Similar and different: systematic investigation of proteogenomic variation between sexes and its relevance for human diseases

- 3 Mine Koprulu¹, Eleanor Wheeler¹, Nicola D. Kerrison¹, Spiros Denaxas^{2,3,4,5}, Julia Carrasco-
- Zanini ^{1,7}, Chloe M. Orkin ^{8,9}, Harry Hemingway ^{2,3,5}, Nicolas J. Wareham ¹, Maik Pietzner
 ^{1,6,7}, Claudia Langenberg ^{1,6,7}
- 1, MRC Epidemiology Unit, University of Cambridge School of Clinical Medicine, Institute of
 Metabolic Science, Cambridge, CB2 0QQ, UK.
- 8 2. Institute of Health Informatics, University College London, London, UK.
- 9 3. Health Data Research UK, London, UK.
- 10 4. British Heart Foundation Data Science Centre, London, UK.
- 5. National Institute of Health Research University College London Hospitals Biomedical
 Research Centre.
- 6. Computational Medicine, Berlin Institute of Health at Charité-Universitätsmedizin Berlin,
 10117 Berlin, Germany.
- 7. Precision Healthcare University Research Institute, Queen Mary University of London,London, UK.
- 8. Blizard Institute and SHARE Collaborative, Queen Mary University of London, London,UK.
- 19 9. Department of Infection and Immunity, Barts Health NHS Trust, London, E1 2AT, UK.
- 20
- 21 Correspondence to:
- 22 Prof Claudia Langenberg (claudia.langenberg@qmul.ac.uk, Precision Healthcare University
- 23 Research Institute, Queen Mary University of London, London, UK)

25 Abstract

To better understand sex differences in human health and disease, we conducted a systematic, large-scale investigation of sex differences in the genetic regulation of the plasma proteome (>5,000 targets), including their disease relevance.

29 Plasma levels of two-thirds of protein targets differed significantly by sex. In contrast, genetic effects on protein targets were remarkably similar, with very few protein quantitative loci 30 (pQTLs, n=74) showing significant sex-differential effects (for 3.9% and 0.3% of protein targets 31 32 from antibody- and aptamer-based platforms, respectively). Most of these 74 pQTLs 33 represented directionally concordant effects significant in both sexes, with only 21 pQTLs 34 showing evidence of sexual dimorphism, i.e. effects restricted to one sex (n=20) or with opposite directions between sexes (n=1 for CDH15). None of the sex-differential pQTLs 35 translated into sex-differential disease risk. 36

Our results demonstrate strong similarity in the genetic regulation of the plasma proteome
between sexes with important implications for genetically guided drug target discovery and
validation.

- 40
- 41 Word count: 153
- 42

43 Main text

Many aspects of human development and health, including the age of onset, prevalence, and severity of many diseases differ between sexes (1-6), but the underlying mechanisms or extent to which genetic factors contribute to any differences remain largely unknown (7-9). Understanding genetically driven sex-differences at the molecular level, specifically proteins as the biologically active entity between the genome and the phenome, is important for basic and translational genetic research, including genetically anchored drug target discovery and validation.

Here, we examined the sex-specific genetic regulation of the plasma proteome across two 51 cohorts and measurements from two different technologies. We contrast sex-differential 52 protein abundance with sex-specific genetic regulation for 4,775 unique proteins, targeted by 53 4.979 unique aptamers in 4.403 females and 3.945 males (aged 29-64) from the Fenland 54 55 study (10) and 1,463 unique proteins, targeted by 1,463 unique antibody assays among 25,904 females and 22,113 males (aged 49-60) from UK Biobank (11) (Supplementary Table 56 57 1) with 1,101 unique proteins being targeted by both platforms. We defined 'female' and 'male' sex by matching the recorded sex and sex chromosomes (XX for females and XY for males) 58 59 for both studies. The recorded sex contained a mixture of self-reported sex and sex through medical records and it was not possible to distinguish sex from gender. We acknowledge the 60 importance of distinguishing between sex and gender in research and that chromosomal 61 make-up does not always align with self-identified gender. 62

Most protein targets (n=3,457 unique proteins out of 5,100 included in this study, 67.8%) 63 showed significant sex differences (Bonferroni corrected p-value threshold for the number of 64 assays measured by each technology, see Methods) in their plasma abundance in at least 65 66 one cohort, including 521 (47.3%) overlapping targets with significant and directionally 67 concordant effects (Fig. 1, Supplementary Table 2). Results exemplified large differences between the sexes, with a slightly larger number of protein targets showing higher levels in 68 males compared to females across both technologies (Fig. 1, Supplementary Table 2). 69 70 Adjustment for hormone replacement therapy/oral contraception or known sex-differential 71 traits such as body mass index, low density lipoprotein cholesterol (LDL) levels, alanine transaminase (ALT) levels, smoking status and the frequency of alcohol consumption 72 73 impacted only a moderate number of significant differences (9.34% and 18.63%, respectively; 74 **Supplementary Table 2**). Proteins with the largest differences reflected sex-specific biology, e.g., specific expression in female- or male- specific tissues, reflected in plasma level 75 76 differences of prostate-specific antigen (12) (PSA, beta [95% confidence interval (CI)]=1.56 [1.53-1.58], p=2.31x10⁻²⁸²³, UniProt: P07288), prokineticin 1 (beta [95% CI]= 1.25 [1.24 -77

1.26], p=1.97x10⁻⁷⁰¹⁹, UniProt: P58294), or follicle stimulating hormone (FSH, beta [95% CI]=
-1.22 [-1.21 - -1.23], p=1.7x10⁻⁶⁸⁶², UniProt: P01215), while others likely reflect the effect of
sex-differences in body composition on plasma abundance of specific protein targets, such as
leptin or adiponectin. We also observed strong sex-differences in established cardiovascular
diagnostic markers such as NT-proBNP (beta [95% CI]= -0.78 [-0.74 - -0.82], p=3.33x10⁻³³⁸,
UniProt: P16860) and troponin T (beta [95% CI]= 0.83 [0.79 - 0.87], p=9.86x10⁻³⁸⁸, UniProt:
P45379).

Males and females do not only differ by disease onset and severity, but also in drug response and a higher frequency of adverse drug reactions is observed in females compared to males (13). We identified a total of 92 proteins that are the targets of already approved drugs or drugs in early clinical trials (14) and showed significant differences between sexes in plasma abundance that were directionally consistent across cohorts. While plasma protein levels are not the primary target for most of those drugs, our results can potentially help understanding sex-differential drug effects.



94

95 Figure 1: Sex differences in the abundance of 5,100 unique proteins measured by 4,979 unique 96 aptamers and 1,463 unique antibody assays. The protein targets were ordered by their effect size in 97 males. Top panel: The top panel shows the proteins for which the plasma abundance significantly 98 differed by sex in at least one technology (p_{het} <1.01x10⁻⁵ for aptamer-based and p_{het} <3.42x10⁻⁵ for antibody-based technology were used as Bonferroni-corrected thresholds respectively). The proteins 99 100 were coloured blue if they had significantly higher levels in males and red if they had higher levels in females. If the protein target was significant in both of the technologies, the effect size estimate from 101 the more significant study was displayed. The dark grey vertical lines represent the 95% confidence 102 103 intervals for the effect size estimates. Bottom panel: The bars in the bottom panel represent the proteins 104 which were targeted by both aptamer-based and antibody-based platforms. The lines were coloured lighter green if the finding was significant and directionally consistent in both technologies, darker green 105 106 if the finding was significant but not directionally consistent across technologies, lilac if the finding was 107 only significant in one of the technologies and black if the finding was not significant in any of the technologies. Results can be found in Supplementary Table 2. 108

109

We next performed sex-specific genome-proteome-wide association studies (15) to systematically identify sex-differential protein quantitative trait loci - 'sd-pQTLs' (**Supplementary Table 3 and 4**). Despite the large number of pQTLs ($p<5x10^{-8}$) identified in each sex ($n_{females}$ = 4,019, n_{males} = 3,540 pQTLs aptamer-based and $n_{females}$ =13,013, n_{males} =10,136 pQTLs for antibody-based technology (**Supplementary Figure 1**)), only very few pQTLs showed significant differences in effects between males and females (i.e. sexdifferential effects), with 15 ($p_{het}<1.01x10^{-11}$) and 59 ($p_{het}<3.42x10^{-11}$) sd-pQTLs being identified for aptamer and antibody-based platforms, respectively (Supplementary Table 3
and 4). Most sd-pQTLs reside close to the cognate gene (cis-pQTLs; n=70.3%). We observed
that the majority of sd-pQTLs (n=53 across technologies) showed effects that are directionally
consistent and significant in both sexes but with significant, small effect size differences (i.e.
same effect direction but different association strength between sexes). In other words,
identified examples were predominantly sex differential rather than sex-dimorphic (i.e. only
evident in one sex or different effect directions between sexes).

- 124 We observed no enrichment of sd-pQTLs on the X-chromosome or among druggable targets (p>0.05). We did not observe a clear bias towards protein encoding gene expression explicitly 125 in reproductive tissues or breast for the proteins for which at least one cis or trans sd-pQTL 126 was identified (16). Overall, 21 sd-pQTLs showed sex-dimorphic effects, with strong evidence 127 of effect in one but not the other sex $(p>5x10^{-8})$ for all, except for CDH15 where the sd-pQTL 128 was significant in both males and females yet showed opposite effect directions. Some of the 129 sex-dimorphic pQTLs mapped to proteins with established roles in only one of the sexes. For 130 131 example, rs10843036 was specifically associated with pregnancy zone protein (PZP) in 132 females (Fig. 2, Supplementary Table 3). Likewise, two protein targets with sd-pQTLs that 133 were significant in males only (prostate and testis expressed protein 4 [PATE4, rs499684] and Kunitz-type protease inhibitor 3 [SPIT3, rs6032259]) have been reported to be involved in 134 male fertility (Fig. 2, Supplementary Table 3). PATE4 has a reported function as a factor 135 contributing to the copulatory plug formation in male fecundity in mouse models (17) and is 136 predominantly expressed in prostate and testis (18). Similarly, SPIT3, encoded by SPINT3, is 137 138 reported to be predominantly expressed in epididymis although Spint3 was reported to be dispensable for mouse fertility (19, 20). Although its sex-specific biological function is not clear, 139 the male-specific sd-pQTLs for neural cell adhesion molecule 1 (NCAM-1) was replicated 140 across both platforms (Fig. 2, Supplementary Table 3 and 4). 141
- For two proteins, we identified more than one sd-pQTL, further supporting a sex-differential or 142 dimorphic genetic regulation (Supplementary Table 4). While the trans-sd-pQTL for cadherin 143 -15 (CDH15) was male-specific, the cis-sd-pQTL (rs113693994, beta_{females}[95% CI]= -0.16 [-144 145 0.12 - 0.19], p_{females}=9.8x10⁻²⁰, beta_{males}[95% CI]= 0.25 [0.21 - 0.28], p_{males}=9.32x10⁻³⁵) for CDH15 was the only example observed in this study where a pQTL was significant in both 146 147 sexes but with opposite effect directions. It is therefore a pQTL that has not been identified in 148 a sex-combined study (beta_{sex_combined}[95% CI]= 0.02 [-0.01 - 0.04], p_{sex_combined}=0.23). CDH15 149 acts a cell adhesion molecule that is involved in facilitating cell-cell adhesion and preserving tissue integrity and is highly expressed in brain and muscle. Similarly, carboxypeptidase E 150 (CPE) which acts as an exopeptidase essential for the activation of peptide hormones (e.g. 151 152 insulin) and neurotransmitters had both a cis and a trans sd-pQTL with both sd-pQTLs having

- 153 stronger effects in males compared to females (Supplementary Table 4). Interestingly CPE
- has been implicated to have a role in osteoclast differentiation and CPE knockout mice
- displayed low bone mineral diversity and increased osteoclastic activity as well as being obese
- and displaying a diabetic phenotype (21, 22). However, neither CDH15 nor CPE have a clearly
- 157 established sex-specific function or disease associations to date, although the fact that these
- 158 two proteins have both cis and trans sex-specific genetic regulation might suggest their
- 159 potential involvement in a sex-specific biological function.





162

Figure 2: Forest plot of all identified sex-differential protein quantitative trait loci (sd-163 pQTLs) from both aptamer- (p_{het}1.01x10⁻¹¹) and antibody-based (p_{het}3.42x10⁻¹¹) 164 technologies. The bottom panel presents sex-dimorphic pQTLs (not significant (p<5x10⁻⁸) in 165 166 one sex or has opposing effect directions), whereas the top panel presents the remaining sexdifferential pQTLs. The significant (p<5x10⁻⁸) pQTLs in each sex are represented by filled 167 circles and non-significant ones are represented by hollow circles. Horizontal lines represent 168 169 95% confidence interval of each finding. Proteins with an asterisk(*) were measured using the aptamer-based technology, otherwise using antibody-based technology. Abbreviations: MAF, 170 171 minor allele frequency.





Figure 3: Phenotypic follow up of the identified sd-pQTLs with 365 disease outcomes 174 with more than 2,500 cases in UK Biobank. A. Miami plot of association of 74 sex-175 differential pQTLs (sd-pQTLs) with 365 disease outcomes among females on the top 176 and among males at the bottom panel. The disease categories have been coloured by 177 178 disease categories. The horizontal dashed line represents a suggestive significance threshold of (p<1x10⁻⁵). B. Comparison of odds ratios for the sd-pQTL – disease associations 179 which meet the suggestive significance threshold (p<1x10⁻⁵) in males or females. The 180 diagonal dashed line represents the equality line (x=y). C. Comparison of $-\log_{10}(\log_{10})$ 181

transformed Pvalues for the sd-pQTL – disease associations which meet the suggestive
 significance threshold (p<1x10⁻⁵) in males or females. The diagonal dashed line
 represents the equality line (x=y).

185

We next attempted to identify any potential sex-differential phenotypic consequences of sdpQTLs (n=74) across 365 diseases with more than 2,500 cases in UK Biobank. Despite 74 significant associations (Bonferroni corrected significance threshold of $p<2.32x10^{-6}$ and $p<9.31x10^{-6}$ for antibody- and aptamer-based technologies, respectively) between sd-pQTLs and disease risk in at least one sex, none of the associations showed evidence of sexdifferential effects (**Supplementary Table 5 and 6**). This leaves the downstream physiological or pathological consequences of the identified sd-pQTLs yet to be determined (**Fig. 3**).

In summary, we identified sd-pQTLs for only a small proportion of the protein targets tested (0.3% and 3.9% of unique assays respectively for aptamer- and antibody-based technologies), a finding in line with previous sex-stratified analyses of tissue-specific gene-expression (16), and further confirmed that most sd-pQTLs act in a sex-differential rather than sex-dimorphic manner. This is in line with a recent study that reported trait variance difference between sexes can predominantly be explained by sex-differential 'amplification effect' (i.e. same effect direction yet different magnitudes of strength between sexes) (23).

200 Our study highlights two important conclusions. Firstly, the fact that we observe sd-pQTLs for 201 a very small percentage of protein targets despite large differences in plasma protein levels emphasizes that the observed sex-differences are likely to be a result of intrinsic (e.g. hormone 202 profiles) and extrinsic mechanisms (e.g. sex-differential lifestyle and risk-factor profiles) other 203 204 than genetic regulation. This finding is in line with what has been reported for sex-differences 205 observed for complex diseases, sex-differential genetic signals being identified for only a small 206 proportion of common diseases (8). Secondly, our results suggest that the use of pQTLs in 207 biomedical research, specifically for drug target discovery and causal inference will - with few 208 exceptions - likely generate findings that are generalisable across sexes for the studied protein targets. However, larger studies should continue to evaluate sex differences as increased 209 power could potentially uncover additional examples with biologically relevant sex-dimorphic 210 effects. 211

Although only few, we did identify some sex-dimorphic genetic effects, with some reflecting sex-specific biology (e.g. PATE4, SPIT3, pregnancy zone protein) (24) acting possibly via steroid hormone responsive elements. Some of the other effects might possibly the result from differential environmental exposures between the sexes, as suggested for genetic variants affecting the risk for gout that may act through differential alcohol consumption as shown previously (25). 218 The restriction to proteins measured in plasma represents a notable limitation of our study, as 219 sex-differential proteogenomic effects within tissues will unlikely be systematically reflected in 220 plasma via secretion, natural cell turnover, or leakage. We obtained some evidence that larger sample sizes can identify a greater number of significant sd-pQTLs, but most act as weak 221 modifiers of strong overall effects at protein encoding loci, and larger studies may possibly 222 reveal even more subtle differences in regulation in trans for the previously targeted protein 223 targets. Given the incomplete proteomic coverage (n=4,775 and n=1,463 unique proteins 224 targeted by aptamer- and antibody-based platforms respectively, within a spectrum of over 225 226 20,000 proteins without taking post-translational modifications or different isoforms into 227 account) as well as limited coverage of the genomic variant spectrum (i.e. rare variants or 228 potentially ancestry-specific effects that we were not able to investigate), future studies might 229 uncover new sd-pQTL signals as genomic and proteomic coverage continues to improve.

230 Main text word count: 1753

231 Methods

232 Study Participants

The Fenland study (26) is a population-based cohort of 12,435 participants of generally white-233 European ancestry, born between 1950 and 1975 who underwent detailed phenotyping at the 234 235 baseline visit from 2005 to 2015. Participants were recruited from general practice surgeries in the Cambridgeshire region in the UK. The participants were excluded from the study if they 236 were (i) clinically diagnosed with diabetes mellitus or a psychotic disorder, or (ii) pregnant or 237 lactating, (iii) unable to walk unaided, or (iv) had a terminal illness. The study was approved 238 by the Cambridge Local Research Ethics Committee (NRES Committee - East of England, 239 Cambridge Central, ref. 04/Q0108/19) and all participants provided written informed consent. 240

241 This study used the largest subset of individuals from the Fenland study (Supplementary 242
 Table 1). 8,348 samples with both genotype information and proteomics measurements were
 taken forward for analyses after excluding ancestry outliers, related individuals or samples 243 which have failed proteomics QC. The samples were well-balanced in terms of the participants 244 from each sex: 4,403 (52.7%) females and 3,945 (47.3%) males were included in the study. 245 246 Sex variable in Fenland study was based on general practitioners (GP) records. We only included participants with matching entries for the recorded sex and sex chromosomes (XX 247 for females and XY for males). Individuals without matching entries were excluded from the 248 study as a part of quality control as a mismatch can sometimes also be indicative of issues 249 250 with genotyping protocol.

UK Biobank is a large-scale, population-based cohort with deep genetic and phenotypic data 251 with the full cohort consisting of approximately 500,000 participants (11). The participants were 252 253 recruited across centres in United Kingdom and were aged 40 to 69 years at the time of 254 recruitment (11). Ethics approval for the UK Biobank study was obtained from the North West Centre for Research Ethics Committee (11/NW/0382) (11) and all participants provided 255 informed consent. This study used the subset of European-ancestry individuals from UK 256 Biobank where both genotype and proteomics measurements were available after excluding 257 ancestry outliers or samples which have failed genomic or proteomics QC (n=48,017). 25,904 258 (53.9%) females and 22,113 (46.1%) males were included in the study (Supplementary Table 259 260 1). Sex in UK Biobank had two definitions, one was based on sex chromosomes (field 22001) and the other was contained a mixture of the sex the NHS had recorded for the participant 261 and self-reported sex (field 31). We only included participants with matching entries for the 262 recorded sex (from medical records or self-reported) and sex chromosomes (XX for females 263 and XY for males). 264

266

267 Genotyping and imputation

The Fenland-OMICS samples have been genotyped using the Affymetrix UK Biobank Axiom 268 array. Sample-level and variant level QC criteria were applied as described elsewhere (26). In 269 summary, the genotyped data was imputed to the HRC (r1) panel (27) using IMPUTE4 270 (https://jmarchini.org/software/) for the autosomes and Sanger Imputation Server for 271 chromosome X (https://imputation.sanger.ac.uk/). The data was also imputed to the UK10K 272 and 1000 Genomes Project 3 panels using and Sanger Imputation Server for both autosomes 273 and chromosome X (28). Additional variants gained from the UK10Kp+1KGp3 imputation were 274 added to the HRC imputed dataset. For basic quality control, variants were filtered for minor 275 allele count (MAC) ≥3 using BCFtools (29) and INFO≥0.4 using QCTOOL v2.0.2 276 277 (https://www.well.ox.ac.uk/~gav/gctool v2/) to eliminate variants with low imputation guality.

The UK Biobank samples were genotyped using the Affymetrix UK BiLEVE or the Affymetrix UK Biobank Axiom arrays. The following QC criteria was applied to the genotyping data (a) routine quality checks carried out during the process of sample retrieval, DNA extraction, and genotype calling; (b) checks and filters for genotype batch effects, plate effects, departures from Hardy Weinberg equilibrium, sex effects, array effects, and discordance across control replicates; and (c) individual and genetic variant call rate filters as previously described (11).

284 Genomic build GRCh37 was used throughout this study.

285

286 Proteomic measurements

287 <u>Aptamer-based platform</u>

Fasting proteomic profiling of EDTA samples from Fenland study participants was performed 288 by SomaLogic Inc. using the SOMAscan proteomic assay (v4). Relative protein abundances 289 of 4,775 human protein targets were measured by 4,979 aptamers (SomaLogic V4). The 290 quality control of the proteomic measurements has been described in detail previously (26). 291 292 Briefly, hybridization control probes were used to generate a hybridization scale factor to account for variation in hybridization within runs. A ratio between each aptamer's measured 293 value and a reference value were computed to control for total signal differences between 294 samples due to variation in overall protein concentration or technical factors. The median of 295 296 these ratios was computed and applied to each dilution set (40%, 1% and 0.005%). Samples 297 were removed if they were deemed by SomaLogic to have failed or did not meet our 298 acceptance criteria of 0.25-4 for all scaling factors. In addition to passing SomaLogic QC,

aptamers were filtered to only include human protein targets for subsequent analysis
 (n=4,979). Aptamers' target annotation and mapping to UniProt (30) accession numbers as
 well as Entrez gene identifiers (31) were provided by SomaLogic and these were used those
 to determine genomic positions of protein encoding genes.

303 Antibody-based platform

304 The UK Biobank proteomic measurements were conducted by antibody-based Olink technology, Explore 1536 platform which uses Proximity Extension Assay (32). In summary, 305 306 each protein is targeted by two unique antibodies with unique complimentary oligonucleotides. which only hybridize when they come into close proximity. This is subsequently quantified by 307 308 next-generation sequencing. Normalized protein expression (NPX) units, which are reported on a log2 scale, are generated by normalization of the extension control and further 309 310 normalization of the plate control. Further details about antibody-based proteomic 311 measurements and QC have been described elsewhere, including the exclusion of samples due to poor quality and selective measurements with assay warnings (33). 312

313

314 Sex-differences in protein abundances

We assessed the differential abundance levels of the 4,979 SomaLogic V4 aptamers between sexes in Fenland study. To estimate the effect of sex, a linear regression model was for implemented in R 3.6, using the inverse rank normalized proteomic values and including covariates age and test site in the model. A stringent Bonferroni-corrected threshold (corrected for n=4,979 aptamers; $p<1.01x10^{-5}$) was applied.

320 1,463 protein targets from Olink Explore 1536 platform in UK Biobank were inverse rank 321 normalized and subsequently restricted cubic splines function through 'rsc' function of 'Hmisc' package was applied to regress out technical covariates such as month of the blood draw, 322 time that blood was drawn, fasting status and sample age in R v4.2.2. Similarly, the effect of 323 324 sex assessed in abundance levels of the 1,463 protein targets from Olink Explore 1536 platform were assessed in a linear regression model using the inverse rank normalized 325 residuals and including covariates age, age² and proteomic batch in R v4.2.2. A stringent 326 Bonferroni-corrected threshold (corrected for n=1,463 aptamers; $p<3.42x10^{-5}$) was applied. 327

Sensitivity analyses was performed by including (a) participants who have undergone hormone replacement therapy or use oral contraception, or (b) for known sex-differential participant characteristics which were body mass index (BMI), low density lipoprotein (LDL) cholesterol levels, alanine transaminase (ALT) levels, smoking status and the frequency of alcohol consumption **(Supplementary Table 1)** as additional covariates in the analyses. The

continuous variables (BMI, LDL and ALT) were inverse rank normalized before being includedas covariates.

335

336 Sex-stratified protein genome-wide association analysis (pGWASs)

For the aptamer-based platform, the protein abundances for 4,979 aptamers measured in Fenland study were inverse rank normalized and regressed for covariates age, test site and the first 10 genetic principal components in R v3.6. The residuals for each sex were used in the subsequent association analyses.

The fastGWA software (34) linear regression analysis was performed through GCTA version 1.93.2 for the sex-stratified genome-wide association analysis in each sex. Further variant level QC was also applied and only variants with MAC \geq 3, INFO \geq 0.4, genotype missingness rate < 5% and MAF>1% were included in the downstream analyses.

345 For the antibody-based platform, the protein abundances for 1,463 assays measured in UK 346 Biobank, the same residuals from the analyses of sex-differences in protein abundances (i.e. inverse rank normalized and technical covariates regressed out) were taken forward. Sex-347 stratified GWASs were performed using REGENIE v.3.4.1 (35) through performing two steps, 348 349 as implemented by the software. In the first step, a whole-genome regression model is fitted for each phenotype to generate a covariate, which is subsequently included in the second step 350 to allow for computationally-efficient analyses of a large number of phenotypes. For the first 351 step, only high-quality variants passing the stringent QC criteria of MAF>1%, MAC>100, 352 Hardy-Weinberg equilibrium p-value< 1×10^{-15} and genotype missingness rate < 10% were 353 354 used and SNPs were pruned for linkage-disequilibrium (LD), specified for 1000 variant 355 windows, 100 sliding windows and r2<0.8 through Plink v.1.9. Subsequently, step 2 was applied to conduct sex-stratified the genome-wide association analyses for 1,463 protein 356 targets with additional per-marker QC filters of MAC>50, MAF>1% and INFO>0.4. 357

358

359 Heterogeneity analysis

We performed an inverse-variance fixed effects meta-analysis for each protein target using female-only and male-only summary statistics through METAL (v.2011-03-25) (34) to assess the heterogeneity in the genetic associations between sexes for each platform. We used a proteome and genome-wide Bonferroni corrected significance threshold (p_{het} <1.01x10⁻¹¹ and p_{het} <3.42x10⁻¹¹ respectively for aptamer- and antibody-based platforms) for heterogeneity pvalue to define sex-differential protein quantitative trait loci (i.e. sd-pQTLs).

Significant genomic regions were defined by 1 Mb regions (±500 Kb on either side) around any variant with significant heterogeneity. The MHC region (chr6:25.5–34.0Mb) was treated as a single region. The regional sentinel variant for each genomic loci was defined as the most significant variant within the region. Variants were defined as cis-pQTLs if they were within the 1 Mb window (±500 Kb on either side) of the protein encoding gene and defined as transpQTLs if they were not within the 1 Mb window.

372

373 Phenome-wide Association Study (PheWAS)

374 We have tested whether any of the significant sd-pQTLs showed heterogeneity between sexes 375 in terms of their disease associations across the phenome. For this purpose, in each sex, we 376 tested the association of sd-pQTLs with 365 binary diseases with more than 2,500 cases in 377 UK Biobank. The binary disease categories were collated through clinical entities named 378 'phecodes' in UK Biobank, which were defined using the International Classification of Diseases, 10th Revision (ICD-10) and the International Classification of Diseases, 10th 379 Revision, Clinical Modification (ICD-10-CM) codes from electronic health records, available in 380 381 UK Biobank (36). We tested the association of each sd-pQTL with each phecode in each sex, using a logistic regression model in R v3.6 and adjusting for age, genotype batch, test centre, 382 and the first ten genetic principal components in unrelated European participants. We have 383 subsequently meta-analysed the female-only and male-only summary statistics using a fixed-384 effects meta-analyses through metafor package in R v3.6 to assess the heterogeneity of the 385 association between sexes. To correct for multiple testing, p-value threshold for PheWAS was 386 defined as p_{het} < 9.31x10⁻⁶ and p_{het} < 2.32x10⁻⁶ for aptamer- and antibody-based platforms 387 respectively, which were corrected for the number of sd-pQTLs and number of phenotypes 388 389 tested (n=365) in each platform.

391 **References**

1. Mauvais-Jarvis F, Bairey Merz N, Barnes PJ, Brinton RD, Carrero JJ, DeMeo DL, et

al. Sex and gender: modifiers of health, disease, and medicine. Lancet.

394 2020;396(10250):565-82.

Zein JG, Denson JL, Wechsler ME. Asthma over the Adult Life Course: Gender and
Hormonal Influences. Clin Chest Med. 2019;40(1):149-61.

397 3. Siegel RL, Miller KD, Fuchs HE, Jemal A. Cancer Statistics, 2021. CA Cancer J Clin.
398 2021;71(1):7-33.

Wagner AD, Oertelt-Prigione S, Adjei A, Buclin T, Cristina V, Csajka C, et al. Gender
medicine and oncology: report and consensus of an ESMO workshop. Ann Oncol.

401 2019;30(12):1914-24.

402 5. Clayton JA. Studying both sexes: a guiding principle for biomedicine. FASEB J.
403 2016;30(2):519-24.

Klein SL, Schiebinger L, Stefanick ML, Cahill L, Danska J, de Vries GJ, et al.
Opinion: Sex inclusion in basic research drives discovery. Proc Natl Acad Sci U S A.
2015;112(17):5257-8.

407 7. Khramtsova EA, Davis LK, Stranger BE. The role of sex in the genomics of human
408 complex traits. Nat Rev Genet. 2019;20(3):173-90.

8. Bernabeu E, Canela-Xandri O, Rawlik K, Talenti A, Prendergast J, Tenesa A. Sex
differences in genetic architecture in the UK Biobank. Nat Genet. 2021;53(9):1283-9.

411 9. Lagou V, Mägi R, Hottenga JJ, Grallert H, Perry JRB, Bouatia-Naji N, et al. Sex412 dimorphic genetic effects and novel loci for fasting glucose and insulin variability. Nat
413 Commun. 2021;12(1):24.

10. Pietzner M, Wheeler E, Carrasco-Zanini J, Cortes A, Koprulu M, Wörheide MA, et al.

415 Mapping the proteo-genomic convergence of human diseases. Science.

416 2021;374(6569):eabj1541.

417 11. Bycroft C, Freeman C, Petkova D, Band G, Elliott LT, Sharp K, et al. The UK Biobank
418 resource with deep phenotyping and genomic data. Nature. 2018;562(7726):203-9.

419 12. Nordström T, Akre O, Aly M, Grönberg H, Eklund M. Prostate-specific antigen (PSA)

420 density in the diagnostic algorithm of prostate cancer. Prostate Cancer Prostatic Dis.

421 2018;21(1):57-63.

Madla CM, Gavins FKH, Merchant HA, Orlu M, Murdan S, Basit AW. Let's talk about
sex: Differences in drug therapy in males and females. Adv Drug Deliv Rev.

424 2021;175:113804.

425 14. Finan C, Gaulton A, Kruger FA, Lumbers RT, Shah T, Engmann J, et al. The

druggable genome and support for target identification and validation in drug development.Sci Transl Med. 2017;9(383).

Pietzner M, Wheeler E, Carrasco-Zanini J, Cortes A, Koprulu M, Wörheide MA, et al.
Mapping the proteo-genomic convergence of human diseases. 2021.

430 16. Oliva M, Muñoz-Aguirre M, Kim-Hellmuth S, Wucher V, Gewirtz ADH, Cotter DJ, et

al. The impact of sex on gene expression across human tissues. Science. 2020;369(6509).

432 17. Noda T, Fujihara Y, Matsumura T, Oura S, Kobayashi S, Ikawa M. Seminal vesicle

433 secretory protein 7, PATE4, is not required for sperm function but for copulatory plug
434 formation to ensure fecundity⁺. Biol Reprod. 2019;100(4):1035-45.

435 18. Consortium G. The GTEx Consortium atlas of genetic regulatory effects across
436 human tissues. Science. 2020;369(6509):1318-30.

19. Clauss A, Persson M, Lilja H, Lundwall Å. Three genes expressing Kunitz domains in
the epididymis are related to genes of WFDC-type protease inhibitors and semen coagulum
proteins in spite of lacking similarity between their protein products. BMC Biochem.
2011;12:55.

20. Robertson MJ, Kent K, Tharp N, Nozawa K, Dean L, Mathew M, et al. Large-scale
discovery of male reproductive tract-specific genes through analysis of RNA-seq datasets.
BMC Biol. 2020;18(1):103.

444 21. Cawley NX, Yanik T, Woronowicz A, Chang W, Marini JC, Loh YP. Obese

445 carboxypeptidase E knockout mice exhibit multiple defects in peptide hormone processing

446 contributing to low bone mineral density. Am J Physiol Endocrinol Metab. 2010;299(2):E189-447 97.

Kim HJ, Hong J, Yoon HJ, Yoon YR, Kim SY. Carboxypeptidase E is a novel
modulator of RANKL-induced osteoclast differentiation. Mol Cells. 2014;37(9):685-90.

Zhu C, Ming MJ, Cole JM, Edge MD, Kirkpatrick M, Harpak A. Amplification is the
primary mode of gene-by-sex interaction in complex human traits. Cell Genomics. 2023;3(5).

452 24. Gegenhuber B, Tollkuhn J. Signatures of sex: Sex differences in gene expression in
453 the vertebrate brain. Wiley Interdiscip Rev Dev Biol. 2020;9(1):e348.

25. Zhou W, Kanai M, Wu KH, Rasheed H, Tsuo K, Hirbo JB, et al. Global Biobank Metaanalysis Initiative: Powering genetic discovery across human disease. Cell Genom.
2022;2(10):100192.

Pietzner M, Wheeler E, Carrasco-Zanini J, Raffler J, Kerrison ND, Oerton E, et al.
Genetic architecture of host proteins involved in SARS-CoV-2 infection. Nat Commun.

459 2020;11(1):6397.

460 27. McCarthy S, Das S, Kretzschmar W, Delaneau O, Wood AR, Teumer A, et al. A

reference panel of 64,976 haplotypes for genotype imputation. Nat Genet.

462 2016;48(10):1279-83.

463 28. Huang J, Howie B, McCarthy S, Memari Y, Walter K, Min JL, et al. Improved

imputation of low-frequency and rare variants using the UK10K haplotype reference panel.

465 Nat Commun. 2015;6:8111.

- 29. Danecek P, Bonfield JK, Liddle J, Marshall J, Ohan V, Pollard MO, et al. Twelve
 years of SAMtools and BCFtools. Gigascience. 2021;10(2).
- 468 30. The UniProt Consortium. UniProt: the universal protein knowledgebase. Nucleic
 469 Acids Res. 2017;45(D1):D158-D69.
- 470 31. Maglott D, Ostell J, Pruitt KD, Tatusova T. Entrez Gene: gene-centered information
 471 at NCBI. Nucleic Acids Res. 2007;35(Database issue):D26-31.

472 32. Assarsson E, Lundberg M, Holmquist G, Björkesten J, Thorsen SB, Ekman D, et al.

Homogenous 96-plex PEA immunoassay exhibiting high sensitivity, specificity, and excellent
scalability. PLoS One. 2014;9(4):e95192.

- 475 33. Sun BB, Chiou J, Traylor M, Benner C, Hsu YH, Richardson TG, et al. Plasma
- 476 proteomic associations with genetics and health in the UK Biobank. Nature.
- 477 2023;622(7982):329-38.
- 478 34. Jiang L, Zheng Z, Qi T, Kemper KE, Wray NR, Visscher PM, et al. A resource-
- 479 efficient tool for mixed model association analysis of large-scale data. Nat Genet.

480 2019;51(12):1749-55.

- 481 35. Mbatchou J, Barnard L, Backman J, Marcketta A, Kosmicki JA, Ziyatdinov A, et al.
 482 Computationally efficient whole-genome regression for quantitative and binary traits. Nat
- 483 Genet. 2021;53(7):1097-103.
- 484 36. Wu P, Gifford A, Meng X, Li X, Campbell H, Varley T, et al. Mapping ICD-10 and
- 485 ICD-10-CM Codes to Phecodes: Workflow Development and Initial Evaluation. JMIR Med
- 486 Inform. 2019;7(4):e14325.
- 487

489 Acknowledgements

We are grateful to all Fenland volunteers and to the General Practitioners and practice staff for assistance with recruitment. We thank the Fenland Study Investigators, Fenland Study Coordination team and the Epidemiology Field, Data and Laboratory teams. SomaLogic proteomic measurements were supported and governed by a collaboration agreement between the University of Cambridge and SomaLogic.

- The Fenland Study (DOI 10.22025/2017.10.101.00001) is funded by the Medical Research 495 496 Council (MC_UU_12015/1). We further acknowledge support for genomics from the Medical Research Council (MC_PC_13046). This work is supported by the Medical Research Council 497 (MC UU 00006/1 - Etiology and Mechanisms) (C.L., E.W., M.P., N.K., and N.J.W.). M.K. is 498 supported by Gates Cambridge Trust. H.H. is supported by Health Data Research UK and the 499 NIHR University College London Hospitals Biomedical Research Centre. S.D. is supported by 500 a) the BHF Data Science Centre led by HDR UK (grant SP/19/3/34678), b) BigData@Heart 501 502 Consortium, funded by the Innovative Medicines Initiative-2 Joint Undertaking under grant 503 agreement 116074, c) the NIHR Biomedical Research Centre at University College London Hospital NHS Trust (UCLH BRC), d) a BHF Accelerator Award (AA/18/6/24223), e) the CVD-504 505 COVID-UK/COVID-IMPACT consortium and f) the Multimorbidity Mechanism and Therapeutic 506 Research Collaborative (MMTRC, grant number MR/V033867/1). J.C.Z. was supported by a 4-year Wellcome Trust PhD Studentship and the Cambridge Trust. 507
- 508

509 Competing interests

510 E.W. is now an employee of AstraZeneca.

511

512 Author contributions

M.K., M.P and C.L. designed the analysis and drafted the manuscript. M.K., E.W., S.D., N.K.,
J.C.Z, H.H. and M.P. have performed the quality control, data preparation or the bioinformatics
analyses. N.J.W. is PI of the Fenland study. C.M.O. contributed to defining sex in this study
and provided insights into broader concepts of sex and gender in research. All authors
contributed to the interpretation of the results and critically reviewed the manuscript.

518

519 Data availability

- 520 Data from the Fenland cohort can be requested by bona fide researchers for specified 521 scientific purposes via the study website (www.mrc-
- 522 epid.cam.ac.uk/research/studies/fenland/information-for-researchers/). Sex-stratified

523 summary statistics will be made available upon publication.

- 524 Access to the UK Biobank genomic, proteomic and phenotype data is open to all approved
- 525 health researchers (http://www.ukbiobank.ac.uk/). This research has been conducted using
- the UK Biobank resource under the application 44448.

527