

1 **Proteomic signatures for identification of impaired glucose tolerance**

2

3 Julia Carrasco-Zanini¹, Maik Pietzner^{1,2}, Joni V Lindbohm^{3,4}, Eleanor Wheeler¹, Erin Oerton¹, Nicola
4 Kerrison¹, Missy Simpson⁵, Matthew Westacott⁵, Dan Drolet⁵, Mika Kivimaki^{3,4}, Rachel Ostroff⁵,
5 Stephen A Williams⁵, Nicholas J Wareham¹, Claudia Langenberg^{1,2,6*}

6

7 **Affiliations**

8 1. MRC Epidemiology Unit, Institute of Metabolic Science, University of Cambridge School of Clinical
9 Medicine, Cambridge, UK

10 2. Computational Medicine, Berlin Institute of Health at Charité – Universitätsmedizin Berlin, Germany

11 3. Clincum, Department of Public Health, University of Helsinki, P.O. Box 41, FI-00014 Helsinki, Finland

12 4. Department of Epidemiology and Public Health, University College London, UK

13 5. SomaLogic, Inc., Boulder, CO, USA

14 6. Precision Healthcare University Research Institute, Queen Mary University of London, UK

15

16 *Correspondence to Dr Claudia Langenberg (claudia.langenberg@mrc-epid.cam.ac.uk)

17

18

19

20 **Abstract**

21 The implementation of recommendations for type 2 diabetes (T2D) screening and diagnosis focus on
22 measurement of HbA1c and fasting glucose. This approach leaves a large number of individuals with
23 isolated impaired glucose tolerance (iIGT), who are only detectable through oral glucose tolerance
24 tests (OGTTs), at risk of diabetes and its severe complications. We applied machine learning to
25 proteomic profiles of a single fasted sample from 11,546 participants of the Fenland study to test
26 discrimination of iIGT defined using gold standard OGTTs. We observed significantly improved
27 discriminative performance by adding only three proteins (RTN4R, CBPM, and GHR) to the best clinical
28 model (0.80 (0.79-0.86), $p=0.004$), which we validated in an external cohort. Increased plasma levels
29 of these candidate proteins were associated with an increased risk for future T2D in an independent
30 cohort and were also increased in individuals genetically susceptible to impaired glucose homeostasis
31 and T2D. Assessment of a limited number of proteins can identify individuals likely to be missed by
32 current diagnostic strategies and at high risk of T2D and its complications.

33

34 **Introduction**

35 Current clinical guidelines for type 2 diabetes (T2D) screening and diagnosis are based on glycated
36 haemoglobin (HbA1c) and fasting glucose (FG) levels for reasons of practicality, however alternative
37 tests can be used^{1,2}. Globally, over 7.5% of adults have impaired glucose tolerance (IGT)³ with
38 increased prevalence reported in older individuals⁴ and specific ethnic groups, such as people from
39 Southeast Asia⁵. A substantial proportion of people with IGT (28 – 86%)⁶⁻⁸ can only be identified
40 through oral glucose tolerance tests (OGTTs), which are inconvenient and time-consuming. Individuals
41 with isolated IGT (iIGT), that is, 2-hour plasma glucose (2hPG) ≥ 7.8 and < 11.1 mmol/L but normal
42 HbA1c and fasting glucose, remain undetected by current T2D detection strategies⁹⁻¹² but are at very
43 high risk of developing diabetes (annualized T2D relative risk of 5.5 compared to normoglycemic
44 individuals)¹³ and presenting with its severe micro- and macrovascular complications^{9-12,14}. Compared
45 to individuals with fasting hyperglycaemia, mortality is twice as high in the iIGT group over a period of
46 5 to 12 years^{15,16}.

47 Small proof-of-concept studies in cohorts of high-risk individuals have demonstrated the value of deep
48 molecular profiling for early identification of pathways that are differentially regulated between
49 individuals with and without insulin resistance^{17,18} and to guide its prediction¹⁹. Deep profiling of the
50 plasma proteome at population scale has become possible through aptamer-based affinity assays²⁰.
51 The systematic study of the circulating proteome promises to improve strategies for prediction and
52 diagnosis¹⁸ as well as aetiological understanding, including identification of novel pathways leading to
53 T2D and refinement of aetiological subtypes.

54 Because of the high global prevalence of IGT and iIGT, their severe complications, and the currently
55 unmet need of screening strategies that can identify iIGT without a challenge test, we used machine
56 learning to test whether large-scale proteomic profiling of a single fasted sample could identify
57 individuals with iIGT and improve current clinical models. We then tested whether the most
58 discriminatory proteins were affected by fasting status, to assess the feasibility of using non-fasted
59 samples to identify iIGT. To gain insights into IGT and iIGT aetiology, we 1) identified and characterised
60 biochemical, phenotypic, and anthropometric drivers of discriminatory proteins, 2) investigated
61 whether their plasma levels were associated with the risk of future T2D in an independent prospective
62 cohort with 521 incident T2D cases, and 3) tested the influence of genetic susceptibility to T2D or
63 related phenotypes on protein levels.

64 **Results**

65 We used an aptamer-based assay to target 4,775 distinct fasting plasma proteins by 4,979 aptamers
66 in 11,546 participants (5,389 men and 6,157 women) without diagnosed diabetes from the

67 contemporary Fenland study²¹ (baseline visit in 2005-2015, mean age 48.5 years (7.5 s.d.),
68 **Supplementary Table 1**), as previously described¹⁸ (Methods). Participants completed a 75-gram
69 OGTT (**Figure 1a**). We defined isolated post challenge hyperglycaemia as 2hPG ≥ 7.8 mmol/L but HbA1c
70 < 42 mmol/mol and FG < 6.1 mmol/L. This definition captured all participants with isolated IGT (2hPG
71 $7.8-11.1$ mmol/L but HbA1c < 42 mmol/mol and FG < 6.1 mmol/L) as well as participants with isolated
72 post-challenge hyperglycaemia in the diabetic range (2hPG ≥ 11.1 mmol/L but HbA1c < 42 mmol/mol
73 and FG < 6.1 mmol/L, N=117), i.e. high-risk individuals missed by standard FG and HbA1c testing. For
74 simplicity, we refer from here on to IGT (or iIGT) for all individuals with 2hPG ≥ 7.8 mmol/L, without
75 specifically distinguishing post-challenge hyperglycaemia ≥ 11.1 mmol/L. We used a least absolute
76 shrinkage and selection operator (LASSO) regression framework implemented as a three-step
77 approach, including independent feature selection (50% sample size), optimization (25%) and
78 validation (25%) to discriminate IGT (prevalence 6.7%) and iIGT (3.9%) based on fasting assessment of
79 4,775 proteins (targeted by 4,979 aptamers) (**Figure 1b**). We defined highly discriminatory proteins as
80 those selected in $> 80\%$, 90% , or 95% of random subsamples of the study population during feature
81 selection (**Extended Data Figure 1**).

82

83 ***Proteomic signatures to discriminate IGT and IIGT***

84 We identified 65 and 68 proteins, respectively, that achieved an area under the receiver operating
85 characteristic curve (AUROC) (95% confidence interval) of 0.83 (0.80 – 0.86) and 0.77 (0.72 – 0.81) for
86 discrimination of IGT and iIGT in the independent validation set (**Extended Data Figure 2,**
87 **Supplementary Tables 2 and 3**). This represented a significantly better predictor when compared to
88 the performance of a T2D genetic risk score (T2D-GRS, AUROC_{IGT} 0.58 (0.52 – 0.63), AUROC_{iIGT} 0.54
89 (0.49 – 0.60)) (**Figure 2a and b,** and **Extended Data Figure 3**). Protein-based models further
90 outperformed the standard patient information-based model (based on the Cambridge Diabetes Risk
91 Score including age, sex, family history of diabetes, smoking status, prescription of steroid or
92 antihypertensive medication and body mass index (BMI))²² (AUROC_{IGT} = 0.71 (0.67 – 0.75); AUROC_{iIGT}
93 = 0.71 (0.66 – 0.76)) and the standard clinical model that additionally included blood test results, that
94 is, FPG and HbA1c (AUROC_{IGT}= 0.78 (0.74 – 0.82); AUROC_{iIGT}=0.75 (0.70 – 0.80)) (**Figure 2a and b,**
95 **Supplementary Table 4**).

96 Considering a limited set of the most informative proteins that were identified by the feature selection
97 framework (Methods), discrimination was still superior to the standard clinical model adding only 8
98 proteins for IGT (AUROC_{IGT} 0.83 (0.80 – 0.86), p-value = 4.13×10^{-5} , **Figure 2a, Supplementary Table**
99 **2**) and 3 proteins for iIGT (AUROC_{iIGT} 0.80 (0.76 – 0.85), p-value = 0.004, **Figure 2b, Supplementary**
100 **Table 3**), including 2 proteins (Reticulon-4 receptor, Carboxypeptidase M) selected for both (**Figure**
101 **2c, Supplementary Table 5**²³⁻³³). The weights for the variables included in these final models are
102 available in **Supplementary Table 6**. We observed significant improvement over and above the clinical
103 model of similar magnitude in the independent Whitehall II (WHII) study (**Supplementary Table 7 and**
104 **8, Extended Data Figure 4**).

105 To identify participants with iIGT and IGT, respectively, we choose a cut-off for the clinical + protein
106 model that optimized sensitivity (recall) at 0.70 and 0.71, which yielded a positive predicted value
107 (precision) of 0.20 and 0.13, respectively. The net reclassification index was higher for the final iIGT
108 model (14.5%) compared to IGT (6.5%), consistent with the current lack of informative predictors.

109 Of the 9 distinct proteins included in the 2 final models, 8 were not significantly affected by fasting
110 status (**Methods**) with maximum postprandial fold changes ranging between 0.07 and 0.16; only
111 HTRA1 showed some evidence of a post-prandial increase (maximum fold change= 0.15, p-
112 value=0.004, **Supplementary Table 9**).

113 Finally, we tested model performance de novo omitting the 3 most informative proteins to predict
114 iIGT. The novel model included 7 proteins and still performed significantly better than the best clinical
115 model (AUROC = 0.78 (0.73 – 0.83), p-value = 0.04, **Extended Data Figure 5**). This finding illustrates

116 redundancy in the protein biomarkers available to select from for iIGT prediction, providing practical
117 benefits for clinical implementation, for example with regard to flexibility of prioritising choice of
118 proteins more easily targeted by clinical chemistry assays, least affected by fasting status or sample
119 handling.

120 *Proteomically informed screening strategies*

121 We calculated the numbers needed to screen (NNS) to determine how many OGTTs would need to be
122 performed to identify one participant with iIGT using a three-stage screening approach (**Figure 3**). We
123 stratified all Fenland individuals based on the patient-derived information model in the first instance
124 and based on their HbA1c levels and the 3-protein iIGT model in the second instance (**Methods**).
125 According to current guidelines², individuals at high predicted risk based on the patient-derived
126 information model, but HbA1c levels below cut-offs for prediabetes or T2D² would not be considered
127 for further testing (N Fenland = 4163, NNS = 14, **Figure 3**). Applying the clinical + 3-protein iIGT model
128 on this group enabled identification of a high-risk subgroup (N = 1739) in which application of an OGTT
129 should be considered, since the NNS was only 7 to identify one additional individual with iIGT (**Figure**
130 **3, Supplementary Table 10**). Hence, our proposed approach identified an additional >30% of
131 individuals that would be reclassified (as having prediabetes) and could be offered preventative
132 interventions, that is, a substantial proportion of high-risk individuals that would otherwise be missed
133 by current strategies. To test for potential bias in the NNS estimates arising from overfitting, we
134 applied the same screening algorithm in the test set only, which provided internal validation for the
135 estimates and results from the entire Fenland set (**Extended Data Figure 6**).

136 *Characterisation of discriminatory proteins*

137 To investigate whether increased genetic risk of diabetes and related metabolic risk factors affect
138 abundances of the identified proteins, we compared their differences in individuals with higher versus
139 lower genetic risk based on genetic risk scores (GRS) for T2D and related endophenotypes, including
140 fasting glucose³⁴, fasting insulin³⁴, 2hPG³⁴, body mass index (BMI)³⁵ and T2D³⁶, using linear regression
141 models. We found evidence of significant, directional concordant associations between genetic
142 susceptibility to these phenotypes and plasma abundances for 4 of the 9 most predictive IGT and iIGT
143 proteins, (p-value < 0.001, **Figure 4c**). Plasma abundances of Growth hormone receptor (GHR),
144 Reticulon-4 receptor (RTN4R), Carboxypeptidase M (CBPM) and Serine protease HTRA1 (HTRA1) were
145 associated with genetic susceptibility to more than one of these phenotypes, including fasting insulin,
146 T2D and BMI.

147 The 3 most predictive iIGT proteins and 6 of the 8 most predictive IGT proteins were significantly
148 associated with higher measured concentrations of fasting and 2-hour glucose, and insulin.

149 Chondroadherin (CHAD) was the only protein inversely associated with all 4 measures. From the
150 remaining two IGT predictor proteins only Cartilage intermediate layer protein 2 (CILP2) was
151 significantly inversely associated with fasting glucose (p -values <0.001 , **Figure 3a**). In the independent
152 prospective WHII cohort ($N = 1,492$, including 521 incident T2D cases, **Supplementary Table 11**), all
153 proteins were significantly associated with an increased risk of developing future T2D, except for
154 CHAD, which was inversely associated (p -value < 0.006 , **Figure 4b**), and CILP2, which showed no
155 significant association. Effect sizes ranged from 0.88–1.51 (hazard ratio for T2D per s.d. difference in
156 the protein target) adjusting for age, sex, and BMI. Associations for HTRA1, GHR, and CBPM remained
157 significant even upon additional adjustment for fasting glucose, total triglycerides, HDL-cholesterol,
158 and lipid lowering medication (**Supplementary Table 12**).

159 Informative biomarkers are not only relevant to improve screening strategies but can inform
160 understanding of the separate and shared aetiologies of IGT and iIGT. Comparison of protein ranking
161 from IGT as opposed to iIGT feature selection revealed that most discriminatory proteins differed
162 strongly between the IGT and iIGT selections (**Extended Data Figure 7**) with only eleven proteins
163 achieving similarly high rankings for both outcomes, that is, being selected in $>80\%$ across random
164 subsets of the study population. The top two biological GO term processes differed between the 65-
165 IGT protein signature (“proteolysis” and “cytokine-mediated signalling pathway”, **Supplementary**
166 **Table 13**) and the 68-iIGT protein signature (“cartilage development”, “collagen fibril organization”,
167 **Supplementary Table 14**), however none were significantly enriched following Bonferroni adjustment
168 for multiple comparisons.

169 To identify potential differences in factors influencing these IGT and iIGT protein signatures, we
170 computed the proportion of variance in the first principal component of the 65-IGT and 68-iIGT protein
171 signatures explained by 24 biochemical, phenotypic, and anthropometric factors. Both signatures had
172 similarly large proportions of explained variance by glycaemic (5.2 – 37.8%) and anthropometric (25.1
173 – 40.9%) measures, blood lipids (2.7 – 33.1%), or an ultrasound-based score for hepatic steatosis (22.4
174 – 24.5%) (**Methods**). Differences included the higher proportion of variance explained by C-reactive
175 protein and the lower proportion explained by ALT (a biomarker of liver injury) for the 65-IGT
176 compared to the 68-iIGT protein signature (CRP 30.2% vs 20.3% and ALT 14.7% vs 23.2%, respectively,
177 **Extended Data Figure 8**). Measures related to glucose metabolism (explaining up to 23.8% of the
178 variance) and adiposity (explaining up to 26.9 % of the variance) were identified as the main factors
179 explaining variance in the 9 predictive IGT or iIGT proteins included in the final prediction models.
180 Other protein specific factors included total triglycerides (explained up to 22.6% of GHR), HDL-
181 cholesterol (up to 13.6% of RTN4R), measures of hepatic steatosis (liver score explained up to 15% of

182 GHR) and inflammation (up to 27.2% of HTRA1), as well as genetic variants in proximity of the relevant
183 protein-encoding gene (up to 11.3% of RTN4R) (**Extended Data Figure 8**).

184

185 ***Long-term health outcomes associated with predicted iIGT***

186 To explore the clinical consequences of isolated impaired glucose tolerance in the absence of an OGTT,
187 we performed an exploratory analysis in a random sub-cohort of the prospective EPIC-Norfolk study³⁷
188 (N=753). We evaluated associations between predicted probabilities based on 1) the final clinical + 3-
189 protein model, 2) the 3-protein model only, and 3) the 68-protein iIGT model with the onset of eight
190 cardiometabolic diseases based on electronic-health record linkage³⁸ (N incident cases 30-235; follow-
191 up time between 18 – 19 years; **Supplementary Table 15 - 16**). All scores were significantly associated
192 with a greater risk of future T2D (52 incident T2D cases) at 5% false discovery rate (FDR). The iIGT final
193 clinical+3-protein score was further associated with cataracts and renal disease, possibly reflecting the
194 known association between chronically elevated 2hPG levels and micro- or macrovascular
195 complications. Predicted probabilities from the best performing 68-protein-based iIGT-model, showed
196 a nominally significant association for coronary artery disease (HR = 1.22, p-value = 0.03, CAD) and
197 peripheral artery disease (HR = 1.27, p-value = 0.04, PAD), T2D-related complications, although these
198 did not reach statistical significance when adjusting for multiple testing given the small number of
199 incident cases in this small exploratory cohort. We observed significant associations for individual
200 proteins with the risk of future T2D, with effect sizes comparable to those in the WHII study³⁹ (**Figure**
201 **5**).

202 We used proteomic measures done with a distinct proteomic technique, the Olink Explore panel⁴⁰ in
203 an independent study (random sub-cohort of the prospective EPIC-Norfolk study, N=602) to test
204 correlation of overlapping protein predictors and to validate some of our findings using an orthogonal
205 technique. We observed a high correlation between the SomaScan and Olink measurements for the
206 top three selected proteins (N=50, Spearman's r: GHR = 0.80, RTN4R = 0.70 and CBPM = 0.87,
207 Pearson's r: GHR = 0.80, RTN4R = 0.72 and CBPM = 0.82). In line with this, we replicated the previously
208 observed associations with an increased risk of incident T2D, including comparable effect sizes, and
209 further observed significant associations between the final clinical + 3-protein model and incident
210 cataracts, heart failure, and coronary heart disease (**Extended Data Figure 9**). These findings suggest
211 cross-platform transferability of our results.

212

213

214

215 **Discussion**

216 Behavioural interventions in individuals with IGT have been shown to delay progression to T2D and
217 reduce the risk of long term microvascular and macrovascular complications⁴¹. However, individuals
218 with iIGT are likely to remain undiagnosed because the current implementation of recommendations
219 for screening and diagnosing T2D does not focus on OGTTs, for reasons of practicality. People with
220 iIGT are at high risk of developing T2D and its associated complications, and failure to identify them
221 can lead to the development of severe and potentially irreversible complications of their unmanaged
222 hyperglycaemia¹⁶.

223 By combining deep plasma proteomic profiling with machine learning, we developed models for
224 improved identification of IGT and iIGT and demonstrated that as few as 8 and 3 proteins, respectively,
225 provided significant improvement over established clinical predictors²². We provided external
226 validation of the significant and substantial improvement achieved by the selected proteins over and
227 above the stringent benchmark provided by the best clinical model, something rarely done in genomic
228 or other 'omic prediction studies. The improvement observed in our independent replication study
229 was slightly greater than what was originally observed, and we note that the lack of HbA1c
230 measurements and other differences in study design (previous phases including OGTT screening) and
231 participant characteristics (older and more males on average) of the Whitehall II cohort³⁹ are likely to
232 have contributed to this, leading to a lower AUROC for the clinical model and/ or potential
233 misclassification of iIGT.

234 We propose a 3-step screening strategy, in line with the current UK Diabetes Prevention
235 Programmes⁴², involving risk assessment by 1) a patient-derived information model, 2) measuring
236 HbA1c levels and only 3 additional proteins from a single spot blood sample, and 3) an OGTT for
237 eventual diagnosis. Implementation of this proposed screening strategy, could lead to a large
238 proportion of individuals with iIGT to be additionally identified with a lower NNS, compared to the
239 currently recommended 2-stage approach⁴². Our findings illustrate how the identified proteins could
240 most efficiently be integrated into existing screening approaches to identify individuals with iIGT, who
241 are at high risk of T2D and its complications but are currently being missed. Behavioural interventions
242 have shown to be effective at reversing post-load hyperglycaemia independently of fasting glucose
243 levels^{43,44}, emphasising the value of identifying individuals with iIGT who would benefit the most from
244 these interventions. We further provided evidence of a link between our developed iIGT predictive
245 scores with incident T2D and several known cardiometabolic comorbidities resulting from chronically
246 elevated 2hPG. These finding highlight the potential of applying such a predictive risk score not only

247 for cross-sectional identification of iIGT, but for monitoring future risk for associated comorbidities
248 that impact patients' quality of life.

249 We showed that the identified proteins are not strongly affected by fasting status, suggesting that
250 they could enable a simple and convenient strategy to better identify individuals with IGT and iIGT,
251 compared to an OGTT, which requires repeated blood draws conveying additional costs¹⁸. Protein
252 assessment could substantially improve the feasibility and acceptability of an improved strategy to
253 identify iIGT, more so than alternative strategies that have been proposed such as a 1-hour OGTT⁴⁵,
254 and hence brings it in line with existing strategies for the screening and diagnosis of T2D. Since HbA1c
255 testing requires anticoagulated whole-blood, usually EDTA, a subset of the same sample type could
256 be processed for plasma preparation to measure discriminatory proteins, avoiding the need for
257 additional blood sampling.

258 This study provided insights into aetiological differences between iIGT and IGT. Our results suggested
259 a stronger low-grade inflammatory component⁴⁶⁻⁴⁹ among proteins discriminatory for IGT compared
260 to those for iIGT. These proteins might represent refined biomarkers of low-grade inflammation, as
261 they were highlighted as being predictive over and above established inflammatory markers also
262 covered in our proteomic study, such as C-reactive protein. At an individual biomarker level, we
263 identified a number of proteins shared or distinctly associated with these metabolic disturbances,
264 including GHR, RTN4R, HTRA1, CBPM, CHAD, CBLN4, NEU1, CILP2, and S100-A10. We used genetic
265 data to provide evidence that early deregulation of diabetes related pathways is linked to the
266 candidate proteins, most of which were also significantly associated with risk of future development
267 of T2D, providing a novel set of high priority T2D targets for further follow-up and assessment in in
268 more diverse settings and ethnicities.

269 While our model estimated a meaningful decrease in the NNS, there are important consideration for
270 implantation of the proposed strategy. A considerable proportion of individuals with iIGT were missed
271 by being classified low risk in either the first or subsequent screening steps. A further limitation of our
272 study was the lack of orthogonal validation of our protein-based prediction models with an alternative
273 proteomic technology. Technical, genetic, and other biological factors can result in biased protein
274 measurements due to changes in affinity of the aptamer reagents⁵⁰. However, the strong correlations
275 observed with the antibody-based Olink Explore panel suggests cross-platform transferability. We
276 further validated the phenotypic association of the iIGT predictive protein scores with incident
277 cardiometabolic diseases using Olink Explore measurements, providing the possibility of
278 implementing our model with alternative proteomic technologies.

279 In summary, we demonstrated the utility of the plasma proteome to inform strategies for screening
280 of iIGT and for gaining novel aetiological insights into early signatures of impaired glucose tolerance,
281 a globally very common and clinically important metabolic disorder, but one that it is difficult to detect
282 and treat in routine clinical practice.

283

284 **Acknowledgements**

285 The Fenland Study (10.22025/2017.10.101.00001) is funded by the Medical Research Council
286 (MC_UU_12015/1). We are grateful to all the volunteers and to the General Practitioners and practice
287 staff for assistance with recruitment. We thank the Fenland Study Investigators, Fenland Study Co-
288 ordination team and the Epidemiology Field, Data and Laboratory teams. We further acknowledge
289 support for genomics from the Medical Research Council (MC_PC_13046). Proteomic measurements
290 were supported and governed by a collaboration agreement between the University of Cambridge
291 and SomaLogic. We thank Ira von Carlowitz and Kaitlin Soucie for their contributions to the fasting
292 proteome analysis. JCZS is supported by a 4-year Wellcome Trust PhD Studentship and the Cambridge
293 Trust, CL, EW, and NJW are funded by the Medical Research Council (MC_UU_12015/1). NJW is a NIHR
294 Senior Investigator. The Whitehall II study and MK are supported by grants from the Wellcome Trust
295 (221854/Z/20/Z); UK Medical Research Council (R024227); and NIA, NIH (R01AG056477). JVL was
296 supported by Academy of Finland (311492 and 339568) and Helsinki Institute of Life Science (H970)
297 grants paid to employer. The funders had no role in study design, data collection and analysis, decision
298 to publish or preparation of the manuscript.

299 **Author Contributions**

300 JCZS, MP, NJW and CL designed the analysis and drafted the manuscript. JCZS analysed the data, JVL
301 did the replication analyses in Whitehall II study. MS and MW did the analysis for assessing the effect
302 of fasting status on protein levels. NJW is PI of the Fenland cohort and MK is PI of the Whitehall II
303 study. All authors contributed to the interpretation of the results and critically reviewed the
304 manuscript.

305

306 **Competing Interests**

307 MS, MW, DD, RO and SAW are employees of SomaLogic. EW and EO are now employees at AstraZeneca.
308 The remaining authors declare no competing interests.

309

310 **References**

- 311 1. American Diabetes, A. 2. Classification and Diagnosis of Diabetes: Standards of Medical Care
312 in Diabetes-2018. *Diabetes Care* **41**, S13-S27 (2018).
- 313 2. International Expert, C. International Expert Committee report on the role of the A1C assay in
314 the diagnosis of diabetes. *Diabetes Care* **32**, 1327-1334 (2009).
- 315 3. Saeedi, P., *et al.* Global and regional diabetes prevalence estimates for 2019 and projections
316 for 2030 and 2045: Results from the International Diabetes Federation Diabetes Atlas, 9(th)
317 edition. *Diabetes Res Clin Pract* **157**, 107843 (2019).
- 318 4. Meisinger, C., *et al.* Prevalence of undiagnosed diabetes and impaired glucose regulation in
319 35-59-year-old individuals in Southern Germany: the KORA F4 Study. *Diabet Med* **27**, 360-362
320 (2010).
- 321 5. Cheng, Y.J., *et al.* Prevalence of Diabetes by Race and Ethnicity in the United States, 2011-
322 2016. *JAMA* **322**, 2389-2398 (2019).
- 323 6. Richter, B., Hemmingsen, B., Metzendorf, M.I. & Takwoingi, Y. Development of type 2 diabetes
324 mellitus in people with intermediate hyperglycaemia. *Cochrane Database Syst Rev* **10**,
325 CD012661 (2018).
- 326 7. Yip, W.C.Y., Sequeira, I.R., Plank, L.D. & Poppitt, S.D. Prevalence of Pre-Diabetes across
327 Ethnicities: A Review of Impaired Fasting Glucose (IFG) and Impaired Glucose Tolerance (IGT)
328 for Classification of Dysglycaemia. *Nutrients* **9**(2017).
- 329 8. Campbell, M.D., *et al.* Benefit of lifestyle-based T2DM prevention is influenced by prediabetes
330 phenotype. *Nat Rev Endocrinol* **16**, 395-400 (2020).
- 331 9. Nichols, G.A., Arondekar, B. & Herman, W.H. Complications of dysglycemia and medical costs
332 associated with nondiabetic hyperglycemia. *The American journal of managed care* **14**, 791-
333 798 (2008).
- 334 10. Cowie, C.C., *et al.* Prevalence of diabetes and high risk for diabetes using A1C criteria in the
335 U.S. population in 1988-2006. *Diabetes Care* **33**, 562-568 (2010).
- 336 11. Cederberg, H., *et al.* Postchallenge glucose, A1C, and fasting glucose as predictors of type 2
337 diabetes and cardiovascular disease: a 10-year prospective cohort study. *Diabetes Care* **33**,
338 2077-2083 (2010).
- 339 12. Balkau, B. The DECODE study. Diabetes epidemiology: collaborative analysis of diagnostic
340 criteria in Europe. *Diabetes Metab* **26**, 282-286 (2000).
- 341 13. Gerstein, H.C., *et al.* Annual incidence and relative risk of diabetes in people with various
342 categories of dysglycemia: a systematic overview and meta-analysis of prospective studies.
343 *Diabetes Res Clin Pract* **78**, 305-312 (2007).
- 344 14. Chen, Y., *et al.* Associations of progression to diabetes and regression to normal glucose
345 tolerance with development of cardiovascular and microvascular disease among people with
346 impaired glucose tolerance: a secondary analysis of the 30 year Da Qing Diabetes Prevention
347 Outcome Study. *Diabetologia* **64**, 1279-1287 (2021).
- 348 15. Shaw, J.E., Hodge, A.M., de Courten, M., Chitson, P. & Zimmet, P.Z. Isolated post-challenge
349 hyperglycaemia confirmed as a risk factor for mortality. *Diabetologia* **42**, 1050-1054 (1999).
- 350 16. Silbernagel, G., *et al.* Isolated post-challenge hyperglycaemia predicts increased
351 cardiovascular mortality. *Atherosclerosis* **225**, 194-199 (2012).
- 352 17. Zhou, W., *et al.* Longitudinal multi-omics of host-microbe dynamics in prediabetes. *Nature*
353 **569**, 663-671 (2019).
- 354 18. Williams, S.A., *et al.* Plasma protein patterns as comprehensive indicators of health. *Nat Med*
355 **25**, 1851-1857 (2019).
- 356 19. Schussler-Fiorenza Rose, S.M., *et al.* A longitudinal big data approach for precision health. *Nat*
357 *Med* **25**, 792-804 (2019).
- 358 20. Gold, L., *et al.* Aptamer-based multiplexed proteomic technology for biomarker discovery.
359 *PLoS One* **5**, e15004 (2010).

- 360 21. Lindsay, T., *et al.* Descriptive epidemiology of physical activity energy expenditure in UK adults
361 (The Fenland study). *International Journal of Behavioral Nutrition and Physical Activity* **16**, 126
362 (2019).
- 363 22. Rahman, M., Simmons, R.K., Harding, A.H., Wareham, N.J. & Griffin, S.J. A simple risk score
364 identifies individuals at high risk of developing Type 2 diabetes: a prospective cohort study.
365 *Fam Pract* **25**, 191-196 (2008).
- 366 23. Deora, A.B., Kreitzer, G., Jacovina, A.T. & Hajjar, K.A. An annexin 2 phosphorylation switch
367 mediates p11-dependent translocation of annexin 2 to the cell surface. *J Biol Chem* **279**,
368 43411-43418 (2004).
- 369 24. Guevara-Aguirre, J., *et al.* Growth hormone receptor deficiency is associated with a major
370 reduction in pro-aging signaling, cancer, and diabetes in humans. *Sci Transl Med* **3**, 70ra13
371 (2011).
- 372 25. Tiaden, A.N., *et al.* Novel Function of Serine Protease HTRA1 in Inhibiting Adipogenic
373 Differentiation of Human Mesenchymal Stem Cells via MAP Kinase-Mediated MMP
374 Upregulation. *Stem Cells* **34**, 1601-1614 (2016).
- 375 26. Haddad, Y. & Couture, R. Kininase 1 As a Preclinical Therapeutic Target for Kinin B1 Receptor
376 in Insulin Resistance. *Front Pharmacol* **8**, 509 (2017).
- 377 27. Klement, J., *et al.* Oxytocin Improves beta-Cell Responsivity and Glucose Tolerance in Healthy
378 Men. *Diabetes* **66**, 264-271 (2017).
- 379 28. Zhong, C., *et al.* Cbln1 and Cbln4 Are Structurally Similar but Differ in GluD2 Binding
380 Interactions. *Cell Rep* **20**, 2328-2340 (2017).
- 381 29. Weingarten, M.F.J., *et al.* Circulating Oxytocin Is Genetically Determined and Associated With
382 Obesity and Impaired Glucose Tolerance. *J Clin Endocrinol Metab* **104**, 5621-5632 (2019).
- 383 30. Wu, T., *et al.* CILP-2 is a novel secreted protein and associated with insulin resistance. *J Mol*
384 *Cell Biol* **11**, 1083-1094 (2019).
- 385 31. Slieker, R.C., *et al.* Novel biomarkers for glycaemic deterioration in type 2 diabetes: an IMI
386 RHAPSODY study. *medRxiv*, 2021.2004.2022.21255625 (2021).
- 387 32. Shen, Z., Gantcheva, S., Mansson, B., Heinegard, D. & Sommarin, Y. Chondroadherin
388 expression changes in skeletal development. *Biochem J* **330 (Pt 1)**, 549-557 (1998).
- 389 33. Hesse, L., *et al.* The skeletal phenotype of chondroadherin deficient mice. *PLoS One* **8**, e63080
390 (2014).
- 391 34. Scott, R.A., *et al.* Large-scale association analyses identify new loci influencing glycemic traits
392 and provide insight into the underlying biological pathways. *Nat Genet* **44**, 991-1005 (2012).
- 393 35. Lotta, L.A., *et al.* Association of Genetic Variants Related to Gluteofemoral vs Abdominal Fat
394 Distribution With Type 2 Diabetes, Coronary Disease, and Cardiovascular Risk Factors. *JAMA*
395 **320**, 2553-2563 (2018).
- 396 36. Mahajan, A., *et al.* Fine-mapping type 2 diabetes loci to single-variant resolution using high-
397 density imputation and islet-specific epigenome maps. *Nat Genet* **50**, 1505-1513 (2018).
- 398 37. Day, N., *et al.* EPIC-Norfolk: study design and characteristics of the cohort. European
399 Prospective Investigation of Cancer. *Br J Cancer* **80 Suppl 1**, 95-103 (1999).
- 400 38. Pietzner, M., *et al.* Plasma metabolites to profile pathways in noncommunicable disease
401 multimorbidity. *Nat Med* **27**, 471-479 (2021).
- 402 39. Marmot, M. & Brunner, E. Cohort Profile: the Whitehall II study. *Int J Epidemiol* **34**, 251-256
403 (2005).
- 404 40. Zhong, W., *et al.* Next generation plasma proteome profiling to monitor health and disease.
405 *Nat Commun* **12**, 2493 (2021).
- 406 41. Gong, Q., *et al.* Morbidity and mortality after lifestyle intervention for people with impaired
407 glucose tolerance: 30-year results of the Da Qing Diabetes Prevention Outcome Study. *Lancet*
408 *Diabetes Endocrinol* **7**, 452-461 (2019).

- 409 42. Barron, E., Clark, R., Hewings, R., Smith, J. & Valabhji, J. Progress of the Healthier You: NHS
410 Diabetes Prevention Programme: referrals, uptake and participant characteristics. *Diabet Med*
411 **35**, 513-518 (2018).
- 412 43. Gong, Q., *et al.* Efficacy of lifestyle intervention in adults with impaired glucose tolerance with
413 and without impaired fasting plasma glucose: A post hoc analysis of Da Qing Diabetes
414 Prevention Outcome Study. *Diabetes Obes Metab* **23**, 2385-2394 (2021).
- 415 44. Knowler, W.C., *et al.* Reduction in the incidence of type 2 diabetes with lifestyle intervention
416 or metformin. *N Engl J Med* **346**, 393-403 (2002).
- 417 45. Bergman, M., *et al.* Lessons learned from the 1-hour post-load glucose level during OGTT:
418 Current screening recommendations for dysglycaemia should be revised. *Diabetes Metab Res*
419 *Rev* **34**, e2992 (2018).
- 420 46. Pham, C.T. Neutrophil serine proteases: specific regulators of inflammation. *Nat Rev Immunol*
421 **6**, 541-550 (2006).
- 422 47. Wiedow, O. & Meyer-Hoffert, U. Neutrophil serine proteases: potential key regulators of cell
423 signalling during inflammation. *J Intern Med* **257**, 319-328 (2005).
- 424 48. Donath, M.Y. & Shoelson, S.E. Type 2 diabetes as an inflammatory disease. *Nat Rev Immunol*
425 **11**, 98-107 (2011).
- 426 49. de Vries, M.A., *et al.* Glucose-dependent leukocyte activation in patients with type 2 diabetes
427 mellitus, familial combined hyperlipidemia and healthy controls. *Metabolism* **64**, 213-217
428 (2015).
- 429 50. Pietzner, M., *et al.* Synergistic insights into human health from aptamer- and antibody-based
430 proteomic profiling. *Nat Commun* **12**, 6822 (2021).

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447 **Figure Legends**

448 **Figure 1. Study design.** **a**, Proteomic profiling was done in fasting plasma samples from participants from the
449 Fenland cohort that had undergone an OGTT. **b**, 3-step modelling framework for IGT and iIGT classification. *For
450 iIGT prediction individuals with non-isolated IGT were excluded. **c**, Association of top discriminatory proteins
451 with incident type 2 diabetes was assessed in the Whitehall II study **d**, Association of iIGT protein scores with 8
452 incident cardiometabolic diseases was assessed in a sub-cohort of the EPIC-Norfolk study. OGTT: oral glucose
453 tolerance test, IGT: impaired glucose tolerance, iIGT: isolated impaired glucose tolerance.

454

455 **Figure 2. Performance of LASSO trained models for impaired glucose tolerance (a) and isolated impaired**
456 **glucose tolerance (b) discrimination in the internal validation test set.** **a**, IGT discrimination performance in the
457 independent internal validation test set (N=2881, 192 IGT individuals) for the standard clinical model (Cambridge
458 T2D risk Score + FPG + HbA1c), a 65-protein model and a clinical + 8 protein model. **b**, iIGT discrimination
459 performance in the independent internal validation test set (N=2795, 111 iIGT individuals) for the standard
460 clinical model, a 68-protein model and a clinical + 3 protein model. **c**, Comparison of protein ranking during
461 feature selection for iIGT (N=2795, 111 iIGT individuals) and IGT (N=2881, 192 IGT individuals) top discriminatory
462 proteins. IGT: impaired glucose tolerance, iIGT: isolated impaired glucose tolerance, FPG: fasting plasma glucose,
463 HbA1c: glycated haemoglobin.

464

465 **Figure 3. Proposed 3-stage screening strategy.** In the first stage, individuals in the whole of Fenland were divided
466 into low and high risk according to the Cambridge T2D risk score. The high risk group would undergo a second
467 stage involving measurement of HbA1c and of the 3 iIGT proteins. Individuals with HbA1c levels within the T2D
468 or prediabetic range would be referred for intervention and lifestyle modifications. Individuals with HbA1c below
469 the prediabetic range, would be further stratified using the final clinical + 3 iIGT protein model to identify a high
470 risk group, which on a third stage would be taken forward for OGTT testing to identify iIGT cases that would
471 have otherwise been missed by current screening guidelines. Figure was designed with biorender.com.

472

473 **Figure 4. Characterization of the association between top impaired glucose tolerance and isolated impaired**
474 **glucose tolerance discriminatory proteins and glycaemic traits, future T2D risk and genetic predisposition to**
475 **metabolic phenotypes.** **a**, Association of top IGT and iIGT discriminatory proteins with fasting and 2-hour glucose
476 and insulin in the Fenland study (N = 10259 individuals). Beta estimates with 95% confidence intervals are shown.
477 **b**, Association of top IGT and iIGT discriminatory proteins with incident T2D in the Whitehall II study (N = 1492,
478 521 incident T2D cases). Hazard ratios (HR) with 95% confidence intervals are shown. **c**, Association of genetic
479 risk scores for fasting glucose, fasting insulin, 2-hour plasma glucose, type 2 diabetes and body mass index with
480 top IGT and iIGT discriminatory proteins in the Fenland study (N = 7973 individuals). Beta estimates with a 95%
481 confidence interval are shown. FG: fasting glucose, FI: fasting insulin, 2hPG: 2-hour plasma glucose, 2hPI: 2-hour
482 plasma insulin, T2D: type 2 diabetes, BMI: body mass index.

483

484 **Figure 5. Association of iIGT protein scores with incident cardiometabolic diseases.** Association of iIGT
485 prediction scores (left panel) or individual top iIGT proteins (right panel) with 8 cardiometabolic disease
486 outcomes in a sub-cohort the EPIC-Norfolk study (N=753 individuals). Hazard ratios (HR) with 95% confidence
487 intervals are shown.

488

489

490

491

492

493 **Methods**

494 ***Study Samples***

495 The Fenland study²¹ is a population-based cohort of 12,435 men and women born between 1950 and
496 1975 who underwent detailed phenotyping at the baseline visit from 2005-2015. Participants were
497 recruited from general practice surgeries in Cambridge, Ely and Wisbech (UK). Exclusion criteria of the
498 Fenland study included pregnancy, prevalent diabetes, an inability to walk unaided, psychosis, or
499 terminal illness. The study was approved by the Cambridge Local Research Ethics Committee (NRES
500 Committee – East of England Cambridge Central, ref. 04/Q0108/19) and all participants provided
501 written informed consent. The consent covered measurements made from blood samples as well as
502 extends beyond the baseline examination as described previously²¹.

503 ***Clinical assessment***

504 All participants completed a 2-hour 75 g OGTT following an overnight fast. Blood samples were
505 collected at fasting and 2-hour post glucose load in EDTA tubes for plasma separation by
506 centrifugation. Samples were kept at -80°C until further analysis. Glucose (assayed in a Dade Behring
507 Dimension RxL analyser) and insulin (DELFI[®] immunoassay, Perkin Elmer) concentrations were
508 measured at fasting and 2-hours, as well as lipid profiles (triglycerides, HDL and total cholesterol),
509 alanine aminotransferase (ALT), alkaline phosphatase (ALP), C-reactive protein (CRP) and serum
510 creatinine (assayed in a Dade Behring Dimension RxL analyser) at fasting, and HbA1c (Tosoh
511 Bioscience, TOSOH G7 analyser).

512 IGT and T2D were defined by 2-hour glucose according to IEC diagnosis criteria² as glucose levels
513 between 7.8 and < 11.1 mmol/L (141 and < 199 mg/dL) and ≥ 11.1 mmol/L (≥ 199 mg/dL), respectively.
514 IGT was defined as 2hPG ≥ 7.8 mmol/L and <11.1 mmol/L, post-challenge hyperglycaemia as 2hPG
515 ≥ 11.1 mmol/L, iIGT as individuals with IGT but HbA1c <42mmol/mol (6%) and FG <6.1 mmol/L
516 (<110mg/dL), and isolated post-challenge hyperglycaemia as individuals with post-challenge
517 hyperglycaemia but HbA1c <42mmol/mol and FG <6.1 mmol/L. The number of individuals with post-
518 challenge hyperglycaemia in the diabetic range (i.e., 2hPG ≥ 11.1 mmol/L) was too low to investigate
519 the performance of our models to identify this group of people with undiagnosed T2D biochemically
520 defined solely due to elevated 2-hour glucose. These individuals would still be missed and remain
521 undiagnosed by FG and HbA1c testing. We therefore used the terms IGT and iIGT to refer to all
522 individuals with 2hPG ≥ 7.8 mmol/L throughout text and in order to develop a model that captures all
523 individuals that would remain undiagnosed by current strategies. We note that the thresholds to
524 define glycaemic categories vary across the American Diabetes Association (ADA) , WHO and the
525 International Expert Committee (IEC)⁵¹. We use the IEC HbA1c and FG thresholds to reflect current

526 clinical practice in the UK. We note that using ADA thresholds will likely results in lower case numbers
527 for IGT and iIGT at the cost of a substantially higher false-positive rate. Body mass index (BMI) was
528 calculated as weight (kg) / square of height (m²). Additionally, the homeostasis model assessment of
529 insulin resistance (HOMA-IR) was calculated as FI (μIU/mL) × fasting glucose (mmol/mL)/22.5⁵².
530 Estimated glomerular filtration rate (eGFR) was calculated by the CKD-EPI equation using serum
531 creatinine⁵³.

532 Hepatic steatosis was evaluated by an abdominal ultrasound and images were scored by two trained
533 operators. Criteria used for scoring included: increased echotexture of the liver parenchyma,
534 decreased visualisation of the intra-hepatic vasculature and attenuation of ultrasound beam. A normal
535 liver was considered as a score from 3 – 4, mild steatosis from 5 – 7, moderate steatosis from 8 – 10
536 and sever steatosis ≥ 11 ⁵⁴.

537 Participants completed DEXA scan measurements using a Lunar Prodigy advanced fan beam scanner
538 (GE Healthcare) performed by trained operators using standard imaging, positioning protocols and
539 manually processed according to a standardized procedure described previously³⁵. Abdominal visceral
540 and subcutaneous fat mass was estimated using the DEXA software.

541 Differences in clinical characteristics were evaluated by ANOVA followed by posthoc Tukey test, or χ²
542 for categorical variables. Non-normally distributed variables were log transformed when appropriate.

543 ***Proteomic profiling of the Fenland cohort***

544 Proteomic profiling was done using an aptamer-based technology (SomaScan proteomic assay).
545 Fasting proteomic profiling was done in participants from the Fenland cohort at baseline, from which
546 relative abundancies of 4,775 unique protein targets (evaluated by 4,979 SOMAmer reagents,
547 SomaLogic v4)^{18,55} was evaluated in EDTA plasma. Briefly, proteins are targeted by modified single
548 stranded DNA sequences (aptamers). Concentration is then approximated as relative fluorescence
549 units using a DNA microarray ⁵⁶.

550 To account for variation in hybridization within runs, hybridization control probes are used to generate
551 a hybridization scale factor for each sample. To control for total signal differences between samples
552 due to variation in overall protein concentration or technical factors such as reagent concentration,
553 pipetting or assay timing, we used the adaptive median normalisation (AMN), unless stated otherwise.
554 Briefly, a ratio between each aptamer's measured value and a reference value from an external
555 reference population is computed, and the median of these ratios is computed for each of the three
556 dilution sets (20%, 1% and 0.005%) and applied to each dilution set to shift the intrapersonal
557 distribution of protein intensities accordingly to match the reference population. We removed
558 samples if they did not meet an acceptance criterion for scaling factors with values outside of the

559 recommend range (0.25-4) or were flagged as technical failures (n=19). Detailed SomaLogic's
560 normalization, calibration data, and quality control processes have been previously described in
561 detail¹⁸. At a protein level, we took only human protein targets forward for subsequent analysis (4,979
562 out of the 5284 aptamers). Intraassay coefficients of variation (calculated based on raw fluorescence
563 units) had a median of 4.98% (interquartile range 3.87% - 6.99%) suggesting good quality measures
564 for the vast majority of protein targets. We decided to not apply any other filters to individual protein
565 qualities given that even poorly measured proteins might be informative and left it to the restrictive
566 feature selection approach applied to drop uninformative proteins, including possibly poorly
567 measured once. Aptamers' target annotation and mapping to UniProt accession numbers as well as
568 Entrez gene identifiers were provided by SomaLogic and we used those to obtain genomic positions
569 of protein encoding genes.

570 ***Genome wide genotyping and imputation***

571 Fenland participants were genotyped using three genotyping arrays: the Affymetrix UK Biobank Axiom
572 array (OMICs, N=8994), Illumina Infinium Core Exome 24v1 (Core-Exome, N=1060) and Affymetrix
573 SNP5.0 (GWAS, N=1402). Samples were excluded for the following reasons: 1) failed channel contrast
574 (DishQC <0.82); 2) low call rate (<95%); 3) gender mismatch between reported and genetic sex; 4)
575 heterozygosity outlier; 5) unusually high number of singleton genotypes or 6) impossible identity-by-
576 descent values. Single nucleotide polymorphisms (SNPs) were removed if: 1) call rate < 95%; 2) clusters
577 failed Affymetrix SNPolisher standard tests and thresholds; 3) MAF was significantly affected by plate;
578 4) SNP was a duplicate based on chromosome, position, and alleles (selecting the best probe set
579 according to Affymetrix SNPolisher); 5) Hardy-Weinberg equilibrium $p < 10^{-6}$; 6) did not match the
580 reference or 7) MAF=0.

581 Autosomes for the OMICS and GWAS subsets were imputed to the HRC (r1) panel using IMPUTE4, and
582 the Core-Exome subset and the X-chromosome (for all subsets) were imputed to HRC.r1.1 using the
583 Sanger imputation server⁵⁷. All three arrays subsets were also imputed to the UK10K+1000Gphase3⁵⁸
584 panel using the Sanger imputation server in order to obtain additional variants that do not exist in the
585 HRC reference panel. Variants with MAF < 0.001, imputation quality (info) < 0.4 or Hardy Weinberg
586 Equilibrium $p < 10^{-7}$ in any of the genotyping subsets were excluded from further analyses.

587 **Statistical Analyses**

588 ***Classification of IGT and iIGT from the fasting proteome***

589 To identify and validate a proteomic signature able to discriminate IGT and iIGT (as a binary outcome),
590 the entire Fenland study (N=11,546 without missing data for 2hPG), was divided into three subsets:
591 for feature selection (50%, N = 5773), parameter optimization (25%, N=2887) and validation (25%,

592 N=2881). IGT and iIGT cases were split equally into 50% for training ($N_{IGT} = 387$, $N_{iIGT} = 222$), 25 % for
593 optimization ($N_{IGT} = 194$, $N_{iIGT} = 111$) and 25% for testing ($N_{IGT} = 193$, $N_{iIGT} = 111$) sets. For these
594 analyses, SOMAmer RFUs were \log_{10} -transformed. Feature selection was carried out by least absolute
595 shrinkage and selection operator (LASSO) regression. We chose to use LASSO because it was the most
596 suitable model to 1) identify the smallest possible set of independent predictors, 2) it is
597 computationally efficient, which allowed us to implement a robust framework using bootstrap
598 resampling to identify a core set of most informative predictors and 3) it is less prone to overfitting.
599 To address case-control imbalance we used the ROSE R package⁵⁹, which implements down-sampling
600 of the majority class (controls) along with synthetic new data points for the minority class (IGT or iIGT).
601 A nested 10-fold cross-validation (inner loop to determine regularization parameter, λ) was done over
602 100 bootstrap samples (outer loop) drawn from the feature selection set. Each protein received a
603 score that was generated by counting the number of times it was included in the final model from
604 each of the 100 bootstrap samples, that is, the score was between 0 (for proteins that were never
605 selected in the final model) and 100 (for proteins that were selected in the final model in all bootstrap
606 samples). We ranked the proteins based on their score to identify the most informative set of features
607 (i.e. with a higher score) (**Supplementary Fig. 1**). This was implemented by the use of the R packages
608 *caret*⁶⁰ and *glmnet*⁶¹. Proteins selected in the final model in more than 80%, 90%, and 95% of the
609 bootstrap samples, were tested as predictors and taken forward for parameter optimization by 10-
610 fold cross validation of the model by LASSO regression in the optimization set. Additional models were
611 optimized by LASSO regression, such as a standard patient information-based model using the
612 variables from the Cambridge Diabetes Risk Score (age, sex, family history of diabetes, smoking status,
613 prescription of steroid or antihypertensive medication and BMI)²², a standard clinical model (including
614 the variables from the Cambridge Diabetes risk Score, FG and HbA1c) and a standard clinical plus the
615 selected proteins model. Clinical predictors were forced to be kept in the clinical plus proteins model
616 by setting the penalty factors of these variables to 0. For comparison, ridge regression (which will keep
617 all proteins in the final model) was used to build a prediction model using all the 4979 proteins as
618 predictors.

619 Performance of the classification models were evaluated in the internal independent validation set,
620 which was never used for training and optimization. The prediction models' discriminatory power was
621 assessed by computing the area under the receiver operating curve (AUROC). Confidence intervals
622 and p-values (using the deLong method implemented by the R package pROC⁶²) were computed for
623 the comparison between the ROC curves for the standard clinical model and clinical with added
624 proteins model. Additionally, models' net reclassification index was evaluated using the R package
625 PredictABEL⁶³.

626 Using an analogous machine learning strategy, we developed models for iIGT discrimination. For these
627 analyses, all individuals with non-isolated IGT (2hPG > 7.8 mmol/L, FPG > 6.1 mmol/L and HbA1c > 42
628 mmol/mol) were excluded from the cohort (leaving N = 11,281), which was subsequently divided into
629 feature selection (50%, N = 5591), parameter optimization (25%, N=2796) and validation (25%,
630 N=2795). Feature selection, optimization and testing were carried out as described for IGT models. To
631 achieve comparable model performance with the minimal number of predictors, we used recursive
632 feature elimination on the set of proteins selected in >95% of boots during feature selection. As a
633 sensitivity analysis, we performed the same framework described above, that is, feature selection,
634 parameter optimization and validation to assess model performance when using protein data
635 reversing the final normalisation step that is unique to the SomaScan platform. We note that using
636 'non-normalised' proteomic data led to broadly comparable results, which are well in the margins of
637 random variation of protein measurements in general, albeit with some difference in the proteins
638 selected as the most predictive markers in the final models (**Supplementary Table 17**).

639 Calibration of the final models was assessed in the internal validation set by computing the calibration
640 slope, which evaluates the spread of the estimated risks and has a target value of one. Calibration
641 slopes less than 1 indicate extreme estimated risks while slopes greater than 1 indicate very moderate
642 risk estimates. Calibration slopes were computed using the R package rms⁶⁴.

643 The number needed to screen (NNS) was calculated using a staged screening scenario. Firstly,
644 participants from the Fenland study were stratified by predicted probabilities from the Cambridge T2D
645 risk Score, that is, non-invasive risk factors that could be obtained by interviewing the patient. The
646 threshold used to stratify individuals into "high" and "low" risk strata according to their predicted
647 probabilities was set to optimize a balance between the total number of individuals that would be
648 needed to screen and sensitivity (as would be appropriate for such a screening setting), which was
649 achieved at 0.7, regardless of specificity. On second instance, participants within the high-risk group
650 were further stratified by HbA1c levels, using IEC cut-offs (normoglycaemic : HbA1c < 42 mmol/mol,
651 prediabetic criteria: HbA1c >= 42 mmol/mol and < 48 mmol/mol, T2D criteria : HbA1c >= 48
652 mmol/mol)⁵¹. On third instance, participants whose HbA1c did not meet the criteria for T2D or
653 prediabetes (that is, normoglycaemic as defined above), were further stratified according to the
654 clinical + 3- iIGT protein model. Similarly, a threshold that optimized testing as few individuals as
655 possible while retaining good sensitivity of 0.7 was set for this model (**Supplementary Table 10**). We
656 estimated the NNS within this stratum compared to the NNS within the full set of individuals with
657 HbA1c in the normoglycaemic range. The NNS was calculated as the total number of individuals within
658 the group divided by the number iIGT cases within the same group and refers to the number of OGTTs

659 that would need to be done to identify one iIGT case within the group of interest. We additionally
660 estimated the NNS in the test set only, as a sensitivity analysis.

661 ***IGT/iIGT model validation and follow-up analyses in the WHII study***

662 The Whitehall II study is a longitudinal, prospective cohort study³⁹ that was approved by the joint
663 University College London / University College London Hospital's Committees on the Ethics of Human
664 Research. Proteomic profiling of fasting EDTA-plasma samples was done for all individuals at phase 5
665 (from 1997 - 1999) with the SomaScan v4.1 proteomic assay. We performed validation of the IGT and
666 iIGT clinical + protein models at phase 5 (from 1997 - 1999) of the study, where proteomic profiling and
667 OGTT values were available. Since HbA1c was not measured at phase 5 of the study, we defined iIGT
668 as 2hPG > 7.8 mmol/L and FPG > 6.1 mmol/L. We used the weights from the models trained in Fenland
669 to evaluate their performance in WHII phase 5 (total sample size = 5058, N_{IGT}= 693, N_{iIGT}=617) for the
670 baseline clinical model (Cambridge T2D risk score + FG) and the baseline clinical + protein iIGT and IGT
671 models (3 and 8 proteins respectively).

672 For the association between top discriminatory proteins and incident T2D in the Whitehall II study
673 individuals were selected as a nested case-control study design in which proteomic profiling of fasting
674 EDTA-plasma samples was done at phase 5 (from 1997 - 1999) with the SomaScan v4 proteomic assay.
675 Incident T2D occurrence was assessed in repeated clinical examinations in 1997-1999, 2002-2004,
676 2007-2009, 2012-2013, and 2015-2016, based on FPG above 7 mmol/L, HbA1c>6.5%, use of diabetes
677 medication, or reported physician diagnosed diabetes, excluding prevalent T2D cases at baseline from
678 the analysis. Additionally, participants with impaired kidney function (eGFR < 30 mL/min/1.73m²),
679 incident cardiovascular diseases or missing data on T2D at follow-up were excluded. The final sample
680 comprised of 521 cases and 971 controls.

681 Association between fasting candidate proteins and incident T2D was assessed using Cox-proportional
682 hazards regression adjusting for the baseline confounders age, sex and BMI. We tested a second
683 model adjusting for additional baseline confounders including FG, triglycerides, HDL-cholesterol and
684 lipid lowering medication on top of age, sex and BMI to determine whether the association persisted
685 in a more refined model.

686 **Effect of fasting status on plasma levels of IGT and iIGT discriminatory proteins**

687 Fourteen adult participants were recruited to participate in the study and provided informed consent
688 appropriately. Participants were asked to fast overnight for at least 12 hours prior to reporting to the
689 study site. Fasting blood samples were collected from each participant, after which they were given a
690 moderate fat meal consisting of 5-8 ounces of Cheerios with 6 ounces of 2% milk, one egg, one slice

691 of bacon, one slice of toast with margarine, and 4 ounces of orange juice (calories: 450, 16.9 grams of
692 fat, 16 grams of protein, and 59 grams of carbohydrates)⁶⁵.

693 The time for each participant to complete the meal ranged from 7 to 19 minutes (average of 16
694 minutes). Post prandial blood samples were collected at 0.5, 1, and 3 hours following completion of
695 the meal. Since each participant consumed their meals at different rates, the actual blood collection
696 times post meal does vary between participants. Participants were not allowed to eat or drink any
697 further caloric items until after the last blood collection. Twelve participants (6 male and 6 female)
698 completed the study. Two participants were excluded due to unmet fasting requirements and an
699 adverse reaction during the first blood draw.

700 Blood samples were processed to obtain EDTA-plasma by centrifugation and frozen at -80°C until
701 delivered to SomaLogic Sample Management for proteomic profiling using the SomaScan v4 assay.
702 The effect of fasting status on 9 unique SOMAmer reagents included in the final clinical + protein
703 models for IGT or iIGT, was tested by repeated measures ANOVA. Proteins with ANOVA p-values <
704 0.0055 (according to Bonferroni adjustment for 9 comparisons) were deemed to be significantly
705 affected by fasting status.

706 ***Functional annotation of IGT and iIGT-protein signatures***

707 Functional annotation of the 65-IGT and 68-iIGT protein signatures was performed using modified
708 Fisher's exact tests as implemented by the Database for Annotation, Visualization and Integrated
709 Discovery (DAVID, version 6.8) and enrichment of biological process GO terms (GOTERM_BP_DIRECT)
710 was analysed, setting the full list of proteins evaluated by the SomaLogic platform as the background.

711 ***Variance explained in top discriminatory protein levels by clinical, biochemical, anthropometric and*** 712 ***behavioural risk factors***

713 The proportion of variance explained in candidate protein levels by several variables was evaluated in
714 the Fenland cohort using the *variancePartition* R package⁶⁶. Analogously, the proportion of variance
715 explained in the first principal component of the 65-IGT and 68-iIGT discriminatory protein signatures
716 was evaluated. Briefly, this package fits a linear mixed model to assess the effect of each variable on
717 the outcome while correcting for all other variables. Variables evaluated were age, sex, IGT, IPCH, FPG,
718 2hPG, FI, 2hPI, HbA1c, total triglycerides, total cholesterol, HDL-cholesterol, LDL-cholesterol, ALT, ALP,
719 a liver score, BMI, waist-to-hip ratio (WHR), amount of subcutaneous fat, amount of visceral fat, CRP,
720 estimated glomerular filtration rate (eGFR) and intake of statins or antihypertensive medication. FPG,
721 2hPG, FI, 2hPI, HbA1c, total triglycerides, ALT, ALP, CRP, subcutaneous fat and visceral fat were natural
722 log-transformed due to skewed distribution of these variables. We fit separate models for each of the
723 variables evaluated adjusting only for age and sex in the entire Fenland cohort (N=11,546) to avoid

724 bias due to strong collinearity among variables tested. For each of the models, participants with
725 missing data were excluded.

726 ***Protein quantitative trait loci (pQTLs) for candidate proteins***

727 Genetic variants associated with candidate proteins (protein quantitative trait loci or pQTLs) were
728 taken from our genome-wide association studies across all aptamers as described in Pietzner et al,
729 2021⁵⁵.

730 ***Percentage of variance explained in protein levels by cis and trans pQTL scores***

731 Polygenic scores were constructed for pQTLs within the *cis* (within ± 500 kb of the protein-encoding
732 gene) and *trans* regions. Cis-pQTL scores were built using conditionally independent variants. The
733 percentage of variance explained in protein levels by the cis and trans-scores was computed as
734 described in the above section adjusting for age and sex.

735 ***Association between top discriminatory proteins and fasting and 2-hour plasma glucose and insulin***

736 Observational associations between the top selected IGT and iIGT discriminatory proteins and FPG, FI,
737 2hPG and 2hPI were assessed in the entire Fenland cohort at baseline (N=10,259 without missing data)
738 by linear regression models adjusting for age, sex, BMI and test site from the study. The models for
739 2hPG and 2hPI were additionally adjusted by FPG and FPG + FI, respectively. Protein levels were log₁₀-
740 transformed and standardized, and 2hPG and 2hPI values were log-transformed for these analyses.
741 Proteins were considered significant at a Bonferroni threshold (p-values < 0.001, accounting for
742 comparisons between the number of protein and number of traits, as for all further association
743 analyses).

744 ***Association between polygenic risk scores for glycaemic traits and top discriminatory proteins***

745 T2D³⁶, fasting glucose (FG)³⁴, fasting insulin³⁴ (FI score), 2hPG³⁴ (2hPG score) and BMI³⁵ polygenic
746 scores, weighted by genetic effect sizes of previously reported genome-wide significant variants, were
747 computed for 7,973 Fenland participants genotyped with the same array (Affymetrix UK Biobank
748 Axiom Array). Variants not available, with low imputation quality scores < 0.6, or with strand
749 ambiguous alleles were excluded from the scores. Each polygenic score was tested for associations
750 with the plasma abundancies of top IGT and iIGT discriminatory proteins by linear regression models
751 adjusting for age, sex, BMI, the first 10 genetic principal components and test site of the study.

752 ***Association between iIGT scores with incident cardiometabolic diseases in a sub-cohort of the EPIC- 753 Norfolk study***

754 The EPIC-Norfolk study is a cohort of 25,639 middle-aged, individuals from the general population of
755 Norfolk a county in Eastern England which is a component of EPIC³⁷. The EPIC-Norfolk study was

756 approved by the Norfolk Research Ethics Committee (ref. 05/Q0101/191); all participants gave their
757 informed written consent before entering the study. All participants were flagged for mortality at the
758 UK Office of National Statistics and vital status was ascertained for the entire cohort. Death certificates
759 were coded by trained nosologists according to the International Statistical Classification of Diseases
760 and Related Health Problems, 10th Revision (ICD-10). Hospitalization data were obtained using
761 National Health Service numbers through linkage with NHS Digital. Participants were identified as
762 having experienced an event if the corresponding ICD-10 code was registered on the death certificate
763 (as the underlying cause of death or as a contributing factor) or as the cause of hospitalization
764 (**Supplementary Table 15**). Since the long-term follow-up of EPIC-Norfolk comprised the ICD-9 and
765 ICD-10 coding system, codes were consolidated. The current study is based on follow-up to 31 March
766 2016. Information on lifestyle factors and medical history was obtained from questionnaires as
767 reported previously³⁷. The current analysis is based on a random sub-cohort (N=875) of the whole
768 EPIC-Norfolk study population that was selected excluding known prevalent case subjects of diabetes
769 at baseline was using the same definitions as used in the InterAct Project⁶⁷; in which proteomic
770 profiling was done at health check 1 using the SOMAscan v4 platform from citrate-plasma samples
771 stored in liquid nitrogen since the baseline visit.

772 Participants with missing data for any of the variables included in the final prediction models
773 developed in the Fenland study were excluded. The final sample comprised of 753 individuals for
774 which characteristics are presented in **Supplementary Table 16**.

775 Final prediction models trained and optimized for iIGT in the Fenland study were used to calculate the
776 predicted probability of iIGT for each participant at health check 1 in this sub-cohort of the EPIC-
777 Norfolk study. Models tested included: the clinical + 3-proteins iIGT model, 3-protein iIGT model (95%
778 feature selection protein set model), 68-protein iIGT model (80% feature selection protein set model)
779 and the clinical model as a baseline comparison. We then tested the association of the predicted iIGT
780 probability with 8 incident cardiometabolic diseases (or associated T2D comorbidities) including type
781 2 diabetes, coronary heart disease, heart failure, peripheral artery disease, cerebral stroke, liver
782 disease, renal disease and cataracts using cox proportional hazards models adjusting by age at
783 baseline and sex (except for the clinical + 3 protein model, which already accounted for these risk
784 factors within the score). Associations were deemed significant at an 5% FDR accounting for
785 comparison between 8 diseases.

786 We aimed for cross-platform validation in a separate random sub-cohort of the prospective EPIC-
787 Norfolk study (N=771), in which proteomic measures were done with the Olink Explore panel⁴⁰ from
788 serum samples. Participants with missing data for any of the variables included in the final prediction

789 models developed in the Fenland study (except HbA1c which was excluded from the models as it was
790 unavailable in a large proportion of participants from this sub-cohort) were excluded. The final sample
791 comprised of 602 individuals for which characteristics are presented in **Supplementary Table 18**.

792 Final prediction models trained and optimized for iIGT in the Fenland study (using SomaScan) were
793 used to calculate the predicted probability of iIGT for each participant at health check 1 in this sub-
794 cohort of the EPIC-Norfolk study, using the Olink measures for the proteins. Models tested included:
795 the clinical + 3-proteins iIGT model, 3-protein iIGT model (95% feature selection protein set model)
796 and the Cambridge T2D risk Score. We then tested the association of the predicted iIGT probability
797 with the same Cox-model setting and set of disease as in the sub-cohort with available SomaLogic
798 measurements except for liver disease (**Supplementary Table 19**). Associations were deemed
799 significant at an 5% FDR accounting for comparison between 7 diseases.

800 All statistical analyses were performed using R language, and environment for statistical computing
801 (version 3.6.1 and 4.1.0, R Core Team).

802

803 **Data availability**

804 Data access for the Fenland and EPIC studies can be requested by bona fide researchers for specified
805 scientific purposes through a simple application process via the study websites below. Data will either
806 be shared through an institutional data sharing agreement or arrangements will be made for analyses
807 to be conducted remotely without the necessity for data transfer.

808 Fenland: <https://www.mrc-epid.cam.ac.uk/research/studies/fenland/information-for-researchers>

809 EPIC-Norfolk: <https://www.mrc-epid.cam.ac.uk/research/studies/epic-norfolk>

810

811 **Code availability**

812 The code employed for the machine learning developed framework has been deposited in the
813 following repository: https://github.com/MRC-Epid/iigt_prediction_proteomics.

814

815

816

817

818 **Methods-only references**

819

- 820 51. Lee, C.M.Y., *et al.* Comparing different definitions of prediabetes with subsequent risk of
821 diabetes: an individual participant data meta-analysis involving 76 513 individuals and 8208
822 cases of incident diabetes. *BMJ Open Diabetes Res Care* **7**, e000794 (2019).
- 823 52. Fukagawa, N.K., *et al.* Insulin-mediated reduction of whole body protein breakdown. Dose-
824 response effects on leucine metabolism in postabsorptive men. *J Clin Invest* **76**, 2306-2311
825 (1985).
- 826 53. Inker, L.A., *et al.* Estimating glomerular filtration rate from serum creatinine and cystatin C.
827 *The New England journal of medicine* **367**, 20-29 (2012).
- 828 54. Mehta, S.R., Thomas, E.L., Bell, J.D., Johnston, D.G. & Taylor-Robinson, S.D. Non-invasive
829 means of measuring hepatic fat content. *World J Gastroenterol* **14**, 3476-3483 (2008).
- 830 55. Pietzner, M., *et al.* Mapping the proteo-genomic convergence of human diseases. *Science*,
831 eabj1541 (2021).
- 832 56. Rohloff, J.C., *et al.* Nucleic Acid Ligands With Protein-like Side Chains: Modified Aptamers and
833 Their Use as Diagnostic and Therapeutic Agents. *Mol Ther Nucleic Acids* **3**, e201 (2014).
- 834 57. McCarthy, S., *et al.* A reference panel of 64,976 haplotypes for genotype imputation. *Nat*
835 *Genet* **48**, 1279-1283 (2016).
- 836 58. Huang, J., *et al.* Improved imputation of low-frequency and rare variants using the UK10K
837 haplotype reference panel. *Nat Commun* **6**, 8111 (2015).
- 838 59. Nicola Lunardon, G.M., Nicola Torelli. ROSE: a Package for Binary Imbalanced Learning. *The R*
839 *Journal* **6**, 79-89 (2014).
- 840 60. Kuhn, M. Building Predictive Models in R Using the caret Package. *J Stat Softw* **28**, 1-26 (2008).
- 841 61. Friedman, J., Hastie, T. & Tibshirani, R. Regularization Paths for Generalized Linear Models via
842 Coordinate Descent. *J Stat Softw* **33**, 1-22 (2010).
- 843 62. Robin, X., *et al.* pROC: an open-source package for R and S+ to analyze and compare ROC
844 curves. *BMC Bioinformatics* **12**, 77 (2011).
- 845 63. Kundu, S., Aulchenko, Y.S., van Duijn, C.M. & Janssens, A.C. PredictABEL: an R package for the
846 assessment of risk prediction models. *European journal of epidemiology* **26**, 261-264 (2011).
- 847 64. Jr, F.E.H. rms: Regression Modeling Strategies. R package version 5.1-1. (2017).
- 848 65. *Pharmacokinetics in Drug Development: Clinical Study Design and Analysis* (American
849 Association of Pharmaceutical Scientists, Arlington, 2004).
- 850 66. Hoffman, G.E. & Schadt, E.E. variancePartition: interpreting drivers of variation in complex
851 gene expression studies. *BMC Bioinformatics* **17**, 483 (2016).
- 852 67. InterAct, C., *et al.* Design and cohort description of the InterAct Project: an examination of the
853 interaction of genetic and lifestyle factors on the incidence of type 2 diabetes in the EPIC
854 Study. *Diabetologia* **54**, 2272-2282 (2011).

855

1 **Proteomic signatures for identification of impaired glucose tolerance**

2

3 Julia Carrasco-Zanini¹, Maik Pietzner^{1,2}, Joni V Lindbohm^{3,4}, Eleanor Wheeler¹, Erin Oerton¹, Nicola
4 Kerrison¹, Missy Simpson⁵, Matthew Westacott⁵, Dan Drolet⁵, Mika Kivimaki^{3,4}, Rachel Ostroff⁵,
5 Stephen A Williams⁵, Nicholas J Wareham¹, Claudia Langenberg^{1,2,6*}

6

7 **Affiliations**

8 1. MRC Epidemiology Unit, Institute of Metabolic Science, University of Cambridge School of Clinical
9 Medicine, Cambridge, UK

10 2. Computational Medicine, Berlin Institute of Health at Charité – Universitätsmedizin Berlin, Germany

11 3. Clincum, Department of Public Health, University of Helsinki, P.O. Box 41, FI-00014 Helsinki, Finland

12 4. Department of Epidemiology and Public Health, University College London, UK

13 5. SomaLogic, Inc., Boulder, CO, USA

14 6. Precision Healthcare University Research Institute, Queen Mary University of London, UK

15

16 *Correspondence to Dr Claudia Langenberg (claudia.langenberg@mrc-epid.cam.ac.uk)

17

18

19

20 **Abstract**

21 The implementation of recommendations for type 2 diabetes (T2D) screening and diagnosis focus on
22 measurement of HbA1c and fasting glucose. This approach leaves a large number of individuals with
23 isolated impaired glucose tolerance (iIGT), who are only detectable through oral glucose tolerance
24 tests (OGTTs), at risk of diabetes and its severe complications. We applied machine learning to
25 proteomic profiles of a single fasted sample from 11,546 participants of the Fenland study to test
26 discrimination of iIGT defined using gold standard OGTTs. We observed significantly improved
27 discriminative performance by adding only three proteins (RTN4R, CBPM, and GHR) to the best clinical
28 model (0.80 (0.79-0.86), $p=0.004$), which we validated in an external cohort. Increased plasma levels
29 of these candidate proteins were associated with an increased risk for future T2D in an independent
30 cohort and were also increased in individuals genetically susceptible to impaired glucose homeostasis
31 and T2D. Assessment of a limited number of proteins can identify individuals likely to be missed by
32 current diagnostic strategies and at high risk of T2D and its complications.

33

34 **Introduction**

35 Current clinical guidelines for type 2 diabetes (T2D) screening and diagnosis are based on glycated
36 haemoglobin (HbA1c) and fasting glucose (FG) levels for reasons of practicality, however alternative
37 tests can be used^{1,2}. Globally, over 7.5% of adults have impaired glucose tolerance (IGT)³ with
38 increased prevalence reported in older individuals⁴ and specific ethnic groups, such as people from
39 Southeast Asia⁵. A substantial proportion of people with IGT (28 – 86%)⁶⁻⁸ can only be identified
40 through oral glucose tolerance tests (OGTTs), which are inconvenient and time-consuming. Individuals
41 with isolated IGT (iIGT), that is, 2-hour plasma glucose (2hPG) ≥ 7.8 and < 11.1 mmol/L but normal
42 HbA1c and fasting glucose, remain undetected by current T2D detection strategies⁹⁻¹² but are at very
43 high risk of developing diabetes (annualized T2D relative risk of 5.5 compared to normoglycemic
44 individuals)¹³ and presenting with its severe micro- and macrovascular complications^{9-12,14}. Compared
45 to individuals with fasting hyperglycaemia, mortality is twice as high in the iIGT group over a period of
46 5 to 12 years^{15,16}.

47 Small proof-of-concept studies in cohorts of high-risk individuals have demonstrated the value of deep
48 molecular profiling for early identification of pathways that are differentially regulated between
49 individuals with and without insulin resistance^{17,18} and to guide its prediction¹⁹. Deep profiling of the
50 plasma proteome at population scale has become possible through aptamer-based affinity assays²⁰.
51 The systematic study of the circulating proteome promises to improve strategies for prediction and
52 diagnosis¹⁸ as well as aetiological understanding, including identification of novel pathways leading to
53 T2D and refinement of aetiological subtypes.

54 Because of the high global prevalence of IGT and iIGT, their severe complications, and the currently
55 unmet need of screening strategies that can identify iIGT without a challenge test, we used machine
56 learning to test whether large-scale proteomic profiling of a single fasted sample could identify
57 individuals with iIGT and improve current clinical models. We then tested whether the most
58 discriminatory proteins were affected by fasting status, to assess the feasibility of using non-fasted
59 samples to identify iIGT. To gain insights into IGT and iIGT aetiology, we 1) identified and characterised
60 biochemical, phenotypic, and anthropometric drivers of discriminatory proteins, 2) investigated
61 whether their plasma levels were associated with the risk of future T2D in an independent prospective
62 cohort with 521 incident T2D cases, and 3) tested the influence of genetic susceptibility to T2D or
63 related phenotypes on protein levels.

64 **Results**

65 We used an aptamer-based assay to target 4,775 distinct fasting plasma proteins by 4,979 aptamers
66 in 11,546 participants (5,389 men and 6,157 women) without diagnosed diabetes from the

67 contemporary Fenland study²¹ (baseline visit in 2005-2015, mean age 48.5 years (7.5 s.d.),
68 **Supplementary Table 1**), as previously described¹⁸ (Methods). Participants completed a 75-gram
69 OGTT (**Figure 1a**). We defined isolated post challenge hyperglycaemia as 2hPG ≥ 7.8 mmol/L but HbA1c
70 < 42 mmol/mol and FG < 6.1 mmol/L. This definition captured all participants with isolated IGT (2hPG
71 7.8-11.1 mmol/L but HbA1c < 42 mmol/mol and FG < 6.1 mmol/L) as well as participants with isolated
72 post-challenge hyperglycaemia in the diabetic range (2hPG ≥ 11.1 mmol/L but HbA1c < 42 mmol/mol
73 and FG < 6.1 mmol/L, N=117), i.e. high-risk individuals missed by standard FG and HbA1c testing. For
74 simplicity, we refer from here on to IGT (or iIGT) for all individuals with 2hPG ≥ 7.8 mmol/L, without
75 specifically distinguishing post-challenge hyperglycaemia ≥ 11.1 mmol/L. We used a least absolute
76 shrinkage and selection operator (LASSO) regression framework implemented as a three-step
77 approach, including independent feature selection (50% sample size), optimization (25%) and
78 validation (25%) to discriminate IGT (prevalence 6.7%) and iIGT (3.9%) based on fasting assessment of
79 4,775 proteins (targeted by 4,979 aptamers) (**Figure 1b**). We defined highly discriminatory proteins as
80 those selected in $> 80\%$, 90% , or 95% of random subsamples of the study population during feature
81 selection (**Extended Data Figure 1**).

82

83 ***Proteomic signatures to discriminate IGT and IIGT***

84 We identified 65 and 68 proteins, respectively, that achieved an area under the receiver operating
85 characteristic curve (AUROC) (95% confidence interval) of 0.83 (0.80 – 0.86) and 0.77 (0.72 – 0.81) for
86 discrimination of IGT and iIGT in the independent validation set (**Extended Data Figure 2,**
87 **Supplementary Tables 2 and 3**). This represented a significantly better predictor when compared to
88 the performance of a T2D genetic risk score (T2D-GRS, AUROC_{IGT} 0.58 (0.52 – 0.63), AUROC_{iIGT} 0.54
89 (0.49 – 0.60)) (**Figure 2a and b,** and **Extended Data Figure 3**). Protein-based models further
90 outperformed the standard patient information-based model (based on the Cambridge Diabetes Risk
91 Score including age, sex, family history of diabetes, smoking status, prescription of steroid or
92 antihypertensive medication and body mass index (BMI))²² (AUROC_{IGT} = 0.71 (0.67 – 0.75); AUROC_{iIGT}
93 = 0.71 (0.66 – 0.76)) and the standard clinical model that additionally included blood test results, that
94 is, FPG and HbA1c (AUROC_{IGT}= 0.78 (0.74 – 0.82); AUROC_{iIGT}=0.75 (0.70 – 0.80)) (**Figure 2a and b,**
95 **Supplementary Table 4**).

96 Considering a limited set of the most informative proteins that were identified by the feature selection
97 framework (Methods), discrimination was still superior to the standard clinical model adding only 8
98 proteins for IGT (AUROC_{IGT} 0.83 (0.80 – 0.86), p-value = 4.13×10^{-5} , **Figure 2a, Supplementary Table**
99 **2**) and 3 proteins for iIGT (AUROC_{iIGT} 0.80 (0.76 – 0.85), p-value = 0.004, **Figure 2b, Supplementary**
100 **Table 3**), including 2 proteins (Reticulon-4 receptor, Carboxypeptidase M) selected for both (**Figure**
101 **2c, Supplementary Table 5**²³⁻³³). The weights for the variables included in these final models are
102 available in **Supplementary Table 6**. We observed significant improvement over and above the clinical
103 model of similar magnitude in the independent Whitehall II (WHII) study (**Supplementary Table 7 and**
104 **8, Extended Data Figure 4**).

105 To identify participants with iIGT and IGT, respectively, we choose a cut-off for the clinical + protein
106 model that optimized sensitivity (recall) at 0.70 and 0.71, which yielded a positive predicted value
107 (precision) of 0.20 and 0.13, respectively. The net reclassification index was higher for the final iIGT
108 model (14.5%) compared to IGT (6.5%), consistent with the current lack of informative predictors.

109 Of the 9 distinct proteins included in the 2 final models, 8 were not significantly affected by fasting
110 status (**Methods**) with maximum postprandial fold changes ranging between 0.07 and 0.16; only
111 HTRA1 showed some evidence of a post-prandial increase (maximum fold change= 0.15, p-
112 value=0.004, **Supplementary Table 9**).

113 Finally, we tested model performance de novo omitting the 3 most informative proteins to predict
114 iIGT. The novel model included 7 proteins and still performed significantly better than the best clinical
115 model (AUROC = 0.78 (0.73 – 0.83), p-value = 0.04, **Extended Data Figure 5**). This finding illustrates

116 redundancy in the protein biomarkers available to select from for iIGT prediction, providing practical
117 benefits for clinical implementation, for example with regard to flexibility of prioritising choice of
118 proteins more easily targeted by clinical chemistry assays, least affected by fasting status or sample
119 handling.

120 *Proteomically informed screening strategies*

121 We calculated the numbers needed to screen (NNS) to determine how many OGTTs would need to be
122 performed to identify one participant with iIGT using a three-stage screening approach (**Figure 3**). We
123 stratified all Fenland individuals based on the patient-derived information model in the first instance
124 and based on their HbA1c levels and the 3-protein iIGT model in the second instance (**Methods**).
125 According to current guidelines², individuals at high predicted risk based on the patient-derived
126 information model, but HbA1c levels below cut-offs for prediabetes or T2D² would not be considered
127 for further testing (N Fenland = 4163, NNS = 14, **Figure 3**). Applying the clinical + 3-protein iIGT model
128 on this group enabled identification of a high-risk subgroup (N = 1739) in which application of an OGTT
129 should be considered, since the NNS was only 7 to identify one additional individual with iIGT (**Figure**
130 **3, Supplementary Table 10**). Hence, our proposed approach identified an additional >30% of
131 individuals that would be reclassified (as having prediabetes) and could be offered preventative
132 interventions, that is, a substantial proportion of high-risk individuals that would otherwise be missed
133 by current strategies. To test for potential bias in the NNS estimates arising from overfitting, we
134 applied the same screening algorithm in the test set only, which provided internal validation for the
135 estimates and results from the entire Fenland set (**Extended Data Figure 6**).

136 *Characterisation of discriminatory proteins*

137 To investigate whether increased genetic risk of diabetes and related metabolic risk factors affect
138 abundances of the identified proteins, we compared their differences in individuals with higher versus
139 lower genetic risk based on genetic risk scores (GRS) for T2D and related endophenotypes, including
140 fasting glucose³⁴, fasting insulin³⁴, 2hPG³⁴, body mass index (BMI)³⁵ and T2D³⁶, using linear regression
141 models. We found evidence of significant, directional concordant associations between genetic
142 susceptibility to these phenotypes and plasma abundances for 4 of the 9 most predictive IGT and iIGT
143 proteins, (p-value < 0.001, **Figure 4c**). Plasma abundances of Growth hormone receptor (GHR),
144 Reticulon-4 receptor (RTN4R), Carboxypeptidase M (CBPM) and Serine protease HTRA1 (HTRA1) were
145 associated with genetic susceptibility to more than one of these phenotypes, including fasting insulin,
146 T2D and BMI.

147 The 3 most predictive iIGT proteins and 6 of the 8 most predictive IGT proteins were significantly
148 associated with higher measured concentrations of fasting and 2-hour glucose, and insulin.

149 Chondroadherin (CHAD) was the only protein inversely associated with all 4 measures. From the
150 remaining two IGT predictor proteins only Cartilage intermediate layer protein 2 (CILP2) was
151 significantly inversely associated with fasting glucose (p -values <0.001 , **Figure 3a**). In the independent
152 prospective WHII cohort ($N = 1,492$, including 521 incident T2D cases, **Supplementary Table 11**), all
153 proteins were significantly associated with an increased risk of developing future T2D, except for
154 CHAD, which was inversely associated (p -value < 0.006 , **Figure 4b**), and CILP2, which showed no
155 significant association. Effect sizes ranged from 0.88–1.51 (hazard ratio for T2D per s.d. difference in
156 the protein target) adjusting for age, sex, and BMI. Associations for HTRA1, GHR, and CBPM remained
157 significant even upon additional adjustment for fasting glucose, total triglycerides, HDL-cholesterol,
158 and lipid lowering medication (**Supplementary Table 12**).

159 Informative biomarkers are not only relevant to improve screening strategies but can inform
160 understanding of the separate and shared aetiologies of IGT and iIGT. Comparison of protein ranking
161 from IGT as opposed to iIGT feature selection revealed that most discriminatory proteins differed
162 strongly between the IGT and iIGT selections (**Extended Data Figure 7**) with only eleven proteins
163 achieving similarly high rankings for both outcomes, that is, being selected in $>80\%$ across random
164 subsets of the study population. The top two biological GO term processes differed between the 65-
165 IGT protein signature (“proteolysis” and “cytokine-mediated signalling pathway”, **Supplementary**
166 **Table 13**) and the 68-iIGT protein signature (“cartilage development”, “collagen fibril organization”,
167 **Supplementary Table 14**), however none were significantly enriched following Bonferroni adjustment
168 for multiple comparisons.

169 To identify potential differences in factors influencing these IGT and iIGT protein signatures, we
170 computed the proportion of variance in the first principal component of the 65-IGT and 68-iIGT protein
171 signatures explained by 24 biochemical, phenotypic, and anthropometric factors. Both signatures had
172 similarly large proportions of explained variance by glycaemic (5.2 – 37.8%) and anthropometric (25.1
173 – 40.9%) measures, blood lipids (2.7 – 33.1%), or an ultrasound-based score for hepatic steatosis (22.4
174 – 24.5%) (**Methods**). Differences included the higher proportion of variance explained by C-reactive
175 protein and the lower proportion explained by ALT (a biomarker of liver injury) for the 65-IGT
176 compared to the 68-iIGT protein signature (CRP 30.2% vs 20.3% and ALT 14.7% vs 23.2%, respectively,
177 **Extended Data Figure 8**). Measures related to glucose metabolism (explaining up to 23.8% of the
178 variance) and adiposity (explaining up to 26.9 % of the variance) were identified as the main factors
179 explaining variance in the 9 predictive IGT or iIGT proteins included in the final prediction models.
180 Other protein specific factors included total triglycerides (explained up to 22.6% of GHR), HDL-
181 cholesterol (up to 13.6% of RTN4R), measures of hepatic steatosis (liver score explained up to 15% of

182 GHR) and inflammation (up to 27.2% of HTRA1), as well as genetic variants in proximity of the relevant
183 protein-encoding gene (up to 11.3% of RTN4R) (**Extended Data Figure 8**).

184

185 ***Long-term health outcomes associated with predicted iIGT***

186 To explore the clinical consequences of isolated impaired glucose tolerance in the absence of an OGTT,
187 we performed an exploratory analysis in a random sub-cohort of the prospective EPIC-Norfolk study³⁷
188 (N=753). We evaluated associations between predicted probabilities based on 1) the final clinical + 3-
189 protein model, 2) the 3-protein model only, and 3) the 68-protein iIGT model with the onset of eight
190 cardiometabolic diseases based on electronic-health record linkage³⁸ (N incident cases 30-235; follow-
191 up time between 18 – 19 years; **Supplementary Table 15 - 16**). All scores were significantly associated
192 with a greater risk of future T2D (52 incident T2D cases) at 5% false discovery rate (FDR). The iIGT final
193 clinical+3-protein score was further associated with cataracts and renal disease, possibly reflecting the
194 known association between chronically elevated 2hPG levels and micro- or macrovascular
195 complications. Predicted probabilities from the best performing 68-protein-based iIGT-model, showed
196 a nominally significant association for coronary artery disease (HR = 1.22, p-value = 0.03, CAD) and
197 peripheral artery disease (HR = 1.27, p-value = 0.04, PAD), T2D-related complications, although these
198 did not reach statistical significance when adjusting for multiple testing given the small number of
199 incident cases in this small exploratory cohort. We observed significant associations for individual
200 proteins with the risk of future T2D, with effect sizes comparable to those in the WHII study³⁹ (**Figure**
201 **5**).

202 We used proteomic measures done with a distinct proteomic technique, the Olink Explore panel⁴⁰ in
203 an independent study (random sub-cohort of the prospective EPIC-Norfolk study, N=602) to test
204 correlation of overlapping protein predictors and to validate some of our findings using an orthogonal
205 technique. We observed a high correlation between the SomaScan and Olink measurements for the
206 top three selected proteins (N=50, Spearman's r: GHR = 0.80, RTN4R = 0.70 and CBPM = 0.87,
207 Pearson's r: GHR = 0.80, RTN4R = 0.72 and CBPM = 0.82). In line with this, we replicated the previously
208 observed associations with an increased risk of incident T2D, including comparable effect sizes, and
209 further observed significant associations between the final clinical + 3-protein model and incident
210 cataracts, heart failure, and coronary heart disease (**Extended Data Figure 9**). These findings suggest
211 cross-platform transferability of our results.

212

213

214

215 **Discussion**

216 Behavioural interventions in individuals with IGT have been shown to delay progression to T2D and
217 reduce the risk of long term microvascular and macrovascular complications⁴¹. However, individuals
218 with iIGT are likely to remain undiagnosed because the current implementation of recommendations
219 for screening and diagnosing T2D does not focus on OGTTs, for reasons of practicality. People with
220 iIGT are at high risk of developing T2D and its associated complications, and failure to identify them
221 can lead to the development of severe and potentially irreversible complications of their unmanaged
222 hyperglycaemia¹⁶.

223 By combining deep plasma proteomic profiling with machine learning, we developed models for
224 improved identification of IGT and iIGT and demonstrated that as few as 8 and 3 proteins, respectively,
225 provided significant improvement over established clinical predictors²². We provided external
226 validation of the significant and substantial improvement achieved by the selected proteins over and
227 above the stringent benchmark provided by the best clinical model, something rarely done in genomic
228 or other 'omic prediction studies. The improvement observed in our independent replication study
229 was slightly greater than what was originally observed, and we note that the lack of HbA1c
230 measurements and other differences in study design (previous phases including OGTT screening) and
231 participant characteristics (older and more males on average) of the Whitehall II cohort³⁹ are likely to
232 have contributed to this, leading to a lower AUROC for the clinical model and/ or potential
233 misclassification of iIGT.

234 We propose a 3-step screening strategy, in line with the current UK Diabetes Prevention
235 Programmes⁴², involving risk assessment by 1) a patient-derived information model, 2) measuring
236 HbA1c levels and only 3 additional proteins from a single spot blood sample, and 3) an OGTT for
237 eventual diagnosis. Implementation of this proposed screening strategy, could lead to a large
238 proportion of individuals with iIGT to be additionally identified with a lower NNS, compared to the
239 currently recommended 2-stage approach⁴². Our findings illustrate how the identified proteins could
240 most efficiently be integrated into existing screening approaches to identify individuals with iIGT, who
241 are at high risk of T2D and its complications but are currently being missed. Behavioural interventions
242 have shown to be effective at reversing post-load hyperglycaemia independently of fasting glucose
243 levels^{43,44}, emphasising the value of identifying individuals with iIGT who would benefit the most from
244 these interventions. We further provided evidence of a link between our developed iIGT predictive
245 scores with incident T2D and several known cardiometabolic comorbidities resulting from chronically
246 elevated 2hPG. These finding highlight the potential of applying such a predictive risk score not only

247 for cross-sectional identification of iIGT, but for monitoring future risk for associated comorbidities
248 that impact patients' quality of life.

249 We showed that the identified proteins are not strongly affected by fasting status, suggesting that
250 they could enable a simple and convenient strategy to better identify individuals with IGT and iIGT,
251 compared to an OGTT, which requires repeated blood draws conveying additional costs¹⁸. Protein
252 assessment could substantially improve the feasibility and acceptability of an improved strategy to
253 identify iIGT, more so than alternative strategies that have been proposed such as a 1-hour OGTT⁴⁵,
254 and hence brings it in line with existing strategies for the screening and diagnosis of T2D. Since HbA1c
255 testing requires anticoagulated whole-blood, usually EDTA, a subset of the same sample type could
256 be processed for plasma preparation to measure discriminatory proteins, avoiding the need for
257 additional blood sampling.

258 This study provided insights into aetiological differences between iIGT and IGT. Our results suggested
259 a stronger low-grade inflammatory component⁴⁶⁻⁴⁹ among proteins discriminatory for IGT compared
260 to those for iIGT. These proteins might represent refined biomarkers of low-grade inflammation, as
261 they were highlighted as being predictive over and above established inflammatory markers also
262 covered in our proteomic study, such as C-reactive protein. At an individual biomarker level, we
263 identified a number of proteins shared or distinctly associated with these metabolic disturbances,
264 including GHR, RTN4R, HTRA1, CBPM, CHAD, CBLN4, NEU1, CILP2, and S100-A10. We used genetic
265 data to provide evidence that early deregulation of diabetes related pathways is linked to the
266 candidate proteins, most of which were also significantly associated with risk of future development
267 of T2D, providing a novel set of high priority T2D targets for further follow-up and assessment in in
268 more diverse settings and ethnicities.

269 While our model estimated a meaningful decrease in the NNS, there are important consideration for
270 implantation of the proposed strategy. A considerable proportion of individuals with iIGT were missed
271 by being classified low risk in either the first or subsequent screening steps. A further limitation of our
272 study was the lack of orthogonal validation of our protein-based prediction models with an alternative
273 proteomic technology. Technical, genetic, and other biological factors can result in biased protein
274 measurements due to changes in affinity of the aptamer reagents⁵⁰. However, the strong correlations
275 observed with the antibody-based Olink Explore panel suggests cross-platform transferability. We
276 further validated the phenotypic association of the iIGT predictive protein scores with incident
277 cardiometabolic diseases using Olink Explore measurements, providing the possibility of
278 implementing our model with alternative proteomic technologies.

279 In summary, we demonstrated the utility of the plasma proteome to inform strategies for screening
280 of iIGT and for gaining novel aetiological insights into early signatures of impaired glucose tolerance,
281 a globally very common and clinically important metabolic disorder, but one that it is difficult to detect
282 and treat in routine clinical practice.

283

284 **Acknowledgements**

285 The Fenland Study (10.22025/2017.10.101.00001) is funded by the Medical Research Council
286 (MC_UU_12015/1). We are grateful to all the volunteers and to the General Practitioners and practice
287 staff for assistance with recruitment. We thank the Fenland Study Investigators, Fenland Study Co-
288 ordination team and the Epidemiology Field, Data and Laboratory teams. We further acknowledge
289 support for genomics from the Medical Research Council (MC_PC_13046). Proteomic measurements
290 were supported and governed by a collaboration agreement between the University of Cambridge
291 and SomaLogic. We thank Ira von Carlowitz and Kaitlin Soucie for their contributions to the fasting
292 proteome analysis. JCZS is supported by a 4-year Wellcome Trust PhD Studentship and the Cambridge
293 Trust, CL, EW, and NJW are funded by the Medical Research Council (MC_UU_12015/1). NJW is a NIHR
294 Senior Investigator. The Whitehall II study and MK are supported by grants from the Wellcome Trust
295 (221854/Z/20/Z); UK Medical Research Council (R024227); and NIA, NIH (R01AG056477). JVL was
296 supported by Academy of Finland (311492 and 339568) and Helsinki Institute of Life Science (H970)
297 grants paid to employer. The funders had no role in study design, data collection and analysis, decision
298 to publish or preparation of the manuscript.

299 **Author Contributions**

300 JCZS, MP, NJW and CL designed the analysis and drafted the manuscript. JCZS analysed the data, JVL
301 did the replication analyses in Whitehall II study. MS and MW did the analysis for assessing the effect
302 of fasting status on protein levels. NJW is PI of the Fenland cohort and MK is PI of the Whitehall II
303 study. All authors contributed to the interpretation of the results and critically reviewed the
304 manuscript.

305

306 **Competing Interests**

307 MS, MW, DD, RO and SAW are employees of SomaLogic. EW and EO are now employees at AstraZeneca.
308 The remaining authors declare no competing interests.

309

310 **References**

- 311 1. American Diabetes, A. 2. Classification and Diagnosis of Diabetes: Standards of Medical Care
312 in Diabetes-2018. *Diabetes Care* **41**, S13-S27 (2018).
- 313 2. International Expert, C. International Expert Committee report on the role of the A1C assay in
314 the diagnosis of diabetes. *Diabetes Care* **32**, 1327-1334 (2009).
- 315 3. Saeedi, P., *et al.* Global and regional diabetes prevalence estimates for 2019 and projections
316 for 2030 and 2045: Results from the International Diabetes Federation Diabetes Atlas, 9(th)
317 edition. *Diabetes Res Clin Pract* **157**, 107843 (2019).
- 318 4. Meisinger, C., *et al.* Prevalence of undiagnosed diabetes and impaired glucose regulation in
319 35-59-year-old individuals in Southern Germany: the KORA F4 Study. *Diabet Med* **27**, 360-362
320 (2010).
- 321 5. Cheng, Y.J., *et al.* Prevalence of Diabetes by Race and Ethnicity in the United States, 2011-
322 2016. *JAMA* **322**, 2389-2398 (2019).
- 323 6. Richter, B., Hemmingsen, B., Metzendorf, M.I. & Takwoingi, Y. Development of type 2 diabetes
324 mellitus in people with intermediate hyperglycaemia. *Cochrane Database Syst Rev* **10**,
325 CD012661 (2018).
- 326 7. Yip, W.C.Y., Sequeira, I.R., Plank, L.D. & Poppitt, S.D. Prevalence of Pre-Diabetes across
327 Ethnicities: A Review of Impaired Fasting Glucose (IFG) and Impaired Glucose Tolerance (IGT)
328 for Classification of Dysglycaemia. *Nutrients* **9**(2017).
- 329 8. Campbell, M.D., *et al.* Benefit of lifestyle-based T2DM prevention is influenced by prediabetes
330 phenotype. *Nat Rev Endocrinol* **16**, 395-400 (2020).
- 331 9. Nichols, G.A., Arondekar, B. & Herman, W.H. Complications of dysglycemia and medical costs
332 associated with nondiabetic hyperglycemia. *The American journal of managed care* **14**, 791-
333 798 (2008).
- 334 10. Cowie, C.C., *et al.* Prevalence of diabetes and high risk for diabetes using A1C criteria in the
335 U.S. population in 1988-2006. *Diabetes Care* **33**, 562-568 (2010).
- 336 11. Cederberg, H., *et al.* Postchallenge glucose, A1C, and fasting glucose as predictors of type 2
337 diabetes and cardiovascular disease: a 10-year prospective cohort study. *Diabetes Care* **33**,
338 2077-2083 (2010).
- 339 12. Balkau, B. The DECODE study. Diabetes epidemiology: collaborative analysis of diagnostic
340 criteria in Europe. *Diabetes Metab* **26**, 282-286 (2000).
- 341 13. Gerstein, H.C., *et al.* Annual incidence and relative risk of diabetes in people with various
342 categories of dysglycemia: a systematic overview and meta-analysis of prospective studies.
343 *Diabetes Res Clin Pract* **78**, 305-312 (2007).
- 344 14. Chen, Y., *et al.* Associations of progression to diabetes and regression to normal glucose
345 tolerance with development of cardiovascular and microvascular disease among people with
346 impaired glucose tolerance: a secondary analysis of the 30 year Da Qing Diabetes Prevention
347 Outcome Study. *Diabetologia* **64**, 1279-1287 (2021).
- 348 15. Shaw, J.E., Hodge, A.M., de Courten, M., Chitson, P. & Zimmet, P.Z. Isolated post-challenge
349 hyperglycaemia confirmed as a risk factor for mortality. *Diabetologia* **42**, 1050-1054 (1999).
- 350 16. Silbernagel, G., *et al.* Isolated post-challenge hyperglycaemia predicts increased
351 cardiovascular mortality. *Atherosclerosis* **225**, 194-199 (2012).
- 352 17. Zhou, W., *et al.* Longitudinal multi-omics of host-microbe dynamics in prediabetes. *Nature*
353 **569**, 663-671 (2019).
- 354 18. Williams, S.A., *et al.* Plasma protein patterns as comprehensive indicators of health. *Nat Med*
355 **25**, 1851-1857 (2019).
- 356 19. Schussler-Fiorenza Rose, S.M., *et al.* A longitudinal big data approach for precision health. *Nat*
357 *Med* **25**, 792-804 (2019).
- 358 20. Gold, L., *et al.* Aptamer-based multiplexed proteomic technology for biomarker discovery.
359 *PLoS One* **5**, e15004 (2010).

- 360 21. Lindsay, T., *et al.* Descriptive epidemiology of physical activity energy expenditure in UK adults
361 (The Fenland study). *International Journal of Behavioral Nutrition and Physical Activity* **16**, 126
362 (2019).
- 363 22. Rahman, M., Simmons, R.K., Harding, A.H., Wareham, N.J. & Griffin, S.J. A simple risk score
364 identifies individuals at high risk of developing Type 2 diabetes: a prospective cohort study.
365 *Fam Pract* **25**, 191-196 (2008).
- 366 23. Deora, A.B., Kreitzer, G., Jacovina, A.T. & Hajjar, K.A. An annexin 2 phosphorylation switch
367 mediates p11-dependent translocation of annexin 2 to the cell surface. *J Biol Chem* **279**,
368 43411-43418 (2004).
- 369 24. Guevara-Aguirre, J., *et al.* Growth hormone receptor deficiency is associated with a major
370 reduction in pro-aging signaling, cancer, and diabetes in humans. *Sci Transl Med* **3**, 70ra13
371 (2011).
- 372 25. Tiaden, A.N., *et al.* Novel Function of Serine Protease HTRA1 in Inhibiting Adipogenic
373 Differentiation of Human Mesenchymal Stem Cells via MAP Kinase-Mediated MMP
374 Upregulation. *Stem Cells* **34**, 1601-1614 (2016).
- 375 26. Haddad, Y. & Couture, R. Kininase 1 As a Preclinical Therapeutic Target for Kinin B1 Receptor
376 in Insulin Resistance. *Front Pharmacol* **8**, 509 (2017).
- 377 27. Klement, J., *et al.* Oxytocin Improves beta-Cell Responsivity and Glucose Tolerance in Healthy
378 Men. *Diabetes* **66**, 264-271 (2017).
- 379 28. Zhong, C., *et al.* Cbln1 and Cbln4 Are Structurally Similar but Differ in GluD2 Binding
380 Interactions. *Cell Rep* **20**, 2328-2340 (2017).
- 381 29. Weingarten, M.F.J., *et al.* Circulating Oxytocin Is Genetically Determined and Associated With
382 Obesity and Impaired Glucose Tolerance. *J Clin Endocrinol Metab* **104**, 5621-5632 (2019).
- 383 30. Wu, T., *et al.* CILP-2 is a novel secreted protein and associated with insulin resistance. *J Mol*
384 *Cell Biol* **11**, 1083-1094 (2019).
- 385 31. Slieker, R.C., *et al.* Novel biomarkers for glycaemic deterioration in type 2 diabetes: an IMI
386 RHAPSODY study. *medRxiv*, 2021.2004.2022.21255625 (2021).
- 387 32. Shen, Z., Gantcheva, S., Mansson, B., Heinegard, D. & Sommarin, Y. Chondroadherin
388 expression changes in skeletal development. *Biochem J* **330 (Pt 1)**, 549-557 (1998).
- 389 33. Hesse, L., *et al.* The skeletal phenotype of chondroadherin deficient mice. *PLoS One* **8**, e63080
390 (2014).
- 391 34. Scott, R.A., *et al.* Large-scale association analyses identify new loci influencing glycemic traits
392 and provide insight into the underlying biological pathways. *Nat Genet* **44**, 991-1005 (2012).
- 393 35. Lotta, L.A., *et al.* Association of Genetic Variants Related to Gluteofemoral vs Abdominal Fat
394 Distribution With Type 2 Diabetes, Coronary Disease, and Cardiovascular Risk Factors. *JAMA*
395 **320**, 2553-2563 (2018).
- 396 36. Mahajan, A., *et al.* Fine-mapping type 2 diabetes loci to single-variant resolution using high-
397 density imputation and islet-specific epigenome maps. *Nat Genet* **50**, 1505-1513 (2018).
- 398 37. Day, N., *et al.* EPIC-Norfolk: study design and characteristics of the cohort. European
399 Prospective Investigation of Cancer. *Br J Cancer* **80 Suppl 1**, 95-103 (1999).
- 400 38. Pietzner, M., *et al.* Plasma metabolites to profile pathways in noncommunicable disease
401 multimorbidity. *Nat Med* **27**, 471-479 (2021).
- 402 39. Marmot, M. & Brunner, E. Cohort Profile: the Whitehall II study. *Int J Epidemiol* **34**, 251-256
403 (2005).
- 404 40. Zhong, W., *et al.* Next generation plasma proteome profiling to monitor health and disease.
405 *Nat Commun* **12**, 2493 (2021).
- 406 41. Gong, Q., *et al.* Morbidity and mortality after lifestyle intervention for people with impaired
407 glucose tolerance: 30-year results of the Da Qing Diabetes Prevention Outcome Study. *Lancet*
408 *Diabetes Endocrinol* **7**, 452-461 (2019).

- 409 42. Barron, E., Clark, R., Hewings, R., Smith, J. & Valabhji, J. Progress of the Healthier You: NHS
410 Diabetes Prevention Programme: referrals, uptake and participant characteristics. *Diabet Med*
411 **35**, 513-518 (2018).
- 412 43. Gong, Q., *et al.* Efficacy of lifestyle intervention in adults with impaired glucose tolerance with
413 and without impaired fasting plasma glucose: A post hoc analysis of Da Qing Diabetes
414 Prevention Outcome Study. *Diabetes Obes Metab* **23**, 2385-2394 (2021).
- 415 44. Knowler, W.C., *et al.* Reduction in the incidence of type 2 diabetes with lifestyle intervention
416 or metformin. *N Engl J Med* **346**, 393-403 (2002).
- 417 45. Bergman, M., *et al.* Lessons learned from the 1-hour post-load glucose level during OGTT:
418 Current screening recommendations for dysglycaemia should be revised. *Diabetes Metab Res*
419 *Rev* **34**, e2992 (2018).
- 420 46. Pham, C.T. Neutrophil serine proteases: specific regulators of inflammation. *Nat Rev Immunol*
421 **6**, 541-550 (2006).
- 422 47. Wiedow, O. & Meyer-Hoffert, U. Neutrophil serine proteases: potential key regulators of cell
423 signalling during inflammation. *J Intern Med* **257**, 319-328 (2005).
- 424 48. Donath, M.Y. & Shoelson, S.E. Type 2 diabetes as an inflammatory disease. *Nat Rev Immunol*
425 **11**, 98-107 (2011).
- 426 49. de Vries, M.A., *et al.* Glucose-dependent leukocyte activation in patients with type 2 diabetes
427 mellitus, familial combined hyperlipidemia and healthy controls. *Metabolism* **64**, 213-217
428 (2015).
- 429 50. Pietzner, M., *et al.* Synergistic insights into human health from aptamer- and antibody-based
430 proteomic profiling. *Nat Commun* **12**, 6822 (2021).

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447 **Figure Legends**

448 **Figure 1. Study design.** **a**, Proteomic profiling was done in fasting plasma samples from participants from the
449 Fenland cohort that had undergone an OGTT. **b**, 3-step modelling framework for IGT and iIGT classification. *For
450 iIGT prediction individuals with non-isolated IGT were excluded. **c**, Association of top discriminatory proteins
451 with incident type 2 diabetes was assessed in the Whitehall II study. **d**, Association of iIGT protein scores with 8
452 incident cardiometabolic diseases was assessed in a sub-cohort of the EPIC-Norfolk study. OGTT: oral glucose
453 tolerance test, IGT: impaired glucose tolerance, iIGT: isolated impaired glucose tolerance.

454

455 **Figure 2. Performance of LASSO trained models for impaired glucose tolerance (a) and isolated impaired**
456 **glucose tolerance (b) discrimination in the internal validation test set.** **a**, IGT discrimination performance in the
457 independent internal validation test set (N=2881, 192 IGT individuals) for the standard clinical model (Cambridge
458 T2D risk Score + FPG + HbA1c), a 65-protein model and a clinical + 8 protein model. **b**, iIGT discrimination
459 performance in the independent internal validation test set (N=2795, 111 iIGT individuals) for the standard
460 clinical model, a 68-protein model and a clinical + 3 protein model. **c**, Comparison of protein ranking during
461 feature selection for iIGT (N=2795, 111 iIGT individuals) and IGT (N=2881, 192 IGT individuals) top discriminatory
462 proteins. IGT: impaired glucose tolerance, iIGT: isolated impaired glucose tolerance, FPG: fasting plasma glucose,
463 HbA1c: glycated haemoglobin.

464

465 **Figure 3. Proposed 3-stage screening strategy.** In the first stage, individuals in the whole of Fenland were divided
466 into low and high risk according to the Cambridge T2D risk score. The high risk group would undergo a second
467 stage involving measurement of HbA1c and of the 3 iIGT proteins. Individuals with HbA1c levels within the T2D
468 or prediabetic range would be referred for intervention and lifestyle modifications. Individuals with HbA1c below
469 the prediabetic range, would be further stratified using the final clinical + 3 iIGT protein model to identify a high
470 risk group, which on a third stage would be taken forward for OGTT testing to identify iIGT cases that would
471 have otherwise been missed by current screening guidelines. Figure was designed with biorender.com.

472

473 **Figure 4. Characterization of the association between top impaired glucose tolerance and isolated impaired**
474 **glucose tolerance discriminatory proteins and glycaemic traits, future T2D risk and genetic predisposition to**
475 **metabolic phenotypes.** **a**, Association of top IGT and iIGT discriminatory proteins with fasting and 2-hour glucose
476 and insulin in the Fenland study (N = 10259 individuals). Beta estimates with 95% confidence intervals are shown.
477 **b**, Association of top IGT and iIGT discriminatory proteins with incident T2D in the Whitehall II study (N = 1492,
478 521 incident T2D cases). Hazard ratios (HR) with 95% confidence intervals are shown. **c**, Association of genetic
479 risk scores for fasting glucose, fasting insulin, 2-hour plasma glucose, type 2 diabetes and body mass index with
480 top IGT and iIGT discriminatory proteins in the Fenland study (N = 7973 individuals). Beta estimates with a 95%
481 confidence interval are shown. FG: fasting glucose, FI: fasting insulin, 2hPG: 2-hour plasma glucose, 2hPI: 2-hour
482 plasma insulin, T2D: type 2 diabetes, BMI: body mass index.

483

484 **Figure 5. Association of iIGT protein scores with incident cardiometabolic diseases.** Association of iIGT
485 prediction scores (left panel) or individual top iIGT proteins (right panel) with 8 cardiometabolic disease
486 outcomes in a sub-cohort the EPIC-Norfolk study (N=753 individuals). Hazard ratios (HR) with 95% confidence
487 intervals are shown.

488

489

490

491

492

493 **Methods**

494 ***Study Samples***

495 The Fenland study²¹ is a population-based cohort of 12,435 men and women born between 1950 and
496 1975 who underwent detailed phenotyping at the baseline visit from 2005-2015. Participants were
497 recruited from general practice surgeries in Cambridge, Ely and Wisbech (UK). Exclusion criteria of the
498 Fenland study included pregnancy, prevalent diabetes, an inability to walk unaided, psychosis, or
499 terminal illness. The study was approved by the Cambridge Local Research Ethics Committee (NRES
500 Committee – East of England Cambridge Central, ref. 04/Q0108/19) and all participants provided
501 written informed consent. The consent covered measurements made from blood samples as well as
502 extends beyond the baseline examination as described previously²¹.

503 ***Clinical assessment***

504 All participants completed a 2-hour 75 g OGTT following an overnight fast. Blood samples were
505 collected at fasting and 2-hour post glucose load in EDTA tubes for plasma separation by
506 centrifugation. Samples were kept at -80°C until further analysis. Glucose (assayed in a Dade Behring
507 Dimension RxL analyser) and insulin (DELFIA® immunoassay, Perkin Elmer) concentrations were
508 measured at fasting and 2-hours, as well as lipid profiles (triglycerides, HDL and total cholesterol),
509 alanine aminotransferase (ALT), alkaline phosphatase (ALP), C-reactive protein (CRP) and serum
510 creatinine (assayed in a Dade Behring Dimension RxL analyser) at fasting, and HbA1c (Tosoh
511 Bioscience, TOSOH G7 analyser).

512 IGT and T2D were defined by 2-hour glucose according to IEC diagnosis criteria² as glucose levels
513 between 7.8 and < 11.1 mmol/L (141 and < 199 mg/dL) and ≥ 11.1 mmol/L (≥ 199 mg/dL), respectively.
514 IGT was defined as 2hPG ≥ 7.8 mmol/L and <11.1 mmol/L, post-challenge hyperglycaemia as 2hPG
515 ≥ 11.1 mmol/L, iIGT as individuals with IGT but HbA1c <42mmol/mol (6%) and FG <6.1 mmol/L
516 (<110mg/dL), and isolated post-challenge hyperglycaemia as individuals with post-challenge
517 hyperglycaemia but HbA1c <42mmol/mol and FG <6.1 mmol/L. The number of individuals with post-
518 challenge hyperglycaemia in the diabetic range (i.e., 2hPG ≥ 11.1 mmol/L) was too low to investigate
519 the performance of our models to identify this group of people with undiagnosed T2D biochemically
520 defined solely due to elevated 2-hour glucose. These individuals would still be missed and remain
521 undiagnosed by FG and HbA1c testing. We therefore used the terms IGT and iIGT to refer to all
522 individuals with 2hPG ≥ 7.8 mmol/L throughout text and in order to develop a model that captures all
523 individuals that would remain undiagnosed by current strategies. We note that the thresholds to
524 define glycaemic categories vary across the American Diabetes Association (ADA) , WHO and the
525 International Expert Committee (IEC)⁵¹. We use the IEC HbA1c and FG thresholds to reflect current

526 clinical practice in the UK. We note that using ADA thresholds will likely results in lower case numbers
527 for IGT and iIGT at the cost of a substantially higher false-positive rate. Body mass index (BMI) was
528 calculated as weight (kg) / square of height (m²). Additionally, the homeostasis model assessment of
529 insulin resistance (HOMA-IR) was calculated as FI (μIU/mL) × fasting glucose (mmol/mL)/22.5⁵².
530 Estimated glomerular filtration rate (eGFR) was calculated by the CKD-EPI equation using serum
531 creatinine⁵³.

532 Hepatic steatosis was evaluated by an abdominal ultrasound and images were scored by two trained
533 operators. Criteria used for scoring included: increased echotexture of the liver parenchyma,
534 decreased visualisation of the intra-hepatic vasculature and attenuation of ultrasound beam. A normal
535 liver was considered as a score from 3 – 4, mild steatosis from 5 – 7, moderate steatosis from 8 – 10
536 and sever steatosis ≥ 11 ⁵⁴.

537 Participants completed DEXA scan measurements using a Lunar Prodigy advanced fan beam scanner
538 (GE Healthcare) performed by trained operators using standard imaging, positioning protocols and
539 manually processed according to a standardized procedure described previously³⁵. Abdominal visceral
540 and subcutaneous fat mass was estimated using the DEXA software.

541 Differences in clinical characteristics were evaluated by ANOVA followed by posthoc Tukey test, or χ²
542 for categorical variables. Non-normally distributed variables were log transformed when appropriate.

543 ***Proteomic profiling of the Fenland cohort***

544 Proteomic profiling was done using an aptamer-based technology (SomaScan proteomic assay).
545 Fasting proteomic profiling was done in participants from the Fenland cohort at baseline, from which
546 relative abundancies of 4,775 unique protein targets (evaluated by 4,979 SOMAmer reagents,
547 SomaLogic v4)^{18,55} was evaluated in EDTA plasma. Briefly, proteins are targeted by modified single
548 stranded DNA sequences (aptamers). Concentration is then approximated as relative fluorescence
549 units using a DNA microarray ⁵⁶.

550 To account for variation in hybridization within runs, hybridization control probes are used to generate
551 a hybridization scale factor for each sample. To control for total signal differences between samples
552 due to variation in overall protein concentration or technical factors such as reagent concentration,
553 pipetting or assay timing, we used the adaptive median normalisation (AMN), unless stated otherwise.
554 Briefly, a ratio between each aptamer's measured value and a reference value from an external
555 reference population is computed, and the median of these ratios is computed for each of the three
556 dilution sets (20%, 1% and 0.005%) and applied to each dilution set to shift the intrapersonal
557 distribution of protein intensities accordingly to match the reference population. We removed
558 samples if they did not meet an acceptance criterion for scaling factors with values outside of the

559 recommend range (0.25-4) or were flagged as technical failures (n=19). Detailed SomaLogic's
560 normalization, calibration data, and quality control processes have been previously described in
561 detail⁴⁸. At a protein level, we took only human protein targets forward for subsequent analysis (4,979
562 out of the 5284 aptamers). Intraassay coefficients of variation (calculated based on raw fluorescence
563 units) had a median of 4.98% (interquartile range 3.87% - 6.99%) suggesting good quality measures
564 for the vast majority of protein targets. We decided to not apply any other filters to individual protein
565 qualities given that even poorly measured proteins might be informative and left it to the restrictive
566 feature selection approach applied to drop uninformative proteins, including possibly poorly
567 measured once. Aptamers' target annotation and mapping to UniProt accession numbers as well as
568 Entrez gene identifiers were provided by SomaLogic and we used those to obtain genomic positions
569 of protein encoding genes.

570 ***Genome wide genotyping and imputation***

571 Fenland participants were genotyped using three genotyping arrays: the Affymetrix UK Biobank Axiom
572 array (OMICs, N=8994), Illumina Infinium Core Exome 24v1 (Core-Exome, N=1060) and Affymetrix
573 SNP5.0 (GWAS, N=1402). Samples were excluded for the following reasons: 1) failed channel contrast
574 (DishQC <0.82); 2) low call rate (<95%); 3) gender mismatch between reported and genetic sex; 4)
575 heterozygosity outlier; 5) unusually high number of singleton genotypes or 6) impossible identity-by-
576 descent values. Single nucleotide polymorphisms (SNPs) were removed if: 1) call rate < 95%; 2) clusters
577 failed Affymetrix SNPolisher standard tests and thresholds; 3) MAF was significantly affected by plate;
578 4) SNP was a duplicate based on chromosome, position, and alleles (selecting the best probe set
579 according to Affymetrix SNPolisher); 5) Hardy-Weinberg equilibrium $p < 10^{-6}$; 6) did not match the
580 reference or 7) MAF=0.

581 Autosomes for the OMICS and GWAS subsets were imputed to the HRC (r1) panel using IMPUTE4, and
582 the Core-Exome subset and the X-chromosome (for all subsets) were imputed to HRC.r1.1 using the
583 Sanger imputation server⁵⁷. All three arrays subsets were also imputed to the UK10K+1000Gphase3⁵⁸
584 panel using the Sanger imputation server in order to obtain additional variants that do not exist in the
585 HRC reference panel. Variants with MAF < 0.001, imputation quality (info) < 0.4 or Hardy Weinberg
586 Equilibrium $p < 10^{-7}$ in any of the genotyping subsets were excluded from further analyses.

587 **Statistical Analyses**

588 ***Classification of IGT and iIGT from the fasting proteome***

589 To identify and validate a proteomic signature able to discriminate IGT and iIGT (as a binary outcome),
590 the entire Fenland study (N=11,546 without missing data for 2hPG), was divided into three subsets:
591 for feature selection (50%, N = 5773), parameter optimization (25%, N=2887) and validation (25%,

592 N=2881). IGT and iIGT cases were split equally into 50% for training ($N_{IGT} = 387$, $N_{iIGT} = 222$), 25 % for
593 optimization ($N_{IGT} = 194$, $N_{iIGT} = 111$) and 25% for testing ($N_{IGT} = 193$, $N_{iIGT} = 111$) sets. For these
594 analyses, SOMAmer RFUs were \log_{10} -transformed. Feature selection was carried out by least absolute
595 shrinkage and selection operator (LASSO) regression. We chose to use LASSO because it was the most
596 suitable model to 1) identify the smallest possible set of independent predictors, 2) it is
597 computationally efficient, which allowed us to implement a robust framework using bootstrap
598 resampling to identify a core set of most informative predictors and 3) it is less prone to overfitting.
599 To address case-control imbalance we used the ROSE R package⁵⁹, which implements down-sampling
600 of the majority class (controls) along with synthetic new data points for the minority class (IGT or iIGT).
601 A nested 10-fold cross-validation (inner loop to determine regularization parameter, λ) was done over
602 100 bootstrap samples (outer loop) drawn from the feature selection set. Each protein received a
603 score that was generated by counting the number of times it was included in the final model from
604 each of the 100 bootstrap samples, that is, the score was between 0 (for proteins that were never
605 selected in the final model) and 100 (for proteins that were selected in the final model in all bootstrap
606 samples). We ranked the proteins based on their score to identify the most informative set of features
607 (i.e. with a higher score) (**Supplementary Fig. 1**). This was implemented by the use of the R packages
608 *caret*⁶⁰ and *glmnet*⁶¹. Proteins selected in the final model in more than 80%, 90%, and 95% of the
609 bootstrap samples, were tested as predictors and taken forward for parameter optimization by 10-
610 fold cross validation of the model by LASSO regression in the optimization set. Additional models were
611 optimized by LASSO regression, such as a standard patient information-based model using the
612 variables from the Cambridge Diabetes Risk Score (age, sex, family history of diabetes, smoking status,
613 prescription of steroid or antihypertensive medication and BMI)²², a standard clinical model (including
614 the variables from the Cambridge Diabetes risk Score, FG and HbA1c) and a standard clinical plus the
615 selected proteins model. Clinical predictors were forced to be kept in the clinical plus proteins model
616 by setting the penalty factors of these variables to 0. For comparison, ridge regression (which will keep
617 all proteins in the final model) was used to build a prediction model using all the 4979 proteins as
618 predictors.

619 Performance of the classification models were evaluated in the internal independent validation set,
620 which was never used for training and optimization. The prediction models' discriminatory power was
621 assessed by computing the area under the receiver operating curve (AUROC). Confidence intervals
622 and p-values (using the deLong method implemented by the R package pROC⁶²) were computed for
623 the comparison between the ROC curves for the standard clinical model and clinical with added
624 proteins model. Additionally, models' net reclassification index was evaluated using the R package
625 PredictABEL⁶³.

626 Using an analogous machine learning strategy, we developed models for iIGT discrimination. For these
627 analyses, all individuals with non-isolated IGT (2hPG > 7.8 mmol/L, FPG > 6.1 mmol/L and HbA1c > 42
628 mmol/mol) were excluded from the cohort (leaving N = 11,281), which was subsequently divided into
629 feature selection (50%, N = 5591), parameter optimization (25%, N=2796) and validation (25%,
630 N=2795). Feature selection, optimization and testing were carried out as described for IGT models. To
631 achieve comparable model performance with the minimal number of predictors, we used recursive
632 feature elimination on the set of proteins selected in >95% of boots during feature selection. As a
633 sensitivity analysis, we performed the same framework described above, that is, feature selection,
634 parameter optimization and validation to assess model performance when using protein data
635 reversing the final normalisation step that is unique to the SomaScan platform. We note that using
636 'non-normalised' proteomic data led to broadly comparable results, which are well in the margins of
637 random variation of protein measurements in general, albeit with some difference in the proteins
638 selected as the most predictive markers in the final models (**Supplementary Table 17**).

639 Calibration of the final models was assessed in the internal validation set by computing the calibration
640 slope, which evaluates the spread of the estimated risks and has a target value of one. Calibration
641 slopes less than 1 indicate extreme estimated risks while slopes greater than 1 indicate very moderate
642 risk estimates. Calibration slopes were computed using the R package rms⁶⁴.

643 The number needed to screen (NNS) was calculated using a staged screening scenario. Firstly,
644 participants from the Fenland study were stratified by predicted probabilities from the Cambridge T2D
645 risk Score, that is, non-invasive risk factors that could be obtained by interviewing the patient. The
646 threshold used to stratify individuals into "high" and "low" risk strata according to their predicted
647 probabilities was set to optimize a balance between the total number of individuals that would be
648 needed to screen and sensitivity (as would be appropriate for such a screening setting), which was
649 achieved at 0.7, regardless of specificity. On second instance, participants within the high-risk group
650 were further stratified by HbA1c levels, using IEC cut-offs (normoglycaemic : HbA1c < 42 mmol/mol,
651 prediabetic criteria: HbA1c >= 42 mmol/mol and < 48 mmol/mol, T2D criteria : HbA1c >= 48
652 mmol/mol)⁵¹. On third instance, participants whose HbA1c did not meet the criteria for T2D or
653 prediabetes (that is, normoglycaemic as defined above), were further stratified according to the
654 clinical + 3- iIGT protein model. Similarly, a threshold that optimized testing as few individuals as
655 possible while retaining good sensitivity of 0.7 was set for this model (**Supplementary Table 10**). We
656 estimated the NNS within this stratum compared to the NNS within the full set of individuals with
657 HbA1c in the normoglycaemic range. The NNS was calculated as the total number of individuals within
658 the group divided by the number iIGT cases within the same group and refers to the number of OGTTs

659 that would need to be done to identify one iIGT case within the group of interest. We additionally
660 estimated the NNS in the test set only, as a sensitivity analysis.

661 ***IGT/iIGT model validation and follow-up analyses in the WHII study***

662 The Whitehall II study is a longitudinal, prospective cohort study³⁹ that was approved by the joint
663 University College London / University College London Hospital's Committees on the Ethics of Human
664 Research. Proteomic profiling of fasting EDTA-plasma samples was done for all individuals at phase 5
665 (from 1997 - 1999) with the SomaScan v4.1 proteomic assay. We performed validation of the IGT and
666 iIGT clinical + protein models at phase 5 (from 1997 - 1999) of the study, where proteomic profiling and
667 OGTT values were available. Since HbA1c was not measured at phase 5 of the study, we defined iIGT
668 as 2hPG > 7.8 mmol/L and FPG > 6.1 mmol/L. We used the weights from the models trained in Fenland
669 to evaluate their performance in WHII phase 5 (total sample size = 5058, N_{IGT}= 693, N_{iIGT}=617) for the
670 baseline clinical model (Cambridge T2D risk score + FG) and the baseline clinical + protein iIGT and IGT
671 models (3 and 8 proteins respectively).

672 For the association between top discriminatory proteins and incident T2D in the Whitehall II study
673 individuals were selected as a nested case-control study design in which proteomic profiling of fasting
674 EDTA-plasma samples was done at phase 5 (from 1997 - 1999) with the SomaScan v4 proteomic assay.
675 Incident T2D occurrence was assessed in repeated clinical examinations in 1997-1999, 2002-2004,
676 2007-2009, 2012-2013, and 2015-2016, based on FPG above 7 mmol/L, HbA1c>6.5%, use of diabetes
677 medication, or reported physician diagnosed diabetes, excluding prevalent T2D cases at baseline from
678 the analysis. Additionally, participants with impaired kidney function (eGFR < 30 mL/min/1.73m²),
679 incident cardiovascular diseases or missing data on T2D at follow-up were excluded. The final sample
680 comprised of 521 cases and 971 controls.

681 Association between fasting candidate proteins and incident T2D was assessed using Cox-proportional
682 hazards regression adjusting for the baseline confounders age, sex and BMI. We tested a second
683 model adjusting for additional baseline confounders including FG, triglycerides, HDL-cholesterol and
684 lipid lowering medication on top of age, sex and BMI to determine whether the association persisted
685 in a more refined model.

686 **Effect of fasting status on plasma levels of IGT and iIGT discriminatory proteins**

687 Fourteen adult participants were recruited to participate in the study and provided informed consent
688 appropriately. Participants were asked to fast overnight for at least 12 hours prior to reporting to the
689 study site. Fasting blood samples were collected from each participant, after which they were given a
690 moderate fat meal consisting of 5-8 ounces of Cheerios with 6 ounces of 2% milk, one egg, one slice

691 of bacon, one slice of toast with margarine, and 4 ounces of orange juice (calories: 450, 16.9 grams of
692 fat, 16 grams of protein, and 59 grams of carbohydrates)⁶⁵.

693 The time for each participant to complete the meal ranged from 7 to 19 minutes (average of 16
694 minutes). Post prandial blood samples were collected at 0.5, 1, and 3 hours following completion of
695 the meal. Since each participant consumed their meals at different rates, the actual blood collection
696 times post meal does vary between participants. Participants were not allowed to eat or drink any
697 further caloric items until after the last blood collection. Twelve participants (6 male and 6 female)
698 completed the study. Two participants were excluded due to unmet fasting requirements and an
699 adverse reaction during the first blood draw.

700 Blood samples were processed to obtain EDTA-plasma by centrifugation and frozen at -80°C until
701 delivered to SomaLogic Sample Management for proteomic profiling using the SomaScan v4 assay.
702 The effect of fasting status on 9 unique SOMAmer reagents included in the final clinical + protein
703 models for IGT or iIGT, was tested by repeated measures ANOVA. Proteins with ANOVA p-values <
704 0.0055 (according to Bonferroni adjustment for 9 comparisons) were deemed to be significantly
705 affected by fasting status.

706 ***Functional annotation of IGT and iIGT-protein signatures***

707 Functional annotation of the 65-IGT and 68-iIGT protein signatures was performed using modified
708 Fisher's exact tests as implemented by the Database for Annotation, Visualization and Integrated
709 Discovery (DAVID, version 6.8) and enrichment of biological process GO terms (GOTERM_BP_DIRECT)
710 was analysed, setting the full list of proteins evaluated by the SomaLogic platform as the background.

711 ***Variance explained in top discriminatory protein levels by clinical, biochemical, anthropometric and*** 712 ***behavioural risk factors***

713 The proportion of variance explained in candidate protein levels by several variables was evaluated in
714 the Fenland cohort using the *variancePartition* R package⁶⁶. Analogously, the proportion of variance
715 explained in the first principal component of the 65-IGT and 68-iIGT discriminatory protein signatures
716 was evaluated. Briefly, this package fits a linear mixed model to assess the effect of each variable on
717 the outcome while correcting for all other variables. Variables evaluated were age, sex, IGT, IPCH, FPG,
718 2hPG, FI, 2hPI, HbA1c, total triglycerides, total cholesterol, HDL-cholesterol, LDL-cholesterol, ALT, ALP,
719 a liver score, BMI, waist-to-hip ratio (WHR), amount of subcutaneous fat, amount of visceral fat, CRP,
720 estimated glomerular filtration rate (eGFR) and intake of statins or antihypertensive medication. FPG,
721 2hPG, FI, 2hPI, HbA1c, total triglycerides, ALT, ALP, CRP, subcutaneous fat and visceral fat were natural
722 log-transformed due to skewed distribution of these variables. We fit separate models for each of the
723 variables evaluated adjusting only for age and sex in the entire Fenland cohort (N=11,546) to avoid

724 bias due to strong collinearity among variables tested. For each of the models, participants with
725 missing data were excluded.

726 ***Protein quantitative trait loci (pQTLs) for candidate proteins***

727 Genetic variants associated with candidate proteins (protein quantitative trait loci or pQTLs) were
728 taken from our genome-wide association studies across all aptamers as described in Pietzner et al,
729 2021⁵⁵.

730 ***Percentage of variance explained in protein levels by cis and trans pQTL scores***

731 Polygenic scores were constructed for pQTLs within the *cis* (within ± 500 kb of the protein-encoding
732 gene) and *trans* regions. Cis-pQTL scores were built using conditionally independent variants. The
733 percentage of variance explained in protein levels by the cis and trans-scores was computed as
734 described in the above section adjusting for age and sex.

735 ***Association between top discriminatory proteins and fasting and 2-hour plasma glucose and insulin***

736 Observational associations between the top selected IGT and iIGT discriminatory proteins and FPG, FI,
737 2hPG and 2hPI were assessed in the entire Fenland cohort at baseline (N=10,259 without missing data)
738 by linear regression models adjusting for age, sex, BMI and test site from the study. The models for
739 2hPG and 2hPI were additionally adjusted by FPG and FPG + FI, respectively. Protein levels were log₁₀-
740 transformed and standardized, and 2hPG and 2hPI values were log-transformed for these analyses.
741 Proteins were considered significant at a Bonferroni threshold (p-values < 0.001, accounting for
742 comparisons between the number of protein and number of traits, as for all further association
743 analyses).

744 ***Association between polygenic risk scores for glycaemic traits and top discriminatory proteins***

745 T2D³⁶, fasting glucose (FG)³⁴, fasting insulin³⁴ (FI score), 2hPG³⁴ (2hPG score) and BMI³⁵ polygenic
746 scores, weighted by genetic effect sizes of previously reported genome-wide significant variants, were
747 computed for 7,973 Fenland participants genotyped with the same array (Affymetrix UK Biobank
748 Axiom Array). Variants not available, with low imputation quality scores < 0.6, or with strand
749 ambiguous alleles were excluded from the scores. Each polygenic score was tested for associations
750 with the plasma abundancies of top IGT and iIGT discriminatory proteins by linear regression models
751 adjusting for age, sex, BMI, the first 10 genetic principal components and test site of the study.

752 ***Association between iIGT scores with incident cardiometabolic diseases in a sub-cohort of the EPIC-*** 753 ***Norfolk study***

754 The EPIC-Norfolk study is a cohort of 25,639 middle-aged, individuals from the general population of
755 Norfolk a county in Eastern England which is a component of EPIC³⁷. The EPIC-Norfolk study was

756 approved by the Norfolk Research Ethics Committee (ref. 05/Q0101/191); all participants gave their
757 informed written consent before entering the study. All participants were flagged for mortality at the
758 UK Office of National Statistics and vital status was ascertained for the entire cohort. Death certificates
759 were coded by trained nosologists according to the International Statistical Classification of Diseases
760 and Related Health Problems, 10th Revision (ICD-10). Hospitalization data were obtained using
761 National Health Service numbers through linkage with NHS Digital. Participants were identified as
762 having experienced an event if the corresponding ICD-10 code was registered on the death certificate
763 (as the underlying cause of death or as a contributing factor) or as the cause of hospitalization
764 (**Supplementary Table 15**). Since the long-term follow-up of EPIC-Norfolk comprised the ICD-9 and
765 ICD-10 coding system, codes were consolidated. The current study is based on follow-up to 31 March
766 2016. Information on lifestyle factors and medical history was obtained from questionnaires as
767 reported previously³⁷. The current analysis is based on a random sub-cohort (N=875) of the whole
768 EPIC-Norfolk study population that was selected excluding known prevalent case subjects of diabetes
769 at baseline was using the same definitions as used in the InterAct Project⁶⁷; in which proteomic
770 profiling was done at health check 1 using the SOMAscan v4 platform from citrate-plasma samples
771 stored in liquid nitrogen since the baseline visit.

772 Participants with missing data for any of the variables included in the final prediction models
773 developed in the Fenland study were excluded. The final sample comprised of 753 individuals for
774 which characteristics are presented in **Supplementary Table 16**.

775 Final prediction models trained and optimized for iIGT in the Fenland study were used to calculate the
776 predicted probability of iIGT for each participant at health check 1 in this sub-cohort of the EPIC-
777 Norfolk study. Models tested included: the clinical + 3-proteins iIGT model, 3-protein iIGT model (95%
778 feature selection protein set model), 68-protein iIGT model (80% feature selection protein set model)
779 and the clinical model as a baseline comparison. We then tested the association of the predicted iIGT
780 probability with 8 incident cardiometabolic diseases (or associated T2D comorbidities) including type
781 2 diabetes, coronary heart disease, heart failure, peripheral artery disease, cerebral stroke, liver
782 disease, renal disease and cataracts using cox proportional hazards models adjusting by age at
783 baseline and sex (except for the clinical + 3 protein model, which already accounted for these risk
784 factors within the score). Associations were deemed significant at an 5% FDR accounting for
785 comparison between 8 diseases.

786 We aimed for cross-platform validation in a separate random sub-cohort of the prospective EPIC-
787 Norfolk study (N=771), in which proteomic measures were done with the Olink Explore panel⁴⁰ from
788 serum samples. Participants with missing data for any of the variables included in the final prediction

789 models developed in the Fenland study (except HbA1c which was excluded from the models as it was
790 unavailable in a large proportion of participants from this sub-cohort) were excluded. The final sample
791 comprised of 602 individuals for which characteristics are presented in **Supplementary Table 18**.

792 Final prediction models trained and optimized for iIGT in the Fenland study (using SomaScan) were
793 used to calculate the predicted probability of iIGT for each participant at health check 1 in this sub-
794 cohort of the EPIC-Norfolk study, using the Olink measures for the proteins. Models tested included:
795 the clinical + 3-proteins iIGT model, 3-protein iIGT model (95% feature selection protein set model)
796 and the Cambridge T2D risk Score. We then tested the association of the predicted iIGT probability
797 with the same Cox-model setting and set of disease as in the sub-cohort with available SomaLogic
798 measurements except for liver disease (**Supplementary Table 19**). Associations were deemed
799 significant at an 5% FDR accounting for comparison between 7 diseases.

800 All statistical analyses were performed using R language, and environment for statistical computing
801 (version 3.6.1 and 4.1.0, R Core Team).

802

803 **Data availability**

804 Data access for the Fenland and EPIC studies can be requested by bona fide researchers for specified
805 scientific purposes through a simple application process via the study websites below. Data will either
806 be shared through an institutional data sharing agreement or arrangements will be made for analyses
807 to be conducted remotely without the necessity for data transfer.

808 Fenland: <https://www.mrc-epid.cam.ac.uk/research/studies/fenland/information-for-researchers>

809 EPIC-Norfolk: <https://www.mrc-epid.cam.ac.uk/research/studies/epic-norfolk>

810

811 **Code availability**

812 The code employed for the machine learning developed framework has been deposited in the
813 following repository: https://github.com/MRC-Epid/iigt_prediction_proteomics.

814

815

816

817

818 **Methods-only references**

819

- 820 51. Lee, C.M.Y., *et al.* Comparing different definitions of prediabetes with subsequent risk of
821 diabetes: an individual participant data meta-analysis involving 76 513 individuals and 8208
822 cases of incident diabetes. *BMJ Open Diabetes Res Care* **7**, e000794 (2019).
- 823 52. Fukagawa, N.K., *et al.* Insulin-mediated reduction of whole body protein breakdown. Dose-
824 response effects on leucine metabolism in postabsorptive men. *J Clin Invest* **76**, 2306-2311
825 (1985).
- 826 53. Inker, L.A., *et al.* Estimating glomerular filtration rate from serum creatinine and cystatin C.
827 *The New England journal of medicine* **367**, 20-29 (2012).
- 828 54. Mehta, S.R., Thomas, E.L., Bell, J.D., Johnston, D.G. & Taylor-Robinson, S.D. Non-invasive
829 means of measuring hepatic fat content. *World J Gastroenterol* **14**, 3476-3483 (2008).
- 830 55. Pietzner, M., *et al.* Mapping the proteo-genomic convergence of human diseases. *Science*,
831 eabj1541 (2021).
- 832 56. Rohloff, J.C., *et al.* Nucleic Acid Ligands With Protein-like Side Chains: Modified Aptamers and
833 Their Use as Diagnostic and Therapeutic Agents. *Mol Ther Nucleic Acids* **3**, e201 (2014).
- 834 57. McCarthy, S., *et al.* A reference panel of 64,976 haplotypes for genotype imputation. *Nat*
835 *Genet* **48**, 1279-1283 (2016).
- 836 58. Huang, J., *et al.* Improved imputation of low-frequency and rare variants using the UK10K
837 haplotype reference panel. *Nat Commun* **6**, 8111 (2015).
- 838 59. Nicola Lunardon, G.M., Nicola Torelli. ROSE: a Package for Binary Imbalanced Learning. *The R*
839 *Journal* **6**, 79-89 (2014).
- 840 60. Kuhn, M. Building Predictive Models in R Using the caret Package. *J Stat Softw* **28**, 1-26 (2008).
- 841 61. Friedman, J., Hastie, T. & Tibshirani, R. Regularization Paths for Generalized Linear Models via
842 Coordinate Descent. *J Stat Softw* **33**, 1-22 (2010).
- 843 62. Robin, X., *et al.* pROC: an open-source package for R and S+ to analyze and compare ROC
844 curves. *BMC Bioinformatics* **12**, 77 (2011).
- 845 63. Kundu, S., Aulchenko, Y.S., van Duijn, C.M. & Janssens, A.C. PredictABEL: an R package for the
846 assessment of risk prediction models. *European journal of epidemiology* **26**, 261-264 (2011).
- 847 64. Jr, F.E.H. rms: Regression Modeling Strategies. R package version 5.1-1. (2017).
- 848 65. *Pharmacokinetics in Drug Development: Clinical Study Design and Analysis* (American
849 Association of Pharmaceutical Scientists, Arlington, 2004).
- 850 66. Hoffman, G.E. & Schadt, E.E. variancePartition: interpreting drivers of variation in complex
851 gene expression studies. *BMC Bioinformatics* **17**, 483 (2016).
- 852 67. InterAct, C., *et al.* Design and cohort description of the InterAct Project: an examination of the
853 interaction of genetic and lifestyle factors on the incidence of type 2 diabetes in the EPIC
854 Study. *Diabetologia* **54**, 2272-2282 (2011).

855









