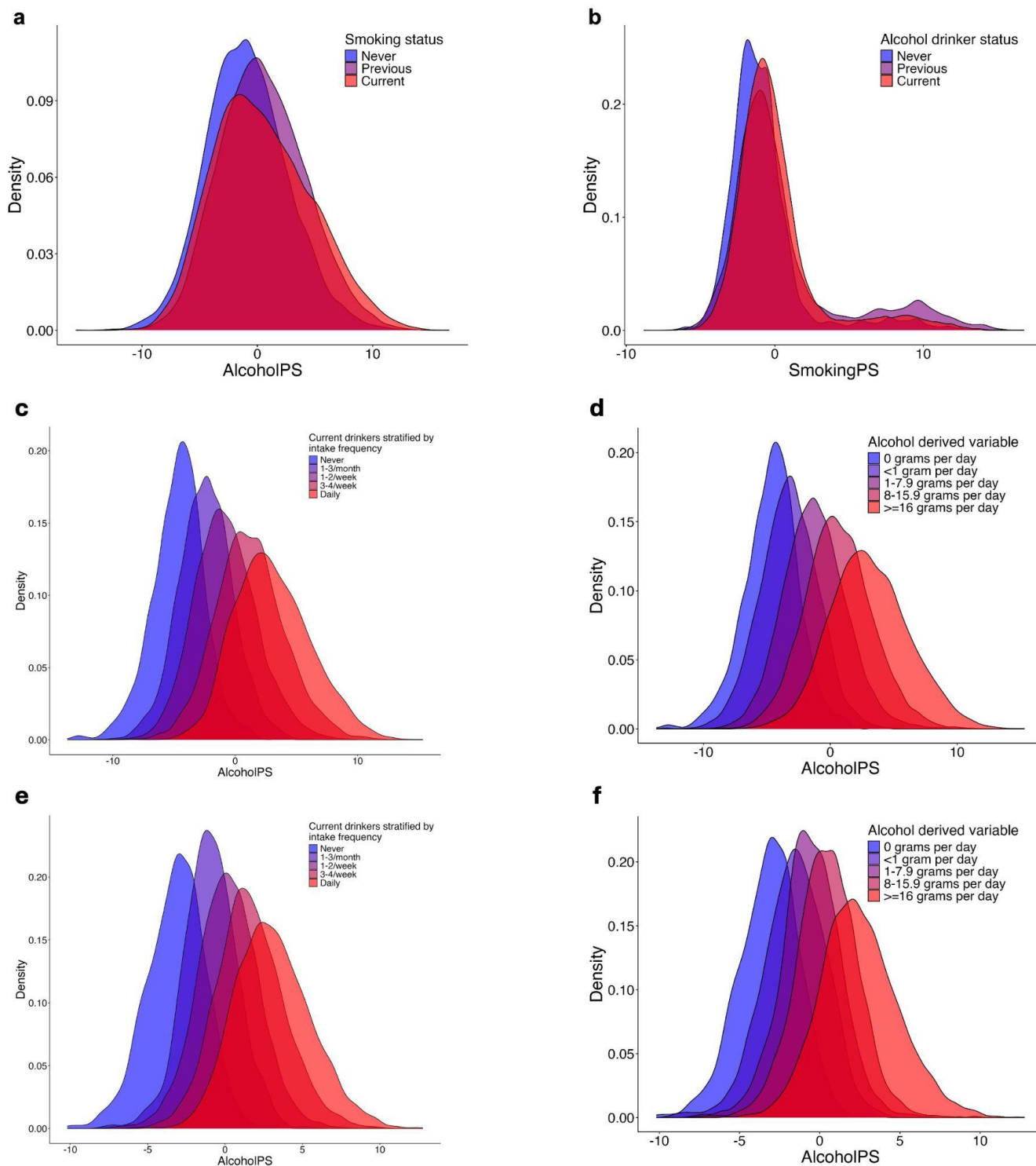
408  
409 Distributions of the AlcoholIPS in self-reported current drinkers stratified by number of drinks per week and month  
410 (**Fig. 5a**). Additionally, we stratified current drinkers by the number of grams of alcohol intake per day, derived  
411 from several alcohol intake variables in the UKB (**Methods**), and similarly observed gradual shifts in the protein  
412 score distribution as the alcohol intake quantities increased (**Fig. 5b**). Given differences in the ratio of drinkers  
413 vs. never-drinkers in males (653 never-drinkers; 5,233 current drinkers) and females (1,498 never-drinkers;  
414 3,891 current drinkers), we performed additional sex-stratified sensitivity analyses. AlcoholIPSs trained  
415 separately in males and females similarly stratified according to alcohol intake frequency and amount (**Extended**  
416 **Data Fig. 8c-f**), indicating that the AlcoholIPS is not merely detecting sex-specific proteins. The AlcoholIPS  
417 derived from matching sample sizes of male and female never and daily drinkers in the training cohort (**Methods**)  
418 had similar performance in the hold-out set (AUC = 0.969), and 293/408 (72%) of the proteins in this AlcoholIPS  
419 overlapped with the AlcoholIPS trained on all individuals without downsampling.

420  
421 We additionally evaluated the ability of the AlcoholIPS to predict clinical biomarkers commonly used to assess  
422 liver function, including alanine aminotransferase (ALT), aspartate aminotransferase (AST), gamma-glutamyl  
423 transferase (GGT), and bilirubin. The AlcoholIPS outperformed self-reported alcohol intake status in daily and  
424 never drinkers for AST (incremental  $R^2$ : 0.005 vs. 0.0009), GGT (0.027 vs. 0.004), and bilirubin (0.02 vs. 0.01),  
425 as well in current and never drinkers (**Extended Data Fig. 9**).



426  
427  
428  
429  
430

**Figure 5. AlcoholPS captures frequency and amount of alcohol intake in current drinkers.** **a**, Density plot of AlcoholPS in current drinkers stratified by self-reported number of drinks per week or month. Self-reported never drinkers shown for reference. **b**, Density plot of AlcoholPS in current drinkers stratified by derived grams of alcohol intake per day.



431

432

433

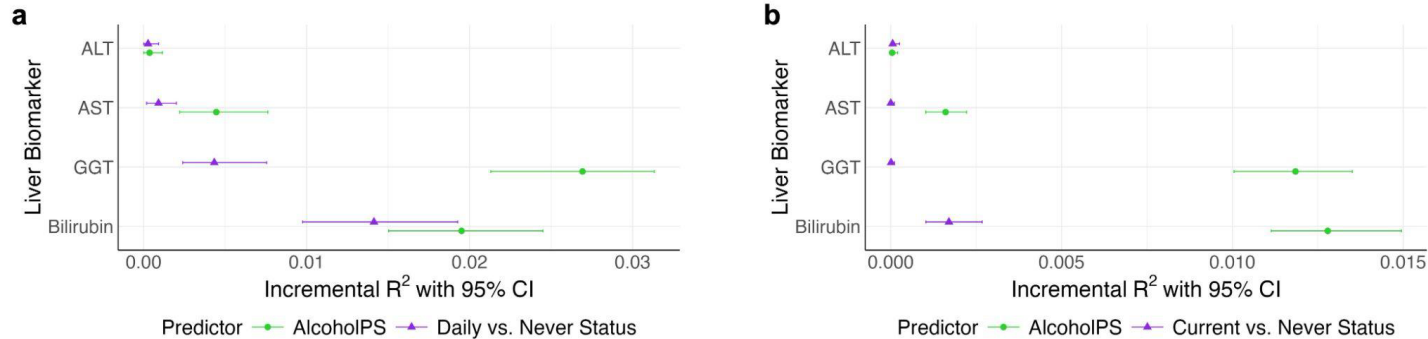
434

435

436

**Extended Data Fig. 8. Comparisons of AlcoholPS and SmokingPS, and AlcoholPS trained in females and males separately.** **a**, Density plot of AlcoholPS in individuals with self-reported smoking information, stratified by smoking status. **b**, Density plot of SmokingPS in individuals with self-reported alcohol intake information, stratified by alcohol drinker status. Density plot of AlcoholPS in current drinkers stratified by self-reported day of drinks per week or month in **c**, females and **e**, males. Self-reported never drinkers shown for reference. Density

437 plot of AlcoholIPS in current drinkers stratified by derived grams of alcohol intake per day in **d**, females and **f**,  
438 males.



439

440

441

442

443

444

445

446

447

**Extended Data Fig. 9. Associations between liver biomarkers and AlcoholIPS.** Comparisons of variance explained in liver biomarkers on *y*-axis by AlcoholIPS and **a**, daily vs. never drinking status or **b**, current vs. never drinking status. Incremental  $R^2$  was estimated as improvement in  $R^2$  with inclusion of either AlcoholIPS or alcohol intake status, comparing the two models: (1) baseline model (biomarker  $\sim$  age + sex + age<sup>2</sup> + age\*sex + age<sup>2</sup>\*sex) and (2) full model (biomarker  $\sim$  AlcoholIPS or alcohol intake status + age + sex + age<sup>2</sup> + age\*sex + age<sup>2</sup>\*sex. (ALT = alanine aminotransferase, AST = aspartate aminotransferase, GGT = gamma-glutamyl transferase)

448

## Discussion

449

450

451

452

453

454

455

456

457

458

459

Our study shows that a substantial proportion of plasma protein-disease associations identified in observational studies are the result of systemic effects of environmental risk factors on protein expression and do not reflect causal relationships between those proteins and disease. Instead, our results suggest that much of the predictive value of the disease-associated proteome stems from its association with environmental impacts, which have been historically challenging to evaluate consistently across studies using objective and quantitative approaches. This has significant implications for the clinical translation of disease risk estimates and biomarkers derived from plasma proteomics, as it indicates that in many cases the proteins may be capturing the effects of an upstream exposure driving changes in protein expression rather than directly contributing to disease onset. Our approach can be applied to numerous diseases and exposure types, helping delineate how plasma proteomic information can inform health interventions.

460

461

462

463

464

465

Associations between protein levels and disease incidence observed in epidemiological studies are susceptible to confounding and other biases. We leveraged MR to determine which subset of protein-disease associations may represent direct, causal effects of proteins on disease onset. Many prior studies have focused on identifying causal proteins for specific diseases<sup>10,13,14</sup>, whereas our study investigated 19 different diseases that span common metabolic diseases (e.g. type 2 diabetes), immune-mediated diseases (e.g. rheumatoid arthritis and IBD), and cancers to discern potential patterns or differences in the proportions of causal proteins identified

466 between traits. Across diseases, we found only 8% of protein-disease pairs tested had any suggestive evidence  
467 for a causal relationship. The small proportion of putatively causal protein drivers identified in this study is  
468 consistent with findings from disease-specific studies, as well as two recent phenome-wide MR studies: one  
469 study found that less than 1% of the tested protein-phenotype pairs across 355 phenotypes were putatively  
470 causal<sup>20</sup>, and another study, which included different diseases and filters, similarly found less than 1% of  
471 significant protein-disease pairs were putatively causal using *cis*-pQTLs<sup>9</sup>. Despite the low numbers of MR-  
472 nominated proteins, we confirm that this group of proteins includes more disease-specific signals, rather than  
473 broadly disease-associated proteins. Though MR-nominated proteins represent a small fraction of the proteome,  
474 their potential causal roles are critical to mechanistically characterize as they may constitute promising  
475 therapeutic targets.

476  
477 In the Cox PH models associating proteins with incident disease, we observed that many proteins were  
478 associated with multiple, often unrelated, diseases. One explanation is that environmental factors that impact  
479 shared pathways underlying many diseases also have widespread effects on the plasma proteome. In this study  
480 we explored this scenario and focused on smoking as an example of a known environmental risk factor for many  
481 diseases with persistent biological effects<sup>37</sup>. By integrating the MR and smoking analyses, we determined that  
482 for many proteins tested for causal effects on disease, their associations with incident disease outcomes are  
483 merely coincidental due to their independent associations with smoking. Our partitioning framework illustrates a  
484 strategy for honing in on the proteins that may be most relevant to disease – whether as putative therapeutic  
485 targets for pharmacological intervention (i.e. causal drivers), as quantitative sensors of environmental influences  
486 on disease (i.e. exposure-associated predictors), or alternatively as biomarkers of disease-specific processes  
487 beginning before diagnosis (i.e. non-causal predictors without smoking or other environmental exposure  
488 association). We note that our study is not an exhaustive investigation of potential environmental perturbations,  
489 and thus future work applying our approach to other environmental risk factors will enable further refinement of  
490 these groups. Based on our finding that disease-based proteomic scores for lung cancer and COPD were largely  
491 composed of smoking-associated predictors, other disease-based proteomic scores may similarly include non-  
492 specific signals of other exposures that have extensive effects on the plasma proteome.

493  
494 We developed a proteomic score for smoking that demonstrates the potential of the plasma proteome to provide  
495 precise biomarkers of environmental risk factors. Distributions of the SmokingPS among previous and current  
496 smokers demonstrated the dynamic nature of the plasma proteome, sensitive to not only smoking status but also  
497 intensity and duration of smoking. This points to the ability of the plasma proteome to serve as an accurate and  
498 detailed proxy for smoking behavior and history in addition to, or even instead of, questionnaire data. For  
499 example, we identified a group of individuals in UKB who were labeled as current smokers but had a SmokingPS

500 distribution overlapping that of previous smokers; other questionnaire data in the UKB about smoking cessation  
501 revealed that they were most likely previous smokers. Potentially of even greater utility is the ability of the  
502 SmokingPS to capture and quantify heterogeneity among current smokers in terms of number of cigarettes  
503 smoked or total years of smoking, particularly for datasets that do not have questionnaire data as detailed as in  
504 the UKB.

505  
506 The ability of the SmokingPS to discriminate between current and never smokers (AUC=0.96) was similar  
507 although slightly lower compared to blood-based DNA methylation predictors of smoking, which were able to  
508 discriminate between current and never smokers with an AUC of 0.98<sup>38</sup>. Studies have found that methylation  
509 levels of most cytosine–phosphate–guanine sites (CpGs) revert to that of never smokers within 5 years of  
510 smoking cessation<sup>39,40</sup>, and the SmokingPS of previous smokers similarly resembled that of never smokers after  
511 around 4-5 years of smoking cessation. However, some CpGs have been found to not revert to the levels of  
512 never smokers even after 30 years of smoking cessation. In this study, we observed that proteins with the largest  
513 weights in the SmokingPS returned back to never-smoker levels within 5 years. Given the complexity and slower  
514 pace of epigenetic reprogramming<sup>41</sup> and the greater biological interpretability of protein levels, plasma proteins  
515 may offer a more sensitive and mechanistically informative indicator of smoking. Nevertheless, large-scale  
516 datasets with both types of measurements will be needed to better understand the dynamics between epigenetic  
517 and proteomic biomarkers.

518  
519 The SmokingPS was also associated with incidence of several diseases, particularly diseases for which smoking  
520 is a known risk factor. The addition of the score to the smoking status variable offered modest improvements in  
521 prediction over the smoking status variable alone. This underscores that for diseases like lung cancer and COPD,  
522 the predictive ability of the plasma proteome likely lies in its ability to accurately quantify smoking effects. Sample  
523 sizes were too small in the UKB-PPP cohort to compare disease-based proteomic scores trained in non-smoking  
524 lung cancer and COPD patients, but this approach may reveal insights into the biology of these diseases  
525 unrelated to smoking, or the effects of other non-smoking exposures such as respiratory infections<sup>42</sup> and poor  
526 cooking ventilation<sup>43</sup>.

527  
528 The AlcoholPS demonstrated similar effectiveness in quantifying alcohol intake behaviors, suggesting that the  
529 value of the proteome in this context is not limited to smoking. Self-reported alcohol intake has been shown to  
530 be particularly subject to misreporting and confounding, which has resulted in biased assessments of the effects  
531 of alcohol consumption in observational studies<sup>44</sup>. Thus, there is a critical need for more objective proxies of  
532 alcohol intake, and our study provides evidence that the plasma proteome may fill this gap. However, we note  
533 that we did not observe associations between alcohol status and incident diseases in this study, and thus we did

534 not pursue disease prediction analyses with the AlcoholPS. Alcohol-related disorders, such as alcoholic  
535 pancreatitis (N = 94), alcoholic cardiomyopathy (N = 8), cirrhosis (N = 43), and oropharyngeal cancers (N = 9),  
536 had limited sample sizes of incident outcomes in the UKB-PPP cohort, precluding prediction analyses using the  
537 AlcoholPS.

538  
539 We note several limitations. First, although we obtained the largest GWAS publicly available without UKB data  
540 for estimating the effects of genetic instruments on disease outcomes, limits in sample sizes may have limited  
541 statistical power for estimating effects in MR. Second, we adopted several strategies to ensure the robustness  
542 of the MR results, but violations of model assumptions are still possible; thus, potential causal drivers identified  
543 in this study will require further validation. Third, we utilized data from only European ancestry populations for  
544 the MR analyses, potentially limiting generalizability of the findings to other groups, although causal genetic  
545 effects are largely shared across populations<sup>45,46</sup>. Fourth, the proteins captured by these proteomic assays are  
546 not a completely random subset of the proteome; therefore, proteins involved in certain disease-specific  
547 pathways may not be represented or below limits of detection, and proteins involved in general disease activity  
548 may be disproportionately represented. Finally, we did not explore the possibility that some of the disease-  
549 associated proteins may be biomarkers of early disease processes that begin prior to diagnosis due to limitations  
550 in EHR data. A key challenge for deeply curated clinical cohorts will be distinguishing between effects from  
551 unmeasured environmental factors and those arising from uncharacterized disease mechanisms.

552  
553 In conclusion, we highlight the environment as a prominent component of the reported predictive power of the  
554 disease-associated plasma proteome for patient outcomes. This underscores the challenge of interpreting  
555 proteomic data, given that environmental exposures, alongside inherited genetic variation and ongoing disease  
556 processes, have substantial impacts on plasma proteins. Drawing insights into disease-specific mechanisms  
557 from this data will require systematic characterizations to separate true causal disease signals from more general  
558 reflections of smoking or other factors. This is critical for not only understanding disease biology but also  
559 improving disease prediction, since models with more disease-specific proteomic biomarkers may be more  
560 portable across populations. Especially as the breadth of protein measurements are increasing at decreasing  
561 cost, clarifying the roles of plasma proteomic measurements from these biobank datasets will be an important  
562 step towards clinical translation.

563  
564 Our work also suggests that proteomic assays may open up a path toward measuring the impacts of the  
565 environment on human health and disease. Assessing genetic risk is now largely straightforward, thanks to  
566 assays that enable consistent, reproducible measurements across studies. However, this has not been the case  
567 for studies of environmental risk factors, where data tends to be collected using non-standardized questionnaires

568 and a diversity of other methods. Our findings point to the potential role of proteomic assays as a way to extend  
569 insights from smaller-scale studies on environmental exposures to other studies lacking comparable data, as  
570 well as develop quantitative biomarkers for exposures that may not have existing biological readouts. While our  
571 study focused on evaluating the aggregate predictive power of plasma proteins for disease-related exposures,  
572 complementary studies have focused on partitioning the variance in individual protein levels explained by  
573 exposures vs. genetics<sup>24,25</sup>. Although these studies show that only a small proportion of proteins are more  
574 influenced by non-genetic factors than genetics, we illustrated that in aggregate plasma proteins are highly  
575 accurate readouts of lifestyle factors (AUC = 0.96 for smoking, 0.98 for alcohol intake) because such a large  
576 proportion of the proteins are perturbed by these factors. Well-powered proteomic cohorts with detailed  
577 environmental measures and longitudinal health records are needed to comprehensively disentangle the effects  
578 of the environment vs. other factors like early disease processes on the disease-associated plasma proteome.

## 579 Methods

### 580 Proteomic profiling in the UK Biobank

581 The UK Biobank Pharma Proteomics Project (UKB-PPP) is a precompetitive biopharmaceutical consortium  
582 formed with the goal of collecting and characterizing the plasma proteomic profiles of participants from the UKB<sup>47</sup>,  
583 a population-based cohort comprising approximately 500,000 individuals from the United Kingdom. The UKB  
584 has been described in Bycroft et al.<sup>47</sup> and details on the data available in the UKB can be found at  
585 <https://biobank.ndph.ox.ac.uk/showcase/>. 54,219 participants were selected for the UKB-PPP cohort; we  
586 restricted the sample used for these analyses to participants within the cohort who were randomly sampled from  
587 the main UKB population. The Olink Explore 3072, an antibody-based proximity extension assay, was used to  
588 measure the abundance of protein analytes in each plasma sample. Measurements were provided in the  
589 Normalized Protein eXpression (NPX) values on a log<sub>2</sub> scale. Full details on sample selection, the Olink assay,  
590 and data processing and quality control are described in Sun et al.<sup>6</sup>. We excluded three proteins missing in >10%  
591 of the sample (CTSS, NPM1, and PCOLCE), and imputed protein expression values of the remaining proteins  
592 using the miceforest package in Python. All proteins except those missing in >30% of participants were used as  
593 predictors for the imputation of each protein. We imputed a single dataset using a maximum of five iterations. All  
594 other parameters were left at default values. After imputation, proteomic data were normalized separately within  
595 each cohort by first rescaling values to be between 0 and 1 using MinMaxScaler() from scikit-learn and then  
596 centering on the median. The final dataset consisted of 2,923 proteins measured in 45,438 individuals.

## Disease outcomes and other phenotypes in the UK Biobank

Codes and data used to define prevalent and incident disease in the UKB are detailed in Supplementary Table 20 from our previous publication<sup>48</sup>. Diagnoses and date of first diagnosis for all diseases in the UKB were ascertained using ICD diagnosis codes and corresponding dates of diagnosis taken from linked hospital inpatient, primary care and death register data. If a participant did not have a diagnosis code in hospital inpatient or primary care records, but the code was listed as a primary or secondary cause of death, then they were coded as a case with the date of diagnosis as the date of death. Primary care read codes were converted to corresponding ICD diagnosis codes using the lookup table provided by the UKB. Linked hospital inpatient, primary care and cancer register data were accessed from the UKB data portal on 22 February 2024, with a censoring date of 31 October 2022; 31 August 2022 or 31 May 2022 for participants recruited in England, Scotland or Wales, respectively (8–16 years of follow-up). Detailed information about the linkage procedure national registries for mortality and cause of death information in the UKB is available online (<https://biobank.ctsu.ox.ac.uk/crystal/refer.cgi?id=115559>) with. Mortality data were accessed from the UKB data portal on 22 February 2024, with a censoring date of 30 November 2022 for all participants (12–16 years of follow-up).

We investigated the 23 diseases studied in Gadd et al.<sup>7</sup>, which represent a selection of leading age-related diseases: liver disease, systemic lupus erythematosus, type 2 diabetes, amyotrophic lateral sclerosis, Alzheimer's dementia, endometriosis, COPD, inflammatory bowel disease, rheumatoid arthritis, ischemic stroke, Parkinson's disease, vascular dementia, ischemic heart disease, major depressive disorder, schizophrenia, multiple sclerosis, cystitis, and lung, prostate, breast, gynecological, brain/central nervous system and colorectal cancers. Gynecological and breast cancer, endometriosis and cystitis were female-stratified; prostate cancer was male-stratified. For these diseases, sex was not included as a covariate in the Cox PH models.

For smoking associations, self-reported smoking status (field ID 20116) was coded as a categorical dummy variable (0 = never, 1 = previous, 2 = current). Individuals who were labeled as never smokers in this field but reported being current or previous tobacco smokers (field ID 22506) were excluded from all analyses. Additional smoking variables used included: pack years (field ID 20161), cigarettes smoked per day in current smokers (field ID 3456), cigarettes smoked per day in previous smokers (field ID 2887), age stopped smoking from questionnaire data (field ID 2897), and age stopped smoking from medical records (field ID 22507). Number of years since smoking was computed based on the baseline age and age stopped smoking variables as follows: first, individuals with an age reported in either of the age stopped smoking variables were included; if individuals had ages reported in both fields which differed by 6 or fewer years, they were included and the average age

630 across the two fields was used; finally, years since smoking was calculated by subtracting the age of smoking  
631 cessation from the baseline age.

632  
633 For sensitivity analyses, C-reactive protein (CRP) and immune cell counts were also extracted from UKB (CRP:  
634 field ID 30710; basophil counts: field ID 30160; eosinophil counts: field ID 30150; lymphocyte counts: field ID  
635 30120; monocyte counts: field ID 30130; neutrophil counts: field ID 30140; white blood cell counts: field ID  
636 30000). CRP levels were natural log transformed; the dataset did not have outliers, defined as 4 standard  
637 deviations from the mean. White blood cell counts were transformed to Z-scores and outliers 4 standard  
638 deviations from the mean were excluded (N outliers = 76). Rank-based inverse normal transformation was  
639 performed on all other immune cell count phenotypes.

640  
641 In never smokers, we additionally collated a selection of 71 exposures, covering behavioral (e.g. exercise habits,  
642 media use), lifestyle (e.g. social activities, transportation) and environmental (e.g. pollution measures, traffic  
643 density) domains. The selection was partially curated based on Choi et al.<sup>49</sup>. The full list of variables can be  
644 found in **Supplementary Table 8**.

645  
646 For alcohol analyses, self-reported alcohol status (field ID 20117) was coded as a categorical dummy variable  
647 (0 = never, 1 = previous, 2 = current), and self-reported alcohol intake (field ID 1558) was also coded as a  
648 categorical dummy variable (4 = 1-3 drinks/month, 3 = 1-2 drinks/week, 2 = 3-4 drinks/week, and 1 = daily).  
649 Grams per day of alcohol intake were calculated using a method from a previously published paper in the UKB<sup>50</sup>.  
650 Liver biomarkers were extracted from UKB (ALT: field ID 30620; AST: field ID 30650; GGT: field ID 30730;  
651 Bilirubin: field ID 30840). Outliers, defined as 4 standard deviations from the mean, were excluded (ALT N outliers  
652 = 322; AST N outliers = 269; GGT N outliers = 425; bilirubin N outliers = 396). Rank-based inverse normal  
653 transformation was performed on all liver biomarkers.

## 654 Cox PH models

655 To identify which proteins were significantly associated with incident disease, we ran Cox proportional hazards  
656 (PH) models between each protein and incident disease outcome using the survival package (v3.4-0) in R<sup>51</sup>.  
657 Protein levels were inverse rank normalized for these models. The models were adjusted for age at baseline and  
658 sex; for sex-stratified outcomes, only age was included as a covariate. To assess significance, we used a  
659 Bonferroni-adjusted *P*-value threshold for multiple testing based on the 23 disease outcomes and 2,923 proteins  
660 tested ( $P\text{-value} < 0.05/(23 \times 2923) = 7.41 \times 10^{-7}$ ).

## Mendelian Randomization

To determine the subset of protein-disease associations identified by the Cox PH models that also potentially represent causal relationships, we performed two-sample Mendelian Randomization (MR). We used protein quantitative trait loci (pQTLs) as genetic instruments, which had been previously identified from GWAS of the 2,941 protein analytes measured in 34,557 individuals of European ancestry in the UKB-PPP cohort<sup>6</sup>. We first mapped pQTLs to the GRCh38/hg38 build and selected any genome-wide significant pQTL ( $P$ -value  $< 5 \times 10^{-8}$ ). We then excluded pQTLs located within the human major histocompatibility complex (MHC) region on chromosome 6 (positions 28510120 to 33480577), due to the complex LD structure in this region, and on the sex chromosomes. We additionally excluded pQTLs with  $MAF \leq 0.005$  and an INFO score  $\leq 0.8$ .

We restricted our analyses to the *cis*-pQTLs, defined as those on the same chromosome and within 1Mb of the transcriptional start site of the protein-coding gene. SNPs were allowed to be *cis*-pQTLs for one or more proteins. To identify independent *cis*-pQTLs, we applied clumping in PLINK (v1.9b7-x86\_64) for each protein, using an LD threshold of  $r^2 \leq 0.001$  and a reference panel including the 34,557 individuals of European ancestry from the UKB-PPP GWAS of protein levels<sup>6</sup>. To prevent underflow,  $-\log_{10}(P\text{-values})$  across proteins were scaled to the [0,1] range before clumping.

Genetic association data for the outcomes were selected based on exclusion of UKB data, to avoid overlap with the UKB-PPP cohort. 19 of the 23 disease outcomes had publicly-available GWAS that did not include UKB data and included only individuals of European ancestry, or were available in the FinnGen study<sup>36</sup> (freeze 10). To boost power, we performed a meta-analysis of GWAS from FinnGen and Okada et al.<sup>52</sup> for rheumatoid arthritis (RA) using the METAL software<sup>53</sup>. For outcome GWAS without minor allele frequency (MAF) information, which is required to perform Steiger filtering, we extracted MAF data from the non-Finnish European ancestry group (NFE) in gnomAD v4.1.0<sup>54</sup>. For the rheumatoid arthritis meta-analysis, the MAF used was computed as follows:  $(MAF_{\text{FinnGen}} * N_{\text{FinnGen}} + MAF_{\text{gnomAD\_NFE}} * N_{\text{Okada\_2014}}) / (N_{\text{FinnGen}} + N_{\text{Okada\_2014}})$ , where FinnGen refers to the FinnGen GWAS for RA and Okada\_2014 refers to Okada et al.<sup>52</sup>. The full list of outcome GWAS can be found in **Supplementary Table 2**.

For each outcome GWAS, we then harmonized the *cis*-pQTLs to ensure the effect allele of the SNPs in the exposure and outcome GWAS matched. We used the R package TwoSampleMR<sup>55</sup> to perform the harmonization. We additionally excluded *cis*-pQTLs with  $F$ -statistics  $\leq 10$ , and performed Steiger filtering to exclude variants with larger correlations with the outcome than the exposure. Ultimately, we tested 2,373 protein-disease pairs with genetic instruments that were significant (Bonferroni-adjusted  $P$ -value  $< (19 * 2923) = 8.97 \times 10^{-7}$ ) in the Cox PH associations (**Supplementary Table 3**).

695

696 The Wald ratio test was applied to proteins with only one *cis*-pQTL, and the inverse-variance-weighted (IVW)  
697 method was applied to proteins with 2 or more *cis*-pQTLs. Additionally, for protein-disease pairs with more than  
698 2 *cis*-pQTLs, MR-Egger was applied to test for horizontal pleiotropy. These tests were run using the  
699 MendelianRandomization package in R<sup>56</sup>. Heterogeneity statistics were also extracted from the IVW tests from  
700 the MendelianRandomization package. Scripts to perform genetic instrument harmonization and run the MR  
701 tests were adapted from <https://github.com/globalbiobankmeta/multi-ancestry-pwm><sup>21</sup>.

## 702 Smoking associations

703 Associations between inverse rank normalized protein levels and self-reported smoking status in the UKB were  
704 tested using linear regression using the speedglm package<sup>57</sup>. All models were adjusted for age, sex, age x sex,  
705 age<sup>2</sup>, and age<sup>2</sup> x sex. A Bonferroni-adjusted *P*-value threshold for multiple testing was used (*P*-value < 0.05/2923  
706 = 1.7 x 10<sup>-5</sup>) to assess significance. The same approach was applied to test associations between the inverse  
707 rank normalized protein levels and self-reported alcohol status in the UKB.

708

709 To partition proteins based on the MR analyses and smoking associations, we first extracted all protein-disease  
710 pairs tested in the Cox PH models that were also tested in MR (i.e. had valid genetic instruments). We then took  
711 the subset of protein-disease pairs that reached significance in the Cox PH models based on a Bonferroni-  
712 adjusted *P*-value threshold (*P*-value < 0.05/(19 x 2923) = 8.97 x 10<sup>-7</sup>). The unique proteins across these protein-  
713 disease pairs (N = 782) were then carried forward and combined with the smoking association results to  
714 determine which of these proteins were significantly associated with smoking. Results were visualized and  
715 plotted in an UpSet plot using the ComplexHeatmap package (v2.15.4)<sup>58</sup> in R. All above analyses were  
716 conducted using R v.4.4.0.

## 717 Proteomic Score for smoking and alcohol intake

718 We developed a proteomic score for smoking using LASSO logistic regression in the R package glmnet  
719 (v4.1.8)<sup>59</sup>. We trained this score using the protein levels of UKB-PPP individuals who reported current and never  
720 smoking. We first randomly sampled 50% of this cohort to use for training and performed 10-fold cross-validation  
721 to select protein features and derive their weighting coefficients. We then generated scores in the remaining 50%  
722 of the cohort not used for training by computing the weighted sum of the levels of proteins selected in training.  
723 These scores were evaluated using AUC statistics calculated via the R package pROC (v1.18.5)<sup>60</sup>. We also  
724 evaluated the score, using the LASSO weights of the selected proteins, in participants from the FinnGen study  
725 with plasma protein measurements from Olink Explore and smoking status information (N current/previous  
726 smokers = 850 and N never smokers = 1,013).

727

728 To further evaluate the SmokingPS, we stratified the UKB-PPP cohort by years since smoking cessation in  
729 previous smokers, pack years in current and previous smokers, and cigarettes smoked per day in current  
730 smokers. Current smokers who reported an age stopped smoking or did not report cigarettes smoked per day  
731 were excluded from these analyses. Additionally, outliers with pack years or cigarettes smoked per day 4  
732 standard deviations away from the mean were excluded.

733

734 We tested associations between the SmokingPS and incident disease outcomes (the 19 diseases that were  
735 tested in MR) using Cox PH models using the survival package (v3.4-0) in R<sup>51</sup>. For comparison, seven models  
736 were tested: (1) a baseline model including age, sex, age<sup>2</sup>, age × sex, age<sup>2</sup> × sex, and first 10 genetic PCs,  
737 sourced from the Pan-UKB project<sup>61</sup>; (2) all covariates in the baseline model and smoking status; (3) all  
738 covariates in the baseline model and the SmokingPS; (4) all covariates in the baseline model, smoking status,  
739 and the SmokingPS; (5) all covariates in the baseline model and CRP; (6) all covariates in the baseline model  
740 and immune cell counts; (7) all covariates in the baseline model, smoking status, the SmokingPS, CRP, and  
741 immune cell counts. Results were plotted using the R package forestploter<sup>62</sup>. Associations between the  
742 SmokingPS and 71 environmental exposures were tested in never smokers using the glm function in R.

743

744 For the proteomic score for alcohol intake, we followed the same protocol used to develop and evaluate the  
745 SmokingPS. We trained this score using the protein levels of UKB-PPP individuals who reported never (N =  
746 2,165) and daily (N = 9,185) drinking. We also developed sex-specific AlcoholPS, by first splitting individuals by  
747 sex and then training the scores separately in males (N daily = 5,233; N never = 653) and females (N daily =  
748 3,891; N never = 1,498). We additionally matched the sample sizes of never and daily drinkers to that of the  
749 smaller sex-specific group (i.e. downsampled daily drinkers in males to the sample size of daily drinkers in  
750 females and never drinkers in females to the sample size of never drinkers in males), and evaluated the score  
751 trained on the combined dataset of males and females with matching sample sizes. To evaluate the AlcoholPS,  
752 we stratified the current drinkers by alcohol intake frequency, as well as bins of grams of alcohol intake as defined  
753 by the derived alcohol intake variable.

## 754 Data availability

755 Individual-level data from the UK Biobank can be accessed via application at <https://www.ukbiobank.ac.uk/>. We  
756 accessed UK Biobank data under application 31063. GWAS summary statistics can be accessed as described  
757 in their respective papers in Supplementary Table 2 and in the FinnGen study<sup>36</sup>. Weights for the proteomic scores  
758 will be made available upon publication. Further information and requests for resources should be directed to

and will be fulfilled by the lead contact, Kristin Tsuo ([ktsuo@broadinstitute.org](mailto:ktsuo@broadinstitute.org)), upon reasonable request. Code used for data preparation and analysis will be made available on GitHub before journal publication.

## Author Contributions

Study design, K.T., A.R.M., and M.J.D.; data analysis, K.T., M.A.A., D.G., D.B.; interpretation of results, K.T., A.R.M., M.J.D., M.A.A., M.K., Z.Z., R.E.M., C.F., H.H., B.S., and C.C.; writing, K.T., A.R.M., and M.J.D.

## Acknowledgements

A.R.M is funded by NIH U01HG011719 as well as Broad Institute Next Gen and Merkin awards. K.T. is funded by F31HL167378 and supported by the ECOR Clafin Award to A.R.M.

## References

1. Gudjonsson, A. *et al.* A genome-wide association study of serum proteins reveals shared loci with common diseases. *Nat. Commun.* **13**, 480 (2022).
2. Emilsson, V. *et al.* Co-regulatory networks of human serum proteins link genetics to disease. *Science* **361**, 769–773 (2018).
3. Sun, B. B. *et al.* Genomic atlas of the human plasma proteome. *Nature* **558**, 73–79 (2018).
4. Ferkingstad, E. *et al.* Large-scale integration of the plasma proteome with genetics and disease. *Nat. Genet.* **53**, 1712–1721 (2021).
5. Koprulu, M. *et al.* Proteogenomic links to human metabolic diseases. *Nat. Metab.* **5**, 516–528 (2023).
6. Sun, B. B. *et al.* Plasma proteomic associations with genetics and health in the UK Biobank. *Nature* **622**, 329–338 (2023).
7. Gadd, D. A. *et al.* Blood protein assessment of leading incident diseases and mortality in the UK Biobank. *Nat Aging* (2024) doi:10.1038/s43587-024-00655-7.
8. Carrasco-Zanini, J. *et al.* Proteomic prediction of common and rare diseases. *bioRxiv* (2023) doi:10.1101/2023.07.18.23292811.
9. Deng, Y.-T. *et al.* Atlas of the plasma proteome in health and disease in 53,026 adults. *Cell* **188**, 253–

- 783 271.e7 (2025).
- 784 10. Schuermans, A. *et al.* Integrative proteomic analyses across common cardiac diseases yield new  
785 mechanistic insights and enhanced prediction. *medRxiv* (2023) doi:10.1101/2023.12.19.23300218.
- 786 11. Ahsan, H. Monoplex and multiplex immunoassays: approval, advancements, and alternatives. *Comp. Clin.*  
787 *Path.* **31**, 333–345 (2022).
- 788 12. Bowman, W. S. *et al.* Proteomic biomarkers of progressive fibrosing interstitial lung disease: a multicentre  
789 cohort analysis. *Lancet Respir. Med.* **10**, 593–602 (2022).
- 790 13. Gudmundsdottir, V. *et al.* Circulating protein signatures and causal candidates for type 2 diabetes.  
791 *Diabetes* **69**, 1843–1853 (2020).
- 792 14. Schuermans, A. *et al.* Genetic associations of circulating cardiovascular proteins with gestational  
793 hypertension and preeclampsia. *JAMA Cardiol.* **9**, 209–220 (2024).
- 794 15. Chong, M. *et al.* Novel drug targets for ischemic stroke identified through Mendelian randomization  
795 analysis of the blood proteome. *Circulation* **140**, 819–830 (2019).
- 796 16. Lu, T., Forgetta, V., Greenwood, C. M. T., Zhou, S. & Richards, J. B. Circulating proteins influencing  
797 psychiatric disease: A Mendelian randomization study. *Biol. Psychiatry* **93**, 82–91 (2023).
- 798 17. Gong, W. *et al.* Genomics-driven integrative analysis highlights immune-related plasma proteins for  
799 psychiatric disorders. *J. Affect. Disord.* **370**, 124–133 (2024).
- 800 18. Bourgault, J. *et al.* Proteome-wide Mendelian randomization identifies causal links between blood proteins  
801 and acute pancreatitis. *Gastroenterology* **164**, 953–965.e3 (2023).
- 802 19. Zhang, S. *et al.* IDDF2024-ABS-0203 Causal links between plasma proteome and digestive system  
803 diseases: mendelian randomization analysis. in *Basic Gastroenterology* A143.1–A143 (BMJ Publishing  
804 Group Ltd and British Society of Gastroenterology, 2024).
- 805 20. Su, C.-Y. *et al.* Multi-ancestry proteome-phenome-wide Mendelian randomization offers a comprehensive  
806 protein-disease atlas and potential therapeutic targets. *medRxiv* (2024)  
807 doi:10.1101/2024.10.17.24315553.
- 808 21. Zhao, H. *et al.* Proteome-wide Mendelian randomization in global biobank meta-analysis reveals multi-

- ancestry drug targets for common diseases. *Cell Genom* **2**, None (2022).
22. Guo, Y. *et al.* Plasma proteomic profiles predict future dementia in healthy adults. *Nat. Aging* **4**, 247–260 (2024).
23. Chan, K. H. *et al.* An exposome-wide assessment of 6600 SomaScan proteins with non-genetic factors in Chinese adults. *medRxiv* (2024) doi:10.1101/2024.10.24.24316041.
24. Carrasco-Zanini, J. *et al.* Mapping biological influences on the human plasma proteome beyond the genome. *Nat. Metab.* **6**, 2010–2023 (2024).
25. Isaac, S. *et al.* Human plasma proteomics links modifiable lifestyle exposome to disease risk. *medRxiv* 2025.05.07.25327178 (2025) doi:10.1101/2025.05.07.25327178.
26. Burgess, S. & Thompson, S. G. Interpreting findings from Mendelian randomization using the MR-Egger method. *Eur. J. Epidemiol.* **32**, 377–389 (2017).
27. Bergeron, N., Phan, B. A. P., Ding, Y., Fong, A. & Krauss, R. M. Proprotein convertase subtilisin/Kexin type 9 inhibition: A new therapeutic mechanism for reducing cardiovascular disease risk. *Circulation* **132**, 1648–1666 (2015).
28. Raulin, A.-C. *et al.* ApoE in Alzheimer’s disease: pathophysiology and therapeutic strategies. *Mol. Neurodegener.* **17**, 72 (2022).
29. Penrose, H. M. *et al.* Ulcerative colitis immune cell landscapes and differentially expressed gene signatures determine novel regulators and predict clinical response to biologic therapy. *Sci. Rep.* **11**, 9010 (2021).
30. Feiner, J. *et al.* 221-LB: Identification of PAM as novel monogenic diabetes gene. *Diabetes* **72**, 221–LB (2023).
31. Sun, Z. *et al.* Comprehensive mendelian randomization analysis of plasma proteomics to identify new therapeutic targets for the treatment of coronary heart disease and myocardial infarction. *J. Transl. Med.* **22**, 404 (2024).
32. Jin, H. *et al.* Smoking, COPD, and 3-nitrotyrosine levels of plasma proteins. *Environ. Health Perspect.* **119**, 1314–1320 (2011).

- 835 33. Elisia, I. *et al.* The effect of smoking on chronic inflammation, immune function and blood cell composition.  
836 *Sci. Rep.* **10**, 19480 (2020).
- 837 34. Chan, K. H. *et al.* Tobacco smoking and risks of more than 470 diseases in China: a prospective cohort  
838 study. *Lancet Public Health* **7**, e1014–e1026 (2022).
- 839 35. Zhang, J.-C. *et al.* TGF- $\beta$ /BAMBI pathway dysfunction contributes to peripheral Th17/Treg imbalance in  
840 chronic obstructive pulmonary disease. *Sci. Rep.* **6**, 31911 (2016).
- 841 36. Kurki, M. I. *et al.* FinnGen provides genetic insights from a well-phenotyped isolated population. *Nature*  
842 **613**, 508–518 (2023).
- 843 37. Saint-André, V. *et al.* Smoking changes adaptive immunity with persistent effects. *Nature* **626**, 827–835  
844 (2024).
- 845 38. Chybowska, A. D. *et al.* A blood- and brain-based EWAS of smoking. *Nat. Commun.* **16**, 3210 (2025).
- 846 39. Joehanes, R. *et al.* Epigenetic signatures of cigarette smoking. *Circ. Cardiovasc. Genet.* **9**, 436–447  
847 (2016).
- 848 40. McCartney, D. L. *et al.* Epigenetic signatures of starting and stopping smoking. *EBioMedicine* **37**, 214–220  
849 (2018).
- 850 41. Fang, F., Andersen, A. M., Philibert, R. & Hancock, D. B. Epigenetic biomarkers for smoking cessation.  
851 *Addict. Neurosci.* **6**, 100079 (2023).
- 852 42. Frickmann, H. *et al.* The influence of virus infections on the course of COPD. *Eur. J. Microbiol. Immunol.*  
853 *(Bp.)* **2**, 176–185 (2012).
- 854 43. Mortimer, K. *et al.* Household air pollution and COPD: cause and effect or confounding by other aspects of  
855 poverty? *Int. J. Tuberc. Lung Dis.* **26**, 206–216 (2022).
- 856 44. Xue, A. *et al.* Genome-wide analyses of behavioural traits are subject to bias by misreports and  
857 longitudinal changes. *Nat. Commun.* **12**, 20211 (2021).
- 858 45. Hou, K. *et al.* Causal effects on complex traits are similar for common variants across segments of  
859 different continental ancestries within admixed individuals. *Nat. Genet.* **55**, 549–558 (2023).
- 860 46. Hu, S. *et al.* Fine-scale population structure and widespread conservation of genetic effect sizes between

- 861 human groups across traits. *Nat. Genet.* **57**, 379–389 (2025).
- 862 47. Bycroft, C. *et al.* The UK Biobank resource with deep phenotyping and genomic data. *Nature* **562**, 203–  
863 209 (2018).
- 864 48. Argentieri, M. A. *et al.* Proteomic aging clock predicts mortality and risk of common age-related diseases  
865 in diverse populations. *Nat. Med.* **30**, 2450–2460 (2024).
- 866 49. Choi, K. W. *et al.* An exposure-wide and Mendelian randomization approach to identifying modifiable  
867 factors for the prevention of depression. *Am. J. Psychiatry* **177**, 944–954 (2020).
- 868 50. Evangelou, E. *et al.* New alcohol-related genes suggest shared genetic mechanisms with neuropsychiatric  
869 disorders. *Nat. Hum. Behav.* **3**, 950–961 (2019).
- 870 51. Therneau, T. A package for survival analysis in R. Preprint at [https://cran.r-](https://cran.r-project.org/web/packages/survival/vignettes/survival.pdf)  
871 [project.org/web/packages/survival/vignettes/survival.pdf](https://cran.r-project.org/web/packages/survival/vignettes/survival.pdf) (2024).
- 872 52. Okada, Y. *et al.* Genetics of rheumatoid arthritis contributes to biology and drug discovery. *Nature* **506**,  
873 376–381 (2014).
- 874 53. Willer, C. J., Li, Y. & Abecasis, G. R. METAL: fast and efficient meta-analysis of genomewide association  
875 scans. *Bioinformatics* **26**, 2190–2191 (2010).
- 876 54. Chen, S. *et al.* A genomic mutational constraint map using variation in 76,156 human genomes. *Nature*  
877 **625**, 92–100 (2024).
- 878 55. Hemani, G. *et al.* The MR-Base platform supports systematic causal inference across the human  
879 phenome. *Elife* **7**, (2018).
- 880 56. Burgess, S. & Yavorska, O. MendelianRandomization: Mendelian randomization package. *CRAN:*  
881 *Contributed Packages* The R Foundation <https://doi.org/10.32614/cran.package.mendelianrandomization>  
882 (2016).
- 883 57. Enea, M. *\_speedglm: Fitting Linear and Generalized Linear Models to Large Data Sets\_*. (2023).
- 884 58. Gu, Z., Eils, R. & Schlesner, M. Complex heatmaps reveal patterns and correlations in multidimensional  
885 genomic data. *Bioinformatics* **32**, 2847–2849 (2016).
- 886 59. Friedman, J., Hastie, T. & Tibshirani, R. Regularization paths for generalized linear models via coordinate

- 887 descent. *J. Stat. Softw.* **33**, 1–22 (2010).
- 888 60. Robin, X. *et al.* pROC: an open-source package for R and S+ to analyze and compare ROC curves. *BMC*  
889 *Bioinformatics* **12**, 77 (2011).
- 890 61. Karczewski, K. J. *et al.* Pan-UK Biobank GWAS improves discovery, analysis of genetic architecture, and  
891 resolution into ancestry-enriched effects. (2024) doi:10.1101/2024.03.13.24303864.
- 892 62. Dayimu, A. *\_forestploter: Create a Flexible Forest Plot\_*. *R package version 1.1.2* [https://CRAN.R-](https://CRAN.R-project.org/package=forestploter)  
893 [project.org/package=forestploter](https://CRAN.R-project.org/package=forestploter) (2024).