

# 1 **Blood protein levels predict leading incident diseases and mortality in UK Biobank**

2 Danni A. Gadd<sup>1,2</sup>, Robert F. Hillary<sup>1,2</sup>, Zhana Kuncheva<sup>1,3</sup>, Tasos Mangelis<sup>1,3</sup>, Yipeng Cheng<sup>2</sup>, Manju  
3 Dissanayake<sup>1,3</sup>, Romi Admanit<sup>4</sup>, Jake Gagnon<sup>4</sup>, Tinchu Lin<sup>4</sup>, Kyle Ferber<sup>4</sup>, Heiko Runz<sup>5</sup>, Biogen  
4 Biobank Team, Riccardo E. Marioni<sup>\*1,2</sup>, Christopher N. Foley<sup>\*1,3</sup>, Benjamin B. Sun<sup>5,6\*</sup>

5 <sup>1</sup> Optima Partners, Edinburgh, EH2 4HQ, UK.

6 <sup>2</sup> Centre for Genomic and Experimental Medicine, Institute of Genetics and Cancer, University of  
7 Edinburgh, Edinburgh, EH4 2XU, UK.

8 <sup>3</sup> Bayes Centre, The University of Edinburgh, Edinburgh, EH8 9BT, UK.

9 <sup>4</sup> Global Analytics and Data Science, Research and Development, Biogen Inc., Cambridge, MA, USA

10 <sup>5</sup> Translational Sciences, Research and Development, Biogen Inc. Cambridge, MA, USA.

11 <sup>6</sup> Cardiovascular Epidemiology Unit, Department of Public Health and Primary Care, University of  
12 Cambridge, Cambridge, CB1 8RN, UK.

13 \* Equal contributions.

14 Correspondence: [benjamin.sun@biogen.com](mailto:benjamin.sun@biogen.com) and [riccardo.marioni@ed.ac.uk](mailto:riccardo.marioni@ed.ac.uk).

15

## 16 **Abstract**

17

18 The circulating proteome offers insights into the biological pathways that underlie disease. Here, we  
19 test relationships between 1,468 Olink protein levels and the incidence of 23 age-related diseases and  
20 mortality, over 16 years of electronic health linkage in the UK Biobank (N=47,600). We report 3,201  
21 associations between 961 protein levels and 21 incident outcomes, identifying proteomic indicators of  
22 multiple morbidities. Next, protein-based scores (ProteinScores) are developed using penalised Cox  
23 regression. When applied to test sets, six ProteinScores improve Area Under the Curve (AUC)  
24 estimates for the 10-year onset of incident outcomes beyond age, sex and a comprehensive set of 24  
25 lifestyle factors, clinically-relevant biomarkers and physical measures. Furthermore, the ProteinScore  
26 for type 2 diabetes outperformed a polygenic risk score, a metabolomic score and HbA1c – a clinical  
27 marker used to monitor and diagnose type 2 diabetes. These data characterise early proteomic  
28 contributions to major age-related disease and demonstrate the value of the plasma proteome for risk  
29 stratification.

## 30 **Introduction**

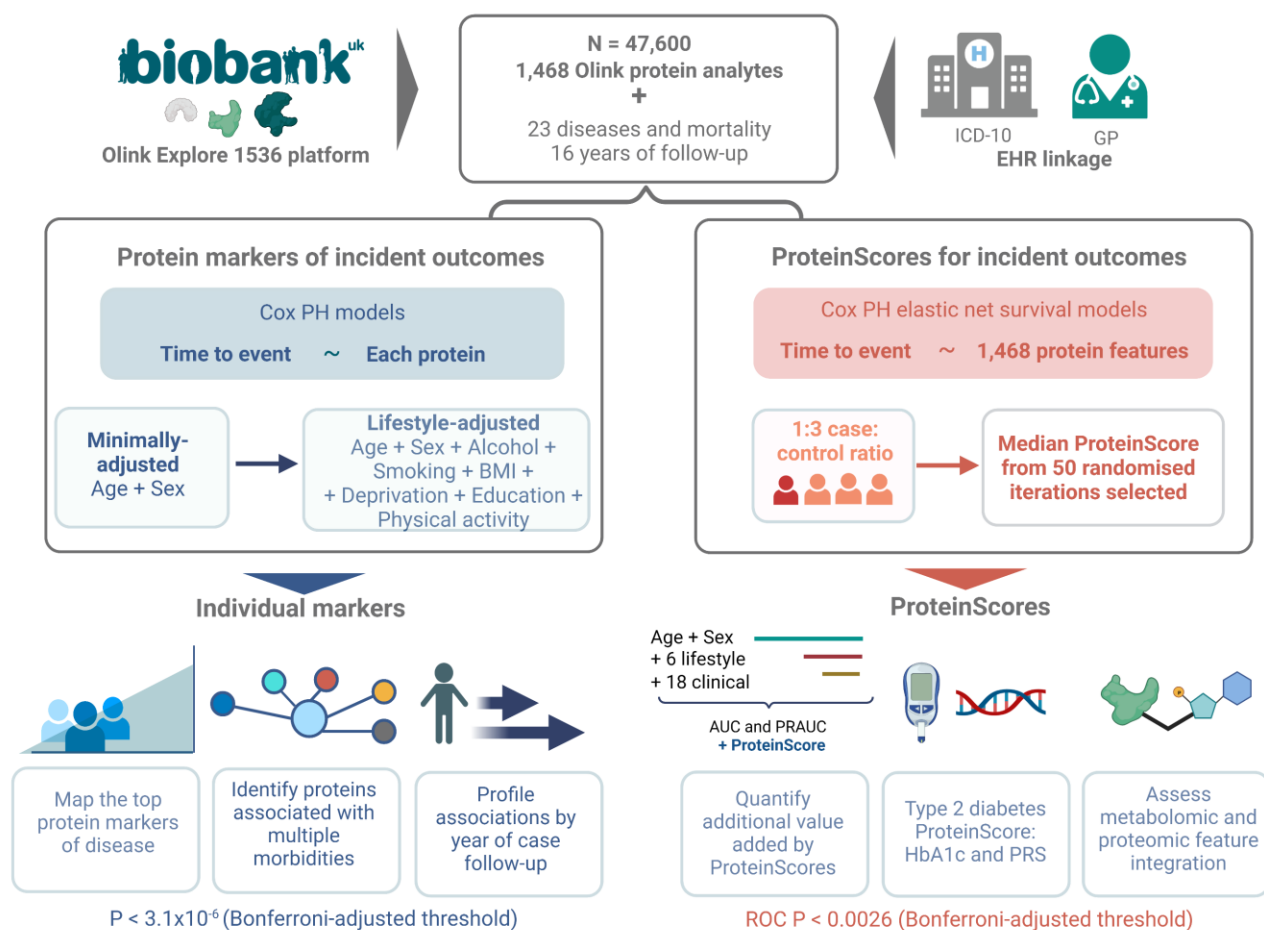
31

32 Omics signatures are increasingly used to hone clinical trial design <sup>1</sup>, while also opening up avenues  
33 for more personalised healthcare <sup>2,3</sup>. Of all the omics layers that can be measured from a single blood  
34 test, proteomics arguably holds the most intrinsic predictive potential, given that proteins are the  
35 intermediary effectors of health maintenance and disease and are often the targets of pharmacological  
36 interventions. Several studies have shown that circulating proteins can discriminate disease cases from  
37 controls and delineate risk of incident diagnoses <sup>4-11</sup>. Screening the proteome against incident  
38 outcomes has been shown to identify sets of individual protein markers – some of which have then  
39 been causally-implicated in disease <sup>8,12-14</sup>. This demonstrates the value protein data have in informing  
40 therapeutic targeting and reflecting the internal processes occurring in the body that precede formal  
41 diagnoses.

42 While singular protein markers offer insight into the mediators of disease, harnessing multiple proteins  
43 simultaneously can be expected to generate predictive tools with even greater clinical utility <sup>15</sup>.  
44 Although cross-sectional case-control studies can inform on the molecular signatures of diagnosed  
45 diseases, longitudinal approaches that assess early biomarker signatures relating to time-to-disease are  
46 more suited to risk stratification. Clinically-available risk profiling scores that rely on lifestyle and  
47 health information such as QRISK and ASSIGN typically profile 10-year onset risk of disease <sup>16,17</sup>.  
48 Scores such as these stratify where individuals lie on the disease-risk continuum for a population, but  
49 do not include omics features. While proteomic and metabolomics scores have been developed for  
50 certain time-to-event outcomes in isolation <sup>9,18-22</sup>, these predictors are rarely developed and tested at  
51 scale. Proteomic predictors have been trained using the SomaScan platform for diabetes and  
52 cardiovascular event risk and multiple lifestyle and health indicators <sup>23</sup>. Metabolomics data have been  
53 recently shown to facilitate incident disease prediction in the UK Biobank <sup>24</sup>. However, no study has  
54 systematically assessed proteomic score generation for multiple incident morbidities.

55 Here, we quantify how large-scale proteomic sampling can identify candidate protein targets and  
56 facilitate the prediction of incident outcomes in the UK Biobank (**Fig. 1**). We use 1,468 Olink plasma  
57 protein measurements in 47,600 individuals available as part of the UK Biobank Pharma Proteomics  
58 Project (UKB-PPP) <sup>25</sup>. First, Cox proportional hazards (PH) models are used to characterise  
59 associations between each protein and 23 incident diseases, ascertained via data linkage to primary  
60 and secondary care records and mortality over 16 years of follow-up. Next, the dataset is randomly  
61 split into training and testing subsets to train proteomic scores (ProteinScores) and assess their utility  
62 for modelling either 5-year or 10-year onset of the 19 incident outcomes that had a minimum of 150  
63 cases available. Type 2 diabetes is taken forward to explore the potential value that ProteinScores may  
64 offer, beyond clinical biomarkers, polygenic risk scores (PRS) and metabolomics measures.

65



66

67

68 **Figure 1. Proteomic assessment of 23 incident diseases and mortality in the UK Biobank**

69 (N=47,600). First, individual Cox proportional hazards (PH) models were used to profile relationships

70 between baseline protein analytes and incident diseases or death, over a maximum of 16 years of

71 electronic health linkage. Associations that had  $P < 3.1 \times 10^{-6}$  (Bonferroni-adjusted threshold) in

72 minimally-adjusted (age and sex) and lifestyle-adjusted models were retained. Proteins associated with

73 multiple morbidities were identified and associations were explored by year of case follow-up. Next,

74 proteomic predictors (ProteinScores) were trained using Cox PH elastic net regression for 19 of the

75 incident outcomes with a minimum of 150 cases. All ProteinScores were developed for 10-year onset

76 of disease, except endometriosis, cystitis and amyotrophic lateral sclerosis that had case distributions

77 that were better-suited to 5-year assessment (80% of cases diagnosed by year 8 of follow-up). Of fifty

78 ProteinScore iterations with randomly sampled train and test populations, the ProteinScore with

79 median improvement in AUC beyond a minimally-adjusted model was selected. Improvements in

80 AUC and PRAUC due to adding the ProteinScores into models with increasingly complex covariate

81 structures were quantified. The type 2 diabetes trait was taken forward as a case study to explore the

82 potential value ProteinScores may offer, in the context of HbA1c (a clinically used biomarker), a

polygenic risk score (PRS) and integration of metabolomics features for scoring.

83

84

## 85 **Results**

### 86 **The UKB-PPP sample**

87 Of the 1,472 protein levels available in the UKB-PPP sample, 1,463 are unique, due to CXCL8, IL6  
88 and TNF having multiple analyte measurements (annotation information provided in **Supplementary**  
89 **Table 1**). After quality control and removal of outliers, measurements for 52,744 individuals were  
90 available. In this study, a total sample of 47,600 individuals with 1,468 protein analytes was used, after  
91 exclusions for related individuals and missing data (**Supplementary Fig. 1, Methods**). The 1,468  
92 analyte measurements correspond to 1,459 unique protein levels. Demographic and phenotypic  
93 information is presented in **Supplementary Table 2**. Principal components analyses indicated that the  
94 first 678 components explained a cumulative variance of 90% in the protein levels (**Supplementary**  
95 **Table 3**).

### 96 **Protein associations with incident outcomes**

97 First, differential plasma protein levels that were associated with the onset of 23 diseases (that included  
98 leading causes of disability, morbidity and reductions in healthy life expectancy)<sup>26–28</sup> were identified,  
99 up to 16 years prior to formal diagnoses. Time-to-mortality was also considered as an outcome (4,446  
100 individuals had died during the 16-year follow-up period). A total of 35,232 associations were tested  
101 (1,468 analytes and 24 outcomes). The number of cases and controls available in Cox PH models, with  
102 mean time-to-onset for cases is presented for each outcome in **Table 1**.

103 In minimally-adjusted (age- or age- and sex-adjusted) models, there were 5,252 associations between  
104 1,209 unique protein analytes and 23 outcomes (Bonferroni-adjusted P threshold =  $3.1 \times 10^{-6}$ )  
105 (**Supplementary Table 4**). Further adjustment for health and lifestyle risk factors (body mass index  
106 (BMI), alcohol consumption, social deprivation, education status, smoking status and physical  
107 activity) led to the attenuation of 2,051 of the minimally-adjusted associations, with 3,201 that

108 remained (Bonferroni-adjusted P threshold =  $3.1 \times 10^{-6}$ ) (**Fig. 2a, Supplementary Table 5**). The 3,201  
109 associations involved 961 unique protein analytes and 21 outcomes, ranging from one association for  
110 amyotrophic lateral sclerosis, cystitis and multiple sclerosis, to 646 and 664 for mortality and liver  
111 disease, respectively. No associations were found for brain/CNS cancer, major depression and  
112 schizophrenia. **Supplementary Table 6** summarises the 961 unique protein analytes selected across  
113 the 3,201 associations by disease and by direction of effect (i.e. 303 associations with Hazard Ratio  
114 (HR) < 1 and 2,898 associations with HR > 1).

### 115 **Proteomic signatures of multimorbidity**

116 Fifty-four proteins had associations with eight or more incident morbidities (**Fig. 2b**); in all instances,  
117 elevated levels of the proteins were associated with the increased incidence of disease or death (i.e.  
118 HR > 1). Of the 54 proteins, GDF15 had the largest number of associations (11 incident outcomes),  
119 followed by IL6 and PLAUR (10 incident outcomes). In logistic regression models run between the  
120 1,468 protein analytes and multimorbidity status (a binary trait defined as individuals that had three or  
121 more diagnoses of the 23 diseases over the 16-year follow-up period), 720 associations had  $P < 3.1 \times 10^{-6}$   
122 (**Supplementary Table 7**). All 54 proteins that were associated with eight or more morbidities in the  
123 Cox PH associations were present in the multimorbidity status associations. GDF15, TNFRSF10B,  
124 WFDC2 and PLAUR had both the largest absolute effect sizes and smallest p-values, which was  
125 consistent with their position as top markers of multimorbidity in the individual Cox PH associations  
126 presented in **Fig. 2b**.

### 127 **Cox PH sensitivity analyses**

128 Understanding whether protein-disease associations are stronger in the near-term of case follow-up is  
129 of interest when considering the clinical use-case for biomarkers. Modelling near-term versus long-  
130 term case follow-up is also important to understand the confidence that can be ascribed to associations  
131 failing the Cox PH assumption (Schoenfeld residual test  $P < 0.05$ ). Therefore, a sensitivity analyses

132 that modelled each of the 35,232 Cox PH associations over increasing yearly case follow-up intervals  
133 was performed (**Supplementary Table 8**). Of the 3,201 protein-disease associations identified over  
134 the maximum 16-year follow-up, 2,915 and 1,957 of these associations remained ( $P < 3.1 \times 10^{-6}$ , the  
135 Bonferroni-adjusted threshold) when restricting cases up to 10-year and 5-year onset, respectively  
136 (**Supplementary Table 9**). Of the 684 failures in the local (protein) Cox PH assumption observed in  
137 the 16-year follow-up analyses, 665 and 410 were observed in the 10-year and 5-year onset analyses.  
138 Relatively minor deviations in magnitude of effect size were observed for these associations by year  
139 of follow-up. These results can be examined visually for each of the 35,232 protein-disease  
140 associations tested in a Shiny app available at: <https://protein-disease-ukb.optima-health.technology>  
141 [Username: ukb\_diseases, Password: UKBshinyapp]. The app also includes an interactive network for  
142 the 3,201 associations that can be manipulated to view multiple proteins and examine their associations  
143 with multiple incident morbidities.

144 A sensitivity analysis was performed to explore the potential impact of medication use on individual  
145 Cox PH associations. A subset of the population with proteomics measures had medication information  
146 available (35,073 of 47,600 individuals). Ischaemic heart disease was chosen, as a range of blood-  
147 pressure lowering medications are used to delay or prevent this disease and these medications were  
148 amongst the most commonly-reported in the population (14,074 of 35,073 individuals reported use of  
149 either statins, antihyperintensives, diuretics, beta blockers, calcium channel blockers  
150 or renin-angiotensin system actors at baseline (**Supplementary Table 10**). In the subset of 35,073  
151 individuals, 370 of the original 403 associations (adjusting for age, sex and six lifestyle factors) for  
152 ischaemic heart disease had  $P < 3.1 \times 10^{-6}$  (**Supplementary Table 11**). With further adjustment for  
153 blood-pressure lowering medication use, 36 of these associations were attenuated, while 344 had  $P <$   
154  $3.1 \times 10^{-6}$ . None of the attenuated associations were present in the top 100 marker associations (ranked  
155 by P-value, or effect size). Adjustment for blood-pressure lowering medication tended to reduce the

156 magnitude of the effect estimate generally across the 370 associations (**Supplementary Fig.2**), but  
157 hazard ratios were nonetheless highly correlated (Pearson correlation  $r = 0.99$ ).

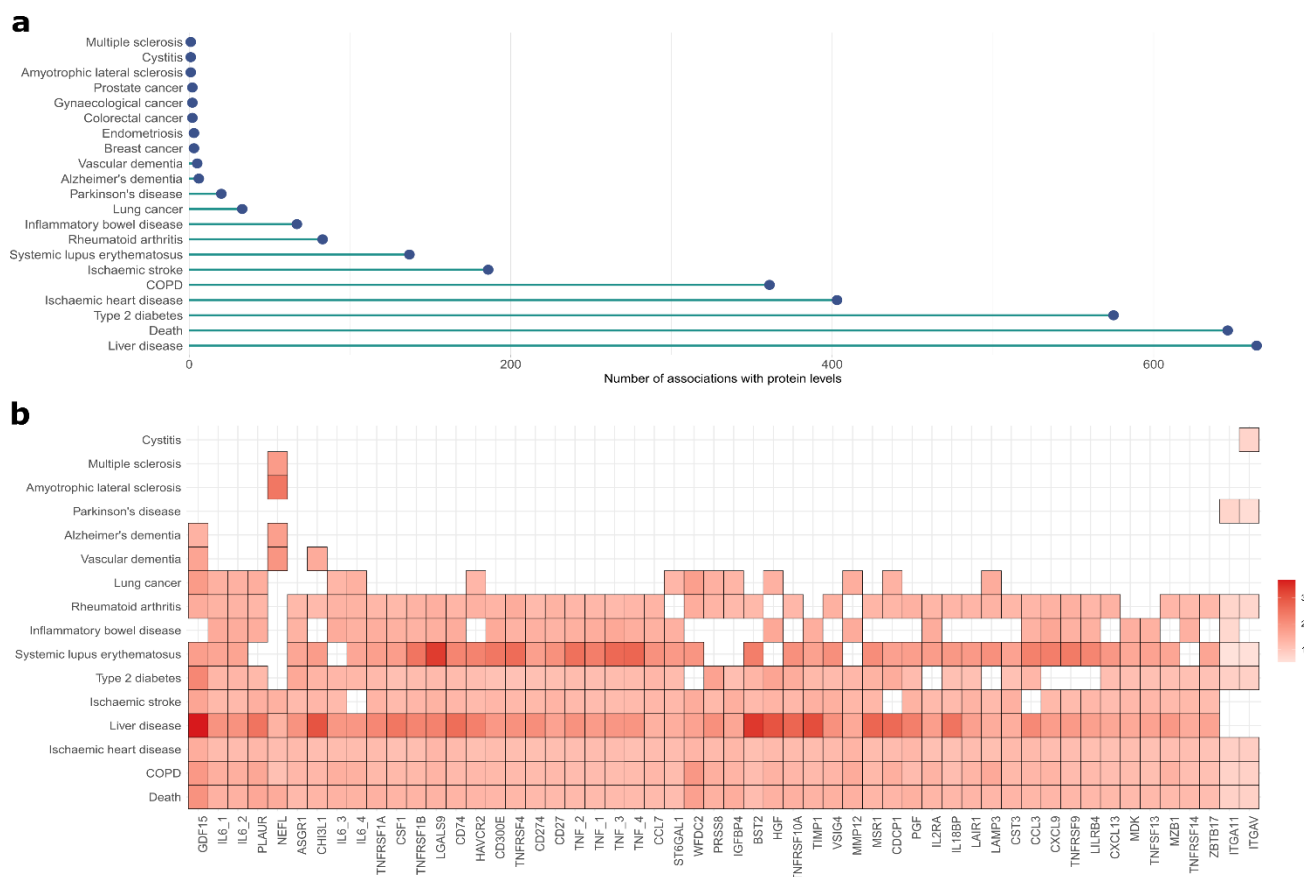


Incident diagnosis	Incident cases (N)	Controls (N)	Mean years to incident case diagnosis (sd)
Schizophrenia	54	47449	6.5 (3.4)
Brain/CNS cancer	82	47507	5.5 (2.8)
Multiple sclerosis	96	47165	5.6 (3.2)
Major depression	111	47229	4.2 (3.1)
Systemic lupus erythematosus	134	47096	5.1 (2.6)
Endometriosis <sup>a</sup>	157	24768	4.8 (3.3)
Vascular dementia <sup>b</sup>	195	33907	8.1 (3)
Gynaecological cancer <sup>a</sup>	256	25185	5 (3)
Amyotrophic lateral sclerosis	264	47269	5.4 (2.7)
Inflammatory bowel disease	275	46727	5.9 (3.3)
Lung cancer	403	47158	5.9 (3.2)
Liver disease	432	47104	7 (3.3)
Alzheimer's dementia <sup>b</sup>	446	33642	7.8 (2.8)
Colorectal cancer	508	46890	5.8 (3.1)
Cystitis <sup>a</sup>	531	24160	4.1 (3)
Rheumatoid arthritis	593	46310	6.8 (3.2)
Parkinson's disease	659	46802	5.4 (3.2)
Ischaemic stroke	765	46657	6.8 (3.4)
Breast cancer <sup>a</sup>	772	24086	5.2 (3.1)
Prostate cancer <sup>a</sup>	1001	20628	5.7 (3.1)
Chronic obstructive pulmonary disease	1998	44948	6.3 (3.4)
Type 2 diabetes	2822	43370	6 (3.3)
Ischaemic heart disease	3338	41341	6.3 (3.4)
Death	4445	43155	7.9 (3.5)

158

159 **Table 1. The 24 incident outcomes profiled over a maximum of 16 years of follow-up in the UK**  
160 **Biobank (N=47,600).** Counts for incident cases and controls are provided, with mean years to  
161 diagnosis for incident cases. These data were used in individual Cox PH models to identify protein  
162 levels that were associated with incident outcomes. <sup>a</sup> Sex-stratified traits. <sup>b</sup> Alzheimer's and vascular  
163 dementia were restricted to individuals aged 65 years or above at the time of diagnosis for cases, or at  
164 the time or censoring for controls. CNS: central nervous system.

165



166

167

168

169

170

171

172

173

174

175

176

177

178

**Figure 2. Individual protein associations with incident outcomes in the UK Biobank (N=47,600).**

**a**, Number of associations between protein analytes and time-to-onset for 20 outcomes that had  $P < 3.1 \times 10^{-6}$  (Bonferroni-adjusted threshold) in both basic and fully-adjusted Cox PH models. There were 3,201 associations in total involving 961 protein analytes. **b**, Hazard ratios (HR) per a one SD increase in levels of the transformed protein analytes are plotted for the 54 protein analytes that were associated with eight or more outcomes in the individual Cox PH models. Each association is represented by a rectangle. Cox PH models were adjusted for age, sex and six lifestyle factors (BMI, alcohol consumption, social deprivation, educational attainment, smoking status and physical activity). Every association identified for these proteins had  $HR > 1$  (red) and associations are shaded based on HR effect size (darkest colouration indicating larger magnitude of effect). The largest HR shown is for the association between GDF15 levels and liver disease ( $HR = 3.67$ ). COPD: chronic obstructive pulmonary disease.

179

## 180 **ProteinScore development**

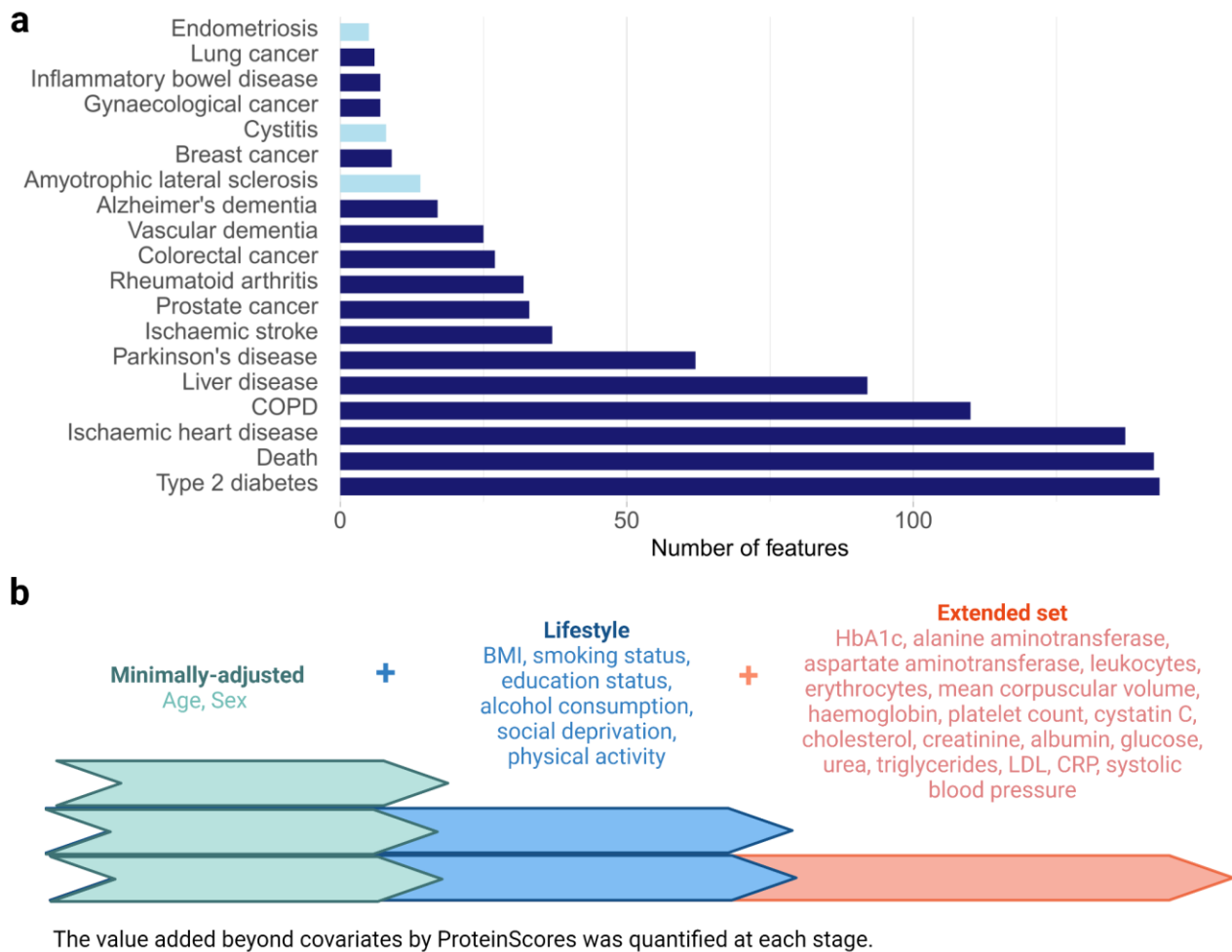
181 ProteinScores for 19 diseases that had a minimum of 150 incident cases available were trained using  
182 Cox PH elastic net regression with cross-validation in a training subset. Cumulative time-to-onset  
183 distributions for cases (**Supplementary Figs. 3-4**) indicated that amyotrophic lateral sclerosis,  
184 endometriosis and cystitis were better-suited to 5-year onset assessments (80% of cases for these traits  
185 were diagnosed by year 8 of follow-up). All remaining ProteinScores were tested in the context of 10-  
186 year onset. Performance was quantified via incremental Cox PH models in the test subset, to obtain  
187 onset probabilities for calculation of AUC and Precision Recall AUC (PRAUC) estimates (**see**  
188 **Methods**). This approach was repeated with fifty randomly sampled train and test subset combinations  
189 for each outcome to assess stability of ProteinScore performance given varied combinations of  
190 individuals in train and test sets. ProteinScores with the median difference in AUC beyond a  
191 minimally-adjusted model were selected for each outcome (**Supplementary Table 12**). Summaries of  
192 protein features selected for the 19 ProteinScores are available in **Supplementary Tables 13-14**,  
193 ranging from five features selected for endometriosis to 143 features selected for type 2 diabetes  
194 (**Fig.3a**).

## 195 **ProteinScore evaluation**

196 Selected ProteinScores were evaluated alongside various combinations of covariates to quantify the  
197 additional improvements in AUC and PRAUC achieved by each score beyond these factors (**Fig.3b**).  
198 Three increasingly complex sets of covariates were considered: 1) age and sex (where traits had not  
199 been sex-stratified), 2) further adjustment for a core set of six lifestyle and health covariates (BMI,  
200 alcohol consumption, social deprivation, educational attainment, smoking status and physical activity)  
201 and 3) further adjustment for an extended set of 18 biochemistry and physical attributes that are  
202 measurable in clinical settings. Performance when using only the ProteinScores was also considered,  
203 to ascertain whether protein information can streamline the signal offered by the set of 26 possible

204 covariates. As these covariates are sourced from a range of physical measures, clinically-used  
205 biomarker assays and self-reporting, they represent a labour and time intensive resource that is rarely  
206 collated for every individual in clinical practice. A tabular summary of both the AUC and PRAUC  
207 statistics for all covariate combinations tested, with ROC P value comparisons comparing models  
208 with/without the addition of the ProteinScore are available in **Supplementary Table 15**. Strikingly, the  
209 singular inclusion of the ProteinScores had either equal or higher performance (as measured by AUC  
210 and PRAUC) than the maximal set of 26 covariates in eight instances (type 2 diabetes, liver disease,  
211 COPD, amyotrophic lateral sclerosis, death, Alzheimer's dementia, ischaemic heart disease and  
212 Parkinson's disease). The difference in AUC resulting from the addition of the ProteinScores into the  
213 three models with increasingly complex sets of covariates are summarised in **Fig.4a**.

214 In tests for significant differences between receiver operating characteristic (ROC) curves for the three  
215 models with increasingly complex sets of covariates with/without the ProteinScores, 10 of the  
216 ProteinScores (type 2 diabetes, liver disease, COPD, lung cancer, death, ischaemic stroke, Alzheimer's  
217 and vascular dementias, ischaemic heart disease and Parkinson's disease) had ROC  $P < 0.0026$  (the  
218 Bonferroni-adjusted P-value threshold) beyond minimally-adjusted covariates. When adding  
219 ProteinScores to models that included both minimally-adjusted and lifestyle covariates, performance  
220 of lung cancer and vascular dementia ProteinScores was attenuated, leaving eight ProteinScores that  
221 had  $P < 0.0026$  in ROC model comparison tests (type 2 diabetes, liver disease, COPD, death, ischaemic  
222 stroke, Alzheimer's dementia, ischaemic heart disease and Parkinson's disease). When assessing  
223 models that further adjusted for an additional 18 clinically-measurable covariates, six of the eight  
224 ProteinScores had  $P < 0.0026$  in model comparisons with/without the ProteinScore, whereas liver disease  
225 and ischaemic stroke were attenuated by the extended covariate set. **Fig.4b** shows the breakdown of  
226 incremental model performance (by AUC) for each of these six ProteinScores: type 2 diabetes, COPD,  
227 death, Alzheimer's dementia, ischaemic heart disease and Parkinson's disease. Models that included  
228 only the ProteinScore are also presented with corresponding AUC performance.



230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

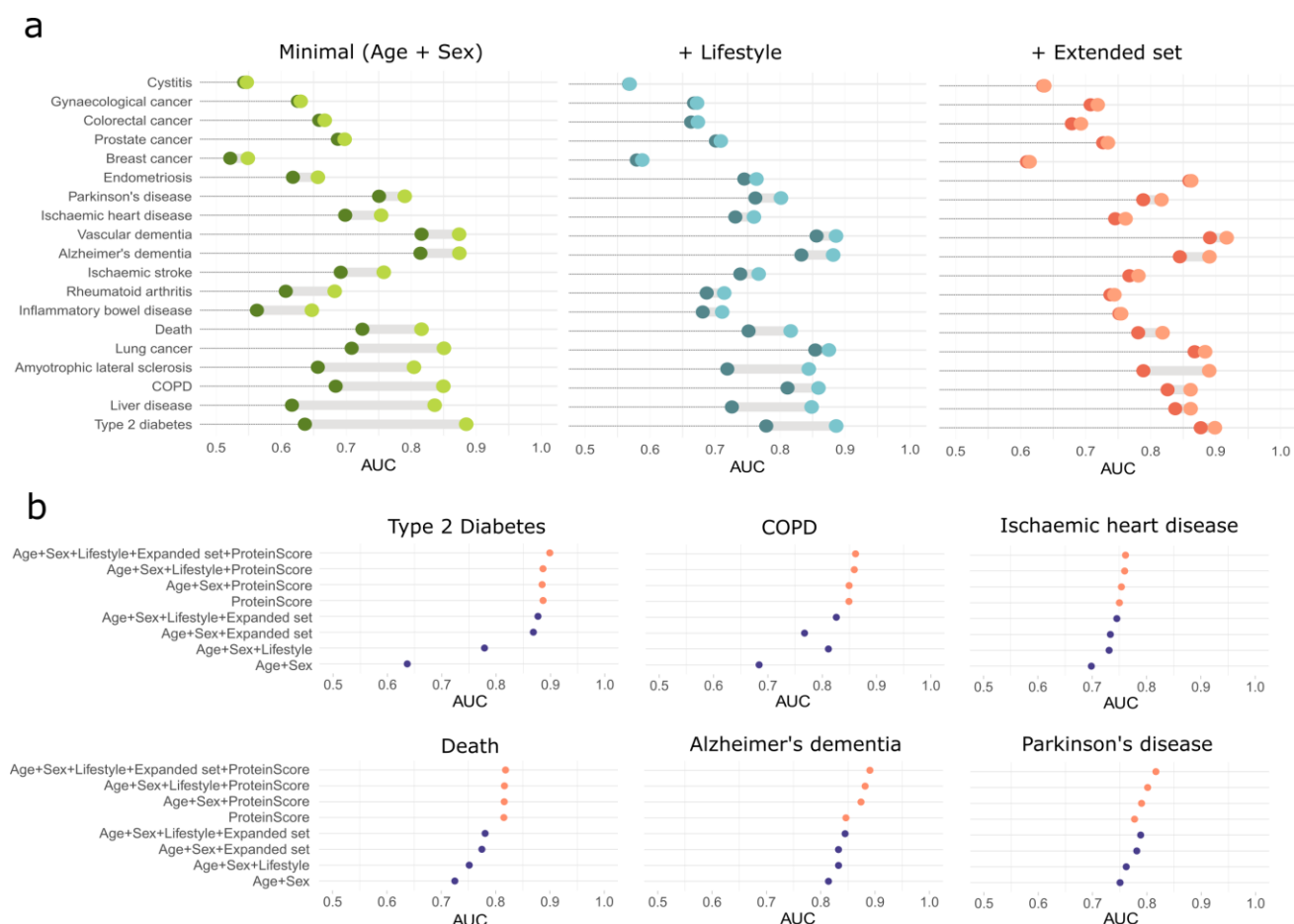
253

254

255

256

257



242

243 **Figure 4. Predictive value offered by ProteinScores for incident outcomes in the UK Biobank. a,**  
 244 **Differences in AUC resulting from the addition of the 19 ProteinScores to models with increasingly**  
 245 **extensive sets of covariates: 1) minimally-adjusted (age and sex where traits were not sex-stratified),**  
 246 **2) minimally-adjusted with the addition of a core set of six lifestyle covariates and 3) further**  
 247 **adjustment for an extended set of 18 covariates that are measured in clinical settings (physical and**  
 248 **biochemical measures). AUC plots are ordered by increasing AUC differences in the minimally-**  
 249 **adjusted models. All ProteinScore performance statistics shown correspond to 10-year onset, except**  
 250 **those for ALS, endometriosis and cystitis that were assessed for 5-year onset. b, A breakdown of the**  
 251 **AUC values achieved by different combinations of risk factors with/out the ProteinScores are shown**  
 252 **for the six incident outcomes whereby the ProteinScore contributed statistically significant beyond a**  
 253 **model including all 24 minimal, lifestyle and extended set variables (ROC P < 0.0026, the Bonferroni-**  
 254 **adjusted threshold). All six of the best-performing ProteinScores shown were assessed for 10-year**  
 255 **onset of disease. ALS: amyotrophic lateral sclerosis. COPD: chronic obstructive pulmonary disease.**

256

257

## 258 Exploration of the type 2 diabetes ProteinScore

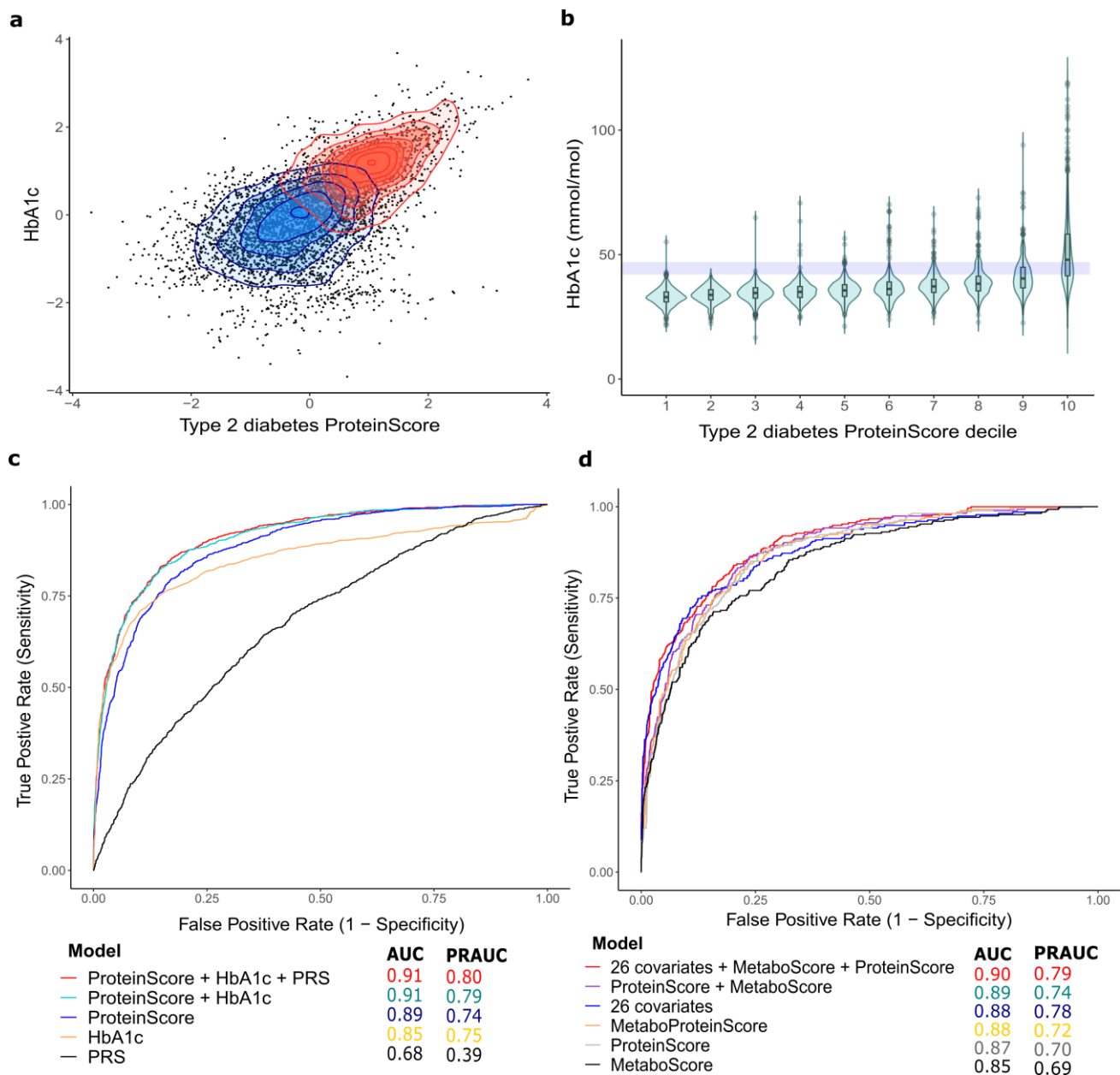
259 To highlight the value that the best-performing ProteinScores may offer, type 2 diabetes was chosen  
260 as a case study for further exploration. Performance of the ProteinScore was assessed in the context of  
261 the current clinically-used biomarker for type 2 diabetes – glycated haemoglobin (HbA1c), in addition  
262 the predictive signals offered by other omics sources (genetic and metabolomic).

263 Given that the ProteinScore for type 2 diabetes added value beyond the extended set of covariates that  
264 included the well-validated biomarker HbA1c, the performance of HbA1c and the ProteinScore was  
265 directly compared in the test sample. As polygenic risk scores are also widely used to quantify the  
266 genetic risk contribution to disease, a polygenic risk score (PRS) for type 2 diabetes was also evaluated.  
267 In the type 2 diabetes test set, 1,105 cases (with mean time to onset 5.4 years [SD 3.0]) and 3,264  
268 controls had all three measures available. HbA1c averages long-term glucose over two to three months  
269 and is widely employed clinically to monitor pre-clinical diabetes risk (42-47mmol/mol) and diagnose  
270 the disease (with two repeated measurements >48mmol/mol) <sup>29,30</sup>. The rank-base inverse normal  
271 transformed levels of the ProteinScore and HbA1c had Pearson  $r=0.50$  and discriminated incident case  
272 and control distributions similarly (**Fig. 5a**). HbA1c levels increased across ProteinScore risk deciles,  
273 with individuals in the upper deciles of the ProteinScore falling within the clinical HbA1c screening  
274 threshold (42-47mmol/mol) for diabetes (**Fig. 5b**). In incremental Cox PH models for the 10-year onset  
275 of type 2 diabetes (**Fig. 5c**) the singular use of the ProteinScore (AUC = 0.89) outperformed both  
276 HbA1c (AUC = 0.85) and the PRS (AUC = 0.68). In ROC model comparisons between HbA1c alone  
277 and HbA1c with the ProteinScore added, a statistically significant improvement due to the  
278 ProteinScore was identified (ROC P < 0.0026). When the PRS was added to this model (including  
279 HbA1c and the ProteinScore), AUC remained unchanged (0.91), whereas PRAUC improved (from  
280 0.79 to 0.80) and a significant difference due to the addition of the PRS was identified (ROC P <

281 0.0026). **Supplementary Table 16** summarises the results from these analyses, which are also  
282 presented in **Fig.5c**.

283 In a preliminary assessment, we aimed to 1) directly compare scores generated with either  
284 metabolomics-only or protein-only features and 2) assess the value added through the consideration of  
285 proteomic and metabolomic features simultaneously. The original type 2 diabetes ProteinScore  
286 populations were subset to training and testing sets that had proteomic and metabolomic measures  
287 available ( $N_{\text{cases}_{\text{train}}} = 377$ ,  $N_{\text{controls}_{\text{train}}} = 1,002$ ,  $N_{\text{cases}_{\text{test}}} = 309$ ,  $N_{\text{controls}_{\text{test}}} = 898$ ). Performance  
288 of a MetaboScore (considering metabolite features), a ProteinScore (considering protein features) and  
289 a MetaboProteinScore (combining metabolomic proteomic features) is summarised for the restricted  
290 population in **Fig. 5d**. The ProteinScore (AUC = 0.87) outperformed the MetaboScore (AUC = 0.85)  
291 The MetaboProteinScore (that considered both omics measures as potential features) had an AUC of  
292 0.88, whereas modelling the independently-trained MetaboScore and ProteinScore together resulted  
293 in an AUC of 0.89. The maximal set of 26 possible covariates (see **Fig.3b**) had an AUC of 0.88, which  
294 rose to a maximal AUC of 0.90 upon the inclusion of the MetaboScore and ProteinScore. The selected  
295 features and weights for each score (MetaboScore = 19 features, ProteinScore = 52 features and  
296 MetaboProteinScore = 37 features) are available in **Supplementary Table 17**, with full AUC and  
297 PRAUC statistics available in **Supplementary Table 18**.





298

299 **Figure 5. Exploration of the type 2 diabetes ProteinScore.** **a**, Case (red) and control (blue)  
300 discrimination for HbA1c and the type 2 diabetes ProteinScore in the test set (1,105 cases, 3,264  
301 controls, mean time to case onset 5.4 years [SD 3.0]). Both markers were rank-based inverse  
302 normalised and scaled to have a mean of 0 and standard deviation of 1. **b**, HbA1c (mmol/mol) per  
303 decile of the type 2 diabetes ProteinScore in the test set. The shaded rectangle indicates the type 2  
304 diabetes HbA1c screening threshold (42-47 mmol/mol). **c**, ROC curves for incremental 10-year onset  
305 models incorporating HbA1c, the type 2 diabetes ProteinScore and a polygenic risk score (PRS) for  
306 type 2 diabetes individually and concurrently. **d**, ROC curves for 10-year onset scores developed in  
307 the subsets of the type 2 diabetes train and test populations that had metabolomics and proteomics  
308 available ( $N_{\text{case}_{\text{Strain}}} = 377$ ,  $N_{\text{control}_{\text{Strain}}} = 1,002$ ,  $N_{\text{case}_{\text{Test}}} = 309$ ,  $N_{\text{control}_{\text{Test}}} = 898$ ). A  
309 Metabolomic score (MetaboScore), ProteinScore, and a joint omics score (MetaboProteinScore) are  
310 modelled individually and concurrently and benchmarked against 26 covariates (age, sex, six lifestyle  
311 factors and the extended set of 18 clinically-relevant covariates).

## 312 Discussion

313 Identifying individuals at risk of a future disease event or death is a priority for prevention-based  
314 medicine during ageing <sup>31</sup>. We report 3,201 associations between 961 circulating proteins and 21  
315 incident outcomes, identifying proteins indicative of multimorbidity. ProteinScores for incident type  
316 2 diabetes, COPD, ischaemic heart disease, Alzheimer’s dementia, Parkinson’s disease and death  
317 demonstrated value beyond a comprehensive set of 26 covariates, offering comparable performance  
318 and minimising the need for extensive recording of lifestyle factors, physical measures and biomarker  
319 assays. Exploration of the type 2 diabetes ProteinScore suggested that while protein information  
320 captures much of the predictive signal, augmenting traditional risk factors with proteomic,  
321 metabolomic and genetic data types may further hone risk classification.

322 The breadth of electronic health data linkage and protein data available in UK Biobank provides a  
323 unique resource for profiling early molecular signatures of age-related disease. This study  
324 demonstrates that for certain diseases, subsets of relatively few circulating proteins can add predictive  
325 value, up to a decade prior to formal diagnoses. As available cases increase, it is likely that the  
326 performance of ProteinScores will be enhanced. Nonetheless, for the best-performing six  
327 ProteinScores, modelling the ProteinScore in isolation resulted in equal or higher AUCs than models  
328 with extensive covariate adjustments. This suggests that ProteinScores for such traits absorb a large  
329 proportion – if not all – of the signal and may offer a streamlined set of metrics to proxy for an  
330 individual’s health status. This often-enhanced predictive quality of the scores presents an exciting  
331 opportunity to reconsider how best to formulate (and maintain) modern clinical prediction models.  
332 This is an important consideration given that self-reported measures are known to be variable in  
333 accuracy and are often misreported <sup>32</sup>. Additionally, while much interest is currently devoted to  
334 employing PRS for disease prediction, they neglect environmental components of disease risk and may  
335 therefore be limited in the context of complex age-related disease <sup>33,34</sup>. Our ProteinScore for type 2

336 diabetes outperformed the PRS, which is likely due to proteins representing an interface that captures  
337 genetic, environmental and lifestyle contributions to disease risk. The improvement in AUC resulting  
338 from concurrent modelling of HbA1c and the type 2 diabetes ProteinScore suggests that the latter may  
339 provide additional predictive value. Similarly, in a subset of the population with metabolomics  
340 measures, the type 2 diabetes ProteinScore outperformed the MetaboScore, with an additive signal  
341 achieved by modelling both scores. While diabetes is typically considered to be a metabolic disease,  
342 the breadth of coverage (249 metabolites, versus 1,468 protein measures) may limit the metabolic score  
343 performance. ProteinScores for multiple diseases within the same individuals may facilitate an  
344 improved understanding of multimorbidity. For example, if an individual falls within the top 5% of  
345 the ProteinScore distributions for type 2 diabetes and Alzheimer's dementia, this information may  
346 enhance personalised intervention plans. The ProteinScore performance for Alzheimer's dementia was  
347 also largely unchanged upon addition of additional covariates. As therapeutic interventions for  
348 neurodegenerative diseases have greater efficacy when implemented earlier in the disease pathogenesis  
349 <sup>35-37</sup>, the ProteinScore for Alzheimer's dementia may hone trial recruitment.

350 The method for ProteinScore generation selects proteins that, in combination, are predictive of  
351 outcomes, but these do not necessarily represent the most probable drivers of disease. It is likely that  
352 a subset of the 3,201 individual protein-disease associations we report represent direct mediators of  
353 disease. The goal of this work was to identify early markers that associate with incident disease and  
354 the markers we identify are therefore useful for risk stratification purposes (even if they are indicative  
355 of underlying morbidities at baseline, or do not represent causal mediators). To delineate the markers  
356 that may be direct mediators of disease, we encourage further exploration of through techniques such  
357 as Mendelian randomisation and colocalisation. Similarly, further modelling that takes into account  
358 multimorbidity trajectories over the lifecourse would also aid in understanding the role of prevalent  
359 diseases and medication use on future disease risk. The largest number of associations and strongest  
360 effect sizes (by magnitude of the absolute log of the hazard ratio) were observed for liver disease in

361 individual Cox PH analyses. For neurological diseases and cancers, where fewer associations were  
362 identified, it is possible that the blood is less able to capture the full spectrum of disease pathogenesis,  
363 which may be localised to distal tissues. Similarly, the panel of proteins available may reflect certain  
364 diseases better than others.

365 All 54 proteins that were associated with eight or more morbidities had associations with hazard ratios  
366 greater than 1, indicating that elevated levels of these proteins may serve as early warning signatures  
367 of disease onset. Elevated growth differentiation factor 15 (GDF15), Interleukin-6 (IL6) and  
368 plasminogen activator urokinase receptor (PLAUR) had the largest number of associations with  
369 incident diseases. This result is in concordance with previous screening of the circulating proteome  
370 against multimorbidity and mortality, which identified GDF15 as the top marker of future  
371 multimorbidity from 1,301 plasma proteins tested <sup>38,39</sup>. Further evidence supporting GDF15 as a  
372 marker of multiple outcomes including heart disease, type 2 diabetes, stroke, dementia and death has  
373 been reported <sup>39-45</sup>. IL6 mediates chronic, low-grade inflammation, is a key biomarker of ageing <sup>46</sup> and  
374 anti-IL6 therapeutics have been developed for a range of inflammation-associated diseases <sup>47,48</sup>. While  
375 less-extensive evidence exists supporting PLAUR as a biomarker of multiple morbidities, it was  
376 associated with incident cancer, cardiovascular disease, diabetes and mortality in previous Cox PH  
377 analyses <sup>49</sup>. Similarly, increased levels of neurofilament light (NEFL) were associated with higher  
378 incidence of multiple neurological traits (Parkinson's disease, Alzheimer's dementia, multiple  
379 sclerosis, amyotrophic lateral sclerosis and ischaemic stroke). These diseases are hallmarked by neuron  
380 degradation and NEFL may therefore be a consequential marker that is released into the blood upon  
381 breakdown of synapses <sup>50,51</sup>. NEFL was also associated with liver disease, COPD and ischaemic heart  
382 disease, which may reflect the presence of underlying synaptic and neuronal dysfunction, or the  
383 presence of comorbidities in individuals with these diagnoses.

384 Across the 16-year window of follow-up in individual Cox PH models, a subset of associations  
385 violated the Cox PH assumption at the local (protein) level. Our Shiny app  
386 <https://protein-disease-ukb.optima-health.technology> [Username: ukb\_diseases, Password:  
387 UKBshinyapp] provides visualisations for sensitivity analyses run across cases over successive years  
388 of case follow up, allowing for interrogation of the stability of individual protein-disease relationships.  
389 This information on near-term versus long-term case follow-up is often of importance to clinicians and  
390 patients for behaviour change and intervention strategies. The Shiny app also visualises the 3,201  
391 fully-adjusted associations in a network view, allowing users to view overlapping signatures between  
392 multiple proteins and the onset of multiple diseases.

393 This study has several limitations. First, a subset of 6,385 individuals in the UKB-PPP sample were  
394 selected by consortium members for enrichment of certain diagnoses and this non-random selection  
395 can introduce biases. Second, as UK Biobank currently represents the largest population with  
396 comprehensive Olink proteomics and electronic health data linkage, it was not possible to source an  
397 external test set for the ProteinScores. Third, variation in protein analyte levels across measurement  
398 technologies has been reported <sup>52</sup>. Results should therefore be corroborated across panels in future.  
399 Fourth, the protein measured were recorded in relative scale, rather than absolute quantification. This  
400 limits direct translation of the ProteinScores for direct prediction in new populations. However, the  
401 early markers of incident disease that are identified may still replicate when datasets become available  
402 to facilitate replication analyses. Fifth, the UK Biobank population is largely comprised of individuals  
403 with European ancestry and a restricted age range (40-71 years, with a mean of 57 years); future studies  
404 in equally well-characterized cohorts will be needed to assess how well ProteinScores translate to other  
405 populations and ethnicities. Sixth, non-linear trajectories of blood-based protein signatures are known  
406 to exist across the life course in the context of ageing <sup>53</sup>. These factors should be considered in disease-  
407 specific analyses in future. Seventh, death was treated as a censoring event; competing risks and multi-  
408 state modelling approaches may be used for disease-protein associations in future to resolve the impact

409 of death as a competing risk for disease onset. Finally, although a comprehensive set of major age-  
410 related morbidities were studied, many diseases were not included in this work. Continued linkage and  
411 proteomic sampling will expand the applications of ProteinScores to further diseases.

412 In conclusion, this study quantified circulating proteome signatures that are reflective of multiple  
413 individual disease states across mid-to-later life. ProteinScores for the incidence of six incident  
414 outcomes significantly improved AUCs for 10-year onset beyond 26 demographic, lifestyle and  
415 clinically-relevant covariates. The type 2 diabetes ProteinScore offered additional value beyond  
416 HbA1c, a PRS and a metabolomic score. A total of 3,123 individual protein-disease associations were  
417 also profiled across the 16-year follow-up period, identifying candidate targets for multimorbidity  
418 prevention. These data suggest that proteomic features are powerful tools for honing risk stratification.

## 419 **Methods**

### 420 **The UK Biobank sample population**

421 UK Biobank (UKB) is a population-based cohort of around 500,000 individuals aged between 40-69  
422 years that were recruited between 2006 and 2010. Genome-wide genotyping, exome sequencing,  
423 electronic health record linkage, whole-body magnetic resonance imaging, blood and urine biomarkers  
424 and physical and anthropometric measurements are available. More information regarding the full  
425 measurements can be found at: <https://biobank.ndph.ox.ac.uk/showcase/>. The UK Biobank Pharma  
426 Proteomics Project (UKB-PPP) is a precompetitive consortium of 13 biopharmaceutical companies  
427 funding the generation of blood-based proteomic data from UKB volunteer samples.

### 428 **Proteomics in the UK Biobank**

429 The UKB-PPP sample includes 54,306 UKB participants and 1,474 protein analytes measured across  
430 four Olink panels (Cardiometabolic, Inflammation, Neurology and Oncology: annotation information  
431 provided in **Supplementary Table 1**)<sup>25</sup>. A randomised subset of 46,673 individuals were selected

432 from baseline UKB, with 6,385 individuals selected by the UKB-PPP consortium members and 1,268  
433 individuals included that participated in a COVID-19 study. The randomised samples have been shown  
434 to be highly representative of the wider UKB population, whereas the consortium-selected individuals  
435 were enriched for 122 diseases<sup>25</sup>. Details on sample selection for UKB-PPP, in addition to processing  
436 and quality control information for the Olink assay are provided in **Supplementary Information**. Of  
437 54,309 individuals that had protein data measured, there were 52,744 that were available after quality  
438 control exclusions with 1,474 Olink protein analytes measured (annotations in **Supplementary Table**  
439 **1**)<sup>25</sup>. The sample is predominantly white/European (93%), but also has individuals with black/black  
440 British, Asian/Asian British, Chinese, mixed, other and missing ethnic backgrounds (7%).

441 **Supplementary Fig. 1** summarises the processing steps applied to this dataset to derive a complete  
442 set of measurements for use. Briefly, of 107,161 related pairs of individuals (calculated through kinship  
443 coefficients  $> 0$  across the full UKB cohort), 1,276 pairs were present in the 52,744 individuals. After  
444 exclusion of 104 individuals in multiple related pairs, in addition to one individual randomly selected  
445 from each of the remaining pairs, there were 51,562 individuals. A further 3,962 individuals were  
446 excluded due to having  $>10\%$  missing protein measurements. Four proteins that had  $>10\%$  missing  
447 measurements (CTSS.P25774.OID21056.v1 and NPM1.P06748.OID20961.v1 from the neurology  
448 panel, PCOLCE.Q15113.OID20384.v1 from the cardiometabolic panel and  
449 TACSTD2.P09758.OID21447.v1 from the oncology panel) were then excluded. The remaining 1% of  
450 missing protein measurements were imputed by K-nearest neighbour ( $k=10$ ) imputation using the  
451 impute R package (Version 1.60.0)<sup>54</sup>. The final dataset consisted of 47,600 individuals and 1,468  
452 protein analytes. Assessments of protein batch, study centre and genetic principal components  
453 suggested that these factors had minimal effects on protein levels (lowest correlation between protein  
454 levels and residuals of 0.94) (**Supplementary Information**). Therefore, protein levels were not  
455 adjusted for these factors.

## 456 **Phenotypes in the UK Biobank**

457 Demographic and phenotypic information for the 47,600 individuals with complete protein data for  
458 1,468 analytes are available in **Supplementary Table 2**. Lifestyle covariates included: BMI (weight  
459 in kilograms divided by height in metres squared), alcohol intake frequency (1 = Daily or almost daily,  
460 2 = Three-Four times a week, 3 = Once or twice a week, 4 = One-Three times a month, 5 = Special  
461 occasions only, 6 = Never), the Townsend index of deprivation (higher score representing greater  
462 levels of deprivation) and smoking status (0 = Never, 1 = Previous, 2 = Current), physical activity (0  
463 = between 0-2 days/week of moderate physical activity, 1 = between 3-4 days/week of moderate  
464 physical activity, 2 = between 5-7 days/week of moderate physical activity) and education status (1 =  
465 college/university educated, 0 = all other education). Of the 47,600 individuals with complete protein  
466 data, there were 52, 52, 236, 56 and 59 missing entries for alcohol, smoking, BMI, physical activity  
467 and deprivation, respectively. No imputation of missing data was performed for the inclusion of these  
468 variables in individual Cox PH analyses. There were an additional 2,556, 188 and 59 individuals that  
469 answered ‘prefer not to answer’ and were excluded from physical activity, smoking and alcohol  
470 variables, respectively.

## 471 **Electronic health data linkage in the UK Biobank**

472 Electronic health linkage to NHS records was used to collate incident diagnoses. Death information  
473 was sourced from the death registry data available through the UK Biobank. Cancer outcomes were  
474 sourced from the cancer registry (ICD codes), whereas non-cancer diseases were sourced from first  
475 occurrence traits available in the UK Biobank. The first occurrence traits integrate GP (read2/3), ICD  
476 (9/10) with self-report and ICD codes present on the death registry to identify the earliest date of  
477 diagnoses. These data sources are linked to 3-digit ICD trait codes. A summary of codes used to extract  
478 each of the outcomes included in the present study are detailed in **Supplementary Information**. The  
479 following 23 diseases were included: liver disease, systemic lupus erythematosus, type 2 diabetes,



480 amyotrophic lateral sclerosis, Alzheimer’s dementia, endometriosis, chronic obstructive pulmonary  
481 disease (COPD), inflammatory bowel disease, rheumatoid arthritis, ischaemic stroke, Parkinson’s  
482 disease, vascular dementia, ischaemic heart disease, major depressive disorder, schizophrenia,  
483 multiple sclerosis, cystitis and lung, prostate, breast, gynaecological, brain/CNS and colorectal  
484 cancers. These represent a selection of leading age-related causes of morbidity, mortality and  
485 disability. In all analyses involving sex-specific diseases, the population was stratified to males or  
486 females and sex was not included as a covariate in incremental Cox PH assessments. Traits that were  
487 stratified included gynaecological cancer, breast cancer, endometriosis and cystitis (all female-  
488 stratified) and prostate cancer (male-stratified).

#### 489 **Incident disease calculation in the UK Biobank**

490 Dates of diagnoses for each disease were ascertained through electronic health linkage. Using the date  
491 of baseline appointment, time-to-first-onset for each diagnoses in years was calculated. Time-to-onset  
492 for controls was defined as the time from baseline to censoring date (**Supplementary Information**).  
493 Death was treated as a censoring event. Time-to-censor date was calculated for the controls that  
494 remained alive, whereas if a control individual had died during follow-up time-to-death was taken  
495 forward for Cox PH models. Any cases that were prevalent at baseline were excluded. Alzheimer’s  
496 and vascular dementias were restricted to age at onset (or censoring) of 65 years or older in all analyses.  
497 Sex-specific traits were stratified across all analyses.

#### 498 **Individual Cox proportional hazards analyses**

499 Cox proportional hazards models were run between each protein and each incident disease using the  
500 ‘survival’ package (Version 3.4-0)<sup>55</sup> in R (Version 4.2.0)<sup>56</sup>. Protein levels were rank-based inverse  
501 normalised and scaled to have a mean of 0 and standard deviation of 1 prior to analyses. Minimally-  
502 adjusted Cox PH models for sex-stratified traits included age at baseline as a covariate, whereas the  
503 remaining models adjusted for age and sex. Lifestyle-adjusted models further controlled for education

504 status, BMI, smoking status, social deprivation rank, physical activity and alcohol intake frequency. A  
505 Bonferroni-adjusted P-value threshold for multiple testing based on the 678 components that explained  
506 90% of the cumulative variance in the 1,468 protein analyte levels (**Supplementary Table 3**) and 24  
507 outcomes tested was applied across all Cox PH models ( $P < 0.05/(678 \times 24) = 3.1 \times 10^{-6}$  used as the  
508 Bonferroni-adjusted P-value threshold). Proportional hazards assumptions were checked through  
509 examination of protein-level Schoenfeld residuals.

510 A sensitivity analysis was performed for each of the 35,232 fully-adjusted associations tested,  
511 restricting cases to successive years of follow-up. These sensitivity analyses were visualised using the  
512 Shiny package (Version 1.7.3)<sup>57</sup> in R. The magnitude of change in hazard ratios for individual  
513 associations can be examined by year of case follow-up to assess consistency of effect sizes. Whether  
514 marker associations are stronger or weaker when restricting to cases occurring in the near-term (1-5  
515 years of follow-up) can also be examined. A network visualisation was also created within the Shiny  
516 interface to highlight the fully-adjusted associations that had  $P < 3.1 \times 10^{-6}$  using networkD3 (Version  
517 3.0.4)<sup>58</sup> and igraph (Version 1.3.5)<sup>59</sup> R packages. To further verify the markers of multiple morbidities  
518 identified in individual Cox PH analyses, logistic regression models were also run between each of the  
519 1,468 protein analyte levels and multimorbidity status (defined as 1,454 individuals that received 3 or  
520 more of the 23 disease diagnoses over the 16-year follow-up period). A sensitivity analyses was also  
521 run for ischaemic heart disease associations with/out adjustment for blood-pressure lowering  
522 medication reported at baseline in a subset of individuals (35,073 of 47,600) that had medication  
523 information available. **Supplementary Information** provides details on the classification of  
524 medications as per the anatomical therapeutic chemical (ATC) classification categories. A total of  
525 14,074 individuals (of the 35,073) indicated they were taking one or more of the above blood-pressure  
526 lowering medications at baseline. This was treated as a binary variable and the comparison with/out  
527 adjustment for this variable was performed for ischaemic heart disease Cox PH associations in the

528 subset of 35,073 individuals. Adjustments for age, sex and six lifestyle factors were included in both  
529 sets of analyses, with 2,456 cases, 27,468 controls.

### 530 **ProteinScore development**

531 MethylPipeR<sup>60</sup> is an R package with accompanying user interface that we have previously developed  
532 for systematic and reproducible development of incident disease predictors. Using MethylPipeR,  
533 ProteinScores that considered 1,468 Olink protein levels were trained using Cox PH elastic net  
534 regression via the R package Glmnet (Version 4.1-4)<sup>61</sup>. Penalised regression minimises overfitting by  
535 the use of a regularisation penalty and the best shrinkage parameter ( $\lambda$ ) was chosen by cross-fold  
536 validation with alpha fixed to 0.5. Of the 24 outcomes featured in the individual Cox PH analyses, 19  
537 that had a minimum case count of 150 were selected for ProteinScore development. The chosen  
538 strategy for ProteinScore development included training ProteinScores for each trait across fifty  
539 randomised iterations (with each iteration including a different combination of cases and controls in  
540 train and test sets). This strategy quantifies the stability of the ProteinScore performance, which is  
541 critical given that unobserved confounders that may be enriched during random selection of individuals  
542 from the wider population. The ProteinScore training strategy is summarised in **Supplementary Fig.**  
543 **5**. Briefly, 50 iterations of each ProteinScore were performed that randomised sample selection by 50  
544 randomly sampled seeds (values between 1 and 5000). For each iteration, cases and controls were  
545 randomly split into 50% groups for training and testing. From the 50% training control population, a  
546 subset of controls were then randomly sampled to give a case:control ratio of 1:3 in order to balance  
547 the datasets. For traits with over 1000 cases in training samples 10 folds were used. For traits with  
548 between 500 and 1000 cases in training, five folds were used. Three folds were used when there were  
549 fewer than 500 cases in the training sample. Protein levels were rank-based inverse normalised and  
550 scaled to have a mean of 0 and standard deviation of 1 in the training set. The linear combination of  
551 weighting coefficients for selected protein features from cross-validation within the folds of the

552 training set were then used to generate a ProteinScore for each individual in the test samples. Of the  
553 50 training iterations tested, models that had no features selected were documented (**Supplementary**  
554 **Table 12**).

### 555 **Assessment of ProteinScore performance**

556 Cumulative time-to-onset distributions for cases (**Supplementary Figs. 3-4**) indicated that  
557 amyotrophic lateral sclerosis, endometriosis and cystitis were better-suited to 5-year onset assessments  
558 in the test sample (80% of cases were diagnosed at 8-years post-baseline). All remaining ProteinScores  
559 were tested in the context of 10-year onset (80% of cases were not diagnosed 8-years post-baseline).  
560 Across the 50 ProteinScore iterations for each trait, 50% of cases and controls that were not randomly  
561 selected for training were reserved for testing. For a visualisation of the test set sampling and  
562 assessment strategy, see **Supplementary Fig. 5**. In the test set, cases that had time-to-event up to or  
563 including the 5-year or 10-year thresholds used for onset prediction were selected, while cases beyond  
564 the threshold were placed with the control population, which was then randomly sampled in a 1:3 ratio.  
565 Weighting coefficients for features selected during ProteinScore training were used to project scores  
566 into the test sample. Incremental Cox PH models were run in the test sample to obtain cumulative  
567 baseline hazard and onset probabilities, which were used to derive AUC and PRAUC estimates. The  
568 test set sampling strategy ensured that while the majority of cases occurred up to the onset threshold,  
569 there were a small proportion (~3%) of cases included in Cox PH models with onset times after the  
570 10- or 5-year threshold, to simulate a real-world scenario for risk stratification. If cases fell beyond the  
571 5-year or 10-year threshold for onset, they were recoded as controls in the AUC calculation.  
572 Cumulative baseline hazard probabilities were calculated using the Breslow estimator available in the  
573 ‘gbm’ R package (Version 2.1.8.1)<sup>62</sup>. Survival probabilities were then generated through taking the  
574 exponential of the negative cumulative baseline hazard at 5 or 10 years to the power of the Cox PH  
575 prediction probabilities. ProteinScore onset probabilities were calculated as one minus these survival

576 probabilities. AUC, PRAUC and ROC statistics were extracted for the survival probabilities using the  
577 calibration function from the ‘caret’ R package (Version 6.0-94)<sup>63</sup> and the evalmod function from the  
578 ‘MLmetrics’ R package (Version 1.1.1)<sup>64</sup>.

579 ProteinScores that yielded the median incremental difference to the AUC of a minimally-adjusted  
580 model (adjusting for age- or age- and sex) were selected from the fifty possible ProteinScores for each  
581 trait. If no features were selected during training, models were weighted as performance of 0 in the  
582 median model selection. In some instances, features were selected during training and incremental Cox  
583 PH models were run successfully, but the random sampling of the test set did not include a case with  
584 time-to-event at or after the 5-year or 10-year onset threshold. Therefore, these models were excluded  
585 as cumulative baseline hazard distributions did not reach the onset threshold and could not be extracted  
586 for AUC and PRAUC calculations. The number of models, with minimum and maximum performance  
587 was documented (**Supplementary Table 12**). Taking this approach mitigated against the presence of  
588 extreme case:control profiles driving ProteinScore performance and minimised the possibility of bias  
589 being introduced by selecting train and test samples based on matching for specific population  
590 characteristics.

591 Selected ProteinScores for each trait were then evaluated to quantify the additional value (in terms of  
592 increases in AUC and PRAUC) that resulted from the addition of ProteinScores. Minimally-adjusted  
593 models included age and sex (if traits were not sex-stratified). Lifestyle-adjusted models then further  
594 accounted for common lifestyle covariates (education status, BMI, smoking status, social deprivation  
595 rank, physical activity and alcohol intake frequency). Finally, models included covariates from the  
596 minimally-adjusted, lifestyle-adjusted and an extended set of clinically-measured variables were then  
597 assessed (**see Fig.3b**). In each case, the difference in AUC and PRAUC resulting from the addition of  
598 the ProteinScore was reported. ROC P-value tests were used to ascertain whether the improvements  
599 offered by selected ProteinScores for each outcome were statistically significant, beyond each set of

600 increasingly saturated covariates. A Bonferroni-adjusted P-value threshold for ROC P tests was used  
601 based on the 19 ProteinScore traits ( $P < 0.05/19 = 0.0026$ ). The ‘precrec’R package (Version 0.12.9)  
602 <sup>65</sup> was used to generate ROC and Precision-Recall curves for each ProteinScore. A series of models  
603 that included only the ProteinScore were also considered for each outcome, to quantify whether protein  
604 data alone could absorb much of the predictive performance achieved by the covariates.

605 A set of 26 possible covariates used across the minimally-adjusted, lifestyle-adjusted and extended set  
606 analyses were assessed for missingness, imputed (where missingness was  $< 10\%$ ) and utilised in  
607 ProteinScore evaluation as a maximal, extended set of covariates. Further details of variable selection  
608 and preparation are supplied in **Supplementary Information**. Additional covariates (considered in  
609 addition to age, sex and six lifestyle traits that were used in individual Cox PH analyses) included:  
610 leukocyte counts ( $10^9$  cells/Litre), erythrocyte counts ( $10^{12}$  cells/Litre), haemoglobin concentration  
611 (grams/decilitre), mean corpuscular volume (femtolitres), platelet count ( $10^9$  cells/Litre), cystatin C  
612 (mg/L), cholesterol (mmol/L), alanine aminotransferase (U/L), creatinine (umol/L), urea (mmol/L),  
613 triglycerides (mmol/L), LDL (mmol/L), CRP (mg/L), aspartate aminotransferase (U/L), glycated  
614 haemoglobin – HbA1c – (mmol/mol), albumin (g/L), glucose (mmol/L) and systolic blood pressure  
615 (mmHg). After covariate processing steps were complete, a population of 43,437 individuals was  
616 available with complete information for ProteinScore testing. Phenotypic summaries of the additional  
617 covariates for this population are summarised in **Supplementary Table 2**.

## 618 **Further assessment of the type 2 diabetes ProteinScore**

619 Glycated Haemoglobin (HbA1c) is a blood-based measure of chronic glycemia that is highly predictive  
620 of type 2 diabetes events and is recommended as a test of choice for the monitoring and diagnosis of  
621 type 2 diabetes <sup>29,30</sup>. HbA1c (mmol/mol) measurements (fieldID 30750) and the type 2 diabetes  
622 polygenic risk score (PRS) available in UK Biobank (fieldID 26285) were extracted. A contour plot  
623 showing both variables grouped by those who went on to be diagnosed with type 2 diabetes over a 10-

624 year period was created. HbA1c levels were also plotted against ProteinScore risk deciles. HbA1c and  
625 the ProteinScore levels were rank-based inverse normalised and assessed individually and concurrently  
626 in incremental models for 10-year onset of type 2 diabetes in the ProteinScore test set. A Pearson  
627 correlation coefficient ( $r$ ) between the transformed HbA1c and ProteinScore levels was calculated.  
628 The 10-year incremental Cox PH models were used to derive onset probabilities for calculation of  
629 AUCs and PRAUCs after adding the ProteinScore to models adjusting for HbA1c and the type 2  
630 diabetes PRS. Model comparisons were used (test of the difference in ROC curves) to quantify the  
631 value added by the ProteinScore beyond the PRS and HbA1c.

632 Metabolomics measures were available in 12,050 of the 47,600 individuals with proteomic data  
633 included in the study (see **Supplementary Information** for details on data preparation). Type 2  
634 diabetes was chosen as a case study for further exploration, as it is typically well-reflected by  
635 circulating levels of both protein and metabolomic markers and is considered as a metabolic disease.  
636 The train and test sets used to develop the main type 2 diabetes ProteinScore were subset to those with  
637 metabolomics available ( $N_{\text{cases}_{\text{Strain}}} = 377$ ,  $N_{\text{controls}_{\text{Strain}}} = 1,002$ ,  $N_{\text{cases}_{\text{Test}}} = 309$ ,  $N_{\text{controls}_{\text{Test}}} =$   
638  $898$ ). Scores that considered only metabolomic features (MetaboScore), only proteomic features  
639 (ProteinScore) and joint omics features (MetaboProteinScore) were trained and tested in these  
640 populations. There were 249 metabolite levels and 1,468 protein levels considered as potentially-  
641 informative features. Performance was evaluated for 10-year onset of type 2 diabetes in the test sample,  
642 modelling the scores individually, concurrently and benchmarking them against the maximal set of 26  
643 possible covariates (see **Fig.3**).

#### 644 **Ethics declarations**

645 All participants provided informed consent. This research has been conducted using the UK Biobank  
646 Resource under approved application numbers 65851, 20361, 26041, 44257, 53639, 69804.

647 **Data availability**

648 Datasets generated in this study are made available in **Supplementary Tables**. Proteomics data is  
649 available in the UK Biobank under Category 1838  
650 at: <https://biobank.ndph.ox.ac.uk/ukb/label.cgi?id=1838>.

651 **Code availability**

652 Code is available with open access at the following Github repository:

653 [https://github.com/DanniGadd/Blood\\_protein\\_levels\\_and\\_incident\\_disease\\_UK\\_Biobank](https://github.com/DanniGadd/Blood_protein_levels_and_incident_disease_UK_Biobank)

654 **Acknowledgements**

655 **This research was funded in whole, or in part, by the Wellcome Trust [108890/Z/15/Z]. For the**  
656 **purpose of open access, the author has applied a CC BY public copyright licence to any Author**  
657 **Accepted Manuscript version arising from this submission.**

658 We thank the participants, contributors, and researchers of UK Biobank for making data available for  
659 this study – with special thanks to Lauren Carson, John Busby, Naomi Allen and Rory Collins for  
660 making the study possible. We are grateful to the research & development leadership teams at the  
661 thirteen participating UKB-PPP member companies (Alnylam Pharmaceuticals, Amgen, AstraZeneca,  
662 Biogen, Bristol-Myers Squibb, Calico, Genentech, Glaxo Smith Klein, Janssen Pharmaceuticals, Novo  
663 Nordisk, Pfizer, Regeneron, and Takeda) for funding the study. We thank the Legal and Business  
664 Development teams at each company for overseeing the contracting of this complex, precompetitive  
665 collaboration – with particular thanks to Erica Olson of Amgen, Andrew Walsh of GSK, and Fiona  
666 Middleton of AstraZeneca. The Biogen team is especially thankful to Helen McLaughlin for her  
667 project management support. Finally, we thank the team at Olink Proteomics (Philippa Pettingell, Klev  
668 Diamanti, Cindy Lawley, Linda Jung, Sara Ghalib, Ida Grundberg and Jon Heimer) for their consistent



669 logistic support throughout the project – with special thanks to Evan Mills for co-championing the  
670 project and leading internal activities at Olink.

671 R.E.M. is supported by Alzheimer’s Society major project grant AS-PG-19b-010. R.F.H is supported  
672 by a MRC IEU Fellowship. D.A.G. is supported by the Wellcome Trust Translational Neuroscience  
673 programme [108890/Z/15/Z].

#### 674 **Author contributions**

675 D.A.G., R.F.H., R.E.M., B.S., C.F., and Z.K., conceptualised the study design and consulted on  
676 methods and results. D.A.G., carried out all analyses. D.A.G., R.F.H., B.B.S., and R.E.M., drafted the  
677 article. R.A., and J.G., conducted preliminary analyses. T.L., and K.F., performed quality control on  
678 the proteomics dataset. Y.C., was consulted on methodology. M.D. contributed to Shiny app  
679 integration of results. All authors reviewed and approved of the manuscript.

#### 680 **Competing interests**

681 B.B.S., R.A., J.G., T.L., K.F., and H.R., are employed by Biogen. C.N.F., Z.K., D.A.G., M.D., and  
682 T.M., are employed by Optima partners. D.A.G., R.F.H., and R.E.M., have received consultancy fees  
683 from Optima Partners. R.E.M. is an advisor to the Epigenetic Clock Development Foundation. R.F.H.,  
684 has received consultant fees from Illumina. All other authors declare no competing interests.

#### 685 **Materials and correspondence**

686 Correspondence and material requests should be sent to Dr Benjamin Sun at  
687 [benjamin.sun@biogen.com](mailto:benjamin.sun@biogen.com) or Prof Riccardo Marioni at [riccardo.marioni@ed.ac.uk](mailto:riccardo.marioni@ed.ac.uk).

## References

1. Fountzilias, E., Tsimberidou, A. M., Vo, H. H. & Kurzrock, R. Clinical trial design in the era of precision medicine. *Genome Med.* **14**, 101 (2022).
2. Duarte, T. T. & Spencer, C. T. Personalized Proteomics: The Future of Precision Medicine. *Proteomes* **4**, 29 (2016).
3. Al-Nesf, M. A. Y. *et al.* Prognostic tools and candidate drugs based on plasma proteomics of patients with severe COVID-19 complications. *Nat. Commun.* **13**, 946 (2022).
4. Yao, C. *et al.* Genome-wide mapping of plasma protein QTLs identifies putatively causal genes and pathways for cardiovascular disease. *Nat. Commun.* **9**, 3268 (2018).
5. Ferkingstad, E. *et al.* Large-scale integration of the plasma proteome with genetics and disease. *Nat. Genet.* **2021 5312** **53**, 1712–1721 (2021).
6. Pietzner, M. *et al.* Mapping the proteo-genomic convergence of human diseases. *Science* eabj1541 (2021).
7. Sun, B. B. *et al.* Genomic atlas of the human plasma proteome. *Nature* **558**, 73–79 (2018).
8. Gudmundsdottir, V. *et al.* Circulating protein signatures and causal candidates for type 2 diabetes. *Diabetes* **69**, 1843–1853 (2020).
9. Nurmohamed, N. S. *et al.* Targeted proteomics improves cardiovascular risk prediction in secondary prevention. *Eur. Heart J.* **43**, 1569–1577 (2022).
10. Huth, C. *et al.* Protein markers and risk of type 2 diabetes and prediabetes: a targeted proteomics approach in the KORA F4/FF4 study. *Eur. J. Epidemiol.* **34**, 409–422 (2019).
11. LaFramboise, W. A. *et al.* Serum protein profiles predict coronary artery disease in symptomatic patients referred for coronary angiography. *BMC Med.* **10**, 157 (2012).
12. Mendelian Randomization Studies in Stroke: Exploration of Risk Factors and Drug Targets With Human Genetic Data | Stroke. <https://www.ahajournals.org/doi/full/10.1161/STROKEAHA.120.032617>.
13. Ritchie, S. C. *et al.* Integrative analysis of the plasma proteome and polygenic risk of cardiometabolic diseases. *Nat. Metab.* **3**, 1476–1483 (2021).

14. Sathyan, S. *et al.* Plasma proteomic profile of age, health span, and all-cause mortality in older adults. *Aging Cell* **19**, e13250 (2020).
15. Borrebaeck, C. A. K. Precision diagnostics: moving towards protein biomarker signatures of clinical utility in cancer. *Nat. Rev. Cancer* **17**, 199–204 (2017).
16. Woodward, M., Brindle, P. & Tunstall-Pedoe, H. Adding social deprivation and family history to cardiovascular risk assessment: the ASSIGN score from the Scottish Heart Health Extended Cohort (SHHEC). *Heart* **93**, 172–176 (2007).
17. Hippisley-Cox, J., Coupland, C. & Brindle, P. Development and validation of QRISK3 risk prediction algorithms to estimate future risk of cardiovascular disease: prospective cohort study. *BMJ* **357**, j2099 (2017).
18. Ganz, P. *et al.* Development and validation of a protein-based risk score for cardiovascular outcomes among patients with stable coronary heart disease. *JAMA* **315**, 2532–2541 (2016).
19. Wang, Z. *et al.* Metabolomic Pattern Predicts Incident Coronary Heart Disease. *Arterioscler. Thromb. Vasc. Biol.* **39**, 1475–1482 (2019).
20. Machado-Fragua, M. D. *et al.* Circulating serum metabolites as predictors of dementia: a machine learning approach in a 21-year follow-up of the Whitehall II cohort study. *BMC Med.* **20**, 334 (2022).
21. Eiriksdottir, T. *et al.* Predicting the probability of death using proteomics. *Commun. Biol.* **4**, 758 (2021).
22. Lind, L. *et al.* Large-Scale Plasma Protein Profiling of Incident Myocardial Infarction, Ischemic Stroke, and Heart Failure. *J. Am. Heart Assoc.* **10**, e023330 (2021).
23. Williams, S. A. *et al.* Plasma protein patterns as comprehensive indicators of health. *Nat. Med.* **25**, 1851–1857 (2019).
24. Buerge, T. *et al.* Metabolomic profiles predict individual multidisease outcomes. *Nat. Med.* 1–12 (2022) doi:10.1038/s41591-022-01980-3.
25. Sun, B. B. *et al.* Genetic regulation of the human plasma proteome in 54,306 UK Biobank participants. *bioRxiv* **20**, 2022.06.17.496443 (2022).

26. Kyu, H. H. *et al.* Global, regional, and national disability-adjusted life-years (DALYs) for 359 diseases and injuries and healthy life expectancy (HALE) for 195 countries and territories, 1990–2017: a systematic analysis for the Global Burden of Disease Study 2017. *The Lancet* **392**, 1859–1922 (2018).
27. James, S. L. *et al.* Global, regional, and national incidence, prevalence, and years lived with disability for 354 Diseases and Injuries for 195 countries and territories, 1990-2017: A systematic analysis for the Global Burden of Disease Study 2017. *The Lancet* **392**, 1789–1858 (2018).
28. Feigin, V. L. *et al.* Global, regional, and national burden of neurological disorders, 1990–2016: a systematic analysis for the Global Burden of Disease Study 2016. *Lancet Neurol.* **18**, 459 (2019).
29. Sherwani, S. I., Khan, H. A., Ekhzaimy, A., Masood, A. & Sakharkar, M. K. Significance of HbA1c Test in Diagnosis and Prognosis of Diabetic Patients. *Biomark. Insights* **11**, 95–104 (2016).
30. WHO. Use of Glycated Haemoglobin (HbA1c) in the Diagnosis of Diabetes Mellitus. Abbreviated Report of a WHO Consultation. WHO/NMH/CHP/CPM/11.1.
31. Next Steps For Risk Stratification in the NHS. NHS England. Available at:  
<https://www.england.nhs.uk/wp-content/uploads/2015/01/nxt-steps-risk-strat-glewis.pdf>.
32. Bollepalli, S., Korhonen, T., Kaprio, J., Anders, S. & Ollikainen, M. EpiSmokEr: A robust classifier to determine smoking status from DNA methylation data. *Epigenomics* **11**, 1469–1486 (2019).
33. Li, R., Chen, Y., Ritchie, M. D. & Moore, J. H. Electronic health records and polygenic risk scores for predicting disease risk. *Nat. Rev. Genet.* **21**, 493–502 (2020).
34. Lewis, C. M. & Vassos, E. Polygenic risk scores: from research tools to clinical instruments. *Genome Med.* **12**, 44 (2020).
35. Barnett, J. H., Lewis, L., Blackwell, A. D. & Taylor, M. Early intervention in Alzheimer’s disease: a health economic study of the effects of diagnostic timing. *BMC Neurol.* **14**, 101 (2014).
36. Crous-Bou, M., Minguillón, C., Gramunt, N. & Molinuevo, J. L. Alzheimer’s disease prevention: from risk factors to early intervention. *Alzheimers Res. Ther.* **9**, 71 (2017).
37. Foster, L. A. & Salajegheh, M. K. Motor Neuron Disease: Pathophysiology, Diagnosis, and Management. *Am. J. Med.* **132**, 32–37 (2019).

38. Tanaka, T. *et al.* Plasma proteomic biomarker signature of age predicts health and life span. *eLife* **9**, 1–24 (2020).
39. Bao, X. *et al.* Growth differentiation factor-15 is a biomarker for all-cause mortality but less evident for cardiovascular outcomes: A prospective study. *Am. Heart J.* **234**, 81–89 (2021).
40. Wu, P. F. *et al.* Growth Differentiation Factor 15 Is Associated With Alzheimer’s Disease Risk. *Front. Genet.* **12**, 1500 (2021).
41. McGrath, E. R. *et al.* Growth Differentiation Factor 15 and NT-proBNP as Blood-Based Markers of Vascular Brain Injury and Dementia. *J. Am. Heart Assoc.* **9**, (2020).
42. Myhre, P. L. *et al.* Growth Differentiation Factor 15 Provides Prognostic Information Superior to Established Cardiovascular and Inflammatory Biomarkers in Unselected Patients Hospitalized With COVID-19. *Circulation* **142**, 2128 (2020).
43. Wang, Z. *et al.* The impact of growth differentiation factor 15 on the risk of cardiovascular diseases: two-sample Mendelian randomization study. *BMC Cardiovasc. Disord.* **20**, 1–7 (2020).
44. Bidadkosh, A. *et al.* Predictive Properties of Biomarkers GDF-15, NTproBNP, and hs-TnT for Morbidity and Mortality in Patients With Type 2 Diabetes With Nephropathy. *Diabetes Care* **40**, 784–792 (2017).
45. Lemmelä, S. *et al.* Integrated analyses of growth differentiation factor-15 concentration and cardiometabolic diseases in humans. *eLife* **11**, e76272 (2022).
46. Zhang, X. *et al.* Association of a blood-based aging biomarker index with death and chronic disease: Cardiovascular Health Study. *J. Gerontol. Ser. A* glad172 (2023) doi:10.1093/gerona/glad172.
47. Choy, E. H. *et al.* Translating IL-6 biology into effective treatments. *Nat. Rev. Rheumatol.* **16**, 335–345 (2020).
48. Ridker, P. M. & Rane, M. Interleukin-6 Signaling and Anti-Interleukin-6 Therapeutics in Cardiovascular Disease. *Circ. Res.* **128**, 1728–1746 (2021).
49. Eugen-Olsen, J. *et al.* Circulating soluble urokinase plasminogen activator receptor predicts cancer, cardiovascular disease, diabetes and mortality in the general population. *J. Intern. Med.* **268**, 296–308 (2010).

50. Alirezaei, Z. *et al.* Neurofilament light chain as a biomarker, and correlation with magnetic resonance imaging in diagnosis of CNS-related disorders. *Mol. Neurobiol.* **57**, 469–491 (2020).
51. Wu, J. *et al.* Plasma neurofilament light chain: A biomarker predicting severity in patients with acute ischemic stroke. *Medicine (Baltimore)* **101**, e29692 (2022).
52. Pietzner, M. *et al.* Synergistic insights into human health from aptamer- and antibody-based proteomic profiling. *Nat. Commun.* **2021 121 12**, 1–13 (2021).
53. Lehallier, B. *et al.* Undulating changes in human plasma proteome profiles across the lifespan. *Nat. Med.* **25**, 1843–1850 (2019).
54. Hastie, T., Tibshirani, R., Narasimhan, B. & Chu, G. Package ‘impute’ Title impute: Imputation for microarray data. R package version 1.60.0. (2022).
55. Therneau, T. M. A Package for Survival Analysis in R. R package version 3.2-7, <https://CRAN.R-project.org/package=survival>. Accessed April 2021. (2020).
56. (2017), R. C. T. R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria.
57. Chang W, Cheng J, Allaire J, Sievert C, Schloerke B, Xie Y, Allen J, McPherson J, Dipert A, Borges B. shiny: Web Application Framework for R. R package version 1.7.3.9002, <https://shiny.rstudio.com/>.
58. J.J. Allaire, Christopher Gandrud, Kenton Russell and CJ Yetman. networkD3: D3 JavaScript Network Graphs from R. R package. <https://CRAN.R-project.org/package=networkD3>. (2017).
59. Csardi G, Nepusz T. The igraph software package for complex network research, *InterJournal, Complex Systems* 1695. <https://igraph.org>. (2006).
60. Cheng, Y. *et al.* DNA Methylation scores augment 10-year risk prediction of diabetes. *medRxiv* 2021.11.19.21266469 (2021) doi:10.1101/2021.11.19.21266469.
61. Simon, N., Friedman, J., Hastie, T. & Tibshirani, R. Regularization Paths for Cox’s Proportional Hazards Model via Coordinate Descent. *J. Stat. Softw.* **39**, (2011).
62. Greenwell B, Boehmke B, Cunningham J, Developers G. gbm: Generalized Boosted Regression Models. R package version 2.1.8.1. (2022).

63. Max Kuhn. Contributions from Jed Wing, Steve Weston, Andre Williams, Chris Keefer, Allan Engelhardt, Tony Cooper, Zachary Mayer, Brenton, Kenkel, the R Core Team, Michael Benesty, Reynald Lescarbeau, Andrew, & Ziem, Luca Scrucca, Yuan Tang and Can Candan. caret: Classification and Regression Training. R package version 6.0-71. (2016).
64. Yan, Y. MLmetrics: Machine Learning Evaluation Metrics. R package version 1.1.1. (2016).
65. Saito, T. & Rehmsmeier, M. Precrec: fast and accurate precision–recall and ROC curve calculations in R. *Bioinformatics* **33**, 145–147 (2017).