

An atlas of genetic associations in UK Biobank

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

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

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

The unprecedented size of this cohort has raised a number of analytical challenges³. First, storing, managing and analyzing ~90 million genetic variants for around half a million individuals is, in itself, a substantial endeavor. Second, the collection of samples at this scale has brought up an analytical challenge, as the cohort is structured by familial relationships and ancestry. For instance, many relatives were unintentionally collected in the cohort, and removing them from the analyses as traditionally done in GWAS would entail a substantial loss of statistical power. Third, although recent developments have reduced the computational costs⁴, fitting a linear mixed model (LMM), the standard analytical technique to perform GWAS when there is population or familial structure, at this scale and for this number of traits, entails a computational burden that may be beyond the means of many research labs.

The objective of the current study was to perform GWAS for 778 traits in UK Biobank, adjusting for the effect of relatedness to minimize the loss of statistical power while reducing false positives due to familial and population structure, in individuals of European ancestry and to make a searchable atlas of genetic associations in UK Biobank for the benefit of the research community.

Results

Data overview. In July 2017, the UK Biobank released genotyped data from approximately 490,000 individuals of largely European descent genotyped for 805,426 genetic variants. We performed

GWAS analyses for 660 binary traits and 118 non-binary traits, the latter including continuous traits and traits with multiple ordered categories (Supplementary Table 1). For each of these traits, we fitted LMMs to test for association with 623,944 genotyped and 30,798,054 imputed genetic polymorphisms imputed using the Haplotype Reference Consortium⁵ as reference panel, as well as 310 imputed HLA alleles. All successfully tested polymorphisms are included in the database (GeneATLAS, <http://geneatlas.roslin.ed.ac.uk>) or associated downloadable files to allow individual researchers to apply their own quality control thresholds. The summary results presented here are based on the quality-controlled imputed polymorphisms (9,113,133 variants after filtering) of 452,264 individuals (Online Methods).

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

To demonstrate the power of using large datasets, we first explored how the analysis of increasingly large sample sizes enables new discoveries and reduces bias when estimating the effect sizes of GWAS hits (Fig. 1 and Supplementary Note). Our results show that the number of GWAS hits increased linearly with the sample size with no sign of saturation, thus suggesting that increasing the size of cohorts like UK Biobank would continue to yield new discoveries. We also observed that the estimated allelic effects of GWAS hits obtained from decreasing sample sizes were generally larger, which is in agreement with a 'winner's curse' effect⁶ (Fig. 1).

Distribution of GWAS hits among non-binary trait. Just under 5 million of the ~1 billion tests performed across 118 non-binary traits were significant at a conventional genome-wide threshold ($P < 10^{-8}$) (Supplementary Table 2), and 3,117,904 were significant after Bonferroni correction ($P < 0.05/9,113,133 \times 118$). The significant associations were distributed across 74,471 lead

¹The Roslin Institute, Royal (Dick) School of Veterinary Studies, The University of Edinburgh, Midlothian, UK. ²MRC Human Genetics Unit at the MRC Institute of Genetics and Molecular Medicine, University of Edinburgh, Western General Hospital, Edinburgh, UK. ³These authors contributed equally: Oriol Canela-Xandri, Konrad Rawlik, Albert Tenesa. *e-mail: albert.tenesa@ed.ac.uk

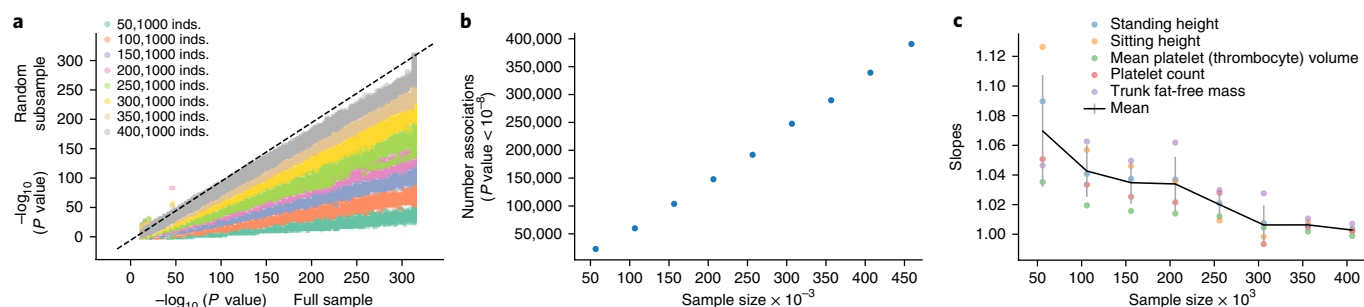


Fig. 1 | The effect of sample size on the number of GWAS hits and their estimated effects. **a**, Comparison between the P values (two-sided t -test) obtained using the whole cohort (452,264 individuals) and random subsamples of increasing sizes. The plot shows only the results for the genetic variants associated with $P < 10^{-8}$ in the whole cohort. The dashed line indicates the diagonal line and is shown for reference. **b**, Total number of detected associated variants (two-sided t -test) at a threshold of $P < 10^{-8}$ as a function of the sample size. **c**, Slope of the effect sizes of the GWAS hits obtained in random subsamples of increasing size versus the same effect sizes estimated in the whole cohort. Slopes larger than one indicate an inflation on the effect estimates in the smaller sample. The black line joins the mean at each sample size shown. Error bars indicate the standard deviation.

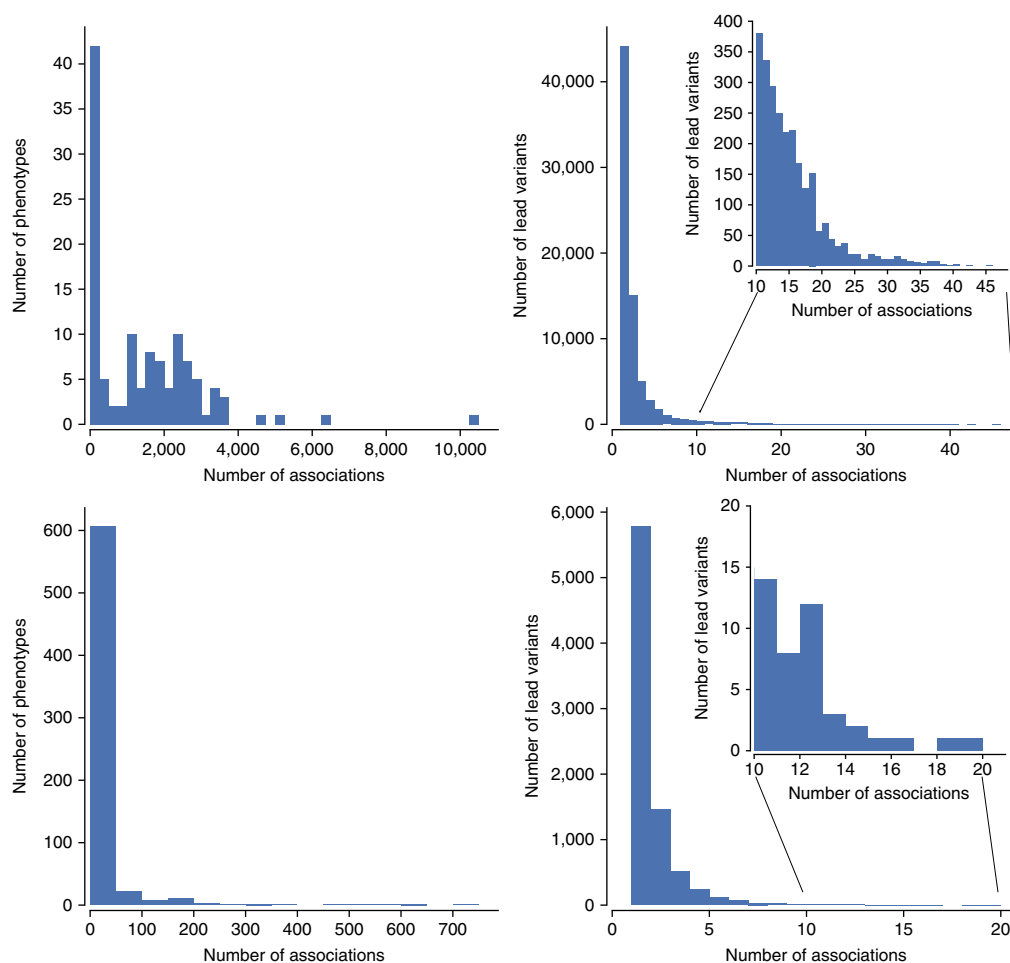


Fig. 2 | Histograms of numbers of significant associations (two-sided t -test, $P < 10^{-8}$). The panels show results for each phenotype (left) and independent lead variant (right) for non-binary (top) and binary (bottom) phenotypes.

polymorphisms mapping to 38,651 independent loci (Online Methods, Fig. 2, and Supplementary Table 3). A substantial proportion of these associations (13.0%) were within the HLA region (Supplementary Table 2).

About 9.5% of the tested polymorphisms reached genome-wide significance ($P < 10^{-8}$) for at least one of the 118 tested traits, while

82% of the tested polymorphisms were associated with at least one of these 118 traits at a P value of 10^{-2} (Supplementary Table 4). There were 20,393 genetic variants each associated with more than 30 of the tested non-binary traits (Figs. 2 and 3, and Supplementary Fig. 1). A cluster of nine variants in a 9-kb region including the genotyped intronic variant rs1421085 within the *FTO* gene had

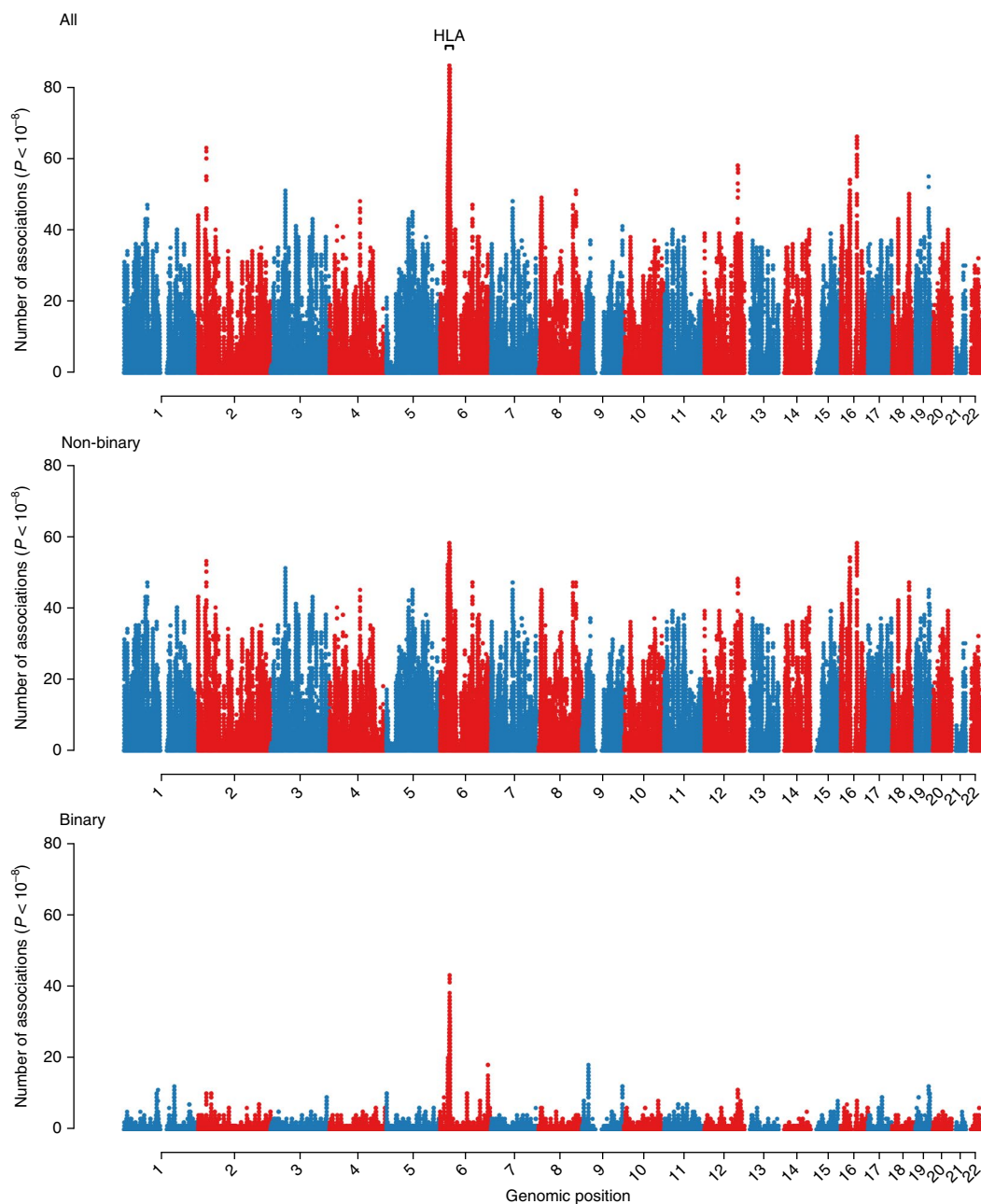


Fig. 3 | Number of significant associations (two-sided t-test, $P < 10^{-8}$). The panels show the number of significant associations at each tested genetic variant for all traits, non-binary and binary phenotypes. The HLA region (+/−10 Mb) is indicated.

the largest number of genome-wide significant associations outside the HLA region, with all nine variants found to be associated with 58 traits (Fig. 3 and Supplementary Fig. 1). The genotyped variant rs1421085 at the *FTO* locus also had the largest average significance across non-binary traits ($P < 10^{-74}$) (Supplementary Fig. 2), which was largely contributed by the associations to anthropometric traits such as body mass index (BMI) and weight, which showed some of the strongest associations ($P < 10^{-300}$). The HLA region contained 362 genetic variants that were significantly ($P < 10^{-8}$) associated with 50 or more of the non-binary traits compared with only 128 such variants in the remaining autosomal variants. About 36% of the analyzed imputed HLA alleles were significant ($P < 10^{-8}$) for at least one trait (Supplementary Fig. 3). Six traits ('standing height', 'sitting height', 'platelet count', 'mean platelet (thrombocyte) volume', 'trunk predicted mass', and 'trunk

fat-free mass') had over 100,000 significant associations ($P < 10^{-8}$) each distributed across 25,352 different independent lead genetic variants (Online Methods). Over 94% of the non-binary traits had more than 100 genome-wide significant hits distributed in 74,442 different leading genetic variants.

Considering the criteria for inclusion of genetic polymorphisms on the genotyping array (Supplementary Table 5), the HLA polymorphisms were the most enriched for associations with at least one non-binary trait (88% had $P < 10^{-8}$), followed by the cardiometabolic, autoimmune/inflammatory and apolipoprotein E (ApoE) criteria, while the lowest enrichment was for two low frequency variants categories ('genome-wide coverage for low frequency variants' and 'rare, possibly disease causing, mutations'). Fewer than 8 in 100 of these polymorphisms were associated with any non-binary trait (Supplementary Table 5).

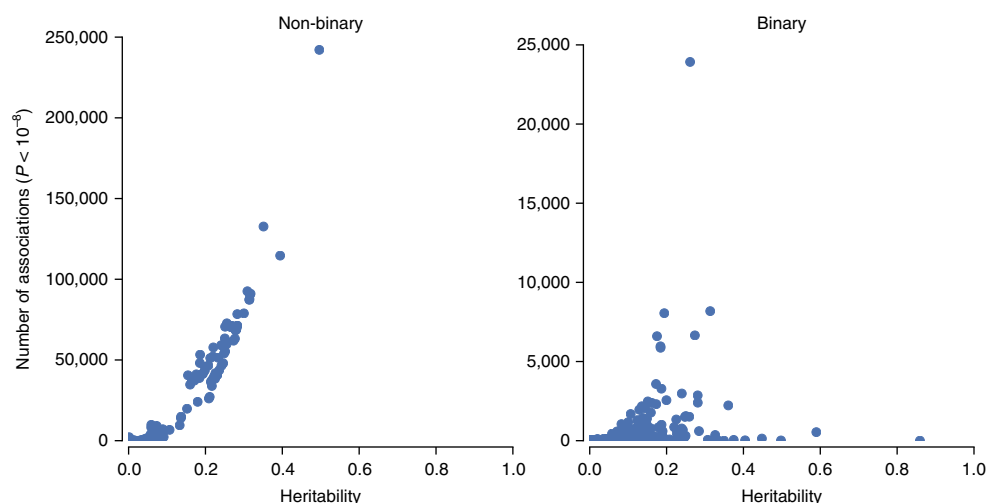


Fig. 4 | Relationship between estimated SNP heritability and numbers of genome-wide significant associations (two-sided t -test, $P < 10^{-8}$). HLA and surrounding 10-Mb region were excluded for non-binary and binary phenotypes, respectively.

We found a significant correlation ($r=0.93$, $P < 10^{-51}$) between the number of hits and the SNP heritability of the traits, suggesting that the number of loci affecting a trait might be proportional to the heritability of the trait (Fig. 4 and Supplementary Fig. 4). Consistent with this model and variation in the distribution of linkage disequilibrium across the genome, the correlation of the SNP heritability with the number of identified independent lead variants was similarly high ($r=0.88$, $P < 10^{-38}$). The number of hits ($P < 10^{-8}$) per chromosome was highly correlated ($r=0.86$) with the length of the chromosome covered by the genotyped SNPs (Supplementary Fig. 5 and Supplementary Table 6). Although this correlation could arise under a polygenic model where the length of the chromosome is correlated with the number of possible variants affecting the traits, the simplest explanation is that it arises as a consequence of the correlation of chromosomal length and number of tested variants per chromosome. We considered a model explaining the number of hits per chromosome as a function of the number of tested genetic variants and the length of the chromosome, and the two nested models including only one of the two factors. The results were consistent with the number of GWAS hits per chromosome correlating with the length of the chromosome, rather than the number of tested variants (Online Methods).

Standing height was the trait with the largest number of hits (Fig. 5), with 261,908 significantly associated variants distributed across 10,374 independent lead variants. We estimated that the lead polymorphisms across the 118 traits studied are distributed among 38,651 independent loci; therefore, 27% of these independent loci contribute to variation in height, as expected for a highly polygenic trait⁷. We also computed the proportion of tested genetic variants associated with at least one disease ($P < 10^{-8}$) that are also associated with height and BMI at different thresholds (Supplementary Table 7). At a threshold of 10^{-8} , ~28% and ~7% of the genetic variants associated with at least one disease, were also associated for height and BMI, respectively. This is important for the interpretation of Mendelian randomization studies as it is likely that one of the critical assumptions to demonstrate causality, that is, that there is no pleiotropy between the exposure and the outcome, may be broken for many exposure-outcome pairs.

Distribution of GWAS hits among binary traits. The binary trait with the largest number of cases was self-reported hypertension, with an average across binary traits of 6,593 cases (Supplementary Table 1). Of the 660 binary phenotypes, 86 were specific to one sex (Supplementary Table 1). Individuals of the unaffected sex were

excluded from the analysis for these phenotypes (Online Methods). Consistent with the reduced statistical power to detect association with binary phenotypes (mainly diseases) compared to non-binary traits, we detected 393,023 associations at $P < 10^{-8}$ (Supplementary Table 2), 61% of those were within the HLA region. Similarly, almost half (that is 48%) of the analyzed imputed HLA alleles were significant ($P < 10^{-8}$) for at least one binary trait (Supplementary Fig. 3). Approximately 1 in 15,000 of the genotype-phenotype pairs were genome-wide significant ($P < 10^{-8}$) for binary traits, while approximately 1 in 200 genotype-phenotype pairs were significant ($P < 10^{-8}$) for non-binary traits. Among the tested genetic variants, 1 in ~80 was associated with at least one binary trait, while 1 in ~10 was associated with one non-binary trait. Only genetic variants within the HLA region were associated with more than 20 binary traits each (Fig. 3 and Supplementary Figs. 1 and 6).

We found a positive correlation ($r=0.64$, $P < 10^{-76}$ in the observed scale, $r=0.56$, $P < 10^{-53}$ in the liability scale) between the heritability of the binary trait and the number of genome-wide significant variants, albeit of smaller magnitude to that found for the non-binary traits (Fig. 4). Some of these traits were obvious outliers as they had large heritabilities but few significantly associated variants. The three largest heritabilities for binary traits were for three autoimmune diseases (ankylosing spondylitis, celiac disease and seropositive rheumatoid arthritis), but few significant variants were found outside the HLA region for these traits. For instance, 5,704 out of 5,706 genome-wide significant associations for ankylosing spondylitis were within the HLA region.

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

We show three examples of Manhattan plots for binary traits (Fig. 5). The first example shows associations with skin cancer (that is melanoma and other malignant neoplasms of the skin). There are 4,795 variants associated ($P < 10^{-8}$) with skin cancer distributed among 172 independent lead variants (Supplementary Table 3). We found associations in genetic variants in or around known susceptibility genes (for example *MC1R*, *IRF4*, *TERT*, *TYR*) for melanoma⁸, but also genes like *FOXP1* (rs13316357, $P = 1.5 \times 10^{-15}$) associated with basal cell carcinoma⁹. The other two examples show the similarity between the results of one of the self-reported and clinically defined traits available in UK Biobank. The Manhattan plots for

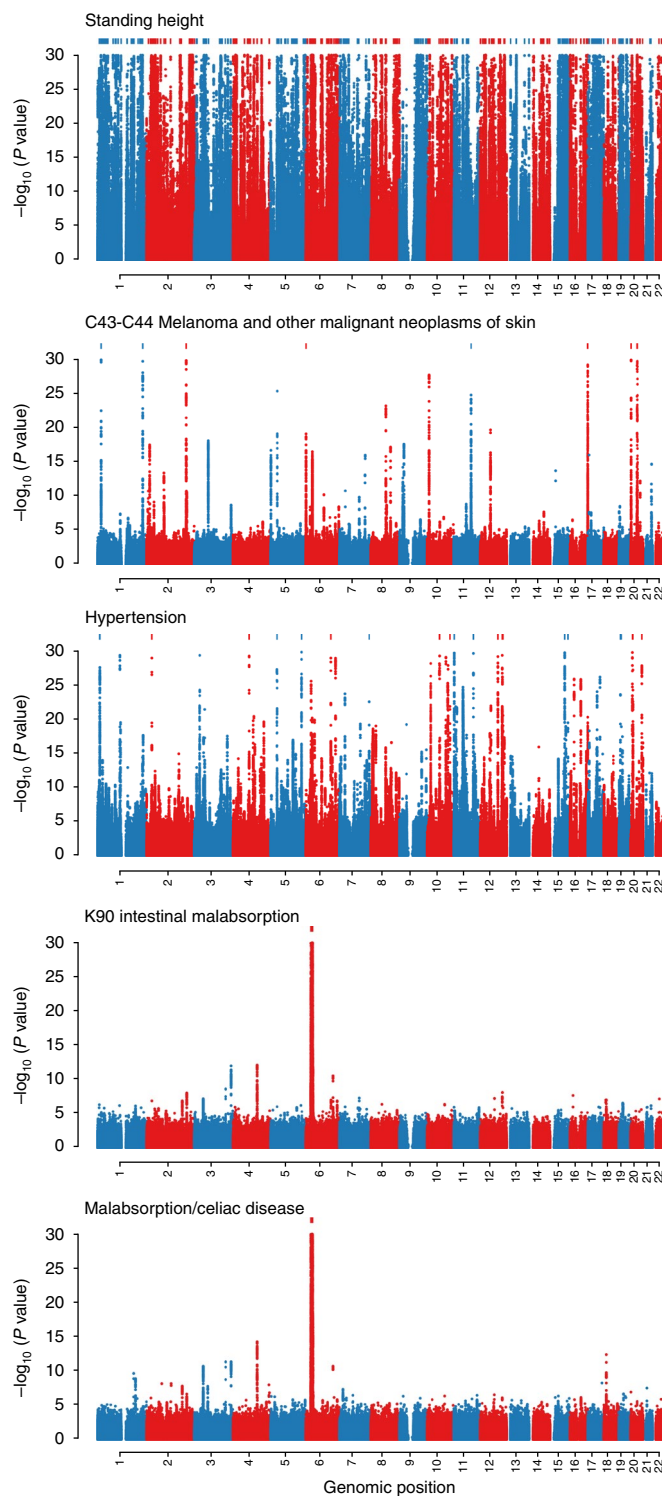


Fig. 5 | Manhattan plots for selected phenotypes. Manhattan plots for the phenotypes with the largest number of genome-wide significant associations (two-sided t -test, $P < 10^{-8}$) within each of these categories: non-binary phenotypes, cancer registry, self-reported non-cancer illness, clinically defined disease from hospital episode statistics, and matching self-reported disease to the clinically defined disease from hospital episode statistics. From top to bottom: non-binary phenotypes (standing height), cancer registry (melanoma and other malignant neoplasms of skin), self-reported non-cancer illness (hypertension), clinically defined malabsorption, and self-reported malabsorption. Genetic variants with $P < 10^{-30}$ are indicated by marks along the top of each plot.

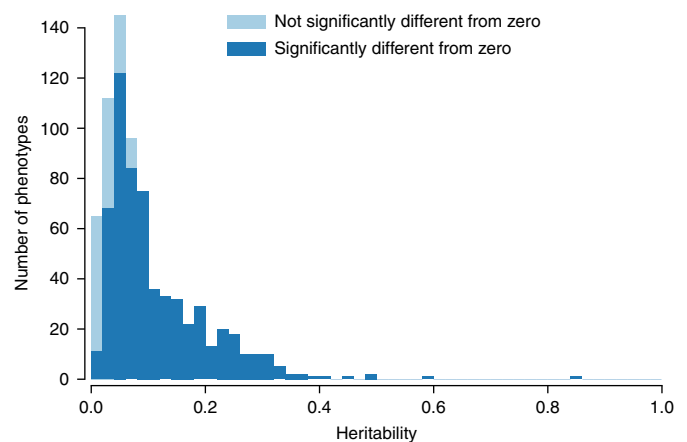


Fig. 6 | Numbers of phenotypes of different SNP heritability. Colors indicate the fraction of phenotypes with heritability significantly ($P < 0.05$, chi-squared test; see Methods for details) different from zero in each bin.

self-reported and clinically defined celiac disease are very similar but not identical, which suggests that generally there will be benefit in analyzing both clinically and self-reported traits.

Heritability estimates. Heritability estimates inform about the contribution of genetics to the observed phenotypic variation. The heritability of many of the 778 traits analyzed here has never been reported, but in any case, it is useful to know how much phenotypic variation is captured by genetic variants in a cohort of the size and interest of UK Biobank. The majority (78%) of the traits analyzed had a significant SNP-heritability ($P < 0.05$; Fig. 6), with the largest SNP-heritability being for ankylosing spondylitis, which was 0.86 on the liability scale. The mean and median heritability among those estimates that were significant were 0.12 and 0.08, respectively. Mean heritabilities were significantly different for binary and non-binary traits ($h^2_{\text{Non-binary}} = 0.17$; $h^2_{\text{Binary}} = 0.10$; $P = 4 \times 10^{-12}$). A total of 36 traits, all binary, had a heritability estimate close to zero ($h^2_{\text{Liability}} < 10^{-4}$). Only 7 of those 36 traits had no genome-wide significant hits ($P < 10^{-8}$), with 9 having more than 10 significant hits, and self-reported gastritis having the largest number of hits with 41. This scenario could arise for monogenic and oligogenic traits for which the model assumptions do not hold or because of false positives. The Manhattan plots for the traits that had the largest numbers of hits seem more consistent with these hits being false positives or perhaps lack of power to detect heritability than with the violation of the model assumptions (Supplementary Fig. 7).

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

Phenotypic prediction from genetic markers. We computed genomic predictions (that is, models of phenotypic prediction based on genetic markers) for all 692 non-gender-dependent traits using Genomic Best Linear Predictions (GBLUP)¹⁰ (Online Methods). GBLUP estimates polygenic risk scores assuming that all fitted variants have an effect. It has been argued that this method has several advantages to traditional polygenic risk scores from GWAS hits^{10,11}. Some of the traits for which we developed GBLUP models

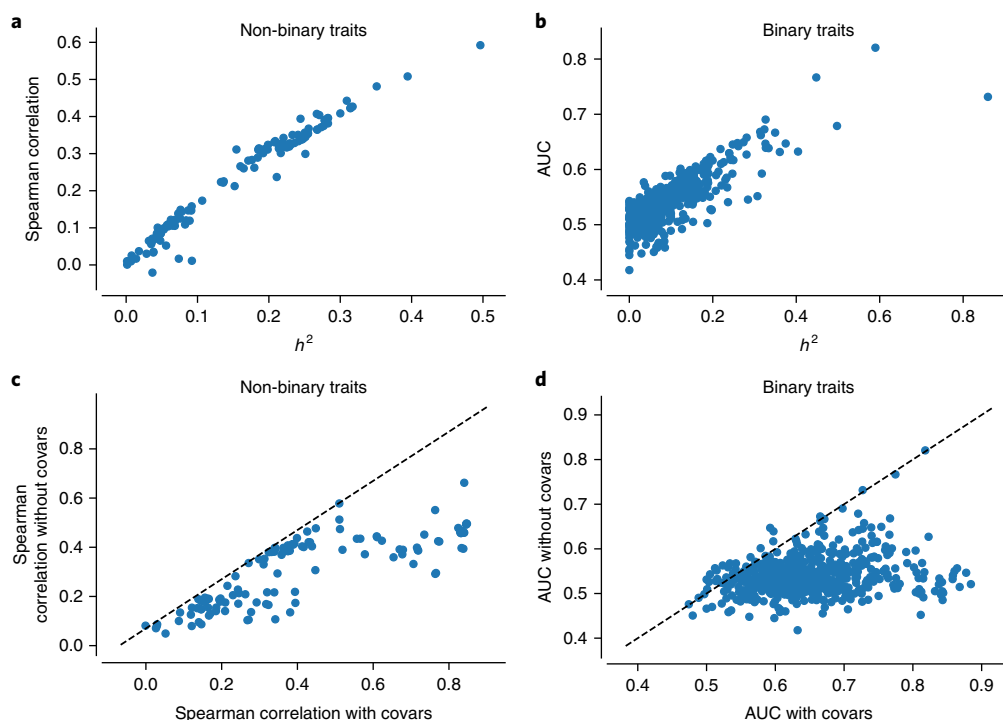


Fig. 7 | Phenotypic prediction accuracy from genetic markers. a,b, Accuracy of phenotypic prediction as a function of the estimated SNP-heritability for non-binary traits (**a**) and binary traits (**b**) when no covariates (covars) were used for prediction. **c,d,** Comparison between prediction accuracy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

did indeed reach large prediction accuracies (Fig. 7), which was further increased when we used additional covariates such as gender or sex. The largest prediction accuracy for a non-binary trait was for height, which was 0.59, while the largest discriminative ability for a binary trait was 0.82 for self-reported malabsorption/celiac disease. We observed a large correlation between the prediction accuracy and the trait heritability (Fig. 7 and Supplementary Table 8). Furthermore, we previously developed a model that predicted the benefit of having increasingly large training datasets for prediction of complex traits in UK Biobank^{11,12}. Our current accuracy of prediction for anthropomorphic traits is very similar to the ones we previously predicted we would achieve with this training set¹¹ (Supplementary Fig. 9).

Discussion

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

URLs. GeneATLAS, <http://geneatlas.roslin.ed.ac.uk/>; UK Biobank, <http://www.ukbiobank.ac.uk/>; ARCHERUK National Supercomputing

Service, <http://www.archer.ac.uk/>; DISSECT, <https://www.dissect.ed.ac.uk/>; GWAS catalog, <https://www.ebi.ac.uk/gwas/>; Affymetrix array, <https://affymetrix.app.box.com/s/6gc2mcw2s6a7zbb7wjn>; PLINK, <http://zzz.bwh.harvard.edu/plink/> and <http://www.cog-genomics.org/plink/1.9/>; BGENIX and BGEN reference implementation, <https://bitbucket.org/gavinband/bgen>.

Online content

Any methods, additional references, Nature Research reporting summaries, source data, statements of data availability and associated accession codes are available at <https://doi.org/10.1038/s41588-018-0248-z>.

Received: 18 September 2017; Accepted: 29 August 2018;

Published online: 22 October 2018

References

- Falconer, D. S. & Mackay, T. F. C. *Introduction to Quantitative Genetics* (Longman, Harlow, 1996).
- 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).
- 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).
- 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).
- McCarthy, S. et al. A reference panel of 64,976 haplotypes for genotype imputation. *Nat. Genet.* **48**, 1279–1283 (2016).
- 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).
- Yang, J. et al. Common SNPs explain a large proportion of the heritability for human height. *Nat. Genet.* **42**, 565–569 (2010).
- 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).

Acknowledgements

This research has been conducted using the UK Biobank Resource under project 788. The work was funded by the Roslin Institute Strategic Programme Grant from the BBSRC (BB/P013732/1) and MRC grant (MR/N003179/1) granted to A.T. A.T. also acknowledges funding from the Medical Research Council and O.C.-X. from MRC fellowship MR/R025851/1. Analyses were performed using the ARCHER UK National Supercomputing Service.

Author contributions

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

Competing interests

The authors declare no competing interests.

Additional information

Supplementary information is available for this paper at <https://doi.org/10.1038/s41588-018-0248-z>.

Reprints and permissions information is available at www.nature.com/reprints.

Correspondence and requests for materials should be addressed to A.T.

Publisher's note: Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

© The Author(s), under exclusive licence to Springer Nature America, Inc. 2018

Methods

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

Phenotypes. In total, we analyzed 778 phenotypes in UK Biobank participants of European ancestry. These included 657 binary phenotypes generated from self-reported disease status (UK Biobank field 20002), ICD10 codes from hospitalization events (UK Biobank fields 41202 and 41204), and ICD10 codes from cancer registries (UK Biobank fields 40006), as well as a further 3 binary and 118 non-binary (comprising continuous and ordered integral measures) phenotypes from across the UK Biobank. Among the 660 binary phenotypes, 86 exhibited either a complete lack of cases in one sex or a strong imbalance in prevalence in the two sexes, that is, the ratio between the smaller and larger prevalence was <0.02 . Of these 86 phenotypes, 72 were specific to women. We only included individuals of the appropriate sex, that is, the sex with higher prevalence, in the analysis of these sex-specific phenotypes. A description of each phenotype, its category and the relevant UK Biobank fields can be found in Supplementary Table 1 and the GeneAtlas website. The non-binary phenotypes were not scale transformed, so the units of the effect sizes are in the units reported in the UK Biobank database. The phenotypes for individuals with negative coding were replaced with the corresponding value (Supplementary Table 9). We also ordered the keys for the ordinal phenotypes with unordered keys in the UK Biobank database (Supplementary Table 10). The individuals with a phenotype departing 10 standard deviations from their gender mean were set as missing for traits with a value type defined as 'Integer' or 'Continuous' by UK Biobank. The exceptions to this were number of self-reported cancers (134-0.0), number of self-reported non-cancer illnesses (135-0.0), nucleated red blood cell percentage (30230-0.0), nucleated red blood cell count (30170-0.0), and frequency of solarium/sunlamp use (2277-0.0), which were left as reported by UK Biobank. Some of the traits analyzed have some redundancy that has been left for completeness. That is, some of these traits were measured in different ways during the study (for example, weight) or are analyzed as self-reported traits and clinical traits (for example, malabsorption). For disease traits, all individuals reporting a disease code were coded as cases, with all other individuals considered controls. Only non-disease phenotypes with missing data rate $<5\%$ were selected for analysis. For these phenotypes, missing values were imputed to the age- and sex-specific mean in the study cohort.

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

Genotypes. The genotypes of the UK Biobank participants were assayed using either of two genotyping arrays, the Affymetrix UK BiLEVE Axiom or Affymetrix UK Biobank Axiom array. These arrays were augmented by imputation of ~90 million genetic variants from the Haplotype Reference Consortium⁵, 1000 Genomes¹³ and UK10K¹³ projects. Full details regarding these data have been published elsewhere¹⁴.

We excluded individuals who were identified by the UK Biobank as outliers based on either genotyping missingness rate or heterogeneity, whose sex inferred from the genotypes did not match their self-reported sex, and who were not of European ancestry (based on both self-reported ethnicity and those from whom one of the two first genomic principal components did not fall within 5 standard deviations from the mean). Finally, we removed individuals with a missingness $>5\%$ across variants which passed our quality control procedure and those that have a missing phenotype for 40 or more traits. The resulting study cohort comprised 452,264 individuals.

From the genotyped data, we only retained biallelic autosomal variants that were assayed by both genotyping arrays employed by UK Biobank. We furthermore excluded variants which had failed UK Biobank quality control procedures in any of the genotyping batches. Additionally, for imputed and genotyped variants, we excluded variants with $P < 10^{-50}$ for departure from Hardy–Weinberg equilibrium, computed on a subset of 344,057 unrelated (Kinship coefficient <0.0442) individuals in the White British subset of the study cohort, and with a missingness

rate $>2\%$ in the study cohort. Although we analyzed all imputed variants and all genotyped variants with $MAF > 10^{-4}$ (all results available at the GeneAtlas website), only imputed variants with (minor allele frequency) $MAF > 10^{-3}$ in the study cohort and imputation score larger than 0.9 were used for the summary results presented here. This cut-off corresponds to less than 905 occurrences of the minor allele in the study cohort. We also filtered the HLA imputed alleles that were present in fewer than 10 individuals.

GWAS analysis. To test each genetic variant while taking into account population structure in UK Biobank (for example presence of related individuals or local structure), we used a LMM. Specifically, the model takes the form:

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

where \mathbf{y} is the vector of phenotypes, \mathbf{X} is the matrix of fixed effects, and $\boldsymbol{\beta}$ the effect sizes of these effects. We included as fixed effects sex, array batch, UK Biobank Assessment Center, age, age², and the leading 20 genomic principal components as computed by UK Biobank. \mathbf{g} is the polygenic effect that captures the population structure, fitted as a random effect. It follows the distribution $\mathbf{g} \sim \text{Normal}(\mathbf{0}, \mathbf{A}\sigma_g^2)$, with \mathbf{A} the Genomic Relationship Matrix (GRM), and σ_g^2 the variance explained by the additive genetic effects. The GRM was computed using common ($MAF > 5\%$) genotyped variants that passed quality control. Finally, $\boldsymbol{\epsilon} \sim \text{Normal}(\mathbf{0}, \mathbf{I}\sigma_e^2)$, with \mathbf{I} being the identity matrix, is a residual effect not accounted for by the fixed and random effects. Under this model, the phenotype vector \mathbf{y} , follows the distribution: $\text{Normal}(\mathbf{X}\boldsymbol{\beta}, \mathbf{A}\sigma_g^2 + \mathbf{I}\sigma_e^2)$.

Fitting one instance of such a LMM model is computationally very demanding. Following a naive approach, the required computational time increasing with the cube of the sample size (N), $\sim O(N^3)$, and the memory requirements with the square of the sample size, $\sim O(N^2)$. Consequently, fitting a single model on a cohort of the size of UK Biobank is challenging, and fitting millions of these models, one for each analyzed genetic variant and phenotype, is not feasible with standard computational and statistical approaches. To address this problem, we took advantage of three different tools. First, we used a large supercomputer, and DISSECT³ to speed up the calculations (for example, computing the GRM eigen-decomposition required 5,040 processor cores working together for ~10 h, and using ~5 TB of memory). Second, we computed the full eigen decomposition of the GRM, $\mathbf{A} = \boldsymbol{\Sigma}\boldsymbol{\Lambda}\boldsymbol{\Lambda}^T$, where $\boldsymbol{\Lambda}$ is the matrix of eigenvectors, and $\boldsymbol{\Sigma}$ is a diagonal matrix containing the eigenvalues. This allowed us to transform all the other model matrices, \mathbf{y} , \mathbf{X} , and $\boldsymbol{\epsilon}$ to the new space where the GRM is diagonal. Although the eigen-decomposition is a computationally intensive process, once diagonalized, the computational time of fitting a model is reduced considerably to $\sim O(N)$, thus enabling us to perform several tests using LMMs on a cohort of hundreds of thousands of individuals. Finally, we performed over 23 billion tests using a two-step approximation that optimizes the computational resources¹⁵. The first step of the approximation fits an LMM that adjusts by the relevant fix (for example, age, sex, and so on) and random effects (genetic effects) to each trait, the second step uses the residuals of LMM to test (two-tailed t -test on effect sizes) all available genetic markers for significance in a linear model. We corrected for the polygenic effect using a leave-one-chromosome-out (LOCO) approach¹⁶.

HLA region. We defined the HLA region as the region of chromosome 6 that spans base pairs 28,866,528 to 33,775,446. Throughout all analyses, we included 10 Mb either side of the above HLA region to account for LD with variants outside this region.

The imputed HLA alleles were tested using the same GWAS model described above, where the independent variable is the best guess allele reported dosage from the HLA imputed values (UK Biobank field 22182). We tested the alleles using two models: a model where the number of copies of each HLA allele for each locus was tested independently as a fixed effect, and a second model where the number of copies of all alleles in a given locus were tested together as fixed effects in the same model (that is, an omnibus test)¹⁷.

Estimation of genetic parameters. To estimate heritabilities and genetic correlations, we fitted LMMs for each trait with a GRM containing all common ($MAF > 5\%$) autosomal genetic variants that passed quality control. The heritability was estimated as $h_g^2 = \sigma_g^2 / (\sigma_g^2 + \sigma_e^2)$, where σ_g^2 and σ_e^2 are the estimates of the genetic and residual variance and the P values were obtained using a chi-squared test following the method described previously^{18,19}. For all binary outcomes, we transformed heritabilities on the observed scaled to the liability scale using the population prevalence of the disease. We provide sex-specific prevalences to allow sex-specific transformations (Supplementary Table 1). Using the model fits, we computed best linear unbiased predictor estimates of genetic additive values for each individual. The genetic correlations were estimated by computing correlations between these additive genetic values. Environmental correlations 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 phenotypic and genetic correlations for traits i and j .

Lead variants and independent loci. We clustered GWAS results into independent lead variants using the `--clump` option of the PLINK 1.9 software^{20,21}. Specifically,

for each trait individually, we clustered GWAS results by selecting genome-wide significant variants as lead variants and assigning to them unassigned variants within 10 Mb that have $P < 10^{-2}$ and $r^2 > 0.3$ with the lead variant. To compute the total number of independent loci across all traits, we performed the same clustering on the lead variants across all traits, choosing the lowest P -value for variants that were lead variants in different traits.

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

Phenotypic prediction. The effect of all common genetic variants ($MAF > 0.05$) were estimated together as a random effect using the model:

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

where μ is the mean term and e_i the residual for individual i . L is the number of fixed effects, x_{il} being the value for the fixed effect l at individual i and β_l the estimated effect of the fixed effect l . We fitted the same covariates as in the GWAS analyses. M is the number of markers and z_{ij} is the standardized genotype of individual i at marker j . The vector of effects of random common genetic variants \mathbf{a} is distributed as $\text{Normal}(0, I\sigma_a^2)$. The vector of environmental effects \mathbf{e} is distributed as $\text{Normal}(0, I\sigma_e^2)$. Defining $\sigma_g^2 = M\sigma_a^2$, the heritabilities were estimated as $\sigma_g^2 / (\sigma_e^2 + \sigma_g^2)$.

The prediction of the phenotype y_i for the individual i was computed as a sum of the product of the SNP effects and the number of reference alleles of the corresponding SNPs:

$\hat{y}_i = \sum_{j=1}^M \frac{(s_{ij} - \mu_j^*)}{\sigma_j^*} a_j$, where s_{ij} is the number of copies of the reference allele at marker j of individual i , M is the number of markers used for the prediction, and a_j the effect of marker j . μ_j^* and σ_j^* are the mean and the standard deviation of the effect allele in the training population.

We split European ancestry individuals into 407,669 genetically confirmed British individuals to train the models and the remaining 44,595 individuals to validate the models. We restricted this analysis to the 692 non-gender specific

phenotypes. Prediction accuracies for non-binary traits were computed as the Spearman correlation between the predicted and the real phenotype of individuals of European but non-British descent after correcting by the estimated effect of the used covariates. Prediction accuracies for binary traits were computed as the area under the curve of a receiver operating characteristic curve using the predicted and the real phenotypes of individuals of non-British descent.

Accession codes. This research has been conducted using the UK Biobank Resource under project 788.

Reporting Summary. Further information on research design is available in the Nature Research Reporting Summary linked to this article.

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

Data availability

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

References

- 1000 Genomes Project Consortium. et al. An integrated map of genetic variation from 1,092 human genomes. *Nature* **491**, 56–65 (2012).
- Bycroft, C. F. et al. Genome-wide genetic data on ~500,000 UK Biobank participants. Preprint at bioRxiv <https://doi.org/10.1101/166298> (2017).
- Aulchenko, Y. S., de Koning, D. J. & Haley, C. Genomewide rapid association using mixed model and regression: a fast and simple method for genomewide pedigree-based quantitative trait loci association analysis. *Genetics* **177**, 577–585 (2007).
- Yang, J., Zaitlen, N. A., Goddard, M. E., Visscher, P. M. & Price, A. L. Advantages and pitfalls in the application of mixed-model association methods. *Nat. Genet.* **46**, 100–106 (2014).
- Patsopoulos, N. A. et al. Fine-mapping the genetic association of the Major Histocompatibility Complex in multiple sclerosis: HLA and non-HLA Effects. *PLoS Genet.* **9**, e1003926 (2013).
- Stram, D. O. & Lee, J. W. Variance components testing in the longitudinal mixed effects model. *Biometrics* **50**, 6 (1994).
- Visscher, P. M. A note on the asymptotic distribution of likelihood ratio tests to test variance components. *Twin. Res. Hum. Genet.* **9**, 490–495 (2012).
- Purcell, S. et al. PLINK: a tool set for whole-genome association and population-based linkage analyses. *Am. J. Hum. Genet.* **81**, 559–575 (2007).
- Chang, C. C. et al. Second-generation PLINK: rising to the challenge of larger and richer datasets. *Gigascience* **4**, 1–16 (2015).

Life Sciences Reporting Summary

Nature Research wishes to improve the reproducibility of the work that we publish. This form is intended for publication with all accepted life science papers and provides structure for consistency and transparency in reporting. Every life science submission will use this form; some list items might not apply to an individual manuscript, but all fields must be completed for clarity.

For further information on the points included in this form, see [Reporting Life Sciences Research](#). For further information on Nature Research policies, including our [data availability policy](#), see [Authors & Referees](#) and the [Editorial Policy Checklist](#).

► Experimental design

1. Sample size

Describe how sample size was determined.

Our study is an analysis of the UK Biobank available data. We did not determine the sample size.

2. Data exclusions

Describe any data exclusions.

We excluded individuals who were identified by the UK Biobank as outliers based on either genotyping missingness rate or heterogeneity, whose sex inferred from the genotypes did not match their self-reported sex and who were not of white ancestry (based on both, self-reported ethnicity and those from whom one of the two first genomic principal components did not fall within 5 standard deviations from the mean). Finally, we removed individuals with a missingness >5% across variants which passed our quality control procedure and those that have a missing phenotype for 40 or more traits. The resulting study cohort comprised 452,264 individuals. The exclusion criteria was pre-established before performing the analyses following best GWAS practices.

3. Replication

Describe whether the experimental findings were reliably reproduced.

Experimental replication was not attempted.

4. Randomization

Describe how samples/organisms/participants were allocated into experimental groups.

Not relevant because our study was not experimental.

5. Blinding

Describe whether the investigators were blinded to group allocation during data collection and/or analysis.

The authors of this manuscript did not participate in the recruitment. The data collection, as described from UK Biobank:
Recruitment were via centrally coordinated identification and invitation from population-based registers (such as those held by the NHS) of potentially eligible people living within a reasonable travelling distance of an assessment centre (located around the UK). This central recruitment strategy will allow invitations to be targeted to enhance generalisability and to make allowance for the impact on participation rates of various factors (e.g. age, sex, ethnicity, socioeconomic status). Each assessment centre will aim to recruit as many as possible of the nearby target population during a period of about six months to one year (depending on the local population density and transport links), and will then be relocated in order to achieve recruitment across most of the UK.

Note: all studies involving animals and/or human research participants must disclose whether blinding and randomization were used.

6. Statistical parameters

For all figures and tables that use statistical methods, confirm that the following items are present in relevant figure legends (or in the Methods section if additional space is needed).

n/a Confirmed

- ☐ ☒ The exact sample size (n) for each experimental group/condition, given as a discrete number and unit of measurement (animals, litters, cultures, etc.)
- ☒ ☐ A description of how samples were collected, noting whether measurements were taken from distinct samples or whether the same sample was measured repeatedly
- ☒ ☐ A statement indicating how many times each experiment was replicated
- ☐ ☒ The statistical test(s) used and whether they are one- or two-sided (note: only common tests should be described solely by name; more complex techniques should be described in the Methods section)
- ☐ ☒ A description of any assumptions or corrections, such as an adjustment for multiple comparisons
- ☐ ☒ The test results (e.g. P values) given as exact values whenever possible and with confidence intervals noted
- ☐ ☒ A clear description of statistics including central tendency (e.g. median, mean) and variation (e.g. standard deviation, interquartile range)
- ☐ ☒ Clearly defined error bars

See the web collection on [statistics for biologists](#) for further resources and guidance.

► Software

Policy information about [availability of computer code](#)

7. Software

Describe the software used to analyze the data in this study.

We used DISSECT (v1.15.2c, May 24, 2018) for the main analyses, which is publicly available at <http://www.dissect.ed.ac.uk/> under GNU Lesser General Public License v3. We also used PLINK (v1.9 and v2.0, freely available online) and BGENIX (v1.0) freely available online (<https://bitbucket.org/gavinband/bgen>) for preparing the data.

For manuscripts utilizing custom algorithms or software that are central to the paper but not yet described in the published literature, software must be made available to editors and reviewers upon request. We strongly encourage code deposition in a community repository (e.g. GitHub). *Nature Methods* [guidance for providing algorithms and software for publication](#) provides further information on this topic.

► Materials and reagents

Policy information about [availability of materials](#)

8. Materials availability

Indicate whether there are restrictions on availability of unique materials or if these materials are only available for distribution by a for-profit company.

No unique materials were used in the study.

9. Antibodies

Describe the antibodies used and how they were validated for use in the system under study (i.e. assay and species).

No antibodies were used in the study.

10. Eukaryotic cell lines

a. State the source of each eukaryotic cell line used.

No eukaryotic cell lines were used in the study.

b. Describe the method of cell line authentication used.

No eukaryotic cell lines were used in the study.

c. Report whether the cell lines were tested for mycoplasma contamination.

No eukaryotic cell lines were used in the study.

d. If any of the cell lines used are listed in the database of commonly misidentified cell lines maintained by [ICLAC](#), provide a scientific rationale for their use.

No eukaryotic cell lines were used in the study.

► Animals and human research participants

Policy information about [studies involving animals](#); when reporting animal research, follow the [ARRIVE guidelines](#)

11. Description of research animals

Provide details on animals and/or animal-derived materials used in the study.

No animals were used in this study.

Policy information about [studies involving human research participants](#)

12. Description of human research participants

Describe the covariate-relevant population characteristics of the human research participants.

UK Biobank recruited ~500,000 people aged between 40-69 years in 2006-2010 from across the country to take part in this project. In our analysis we included the individuals of white ancestry, which were not genetic outliers, that passed the quality control tests.