

# 1 **Accurate genetic profiling of anthropometric traits using a** 2 **big data approach**

3

## 4 **Authors:**

5 Oriol Canela-Xandri<sup>1,\*</sup>, Konrad Rawlik<sup>1,\*</sup>, John A. Woolliams<sup>1</sup>, Albert Tenesa<sup>1,2,3,\*</sup>

6

## 7 **Affiliations:**

8 <sup>1</sup>The Roslin Institute, Royal (Dick) School of Veterinary Studies, The University of Edinburgh,  
9 Easter Bush Campus, Midlothian, EH25 9RG. Scotland. UK.

10 <sup>2</sup>MRC HGU at the MRC IGMM, University of Edinburgh, Western General Hospital, Crewe  
11 Road South, Edinburgh. EH4 2XU. UK

12

13 \*These authors contributed equally to this work

14

15

16

17 <sup>3</sup> Corresponding author

18 Dr Albert Tenesa

19 The Roslin Institute

20 The University of Edinburgh

21 Easter Bush

22 Roslin, Midlothian

23 EH25 9RG

24 Scotland

25 Tel: 0044 (0)131 651 9100

26 Fax: 0044 (0)131 651 9220

27 Email: [Albert.Tenesa@ed.ac.uk](mailto:Albert.Tenesa@ed.ac.uk)

28

29

30

31 **Genome-wide association studies (GWAS) promised to translate their findings into**  
32 **clinically beneficial improvements of patient management by tailoring disease**  
33 **management to the individual through the prediction of disease risk<sup>1,2</sup>. However, the**  
34 **ability to translate genetic findings from GWAS into predictive tools that are of clinical**  
35 **utility and which may inform clinical practice has, so far, been encouraging but**  
36 **limited<sup>1,2</sup>. Here we propose to use a more powerful statistical approach that enables the**  
37 **prediction of multiple medically relevant phenotypes without the costs associated with**  
38 **developing a genetic test for each of them. As a proof of principle, we used a common**  
39 **panel of 319,038 SNPs to train the prediction models in 114,264 unrelated White-British**  
40 **for height and four obesity related traits (body mass index, basal metabolic rate, body**  
41 **fat percentage, and waist-to-hip ratio). We obtained prediction accuracies that ranged**  
42 **between 46% and 75% of the maximum achievable given their explained heritable**  
43 **component. This represents an improvement of up to 75% over the phenotypic variance**  
44 **explained by the predictors developed through large collaborations<sup>3</sup>, which used more**  
45 **than twice as many training samples. Across-population predictions in White non-**  
46 **British individuals were similar to those of White-British whilst those in Asian and Black**  
47 **individuals were informative but less accurate. The genotyping of circa 500,000 UK**  
48 **Biobank<sup>4</sup> participants will yield predictions ranging between 66% and 83% of the**  
49 **maximum. We anticipate that our models and a common panel of genetic markers,**  
50 **which can be used across multiple traits and diseases, will be the starting point to tailor**  
51 **disease management to the individual. Ultimately, we will be able to capitalise on whole-**  
52 **genome sequence and environmental risk factors to realise the full potential of genomic**  
53 **medicine.**

54

55 Phenotypic prediction of complex traits from genomic data could transform clinical practice by  
56 enabling tailored treatment and targeted disease screening programs based on the genetic

57 make-up of the individual, and by facilitating more efficient allocation of resources within the  
58 health systems<sup>5-7</sup>. Ultimately, it would help to understand the underlying disease mechanisms  
59 and open the targeted search of specific solutions based on this knowledge. With this in mind,  
60 large efforts and investments in the past years have been directed towards generating  
61 genotypic and phenotypic data for identifying individual genetic variants associated with  
62 different traits through genome-wide association studies (GWAS)<sup>8</sup>. Although using this  
63 approach a large number of susceptibility variants for many diseases have been identified, the  
64 strategy has several limitations. First, the accuracy of prediction has been disappointingly low  
65 for traits affected by a large numbers of susceptibility variants<sup>7</sup>. Second, the approach of  
66 identifying one single nucleotide polymorphism (SNP) at a time and including such newly  
67 identified SNPs in the prediction models as and when they are identified is unpractical if one  
68 wishes to use genetic tests for multiple traits because the composition of each trait's genetic  
69 test would need to be continuously updated and each trait would require its own SNP panel.  
70 Third, statistical considerations and simulation studies have shown that the accuracy of  
71 prediction for complex traits increases by modelling all available SNPs simultaneously<sup>9</sup>.

72

73 Recent studies have shown that SNP arrays containing common genetic variants capture a  
74 substantial amount of the genetic variation for each trait and that the contributing SNPs have  
75 effects generally too small to be detected with current GWAS sample sizes due to the stringent  
76 genome-wide significance levels applied<sup>3,10,11</sup>. Furthermore, we have previously shown  
77 through simulations that the size of the studies that have estimated heritability from SNP  
78 arrays have been too small to properly estimate SNP effects for accurate phenotypic  
79 prediction<sup>12</sup>. However, the availability of large genotyped cohorts for which individual-level  
80 data is available, e.g. the UK Biobank<sup>4,13</sup>, combined with new and powerful computational  
81 tools<sup>12</sup> capable of fitting complex statistical models to big datasets and access to high-  
82 performance computational infrastructure has the potential to provide accurate SNP effects  
83 for genomic prediction.

84

85 We show that modelling individual-level data of circa 120,000 individuals can lead to accurate  
86 predictions across multiple traits by jointly fitting the SNPs of a single array of common SNPs.  
87 These predictions significantly improve on the accuracies of models derived from summary  
88 statistics obtained from large GWAS meta-analyses and in turn would ease clinical  
89 implementation and direct-to-consumer genetic testing, as well as improve the accuracy of the  
90 predictions as the sample sizes of the training datasets increase.

91

92 We first focused on human height, a highly heritable quantitative trait commonly used as a  
93 model in the study of the genetic architecture of complex traits<sup>3,10,14</sup> and one of the traits for  
94 which most contributing loci have been identified to date. To increase the generality of our  
95 findings, we then selected four obesity related traits - BMI, percentage body fat, waist-to-hip  
96 ratio (WHR) and basal metabolic rate (BMR). In order to jointly estimate the additive effects of  
97 all SNPs we fitted them as random effects in a Mixed Linear Model (MLM) on the training  
98 population, with gender and age as fixed effects (Online Methods). As the computational  
99 requirements of MLM fitting rapidly increases with incrementing sample sizes, we used  
100 DISSECT<sup>12</sup> (<https://www.dissect.ed.ac.uk>), a software specifically designed to perform  
101 genomic analysis in large supercomputers. Each analysis required ~1h of computing time on  
102 the ARCHER supercomputer, harnessing the joint power of 1,152 processors. Using the  
103 estimated SNP effects to predict the genetic value of individuals in an independent validation  
104 dataset (Online Methods), we computed the prediction accuracy as the correlation between  
105 these predicted genetic values and the phenotypes corrected for gender and age.

106 For our analyses, we used the 152,736 genotyped individuals available from the UK Biobank  
107 cohort<sup>4</sup>. After applying stringent quality control criteria, we divided our sample into White-  
108 British (123,847 individuals) and non White-British (27,685 individuals), the latter including  
109 individuals from different ethnic backgrounds (Online Methods and Supplementary Fig. 1). We  
110 divided the White-British further into a group of 114,264 unrelated individuals with a  
111 relatedness below 0.0625 (i.e. less related to each other than second cousins once removed),

112 another group of 9,583 individuals that had at least one relationship above 0.0625 with the  
113 unrelated White-British group, and a group of self-reported White-British (Online Methods and  
114 Supplementary Fig. 1). We modelled 319,038 common SNPs, that is, variants with a minor  
115 allele frequency (MAF)  $>0.05$  that passed our genotype quality control.

116 We used the 114,264 unrelated White-British individuals to train the prediction models and  
117 assessed the validity of the within-population predictions using the 9,583 related White-British  
118 individuals and the 12,640 self-reported White-British individuals. Prediction accuracy in the  
119 self-reported White-British ranged from 0.51 (95% CI 0.49-0.52) for height to 0.20 (95% CI 0.19-  
120 0.22) for WHR (Table 1). We evaluated whether prediction accuracies can be further improved  
121 by using more complex models (Online Methods). The accuracy for height improved to 0.55  
122 (95% CI 0.53-0.56) despite small reduction in the estimate of heritability. However, accuracies  
123 for the other traits decreased. The accuracies we obtained represent between 75% and 46%  
124 of the maximum achievable given the estimated SNP-based heritabilities of the traits in  
125 unrelated White-British, i.e., 0.53 (SE = 0.004), 0.26 (SE = 0.005), 0.26 (SE = 0.005), 0.20  
126 (SE = 0.005) and 0.31 (SE = 0.005) for height, body fat percentage, BMI, WHR, and BMR,  
127 respectively (Online Methods). As expected, phenotypic prediction of relatives was more  
128 accurate than that of unrelated people because their phenotypes and genotypes are more  
129 correlated to the training samples. The robustness of our within-population predictions were  
130 confirmed using 10-fold cross-validation<sup>15</sup> within the White-British participants (Supplementary  
131 Table 1). The SNP effects are available as Supplementary Information.

132 We also investigated to what extent across-population predictions were feasible. To this end,  
133 we further subdivided the non White-British subset by self-reported ethnic background.  
134 Excluding ethnicities with less than 1,000 individuals and removing outliers resulted in 7,541  
135 White individuals who did not self-report as White-British, 1,954 Asian or Asian-British  
136 individuals, and 1,591 Black or Black-British individuals (Online Methods). Predictions  
137 obtained in the White cohort (Table 2) were almost as accurate as to those obtained in the  
138 self-reported White-British cohort. Predictions for the other two ethnicities remained

139 considerable but lower than within-population predictions, especially for Black or Black British  
140 as expected from the genetic distance between populations (Supplementary Fig. 2), indicating  
141 that predictions may benefit from within-ethnic group tailored models.

142 Although sample sizes for training the models will increase in the future, it is unlikely that they  
143 will increase indefinitely. Therefore, we argued that it would be useful to know what sample  
144 size would be required to exploit all the genetic variation captured by the SNP array. To gauge  
145 that, we computed prediction accuracies for samples of decreasing size, by randomly  
146 subsampling the unrelated White-British individuals (Online Methods). Our data fitted very well  
147 to a well-known theoretical model<sup>16</sup>. Our predictions suggest that prediction accuracies for  
148 height will reach 0.6 (SE = 0.02) when training the models using ~500,000 individuals (Fig. 1),  
149 the samples planned to be genotyped UK Biobank in the near future. This prediction accuracy  
150 would represent 82% of the maximum accuracy possible given the explained heritability.  
151 Similarly, we estimate that genetic prediction models for BMI, WHR, body fat percentage and  
152 BMR will reach prediction accuracies of 0.36 (SE = 0.05), 0.29 (SE = 0.03), 0.37 (SE = 0.03),  
153 and 0.42 (SE = 0.03) respectively (Supplementary Fig. 3).

154 Our results confirm previous findings that many variants with small effect can explain a large  
155 proportion of the genetic variance. Due to several factors, this part of the genetic variance has  
156 so far remained largely unexploited for phenotypic prediction. These factors include the  
157 statistical methods used, the available sample sizes, and computational software available to  
158 analyse the data. However, as we have shown, predictions which are significantly more  
159 accurate can be obtained by increasing sample sizes and using powerful computational  
160 approaches to jointly estimate all SNP effects. The phenotypic variance explained by our  
161 predictor for height is ~75% larger than that of the largest height meta-analysis to date, which  
162 used a discovery sample size ~250% larger than ours, but that used genome-wide significant  
163 SNPs from the meta-analyses as predictors<sup>3</sup>. Our prediction accuracies are very close to the  
164 maximum achievable given the estimates of the corresponding explained heritable  
165 components, and we predict that they will become even closer when the number of samples

166 increase (e.g. when the UK Biobank is fully genotyped). For BMI, which is affected by a smaller  
167 genetic component, our predictor explains slightly more variance than previous work that used  
168 a discovery sample size ~3 times larger<sup>17</sup>. Finally, we demonstrated that more complex models  
169 have the potential to further improve prediction accuracies, although our results also indicate  
170 that the optimal model may be trait specific. In conclusion, the presented results support our  
171 initial hypothesis and suggest a promising future for genomic prediction of complex traits.

172

## 173 **Methods**

### 174 **Genotype Quality Control**

175 For our analysis, we used the data for the genotyped individuals in phase 1 of the UK Biobank  
176 genotyping program. 49,979 individuals were genotyped by using the Affymetrix UK BiLEVE  
177 Axiom array and 102,750 individuals by using the Affymetrix UK Biobank Axiom array. Details  
178 regarding genotyping procedure and genotype calling protocols are provided elsewhere  
179 (<http://biobank.ctsu.ox.ac.uk/crystal/refer.cgi?id=155580>). From the overlapping markers, we  
180 excluded those which were multi-allelic, their overall missingness rate exceeded 2% or they  
181 exhibited a strong platform specific missingness bias (Fisher's exact test,  $P < 10^{-100}$ ). We also  
182 excluded individuals if they exhibited excess heterozygosity, as identified by UK Biobank  
183 internal QC procedures (<http://biobank.ctsu.ox.ac.uk/crystal/refer.cgi?id=155580>), if their  
184 missingness rate exceeded 5% or if their self-reported sex did not match genetic sex estimated  
185 from X chromosome inbreeding coefficients. These criteria resulted in a reduced dataset of  
186 151,532 individuals. Finally, we only kept the common variants (i.e. with a MAF > 0.05) and  
187 those that did not exhibit departure from Hardy-Weinberg equilibrium ( $P < 10^{-50}$ ) in the  
188 unrelated (subset of individuals with a relatedness below 0.0625) White-British cohort (see  
189 below).

190

## 191 **Ethnicity**

192 The UK Biobank samples are from individuals of diverse ethnicities. To define the White-British  
193 cohort, we performed a Principal Components Analysis (PCA) of all individuals passing  
194 genotypic QC using a linkage disequilibrium (LD) pruned set of 99,101 autosomal markers  
195 (<http://biobank.ctsu.ox.ac.uk/crystal/refer.cgi?id=149744>) that passed our SNP QC protocol.  
196 The related and unrelated White-British individuals were defined as those for whom the  
197 projections onto the leading twenty genomic principal components fell within three standard  
198 deviations of the mean and who identified themselves as White-British. We defined the other  
199 removed White-British as self-reported White-British. The other ethnicities were defined using  
200 the self-identified ethnic background (<http://biobank.ctsu.ox.ac.uk/crystal/field.cgi?id=21000>).  
201 As we did with White-British individuals, we only retained those individuals whose projections  
202 onto the leading twenty genomic principal components fell within three standard deviations of  
203 the ethnicity group mean (Supplementary Fig. 2).

## 204 **Phenotype quality control**

205 We defined outliers as males and females that were outside  $\pm 3$  standard deviations from their  
206 gender mean of all the individuals in the UK Biobank, and removed them from the analyses.

## 207 **Software**

208 The genotype quality control and data filtering was performed using plink<sup>18</sup> (<https://www.cog->  
209 [genomics.org/plink2](https://www.cog-genomics.org/plink2)). The PCA, MLMs fittings for estimating SNP effects and phenotype  
210 predictions were performed using DISSECT (<https://www.dissect.ed.ac.uk>) on the UK National  
211 Supercomputer (ARCHER). DISSECT software is designed to perform genomic analyses on  
212 very large sample sizes without the need to perform mathematical approximations by using  
213 the power of large supercomputers.

## 214 **Phenotype prediction**



215 The effect of all SNPs were estimated together as a random effect using the model,

$$216 \quad y_i = \mu + \sum_{l=1}^L x_{il}\beta_l + \sum_{j=1}^M z_{ij}a_j + e_i,$$

217 where  $\mu$  is the mean term and  $e_i$  the residual for individual  $i$ .  $L$  is the number of fixed effects,  
218  $x_{il}$  being the value for the fixed effect  $l$  at individual  $i$  and  $\beta_l$  the estimated effect of the fixed  
219 effect  $l$ .  $M$  is the number of markers and  $z_{ij}$  is the standardised genotype of individual  $i$  at  
220 marker  $j$ . The vector of random SNP effects  $a$  is distributed as  $N(0, \mathbf{I}\sigma_u^2)$ . The vector of  
221 environmental effects  $e$  is distributed as  $N(0, \mathbf{I}\sigma_e^2)$ . The heritabilities were estimated as  
222  $\sigma_u^2/(\sigma_e^2 + \sigma_u^2)$ .

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

$$225 \quad \hat{y}_i = \sum_{j=1}^M \frac{(s_{ij} - \mu_j^*)}{\sigma_j^*} a_j$$

226 Where  $s_{ij}$  is the number of copies of the reference allele at SNP  $j$  of individual  $i$ ,  $M$  is the  
227 number of SNPs used for the prediction, and  $a_j$  the effect of SNP  $j$ .  $\mu_j^*$  and  $\sigma_j^*$  are the mean  
228 and the standard deviation of the reference allele in the training population.

229 Prediction accuracies were computed as the correlation between the predicted phenotype and  
230 the real one after correcting by the estimated effect of the used covariates (e.g. sex and age).

### 231 **Phenotype prediction using a two variance components model**

232 The MLM of the previous section assumes that all SNP effects follow a Gaussian distribution  
233 with one variance. However this is may not be true. To improve the model we first fitted all  
234 SNPs independently using a standard GWAS model,

$$235 \quad y_i = \mu + \sum_{l=1}^L x_{il}\beta_l + a_j^* + e_i,$$

236 Here, the parameters are the same as in the previous MLM, and the SNP effect size  $a_i^*$  is  
237 estimated independently for each SNP as a fixed effect. We then divided the SNPs into two  
238 groups based on their effect size. Specifically, one group of SNPs in the main distribution and  
239 a group of outliers, which were defined as SNPs with effect sizes more than 3 standard  
240 deviations away from the mean effect across all SNPs. Using these groups, we fit an extended  
241 MLM where we assume the SNP effects were distributed in two different Gaussian  
242 distributions with a different variance each one,

$$243 \quad y_i = \mu + \sum_{l=1}^L x_{il} \beta_l + \sum_{j=1}^M z_{ij}^m a_j^m + \sum_{k=1}^K z_{ik}^t a_k^t + e_i,$$

244 where all parameters are the same as in the simpler MLM, but now M and K are the number  
245 of SNPs in the main distribution and the two tails, respectively, while  $z_{ij}^m$  and  $z_{ik}^t$  are the  
246 corresponding genotypes. We fit independent variances for the two groups of SNPs, so that  
247 the vector of SNP effects in the main distribution,  $\mathbf{a}^m$ , is distributed as  $N(0, \mathbf{I}(\sigma_u^m)^2)$  and the  
248 vector of SNP effects in the tails,  $\mathbf{a}^t$ , is distributed as  $N(0, \mathbf{I}(\sigma_u^t)^2)$ .

249

## 250 **Random subsampling**

251 We computed accuracies for samples of decreasing size, by randomly subsampling 5 of the  
252 10-fold cross-validation subsets used in the within unrelated White-British population  
253 predictions (Supplementary Table 1).

254

255

## 256 **References**

257

- 258 1. Dunlop, M. G. *et al.* Cumulative impact of common genetic variants and other risk  
259 factors on colorectal cancer risk in 42,103 individuals. *Gut* **62**, 871–81 (2013).

- 260 2. Manolio, T. A. Bringing genome-wide association findings into clinical use. *Nat. Rev.*  
261 *Genet.* **14**, 549–58 (2013).
- 262 3. Wood, A. R. *et al.* Defining the role of common variation in the genomic and biological  
263 architecture of adult human height. *Nat. Genet.* **46**, 1173–1186 (2014).
- 264 4. Sudlow, C. *et al.* UK Biobank: An Open Access Resource for Identifying the Causes  
265 of a Wide Range of Complex Diseases of Middle and Old Age. *PLoS Med.* **12**,  
266 e1001779 (2015).
- 267 5. Bowles Biesecker, B. & Marteau, T. M. The future of genetic counselling: an  
268 international perspective. *Nat. Genet.* **22**, 133–7 (1999).
- 269 6. De los Campos, G., Gianola, D. & Allison, D. B. Predicting genetic predisposition in  
270 humans: the promise of whole-genome markers. *Nat. Rev. Genet.* **11**, 880–6 (2010).
- 271 7. Schrod, S. J. *et al.* Genetic-based prediction of disease traits: prediction is very  
272 difficult, especially about the future. *Front. Genet.* **5**, 162 (2014).
- 273 8. Visscher, P. M., Brown, M. A., McCarthy, M. I. & Yang, J. Five years of GWAS  
274 discovery. *Am. J. Hum. Genet.* **90**, 7–24 (2012).
- 275 9. Meuwissen, T. H. E., Hayes, B. J. & Goddard, M. E. Prediction of Total Genetic Value  
276 Using Genome-Wide Dense Marker Maps. *Genetics* **157**, 1819–1829 (2001).
- 277 10. Yang, J. *et al.* Common SNPs explain a large proportion of the heritability for human  
278 height. *Nat. Genet.* **42**, 565–9 (2010).
- 279 11. Yang, J. *et al.* Genome partitioning of genetic variation for complex traits using  
280 common SNPs. *Nat. Genet.* **43**, 519–25 (2011).
- 281 12. Canela-Xandri, O., Law, A., Gray, A., Woolliams, J. A. & Tenesa, A. A new tool called  
282 DISSECT for analyzing large genomic datasets using a Big Data approach. *Nat.*  
283 *Commun.* **In press**, (2015).
- 284 13. Collins, R. What makes UK Biobank special? *Lancet* **379**, 1173–4 (2012).
- 285 14. Lango Allen, H. *et al.* Hundreds of variants clustered in genomic loci and biological  
286 pathways affect human height. *Nature* **467**, 832–8 (2010).
- 287 15. Molinaro, A. M., Simon, R. & Pfeiffer, R. M. Prediction error estimation: a comparison  
288 of resampling methods. *Bioinformatics* **21**, 3301–7 (2005).
- 289 16. Daetwyler, H. D., Villanueva, B. & Woolliams, J. A. Accuracy of predicting the genetic  
290 risk of disease using a genome-wide approach. *PLoS One* **3**, e3395 (2008).
- 291 17. Locke, A. E. *et al.* Genetic studies of body mass index yield new insights for obesity  
292 biology. *Nature* **518**, 197–206 (2015).
- 293 18. Purcell, S. *et al.* PLINK: a tool set for whole-genome association and population-  
294 based linkage analyses. *Am. J. Hum. Genet.* **81**, 559–75 (2007).

## 296 **Supplementary Information**

297 Supplementary Information contains the SNP estimated effect sizes for each trait.

298

## 299 **Acknowledgments**

300

301 This work was mainly supported by the Medical Research Council [grant numbers  
302 MR/K014781/1 and MR/N003179/1] and The Roslin Institute Strategic Grant funding from the  
303 BBSRC. AT also acknowledges funding from the Medical Research Council Human Genetics  
304 Unit. This work used the ARCHER UK National Supercomputing Service  
305 (<http://www.archer.ac.uk>) and the Edinburgh Compute and Data Facility (ECDF)  
306 (<http://www.ecdf.ed.ac.uk/>). This research has been conducted using the UK Biobank  
307 Resource. We also acknowledge the fruitful comments received from Chris Haley, Wendy  
308 Bickmore and Pau Navarro.

309

## 310 **Contributions**

311 All authors contributed equally on this study.

312

## 313 **Ethical approval**

314 The use of the UK Biobank dataset falls within the study's ethical approval from the North  
315 West Medical Research Ethics Committee (Reference 11/NW/0382). An informed consent  
316 was obtained from all subjects.

317

318

## 319 **Competing financial interests**

320 The authors declare no competing financial interests.

321 **URLs**

322 DISSECT and documentation available at: <https://www.dissect.ed.ac.uk>

323 PLINK2 and documentation available at: <https://www.cog-genomics.org/plink2>

324

325

326

327 **Tables**

328

<b>Traits</b>	<b>self-reported White-British (95% CI)</b>	<b>related White-British (95% CI)</b>
<b>Height</b>	0.51 (0.49-0.52)	0.53 (0.52-0.55)
<b>Body fat percentage</b>	0.27 (0.25-0.29)	0.28 (0.26-0.30)
<b>BMI</b>	0.25 (0.24-0.27)	0.27 (0.26-0.29)
<b>WHR</b>	0.20 (0.19-0.22)	0.23 (0.21-0.25)
<b>BMR</b>	0.32 (0.31-0.34)	0.34 (0.32-0.36)

329 **Table 1: Prediction accuracies on related White-British and self-reported White-British.**

330

331

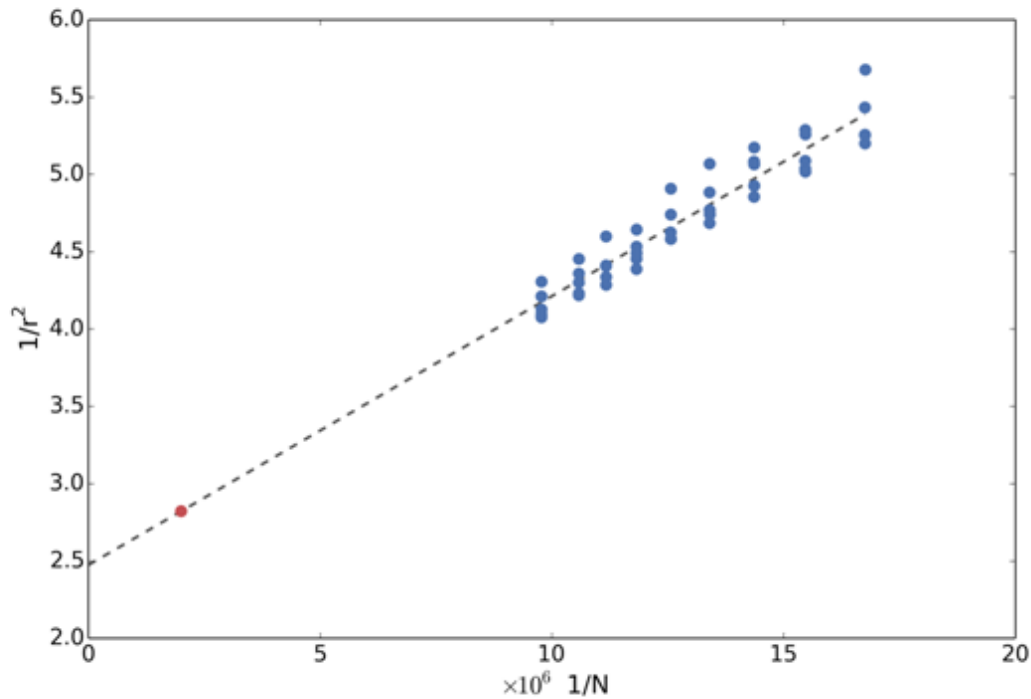
<b>Traits</b>	<b>White non British (95% CI)</b>	<b>Asian/Asian-British (95% CI)</b>	<b>Black/Black-British (95% CI)</b>
<b>Height</b>	0.50 (0.48-0.51)	0.34 (0.30-0.38)	0.18 (0.14-0.23)
<b>Body fat percentage</b>	0.26 (0.24-0.28)	0.21 (0.16-0.25)	0.12 (0.07-0.17)
<b>BMI</b>	0.25 (0.23-0.27)	0.22 (0.18-0.26)	0.11 (0.06-0.16)
<b>WHR</b>	0.21 (0.19-0.23)	0.14 (0.10-0.19)	0.07 (0.02-0.12)
<b>BMR</b>	0.32 (0.30-0.34)	0.22 (0.18-0.26)	0.12 (0.07-0.17)

332 **Table 2: Across-population prediction accuracies.**

333

334

## 335 Figures



336

337 **Figure 1: Prediction accuracy as a function of sample size for height.**

338 Inverse of the square of the prediction accuracy as a function of the inverse of the training  
339 sample size. Blue dots indicate prediction accuracies achieved on several trials. The dashed  
340 straight line shows the linear regression fit to the blue dots. The regression intercept indicates  
341 the maximum accuracy achievable using common variants represented in the array. The red  
342 dot is the expected prediction accuracy with a training sample size of 500,000 individuals.

343

344