



THE UNIVERSITY *of* EDINBURGH

Edinburgh Research Explorer

An atlas of genetic associations in UK Biobank

Citation for published version:

Canela-Xandri, O, Rawlik, K & Tenesa, A 2018, 'An atlas of genetic associations in UK Biobank', *Nature Genetics*, vol. 50, no. 11, pp. 1593-1599. <https://doi.org/10.1038/s41588-018-0248-z>

Digital Object Identifier (DOI):

[10.1038/s41588-018-0248-z](https://doi.org/10.1038/s41588-018-0248-z)

Link:

[Link to publication record in Edinburgh Research Explorer](#)

Document Version:

Peer reviewed version

Published In:

Nature Genetics

General rights

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact openaccess@ed.ac.uk providing details, and we will remove access to the work immediately and investigate your claim.



1 **Title:**

2 An atlas of genetic associations in UK Biobank

3

4 **Authors:**

5 Oriol Canela-Xandri^{1,2,*}, Konrad Rawlik^{1,*}, Albert Tenesa^{1,2,3,*}

6

7

8 **Affiliations:**

9 ¹The Roslin Institute, Royal (Dick) School of Veterinary Studies, The University of
10 Edinburgh, Easter Bush Campus, Midlothian, UK.

11 ² MRC Human Genetics Unit at the MRC Institute of Genetics and Molecular
12 Medicine, University of Edinburgh, Western General Hospital, Edinburgh, UK

13

14 ***contributed equally**

15

16 ³ Corresponding author

17 Dr Albert Tenesa

18 The Roslin Institute

19 The University of Edinburgh

20 Easter Bush

21 Roslin, Midlothian

22 EH25 9RG

23 UK

24 Tel: 0044 (0)131 651 9100

25 Fax: 0044 (0)131 651 9220

26 Email: Albert.Tenesa@ed.ac.uk

27

28

29

30 **ABSTRACT**

31 **Genome-wide association studies have revealed many loci contributing to**
32 **the variation of complex traits, yet the majority of loci that contribute to the**
33 **heritability of complex traits remain elusive. Large study populations with**
34 **sufficient statistical power are required to detect the small effect sizes of**
35 **the yet unidentified genetic variants. However, the analysis of huge**
36 **cohorts, like UK Biobank, is challenging. Here we present an atlas of**
37 **genetic associations for 118 non-binary and 660 binary traits of 452,264 UK**
38 **Biobank participants of white descent. Results are compiled in a publicly**
39 **accessible database that allows querying genome-wide association results**
40 **for 9,113,133 genetic variants, as well as downloading whole GWAS**
41 **summary statistics for over 30 million imputed genetic variants (>23 billion**
42 **phenotype-genotype pairs). Our atlas of associations (GeneATLAS,**
43 **<http://geneatlas.roslin.ed.ac.uk>) will help researchers to query UK Biobank**
44 **results in an easy and uniform way without the need to incur in high**
45 **computational costs.**

46

47

48 **INTRODUCTION**

49 Most human traits are complex and influenced by the combined effect of large
50 numbers of small genetic and environmental effects¹. Genome-wide association
51 studies (GWAS) have identified many genetic variants influencing many complex
52 traits. The largest genetic effects were discovered with modest sample sizes, with
53 researchers subsequently joining efforts to increase the size of the study cohorts,
54 thus allowing them to identify much smaller genetic effects. The UK Biobank², a
55 large prospective epidemiological study comprising approximately 500,000
56 deeply phenotyped individuals from the United Kingdom, has been genotyped
57 using an array that comprises 847,441 genetic polymorphisms, with a view to
58 identifying new genetic variants in a uniformly genotyped and phenotyped cohort
59 of unprecedented size, both in terms of the number of samples and number of
60 traits.

61 The unprecedented size of this cohort has raised a number of analytical
62 challenges³. First, storing, managing and analysing the circa 90 million genetic
63 variants for around half a million individuals is, in itself, a substantial endeavour.
64 Second, the collection of samples at this scale has brought up an analytical
65 challenge, as the cohort is structured by familial relationships and ethnicity. For
66 instance, many relatives were unintentionally collected in the cohort, and
67 removing them from the analyses as traditionally done in GWAS would entail a
68 substantial loss of statistical power. Third, although recent developments have
69 reduced the computational costs⁴, fitting a Linear Mixed Model (LMM), the
70 standard analytical technique to perform GWAS when there is population or
71 familial structure, at this scale and for this number of traits, entails a
72 computational burden which may be beyond the means of many research labs.

73 The objective of the current study was to perform GWAS for 778 traits in UK
74 Biobank, adjusting for the effect of relatedness to minimise the loss of statistical
75 power whilst reducing false positives due to familial and population structure, in
76 individuals of white ancestry and to make a searchable atlas of genetic
77 associations in UK Biobank for the benefit of the research community.

78 **RESULTS**

79 **Data overview**

80 In July 2017, the UK Biobank released genotyped data from circa 490,000
81 individuals of largely white descent genotyped for 805,426 genetic variants. We
82 performed GWASs for 660 binary traits and 118 non-binary traits, the latter
83 including continuous traits and traits with multiple ordered categories
84 (Supplementary Table 1). For each of these traits we fitted LMMs to test for
85 association with 623,944 genotyped and 30,798,054 imputed genetic
86 polymorphisms imputed using the Haplotype Reference Consortium⁵ as
87 reference panel, as well as 310 imputed HLA alleles. All successfully tested
88 polymorphisms are shown in the database (GeneATLAS,
89 <http://geneatlas.roslin.ed.ac.uk>) or associated downloadable files to allow
90 individual researchers to apply their own quality control thresholds. The summary
91 results presented here are based on the quality controlled imputed
92 polymorphisms (9,113,133 variants after filtering) of 452,264 individuals
93 (Methods).

94

95 The phenotypes selected comprise a mix of baseline measurements (e.g. height),
96 self-reported traits at recruitment (e.g. self-reported depression), and Hospital
97 Episode Statistics (i.e. data collected during hospital admissions) as well as
98 cancer diagnoses from the appropriate UK Cancer Registry. Since UK Biobank
99 is a recently established prospective cohort, we allowed for potential differences in
100 statistical power among binary and non-binary traits by splitting the presentation
101 of the data into non-binary and binary traits.

102

103 To demonstrate the power of using large datasets (so called, Big Data), we first
104 explored how the analysis of increasingly large sample sizes enable new
105 discoveries, and reduce bias when estimating the effect sizes of GWAS hits (Fig.
106 1 and Supplementary Note). Our results show that the number of GWAS hits
107 increased linearly with the sample size with no sign of saturation, thus suggesting
108 that increasing the size of cohorts like UK Biobank would continue to yield new
109 discoveries. We also observed that the estimated allelic effects of GWAS hits
110 obtained from decreasing sample sizes were generally larger, which is in
111 agreement with a Winner's Curse effect⁶ (Fig. 1).

112

113 **Distribution of GWAS hits among non-binary trait**

114 Just below 5 million of the circa 1 billion tests performed across 118 non-binary
115 traits were significant at a conventional genome wide threshold ($P < 10^{-8}$)
116 (Supplementary Table 2), and 3,117,904 were significant after Bonferroni
117 correction ($P < 0.05/9,113,133 \times 118$). The significant associations were
118 distributed across 74,471 leading polymorphisms mapping to 38,651
119 independent loci (Methods, Fig. 2, Supplementary Table 3). A substantial
120 proportion of these associations (13.0%) were within the HLA region
121 (Supplementary Table 2).

122

123 About 9.5% of the tested polymorphisms reached genome-wide significant
124 thresholds ($P < 10^{-8}$) for at least one of the 118 tested traits, whilst 82% of the
125 tested polymorphisms were associated with at least one of these 118 traits at a
126 significance level of 10^{-2} (Supplementary Table 4). There were 20,393 genetic
127 variants each associated with more than 30 of the tested non-binary traits (Figs.
128 2 and 3, Supplementary Fig. 1). A cluster of nine variants in a 9kb region including
129 the genotyped intronic variant rs1421085 within the *FTO* gene had the largest
130 number of genome-wide significant associations outside the HLA region, all nine
131 variants being found to be associated with 58 traits (Fig. 3 and Supplementary
132 Fig. 1). The genotyped variant rs1421085 at the *FTO* locus also had the largest
133 average significance across non-binary traits ($P < 10^{-74}$) (Supplementary Fig. 2),
134 which was largely contributed by the associations to anthropometric traits such
135 as BMI and Weight which showed some of the strongest associations ($P < 10^{-300}$).
136 The HLA region contained 362 genetic variants which were significantly ($P < 10^{-8}$)
137 associated with 50 or more of the non-binary traits compared to only 128 such
138 variants in the remaining autosomal variants. About 36% of the analyzed imputed
139 HLA alleles were significant ($P < 10^{-8}$) for at least one trait (Supplementary Fig. 3).
140 Six traits ('Standing height', 'Sitting height', 'Platelet count', 'Mean platelet
141 (thrombocyte) volume', 'Trunk predicted mass', 'Trunk fat-free mass') had over
142 100,000 significant associations ($P < 10^{-8}$) each distributed across 25,352 different
143 independent lead genetic variants (Methods). Over 94% of the non-binary traits
144 had more than 100 genome-wide significant hits distributed in 74,442 different
145 leading genetic variants.

146

147 Considering the criteria for inclusion of genetic polymorphisms on the genotyping
148 array (Supplementary Table 5), the HLA polymorphisms were the most enriched
149 for associations with at least one non-binary trait (88% had a $P < 10^{-8}$), followed
150 by the Cardiometabolic, Autoimmune/Inflammatory and ApoE criteria, whilst the
151 lowest enrichment was for two low frequency variants categories (“Genome-wide
152 coverage for low frequency variants” and “Rare, possibly disease causing,
153 mutations”). Less than 8 in 100 of these polymorphisms were associated with any
154 non-binary trait (Supplementary Table 5).

155

156 We found a significant correlation ($r=0.93$, $P < 10^{-51}$) between the number of hits
157 and the SNP heritability of the traits, suggesting that the number of loci affecting
158 a trait might be proportional to the heritability of the trait (Fig. 4, Supplementary
159 Fig. 4). Consistent with this model and variation in the distribution of linkage
160 disequilibrium across the genome, the correlation of the SNP heritability with the
161 number of identified independent lead variants was similarly high ($r=0.88$, $P < 10^{-$
162 $38}$). The number of hits ($P < 10^{-8}$) per chromosome was highly correlated ($r=0.86$)
163 with the length of the chromosome covered by the genotyped SNPs
164 (Supplementary Fig. 5, Supplementary Table 6). Although this correlation could
165 arise under a polygenic model where the length of the chromosome is correlated
166 with the number of possible variants affecting the traits, the simplest explanation
167 is that it arises as a consequence of the correlation of chromosomal length and
168 number of tested variants per chromosome. Comparing the fit of two nested
169 models to explain the number of hits per chromosome as a function of number of
170 tested genetic variants and length of the chromosome or just the number of
171 genetic variants was consistent with the number of GWAS hits per chromosome
172 correlating with the length of the chromosome rather than the number of tested
173 variants (Methods).

174

175 Standing height was the trait with the largest number of hits (Fig. 5) with 261,908
176 significantly associated variants distributed across 10,374 independent lead
177 variants. We estimated that the leading polymorphisms across the 118 traits
178 studied are distributed among 38,651 independent loci, therefore 27% of these
179 independent loci contribute to the variation of height, as expected by a highly
180 polygenic trait⁷. We also computed the proportion of tested genetic variants

181 associated with at least one disease ($P < 10^{-8}$) that are also associated with height
182 and BMI at different thresholds (Supplementary Table 7). At a threshold of 10^{-8} ,
183 ~28% and ~7% of the genetic variants associated for height and BMI,
184 respectively, were also associated with at least one disease. This is important for
185 the interpretation of Mendelian Randomisation studies as it is likely that one of
186 the critical assumptions to demonstrate causality, that is, that there is no
187 pleiotropy between the exposure and the outcome, may be broken for many
188 exposure-outcome pairs.

189

190 **Distribution of GWAS hits among binary traits**

191 The binary trait with the largest number of cases was self-reported hypertension,
192 with an average across binary traits of 6,593 cases (Supplementary Table 1). Of
193 the 660 binary phenotypes 86 were specific to one sex (Supplementary Table 1).
194 Individuals of the unaffected sex were excluded from the analysis for these
195 phenotypes (Methods). Consistent with the reduced statistical power to detect
196 association with binary phenotypes (mainly diseases) compared to non-binary
197 traits we detected 393,023 associations at a $P < 10^{-8}$ (Supplementary Table 2),
198 61% of those were within the HLA region. Similarly, almost half (i.e. 48%) of the
199 analyzed imputed HLA alleles were significant ($P < 10^{-8}$) for at least one binary trait
200 (Supplementary Fig. 3). Approximately 1 in 15,000 of the genotype-phenotype
201 pairs was genome-wide significant ($P < 10^{-8}$) for binary traits, whilst approximately
202 1 in 200 genotype-phenotype pairs were significant ($P < 10^{-8}$) for non-binary traits.
203 Among the tested genetic variants, one in ~80 was associated with at least one
204 binary trait, whilst one in ~10 was associated with one non-binary trait. Only
205 genetic variants within the HLA region were associated with more than 20 binary
206 traits each (Figs. 3, Supplementary Fig. 1 and 6).

207

208 We found a positive correlation ($r=0.64$, $P < 10^{-76}$ in the observed scale, $r=0.56$,
209 $P < 10^{-53}$ in the liability scale) between the heritability of the binary trait and the
210 number of genome-wide significant variants, albeit of smaller magnitude to that
211 found for the non-binary traits (Fig. 4). Some of these traits were obvious outliers
212 as they had large heritabilities but few significantly associated variants. The three
213 largest heritabilities for binary traits were for three autoimmune diseases

214 (ankylosing spondylitis, coeliac disease and seropositive rheumatoid arthritis) but
215 few significant variants were found outside the HLA region for these traits. For
216 instance, 5,704 out of 5,706 genome-wide significant associations for ankylosing
217 spondylitis were within the HLA region.

218

219 Among the categories for inclusion of genetic variants in the genotyping array
220 there was a substantial enrichment for HLA (79%), ApoE (48%), and Cancer
221 common variants (40%). The categories with the lowest enrichment were
222 genome-wide coverage for low frequency variants (0.15%) and tags for
223 Neanderthal ancestry (0.8%) (Supplementary Table 5).

224

225 We show three examples of Manhattan plots for binary traits (Fig. 5). The first
226 example shows where there are associations with skin cancer (i.e melanoma and
227 other malignant neoplasms of the skin). There are 4795 variants associated
228 ($P < 10^{-8}$) with skin cancer distributed among 172 independent lead variants
229 (Supplementary Table 3). We found associations in genetic variants in or around
230 known susceptibility genes (e.g. MC1R, IRF4, TERT, TYR) for melanoma⁸, but
231 also genes like FOXP1 (rs13316357, $P = 1.5 \times 10^{-15}$) associated with basal cell
232 carcinoma⁹. The other two examples show the similarity between the results of
233 one of the self-reported and clinically defined traits available in UK Biobank. The
234 Manhattan plots for self-reported and clinically defined coeliac disease are very
235 similar but not identical, which suggests that generally there will be benefit in
236 analyzing both clinically and self-reported traits.

237

238 **Heritability Estimates**

239 Heritability estimates inform about the contribution of genetics to the observed
240 phenotypic variation. The heritability of many of the 778 traits analysed here has
241 never been reported, but even if they have been reported it is useful to know how
242 much phenotypic variation is captured by genetic variants in a cohort of the size
243 and interest of UK Biobank. The majority (78%) of the traits analyzed had a
244 significant SNP-heritability ($P < 0.05$; Fig. 6), with the largest SNP-heritability being
245 for ankylosing spondylitis, which was 0.86 on the liability scale. The mean and
246 median heritability among those estimates that were significant were 0.12 and

247 0.08, respectively. Mean heritabilities were significantly different for binary and
248 non-binary traits ($h^2_{\text{Non-binary}}=0.17$; $h^2_{\text{Binary}}=0.10$; $P=4\times 10^{-12}$). A total of thirty-six
249 traits, all binary, had a heritability estimate close to zero ($h^2_{\text{Liability}} < 10^{-4}$). Only
250 seven of those thirty-six traits had no genome-wide significant hits ($P < 10^{-8}$), with
251 nine having more than ten significant hits, self-reported gastritis having the largest
252 number of hits with 41. This scenario could arise for monogenic and oligogenic
253 traits for which the model assumptions do not hold or because of false positives.
254 The Manhattan plots for the traits that had the largest numbers of hits seem more
255 consistent with these hits being false positives or perhaps lack of power to detect
256 heritability than with the violation of the model assumptions (Supplementary Fig.
257 7).

258

259 Estimates of genetic and environmental correlations show that for 15% of the
260 pairs of non-binary traits the genetic and environmental correlation changes sign
261 (Supplementary Fig. 8, GeneATLAS web page). Across all pairs of non-binary
262 traits for which the genetic and environmental correlation had the same sign the
263 absolute value of the genetic correlation was smaller in 31% of the cases. Overall,
264 taking into account the size of observed heritabilities, this suggests that the
265 phenotypic covariance of many of these traits is likely driven by the environment
266 and not genetics (average $(\text{cov}_g/\text{cov}_e)=0.24$, among traits where cov_g and cov_e
267 have the same sign).

268

269 **Phenotypic prediction from genetic markers**

270 We computed genomic predictions (that is, models of phenotypic prediction
271 based on genetic markers) for all 692 non-gender dependent traits using
272 Genomic Best Linear Predictions (GBLUP)¹⁰ (Methods). GBLUP estimates
273 polygenic risk scores assuming that all fitted variants have an effect. It has been
274 argued that this method has several advantages to traditional polygenetic risk
275 scores from GWAS hits^{10,11}. Some of the traits for which we developed GBLUP
276 models did indeed reach large prediction accuracies (Fig. 7), which was further
277 increased when we used additional covariates such as gender or sex. The largest
278 prediction accuracy for a non-binary trait was for height which was 0.59, whilst
279 the largest discriminative ability for a binary trait was 0.82 for self-reported

280 malabsorption/coeliac disease. We observed a large correlation between the
281 prediction accuracy and the trait heritability (Fig. 7 and Supplementary Table 8).
282 Furthermore, we previously developed a model that predicted the benefit of
283 having increasingly large training datasets for prediction of complex traits in UK
284 Biobank^{11,12}. Our current accuracy of prediction for anthropomorphic traits is very
285 similar to the ones we previously predicted we would achieve with this training
286 set¹¹ (Supplementary Fig. 9).

287 **DISCUSSION**

288 We used circa 452,000 related and unrelated UK Biobank participants of white
289 ethnicity to build the largest atlas of genetic associations to date. Summary
290 statistics for 778 traits will be available to the research community to help them
291 gain further insight into the genetic architecture of complex traits. Unlike other
292 currently available databases, like the GWAS catalog (which contains
293 ~39,366 unique SNP-trait associations), our database includes significant and
294 non-significant associations, thus providing an unbiased view of phenotype-
295 genotype associations across a large number of traits within a single cohort. In
296 addition, the database contains 182,266 independent genotype-phenotype
297 associations, genetic and environmental correlations, and estimates of SNP
298 heritability to allow researchers to perform their own filters on what a meaningful
299 association or heritability is. We hope this database will be useful to those
300 working on complex traits genetics, but also to those that have not got the
301 expertise or capabilities to perform analyses at this scale.

302

303

304 **ACCESSION CODES**

305 This research has been conducted using the UK Biobank Resource under project
306 788.

307

308 **ACKNOWLEDGEMENTS**

309 This research has been conducted using the UK Biobank Resource under project
310 788. The work was funded by the Roslin Institute Strategic Programme Grant
311 from the BBSRC (BB/P013732/1) and MRC grant (MR/N003179/1) granted to
312 AT. AT also acknowledge funding from the Medical Research Council and OCX

313 from MRC fellowship MR/R025851/1. Analyses were performed using the
314 ARCHER UK National Supercomputing Service.

315

316 **AUTHOR CONTRIBUTIONS**

317 All authors contributed equally to the design, running of the analyses, and writing
318 of the manuscript.

319 **COMPETING INTEREST STATEMENT**

320 The authors declare no competing financial interests.

321

322 **ETHICAL COMPLIANCE**

323 The UK Biobank project was approved by the National Research Ethics Service
324 Committee North West-Haydock (REC reference: 11/NW/0382). An electronic
325 signed consent was obtained from the participants.

326

327 **URLs**

328 GeneATLAS, <http://geneatlas.roslin.ed.ac.uk>; UK Biobank,
329 <http://www.ukbiobank.ac.uk>; ARCHER UK National Supercomputing Service,
330 <http://www.archer.ac.uk>; DISSECT, <https://www.dissect.ed.ac.uk>; GWAS catalog
331 <https://www.ebi.ac.uk/gwas/>; Affymetrix array
332 <https://affymetrix.app.box.com/s/6gc2mcw2s6a7zbb7wijn>; PLINK,
333 <http://zzz.bwh.harvard.edu/plink/> and <http://www.cog-genomics.org/plink/1.9/>).

334 BGENIX and BGEN reference implementation, <https://bitbucket.org/gavinband/bgen>

335 .

336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373

REFERENCES

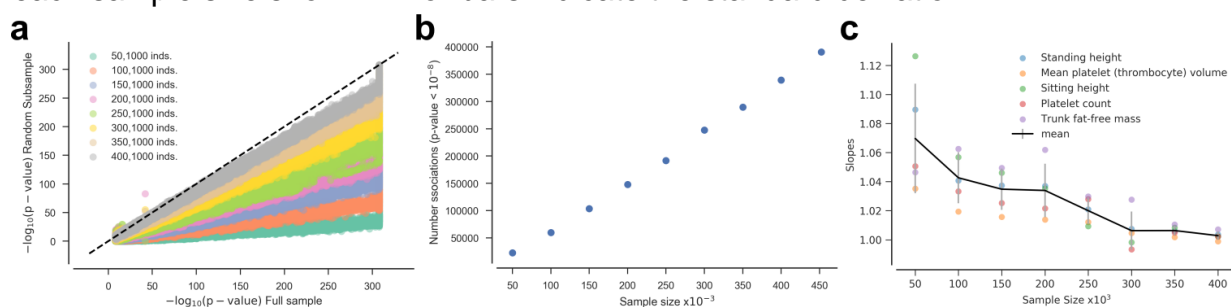
1. Falconer, D.S. & Mackay, T.F.C. *Introduction to Quantitative Genetics*, (Longman, 1996).
2. Sudlow, C. *et al.* UK biobank: an open access resource for identifying the causes of a wide range of complex diseases of middle and old age. *PLoS Med.* **12**, e1001779 (2015).
3. Canela-Xandri, O., Law, A., Gray, A., Woolliams, J.A. & Tenesa, A. A new tool called DISSECT for analysing large genomic data sets using a Big Data approach. *Nat. Commun.* **6**, 10162 (2015).
4. Loh, P.-R., Kichaev, G., Gazal, S., Schoech, A.P. & Price, A.L. Mixed-model association for biobank-scale datasets. *Nat. Genet.* **50**, 906-908 (2018).
5. McCarthy, S. *et al.* A reference panel of 64,976 haplotypes for genotype imputation. *Nat. Genet.* **48**, 1279-83 (2016).
6. Palmer, C. & Pe'er, I. Statistical correction of the Winner's Curse explains replication variability in quantitative trait genome-wide association studies. *PLOS Genet.* **13**, e1006916 (2017).
7. Yang, J. *et al.* Common SNPs explain a large proportion of the heritability for human height. *Nat. Genet.* **42**, 565-9 (2010).
8. Ransohoff, K.J. *et al.* Two-stage genome-wide association study identifies a novel susceptibility locus associated with melanoma. *Oncotarget* **8**, 17586-17592 (2017).
9. Chahal, H.S. *et al.* Genome-wide association study identifies 14 novel risk alleles associated with basal cell carcinoma. *Nat. Commun.* **7**, 12510 (2016).
10. Meuwissen, T., Hayes, B. & Goddard, M. Prediction of total genetic value using genome-wide dense marker maps. *Genetics* **157**, 1819-1829 (2001).
11. Canela-Xandri, O., Rawlik, K., Woolliams, J.A. & Tenesa, A. Improved Genetic Profiling of Anthropometric Traits Using a Big Data Approach. *PloS one* **11**, e0166755 (2016).
12. Daetwyler, H. D., Villanueva, B. & Woolliams, J. A. Accuracy of predicting the genetic risk of disease using a genome-wide approach. *PLoS One* **3**, e3395 (2008).

374 **Figures**

375

376

377 **Figure 1: The effect of sample size on the number of GWAS hits and their**
378 **estimated effects. (a)** Comparison between the p-values (two-sided t-test)
379 obtained using the whole cohort (452,264 individuals) and random subsamples
380 of increasing sizes. The plot shows only the results for the genetic variants
381 associated with a p-value $< 10^{-8}$ in the whole cohort. **(b)** Total number of
382 detected associated variants (two-sided t-test) at a threshold of p-value $< 10^{-8}$
383 as a function of the sample size. **(c)** Slope of the effect sizes of the GWAS hits
384 obtained in random subsamples of increasing size vs the same effect sizes
385 estimated in the whole cohort. Slopes larger than one indicate an inflation on
386 the effect estimates in the smaller sample. The black line joints the mean at
387 each sample size shown. Error bars indicate the standard deviation.



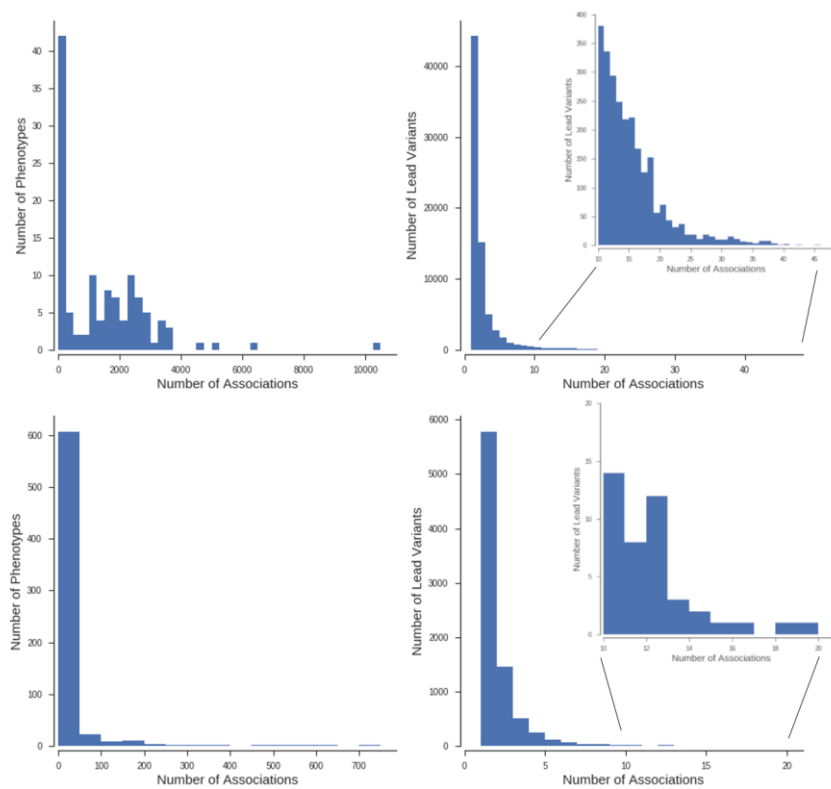
388

389

390

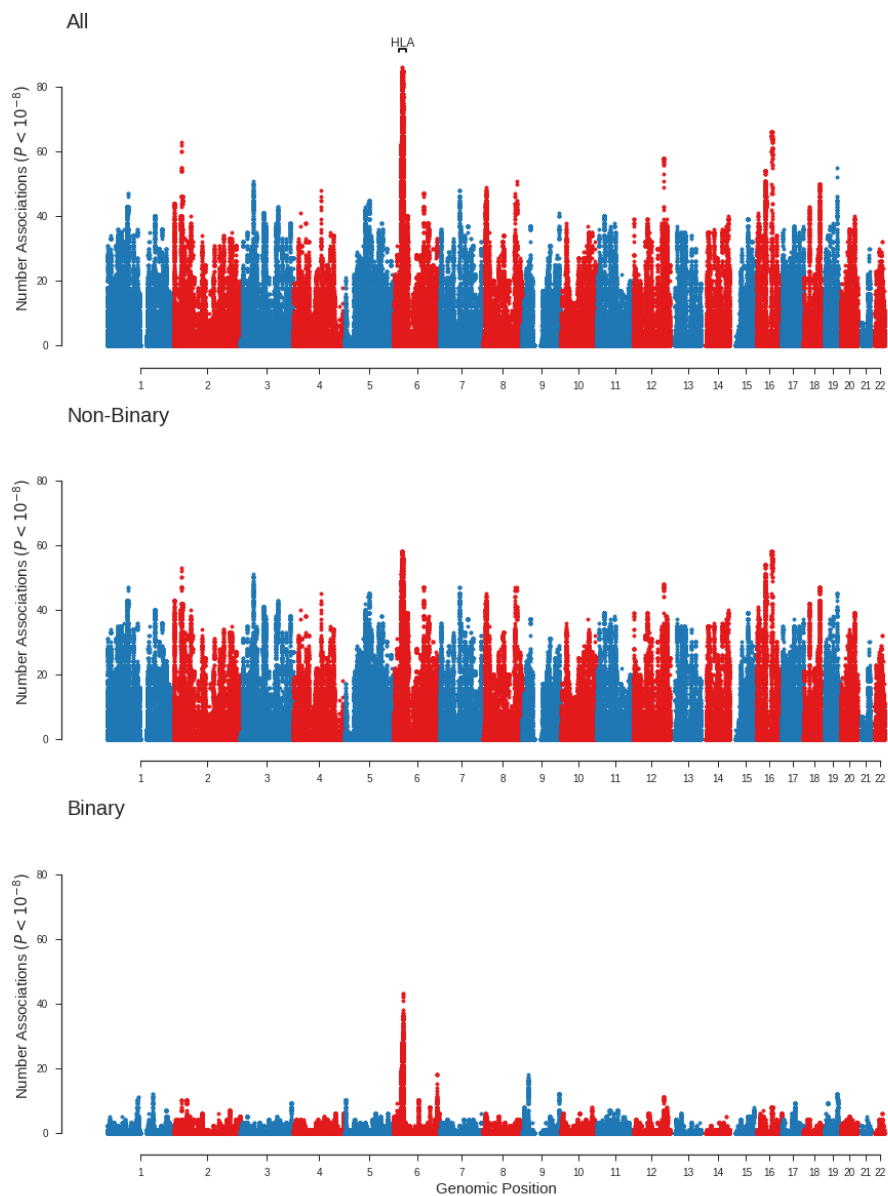
391

392 **Figure 2. Histograms of numbers of significant associations (two-sided t-**
393 **test, $P < 10^{-8}$).** The panels show results for each phenotype (left) and
394 independent lead variant (right) for non-binary (top) and binary (bottom)
395 phenotypes.
396



397
398
399
400
401
402
403
404
405
406
407
408

409 **Figure 3. Number of significant associations (two-sided t-test, $P < 10^{-8}$).**
410 The panels show the number of significant associations at each tested genetic
411 variant for all traits, non-binary and binary phenotypes. The HLA region
412 ($\pm 10\text{Mb}$) is indicated.
413
414
415



416

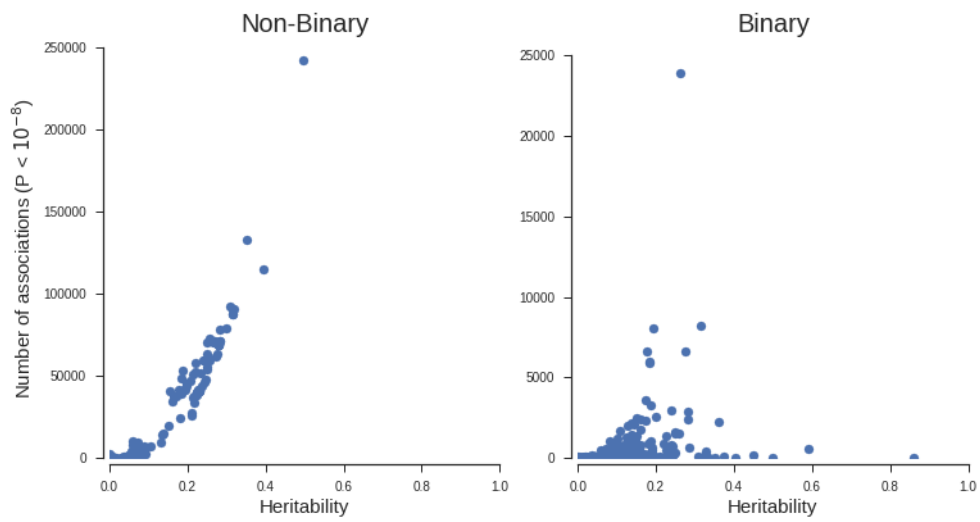
417 **Figure 4. Relationship between estimated SNP heritability and numbers**
418 **of genome wide significant associations (two-sided t-test, $P < 10^{-8}$).** HLA
419 and surrounding 10Mb region were excluded for non-binary and binary
420 phenotypes respectively.

421

422

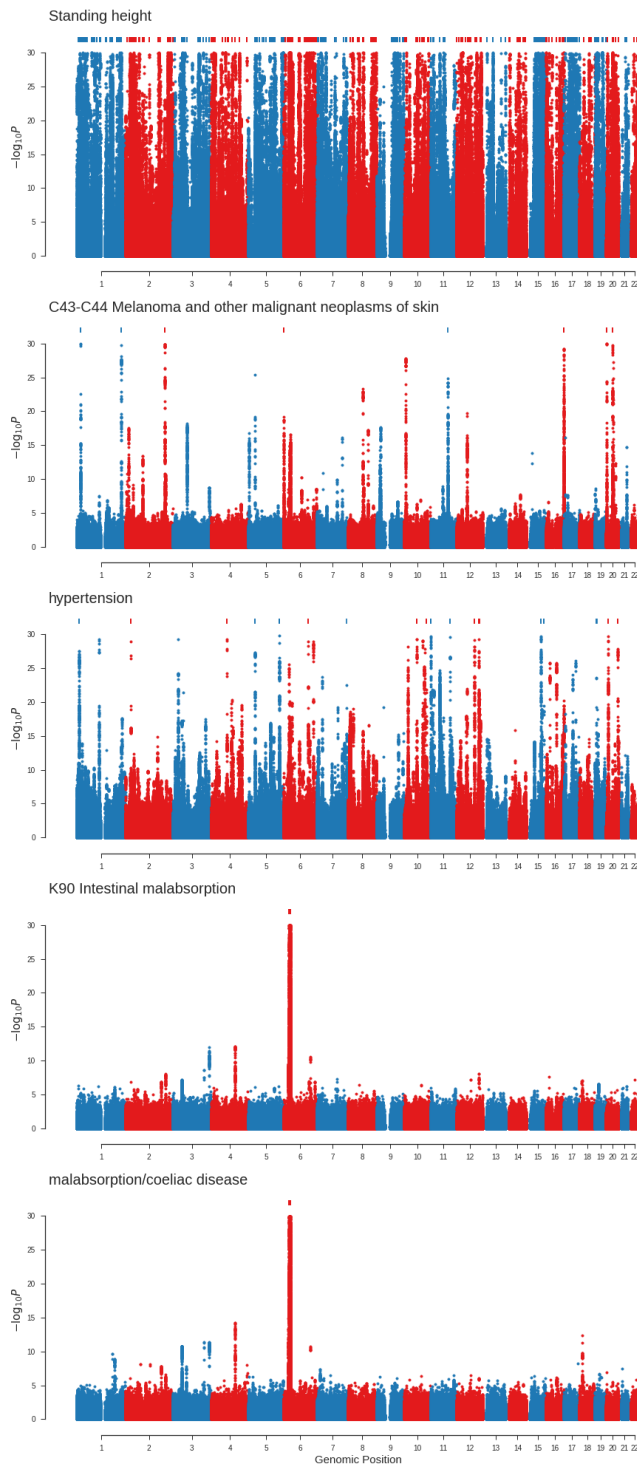
423

424



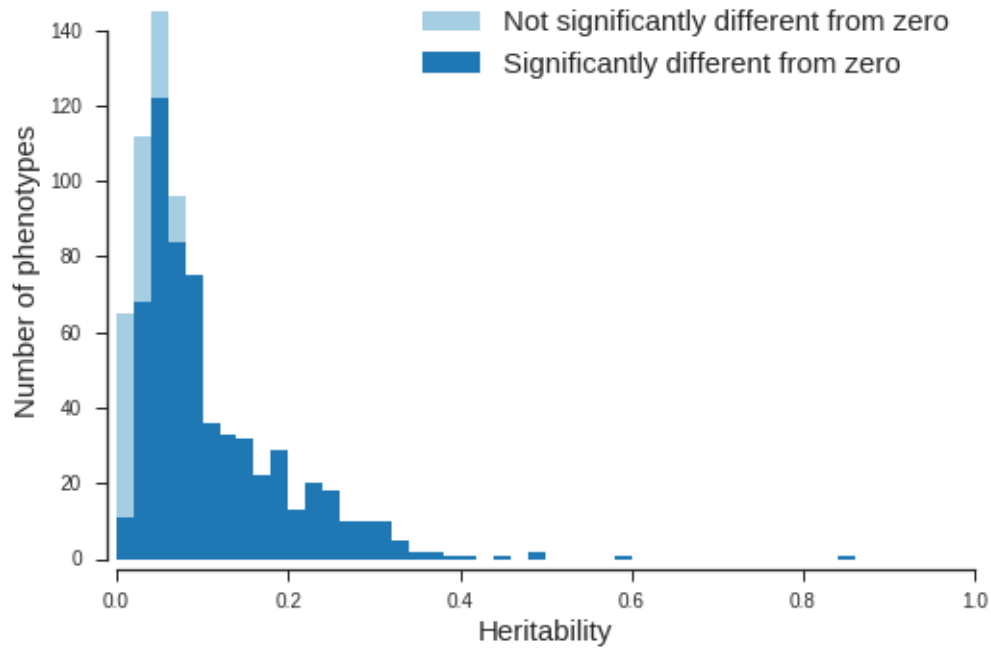
425

426 **Figure 5. Manhattan plots for selected phenotypes.** Manhattan plots for the
427 phenotypes with the largest number of genome wide significant associations
428 (two-sided t-test, $P < 10^{-8}$) within each of these categories: non-binary
429 phenotypes, cancer registry, self-reported non-cancer illness, clinically defined
430 disease from hospital episode statistics and matching self-reported disease to
431 the clinically defined disease from hospital episode statistics. From top to
432 bottom: non-binary phenotypes (Standing height), cancer registry (Melanoma
433 and other malignant neoplasms of skin), self-reported non-cancer illness
434 (hypertension), clinically defined malabsorption, and self-reported
435 malabsorption. Genetic variants with $P < 10^{-30}$ are indicated by marks along the
436 top of each plot.



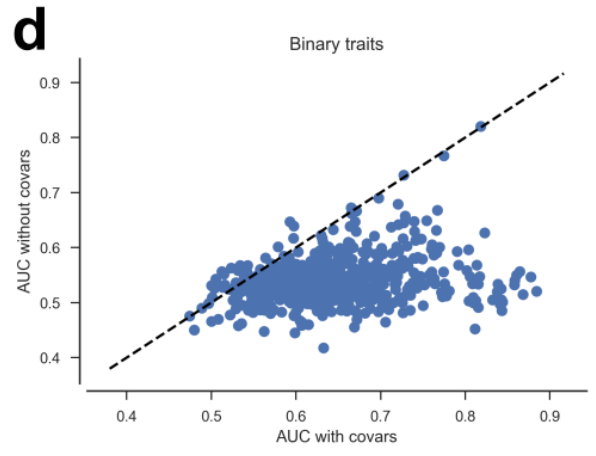
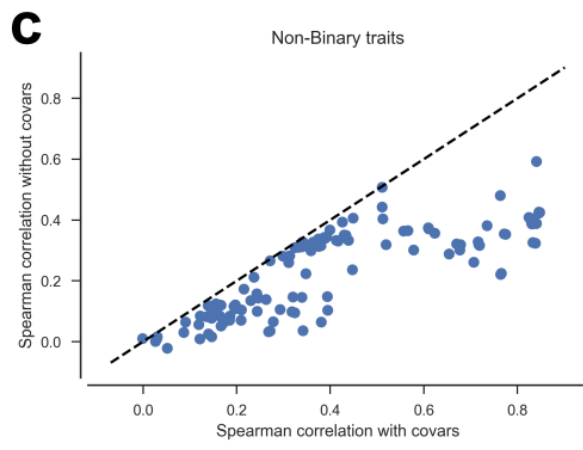
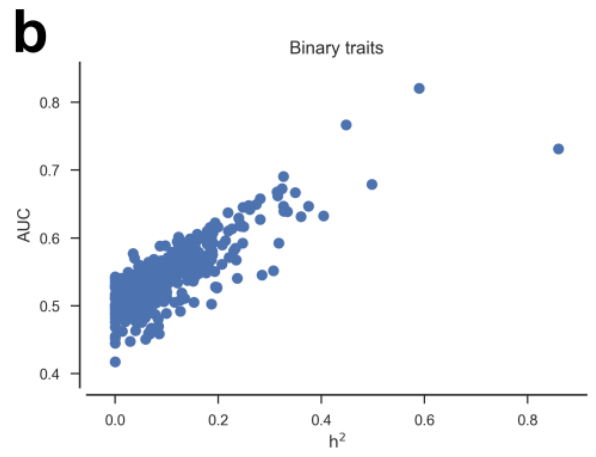
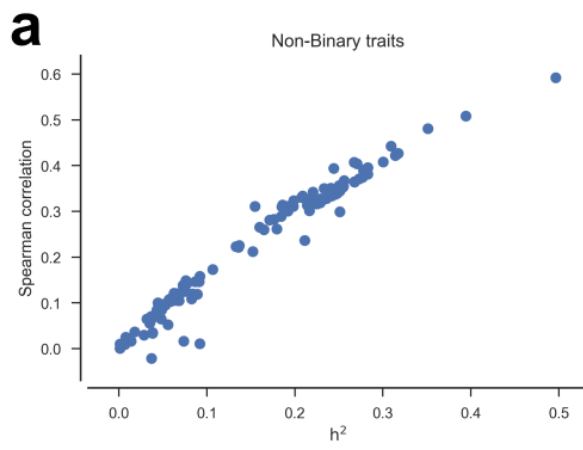
437
 438
 439
 440
 441
 442
 443

Figure 6. Numbers of phenotypes of different SNP heritability. Colours indicate the fraction of phenotypes with heritability significantly ($P < 0.05$, Chi-squared test, see Online Methods for details) different from zero in each bin.



444
 445
 446
 447
 448
 449
 450
 451
 452
 453
 454
 455
 456
 457
 458
 459
 460
 461
 462
 463
 464
 465
 466
 467
 468
 469
 470

Figure 7. Phenotypic prediction accuracy from genetic markers. Accuracy of phenotypic prediction as a function of the estimated SNP-heritability for (a) non-binary traits and (b) binary traits when no covariates were used for prediction. Comparison between prediction accuracy when covariates are included or not included for (c) non-binary traits and (d) binary traits.



471
472

473 **ONLINE METHODS**

474 **Phenotypes**

475 In total we analysed 778 phenotypes in UK Biobank participants of white
476 ethnicity. These included 657 binary phenotypes generated from self-reported
477 disease status (UK Biobank field 20002), ICD10 codes from hospitalization
478 events (UK Biobank fields 41202 and 41204), and ICD10 codes from cancer
479 registries (UK Biobank fields 40006), as well as a further 3 binary and 118 non-
480 binary (comprising continuous and ordered integral measures) phenotypes
481 from across the UK Biobank. Amongst the 660 binary phenotypes 86 exhibited
482 either a complete lack of cases in one sex or a strong imbalance in prevalence
483 in the two sexes, i.e, the ratio between the smaller and larger prevalence was
484 <0.02 . Of these 86 phenotypes 72 were specific to women. We only included
485 individuals of the appropriate sex, i.e., the sex with higher prevalence, in the
486 analysis of these sex specific phenotypes. A description of each phenotype, its
487 category and the relevant UK Biobank fields can be found in Supplementary
488 Table 1 and Gene ATLAS website. The non-binary phenotypes were not scale
489 transformed, so the units of the effect sizes are in the units reported in the UK
490 Biobank database. The phenotypes for individuals with negative coding were
491 replaced with the corresponding value (Supplementary Table 9). We also
492 ordered the keys for the ordinal phenotypes with unordered keys in the UK
493 Biobank database (Supplementary Table 10). The individuals with a phenotype
494 departing 10 standard deviations from their gender mean were set as missing
495 for traits with a value type defined as "Integer" or "Continuous" by UK Biobank.
496 The exceptions to this were Number of self-reported cancers (134-0.0), Number
497 of self-reported non-cancer illnesses (135-0.0), Nucleated red blood cell
498 percentage (30230-0.0), Nucleated red blood cell count (30170-0.0), and
499 Frequency of solarium/sunlamp use (2277-0.0) which were left as reported by
500 UK Biobank. Some of the traits analysed have some redundancy that has been
501 left for completeness. That is, some of these traits were measured in different
502 ways during the study (e.g. weight) or are analysed as self-reported traits and
503 clinical traits (e.g. malabsorption). For disease traits all individuals reporting a
504 disease code were coded as cases with all other individuals considered
505 controls. Only non-disease phenotypes with missing data rate $< 5\%$ were

506 selected for analysis. For these phenotypes missing values were imputed to
507 the age and sex specific mean in the study cohort.

508

509 **Analysis Checks**

510 Extensive validation steps were performed to ensure the reliability of the data
511 (Supplementary Material). These steps included, for instance, a comparison of
512 effect sizes with previous results from GWAS published in GWAS Catalog
513 (Supplementary Figs. 10-18), the investigation of how the polygenicity of the
514 traits drive inflation factors in GWAS (Supplementary Fig. 19), and comparisons
515 with repeated analyses where the non-binary phenotypes containing at least
516 500 different values were transformed using a rank-based normal
517 transformation (Supplementary Note, Supplementary Table 11, and
518 Supplementary Fig. 20). The results are in good agreement. Since the
519 statistical power may be different in some cases, the results are available at the
520 GeneATLAS web. Furthermore, the comparison between our heritability
521 estimations with previously published heritabilities showed a good agreement
522 (Supplementary Fig. 21 and Supplementary Table 12) when comparing ten
523 traits. In addition, we computed the Q-Q plots (Supplementary Fig. 22, and
524 summary plots in GeneATLAS website). We also checked whether there were
525 any areas depleted of associations, that is, that showed few significant
526 associations (Supplementary Fig. 23 and 24). Finally, we compared the
527 coherence of the effect size directions estimated with the whole cohort and
528 subsets of it of different sizes (Supplementary Table 13).

529

530 **Genotypes**

531 The genotypes of the UK Biobank participants were assayed using either of two
532 genotyping arrays, the Affymetrix UK BiLEVE Axiom or Affymetrix UK Biobank
533 Axiom array. These arrays were augmented by imputation of ~90 million
534 genetic variants from the Haplotype Reference Consortium⁵, the thousand
535 genomes¹³ and the UK 10K¹³ projects. Full details regarding these data have
536 been published elsewhere¹⁴.

537

538 We excluded individuals who were identified by the UK Biobank as outliers
539 based on either genotyping missingness rate or heterogeneity, whose sex

540 inferred from the genotypes did not match their self-reported sex and who were
541 not of white ancestry (based on both, self-reported ethnicity and those from
542 whom one of the two first genomic principal components did not fall within 5
543 standard deviations from the mean). Finally, we removed individuals with a
544 missingness >5% across variants which passed our quality control procedure
545 and those that have a missing phenotype for 40 or more traits. The resulting
546 study cohort comprised 452,264 individuals.

547

548 From the genotyped data we only retained bi-allelic autosomal variants which
549 were assayed by both genotyping arrays employed by UK Biobank. We
550 furthermore excluded variants which had failed UK Biobank quality control
551 procedures in any of the genotyping batches. Additionally, for imputed and
552 genotyped variants, we excluded variants with $P < 10^{-50}$ for departure from
553 Hardy-Weinberg, computed on a subset of 344,057 unrelated (Kinship
554 coefficient < 0.0442) individuals in the White-British subset of the study cohort,
555 and with a missingness rate > 2% in the study cohort. Although we analysed all
556 imputed variants and all genotyped variants with $MAF > 10^{-4}$ (all results
557 available on the GeneATLAS website), only imputed variants with $MAF > 10^{-3}$
558 in the study cohort and imputation score larger than 0.9 were used for the
559 summary results presented here. This cut-off corresponds to less than 905
560 occurrences of the minor allele in the study cohort. We also filtered the HLA
561 imputed alleles that were present in fewer than 10 individuals.

562

563 **GWAS Analysis**

564 To test each genetic variant whilst taking into account population structure in
565 UK Biobank (e.g. presence of related individuals or local structure), we used a
566 Linear Mixed Model. Specifically, the model takes the form

567

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{g} + \boldsymbol{\epsilon},$$

568 where \mathbf{y} is the vector of phenotypes, \mathbf{X} , is the matrix of fixed effects, and $\boldsymbol{\beta}$ the
569 effect size of these effects. We included as fixed effects sex, array batch, UK
570 Biobank Assessment Center, age, age², and the leading 20 genomic principal
571 components as computed by UK Biobank. \mathbf{g} is the polygenic effect that
572 captures the population structure, fitted as a random effect. It follows the

573 distribution $\mathbf{g} \sim \mathbf{N}(0, \mathbf{A}\sigma_g^2)$, with \mathbf{A} the Genomic Relationship Matrix (GRM), and
574 σ_g^2 the variance explained by the additive genetic effects. The GRM was
575 computed using common (MAF > 5%) genotyped variants that passed quality
576 control. Finally, $\epsilon \sim \mathbf{N}(0, \mathbf{I}\sigma_\epsilon^2)$ is a residual effect not accounted for by the fixed
577 and random effects. Under this model, the phenotype vector \mathbf{y} , follows the
578 distribution $\mathbf{N}(\mathbf{X}\boldsymbol{\beta}, \mathbf{A}\sigma_g^2 + \mathbf{I}\sigma_\epsilon^2)$.

579

580 Fitting one instance of such a LMM model is computationally very demanding.
581 Following a naïve approach, the required computational time increasing with
582 the cube of the sample size, $\sim O(N^3)$, and the memory requirements with the
583 square of the sample size, $\sim O(N^2)$. Consequently, fitting a single model on a
584 cohort of the size of UK Biobank is challenging, and fitting millions of these
585 models, one for each analysed genetic variant and phenotype is not feasible
586 with standard computational and statistical approaches. To address this
587 problem, we took advantage of three different tools. First, we used a large
588 supercomputer, and DISSECT³ to speed up the calculations (e.g. computing
589 the GRM eigen-decomposition required 5,040 processor cores working
590 together for ~ 10 h, and using ~ 5 TB of memory). Second, we computed the full
591 eigen decomposition of the GRM, $\mathbf{A} = \boldsymbol{\Lambda}\boldsymbol{\Sigma}\boldsymbol{\Lambda}^T$, where $\boldsymbol{\Lambda}$ is the matrix of
592 eigenvectors, and $\boldsymbol{\Sigma}$ is a diagonal matrix containing the eigenvalues. This
593 allowed us to transform all the other model matrices, \mathbf{y} , \mathbf{X} , and ϵ to the new
594 space where the GRM is diagonal. Although the eigen-decomposition is a
595 computationally intensive process, once diagonalized, the computational time
596 of fitting a model is reduced considerably to $\sim O(N)$, thus enabling us to perform
597 several tests using Mixed Linear Models on a cohort of hundreds of thousands
598 of individuals. Finally we performed over 23 billion tests using a two-step
599 approximation that optimizes the computational resources¹⁵. The first step of
600 the approximation fits a LMM that adjusts by the relevant fix (e.g. age, sex, etc.)
601 and random effects (genetic effects) to each trait, the second step uses the
602 residuals of LMM to test (two-tailed t-test on effect sizes) all available genetic
603 markers for significance in a linear model. We corrected for the polygenic effect
604 using a Leave-One-Chromosome-Out (LOCO) approach¹⁶.

605

606 **HLA Region**

607 We defined the *HLA* region as the region of chromosome 6 spanning base pairs
608 28,866,528 to 33,775,446. Throughout all analyses we included 10Mb either
609 side of the above *HLA* region to account for LD with variants outside this region.
610 The imputed HLA alleles were tested using the same GWAS model described
611 above, where the independent variable is the best guess allele reported dosage
612 from the HLA imputed values (UK Biobank field 22182). We tested the alleles
613 using two models. A model where the number of copies of each HLA allele for
614 each locus was tested independently as a fixed effect, and a second model
615 where the number of copies of all alleles in a given locus were tested together
616 as fixed effects in the same model (i.e. an omnibus test)¹⁷.

617

618 **Estimation of Genetic Parameters**

619 In order to estimate heritabilities and genetic correlations we fitted LMMs for
620 each trait with a GRM containing all common (MAF > 5%) autosomal genetic
621 variants which passed QC. The heritability was estimated as $h_g^2 =$
622 $\sigma_g^2 / (\sigma_g^2 + \sigma_e^2)$, where σ_g^2 and σ_e^2 are the estimates of the genetic and residual
623 variance and the p-values were obtained using a Chi-squared test following the
624 method described previously^{18,19}. For all binary outcomes, we transformed
625 heritabilities on the observed scaled to the liability scale using the population
626 prevalence of the disease. We provide sex-specific prevalences to allow sex-
627 specific transformations (Supplementary Table 1). Using the model fits we
628 computed best linear unbiased predictor estimates of genetic additive values
629 for each individual. The genetic correlations were estimated by computing
630 correlations between these additive genetic values. Environmental correlations

631 were estimated as $r_e = (r_y - \sqrt{h_i^2 h_j^2} r_g) / \sqrt{(1 - h_i^2)(1 - h_j^2)}$, where r_y, r_g are the
632 phenotypic and genetic correlations for traits i and j .

633

634 **Lead variants and Independent Loci**

635 We clustered GWAS results into independent lead variants using the `--clump`
636 option of the PLINK 1.9 software^{20,21}. Specifically, for each trait individually, we
637 clustered GWAS results by selecting genome wide significant variants as lead
638 variants and assigning to them unassigned variants within 10Mb, that have

639 $P < 10^{-2}$ and a $r^2 > 0.3$ with the lead variant. To compute the total number of
640 independent loci across all traits, we performed the same clustering on the lead
641 variants across all traits, choosing the lowest p-value for variants which were
642 lead variants in different traits.

643

644 **Relation of number of associations and chromosome length**

645 We regressed the number of significant associations ($P < 10^{-8}$) across traits for
646 each chromosome on the covered length of the chromosome, i.e., distance in
647 base pairs of the first and last tested genetic variants, and the number of genetic
648 variants tested on the chromosome. For chromosome 6 we excluded the HLA
649 region and variants contained therein from the statistics. We compared the full
650 model to one with either the chromosomal length or number of tested genetic
651 variants removed using the likelihood ratio test. The full model was not
652 significantly better than the model containing only chromosomal length
653 ($P = 0.08$) but was significantly better than the model containing only the number
654 of genetic variants ($P = 0.004$). Both reduced models were significant when
655 compared to a null model containing only an intercept.

656

657

658 **Phenotypic prediction**

659 The effect of all common genetic variants ($MAF > 0.05$) were estimated together
660 as a random effect using the model,

$$661 \quad y_i = \mu + \sum_{l=1}^L x_{il}\beta_l + \sum_{j=1}^M z_{ij}a_j + e_i,$$

662 where μ is the mean term and e_i the residual for individual i . L is the number of
663 fixed effects, x_{il} being the value for the fixed effect l at individual i and β_l the
664 estimated effect of the fixed effect l . We fitted the same covariates as in the
665 GWAS analyses. M is the number of markers and z_{ij} is the standardised
666 genotype of individual i at marker j . The vector of effects of random common
667 genetic variants \mathbf{a} is distributed as $N(0, \mathbf{I}\sigma_u^2)$. The vector of environmental

668 effects \mathbf{e} is distributed as $N(0, \mathbf{I}\sigma_e^2)$. Defining $\sigma_g^2 = M\sigma_u^2$, the heritabilities were
669 estimated as $\sigma_g^2 / (\sigma_e^2 + \sigma_g^2)$.

670 The prediction of the phenotype \hat{y}_i for the individual i was computed as a sum
671 of the product of the SNP effects and the number of reference alleles of the
672 corresponding SNPs:

$$673 \quad \hat{y}_i = \sum_{j=1}^M \frac{(s_{ij} - \mu_j^*)}{\sigma_j^*} a_j,$$

674 where s_{ij} is the number of copies of the reference allele at marker j of individual
675 i , M is the number of markers used for the prediction, and a_j the effect of marker
676 j . μ_j^* and σ_j^* are the mean and the standard deviation of the effect allele in the
677 training population.

678 We used 407,669 genetically confirmed white British to train the models and
679 44,595 whites of non-British descent to validate the models. We restricted this
680 analysis to the 692 non-gender specific phenotypes. Prediction accuracies for
681 non-binary traits were computed as the Spearman correlation between the
682 predicted and the real phenotype of white participants of non-British descent
683 after correcting by the estimated effect of the used covariates. Prediction
684 accuracies for binary traits were computed as the Area Under the Curve (AUC)
685 of a Receiver Operating Characteristic (ROC) curve using the predicted and
686 the real phenotypes of white individuals of non-British descent.

687

688 **Reporting Summary**

689 Further information on experimental design is available in the Life Sciences
690 Reporting Summary linked to this article.

691

692 **Code availability**

693 The source code of DISSECT, the tool used for GWAS and heritability
694 estimations, is freely available at <https://www.dissect.ed.ac.uk> under GNU Lesser
695 General Public License v3.

696

697 **Data availability**

698 All summary results from the analyses performed are available at GeneATLAS
699 website, <http://geneatlas.roslin.ed.ac.uk/>.

700

701

702 **ONLINE METHODS REFERENCES**

- 703 13. Genomes Project, C. *et al.* An integrated map of genetic variation from
704 1,092 human genomes. *Nature* **491**, 56-65 (2012).
- 705 14. Bycroft, C.F., C.; Petkova, D.; Band, G.; Elliott, L.T.; Sharp, K.; Motyer,
706 A.; Vukcevic, D.; Delaneau, O.; O'Connell, J.; Cortes, A.; Welsh, S.;
707 McVean, G.; Leslie, S.; Donnelly, P.; Marchini, J. Genome-wide genetic
708 data on ~500,000 UK Biobank participants. *Biorxiv* (2017).
- 709 15. Aulchenko, Y.S., de Koning, D.J. & Haley, C. Genomewide rapid
710 association using mixed model and regression: A fast and simple
711 method for genomewide pedigree-based quantitative trait loci
712 association analysis. *Genetics* **177**, 577-585 (2007).
- 713 16. Yang, J., Zaitlen, N.A., Goddard, M.E., Visscher, P.M. & Price, A.L.
714 Advantages and pitfalls in the application of mixed-model association
715 methods. *Nat. Genet.* **46**, 100 (2014).
- 716 17. Patsopoulos, N.A. *et al.* Fine-Mapping the Genetic Association of the
717 Major Histocompatibility Complex in Multiple Sclerosis: HLA and Non-
718 HLA Effects. *PLOS Genet.* **9**, e1003926 (2013).
- 719 18. Stram, D.O. & Lee, J.W. Variance Components Testing in the
720 Longitudinal Mixed Effects Model. *Biometrics* **50**, 6 (1994).
- 721 19. Visscher, P.M. A Note on the Asymptotic Distribution of Likelihood Ratio
722 Tests to Test Variance Components. *Twin Res. Hum. Genet.* **9**, 490-495
723 (2012).
- 724 20. Purcell, S. *et al.* PLINK: a tool set for whole-genome association and
725 population-based linkage analyses. *Am. J. Hum. Genet.* **81**, 559-75
726 (2007).
- 727 21. Chang, C.C. *et al.* Second-generation PLINK: rising to the challenge of
728 larger and richer datasets. *GigaScience* **4**, 1-16 (2015).
- 729

730

731

732

733