

# Genetic associations of protein-coding variants in human disease

Benjamin B. Sun<sup>\*1,2</sup>, Mitja I. Kurki<sup>3,4,5,6</sup>, Christopher N. Foley<sup>7,8</sup>, Asma Mechakra<sup>9</sup>, Chia-Yen Chen<sup>1</sup>, Eric Marshall<sup>1</sup>, Jemma B. Wilk<sup>1</sup>, Biogen Biobank Team<sup>1</sup>, Mohamed Chahine<sup>10</sup>, Philippe Chevalier<sup>9</sup>, Georges Christé<sup>9</sup>, FinnGen<sup>11</sup>, Aarno Palotie<sup>3,4,5,6</sup>, Mark J. Daly<sup>3,4,5,6</sup>, Heiko Runz<sup>\*1</sup>.

1. Translational Biology, Research & Development, Biogen Inc., Cambridge, MA, US
2. BHF Cardiovascular Epidemiology Unit, Department of Public Health and Primary Care, University of Cambridge, Cambridge, UK
3. Psychiatric & Neurodevelopmental Genetics Unit, Massachusetts General Hospital, Boston, MA, US
4. The Stanley Center for Psychiatric Research, The Broad Institute of MIT and Harvard, Cambridge, MA, US
5. Institute for Molecular Medicine Finland (FIMM), University of Helsinki, Helsinki, Finland
6. Analytic and Translational Genetics Unit, Department of Medicine, Massachusetts General Hospital, Boston, MA, US
7. MRC Biostatistics Unit, School of Clinical Medicine, University of Cambridge, Cambridge, UK
8. Optima Partners, Edinburgh, UK
9. Université Claude Bernard Lyon 1, EA4612 Neurocardiology, Lyon, France
10. CERVO Brain Research Center and Department of Medicine, Faculty of Medicine, Université Laval, Quebec City, QC, Canada
11. FinnGen, Finland

\*Correspondence: [bbsun92@outlook.com](mailto:bbsun92@outlook.com) (B.B.S.), [heiko.runz@gmail.com](mailto:heiko.runz@gmail.com) (H.R.)

## 25 **Abstract**

26 Genome-wide association studies (GWAS) have identified thousands of genetic variants linked  
27 to the risk of human disease. However, GWAS have thus far remained largely underpowered  
28 to identify associations in the rare and low frequency allelic spectrum and have lacked the  
29 resolution to trace causal mechanisms to underlying genes. Here, we combined whole exome  
30 sequencing in 392,814 UK Biobank participants with imputed genotypes from 260,405  
31 FinnGen participants (653,219 total individuals) to conduct association meta-analyses for 744  
32 disease endpoints across the protein-coding allelic frequency spectrum, bridging the gap  
33 between common and rare variant studies. We identified 975 associations, with more than one-  
34 third of our findings not reported previously. We demonstrate population-level relevance for  
35 mutations previously ascribed to causing single-gene disorders, map GWAS associations to  
36 likely causal genes, explain disease mechanisms, and systematically relate disease associations  
37 to levels of 117 biomarkers and clinical-stage drug targets. Combining sequencing and  
38 genotyping in two population biobanks allowed us to benefit from increased power to detect  
39 and explain disease associations, validate findings through replication and propose medical  
40 actionability for rare genetic variants. Our study provides a compendium of protein-coding  
41 variant associations for future insights into disease biology and drug discovery.

42

## 43 **Introduction**

44 Inherited protein-coding and non-coding DNA variations play a role in the risk, onset, and  
45 progression of human disease. Traditionally, geneticists have dichotomized diseases as either  
46 caused by coding mutations in single genes that tend to be rare, highly penetrant, and often  
47 compromise survival and reproduction (often termed “Mendelian” diseases), or alternatively  
48 as common diseases that show a complex pattern of inheritance influenced by the joint  
49 contributions of hundreds of low-impact, typically non-coding genetic variants (often termed  
50 “complex” diseases). For both rare and common conditions, large human cohorts  
51 systematically characterized for a respective trait of interest have enabled the identification of  
52 thousands of disease-relevant variants through either sequencing-based approaches or genome-  
53 wide association studies (GWAS). Nevertheless, the exact causal alleles and mechanisms that  
54 underlie associations of genetic variants to disease have thus far remained largely elusive<sup>1</sup>.

55  
56 In recent years, population biobanks have been added to the toolkit for disease gene discovery.  
57 Biobanks provide the opportunity to simultaneously investigate multiple traits and diseases at  
58 once and uncover relationships between previously unconnected phenotypes. For instance, the  
59 UK Biobank (UKB) is a resource that captures detailed phenotype information matched to  
60 genetic data for over 500,000 individuals and, since its inception, has facilitated biomedical  
61 discoveries at an unprecedented scale<sup>2</sup>. We and others have recently reported on the ongoing  
62 efforts to sequence the exomes of all UKB participants and link genetic findings to a broad  
63 range of phenotypes<sup>3-5</sup>. We also established FinnGen (FG), an academic-industry collaboration  
64 to identify genotype-phenotype correlations in the Finnish founder population with the aim to  
65 better understand how the genome affects health (<https://www.finnngen.fi>). Finland is a well-  
66 established genetic isolate and a unique gene pool distinguishes Finns from other Europeans<sup>6</sup>.  
67 The distinct Finnish haplotype structure is characterized by large blocks of co-inherited DNA

68 in linkage disequilibrium and an enrichment for alleles that are rare in other populations, but  
69 can still be confidently imputed from genotyping data even in the rare and ultra-rare allele  
70 frequency spectrum<sup>7-9</sup>. Through combining imputed genotypes with detailed phenotypes  
71 ascertained through national registries, FG holds the promise to provide particular insights into  
72 the yet little examined allele frequency spectrum between 0.1 and 2% where both sequencing  
73 studies and GWAS have thus far remained largely underpowered to identify associations to  
74 disease. This spectrum includes many coding variants with moderate to large effect sizes that  
75 can help identify causal genes in GWAS loci, provide mechanistic insights into disease  
76 pathologies, and potentially bridge rare and common diseases.

77

78 Here, we have leveraged the combined power of UKB and FG to investigate how rare and low-  
79 frequency variants in protein-coding regions of the genome contribute to the risk for human  
80 traits and diseases. Using data from a total of 653,219 individuals, we tested how ~48,000  
81 coding variants identified in both biobanks through either whole-exome sequencing (WES) or  
82 genotype imputation associate with 744 distinct disease endpoints. Disease associations were  
83 compared against information from rare disease, biomarker and drug target resources and  
84 complemented by deep dives into distinct disease mechanisms of individual genes and coding  
85 variants. Our results showcase the benefits of combining large population cohorts to discover  
86 and replicate novel associations, explain disease mechanisms across a range of common and  
87 rare diseases, and shed light on a substantial gap in the allelic spectrum that neither genotyping  
88 nor sequencing studies have previously been able to address.

89

## 90 **Results**

91 An overview of the study design and basic demographics are provided in **Extended Data**  
92 **Figure 1 and Supplementary Table 1**. In brief, we systematically harmonised disease



93 phenotypes across UKB and FG using Phecode and ICD10 mappings and retained 744 specific  
94 disease endpoints grouped into 580 disease clusters that span a broad range of diseases  
95 (**Methods, Supplementary Table 2**). Disease case counts relative to cohort size showed good  
96 correlations both, overall between UKB and FG (Spearman's  $\rho=0.65$ ,  $p<5.3\times 10^{-90}$ ) and across  
97 distinct disease groups (**Extended Data Figure 2**).

98

## 99 **Coding-wide association analyses in 653,219 individuals across 744 disease** 100 **endpoints identify 975 genetic signals**

101 We performed coding-wide association studies (CWAS) across 744 disease endpoints over a  
102 mean of 48,189 (range: 25,309-89,993, **Methods, Supplementary Table 2**) post-QC coding  
103 variants across the allele frequency spectrum derived from whole-exome sequencing (WES) of  
104 392,814 European ancestry individuals in UKB and meta-analysed these data with summary  
105 results from up to 260,405 individuals in FG (**Methods, Supplementary Table 2**).

106

107 We identified 975 associations (534 variants in 301 distinct regions across 148 disease clusters;  
108 620 distinct region-disease cluster associations) meeting genome-wide significance ( $p<5\times 10^{-8}$ ),  
109 and 717 associations (378 variants in 231 distinct regions across 121 disease clusters; 445  
110 distinct region-disease cluster associations) at a conservative (Bonferroni) multiple testing  
111 threshold of  $p<2\times 10^{-9}$  (correcting for the number of approximate independent tests) (**Methods,**  
112 **Figure 1a, Supplementary Figure 1 (interactive), Supplementary Table 3**). The  
113 distributions of coding variant annotation categories were largely similar for variants with at  
114 least one significant association ( $p<5\times 10^{-8}$ ) relative to all variants tested, with missense variants  
115 showing a higher fraction of significant variants than in-frame indel or predicted loss-of-  
116 function (pLoF) variants (**Extended Data Figure 3**). Inflation was well controlled with a mean

117 genomic inflation factor of 1.04 (5-95 percentiles: 1.00-1.09, **Extended Data Figure 4a**).  
118 Effect sizes were generally well aligned between UKB and FG (Spearman's  $\rho=0.90$ ,  $p<10^{-300}$ ,  
119 **Extended Data Figure 4b**). MAFs of lead variants correlated well overall between UKB and  
120 FG (Spearman's  $\rho=0.97$ ,  $p<10^{-300}$ , **Figure 1b**), especially for variants with  $MAF>1\%$ , yet as  
121 expected<sup>8</sup> from genetic differences between Finns and non-Finnish Europeans (NFEs) was  
122 reduced for variants with  $MAF<1\%$  (Spearman's  $\rho=0.32$ ,  $p=0.023$ ).

123  
124 Across all diseases, we found generally larger effect sizes for low frequency and rare variants  
125 (**Figure 1c**). 387 of the 975 (39.7% at  $p<5\times 10^{-8}$ ; 270/717 (37.7%) at  $p<2\times 10^{-9}$ ) associations  
126 would not have been detected if analysed in UKB (61.5% at  $p>5\times 10^{-8}$ ; 60.1% at  $p>2\times 10^{-9}$ ) or  
127 FG (59.6% at  $p>5\times 10^{-8}$ ; 58.6% at  $p>2\times 10^{-9}$ ) alone. Association testing within UKB and FG  
128 individually would have yielded 318 and 479 associations respectively at  $p<5\times 10^{-8}$   
129 (**Supplementary Tables 4 and 5**). Thus our combined approach utilizing both biobanks  
130 increased the number of significant findings by approximately 3- and 2-fold, respectively. Of  
131 the 318 and 479 significant sentinel variants in UKB and FG, 252 (72.6%) and 258 (53.9%)  
132 replicated at  $p<0.05$  in FG and UKB respectively (**Supplementary Tables 4 and 5**),  
133 highlighting further the strength of our approach to yield results that are more robust through  
134 replication than would be findings derived from just a single biobank.

135  
136 Our study benefits from population enrichment of rare alleles in Finns versus NFEs (and vice  
137 versa) that increases the power for association discovery. Using a combination of theoretical  
138 analyses and empirical simulations, we show that by leveraging population-enriched variants  
139 we could increase inverse-variance weighted meta-analysis Z-scores and hence our ability to  
140 detect underlying associations. The gain in power from enriched alleles was present across a  
141 range of rare MAFs (0.01-1%), with the strongest power gain in the rare and ultra-rare minor

142 allele frequency (MAF) range of 0.01% to 0.25% (**Figure 1d, Supplementary Information,**  
143 **Extended Data Figure 5, Supplementary Figure 2 (interactive)**). Of the sentinel variants,  
144 we found 73 (33 in UKB, 40 in FG) to be enriched by >2-fold and 23 (8 in UKB, 15 in FG)  
145 by >4-fold relative to the respective other biobank (**Figure 1b, Supplementary Table 6**). The  
146 majority of highly population enriched variants are rare (MAF<1%) or low frequency (MAF  
147 1-5%), whereby 20 of 23 variants with >4-fold population enrichment (13 in FG and 7 in UKB)  
148 had MAF <1% (**Table 1, Supplementary Table 6**).

149  
150 We systematically cross-referenced our results with previously described GWAS associations  
151 (via GWAS Catalog<sup>10</sup> and PhenoScanner<sup>11</sup>) and disease relevance as reported in ClinVar<sup>12</sup>  
152 (**Methods**). In total, we found that 216 of 620 (34.8%) distinct region-disease cluster  
153 associations had not previously been reported at  $p < 5 \times 10^{-8}$  (130/445 [29.2%] at  $p < 2 \times 10^{-9}$ ). Of  
154 the 216 distinct loci, 177 (104/130 at  $p < 2 \times 10^{-9}$ ) were in genes not previously mapped to the  
155 respective diseases (**Supplementary Table 3, Figure 1a, Supplementary Figure 1**  
156 **(interactive)**). Of the novel associations at GWAS significance ( $p < 5 \times 10^{-8}$ ), roughly one third  
157 had MAF<5% in either UKB or FG and 15% had MAF<1% (**Supplementary Table 3**).  
158 Importantly, 17% of known (UKB; 19% in FG), but 31% of novel (UKB; 28% in FG)  
159 associations had a MAF<5%. Correspondingly, 5% of known (UKB; 6% in FG) and 15% of  
160 novel (UKB; 10% in FG) associations had a MAF<1%, highlighting the power gained through  
161 our approach especially in the low and rare allele frequency spectrum (**Figure 1e,**  
162 **Supplementary Table 3**).

163  
164 Mapping associations to genes, we found the majority of gene loci (81.2% at  $p < 5 \times 10^{-8}$ , MHC  
165 region counted as one locus) to be associated with a single disease cluster (**Extended Data**  
166 **Figure 6a**). Thirteen loci were associated with  $\geq 5$  trait clusters (at  $p < 5 \times 10^{-8}$ ), including well

167 established pleiotropic regions such as the *MHC*, *APOE*, *PTPN22*, *GCKR*, *SH2B3* and *FUT2*  
168 (**Figure 1a**). For instance, in addition to a known association with breast cancer, we found  
169 variants in *CHEK2* as associated with the risk of colorectal and thyroid cancers, uterine  
170 leiomyoma, benign meningeal tumours and ovarian cysts. Also, in addition to a known  
171 association with prostate hyperplasia, we found an *ODF3* missense variant (rs72878024) to be  
172 associated with risk of uterine leiomyoma, benign meningeal tumour, lipoma and polyps in the  
173 female genital tract (**Supplementary Table 3**).

174  
175 Harnessing the added power of UKB and FG, we were able to detect GWAS associations for  
176 rare variants previously only annotated as causal for single-gene diseases, establishing a  
177 disease relevance for these variants at the population level. Of the 534 distinct variants with  
178 significant disease associations in our study ( $p < 5 \times 10^{-8}$ ), 152 had previously been linked to  
179 diseases in ClinVar. For 45 of these variants, the associated disease cluster matched with a  
180 previously reported phenotype in ClinVar. Notably, only six of these 45 variants (in *GJB2*,  
181 *ABCC6*, *BRCA1*, *SERPINA1*, *FLG*, and *MYOC*) had a previous annotation as either pathogenic  
182 or likely pathogenic (**Supplementary Table 7**), with 15 others annotated as benign. Of the  
183 novel trait cluster associations, 17 had been reported in ClinVar for the same/similar diseases,  
184 with 4 being classified as pathogenic/likely pathogenic and 13 classified either as benign or  
185 having “conflicting interpretation of pathogenicity” for the associated trait (**Supplementary**  
186 **Table 3, Supplementary Table 7**). For instance, we found a rare missense variant annotated  
187 as showing conflicting pathogenicity in ClinVar in *VWF* (rs1800386:C; Tyr1584Cys;  
188 MAF=0.44% [UKB], 0.47% [FG]) to be associated with the risk of von Willebrand disease<sup>12</sup>  
189 ( $\log[\text{OR}] = 2.09$ ,  $p = 8.7 \times 10^{-9}$ ); or a missense variant in *SPINK1* (rs17107315:C; Asn34Ser;  
190 MAF=1.3% [UKB], 1.6% [FG]) annotated as showing conflicting pathogenicity in ClinVar for  
191 chronic pancreatitis to be associated with chronic pancreatitis risk<sup>12</sup> ( $\log[\text{OR}] = 1.16$ ,  $p = 6.9 \times 10^{-$

192 <sup>25</sup>) and acute pancreatitis risk ( $\log[\text{OR}]=0.69$ ,  $p=2.3 \times 10^{-18}$ ). These examples highlight that  
193 population-scale analyses like ours can help refine pathogenicity assignments through  
194 contributing quantitative information on relative disease risks for variant carriers. Likewise, 17  
195 of the 23 genes with highly population-enriched sentinel variants (**Table 1**) were OMIM listed  
196 disease genes. Of these, 10 (*CHEK2*, *DBH*, *SCL24A5*, *CFI*, *FLG*, *XPA*, *F10*, *BRCA1*, *SCN5A*,  
197 *CACNAID*) showed associations with conditions identical or related to the respective  
198 Mendelian disease, unveiling a relevance of the associated variants on the population level. For  
199 instance, we found the missense variant rs77273740 in *DBH* (enriched by >50x in FG), a gene  
200 associated with orthostatic hypotension, to be associated with *reduced* risk of hypertension  
201 ( $\log[\text{OR}]=-0.19$ ,  $p=1.3 \times 10^{-23}$ ), whilst an in-frame deletion (rs1250342280) in *CACNAID*  
202 (enriched by 4.3x in UKB), a gene associated with primary aldosteronism, was associated with  
203 *increased* risk of hypertension ( $\log[\text{OR}]=0.19$ ,  $p=2.0 \times 10^{-8}$ ) (**Table 1**).

204

## 205 **Biomedical insights from coding variant associations**

206 We leveraged the coding variant associations identified in our study to generate biological  
207 insights for a range of distinct genes, pathways and diseases and in the following exemplify  
208 the broad utility of our resource through a set of selected use cases.

209

### 210 **New roles of coagulation pathway proteins in conferring pulmonary embolism risk**

211 We found known and novel associations with pulmonary embolism (PE) risk, including two  
212 rare variant associations (average MAF<1%) in genes encoding components of the coagulation  
213 cascade at the convergent common pathway (**Figure 2a**). For instance, we discovered a rare  
214 missense mutation in *F10*, enriched by ~5-fold in FG (rs61753266:A; Glu142Lys; MAF=0.33%  
215 [UKB], 1.85% [FG]), to be protective against PE ( $\log[\text{OR}]=-0.44$ ,  $p=2.9 \times 10^{-9}$ ). This variant  
216 has been associated with reduced plasma coagulation factor X ( $\beta=-1.12$ ,  $p=2.0 \times 10^{-8}$ ) and

217 factor Xa ( $\beta=-1.54, p=7.9 \times 10^{-15}$ ) levels previously<sup>13</sup>, as well as clinical factor X deficiency<sup>14</sup>.  
218 Deficiencies in coagulation factors, including factor X, are associated with increased bleeding  
219 liability and reduced thrombotic risk. In a similar fashion, we found a previously reported  
220 venous thromboembolism risk-reducing variant (rs4525:C; His865Arg; MAF=27.2% [UKB],  
221 22.3% [FG]) in *F5* that is also protective for PE ( $\log[\text{OR}]=-0.14, p=1.2 \times 10^{-15}$ ) and associated  
222 with reduced plasma F5 levels<sup>15</sup> ( $\beta=-0.25, p=6.0 \times 10^{-7}$ ). This variant acts opposite to the  
223 well-established risk promoting F5 Leiden missense mutation, which leads to increased  
224 resistance to activated protein C cleavage<sup>16</sup> and thromboembolism liability, thus unravelling  
225 that coding variants in *F5* can have opposite effects on PE risk at the population level (**Figure**  
226 **2b**). We performed Mendelian randomisation (MR) using rs4525 and rs61753266 as  
227 instruments to estimate the relative reduction in PE risk due to reduced F5 ( $\beta_{\text{MR}}=0.57,$   
228  $p=1.0 \times 10^{-15}$ ) and F10 levels (F10:  $\beta_{\text{MR}}=0.40, p=2.9 \times 10^{-9}$ ; F10a:  $\beta_{\text{MR}}=0.28, p=2.9 \times 10^{-9}$ )  
229 respectively (**Figure 2b**). MR results support the expected clinical indication of factor X  
230 inhibitors in thromboembolic diseases and the hypothesis that developing drugs inhibiting  
231 factor V will also likely be beneficial for PE. We also found a rare variant in fibrinogen (*FGB*  
232 rs2227434:T; Pro100Ser; MAF=0.13% [UKB], 0.15% [FG], **Figure 2a, Supplementary**  
233 **Table 3**) that associated with increased PE risk at nominal GWAS significance ( $\log[\text{OR}]=1.03,$   
234  $p=1.5 \times 10^{-8}$ ). Missense mutations in *FGB* have previously been linked to both elevated and  
235 reduced fibrinogen levels through GWAS<sup>17,18</sup>, as well as congenital afibrinogenemia<sup>19</sup>.

236

### 237 **Rare variant biomarker associations yield insights into disease mechanisms**

238 We interrogated the sentinel variants identified in this study for associations with 117  
239 quantitative biomarkers spanning eight categories in UKB (**Supplementary Table 8**). At a  
240 multiple testing adjusted threshold of  $p < 1 \times 10^{-6}$ , we found 108 of the biomarkers to be  
241 associated with at least one of 417 sentinel variants across 239 regions (**Figure 3a,**

242 **Supplementary Table 9).** 47 of the regions were associated with 5 or more biomarker  
243 categories (**Extended Data Figure 6b, Supplementary Table 9**), including pleiotropic  
244 disease loci such as *MHC*, *APOE*, *GCKR*, *SH2B3*, *FUT2*, *MC1R*, *ABCG5*.

245

246 Many of the newly discovered associations with biomarkers are biologically plausible. For  
247 example, a low-frequency missense variant in *ADH1B* (rs1229984:T; Arg48His; MAF=2.2%  
248 [UKB], 0.5% [FG]) that is associated with increased enzymatic activity of alcohol  
249 dehydrogenase and reduced alcohol tolerance, is also associated with reduced risk of alcohol-  
250 related disorders (alcoholic liver disease:  $\log[\text{OR}]=-1.08$ ,  $p=1.5 \times 10^{-9}$ ; mental and behavioural  
251 disorders due to alcohol:  $\log[\text{OR}]=-0.82$ ,  $p=1.2 \times 10^{-33}$ ) and increased risk of gout  
252 ( $\log[\text{OR}]=0.39$ ,  $p=3.3 \times 10^{-10}$ ). Notably, the alcohol dependence disorder-promoting *ADH1B*  
253 allele (C) is also associated with reduced IGF-1 ( $\beta=-0.11$ ,  $p=1.5 \times 10^{-51}$ ) and vitamin D levels  
254 ( $\beta=-0.049$ ,  $p=2.6 \times 10^{-10}$ ), increased levels of liver enzymes (alkaline phosphatase:  
255  $\beta=0.087$ ,  $p=3.1 \times 10^{-37}$ ; gamma-glutamyl transferase:  $\beta=0.041$ ,  $p=1.62 \times 10^{-9}$ ) and total  
256 bilirubin ( $\beta=0.031$ ,  $p=5.1 \times 10^{-7}$ ), macrocytosis with increased mean corpuscular volume  
257 ( $\beta=0.047$ ,  $p=4.2 \times 10^{-12}$ ) and mean corpuscular haemoglobin ( $\beta=0.048$ ,  $p=1.4 \times 10^{-12}$ ), as  
258 well as reduced erythrocyte count ( $\beta=-0.034$ ,  $p=3.2 \times 10^{-8}$ ). The gout risk reducing C allele  
259 is associated with reduced urate levels ( $\beta=-0.061$ ,  $p=1.4 \times 10^{-21}$ ).

260

### 261 *A deletion in SLC34A1 is associated with multiple blood and urinary abnormalities*

262 Cross-referencing with biomarkers provided mechanistic insights into novel findings. For  
263 instance, we discovered a novel association between a low frequency in-frame deletion in  
264 *SLC34A1* (rs1460573878, also known as rs876661296; MAF=2.6% [UKB], 2.7% [FG];  
265 p.Val91\_Alal97del) coding for the type II sodium phosphate cotransporter, NPT2a, which is  
266 expressed specifically in renal proximal tubular cells, to be associated with increased risk of



267 renal ( $\log[\text{OR}]=0.24$ ,  $p=4.0 \times 10^{-9}$ ) and urinary tract stones ( $\log[\text{OR}]=0.21$ ,  $p=6.8 \times 10^{-9}$ ). The  
268 deletion has previously been implicated in hypercalciuric renal stones<sup>20,21</sup> and autosomal  
269 recessive idiopathic infantile hypercalcaemia<sup>22</sup> in family studies. The variant is also associated  
270 with increased serum calcium ( $\beta=0.047$ ,  $p=5.4 \times 10^{-11}$ ) and reduced phosphate ( $\beta=-0.075$ ,  
271  $p=3.3 \times 10^{-26}$ ), consistent with a disrupted function/cell surface expression of the transporter<sup>22</sup>  
272 (**Figure 3b**). We further find associations with increased levels of serum urate ( $\beta=0.048$ ,  
273  $p=4.5 \times 10^{-17}$ ), suggesting an increased risk also of uric acid stones. Additionally, we found  
274 associations with increased erythrocyte count ( $\beta=0.035$ ,  $p=4.7 \times 10^{-10}$ ), haemoglobin  
275 concentration ( $\beta=0.033$ ,  $p=7.7 \times 10^{-10}$ ) and haematocrit percentage ( $\beta=0.036$ ,  $p=9.9 \times 10^{-11}$ ),  
276 suggesting increased renal-driven erythropoiesis (**Figure 3b**). Serum creatinine was not  
277 increased in carriers of the deletion ( $\beta=-0.07$ ,  $p=3.6 \times 10^{-33}$ ), suggesting renal function is not  
278 adversely affected in deletion carriers. Amongst 11,114 renal/ureteric and 13,319 urinary tract  
279 stone cases, we identified 735 (renal/ureteric) and 863 (urinary tract) carriers of the deletion  
280 who may benefit from clinical interventions targeting NPT2a related pathways and monitoring  
281 for deranged biochemical and haematological biomarkers.

282

### 283 *A CHEK2 deletion is associated with blood cell counts and haematological malignancies*

284 A frameshift deletion in *CHEK2* (rs555607708; MAF=0.64% [FG], 0.24% [UKB]) that  
285 increases breast cancer risk has been previously implicated also in myeloproliferative  
286 neoplasms through GWAS<sup>23</sup> and lymphoid leukaemia in a candidate variant study<sup>24</sup>.  
287 Consistently, we found nominally significant associations with risks of both, myeloid  
288 ( $\log[\text{OR}]=1.52$ ,  $p=9.5 \times 10^{-8}$ ) and lymphoid ( $\log[\text{OR}]=1.38$ ,  $p=3.1 \times 10^{-7}$ ) leukaemia, but also  
289 multiple myeloma ( $\log[\text{OR}]=1.07$ ,  $p=5.1 \times 10^{-5}$ ) and non-Hodgkin lymphoma ( $\log[\text{OR}]=0.81$ ,  
290  $p=4.7 \times 10^{-4}$ ). Association of rs555607708 with clinical haematology traits showed statistically  
291 significant associations with increased blood cell counts for both, myeloid (leukocytes,



292 neutrophils, platelets at  $p < 1 \times 10^{-6}$ ; monocyte and erythrocytes at  $p < 1 \times 10^{-3}$ ) and lymphoid  
293 (lymphocytes,  $p = 5.7 \times 10^{-17}$ ) lineages (**Figure 3c**). Furthermore, we found associations with  
294 increased mean platelet volume (MPV,  $p = 1.3 \times 10^{-16}$ ) and platelet distribution width (PDW,  
295  $p = 5.2 \times 10^{-13}$ ), consistent with increased platelet activation and previous associations of MPV  
296 and PDW with chronic myeloid leukaemia<sup>25</sup>. We also found associations with decreased mean  
297 corpuscular haemoglobin ( $p = 7.8 \times 10^{-12}$ ) and mean corpuscular volume ( $p = 5.3 \times 10^{-10}$ ),  
298 suggesting predisposition to haematological cancers by loss-of *CHEK2* function is  
299 accompanied by a microcytic red blood cell phenotype (**Figure 3c**).

300

### 301 **Coding variant associations inform drug discovery and development**

302 We cross-referenced genes with significant coding variant associations with drug targets using  
303 the therapeutic targets database<sup>26</sup>. We found 66 genes with trait cluster associations that are the  
304 targets of either approved drugs (26 genes) or drugs currently being tested in clinical trials (40  
305 genes), among these, 14 in phase 3 trials (**Supplementary Table 10**). We found a statistical  
306 enrichment for significant genes in our study to also be approved drug targets (26/482;  
307 compared with a background of 569 approved targets/19,955 genes, OR=1.9,  $p = 0.0024$ ), which  
308 is in line with previous estimates of a higher success rate for drug targets supported by  
309 genetics<sup>27,28</sup>. Sensitivity analyses using more stringent association  $p$ -value thresholds further  
310 increased these probability estimates ( $p = 5 \times 10^{-9}$  [OR 2.3,  $p = 0.00070$ ];  $p = 5 \times 10^{-10}$  [OR 2.5,  
311  $p = 0.00037$ ], supporting previous observations that the stronger the genetic association, the  
312 higher the likelihood of therapeutic success (**Supplementary Table 11**). In addition to  
313 providing further support for well-established drug target associations such as between *PCSK9*  
314 loss-of-function and hypercholesterolaemia, or *F10* loss-of-function and venous  
315 thromboembolism, we also found an association between a common missense variant  
316 (rs231775:G) in *CTLA4* with increased risk of thyrotoxicosis ( $\log[\text{OR}] = 0.12$ ,  $p = 8.5 \times 10^{-13}$ ).

317 Since this variant is also a blood eQTL for decreased *CTLA4* expression<sup>29</sup> (Z-score=-6.91,  
318  $p=5.0 \times 10^{-12}$ ), the association between genetic CTLA4 reduction and thyroid dysfunction might  
319 contribute to the adverse event of hyperthyroidism in cancer patients treated with CTLA4  
320 inhibitors<sup>30</sup>.

321

322 Genetics can inform drug discovery also on alternative indications for repurposing. For  
323 example, *TYK2* inhibitors are being tested in clinical trials for various autoimmune and  
324 psoriatic diseases<sup>31</sup>. Consistent with previous GWAS<sup>10</sup>, we found a missense variant in *TYK2*  
325 (rs34536443:C) to be associated with reduced risk of rheumatoid arthritis and psoriatic diseases  
326 (**Supplementary Table 3**). Our analyses establish this variant to also be associated with  
327 sarcoidosis ( $\log[\text{OR}]=-0.41$ ,  $p=3.6 \times 10^{-8}$ ), proposing sarcoidosis as a new indication for *TYK2*  
328 inhibitors. Similarly, while the pleiotropy of *CHEK2* provides support for exploring CHEK2  
329 inhibitors against a broader spectrum of malignancies, our analyses also highlight a risk for  
330 potential haematological perturbations upon CHEK2 inhibitor treatment.

331

### 332 **Genetic insights into atrial fibrillation**

333 Atrial fibrillation (AF) GWAS have yielded a sizeable number of loci<sup>32,33</sup>. We chose AF to  
334 exemplify how results from our study can further elucidate the genetics and biological basis of  
335 one distinct human trait, with a particular emphasis on how our results might help to  
336 disambiguate AF loci to causal genes and explain the functional significance of coding variant  
337 associations. Indeed, we report several coding variant associations (**Supplementary Table 3**)  
338 where prior GWAS<sup>32,33</sup> had fallen short to resolve GWAS loci to coding genes and explain  
339 disease mechanisms.

340

341 ***A binding motif disrupting missense variant reveals a role for METTL11B methylase in AF***

342 The AF GWAS sentinel variant rs72700114 is an intergenic variant located between  
343 *METTL11B* and *LINC01142*<sup>32-34</sup>. Our study unveiled a low frequency missense variant in  
344 *METTL11B* (rs41272485:G; Ile127Met; MAF=3.9% [UKB], 3.8% [FG]) as associated with  
345 increased AF risk (log[OR]=0.14,  $p=4.0 \times 10^{-11}$ ). *METTL11B* is a N-terminal monomethylase  
346 that methylates target proteins containing an N-terminal [Ala/Pro/Ser]-Pro-Lys motif<sup>35</sup>. The  
347 missense variant Ile127Met (SIFT=0, PolyPhen=1.0) falls within a conserved motif in the  
348 enzyme's S-adenosylmethionine/S-adenosyl-l-homocysteine ligand binding site<sup>36</sup>. *METTL11B*  
349 expression is enriched in heart and skeletal muscles with highest expression in heart muscle,  
350 in particular cardiomyocytes<sup>37,38</sup>. We scanned protein sequences for a presence of the  
351 [Ala/Pro/Ser]-Pro-Lys motif and elevated expression in cardiomyocytes (**Methods,**  
352 **Supplementary Table 12**). We found statistically significant enrichment of genes encoding  
353 [Ala/Pro/Ser]-Pro-Lys motif containing proteins amongst genes with elevated expression in  
354 cardiomyocytes (OR=1.34, 95% CI=[1.16, 1.54],  $p=3.2 \times 10^{-5}$ ), many of which show N-  
355 terminal [Ala/Pro/Ser]-Pro-Lys motifs (OR=1.24, 95% CI=[1.06, 1.44],  $p=5.6 \times 10^{-3}$ ). This  
356 group included several well-established AF genes<sup>39</sup> such as potassium channels (*KCNA5*,  
357 *KCNE4*, *KCNN3*), sodium channels (*SCN5A*, *SCN10A*), *NPPA*, and *TTN*. Our data support  
358 *METTL11B* as the causal gene in this GWAS locus and a relevance for N-terminal  
359 [Ala/Pro/Ser]-Pro-Lys methylation in cardiomyocytes for AF.

360

361 ***Rare variants unveil causal mechanisms in SCN5A-SCN10A and HCN4-REC114 AF loci***

362 Within the *SCN5A-SCN10A* locus, we replicated a common missense variant in *SCN10A*  
363 (rs6795970:A; Ala1073Val; MAF=40.0% [UKB], 44.6% [FG]) that was previously described  
364 to prolong cardiac conduction<sup>40</sup>. Additionally, we found associations with reduced AF risk  
365 (log[OR]=-0.06,  $p=2.1 \times 10^{-12}$ ), reduced pulse rate (beta=-0.02,  $p=4.8 \times 10^{-18}$ ), and a suggestive

366 signal for increased risk of atrioventricular block ( $\log[\text{OR}]=0.10$ ,  $p=1.9\times 10^{-7}$ ). On top of this,  
367 we discovered a rare, FG enriched missense variant (rs45620037:A; Thr220Ile; MAF=0.11%  
368 [UKB], 0.47% [FG]; SIFT=0.03, PolyPhen=0.96) in *SCN5A*, encoding the cardiac sodium  
369 channel  $\text{Na}_v1.5$ , as associated with decreased risk of AF ( $\log[\text{OR}]=-0.65$ ,  $p=1.3\times 10^{-12}$ ). The  
370 variant lies in the S4 voltage sensing segment of the first transmembrane domain of *SCN5A*<sup>41</sup>  
371 and leads to a substitution of the polar Thr to a non-polar Ile residue, most probably causing  
372 loss of function and electrophysiological changes<sup>42,43</sup>. The Thr220Ile variant has been  
373 associated with dilated cardiomyopathy<sup>44</sup> and conduction defects including sick sinus  
374 syndrome and atrial standstill<sup>45</sup> in family studies with bradycardic changes. Consistently, we  
375 found a nominal association with reduced pulse rate ( $\beta=-0.078$ ,  $p=0.023$ ), suggesting that  
376 protective effects of the variant will be most beneficial for the common tachycardic form of  
377 AF through reducing pulse rate. Notably, our results identify both, *SCN5A* and *SCN10A* as  
378 likely causal genes at this AF locus.

379  
380 Another AF GWAS locus is tagged by the common intergenic sentinel variant rs74022964  
381 between *HCN4* and *REC114*<sup>32,33</sup>. We identified a rare, FG enriched variant in *HCN4*  
382 (rs151004999:T; Asp364Asn; MAF=0.045% [UKB], 0.17% [FG]; SIFT=0.05, PolyPhen=0.41)  
383 as associated with increased AF risk ( $\log[\text{OR}]=0.72$ ,  $p=2.8\times 10^{-8}$ ). *HCN4* is a  
384 hyperpolarization-activated ion channel contributing to cardiac pacemaker (funny) currents ( $I_f$ ).  
385 Mutations in *HCN4* have been associated with familial bradycardia (also known as sick sinus  
386 syndrome 2) and Brugada syndrome 8 in family studies<sup>46</sup>. Consistently, in our study we also  
387 found an association with decreased heart rate ( $\beta=-0.49$ ,  $p=3.8\times 10^{-21}$ ).

388

389 ***Genetic variants underlying AF risk differentially modulate pulse rate***

390 To further evaluate the hypothesis that distinct genetic mechanisms underlying AF risk  
391 inversely modulate pulse rate, we applied clustered Mendelian randomization (MR-Clust)<sup>47</sup>  
392 with a slight modification to the mixture-model to better accommodate rare-variants.  
393 Specifically, we related expectation maximization clustering of AF associated variants with  
394 homogenous directional effects on pulse rate (**Methods**). Among the AF sentinel variants from  
395 our coding variant association analyses, we found the two above described clusters of variants  
396 in *SCN10A* (rs6795970) and *HCN4* (rs151004999), suggestive of two genetic components of  
397 AF that can increase and decrease pulse rate, respectively (**Figure 4a, left**). Identifying  
398 components of AF with diverging directionality on pulse rate is not surprising given clinically  
399 AF can be driven by both tachycardia and bradycardia through distinct mechanisms<sup>48</sup>. As  
400 sensitivity analyses, we used sentinel variants from a recent AF GWAS<sup>32</sup> where the candidate  
401 gene sets overlapped those in our study and found concordant patterns (**Figure 4a, left-middle,**  
402 **Supplementary Table 13**). We also included all sentinel variants in the AF GWAS loci<sup>32</sup> and  
403 found additional clusters with differing impact on pulse rate (**Figure 4a, right-middle**).  
404 Conversely, as expected, permuting pulse rate led to no clustering and null associations  
405 between AF and pulse (**Figure 4a, right**). Expectedly, within the AF GWAS loci<sup>32</sup> the two  
406 rare missense alleles in *HCN4* (rs151004999:T, log[OR]=0.72) and *SCN5A* (rs45620037:A,  
407 log[OR]=-0.65) identified in our study had much larger effect sizes on AF risk than the  
408 respective non-coding sentinel GWAS variants (rs74022964:T [*HCN4* locus], log[OR]=0.12;  
409 rs6790396:C [*SCN5A, SCN10A* locus], log[OR]=-0.058) (**Figure 4a**).

410

411 ***PITX2c Pro41Ser increases AF risk through a gain-of-function mechanism***

412 Lastly, we found a rare missense variant in *PITX2* as associated with increased risk of AF  
413 (log[OR]=0.38,  $p=1.1 \times 10^{-9}$ ). This variant is enriched nearly 50-fold in FG (rs143452464:A;

414 Pro41Ser; MAF=0.023% [UKB], 1.1% [FG]) and was independently identified in a French  
415 family with AF (**Supplementary Information, Supplementary Figure Pro41Ser 1**), whilst  
416 GWAS had linked intergenic variants between *PITX2* and *FAM241A* to AF risk. *PITX2* is a  
417 bicoid type homeobox transcription factor previously assumed to play a role in cardiac rhythm  
418 control<sup>49</sup>. The Pro41Ser variant lies in the N-terminal domain that is only present in the *PITX2c*  
419 isoform expressed in cardiac muscle (**Figure 4b left**). We performed reporter assays comparing  
420 the ability of Xpress-*PITX2c* constructs to transactivate a luciferase reporter plasmid  
421 containing a putative *PITX2c* binding element (**Supplementary Information**). A construct  
422 containing the Pro41Ser variant showed a ~2.4-fold higher activation of the reporter than the  
423 wild-type construct ( $p=0.006$ , **Figure 4b middle**). This effect was abrogated upon deletion of  
424 the putative *PITX2c* binding site (**Figure 4b middle**; see also **Supplementary Information,**  
425 **Supplementary Figure Pro41Ser 2** and **Supplementary Figure Pro41Ser 3**). In cultured  
426 cardiac muscle HL-1 cells, Pro41Ser increased the transcription of several presumed *PITX2c*  
427 target genes, specifically *GJA1* (*Cx40*, 1.76-fold,  $p=0.012$ ), *GJA5* (*Cx43*, 1.85-fold,  $p=0.005$ )  
428 and *KCNH2* (1.81-fold,  $p=0.009$ ), while the transcription of other selected genes with putative  
429 roles in AF was not substantially altered (**Figure 4b right, Supplementary Table 14**, see also  
430 **Supplementary Information**). Together, these results are consistent with a putative gain-of-  
431 function mechanism of Pro41Ser on *PITX2c* transactivation potential and AF risk.

432

## 433 **Discussion**

434 Here, we have conducted the to date largest association study of protein-coding genetic variants  
435 against hundreds of disease endpoints ascertained from two massive population biobanks, UK  
436 Biobank and FinnGen. We report novel disease associations, most notably in the rare and low  
437 allelic frequency spectrum, replicate and assign putative causalities to many previously  
438 reported GWAS associations, and leverage the insights gained to elucidate disease

439 mechanisms. In addition to a substantial gain in power over previous studies, our analyses  
440 benefit from replication between two population cohorts, increasing the robustness of our  
441 findings and setting the stage for future similar studies in ethnically more diverse populations.

442

443 Importantly, our study identifies both, pathogenic variants residing in monogenic disease genes  
444 to impact the risk for related complex conditions, as well as new, likely causal sentinel variants  
445 within GWAS loci in genes with known and novel biological roles in the respective GWAS  
446 trait. With this, our study is one of the first to help bridge the gap between common and rare  
447 disease genetics across a broad range of conditions and provides support for the hypothesis that  
448 the genetic architecture of many diseases is continuous<sup>1</sup>. As reflected in a recent schizophrenia  
449 study<sup>50</sup>, GWAS tend to identify association signals primarily for variants with MAF>2%, while  
450 most variants identified through exome sequencing are ultra-rare (MAF<0.01%). Of the 975  
451 associations identified in our study, 145 are driven by unique variants in the yet little  
452 interrogated rare and low allelic frequency spectrum that is hypothesized to contribute  
453 substantially to the “missing heritability” of many human diseases<sup>51</sup>.

454

455 Our approach benefits considerably from the Finnish genetic background where, consistent  
456 with previous observations<sup>6-8</sup>, certain alleles are stochastically enriched to unusually high allele  
457 frequencies, at times exceeding population frequencies in the UK Biobank by >50-fold (such  
458 as for instance the *PITX2*-Pro41Ser variant). Our theoretical and empirical results suggest the  
459 increasing utility of enriched variants for identifying associations quantitatively towards lower  
460 allelic frequencies. Notably, we identify the most prominent relative power gain in the rarest  
461 variant frequency spectrum, highlighting a role for sequencing studies and integrating  
462 additional population cohorts with enriched variants for identifying novel disease associations  
463 at scale. In our study, we identify several alleles with comparatively high effect sizes and a



464 prevalence in the population that warrants follow-up, both experimentally as well as potentially  
465 directly in clinical settings to help improve disease outcomes. For instance, our data propose  
466 that 6.5% of UKB and FG participants with kidney or urinary tract stones, conditions  
467 debilitating >15% of men and 5% of women by 70 years of age<sup>52</sup>, carry a deletion in *SLC34A1*.  
468 Monitoring patients for the clinical biomarkers identified here as associated with this deletion  
469 might help to differentiate aetiologies and guide individualized treatments. Likewise, coding  
470 variant associations identified in our study may serve as an attractive source to generate  
471 hypotheses for drug discovery programs. Our results support previous studies<sup>27,28</sup> that drug  
472 targets supported by human genetics have an increased likelihood of success, which can be  
473 considered particularly high when the genetic effect on a drug target closely mimics that of a  
474 pharmacological intervention<sup>53</sup>.

475

476 We demonstrate the broad utility of our results through numerous examples. For instance, in  
477 the case of AF we show how coding variants can help disambiguate GWAS loci to likely causal  
478 genes and in some cases predict specific changes in a protein's function; how integration of  
479 genetics with intermediate traits (such as slow versus fast heart rate) can unravel different  
480 biological mechanisms underlying a disease entity; or how a putative function of sentinel  
481 coding variants can be further validated through experiments. Our examples highlight that the  
482 step from association to biological insight may be considerably shorter for coding variant  
483 association studies than it has traditionally been for GWAS.

484

485 The results of our study foreshadow the discovery of many additional coding and non-coding  
486 associations from cross-biobank analyses at even larger sample sizes. With the continued  
487 growth of population biobanks with comprehensive health data also in non-European  
488 populations, the emergence of more and more cost-effective technologies for sequencing and



489 genotyping, and computational advances to analyse genetic and non-genetic data at scale,  
490 future studies will be able to assess the genetic contribution to health and disease at even finer  
491 resolution.

492

## 493 **Methods**

### 494 **Samples and participants**

495 UK biobank (UKB) is a UK population study of approximately 500,000 participants aged 40-  
496 69 years at recruitment<sup>2</sup>. Participant data include genomic, electronic health record linkage,  
497 blood, urine and infection biomarkers, physical and anthropometric measurements, imaging  
498 data and various other intermediate phenotypes that are constantly being updated. Further  
499 details are available at <https://biobank.ndph.ox.ac.uk/showcase/>. Analyses in this study were  
500 conducted under UK Biobank Approved Project number 26041.

501  
502 FinnGen (FG) is a public-private partnership project combining electronic health record and  
503 registry data from six regional and three Finnish biobanks. Participant data include genomics  
504 and health records linked to disease endpoints. Further details are available at  
505 <https://www.finngen.fi/>. More details on FG and ethics protocols are provided in  
506 **Supplementary Information**. We used data from FG participants with completed genetic  
507 measurements (R5 data release) and imputation (R6 data release). FinnGen participants  
508 provided informed consent for biobank research. Recruitment protocols followed the biobank  
509 protocols approved by Fimea, the National Supervisory Authority for Welfare and Health. The  
510 Coordinating Ethics Committee of the Hospital District of Helsinki and Uusimaa (HUS)  
511 approved the FinnGen study protocol Nr HUS/990/2017. The FinnGen study is approved by  
512 Finnish Institute for Health and Welfare.

513

### 514 **Disease phenotypes**

515 FG phenotypes were automatically mapped to those used in the Pan UKBB  
516 (<https://pan.ukbb.broadinstitute.org/>) project. Pan UKBB phenotypes are a combination of

517 Phecodes<sup>54</sup> and ICD10 codes. Phecodes were translated to ICD10  
518 ([https://phewascatalog.org/phecodes\\_icd10](https://phewascatalog.org/phecodes_icd10), v.2.1) and mapping was based on ICD-10  
519 definitions for FG endpoints obtained from cause of death, hospital discharge and cancer  
520 registries. For disease definition consistency, we reproduced the same Phecode maps using the  
521 same ICD-10 definitions in UKB. In particular, we expertly curated 15 neurological  
522 phenotypes using ICD10 codes. We retained phenotypes where the similarity score (Jaccard  
523 index:  $ICD10_{FG} \cap ICD10_{UKB} / ICD10_{FG} \cup ICD10_{UKB}$ ) was  $>0.7$  and additionally excluded  
524 spontaneous deliveries and abortions.

525  
526 Phecodes and ICD10 coded phenotypes were first mapped to unified disease names and disease  
527 groups using mappings from Phecode, “PheWAS” and “icd” R packages followed by manual  
528 curation of unmapped traits and diseases groups, mismatched and duplicate entries. Disease  
529 endpoints were mapped to Experimental Factor Ontology (EFO) terms using mappings from  
530 EMBL-EBI and Open Targets based on exact disease entry matches followed by manual  
531 curation of unmapped traits.

532  
533 Disease trait clusters were determined through first calculating the phenotypic similarity via  
534 the cosine similarity, then determining clusters via hierarchical clustering on the distance  
535 matrix (1-similarity) using the Ward algorithm and cutting the hierarchical tree, after inspection,  
536 at height 0.8 to provide the most semantically meaningful clusters.

537

## 538 **Genetic data processing**

### 539 **UKB genetic QC**

540 UKB genotyping and imputation were performed as described previously<sup>2</sup>. WES data for UKB  
541 participants were generated at the Regeneron Genetics Center (RGC) as part of a collaboration

542 between AbbVie, Alnylam Pharmaceuticals, AstraZeneca, Biogen, Bristol-Myers Squibb,  
543 Pfizer, Regeneron and Takeda with the UK Biobank. WES data were processed using the RGC  
544 SBP pipeline as described in <sup>3,55</sup>. RGC generated a QC-passing “Goldilocks” set of genetic  
545 variants from a total of 454,803 sequenced UK Biobank participants for analysis. Additional  
546 QC were performed prior to association analyses as detailed below.

547

#### 548 **FG genetic QC**

549 Samples were genotyped with Illumina (Illumina Inc., San Diego, CA, USA) and Affymetrix  
550 arrays (Thermo Fisher Scientific, Santa Clara, CA, USA). Genotype calls were made with  
551 GenCall and zCall algorithms for Illumina and AxiomGT1 algorithm for Affymetrix data.  
552 Sample, genotyping as well as imputation procedures and QC are detailed in **Supplementary**  
553 **Information**.

554

#### 555 **Coding variant selection**

556 GnomAD v.2.0 variant annotations were used for FinnGen variants<sup>56</sup>. The following gnomAD  
557 annotation categories are included: predicted loss-of-function (pLoF), low-confidence loss-of-  
558 function (LC), in-frame indel, missense, start lost, stop lost, stop gained. Variants have been  
559 filtered to imputation INFO score > 0.6. Additional variant annotations were performed using  
560 variant effect predictor (VEP)<sup>57</sup> with SIFT and PolyPhen scores averaged across the canonical  
561 annotations.

562

#### 563 **Disease endpoint association analyses**

564 For optimized meta-analyses with FG, analyses in UKB were performed in the subset of  
565 exome-sequence UKB participants with white European ancestry for consistency with FG  
566 (n=392,814). We used REGENIE v1.0.6.7 for association analyses via a two-step procedure as

567 detailed in<sup>58</sup>. In brief, the first step fits a whole genome regression model for individual trait  
568 predictions based on genetic data using the leave one chromosome out (LOCO) scheme. We  
569 used a set of high-quality genotyped variants: minor allele frequency (MAF)>5%, minor allele  
570 count (MAC)>100, genotyping rate >99%, Hardy-Weinberg equilibrium (HWE) test  $p>10^{-15}$ ,  
571 <5% missingness and linkage-disequilibrium (LD) pruning (1000 variant windows, 100 sliding  
572 windows and  $r^2<0.8$ ). Traits where the step 1 regression failed to converge due to case  
573 imbalances were subsequently excluded from subsequent analyses. The LOCO phenotypic  
574 predictions were used as offsets in step 2 which performs variant association analyses using  
575 the approximate Firth regression detailed in<sup>58</sup> when the  $p$ -value from the standard logistic  
576 regression score test is below 0.01. Standard errors (SEs) were computed from the effect size  
577 estimate and the likelihood ratio test  $p$ -value. To avoid issues related to severe case imbalance  
578 and extremely rare variants, we limited association test to phenotypes with >100 cases and for  
579 variants with  $MAC\geq 5$  in total samples and  $MAC\geq 3$  in cases and controls. The number of  
580 variants used for analyses varies for different diseases as a result of the MAC cut-off for  
581 different disease prevalence. The association models in both steps also included the following  
582 covariates: age, age<sup>2</sup>, sex, age\*sex, age<sup>2</sup>\*sex, first 10 genetic principle components (PCs).  
583  
584 Association analyses in FG were performed using mixed model logistic regression method  
585 SAIGE v0.39<sup>59</sup>. Age, sex, 10 PCs and genotyping batches were used as covariates. For null  
586 model computation for each endpoint each genotyping batch was included as a covariate for  
587 an endpoint if there were at least 10 cases and 10 controls in that batch to avoid convergence  
588 issues. One genotyping batch need be excluded from covariates to not have them saturated. We  
589 excluded Thermo Fisher batch 16 as it was not enriched for any particular endpoints. For  
590 calculating the genetic relationship matrix, only variants imputed with an INFO score >0.95 in  
591 all batches were used. Variants with >3% missing genotypes were excluded as well as variants

592 with MAF<1%. The remaining variants were LD pruned with a 1Mb window and  $r^2$  threshold  
593 of 0.1. This resulted in a set of 59,037 well-imputed not rare variants for GRM calculation.  
594 SAIGE options for null computation were: “LOCO=false, numMarkers=30,  
595 traceCVcutoff=0.0025, ratioCVcutoff=0.001”. Association tests were performed phenotypes  
596 with case counts >100 and for variants with minimum allele count of 3 and imputation info >0.6  
597 were used.

598

599 We additionally performed sex-specific associations for a subset of gender-specific diseases  
600 (60 female diseases and in 50 disease clusters, 14 male diseases and in 13 disease clusters) in  
601 both FG and UKB using the same approach without inclusion of sex-related covariates  
602 (**Supplementary Table 2**)

603

604 We performed fixed-effect inverse-variance meta-analysis combining summary effect sizes  
605 and standard errors for overlapping variants with matched alleles across FG and UKB using  
606 METAL<sup>60</sup>.

607

## 608 **Definition and refinement of significant regions**

609 To define significance, we used a combination of (1) multiple testing corrected threshold of  
610  $p < 2 \times 10^{-9}$ ,  $0.05 / (\sim 26.8 \times 10^6)$  [sum (mean number of variants tested per disease cluster)], to  
611 account for the fact that some traits are highly correlated disease subtypes, (2) concordant  
612 direction of effect between UKB and FG associations, and (3)  $p < 0.05$  in both UKB and FG.

613

614 We defined independent trait associations through LD-based ( $r^2=0.1$ ) clumping  $\pm 500$ Kb  
615 around the lead variants using PLINK<sup>61</sup>, excluding the HLA region (chr6:25.5-34.0Mb) which  
616 is treated as one region due to complex and extensive LD patterns. We then merged overlapping

617 independent regions ( $\pm 500\text{Kb}$ ) and further restricted each independent variant ( $r^2=0.1$ ) to the  
618 most significant sentinel variant for each unique gene. For defining region associations across  
619 traits, we merged overlapping independent regions for each individual trait.

620

### 621 **Cross reference with known genetic associations**

622 We cross-referenced the sentinel variants and their proxies ( $r^2>0.2$ ) for significant associations  
623 ( $p<5\times 10^{-8}$ ) of mapped Experimental Factor Ontology (EFO) terms and their descendants in  
624 GWAS Catalog<sup>10</sup> and PhenoScanner<sup>11</sup>. To be more conservative with reporting of novel  
625 associations, we also considered whether the most-severe associated gene in our analyses were  
626 reported in GWAS Catalog and PhenoScanner. In addition, we also queried our sentinel  
627 variants in ClinVar<sup>12</sup> to define known associations with rarer genetic diseases and further  
628 manually curated novel associations for previous genome-wide significant ( $p<5\times 10^{-8}$ )  
629 associations.

630

### 631 **Biomarker associations of lead variants**

632 For the lead sentinel variants, we performed association analyses using the two-step REGENIE  
633 approach described above with 117 biomarkers including anthropometric traits, physical  
634 measurements, clinical haematology measurements, blood and urine biomarkers available in  
635 UKB (detailed in **Supplementary Table 8**).

636

### 637 **Drug target mapping and enrichment**

638 We mapped the annotated gene for each sentinel variant to drugs using the therapeutic target  
639 database (TTD)<sup>26</sup>. We retained only drugs which have been approved or are in clinical trial  
640 stages. For enrichment analysis of approved drugs with genetic associations, we used Fisher's  
641 exact test on the proportion of significant genes targeted by approved drug against a

642 background of all approved drugs in TTD<sup>26</sup> (n=595) and 20,437 protein coding genes from  
643 Ensembl annotations<sup>62</sup>.

644

## 645 **Mendelian randomization (MR) analyses**

### 646 **F5 and F10 effect on pulmonary embolism (PE) risk**

647 The missense variants rs4525 and rs61753266 in F5 and F10 genes were taken as genetic  
648 instruments for MR analyses. To assess potential that each factor level is causally associated  
649 with PE we employed two-sample MR using summary statistics, with effect of the variants on  
650 their respective factor levels obtained from previous large scale (protein quantitative trait loci)  
651 pQTL studies<sup>13,15</sup>. Let  $\beta_{XY}$  denote the estimated causal effect of a factor level on PE risk and  $\beta_X$ ,  
652  $\beta_Y$  be the genetic association with a factor level (FV, FX or FXa) and PE risk respectively.  
653 Then, the MR ratio-estimate of  $\beta_{XY}$  is given by:

$$654 \quad \beta_{XY} = \frac{\beta_Y}{\beta_X}$$

655 where the corresponding standard error  $se(\beta_{XY})$ , computed to leading order, is:

$$656 \quad se(\beta_{XY}) = \frac{se(\beta_Y)}{|\beta_X|}$$

657

### 658 **Clustered MR**

659 To assess evidence of several distinct causal mechanisms by which atrial fibrillation (AF) may  
660 influence pulse rate (PR) we used MR-Clust<sup>47</sup>. In brief, MR-Clust is a purpose-built clustering  
661 algorithm for use in univariate MR analyses. It extends the typical MR assumption that a risk  
662 factor can influence an outcome via a single causal mechanism<sup>63</sup> to a framework that allows  
663 one or more mechanisms to be detected. When a risk-factor affects an outcome via several  
664 mechanisms, the set of two-stage ratio-estimates can be divided into clusters, such that variants  
665 within each cluster have similar ratio-estimates. As shown in<sup>47</sup>, two or more variants are



666 members of the same cluster if and only if they affect the outcome via the same distinct causal  
667 pathway. Moreover, the estimated causal effect from a cluster is proportional to the total causal  
668 effect of the mechanism on the outcome. We included variants within clusters where the  
669 probability of inclusion  $>0.7$ . We used MR-Clust algorithm allowing for singletons/outlier  
670 variants to be identified as their own “clusters” to reflect the large but biologically plausible  
671 effect sizes seen with rare and low frequency variants.

672

### 673 **Bioinformatic analyses of motif and expression for *METTL11B***

674 We searched [Ala/Pro/Ser]-Pro-Lys motif containing proteins using the “peptide search”  
675 function on UniProt<sup>64</sup>, filtering for reviewed Swiss-Prot proteins and proteins listed in Human  
676 Protein Atlas (HPA)<sup>38</sup> (n=7,656). We obtained genes with elevated expression in  
677 cardiomyocytes (n=880) from HPA based on the criteria: “cell\_type\_category\_rna:  
678 cardiomyocytes; cell type enriched, group enriched, cell type enhanced” as defined by HPA in  
679 (<https://www.proteinatlas.org/humanproteome/celltype/Muscle+cells#cardiomyocytes>  
680 [accessed 20/03/2021]) with filtering for those with valid UniProt IDs (Swiss-Prot, n=863).  
681 Enrichment test was performed using Fisher’s exact test. Additionally, we performed  
682 enrichment analyses using any Ala/Pro/Ser]-Pro-Lys motif positioned within the N-terminal  
683 half of the protein (n=4,786).

684

## 685 **Acknowledgements**

686 We thank all the participants, contributors and researchers of UK Biobank and FinnGen (and  
687 its participating biobanks) for making data available for this study. We thank the UK Biobank  
688 Exome Sequencing Consortium (AbbVie, Alnylam Pharmaceuticals, AstraZeneca, Biogen,  
689 Bristol-Myers Squibb, Pfizer, Regeneron and Takeda) for generation the whole exome  
690 sequencing data and Regeneron Genetics Centre for initial quality control of the exome  
691 sequencing data. The FinnGen project is funded by two grants from Business Finland (HUS  
692 4685/31/2016 and UH 4386/31/2016) and the following industry partners: AbbVie Inc.,  
693 AstraZeneca UK Ltd, Biogen MA Inc., Celgene Corporation, Celgene International II Sàrl,  
694 Genentech Inc., Merck Sharp & Dohme Corp, Pfizer Inc., GlaxoSmithKline Intellectual  
695 Property Development Ltd., Sanofi US Services Inc., Maze Therapeutics Inc., Janssen Biotech  
696 Inc, and Novartis AG. We thank Susanna Lemmelä for her contribution to FinnGen data  
697 curation. We further thank Yi-Qing Yang, Tim Footz, Michael Walter, Amelia Aránega,  
698 Francisco Hernández-Torres, Elodie Morel and Gilles Millat for their contributions to the  
699 functional characterisation of PITX2c. PITX2 functional work was supported in part by grants  
700 from the National Natural Science Fund of China (81070153), the Personnel Development  
701 Foundation of Shanghai, China (2010019), and the Key Program of Basic Research of  
702 Shanghai, China (10JC1414002), and by the Canadian Institutes of Health Research (grants  
703 MOP-111072 and MOP-130373 to Mohamed Chahine). Asma Mechakra was supported by a  
704 bursary of the French Ministry of Research and Technology (MRT).

705

## 706 **Author contributions**

707 Conceptualization and experimental design: B.B.S., H.R.; methodology: B.B.S., H.R., C.N.F.,  
708 C.C., M.J.D.; analysis: B.B.S., M.I.K., C.N.F., A.M., C.C., E.M., J.B.W., Biogen Biobank

709 Team; experimental work: A.M., G.C., M.C., P.C.; FinnGen protocols and analysis: M.I.K.,  
710 A.P., M.J.D., FinnGen; writing: B.B.S., H.R.; all authors critically reviewed the manuscript.

711

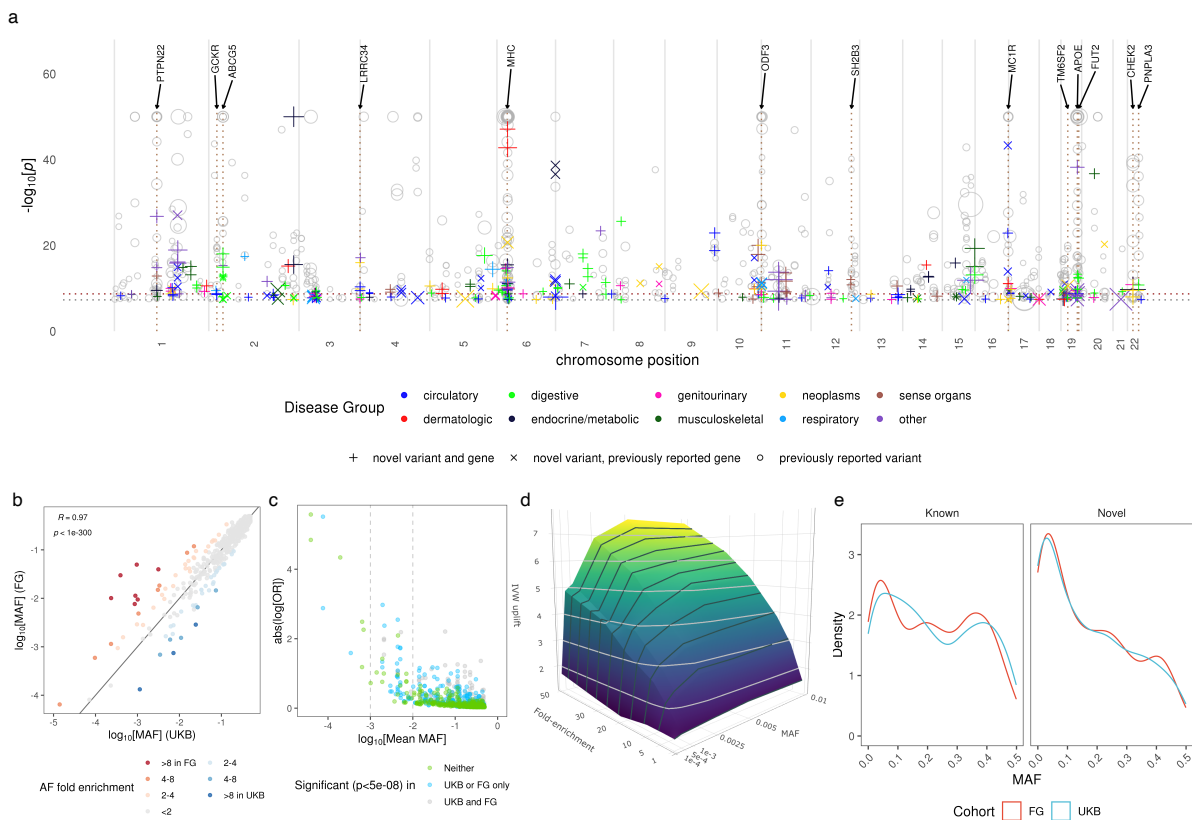
## 712 **Correspondence**

713 Correspondence and requests for materials should be addressed to B.B.S.  
714 ([bbsun92@outlook.com](mailto:bbsun92@outlook.com)) or H.R. ([heiko.runz.@gmail.com](mailto:heiko.runz.@gmail.com)).

715

## 716 Figures and Tables

717 **Figure 1. Coding genetic associations with disease. Manhattan plot for novel associations and**  
 718 **allelic enrichment surface plots are provided as Interactive Supplementary Figures 1 and 2. (a)**  
 719 **Summary of sentinel variant associations.** Size of the point is proportional to effect size.  $-\log_{10}(p)$   
 720 capped at  $-\log_{10}(10^{-50})$ . Labels highlight pleiotropic associations ( $\geq 5$  trait clusters). Colours indicate  
 721 disease groups. Shape indicate novel/known (grey circles) associations. Dotted horizontal lines:  $-\log_{10}(2 \times 10^{-9})$  [brown] and  $-\log_{10}(5 \times 10^{-8})$  [grey]. (b). **Comparison of sentinel variant MAF between**  
 723 **UKB and FG. (c) Effect size against MAF of sentinel variants.** Dashed lines indicate MAF of 0.1% (left) and 1% (right). (d) **Surface plot of effects of cohort specific allele enrichment on inverse**  
 725 **variant weighted meta-analysis z-scores (IVW uplift) across MAFs (up to MAF 1%).** Uplift is  
 726 defined as the ratio of meta-analysed IVW Z-score to the Z-score of an individual study (details in  
 727 **Supplementary Information).** (e) **Density plot of MAF for sentinel variants for known vs novel**  
 728 **associations.**  
 729

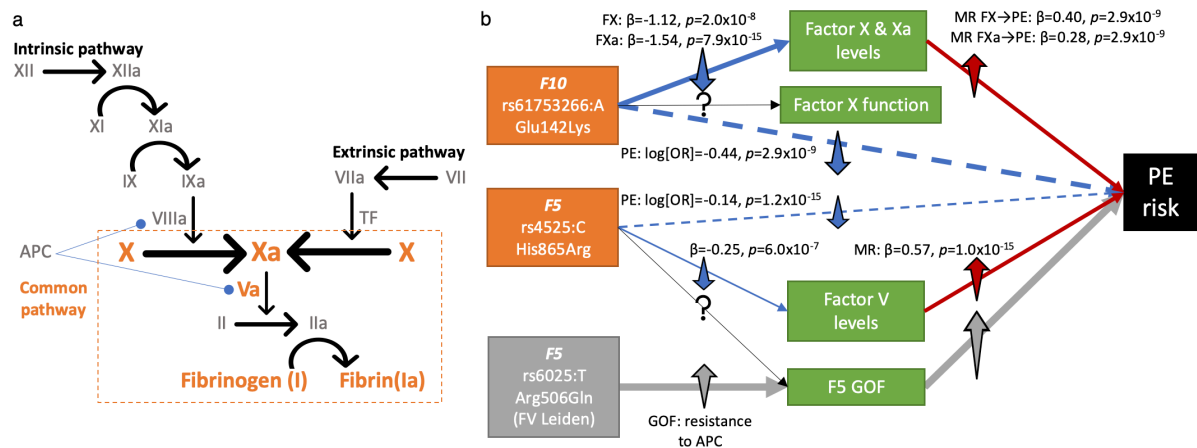


730

731

732

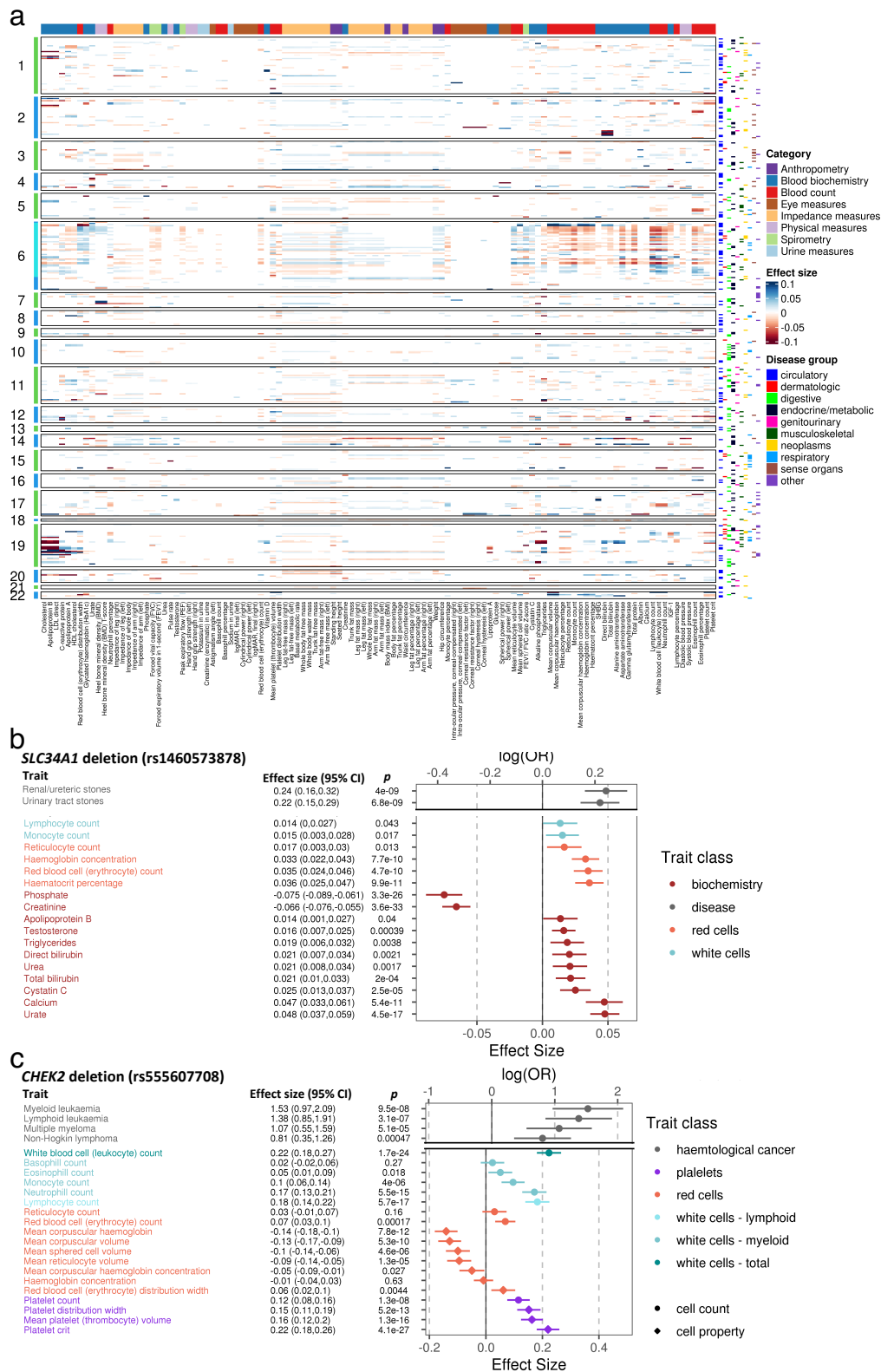
733 **Figure 2. (a) Simplified diagram of the coagulation cascade.** Factors (in roman numerals, “a”  
 734 represents activated) with genetic association with PE highlighted in orange. Blue line (round end)  
 735 indicates inhibitory effect of APC on VIIIa and Va. **(b) Schematic of potential pathway from**  
 736 **missense variants in *F5* and *F10* to PE risk.** Factor V Leiden variant had null associations with *F5*  
 737 levels ( $\beta_{F5 \text{ levels}}=0.21, p=0.091$ ). Dashed blue lines suggest effect of the variants on PE risk which we  
 738 assume under MR framework acts through factor levels (solid blue lines). Grey box and arrows  
 739 represent known pathway for Factor V Leiden mutation.  
 740 *GOF*: Gain of function, *APC*: Activated protein C, *MR*: Mendelian randomisation, *PE*: Pulmonary  
 741 embolism.  
 742



743

744

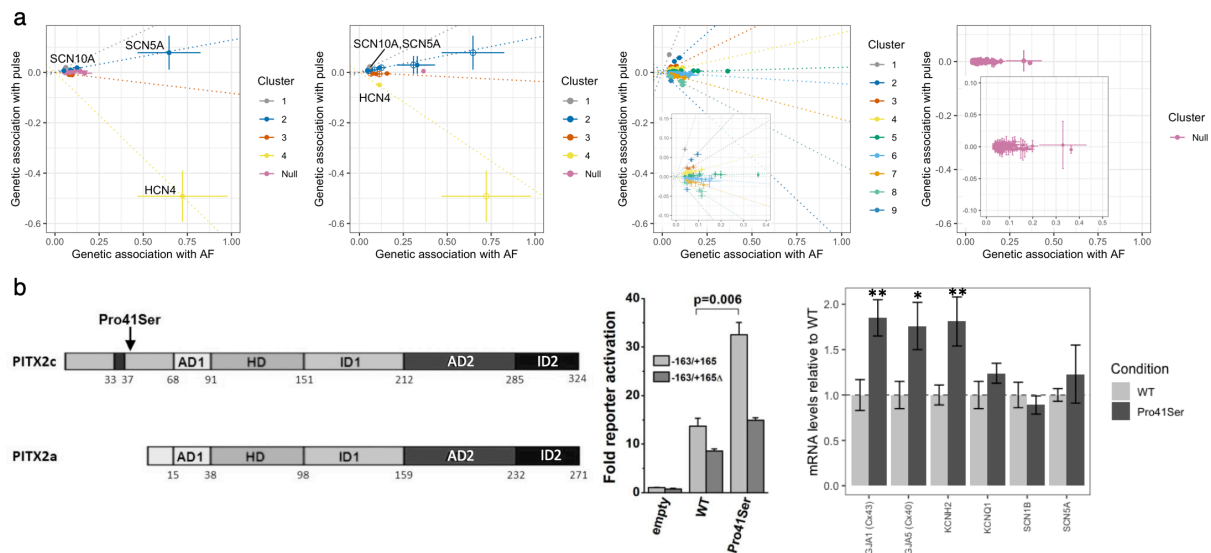
745 **Figure 3. Biomarker associations with sentinel variants.** (a) Heatmap of sentinel associations with  
 746 **biomarkers.** Only significant associations ( $p < 10^{-6}$ ) displayed. Left side indicates chromosome with  
 747 cyan indicating MHC region. Right-side: sentinel association with disease by group (colours). Top  
 748 colours: category of biomarkers. (b) Forest plot of associations between *SLC34A1* deletion  
 749 (rs1460573878) with haematological and biochemistry biomarkers.  $P < 0.05$  displayed. (c) Forest  
 750 plot of associations between *CHEK2* deletion (rs555607708) with haematological biomarkers.



752 **Figure 4. Genetic and functional insights into AF.**

753 **(a) Clustered MR plot of AF loci on pulse rate.** Only variants with cluster inclusion probability > 0.7  
 754 are included. Left to right: CWAS loci (sentinels), Overlapping CWAS and AF GWAS loci, All AF  
 755 GWAS loci from Nielsen *et al.* (with zoomed inset), All AF GWAS loci with permuted pulse (null,  
 756 with zoomed inset). **(b) Functional effect of PITX2c Pro41Ser variant (rs143452464) in vitro.** Left:  
 757 schematic of the location of the Pro41Ser variant in PITX2c as compared to the PITX2a splicing  
 758 alternative. AD1: common sequence, HD: homeodomain, ID1: transcriptional inhibitory domain 1,  
 759 AD2: second common sequence, ID2: transcriptional inhibitory domain 2. Pro41Ser is within the N-  
 760 terminal domain (grey), near to the 5-AA sequence (33 to 37 red, LAMAS) important for transcriptional  
 761 activity of the N-terminal of PITX2c. **Middle: Reporter gene assays in TM-1 cells.** Luciferase values  
 762 from activation of the SLC13A3-reporter plasmids (n=3) were normalized to  $\beta$ -galactosidase  
 763 (expressed from the transfection control plasmid), relative to the ratio for empty expression vector plus  
 764 non-deleted SLC13A3 reporter ("-163/+165"). The reporter plasmid designated as "-163/+165 $\Delta$ "  
 765 contains a deletion of 8bp corresponding to the predicted PITX2 binding site. **Right: qRT-PCR**  
 766 **analysis of HL-1 cells transfected with PITX2c recombinant plasmids.** Effect of Pro41Ser PITX2c  
 767 variant expression on Cx40, Cx43, KCNQ1, KCNH2, SCN1B and SCN5A. \* $p < 0.05$ , \*\* $p < 0.01$ .

768  
769



770  
771  
772



773 **Table 1. Genes with sentinel variants enriched >4 fold in either UKB or FG. All**  
774 **enrichment  $p < 5 \times 10^{-5}$ .**

Gene	rsID (protein change)	CHR	A0/A1	A1 freq UKB/FG%	log2 FE (FG/UKB)	OMIM gene phenotype relationships	CWAS gene phenotype relationships
<i>CHEK2</i>	rs17879961 (I200T) rs555607708 <sup>1</sup> (T410fs)	22	A/G AG/-	0.04/2.99% 0.24/0.64%	6.25 1.42	Cancer (breast, prostate, colorectal, osteosarcoma); Li-Fraumeni syndrome	rs17879961:G – benign meningeal neoplasm rs555607708:del (2.8x FG enriched) – Cancer (breast, thyroid, colorectal (benign)); uterine leiomyoma; ovarian cysts; PCOS
<i>DBH</i>	rs77273740 (R79W)	9	C/T	0.10/4.95%	5.69	Orthostatic hypotension	Hypertension (IA)
<i>PITX2</i>	rs143452464 (P41S)	4	G/A	0.02/1.01%	5.42	Anterior segment dysgenesis; Axenfeld-Rieger syndrome, ring dermoid of cornea	Arrhythmia and AF
<i>SLC24A5</i>	rs1426654 (T111A)	15	A/G	0.09/1.13%	3.67	Skin/hair/eye pigmentation (dark); oculocutaneous albinism	Non-epithelial cancer of skin (other) (IA)
<i>CFHR5</i>	rs565457964 (E163fs)	1	C/CAA	0.32/3.96%	3.66	Nephropathy due to CFHR5 deficiency	Degeneration of macula and posterior pole of retina (IA)
<i>ANKH</i>	rs146886108 (R187Q)	5	C/T	0.72/0.07%	-3.28	Chondrocalcinosis; craniometaphyseal dysplasia	Type 2 diabetes mellitus (IA)
<i>ALDH16A1</i>	rs150414818 (P527R)	19	C/G	0.10/0.95%	3.23	-	Gout
<i>LRRK1</i>	rs41531245 (T967M)	15	C/T	0.09/0.76%	3.15	-	Contracture of palmar fascia; fasciitis; umbilical hernia
<i>CFI</i>	rs141853578 (G119R)	4	C/T	0.11/0.01%	-3.10	Atypical haemolytic uremic syndrome; age-related macular degeneration; CFI deficiency	Retinal disorders (other)
<i>FLG</i>	rs61816761 (R501*) rs138381300 <sup>1</sup> (S761fs)	1	G/A CACTG/-	2.45/0.29% 2.45/1.35%	-3.10 -0.85	Atopic dermatitis; ichthyosis vulgaris	rs61816761:A – dermatitis (other) rs138381300:del (1.8x UKB enriched) – asthma; non-epithelial cancer of skin (other)
<i>SOS2</i>	rs72681869 (P191R)	14	G/C	1.09/0.15%	-2.84	Noonan syndrome	Hypertension (IA)
<i>XPA</i>	rs144725456 (H244R)	9	T/C	0.01/0.06%	2.61	Xeroderma pigmentosum	Non-epithelial cancer of skin (other)
<i>CDC25A</i>	rs146179438 (Q24H)	3	C/A	1.52/8.72%	2.52	-	Kidney and urinary stones (IA)
<i>F10</i>	rs61753266 (E142K)	13	G/A	0.33/1.83%	2.46	Factor X deficiency	PE and pulmonary heart disease (inverse association)
<i>TNXB</i>	rs61745355 (G2848R) rs10947230 <sup>1</sup> (R2704H) rs1150752 <sup>1</sup> (T302A)	6	C/T C/T T/C	2.22/11.86% 5.96/14.75% 13.29/9.17%	2.42 1.31 -0.54	Ehlers-Danlos syndrome; vesicoureteral reflux	rs61745355:T – lymphoma rs10947230:T – lichen planus rs1150752:C – chronic hepatitis; other inflammatory liver diseases; atherosclerosis
<i>SLC39A8</i>	rs13107325 (A391T)	4	C/T	7.40/1.46%	-2.35	Congenital disorder of glycosylation	Shoulder lesions
<i>CLPTM1</i>	rs150484293 (L140F)	19	C/T	0.35/0.07%	-2.33	-	Dementia
<i>ELL2</i>	rs141299831 (S18L)	5	G/A	0.02/0.12%	2.29	-	Benign neoplasm of other and ill-defined parts of digestive system
<i>CASP7</i>	rs141266925 (F214L)	10	T/C	0.31/1.5%	2.29	-	Cataracts
<i>BRCA1</i>	rs80357906 (Q1777fs)	17	T/TG	0.001/0.01%	2.21	Cancer (breast, ovarian, pancreatic); Fanconi anaemia	Breast cancer
<i>SCN5A</i>	rs45620037 (T220I)	3	G/A	0.11/0.49%	2.20	Sudden infant death syndrome; dilated cardiomyopathy; arrhythmia <sup>2</sup>	Arrhythmia and AF
<i>CACNA1D</i>	rs1250342280 (F1943del)	3	CCTT/C	0.60/0.14%	-2.09	Primary aldosteronism, seizures, and neurologic abnormalities; sinoatrial node dysfunction and deafness	Hypertension
<i>WNT10A</i>	rs121908120 (F228I)	2	T/A	2.72/0.65%	-2.06	Odontoonychodermal dysplasia; Schopf-Schulz-Passarge syndrome; selective tooth agenesis	Follicular cysts of skin and subcutaneous tissue (IA)

775  
776 <sup>1</sup>Other sentinel variants in the gene with <4 FE

777 <sup>2</sup>Sudden infant death syndrome; atrial fibrillation; Brugada syndrome; progressive and non-progressive heart block; long QT  
778 syndrome, sick sinus syndrome; ventricular fibrillation

779 FE: fold enrichment; IA: inverse association; PCOS: polycystic ovarian syndrome; PE: pulmonary embolism; AF: atrial  
780 fibrillation

781



## 782 References

- 783 1 Claussnitzer, M. *et al.* A brief history of human disease genetics. *Nature* **577**, 179-  
784 189, doi:10.1038/s41586-019-1879-7 (2020).
- 785 2 Bycroft, C. *et al.* The UK Biobank resource with deep phenotyping and genomic data.  
786 *Nature* **562**, 203-209, doi:10.1038/s41586-018-0579-z (2018).
- 787 3 Van Hout, C. V. *et al.* Exome sequencing and characterization of 49,960 individuals  
788 in the UK Biobank. *Nature* **586**, 749-756, doi:10.1038/s41586-020-2853-0 (2020).
- 789 4 Szustakowski, J. D. *et al.* Advancing Human Genetics Research and Drug Discovery  
790 through Exome Sequencing of the UK Biobank. *medRxiv*, 2020.2011.2002.20222232,  
791 doi:10.1101/2020.11.02.20222232 (2020).
- 792 5 Wang, Q. *et al.* Surveying the contribution of rare variants to the genetic architecture  
793 of human disease through exome sequencing of 177,882 UK Biobank participants.  
794 *bioRxiv*, 2020.2012.2013.422582, doi:10.1101/2020.12.13.422582 (2020).
- 795 6 Peltonen, L., Jalanko, A. & Varilo, T. Molecular genetics of the Finnish disease  
796 heritage. *Hum Mol Genet* **8**, 1913-1923, doi:10.1093/hmg/8.10.1913 (1999).
- 797 7 Lim, E. T. *et al.* Distribution and medical impact of loss-of-function variants in the  
798 Finnish founder population. *PLoS Genet* **10**, e1004494,  
799 doi:10.1371/journal.pgen.1004494 (2014).
- 800 8 Locke, A. E. *et al.* Exome sequencing of Finnish isolates enhances rare-variant  
801 association power. *Nature* **572**, 323-328, doi:10.1038/s41586-019-1457-z (2019).
- 802 9 Hassan, S. *et al.* High-resolution population-specific recombination rates and their  
803 effect on phasing and genotype imputation. *Eur J Hum Genet*, doi:10.1038/s41431-  
804 020-00768-8 (2020).
- 805 10 Buniello, A. *et al.* The NHGRI-EBI GWAS Catalog of published genome-wide  
806 association studies, targeted arrays and summary statistics 2019. *Nucleic acids*  
807 *research* **47**, D1005-D1012, doi:10.1093/nar/gky1120 (2019).
- 808 11 Staley, J. R. *et al.* PhenoScanner: a database of human genotype-phenotype  
809 associations. *Bioinformatics (Oxford, England)* **32**, 3207-3209,  
810 doi:10.1093/bioinformatics/btw373 (2016).
- 811 12 Landrum, M. J. *et al.* ClinVar: improving access to variant interpretations and  
812 supporting evidence. *Nucleic acids research* **46**, D1062-D1067,  
813 doi:10.1093/nar/gkx1153 (2018).
- 814 13 Sun, B. B. *et al.* Genomic atlas of the human plasma proteome. *Nature* **558**, 73-79,  
815 doi:10.1038/s41586-018-0175-2 (2018).
- 816 14 Forberg, E., Huhmann, I., Jimenez-Boj, E. & Watzke, H. H. The impact of  
817 Glu102Lys on the factor X function in a patient with a doubly homozygous factor X  
818 deficiency (Gla14Lys and Glu102Lys). *Thromb Haemost* **83**, 234-238 (2000).
- 819 15 Suhre, K. *et al.* Connecting genetic risk to disease end points through the human  
820 blood plasma proteome. *Nat Commun* **8**, 14357-14357, doi:10.1038/ncomms14357  
821 (2017).
- 822 16 Kujovich, J. L. Factor V Leiden thrombophilia. *Genetics in Medicine* **13**, 1-16,  
823 doi:10.1097/GIM.0b013e3181faa0f2 (2011).
- 824 17 de Vries, P. S. *et al.* A meta-analysis of 120 246 individuals identifies 18 new loci for  
825 fibrinogen concentration. *Hum Mol Genet* **25**, 358-370, doi:10.1093/hmg/ddv454  
826 (2016).
- 827 18 Wassel, C. L. *et al.* Association of genomic loci from a cardiovascular gene SNP  
828 array with fibrinogen levels in European Americans and African-Americans from six

- 829 cohort studies: the Candidate Gene Association Resource (CARE). *Blood* **117**, 268-  
830 275, doi:10.1182/blood-2010-06-289546 (2011).
- 831 19 Simurda, T. *et al.* Genetic Variants in the FGB and FGG Genes Mapping in the Beta  
832 and Gamma Nodules of the Fibrinogen Molecule in Congenital Quantitative  
833 Fibrinogen Disorders Associated with a Thrombotic Phenotype. *Int J Mol Sci* **21**,  
834 doi:10.3390/ijms21134616 (2020).
- 835 20 Lapointe, J. Y. *et al.* NPT2a gene variation in calcium nephrolithiasis with renal  
836 phosphate leak. *Kidney Int* **69**, 2261-2267, doi:10.1038/sj.ki.5000437 (2006).
- 837 21 Halbritter, J. *et al.* Fourteen monogenic genes account for 15% of  
838 nephrolithiasis/nephrocalcinosis. *J Am Soc Nephrol* **26**, 543-551,  
839 doi:10.1681/ASN.2014040388 (2015).
- 840 22 Schlingmann, K. P. *et al.* Autosomal-Recessive Mutations in SLC34A1 Encoding  
841 Sodium-Phosphate Cotransporter 2A Cause Idiopathic Infantile Hypercalcemia. *J Am*  
842 *Soc Nephrol* **27**, 604-614, doi:10.1681/ASN.2014101025 (2016).
- 843 23 Hinds, D. A. *et al.* Germ line variants predispose to both JAK2 V617F clonal  
844 hematopoiesis and myeloproliferative neoplasms. *Blood* **128**, 1121-1128,  
845 doi:10.1182/blood-2015-06-652941 (2016).
- 846 24 Sellick, G. S., Sullivan, K., Catovsky, D. & Houlston, R. S. CHEK2\*1100delC and  
847 risk of chronic lymphocytic leukemia. *Leuk Lymphoma* **47**, 2659-2660,  
848 doi:10.1080/10428190600942462 (2006).
- 849 25 Yan, K. *et al.* Normal platelet counts mask abnormal thrombopoiesis in patients with  
850 chronic myeloid leukemia. *Oncol Lett* **10**, 2390-2394, doi:10.3892/ol.2015.3502  
851 (2015).
- 852 26 Wang, Y. *et al.* Therapeutic target database 2020: enriched resource for facilitating  
853 research and early development of targeted therapeutics. *Nucleic Acids Res* **48**,  
854 D1031-D1041, doi:10.1093/nar/gkz981 (2020).
- 855 27 Nelson, M. R. *et al.* The support of human genetic evidence for approved drug  
856 indications. *Nat Genet* **47**, 856-860, doi:10.1038/ng.3314 (2015).
- 857 28 King, E. A., Davis, J. W. & Degner, J. F. Are drug targets with genetic support twice  
858 as likely to be approved? Revised estimates of the impact of genetic support for drug  
859 mechanisms on the probability of drug approval. *PLoS Genet* **15**, e1008489,  
860 doi:10.1371/journal.pgen.1008489 (2019).
- 861 29 Vösa, U. *et al.* Unraveling the polygenic architecture of complex traits using blood  
862 eQTL metaanalysis. *bioRxiv*, 447367, doi:10.1101/447367 (2018).
- 863 30 Girotra, M. *et al.* The Current Understanding of the Endocrine Effects From Immune  
864 Checkpoint Inhibitors and Recommendations for Management. *JNCI Cancer Spectr*  
865 **2**, pky021, doi:10.1093/jncics/pky021 (2018).
- 866 31 Diogo, D. *et al.* Phenome-wide association studies across large population cohorts  
867 support drug target validation. *Nat Commun* **9**, 4285, doi:10.1038/s41467-018-06540-  
868 3 (2018).
- 869 32 Nielsen, J. B. *et al.* Biobank-driven genomic discovery yields new insight into atrial  
870 fibrillation biology. *Nature Genetics* **50**, 1234-1239, doi:10.1038/s41588-018-0171-3  
871 (2018).
- 872 33 Roselli, C. *et al.* Multi-ethnic genome-wide association study for atrial fibrillation.  
873 *Nature genetics* **50**, 1225-1233, doi:10.1038/s41588-018-0133-9 (2018).
- 874 34 Thorolfsdottir, R. B. *et al.* A Missense Variant in PLEC Increases Risk of Atrial  
875 Fibrillation. *J Am Coll Cardiol* **70**, 2157-2168, doi:10.1016/j.jacc.2017.09.005 (2017).
- 876 35 Petkowski, J. J. *et al.* NRMT2 is an N-terminal monomethylase that primes for its  
877 homologue NRMT1. *Biochem J* **456**, 453-462, doi:10.1042/BJ20131163 (2013).

- 878 36 Dong, C. *et al.* An asparagine/glycine switch governs product specificity of human N-  
879 terminal methyltransferase NTMT2. *Commun Biol* **1**, 183, doi:10.1038/s42003-018-  
880 0196-2 (2018).
- 881 37 Consortium, G. T. The Genotype-Tissue Expression (GTEx) project. *Nat Genet* **45**,  
882 580-585, doi:10.1038/ng.2653 (2013).
- 883 38 Uhlen, M. *et al.* Proteomics. Tissue-based map of the human proteome. *Science* **347**,  
884 1260419, doi:10.1126/science.1260419 (2015).
- 885 39 Feghaly, J., Zakka, P., London, B., MacRae, C. A. & Refaat, M. M. Genetics of Atrial  
886 Fibrillation. *J Am Heart Assoc* **7**, e009884, doi:10.1161/JAHA.118.009884 (2018).
- 887 40 Chambers, J. C. *et al.* Genetic variation in SCN10A influences cardiac conduction.  
888 *Nat Genet* **42**, 149-152, doi:10.1038/ng.516 (2010).
- 889 41 Li, W. *et al.* SCN5A Variants: Association With Cardiac Disorders. *Front Physiol* **9**,  
890 1372, doi:10.3389/fphys.2018.01372 (2018).
- 891 42 Gui, J. *et al.* Multiple loss-of-function mechanisms contribute to SCN5A-related  
892 familial sick sinus syndrome. *PLoS One* **5**, e10985,  
893 doi:10.1371/journal.pone.0010985 (2010).
- 894 43 Benson, D. W. *et al.* Congenital sick sinus syndrome caused by recessive mutations in  
895 the cardiac sodium channel gene (SCN5A). *J Clin Invest* **112**, 1019-1028,  
896 doi:10.1172/JCI18062 (2003).
- 897 44 Olson, T. M. *et al.* Sodium channel mutations and susceptibility to heart failure and  
898 atrial fibrillation. *JAMA* **293**, 447-454, doi:10.1001/jama.293.4.447 (2005).
- 899 45 Zaklyazminskaya, E. & Dzemeshevich, S. The role of mutations in the SCN5A gene  
900 in cardiomyopathies. *Biochim Biophys Acta* **1863**, 1799-1805,  
901 doi:10.1016/j.bbamer.2016.02.014 (2016).
- 902 46 Verkerk, A. O. & Wilders, R. Pacemaker activity of the human sinoatrial node:  
903 effects of HCN4 mutations on the hyperpolarization-activated current. *Europace* **16**,  
904 384-395, doi:10.1093/europace/eut348 (2014).
- 905 47 Foley, C. N., Mason, A. M., Kirk, P. D. W. & Burgess, S. MR-Clust: Clustering of  
906 genetic variants in Mendelian randomization with similar causal estimates.  
907 *Bioinformatics*, doi:10.1093/bioinformatics/btaa778 (2020).
- 908 48 Sidhu, S. & Marine, J. E. Evaluating and managing bradycardia. *Trends Cardiovasc*  
909 *Med* **30**, 265-272, doi:10.1016/j.tcm.2019.07.001 (2020).
- 910 49 Syeda, F., Kirchhof, P. & Fabritz, L. PITX2-dependent gene regulation in atrial  
911 fibrillation and rhythm control. *J Physiol* **595**, 4019-4026, doi:10.1113/JP273123  
912 (2017).
- 913 50 Ripke, S., Walters, J. T. & O'Donovan, M. C. Mapping genomic loci prioritises genes  
914 and implicates synaptic biology in schizophrenia. *medRxiv*,  
915 2020.2009.2012.20192922, doi:10.1101/2020.09.12.20192922 (2020).
- 916 51 Eichler, E. E. *et al.* Missing heritability and strategies for finding the underlying  
917 causes of complex disease. *Nat Rev Genet* **11**, 446-450, doi:10.1038/nrg2809 (2010).
- 918 52 Howles, S. A. & Thakker, R. V. Genetics of kidney stone disease. *Nat Rev Urol* **17**,  
919 407-421, doi:10.1038/s41585-020-0332-x (2020).
- 920 53 Plenge, R. M., Scolnick, E. M. & Altshuler, D. Validating therapeutic targets through  
921 human genetics. *Nat Rev Drug Discov* **12**, 581-594, doi:10.1038/nrd4051 (2013).
- 922 54 Denny, J. C. *et al.* Systematic comparison of phenome-wide association study of  
923 electronic medical record data and genome-wide association study data. *Nat*  
924 *Biotechnol* **31**, 1102-1110, doi:10.1038/nbt.2749 (2013).
- 925 55 Kosmicki, J. A. *et al.* A catalog of associations between rare coding variants and  
926 COVID-19 outcomes. *medRxiv*, 2020.2010.2028.20221804,  
927 doi:10.1101/2020.10.28.20221804 (2021).

- 928 56 Karczewski, K. J. *et al.* The mutational constraint spectrum quantified from variation  
929 in 141,456 humans. *Nature* **581**, 434-443, doi:10.1038/s41586-020-2308-7 (2020).
- 930 57 McLaren, W. *et al.* The Ensembl Variant Effect Predictor. *Genome Biol* **17**, 122,  
931 doi:10.1186/s13059-016-0974-4 (2016).
- 932 58 Mbatchou, J. *et al.* Computationally efficient whole genome regression for  
933 quantitative and binary traits. *bioRxiv*, 2020.2006.2019.162354,  
934 doi:10.1101/2020.06.19.162354 (2020).
- 935 59 Zhou, W. *et al.* Efficiently controlling for case-control imbalance and sample  
936 relatedness in large-scale genetic association studies. *Nat Genet* **50**, 1335-1341,  
937 doi:10.1038/s41588-018-0184-y (2018).
- 938 60 Willer, C. J., Li, Y. & Abecasis, G. R. METAL: fast and efficient meta-analysis of  
939 genomewide association scans. *Bioinformatics* **26**, 2190-2191,  
940 doi:10.1093/bioinformatics/btq340 (2010).
- 941 61 Chang, C. C. *et al.* Second-generation PLINK: rising to the challenge of larger and  
942 richer datasets. *Gigascience* **4**, 7, doi:10.1186/s13742-015-0047-8 (2015).
- 943 62 Yates, A. D. *et al.* Ensembl 2020. *Nucleic Acids Res* **48**, D682-D688,  
944 doi:10.1093/nar/gkz966 (2020).
- 945 63 Burgess, S., Foley, C. N. & Zuber, V. Inferring Causal Relationships Between Risk  
946 Factors and Outcomes from Genome-Wide Association Study Data. *Annu Rev*  
947 *Genomics Hum Genet* **19**, 303-327, doi:10.1146/annurev-genom-083117-021731  
948 (2018).
- 949 64 UniProt, C. UniProt: the universal protein knowledgebase in 2021. *Nucleic Acids Res*  
950 **49**, D480-D489, doi:10.1093/nar/gkaa1100 (2021).  
951