Check for updates

# A cross-population atlas of genetic associations for 220 human phenotypes

Saori Sakaue <sup>(1),2,3,4,5,46</sup> <sup>(2)</sup>, Masahiro Kanai <sup>(1),5,6,7,8,9,46</sup>, Yosuke Tanigawa <sup>(1)</sup>, Juha Karjalainen<sup>5,6,7,9</sup>, Mitja Kurki<sup>5,6,7,9</sup>, Seizo Koshiba<sup>11,12</sup>, Akira Narita <sup>(1)</sup>, Takahiro Konuma<sup>1</sup>, Kenichi Yamamoto<sup>1,13,14</sup>, Masato Akiyama<sup>2,15</sup>, Kazuyoshi Ishigaki <sup>(2),2,4,5</sup>, Akari Suzuki<sup>16</sup>, Ken Suzuki <sup>(6)</sup>, Wataru Obara<sup>17</sup>, Ken Yamaji<sup>18</sup>, Kazuhisa Takahashi<sup>19</sup>, Satoshi Asai<sup>20,21</sup>, Yasuo Takahashi <sup>(2)</sup>, Takao Suzuki<sup>22</sup>, Nobuaki Shinozaki<sup>22</sup>, Hiroki Yamaguchi <sup>(2)</sup>, Shiro Minami<sup>24</sup>, Shigeo Murayama<sup>25</sup>, Kozo Yoshimori<sup>26</sup>, Satoshi Nagayama<sup>27</sup>, Daisuke Obata<sup>28</sup>, Masahiko Higashiyama<sup>29</sup>, Akihide Masumoto<sup>30</sup>, Yukihiro Koretsune<sup>31</sup>, FinnGen<sup>\*</sup>, Kaoru Ito <sup>(3)</sup>, Chikashi Terao <sup>(3)</sup>, Toshimasa Yamauchi<sup>33</sup>, Issei Komuro <sup>(3)</sup>, Takashi Kadowaki<sup>33,35</sup>, Gen Tamiya<sup>11,12,36,37</sup>, Masayuki Yamamoto <sup>(3)</sup>, Yusuke Nakamura<sup>38,39</sup>, Michiaki Kubo <sup>(4)</sup>, Yoshinori Murakami <sup>(3)</sup>, Mark J. Daly <sup>(5,6,7,9)</sup>, Koichi Matsuda <sup>(3)</sup> <sup>(4)</sup> <sup>(2)</sup> and Yukinori Okada <sup>(3)</sup>, <sup>(1)</sup>, <sup>(1)</sup>, <sup>(4)</sup>, <sup></sup>

Current genome-wide association studies do not yet capture sufficient diversity in populations and scope of phenotypes. To expand an atlas of genetic associations in non-European populations, we conducted 220 deep-phenotype genome-wide association studies (diseases, biomarkers and medication usage) in BioBank Japan (n = 179,000), by incorporating past medical history and text-mining of electronic medical records. Meta-analyses with the UK Biobank and FinnGen ( $n_{total} = 628,000$ ) identified ~5,000 new loci, which improved the resolution of the genomic map of human traits. This atlas elucidated the landscape of pleiotropy as represented by the major histocompatibility complex locus, where we conducted HLA fine-mapping. Finally, we performed statistical decomposition of matrices of phenome-wide summary statistics, and identified latent genetic components, which pinpointed responsible variants and biological mechanisms underlying current disease classifications across populations. The decomposed components enabled genetically informed subtyping of similar diseases (for example, allergic diseases). Our study suggests a potential avenue for hypothesis-free re-investigation of human diseases through genetics.

edical diagnosis has been shaped through description of organ dysfunctions and extraction of shared key symptoms, which categorizes a group of individuals into a specific disease to provide an optimal treatment. The earliest physicians in ancient Egypt empirically made disease diagnoses based on clinical symptoms, palpitation and auscultation (~2600 BC)<sup>1</sup>. Since then, physicians have refined the disease diagnosis by empirically categorizing the observed symptoms (for example, cough, sputum and fever) to describe underlying dysfunction (for example, pneumonia, a lung infection). An increased understanding of organ functions and the availability of diagnostic tests including biomarkers and imaging techniques have contributed to the current disease classifications, such as the International Classification of Diseases (ICD)<sup>2</sup> and phecode<sup>3</sup>.

In the past decades, genome-wide association studies (GWASs)<sup>4</sup> and phenome-wide association studies (PheWASs)<sup>5</sup> have provided novel insights into the biological basis underlying disease diagnoses. While disease pathogenesis is quite multifactorial, genetic underpinnings provide us with one way to independently assess the validity of historically defined disease classifications. To this end, a comprehensive catalog of disease genetics is warranted. While previous works have broadly contributed the catalog<sup>6</sup>, current genetic studies are still short of comprehensiveness in three ways: (1) population, in that the vast majority of GWASs have been predominated by European populations<sup>7</sup>; (2) scope of phenotypes, which were mostly limited to predetermined diseases on which participants' recruitment had been performed; and (3) a systematic method to interpret a plethora of summary results for understanding disease pathogenesis. We thus need to promote equity in genetic studies by sharing the results of genetic studies of deep phenotypes from diverse populations.

To expand the atlas of genetic associations, here we conducted 220 deep-phenotype GWASs (that is, 197 electronic medical records (EMR)-based disease and medication records and 23 biomarkers) in BioBank Japan (BBJ), including 108 phenotypes on which GWAS has never been conducted in East Asian populations. We then conducted GWASs for corresponding phenotypes in UK Biobank (UKB) and FinnGen, and performed cross-population meta-analyses ( $n_{total}$  = 628,000). We sought to elucidate the land-scape of pleiotropy and genetic correlation across diseases and populations. Furthermore, we applied DeGAs<sup>8</sup> to perform truncated singular-value decomposition (TSVD) on matrices of GWAS summary statistics of 159 diseases in Japanese and European ancestries, and derived latent components shared across the diseases.

A full list of affiliations appears at the end of the paper.

We interpreted the derived components by (1) functional annotation of genetic variants explaining the component, (2) identification of important cell types where the genes contributing to each component are specifically regulated and (3) projection of GWASs of biomarkers or metabolomes into the component space. The latent components recapitulated the hierarchy of current disease classifications, while different diseases sometimes converged on the same component which implicated shared biological pathways and relevant tissues. We classified a group of similar diseases (for example, allergic diseases) into subgroups based on these components. Analogous to the conventional classification of diseases structured by the shared symptoms, an atlas of genetic studies suggested the latent structure behind human diseases, which can elucidate the genetic variants, genes, organs and biological functions underlying human diseases.

#### Results

GWAS of 220 traits in BBJ and cross-population meta-analysis. An overview of this study is presented in Extended Data Fig. 1. BBJ is a nationwide biobank in Japan, and recruited participants based on the diagnosis of at least 1 of 47 target diseases (Supplementary Notes)<sup>9</sup>. Along with the target disease status, deep-phenotype data, such as past medical history (PMH), drug prescription records (~7 million), text data retrieved from EMR and biomarkers, have been collected. Beyond the collection of case samples based on the predetermined target diseases, the PMH and EMR have provided broader insights into disease genetics, as shown in recently launched biobanks such as UKB<sup>10</sup> and BioVU<sup>11</sup>. We therefore curated the PMH, performed text-mining of the EMR and merged them with 47 target disease statuses<sup>12</sup>. We created individual-level phenotypes on 159 disease endpoints (38 target diseases with median 1.25 times increase in case samples and 121 novel disease endpoints) and 23 categories of medication usage. We then systematically mapped the disease endpoints into (1) phecode<sup>3</sup>, which is a hierarchical grouping system of EHR-based disease codes to conduct PheWAS, and (2) ICD10 (ref.<sup>2</sup>), which is also a medical classification list by the World Health Organization (WHO) and widely used for billing purposes, to enable harmonized GWASs in UKB and FinnGen. We also analyzed a quantitative phenotype of 38 biomarkers in BBJ, of which individual phenotype data are available in UKB<sup>13</sup>. Using genotypes imputed with the 1000 Genomes Project Phase 3 data (n = 2,504) and population-specific whole-genome sequencing data (n=1,037) as a reference panel<sup>14</sup>, we conducted the GWASs of 159 binary disease endpoints, 38 biomarkers and 23 medication usages in ~179,000 individuals in BBJ (Fig. 1a-c and Supplementary Tables 1 and 2 for phenotype summary). To maximize statistical power, we used a linear mixed model implemented in SAIGE (v.0.37)<sup>15</sup> for binary traits and BOLT  $(v.2.3.4)^{16}$  for quantitative traits. By using linkage disequilibrium (LD)-score regression (LDSC)<sup>17</sup>, we confirmed that potential biases were controlled in the GWASs (Supplementary Table 3). In this expanded scope of GWASs in the Japanese population, we identified 519 genome-wide significant loci across 159 disease endpoints, 2,249 across 38 biomarkers, of which 113 and 281 loci were new, respectively ( $P < 5.0 \times 10^{-8}$ ; Methods and Supplementary Table 4). We conducted the initial medication-usage GWASs in East Asian populations and detected 215 genome-wide significant loci across 23 traits (Methods). These signals underscore the value of (1) conducting GWASs in non-Europeans and (2) expanding the scope of phenotypes by incorporating biobank resources such as PMH and EMR. For example, we detected an East Asian-specific variant, rs140780894, at the major histocompatibility complex (MHC) locus in pulmonary tuberculosis (PTB; odds ratio (OR) = 1.2,  $P = 2.9 \times 10^{-23}$ , minor allele frequency  $(MAF)_{FAS} = 0.24$ ; Extended Data Fig. 2a), which was not present in the European population (minor allele count  $(MAC)_{EUR} = 0$ )<sup>18</sup>. PTB is a serious global health burden and relatively endemic in Japan<sup>19</sup> (annual incidence

per 100,000 was 14 in Japan but 8 in the United Kingdom and 3 in the United States in 2018(ref. 20). Because PTB, an infectious disease, can be treatable and remittable, we substantially increased the number of cases by combining the participants with PMH of PTB to the patients with active PTB at the time of recruitment (from 549 (ref. 12) to 7,800 case individuals). We also identified new signals in common diseases that had not been target diseases but were included in the PMH record, such as rs715 at 3'UTR of CPS1 in cholelithiasis (Extended Data Fig. 2b; OR=0.87,  $P=9.6 \times 10^{-13}$ ) and rs2976397 at the PSCA locus in gastric ulcer, gastric cancer and gastric polyp (Extended Data Fig. 2c; OR = 0.86,  $P = 6.1 \times 10^{-24}$  for gastric ulcer). We detected pleiotropic functional variants, such as a deleterious missense variant, rs28362459 (p.Leu20Arg), in FUT3 associated with gall bladder polyp (OR = 1.46,  $P = 5.1 \times 10^{-11}$ ) and cholelithiasis (OR=1.11,  $P=7.3 \times 10^{-9}$ ; Extended Data Fig. 2d), and a splice donor variant, rs56043070 (c.89 + 1 G > A), causing loss of function of GCSAML associated with urticaria (OR=1.24,  $P=6.9\times10^{-12}$ ; Extended Data Fig. 2e), which was previously reported to be associated with platelet and reticulocyte counts<sup>4</sup>. Medication-usage GWASs also provided interesting signals as an alternative perspective for understanding disease genetics<sup>21</sup>. For example, individuals taking HMG CoA reductase inhibitors (C10AA in the Anatomical Therapeutic Chemical Classification (ATC)) were likely to harbor variations at HMGCR (lead variant at rs4704210, OR=1.11,  $P = 2.0 \times 10^{-27}$ ). Prescription of salicylic acids and derivatives (N02BA in ATC) was significantly associated with a rare East Asian missense variant in PCSK9, rs151193009 (p.Arg93Cys; OR=0.75,  $P = 7.1 \times 10^{-11}$ , MAF<sub>EAS</sub> = 0.0089, MAF<sub>EUR</sub> = 0.000; Extended Data Fig. 2f), which might indicate a strong protective effect against thromboembolic diseases in general.

We next conducted GWASs of corresponding phenotypes (that is, disease endpoints and biomarkers) that can be mapped in UKB and FinnGen (196 and 128 traits, respectively; Methods), and collected summary statistics of a medication-usage GWAS conducted in UKB<sup>21</sup> (23 traits; Supplementary Table 5). To confirm that the signals identified in BBJ were validated across populations, we systematically compared the effect sizes of the genome-wide significant variants in BBJ with those in a European dataset across binary and quantitative traits (Methods). The loci identified in BBJ GWASs were successfully validated in the same effect direction (2,171 of 2,305 (94.2%),  $P < 10^{-325}$  in sign test) and with high effect-size correlation (Extended Data Fig. 3). We also note that the genetic correlations encompassing genome-wide polygenic signals were generally high between BBJ and European GWASs (median  $\rho_{ge} = 0.82$ ; Supplementary Table 6 and Methods).

Motivated by the high replicability, we performed cross-population meta-analyses of these 220 harmonized phenotypes across three biobanks (Methods). We identified 1,730 disease-associated, 12,066 biomarker-associated and 1,018 medication-associated loci in total, of which 571, 4,471 and 301 were new, respectively (Fig. 1d and Supplementary Table 7). We note that when we strictly control for multiple testing burden by Bonferroni correction ( $P < 5.0 \times 10^{-8}$ / (220 phenotypes  $\times$  3 populations) = 7.6  $\times$  10<sup>-11</sup>), the number of significantly associated loci was 844, 7,309 and 500, respectively. All these summary statistics of GWASs are openly distributed through the PheWeb.jp website, with interactive Manhattan plots, LocusZoom plots and PheWAS plots based on the PheWeb platform<sup>22</sup>. Together, we successfully expanded the genomic map of human complex traits in terms of populations and scope of phenotypes through conducting deep-phenotype GWASs across global nationwide biobanks.

The regional landscape of pleiotropy. Because human traits are highly polygenic and the observed variations within the human genome are finite in number, pleiotropy, where a single variant affects multiple traits, is pervasive<sup>23</sup>. While pleiotropy has been intensively





8

**Fig. 1** Overview of the identified loci in the cross-population meta-analyses of 220 deep-phenotype GWASs. a-c, The pie charts describe the phenotypes analyzed in this study. The disease endpoints (**a**;  $n_{trait} = 159$ ) were categorized based on the ICD10 classifications (A to Z; Supplementary Table 1a), the biomarkers (**b**;  $n_{trait} = 38$ ; Supplementary Table 1b) were classified into nine categories and medication usage (**c**;  $n_{trait} = 23$ ) was categorized based on the ATC system (A to S; Supplementary Table 1c). **d**, The genome-wide significant loci identified in the cross-population meta-analyses and pleiotropic loci ( $P < 5.0 \times 10^{-8}$ ). The traits (rows) are sorted as shown in the pie chart, and each dot represents a significant locus in each trait. Pleiotropic loci are annotated by lines with a locus symbol.

studied in European populations by compiling previous GWASs<sup>23–26</sup>, the landscape of pleiotropy in non-European populations has been understudied. By leveraging this opportunity for comparing

the genetics of deep phenotypes across populations, we sought to investigate the landscape of regional pleiotropy in both Japanese and European populations. We defined the degree of pleiotropy as

the number of significant associations per variant  $(P < 5.0 \times 10^{-8})^{24}$ . In the Japanese, rs671, a missense variant at the ALDH2 locus, harbored the largest number of genome-wide significant associations (47 traits; Fig. 2a). Following this, rs2523559 at the MHC locus (24 traits) and rs1260326 at the GCKR locus (20 traits) were most pleiotropic. In Europeans, rs9265949 at the MHC locus harbored the largest number of genome-wide significant associations (46 traits; Fig. 2b), followed by rs7310615 at the ATXN2/SH2B3 locus (38 traits), rs1260326 at the GCKR locus (28 traits) and rs2519093 at the ABO locus (28 traits). We note that those pleiotropic loci were not affected when we adjusted for phenotypically closely correlated traits (Extended Data Figs. 4 and 5) or genetically closely related traits (Extended Data Fig. 5 and methods in the Supplementary Notes). Notably, the ALDH2 locus (pleiotropic in Japanese) and the MHC locus (pleiotropic in Japanese and Europeans) are known to be under recent positive selection<sup>27,28</sup>. To systematically assess whether pleiotropic regions in the genome were likely to be under selection pressure in each of the populations, we investigated the enrichment of the signatures of recent positive selection quantified by the metric singleton density score (SDS)<sup>27</sup> values within the pleiotropic loci (Methods). Intriguingly, when compared with those under the null hypothesis, we observed significantly higher values of SDS  $\chi^2$  values within the pleiotropic loci, and this fold change increased as the number of associations increased (that is, more pleiotropic) in both Japanese and Europeans (Fig. 2c,d). To summarize, the cross-population atlas of genetic associations elucidated the broadly shared landscape of pleiotropy, which implied a potential connection to natural selection signatures affecting diverse human populations.

Pleiotropic associations in MHC and ABO locus. Given the high degree of pleiotropy in both populations, we next sought to fine-map the pleiotropic signals within the MHC locus. To this end, we imputed the classical HLA alleles in BBJ and UKB, and performed association tests for 159 disease endpoints and 38 biomarkers (Fig. 3a,b). After the fine-mapping and conditional analyses (Methods), we identified 75 and 129 independent association signals in BBJ and UKB, respectively ( $P < 5.0 \times 10^{-8}$ ; Supplementary Table 8). Among 53 and 63 traits associated with MHC in BBJ and UKB, 2 and 9 traits had never been previously shown to be associated with MHC, respectively. Overall, HLA-B in class I and HLA-DRB1 in class II harbored the largest number of associations in both BBJ and UKB. For example, we fine-mapped the strong signal associated with PTB to HLA-DR $\beta$ 1 Ser57 (OR = 1.20,  $P = 7.1 \times 10^{-19}$ ) in BBJ. This is the third line of evidence showing the robust association of HLA with tuberculosis identified to date<sup>29,30</sup>, and we fine-mapped the signal to HLA-DRB1. Interestingly, HLA-DRβ1 at position 57 also showed pleiotropic associations with other autoimmune and thyroid-related diseases, such as Graves' disease, hyperthyroidism, Hashimoto's disease, Sjögren's syndrome, chronic hepatitis B and atopic dermatitis in BBJ. Of note, the effect direction of the association of HLA-DR $\beta$ 1 Ser57 was the same between hyperthyroid status (OR=1.29,  $P = 2.6 \times 10^{-14}$  in Graves' disease and OR = 1.37,  $P = 1.4 \times 10^{-8}$  in hyperthyroidism) and hypothyroid status (OR = 1.50,  $P = 9.0 \times 10^{-8}$ in Hashimoto's disease and OR = 1.31,  $P = 1.5 \times 10^{-7}$  in hypothyroidism), despite the opposite direction of thyroid hormone abnormality. This association of HLA-DRB1 was also observed in Sjögren's syndrome (OR = 2.04,  $P = 7.9 \times 10^{-12}$ ), which might underlie epidemiological comorbidities of these diseases<sup>31</sup>. Other novel associations in BBJ included HLA-DR $\beta$ 1 Asn197 with sarcoidosis (OR=2.07,  $P = 3.7 \times 10^{-8}$ ), and four independent signals with chronic sinusitis (that is, HLA-DRA, HLA-B, HLA-A and HLA-DQA1).

Another representative pleiotropic locus in the human genome is the *ABO* locus. We performed ABO blood-type PheWAS in BBJ and UKB (Fig. 3c,d). We estimated the ABO blood type from three variants (rs8176747, rs8176746 and rs8176719 at 9q34.2)<sup>32</sup>, and associated them with the risk of diseases and quantitative traits for each blood group. A variety of phenotypes, including common diseases such as myocardial infarction, as well as biomarkers such as blood cell traits and lipids, were associated with the blood types in both biobanks (Supplementary Table 9). We replicated an increased risk of gastric cancer in blood-type A as well as an increased risk of gastric ulcer in blood-type O in BBJ<sup>33</sup>.

Genetic correlation across populations. The interplay between polygenicity and pleiotropy suggests widespread genetic correlations among complex human traits<sup>34</sup>. Genetic relationships among human diseases have contributed to the refinement of disease classifications<sup>35</sup> and elucidation of the biology underlying the epidemiological comorbidity<sup>34</sup>. To obtain insights into the interconnections among human traits and compare them across populations, we computed pairwise genetic correlations  $(r_{a})$  across 106 traits (in Japanese) and 148 traits (in Europeans) with Z-score for  $h^{2}_{SNP}$  (that is, SNP heritability of the trait) > 2, using bivariate LDSC (Methods). We then defined the correlated trait domains by searching for the phenotype blocks with pairwise  $r_g > 0.7$  within 70% of  $r_g$ values in the block on the hierarchically clustered matrix of pairwise  $r_{g}$  values using a greedy algorithm (Methods and Extended Data Fig. 6). We detected domains of tightly correlated phenotypes, such as (1) cardiovascular- acting medications, (2) coronary artery diseases, (3) type 2 diabetes-related phenotypes, (4) allergy-related phenotypes and (5) blood cell phenotypes in BBJ (Extended Data Fig. 6a). These domains implicated the shared genetic backgrounds on the similar diseases and their treatments (for example, (2) diseases of the circulatory system in ICD10 and their treatments) and diagnostic biomarkers (for example, (3) glucose and hemoglobin A1C (HbA1c) in type 2 diabetes). Intriguingly, the corresponding trait domains were mostly identified in UKB as well (Extended Data Fig. 6b). We considered that the current clinical boundaries for human diseases broadly reflect the shared genetic etiology across populations, despite differences in populations and despite potential differences in diagnostic, environmental and prescription practices.

Deconvolution of GWAS statistics provides insights into biol**ogy.** A major challenge in genetic correlation is that the  $r_{g}$  is a scalar value between two traits, which collapses the correlation over the whole genome into an averaged metric<sup>36</sup>. This approach is not straightforward in specifying a set of genetic variants driving the observed correlation, which would pinpoint biological pathways explaining the shared pathogenesis. To address this, genetic association statistics of diverse phenotypes have implicated latent structures underlying genotype-phenotype associations without a prior hypothesis. In particular, matrix decomposition of GWAS statistics is a promising approach<sup>8,37,38</sup>, which derives orthogonal components that explain association variance across multiple traits. This decomposition addresses two challenges in current genetic correlation studies. First, it informs us of genetic variants that explain the shared structure across multiple diseases, thereby enabling functional interpretation of the component. Second, it can be applicable to subsignificant associations, which are important in understanding contribution of common variants in rare diseases<sup>37</sup> or in genetic studies in underrepresented populations where lower statistical power is inevitable.

Therefore, we applied DeGAs<sup>8</sup> on a matrix of the disease GWAS summary statistics in Japanese and European individuals ( $n_{\text{disease}}$ =159; Fig. 4a,b). To interpret the derived latent components, we annotated the genetic variants explaining each component (1) through GREAT genomic region ontology enrichment analysis<sup>39</sup>; (2) through identification of relevant cell types using tissue-specific regulatory DNA (ENCODE3 (ref. <sup>40</sup>)) and expression (GTEx<sup>41</sup>) profiles; and (3) by projecting biomarker GWASs in BBJ and UKB ( $n_{\text{biomarker}}$ =38) or metabolome GWASs in the



**Fig. 2** | Number of significant associations per variant. a,b, The Manhattan-like plots show the number of significant associations ( $P < 5.0 \times 10^{-8}$ ) at each tested genetic variant for all traits ( $n_{trait} = 220$ ) in Japanese (**a**) and in European GWASs (**b**). Loci with a large number of associations were annotated based on the closest genes of each variant. **c**,**d**, The plots indicate the fold change of the sum of SDS  $\chi^2$  within variants with a larger number of significant associations than a given number on the *x* axis compared with those under the null hypothesis in Japanese (**c**) and in Europeans (**d**). We also illustrated a regression line based on local polynomial regression fitting.

East Asian (EAS) cohort of the Tohoku Medical Megabank Organization (ToMMo; Methods) and the European (EUR) cohort<sup>42</sup> ( $n_{\text{metabolite}\_EAS}$  = 206,  $n_{\text{metabolite}\_EUR}$  = 248) into the component space derived from disease genetics (Fig. 4a). We applied TSVD on the sparse Z-score matrix of 22,980 variants, 159 phenotypes each in 2 populations (Japanese and Europeans), and derived 40 components that together explained 36.7% of the variance (Extended Data Fig. 7 and Supplementary Fig. 1).

Globally, similar diseases as defined by the conventional ICD10 classification were explained by the same components, based on DeGAs trait squared cosine score that quantifies component loadings8 (Fig. 4c,d and Extended Data Fig. 8). This would be considered as a hypothesis-free support of the historically defined disease classifications. For example, component 1 explained the genetic association patterns of diabetes (E10 and E11 in ICD10) and component 2 explained those of cardiac and vascular diseases (I00-183), in both populations. Functional annotation of the genetic variants explaining these components showed that component 1 (diabetes component) was associated with abnormal pancreas size (binomial  $P_{\text{enrichment}} = 7.7 \times 10^{-19}$ ) as a human phenotype, whereas component 2 (cardiovascular disease component) was associated with xanthelasma (that is, cholesterol accumulation on the eyelids; binomial  $P_{\text{enrichment}} = 3.0 \times 10^{-10}$ ). Further, the genes comprising component 1 were enriched in genes specifically expressed in the pancreas ( $P_{\text{enrichment}} = 5.5 \times 10^{-4}$ ), and those comprising component 2 were enriched in genes specifically expressed in the aorta  $(P_{\text{enrichment}} = 1.9 \times 10^{-3};$  Extended Data Fig. 9). By projecting the biomarker GWASs in BBJ and UKB and the metabolite GWASs in

independent cohorts of EAS and EUR into this component space, we observed that component 1 represented the genetics of glucose and HbA1c, and component 2 represented the genetics of blood pressure and lipids, which are biologically relevant. This deconvolution-projection analysis suggested the latent genetic structure behind human diseases, which recapitulated the underlying biological functions, relevant tissues and associated markers.

The latent components shared across diseases explained the convergent biology behind etiologically similar diseases. For example, component 10 explained the genetics of cholelithiasis (gall stone), cholecystitis (inflammation of gall bladder) and gall bladder polyp (Fig. 5a). The projection of publicly available EUR metabolite GWASs into the component space identified that component 10 represented the bilirubin metabolism pathway. Component 10 was composed of variants involved in intestinal cholesterol absorption in mouse phenotype (binomial  $P_{\text{enrichment}} = 3.8 \times 10^{-10}$ ). This is biologically relevant, since increased absorption of intestinal cholesterol is a major cause of cholelithiasis, which also causes cholecystitis<sup>43</sup>. The projection of the metabolite GWASs in an independent Japanese cohort of ToMMo showed the connection between component 10 and glycine that conjugates with bile acids<sup>44</sup>.

Some components were further utilized to interpret the underpowered GWAS with the use of the well-powered GWAS, and to identify the contributor of shared genetics between different diseases. For example, we complemented an underpowered varicose vein GWAS in BBJ ( $n_{case} = 474$ , genome-wide significant loci=0) with a more powered GWAS in Europeans ( $n_{case} = 22,037$ , genome-wide significant loci=70). Both GWASs were mostly

### **NATURE GENETICS**



**Fig. 3 | HLA and ABO association PheWAS. a,b**, Significantly associated HLA genes identified by HLA PheWAS in BBJ (**a**) or in UKB (**b**) are plotted ( $P < 5.0 \times 10^{-8}$ ). In addition to primary association signals of the phenotypes, independent associations identified by conditional analyses are also plotted, and the primary association is indicated by the plots with a gray border. The color of each plot indicates two-tailed *P* values calculated with logistic regression (for binary traits) or linear regression (for quantitative traits) as designated in the color bar at the bottom. The bars in green at the top indicate the number of significant associations per gene in each of the populations. The detailed allelic or amino acid position as well as statistics in the association are provided in Supplementary Table 8. **c,d**, Significant associations identified by ABO blood-type PheWAS in BBJ (**c**) or in UKB (**d**) are shown as boxes and colored based on the OR. The size of each box indicates two-tailed *P* values calculated with logistic regression (for binary traits) or linear regression (for quantitative traits).

represented by component 11, which was explained by variants related to abnormal vascular development (binomial  $P_{\text{enrichment}} = 4.2 \times 10^{-7}$ ; Fig. 5b). Another example is component 27, which was shared with rheumatoid arthritis and systemic lupus erythematosus, two distinct but representative autoimmune diseases. Component 27 was explained by the variants associated with interleukin secretion and plasma cell number (binomial  $P_{\text{enrichment}} = 6.1 \times 10^{-10}$  and  $9.3 \times 10^{-10}$ , respectively), and significantly enriched in the DNase I hypersensitive site (DHS) signature of lymphoid tissue ( $P_{\text{enrichment}} = 1.3 \times 10^{-4}$ ; Fig. 5c). This might suggest the convergent etiology of the two autoimmune diseases, which was not elucidated by the genetic correlation alone.

Finally, we aimed at hypothesis-free categorization of diseases based on these components. Historically, hypersensitivity reactions have been classified into four types (for example, types I to IV)<sup>45</sup>, but the clear subcategorization of allergic diseases based on this pathogenesis and whether the categorization can be achieved solely by genetics were unknown. In TSVD results, the allergic diseases (mostly J and L in ICD10) were represented by the four components 3, 16, 26 and 34. By combining these components as axis-1 (components 3 and 16) and axis-2 (components 26 and 34), and comparing the cumulative variance explained by these axes, we defined axis-1 dominant allergic diseases (for example, asthma and allergic rhinitis) and axis-2 dominant allergic diseases (metal allergy, contact dermatitis and atopic dermatitis; Fig. 5d). Intriguingly, the axis-1 dominant diseases etiologically corresponded to type I allergy (that is, immediate hypersensitivity). The variants explaining axis-1 were biologically related to IgE secretion and T helper



**Fig. 4 | The deconvolution analysis of a matrix of summary statistics of 159 diseases across populations. a**, An illustrative overview of deconvolution-projection analysis. Using DeGAs framework, a matrix of summary statistics from two populations (EUR: European and ASN: Asian (BBJ)) was decomposed into latent components, which were interpreted by annotation of a set of genetic variants driving each component and in the context of other GWASs through projection. b, A schematic representation of TSVD applied to decompose a summary statistic matrix **W** to derive latent components. **U, S** and **V** represent resulting matrices of singular values (**S**) and singular vectors (**U** and **V**). **c**, A heatmap representation of DeGAs squared cosine scores of diseases (columns) to components (rows). The components are shown from 1 (top) to 40 (bottom), and diseases are sorted based on the contribution of each component to the disease measured by the squared cosine score (from component 1 to 40). Full results with disease and component labels are in Extended Data Fig. 8. **d**, Results of TSVD of the disease genetics matrix and the projection of biomarker genetics. Diseases (left) and biomarkers (right) are colored based on the ICD10 classification and functional categorization, respectively. The derived components (middle; from 1 to 40) are colored alternately in blue or red. The squared cosine score of each disease with squared cosine score > 0.3 in at least one component are displayed. Anth, anthropometry; BC, blood cell; BP, blood pressure; Ele, electrolytes; Infl, inflammatory; Kidn, kidney-related; Liver, liver-related; Metab, metabolic; Prot, protein.

2 (Th<sub>2</sub>) cells (binomial  $P_{\text{enrichment}} = 9.9 \times 10^{-46}$  and  $2.9 \times 10^{-44}$ , respectively). Furthermore, GWAS of eosinophil count was projected onto axis-1, which recapitulated the biology of type I allergy<sup>46</sup>. By contrast, the axis-2 dominant diseases corresponded to type IV allergy (that is, cell-mediated delayed hypersensitivity). The variants explaining axis-2 were associated with interleukin-13 and interferon secretion (binomial  $P_{\text{enrichment}} = 1.6 \times 10^{-10}$  and  $5.2 \times 10^{-9}$ , respectively), and GWAS of C-reactive protein was projected onto axis-2, which was distinct from axis-1 (ref. <sup>47</sup>). To summarize, our deconvolution approach (1) recapitulated the existing disease classifications, (2) implicated underlying biological mechanisms and relevant tissues shared among related diseases and (3) suggested potential application for genetics-driven categorization of human diseases.

#### Discussion

Here, we performed 220 GWASs of human traits by incorporating the PMH and EMR data in BBJ, substantially expanding the atlas of genotype-phenotype associations in non-Europeans. We note that we additionally discovered 92 loci across 38 disease endpoints of which we had previously conducted GWASs in BBJ<sup>12</sup>, explaining an additional 0.21% of trait heritability on average in the liability threshold model, which highlights the value of curating PMH and EMR in biobanks. We then systematically compared their genetic basis with GWASs of corresponding phenotypes in Europeans. We confirmed the global replication of loci identified in BBJ, and discovered 5,343 new loci through cross-population meta-analyses. The results are openly shared through web resources, which will be a platform to accelerate further research such as functional follow-up studies and drug discovery<sup>48</sup>. Of note, leveraging these well-powered GWASs, we observed that the genes associated with endocrine/metabolic, circulatory and respiratory diseases (E, I and J by ICD10) were systematically enriched in targets of approved medications treating those diseases<sup>49</sup> (Supplementary Fig. 2). This should motivate us to use this expanded resource for genetics-driven novel drug discovery and drug repositioning.

The landscape of regional pleiotropy was globally shared across populations, and pleiotropic regions tended to have been under recent positive selection. One limitation of the current analysis is that we did not conduct statistical fine-mapping for every locus we identified, which might cause a concern over potential effects due to LD tagging for observed pleiotropy. Although we confirmed that the same pleiotropic variants were included within the 95% credible set for representative loci (*ALDH2* and *GCKR*; Supplementary Notes), more comprehensive statistical fine-mapping would further

### **NATURE GENETICS**



**Fig. 5 | Examples of disease-component correspondence and biological interpretation of the components by projection and enrichment analysis using GREAT.** Shown is a representative component explaining a group of diseases based on the contribution score, along with responsible genes, functional enrichment results by GREAT, relevant tissues and relevant biomarkers/metabolites. **a**, The functional annotation of gall bladder-related diseases and the component 10. **b**, The functional annotation of varicose vein and the component 11. **c**, The functional annotation of autoimmune diseases and the component 27. **d**, The characterization of allergic diseases based on the components 3, 16, 26 and 34. The red bars indicate the sum of squared cosine scores of components 26 and 34 (axis 2). We also performed functional characterization of those components by projection analysis and GREAT enrichment analysis. ASN, Asian; EUR, European; GB, gall bladder; RA, rheumatoid arthritis; SLE, systemic lupus erythematosus.

illuminate a global landscape of pleiotropic variants in the future. Moreover, elucidation of pleiotropy in other populations is warranted to replicate our results. Finally, to highlight the utility of deep-phenotype GWASs, we decomposed the cross-population genotype-phenotype association patterns by TSVD. The latent components derived from TSVD showed the convergent biological mechanisms and relevant cell types across diseases, which can be utilized for re-evaluation of existing disease classifications. The incorporation of biomarker and metabolome GWAS summary statistics enabled interpretation of the latent components. Our approach might suggest a potential avenue for restructuring of medical diagnoses through dissecting the shared genetic basis across a spectrum of diseases, as analogous to the current disease classifications historically and empirically shaped through categorization of key symptoms across a spectrum of organ dysfunctions. However, we note that one major challenge in the deconvolution analyses is that the derived components and their order are affected by selection of phenotypes analyzed as input matrices. Thus, external validations are necessary before being generalized and used for refinement of disease categorization.

In conclusion, our study substantially expanded the atlas of genetic associations, supported the historically defined categories of human diseases and should accelerate the discovery of the biological basis contributing to complex human diseases.

#### **Online content**

Any methods, additional references, Nature Research reporting summaries, source data, extended data, supplementary information, acknowledgements, peer review information; details of author contributions and competing interests; and statements of data and code availability are available at https://doi.org/10.1038/ s41588-021-00931-x.

Received: 20 October 2020; Accepted: 4 August 2021; Published online: 30 September 2021

#### References

- 1. Berger, D. A brief history of medical diagnosis and the birth of the clinical laboratory. Part 1—ancient times through the 19th century. *MLO Med. Lab. Obs.* **31**, 28–30 (1999).
- Organización Mundial de la Salud. International Statistical Classification of Diseases and Related Health Problems, 10th revision (ICD-10) (World Health Organization, 2016).
- Denny, J. C. et al. Systematic comparison of phenome-wide association study of electronic medical record data and genome-wide association study data. *Nat. Biotechnol.* 31, 1102–1110 (2013).
- Welter, D. et al. The NHGRI GWAS Catalog, a curated resource of SNP-trait associations. *Nucleic Acids Res.* 42, D1001–D1006 (2014).
- Denny, J. C. et al. PheWAS: demonstrating the feasibility of a phenomewide scan to discover gene-disease associations. *Bioinformatics* 26, 1205–1210 (2010).
- Claussnitzer, M. et al. A brief history of human disease genetics. *Nature* 577, 179–189 (2020).
- Martin, A. R. et al. Clinical use of current polygenic risk scores may exacerbate health disparities. *Nat. Genet.* 51, 584–591 (2019).
- Tanigawa, Y. et al. Components of genetic associations across 2,138 phenotypes in the UK Biobank highlight adipocyte biology. *Nat. Commun.* 10, (2019).
- Nagai, A. et al. Overview of the BioBank Japan project: study design and profile. J. Epidemiol. 27, S2–S8 (2017).
- Bycroft, C. et al. The UK Biobank resource with deep phenotyping and genomic data. *Nature* 562, 203–209 (2018).

### **NATURE GENETICS**

# ARTICLES

- Ritchie, M. D. et al. Robust replication of genotype-phenotype associations across multiple diseases in an electronic medical record. *Am. J. Hum. Genet.* 86, 560–572 (2010).
- Ishigaki, K. et al. Large-scale genome-wide association study in a Japanese population identifies novel susceptibility loci across different diseases. *Nat. Genet.* 52, 669–679 (2020).
- Kanai, M. et al. Genetic analysis of quantitative traits in the Japanese population links cell types to complex human diseases. *Nat. Genet.* 50, 390–400 (2018).
- 14. Akiyama, M. et al. Characterizing rare and low-frequency height-associated variants in the Japanese population. *Nat. Commun.* **10**, 4393 (2019).
- Zhou, W. et al. Efficiently controlling for case-control imbalance and sample relatedness in large-scale genetic association studies. *Nat. Genet.* 50, 1335–1341 (2018).
- Loh, P. R. et al. Efficient Bayesian mixed-model analysis increases association power in large cohorts. *Nat. Genet.* 47, 284–290 (2015).
- Bulik-Sullivan, B. et al. LD score regression distinguishes confounding from polygenicity in genome-wide association studies. *Nat. Genet.* 47, 291–295 (2015).
- Karczewski, K. J. et al. The mutational constraint spectrum quantified from variation in 141,456 humans. *Nature* 581, 434–443 (2020).
- Hagiya, H. et al. Trends in incidence and mortality of tuberculosis in Japan: a population-based study, 1997-2016. *Epidemiol. Infect.* 147, e38 (2019).
- WHO. Global Tuberculosis Report. https://apps.who.int/iris/bitstream/handle/ 10665/336069/9789240013131-eng.pdf (2020).
- Wu, Y. et al. Genome-wide association study of medication-use and associated disease in the UK Biobank. *Nat. Commun.* 10, 1891 (2019).
- Gagliano Taliun, S. A. et al. Exploring and visualizing large-scale genetic associations by using PheWeb. *Nat. Genet.* 52, 550–552 (2020).
- 23. Watanabe, K. et al. A global overview of pleiotropy and genetic architecture in complex traits. *Nat. Genet.* **51**, 1339–1348 (2019).
- Canela-Xandri, O., Rawlik, K. & Tenesa, A. An atlas of genetic associations in UK Biobank. *Nat. Genet.* 50, 1593–1599 (2018).
- 25. Pendergrass, S. A. et al. A phenome-wide association study (PheWAS) in the Population Architecture using Genomics and Epidemiology (PAGE) study reveals potential pleiotropy in African Americans. *PLoS ONE* 14, e0226771 (2019).
- Verma, A. et al. PheWAS and beyond: the landscape of associations with medical diagnoses and clinical measures across 38,662 individuals from Geisinger. Am. J. Hum. Genet. 102, 592–608 (2018).
- 27. Field,  $\bar{Y}$ . et al. Detection of human adaptation during the past 2000 years. *Science* **354**, 760–764 (2016).
- Okada, Y. et al. Deep whole-genome sequencing reveals recent selection signatures linked to evolution and disease risk of Japanese. *Nat. Commun.* 9, 1631 (2018).
- Qi, H. et al. Discovery of susceptibility loci associated with tuberculosis in Han Chinese. Hum. Mol. Genet. 26, 4752–4763 (2017).

- 30. Sveinbjornsson, G. et al. HLA class II sequence variants influence tuberculosis risk in populations of European ancestry. *Nat. Genet.* **48**, 318–322 (2016).
- Baldini, C., Ferro, F., Mosca, M., Fallahi, P. & Antonelli, A. The association of Sjögren syndrome and autoimmune thyroid disorders. *Front. Endocrinol.* 9, 121 (2018).
- 32. Nakao, M. et al. ABO blood group alleles and the risk of pancreatic cancer in a Japanese population. *Cancer Sci.* **102**, 1076–1080 (2011).
- Edgren, G. et al. Risk of gastric cancer and peptic ulcers in relation to ABO blood type: a cohort study. Am. J. Epidemiol. 172, 1280–1285 (2010).
- Bulik-Sullivan, B. et al. An atlas of genetic correlations across human diseases and traits. *Nat. Genet.* 47, 1236–1241 (2015).
- Anttila, V. et al. Analysis of shared heritability in common disorders of the brain. *Science* 360, eaap8757 (2018).
- Shi, H., Mancuso, N., Spendlove, S. & Pasaniuc, B. Local genetic correlation gives insights into the shared genetic architecture of complex traits. *Am. J. Hum. Genet.* **101**, 737–751 (2017).
- Burren, O. S. et al. Informed dimension reduction of clinically-related genome-wide association. Preprint at *bioRxiv* https://www.biorxiv.org/content /10.1101/2020.01.14.905869v3 (2020).
- Chasman, D. I., Giulianini, F., Demler, O. V. & Udler, M. S. Pleiotropy-based decomposition of genetic risk scores: association and interaction analysis for type 2 diabetes and CAD. Am. J. Hum. Genet. 106, 646–658 (2020).
- McLean, C. Y. et al. GREAT improves functional interpretation of cis-regulatory regions. Nat. Biotechnol. 28, 495–501 (2010).
- Meuleman, W. et al. Index and biological spectrum of human DNase I hypersensitive sites. *Nature* 584, 244–251 (2020).
- GTEx Consortium, F. et al. Genetic effects on gene expression across human tissues. Nature 550, 204–213 (2017).
- Shin, S. Y. et al. An atlas of genetic influences on human blood metabolites. Nat. Genet. 46, 543–550 (2014).
- Portincasa, P. & Wang, D. Q. H. Intestinal absorption, hepatic synthesis, and biliary secretion of cholesterol: where are we for cholesterol gallstone formation? *Hepatology* 55, 1313–1316 (2012).
- Vessey, D. A. The biochemical basis for the conjugation of bile acids with either glycine or taurine. *Biochem. J.* 174, 621–626 (1978).
- Coombs, R. R. A. & Gell, P. G. (eds) in *Clinical Aspects of Immunology* 317–337 (Blackwell Science, 1963).
- Stone, K. D., Prussin, C. & Metcalfe, D. D. IgE, mast cells, basophils, and eosinophils. J. Allergy Clin. Immunol. 125, S73 (2010).
- Kobayashi, K., Kaneda, K. & Kasama, T. Immunopathogenesis of delayed-type hypersensitivity. *Microsc. Res. Tech.* 53, 241–245 (2001).
- 48. Okada, Y. et al. Genetics of rheumatoid arthritis contributes to biology and drug discovery. *Nature* **506**, 376–381 (2014).
- Sakaue, S. & Okada, Y. GREP: Genome for REPositioning drugs. Bioinformatics 35, 3821–3823 (2019).

**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. 2021

<sup>1</sup>Department of Statistical Genetics, Osaka University Graduate School of Medicine, Suita, Japan. <sup>2</sup>Laboratory for Statistical and Translational Genetics, RIKEN Center for Integrative Medical Sciences, Yokohama, Japan. <sup>3</sup>Center for Data Sciences, Harvard Medical School, Boston, MA, USA. <sup>4</sup>Divisions of Genetics and Rheumatology, Department of Medicine, Brigham and Women's Hospital, Harvard Medical School, Boston, MA, USA. <sup>5</sup>Program in Medical and Population Genetics, Broad Institute of Harvard and MIT, Cambridge, MA, USA. <sup>6</sup>Analytic and Translational Genetics Unit, Massachusetts General Hospital, Boston, MA, USA. <sup>7</sup>Stanley Center for Psychiatric Research, Broad Institute of Harvard and MIT, Cambridge, MA, USA. <sup>8</sup>Department of Biomedical Informatics, Harvard Medical School, Boston, MA, USA. 9Institute for Molecular Medicine Finland (FIMM), University of Helsinki, Helsinki, Finland. <sup>10</sup>Department of Biomedical Data Science, School of Medicine, Stanford University, Stanford, CA, USA. <sup>11</sup>Tohoku Medical Megabank Organization, Tohoku University, Sendai, Japan. <sup>12</sup>Advanced Research Center for Innovations in Next-Generation Medicine (INGEM), Sendai, Japan. <sup>13</sup>Department of Pediatrics, Osaka University Graduate School of Medicine, Suita, Japan. <sup>14</sup>Laboratory of Statistical Immunology, Immunology Frontier Research Center (WPI-IFReC), Osaka University, Suita, Japan.<sup>15</sup>Department of Ocular Pathology and Imaging Science, Graduate School of Medical Sciences, Kyushu University, Fukuoka, Japan. <sup>16</sup>Laboratory for Autoimmune Diseases, RIKEN Center for Integrative Medical Sciences, Yokohama, Japan. <sup>17</sup>Department of Urology, Iwate Medical University, Iwate, Japan. <sup>18</sup>Department of Internal Medicine and Rheumatology, Juntendo University Graduate School of Medicine, Tokyo, Japan. <sup>19</sup>Department of Respiratory Medicine, Juntendo University Graduate School of Medicine, Tokyo, Japan. <sup>20</sup>Division of Pharmacology, Department of Biomedical Science, Nihon University School of Medicine, Tokyo, Japan. <sup>21</sup>Division of Genomic Epidemiology and Clinical Trials, Clinical Trials Research Center, Nihon University School of Medicine, Tokyo, Japan. 22 Tokushukai Group, Tokyo, Japan. 23 Department of Hematology, Nippon Medical School, Tokyo, Japan.<sup>24</sup>Department of Bioregulation, Nippon Medical School, Kawasaki, Japan.<sup>25</sup>Tokyo Metropolitan Geriatric Hospital and Institute of Gerontology, Tokyo, Japan. <sup>26</sup>Fukujuji Hospital, Japan Anti-Tuberculosis Association, Tokyo, Japan. <sup>27</sup>The Cancer Institute Hospital of the Japanese Foundation for Cancer Research, Tokyo, Japan.<sup>28</sup>Center for Clinical Research and Advanced Medicine, Shiga University of Medical Science, Otsu, Japan. <sup>29</sup>Department of General Thoracic Surgery, Osaka International Cancer Institute, Osaka, Japan. <sup>30</sup>Aso Iizuka Hospital, Fukuoka, Japan. <sup>31</sup>National Hospital Organization Osaka National Hospital, Osaka, Japan. <sup>32</sup>Laboratory for Cardiovascular Genomics and Informatics, RIKEN Center for Integrative Medical Sciences, Yokohama, Japan. <sup>33</sup>Department of Diabetes and Metabolic Diseases, Graduate School of Medicine, The University of Tokyo, Tokyo, Japan. <sup>34</sup>Department of Cardiovascular Medicine, Graduate School of Medicine, The University of Tokyo, Tokyo, Japan. <sup>35</sup>Toranomon Hospital,

Tokyo, Japan. <sup>36</sup>Graduate School of Medicine, Tohoku University, Sendai, Japan. <sup>37</sup>Center for Advanced Intelligence Project, RIKEN, Tokyo, Japan. <sup>38</sup>Human Genome Center, Institute of Medical Science, The University of Tokyo, Tokyo, Japan. <sup>39</sup>Cancer Precision Medicine Center, Japanese Foundation for Cancer Research, Tokyo, Japan. <sup>40</sup>RIKEN Center for Integrative Medical Sciences, Yokohama, Japan. <sup>41</sup>Division of Molecular Pathology, Institute of Medical Sciences, The University of Tokyo, Tokyo, Japan. <sup>41</sup>Division of Molecular Pathology, Institute of Medical Science, The University of Tokyo, Tokyo, Japan. <sup>42</sup>Laboratory of Complex Trait Genomics, Department of Computational Biology and Medical Sciences, Graduate School of Frontier Sciences, The University of Tokyo, Tokyo, Japan. <sup>43</sup>Psychiatric & Neurodevelopmental Genetics Unit, Department of Psychiatry, Analytic and Translational Genetics Unit, Department of Medical Sciences, Graduate school of Frontier Sciences, Unit, Department of Medical Sciences, Graduate school of Frontier Sciences, The University, Boston, MA, USA. <sup>44</sup>Department of Computational Biology and Medical Sciences, Graduate school of Frontier Sciences, The University, Suita, Japan. <sup>45</sup>Integrated Frontier Research for Medical Science Division, Institute for Open and Transdisciplinary Research Initiatives, Osaka University, Suita, Japan. <sup>46</sup>These authors contributed equally: Saori Sakaue, Masahiro Kanai. \*A list of authors and their affiliations appears at the end of the paper. <sup>Ka</sup>e-mail: ssakaue@bwh.harvard.edu; kmatsuda@edu.k.u-tokyo.ac.jp; yokada@sg.med.osaka-u.ac.jp

### FinnGen

### Juha Karjalainen<sup>5,6,7,9</sup>, Mitja Kurki<sup>5,6,7,9</sup>, Aarno Palotie<sup>5,9,43</sup> and Mark J. Daly<sup>5,6,7,9</sup>

A full list of members and their affiliations appears in the Supplementary Information.

#### Methods

GWAS of 220 traits in BBJ. All the participants provided written, informed consent approved by ethics committees of the Institute of Medical Sciences, the University of Tokyo and RIKEN Center for Integrative Medical Sciences. We conducted 220 deep-phenotype GWASs in BBJ. BBJ is a prospective biobank that collaboratively collected DNA and serum samples from 12 medical institutions in Japan and recruited approximately 200,000 participants, mainly of Japanese ancestry (Supplementary Notes). Mean age of participants at recruitment was 63.0 yr old, and 46.3% were female. All study participants had been diagnosed with one or more of 47 target diseases by physicians at the cooperating hospitals. We previously conducted GWASs of 42 of the 47 target diseases12. In this study, we curated the PMH records included in the clinical data, and performed text-mining to retrieve disease records from the free-format EMR as well. For disease phenotyping, the PMH record has already been curated and formatted as a sample × phenotype table. Regarding EMR, we searched for the Japanese term of a given disease diagnosis in the cells designated as the presence of PMH, which was compiled into a sample × phenotype table. We merged both pieces of information with the target disease status, and defined the case status for 159 diseases with a case count >50 (Supplementary Table 2). As controls, we used samples in the cohort without a given diagnosis or related diagnoses, which was systematically defined by using the phecode framework<sup>3</sup> (Supplementary Table 1). For medication-usage phenotyping, we again retrieved information by text-mining of 7,018,972 medication records. Then, we categorized each medication trade name by using the ATC, WHO, which is used for the classification of active ingredients of drugs according to the organ or system on which they act and their therapeutic, pharmacological and chemical properties. For biomarker phenotyping, we used the same processing and quality control method as previously described (Supplementary Table 2 for phenotype summary)13,50. In brief, we generally used the laboratory values measured at the participants' first visit to the recruitment center, and excluded measurements outside three times the interquartile range (IQR) of the upper/lower quartile across participants. For individuals taking anti-hypertensive medications, we added 15 mmHg to systolic blood pressure and 10 mmHg to diastolic blood pressure. For individuals taking a statin, we applied the following correction to the lipid measurements: (1) total cholesterol was divided by 0.8; (2) measured LDL-cholesterol (LDLC) was adjusted as LDLC/0.7; (3) derived LDLC from the Friedewald equation was re-derived as (total cholesterol/0.8) - HDL-cholesterol (HDLC) - (triglyceride/5).

We genotyped participants with the Illumina HumanOmniExpressExome BeadChip or a combination of the Illumina HumanOmniExpress and HumanExome BeadChip. Quality control of participants and genotypes was performed as described elsewhere<sup>14</sup>. In this project, we analyzed 178,726 participants of East Asian ancestry as estimated by the principal component analysis (PCA)-based sample selection criteria. The genotype data were further imputed with 1000 Genomes Project Phase 3 version 5 genotype data (n=2,504) and Japanese whole-genome sequencing data (n=1,037) using Minimac3 software. After this imputation, we excluded variants with an imputation quality of Rsq < 0.7, resulting in 13,530,797 variants analyzed in total.

We conducted GWASs for binary traits (that is, disease endpoints and medication usage) by using a generalized linear mixed model implemented in SAIGE (v.0.37), which had substantial advantages in terms of (1) maximizing the sample size by including genetically related participants and (2) controlling for case-control imbalance<sup>15</sup>, which was the case in many of the disease endpoints in this study. We included adjustments for age, age<sup>2</sup>, sex, age × sex, age<sup>2</sup> × sex and the top 20 principal components as covariates used in step 1. For sex-specific diseases, we alternatively adjusted for age, age<sup>2</sup> and the top 20 principal components as covariates used in step 1. For sex-specific diseases is specific. We conducted GWASs for quantitative traits (that is, biomarkers) by using a linear mixed model implemented in BOLT-LMM (v.2.3.4). We included the same covariates as used in the binary traits above.

Harmonized GWAS of 220 traits in UKB and FinnGen. We conducted the GWASs harmonized with BBJ in UKB and in FinnGen. The UKB project is a population-based prospective cohort that recruited approximately 500,000 people across the United Kingdom (Supplementary Notes). Mean age of participants at recruitment was 56.8 yr old, and 53.8% were female. We defined case and control status of 158 disease endpoints, which were originally retrieved from the clinical information in UKB and mapped to BBJ phenotypes via phecode (Supplementary Table 1). We also analyzed 38 biomarker values provided by the UKB. The genotyping was performed using the Applied Biosystems UK BiLEVE Axiom Array or the Applied Biosystems UK Biobank Axiom Array. The genotypes were further imputed using a combination of the Haplotype Reference Consortium, UK10K and 1000 Genomes Project Phase 3 reference panels by IMPUTE4 software<sup>10</sup>. In this study, we analyzed 361,194 individuals of white British genetic ancestry as determined by the PCA-based sample selection criteria (https://github. com/Nealelab/UK\_Biobank\_GWAS/blob/master/ukb31063\_eur\_selection.R). We excluded the variants with (1) an imputation information metric (INFO score)  $\leq$  0.8; (2) MAF  $\leq$  0.0001 (except for missense and protein-truncating variants annotated by VEP<sup>51</sup>, which were excluded if MAF  $\leq 1 \times 10^{-6}$ ); and (3)  $P_{\rm HWE} \le 1 \times 10^{-10}$ , resulting in 13,791,467 variants analyzed in total. We conducted

# ARTICLES

GWASs for 159 disease endpoints by using SAIGE with the same covariates used in the BBJ GWAS. For biomarker GWASs, we used publicly available summary statistics of UKB biomarker GWASs when available through Neale's lab website: http://www.nealelab.is/uk-biobank/ukbround2announcement, and otherwise performed linear regression using PLINK software with the same covariates, excluding the genetically related individuals (the first, second or third degree)<sup>10</sup>. For medication-usage GWASs, we used publicly available summary statistics of medication usage in UKB<sup>21</sup>, which was organized by the ATC and thus could be harmonized with BBJ GWASs.

FinnGen is a public-private partnership project combining genotype data from Finnish biobanks and digital health record data from Finnish health registries (Supplementary Notes). Mean age of participants at DNA sample collection was 51.8 yr old, and 56.3% were female. For GWASs, we used the summary statistics of FinnGen release 3 data (accessed through https://www.finngen.fi/en/access\_ results). The disease endpoints were mapped to BBJ phenotypes by using ICD10 codes, and we defined 128 of 159 endpoints in BBJ. The genome coordinates in summary statistics were lifted over to hg19, and we analyzed 16,859,359 variants after quality control. We did not conduct biomarker and medication-related GWASs because the availability of these phenotypes was limited.

# **Meta-analysis and annotation of the genome-wide significant variants.** First, we performed intraEuropean meta-analysis when summary statistics of both UKB and FinnGen were available, and then performed cross-population meta-analysis across three or two cohorts in 159 disease endpoints, 38 biomarker values and 23 medication-usage GWASs. We conducted these meta-analyses by using the

inverse-variance method and estimated heterogeneity with Cochran's Q test with METAL software (v.2011-03-25)<sup>52</sup>. In this meta-analysis, we included all variants after quality control in each of the three cohorts. The overlapping variants among the cohorts are summarized in Extended Data Fig. 10. The summary statistics of primary GWASs in BBJ and cross-population meta-analysis GWASs are openly shared without any restrictions.

We adopted the conventional genome-wide significance threshold of  $<5.0 \times 10^{-8}$ , as well as considering the Bonferroni-corrected threshold of  $<7.6 \times 10^{-11}$  ( $5.0 \times 10^{-8}/(220 \text{ phenotypes} \times 3 \text{ populations})$ ) in the context of cross-population meta-analysis. We defined independent genome-wide significant loci on the basis of genomic positions within  $\pm 500 \text{ kb}$  from the lead variant. We considered a trait-associated locus as novel when the locus within  $\pm 1$  megabase (Mb) from the lead variant did not include any variants that were previously reported to be significantly associated with the same disease.

To systematically collect previously reported significant associations  $(5.0 \times 10^{-8})$  as known variants, we (1) exhaustively searched for previous reports of genetic association in a given trait using the GWAS Catalog<sup>4</sup>, since it is currently recognized as a standard and most comprehensive database of genetic associations; (2) systematically searched PubMed when the corresponding trait was not included in the GWAS Catalog; and (3) exceptionally included preprints in case we have collaboratively worked on them, to avoid duplicated publication.

The goal was to comprehensively include only robust and invariant associations. In this way, we included 75,230 associations across 181 traits in 1,792 literatures as of 31 December 2020 (Supplementary Table 10).

We annotated the lead variants using ANNOVAR software, such as rsIDs in the dbSNP database (https://www.ncbi.nlm.nih.gov/snp/), the genomic region and closest genes, and functional consequences. We also supplemented this with the gnomAD database<sup>18</sup>, and also looked for the allele frequencies in global populations as an independent resource.

**Replication of significant associations in BBJ.** For 2,287 lead variants in the genome-wide significant loci of 159 disease endpoints and 38 biomarkers in BBJ, we compared the effect sizes and directions with European-only meta-analysis when available and with UKB-based summary statistics otherwise. Of them, 1,929 variants could be compared with the corresponding European GWASs. Thus, we performed the Pearson's correlation test for these variants' beta values in the association test in BBJ and in European GWASs. We also performed the correlation tests with variants with  $P_{\text{EUR}} < 0.05$  and with those with  $P_{\text{EUR}} < 5.0 \times 10^{-8}$ .

**Cross-population genetic correlation.** To estimate cross-population genetic-effect correlations between BBJ and European GWASs considering polygenic signals, we used Popcorn software  $(v.1.0)^{s3}$ . For this analysis, we restricted the traits to those with (1) heritability Z-score from LDSC>2 (which will be explained later in the Methods), and (2) both BBJ and European heritability calculated by Popcorn >0.01. We excluded the MHC region from the analysis because of its complex LD structure. Using these quality-controlled traits' summary statistics, we calculated the cross-population genetic-effect correlation between EUR and EAS with precomputed cross-population scores for EUR and EAS 1000 Genomes Project populations provided by the authors.

**Evaluation of regional pleiotropy.** We assessed the regional pleiotropy based on each tested genetic variant separately for BBJ GWASs and for European GWASs (that is, intraEuropean meta-analysis when FinnGen GWAS was available and UKB summary statistics otherwise). We quantified the degree of pleiotropy per genetic

variant by aggregating and counting the number of genome-wide significant associations across 220 traits. We then annotated loci from the largest number of associations ( $n_{\text{associations}} \ge 9$  in BBJ and  $\ge 18$  in Europeans) in Fig. 2a,b.

Next, we assessed the recent natural selection signature within the pleiotropic loci separately for Japanese and for Europeans. To do this, we first defined the pleiotropic loci by identifying genetic variants that harbored a larger number of significant associations than a given threshold. We varied this threshold from 1 to 40. Then, at each threshold, we calculated the sum of SDS  $\chi^2$  values within the pleiotropic loci, and compared this with the  $\chi^2$  distribution under the null hypothesis with a degree of freedom equal to the number of variants in the loci. We thus estimated the SDS enrichment within the pleiotropic loci defined by a given threshold as fold change and *P* value. The SDS values in UK10K were provided by the web resource at Pritchard's lab (http://web.stanford.edu/group/pritchardlab/UK10K-SDS-values.zip)<sup>27</sup> and provided by the authors on the Japanese population<sup>28</sup>. The raw SDS values were normalized according to the derived allele frequency as described previously<sup>27</sup>.

**Fine-mapping of the MHC region.** We performed the fine-mapping of MHC associations in BBJ and UKB by HLA imputation<sup>54</sup>. In BBJ, we imputed classical HLA alleles and corresponding amino acid sequences using the reference panel recently constructed from 1,120 individuals of Japanese ancestry by the combination of SNP2HLA software, Eagle and minimac3, as described previously<sup>55</sup>. We applied postimputation quality control to keep the imputed variants with MAF  $\geq$  0.5% and Rsq > 0.7. For each marker dosage that indicated the presence or absence of an investigated HLA allele or an amino acid sequence, we performed an association test with the disease endpoints and biomarkers. We assumed additive effects of the allele dosages on phenotypes in the regression models. We included the same covariates as in the GWAS. In UKB, we imputed classical HLA alleles and corresponding amino acid sequences using the T1DGC reference panel of European ancestry (n = 5,225)<sup>56</sup>. We applied the same postimputation quality control and performed the association tests as in BBJ.

**ABO blood group genotyping and analysis.** We extracted the best guess genotype of three variants (rs8176747, rs8176746 and rs8176719 at 9q34.2)<sup>32</sup>, and inferred blood group status of individuals in BBJ and UKB. We then performed logistic regression for 159 disease endpoints and linear regression for 38 biomarkers to test the association with the blood groups. Blood group-specific ORs or effect sizes (beta) were calculated by making four different groupings as A versus B/AB/O, B versus A/AB/O, AB versus A/B/O and O versus A/AB/B, as described elsewhere<sup>57</sup>. We described the traits with association  $P < 5 \times 10^{-8}$  in at least one of the blood groups in either of the cohorts in Fig. 3c,d.

Heritability and genetic correlation estimation. We performed LDSC by using LDSC software (v.1.0.1; https://github.com/bulik/ldsc) for GWASs of BBJ and Europeans to estimate SNP-based heritability, potential bias and pairwise genetic correlations. Variants in the MHC region (chromosome 6: 25-34 Mb) were excluded. We also excluded variants with  $\chi^2 > 80$ , as recommended previously For heritability estimation, we used the baselineLD model (v.2.2), which included 97 annotations that correct for bias in heritability estimates<sup>59</sup>. We note that we did not report liability-scale heritability, since the population prevalence of 159 diseases in each country was not always available, and the main objective of this analysis was an assessment of bias in GWAS, rather than the accurate estimation of heritability. We calculated the heritability Z-score to assess the reliability of heritability estimation, and reported the LDSC results with Z-score for  $h_{SNP}^2 > 2$  (Supplementary Table 3). For calculating pairwise genetic correlation, we again restricted the target GWASs to those whose Z-score for  $h_{SNP}^2$  is >2, as recommended previously58. In total, we calculated genetic correlation for 106 GWASs in BBJ and 148 European GWASs, which resulted in 5,565 and 10,878 trait pairs, respectively.

To illustrate trait-by-trait genetic correlation, we hierarchically clustered the  $r_{\rm g}$  values with hclust and colored them as a heatmap (Extended Data Fig. 6). To adopt reliable genetic correlations, we restricted the  $r_{\rm g}$  values that had  $P_{\rm cor} < 0.05$ . Otherwise, the  $r_{\rm g}$  values were replaced with 0. We then defined the tightly clustered trait domains by greedily searching for the phenotype blocks with pairwise  $r_{\rm g} > 0.7$  within 70% of  $r_{\rm g}$  values in the block from the top left of the clustered correlation matrix. We manually annotated each trait domain by extracting the characteristics of traits constituting the domain (Extended Data Fig. 6).

**Deconvolution of a matrix of summary statistics by TSVD.** We performed the TSVD on the matrix of genotype–phenotype association Z-scores as described previously as DeGAs framework<sup>8</sup>. In this study, we first focused on 159 disease endpoint GWASs in BBJ and European GWASs (that is, 318 in total) to derive latent components through TSVD. On constructing a Z-score matrix, we conducted variant-level quality control. We removed variants located in the MHC region (chromosome 6: 25–34 Mb), and replaced unreliable Z-score estimates with zero when one of the following conditions was satisfied as in Tanigawa et al.<sup>8</sup>: (1) *P* value of marginal association  $\geq 0.001$  or (2) standard error of beta value  $\geq 0.2$ . Considering that rows and columns with all zeros do not contribute to matrix decomposition, we excluded variants that had all zero Z-scores across 159 traits in

### NATURE GENETICS

either BBJ or Europeans. We then performed LD pruning using PLINK software<sup>60</sup> ('--indep-pairwise 50 5 0.1') with an LD reference of 5,000 randomly selected individuals of white British UKB participants to select LD-independent variant sets, which resulted in a total of 22,980 variants. Thus, we made a Z-score matrix (=**W**) with a size of 318 (*N*: 159 diseases × 2 populations) × 22,980 (*M*: variants). With a predetermined number of *K*, TSVD decomposed **W** into a product of three matrices: **U**, **S** and **V**<sup>T</sup>: **W** = **USV**<sup>T</sup>. **U** = ( $u_{i,k}$ )<sub>*i,k*</sub> is an orthonormal matrix of size *K*×*K* whose columns are phenotype singular vectors, **S** is a diagonal matrix of size *K*×*K* whose clements are singular values and **V** = ( $v_{j,k}$ )<sub>*i,k*</sub> is an orthonormal matrix. This value was determined by experimenting with different values from 20 to 100 and selecting the informative and sufficient threshold. We used the TruncatedSVD module in the sklearn.decomposition library of python for performing TSVD.

To interpret and visualize the results of TSVD, we calculated the squared cosine scores. The phenotype squared cosine score,  $\cos_i^{2^{phe}}(k)$ , is a metric to quantify the relative importance of the *k*th latent component for a given phenotype *i*, and is defined as follows:

$$\cos_{i}^{2^{phe}}\left(k\right) = \frac{\left(f_{i,k}^{p}\right)^{2}}{\sum_{k'}\left(f_{i,k'}^{p}\right)^{2}}$$

where

Annotation of the components by using GREAT and identification of relevant cell types. We calculated the variant contribution score, which is a metric to quantify the contribution of a given variant *j* to a given component *k*, as follows:

 $\mathbf{F}_p = \mathbf{U}\mathbf{S} = \left(f_{ik}^p\right)_{...}$ 

$$contr_{k}^{var}(j) = (v_{i,k})^{2}$$

For each component, we can thus rank the variants based on their contribution to the component and calculate the cumulative contribution score. We defined a set of contributing variants to a given component to include top-ranked variants that had high contribution scores until the cumulative contribution score to the component exceeded 0.5. For these variant sets contributing to the latent components, we performed the GREAT (v.4.0.4) binomial genomic region enrichment analysis<sup>19</sup> based on the size of the regulatory domain of genes and quantified the significance of enrichment in terms of binomial fold enrichment and binomial *P* value to biologically interpret these components. We used the human phenotype and mouse genome informatics phenotype ontology, which contains manually curated knowledge about the hierarchical structure of phenotypes and genotype-phenotype mapping of human and mouse, respectively. The enriched annotation with a false discovery rate < 0.05 is considered significant and displayed in the figures.

For a gene set associated with the contributing variants with a given component (P < 0.05), we sought to identify relevant cell types by integrating two datasets: (1) ENCODE3 DHS regulatory patterns across human tissues from non-negative matrix factorization<sup>40</sup> and (2) specifically expressed genes defined from GTEx data<sup>41</sup>. In brief, a vocabulary (that is, DHS patterns) for regulatory patterns was defined from the non-negative matrix factorization of 3 million DHSs × 733 human biosamples encompassing 438 cell and tissue types. Then, for each regulatory vocabulary, GENCODE genes were assigned based on their overlying DHSs. The gene labeling result was downloaded from the journal website<sup>40</sup>. We also defined genes specifically expressed in 53 tissues from GTEx version 7 data, based on the top 5% of the *t*-statistics in each tissue as described elsewhere<sup>61</sup>. Then, for (1) each regulatory vocabulary wocabulary and (2) each tissue, we performed Fisher's exact tests to investigate whether the genes associated with a given component are significantly enriched in the defined gene set.

**Projection of biomarker and metabolite GWASs into the component space.** To further help interpret the latent components derived from disease-based TSVD, we projected the Z-score matrix of biomarker GWASs and metabolite GWASs into the component space. Briefly, we constructed the Z-score matrices (W') of 38 biomarkers of BBJ and European GWASs (that is, 76 rows) and 248 known metabolites of independent previous GWASs in the European population (http://metabolomics.helmholtz-muenchen.de/gwas/index.php?task=download) (ref. <sup>42</sup>) ×22,980 variants (Supplementary Table 11). Then, using the V from the disease-based TSVD, we calculated the phenotype contribution as follows:

$$\mathbf{F}_{p}^{projection} = \mathbf{W}'\mathbf{V} = \left(f_{i,k}^{projection}\right)_{i,k}$$

We note that for metabolite GWASs, since the GWASs were imputed with the HapMap reference panel, we imputed Z-scores of missing variants using ssimp software<sup>62</sup> (v.0.5.5 –ref 1KG/EUR –impute.maf 0.01), and otherwise we set the missing Z-scores to zero.

### **NATURE GENETICS**

Projection of metabolite GWASs in Japanese into the component space. To investigate whether the projection analysis is applicable to independent datasets, we conducted metabolite GWASs in ToMMo. ToMMo is a community-based biobank that combines medical and genome information from the participants in the Tohoku region of Japan<sup>63</sup>. Detailed cohort description is presented in the Supplementary Notes. In this study, we analyzed a total of 206 metabolites<sup>64</sup> by proton nuclear magnetic resonance (NMR) or liquid chromatography-mass spectrometry (LC-MS) (Supplementary Table 12). For sample quality control, we excluded samples meeting any of the following criteria: (1) genotype call rate < 95%; (2) one individual from each pair of those in close genetic relation (PI\_HAT calculated by PLINK<sup>60</sup>  $\geq$  0.1875) based on call rate; and (3) outliers from Japanese ancestry clustering based on the PCA with samples of 1000 Genomes Project Phase 3 data. For phenotype quality control, we excluded (1) the measurements in pregnant women, (2) measurements that took time from sampling to biobanking  $\geq 2 d$  and (3) phenotypic outliers defined as log-transformed measurements laying more than 4 s.d. from the mean for each metabolite. The participants were genotyped with a custom SNP array for the Japanese population (that is, Japonica Array v.2). For genotype quality control, we excluded variants meeting any of the following criteria: (1) call rate < 98%, (2) *P* value for Hardy–Weinberg equilibrium  $< 1.0 \times 10^{-6}$  and (3) MAF < 0.01. The quality-controlled genotype data were prephased by using SHAPEIT2 software (r837), and imputed by using IMPUTE4 software (r300.3) with a combined reference panel of 1000 Genomes Project Phase 3 (n = 2,504) and population-specific WGS data (that is, 3.5KJPNv2; n = 3,552)<sup>64</sup>. After imputation, we excluded variants with imputation INFO < 0.7.

For GWASs, we obtained the residuals from a linear regression model of each of the log-transformed metabolites adjusted for age, age<sup>2</sup>, sex, time period from sampling to biobanking and top 20 genotype principal components. The residuals were then transformed by rank-based inverse normalization. Association analysis of imputed genotype dosage with the normalized residual of each metabolite was performed using PLINK2 software. We constructed the Z-score matrices (W') of the Japanese metabolites GWASs (that is, 206 rows)  $\times$  22,980 variants, in which we applied the same quality control to the Z-scores and set the missing Z-scores to zero again. We then performed the projection as described above.

**Drug target enrichment analysis.** To investigate whether disease-associated genes are systematically enriched in the targets of the approved drugs for the treatment of those diseases, the Genome for REPositioning drugs (GREP)<sup>49</sup> was used. A list of genes closest to the lead variants from GWAS, which was concatenated based on the alphabetical category of ICD10 (A to N), was used as an input gene set to test the enrichment for the target genes of approved drugs for diseases of a given ICD10 category.

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

#### Data availability

The genotype data of BBJ used in this study are available from the Japanese Genotype-phenotype Archive (JGA) with accession codes JGAS000114/J GAD000123 and JGAS000114/JGAD000220, which can be accessed through application at https://humandbs.biosciencedbc.jp/en/hum0014-latest. The UKB analysis was conducted via application number 47821. The genotype and phenotype data can be accessed through application at https://www.ukbiobank. ac.uk. This study used the FinnGen release 3 data. Summary results can be accessed through application at https://www.finngen.fi/en/access\_results. We provide downloadable full GWAS summary statistics with an interactive visualization of Manhattan, LocusZoom and PheWAS plots at our PheWeb.jp website (https:// pheweb.jp/). The summary statistics of GWASs in this study (BioBank Japan, European and cross-population meta-analyses) are also deposited at the National Bioscience Database Center (NBDC) Human Database (https://humandbs. biosciencedbc.jp/en/) with the accession code hum0197, and the GWAS Catalog (https://www.ebi.ac.uk/gwas/) with the study accession IDs from GCST90018563 (https://www.ebi.ac.uk/gwas/studies/GCST90018563) to GCST90019002 (https:// www.ebi.ac.uk/gwas/studies/GCST90019002) (full IDs are described in the Supplementary Notes). The summary statistics of metabolite GWASs in the Japanese population (Tohoku Medical Megabank Organization) which we used for decomposition-projection analysis are available at https://jmorp.megabank. tohoku.ac.jp/202102/gwas/TGA000005. We used gnomAD database (https:// gnomad.broadinstitute.org/) to refer to the allele frequencies.

#### Code availability

We used publicly available software for the analyses. The software used is listed and described in the Methods section of our manuscript.

#### References

 Sakaue, S. et al. Trans-biobank analysis with 676,000 individuals elucidates the association of polygenic risk scores of complex traits with human lifespan. *Nat. Med.* 26, 542–548 (2020).

- McLaren, W. et al. The Ensembl Variant Effect Predictor. *Genome Biol.* 17, 122 (2016).
- Willer, C. J., Li, Y. & Abecasis, G. R. METAL: fast and efficient meta-analysis of genomewide association scans. *Bioinformatics* 26, 2190–2191 (2010).
- Brown, B. C., Ye, C. J., Price, A. L. & Zaitlen, N. Transethnic genetic-correlation estimates from summary statistics. *Am. J. Hum. Genet.* 99, 76–88 (2016).
- 54. Raychaudhuri, S. et al. Five amino acids in three HLA proteins explain most of the association between MHC and seropositive rheumatoid arthritis. *Nat. Genet.* **44**, 291–296 (2012).
- Hirata, J. et al. Genetic and phenotypic landscape of the major histocompatibility complex region in the Japanese population. *Nat. Genet.* 51, 470–480 (2019).
- Jia, X. et al. Imputing amino acid polymorphisms in human leukocyte antigens. *PLoS ONE* 8, e64683 (2013).
- Severe Covid-19 GWAS Group et al. Genomewide association study of severe Covid-19 with respiratory failure. N. Engl. J. Med. 383, 1522–1534 (2020).
- Zheng, J. et al. LD Hub: a centralized database and web interface to perform LD score regression that maximizes the potential of summary level GWAS data for SNP heritability and genetic correlation analysis. *Bioinformatics* 33, 272–279 (2017).
- Gazal, S. et al. Linkage disequilibrium-dependent architecture of human complex traits shows action of negative selection. *Nat. Genet.* 49, 1421–1427 (2017).
- 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).
- Finucane, H. K. et al. Heritability enrichment of specifically expressed genes identifies disease-relevant tissues and cell types. *Nat. Genet.* 50, 621–629 (2018).
- Rüeger, S., McDaid, A. & Kutalik, Z. Evaluation and application of summary statistic imputation to discover new height-associated loci. *PLoS Genet.* 14, e1007371 (2018).
- Kuriyama, S. et al. The Tohoku Medical Megabank Project: design and mission. J. Epidemiol. 26, 493–511 (2016).
- Tadaka, S. et al. JMorp: Japanese Multi Omics Reference Panel. Nucleic Acids Res. 46, D551–D557 (2018).

#### Acknowledgements

We thank all the participants of BioBank Japan, UK Biobank and FinnGen. We thank K. Watanabe for her input in the analysis of phenotypic correlations and pleiotropy. This research was supported by the Tailor-Made Medical Treatment program (the BioBank Japan Project) of the Ministry of Education, Culture, Sports, Science, and Technology (MEXT), the Japan Agency for Medical Research and Development (AMED). The FinnGen project is funded by two grants from Business Finland (grant nos. HUS 4685/31/2016 and UH 4386/31/2016) and nine industry partners (AbbVie, AstraZeneca, Biogen, Celgene, Genentech, GSK, MSD, Pfizer and Sanofi). The following biobanks are acknowledged for collecting the FinnGen project samples: Auria Biobank (https:// www.auria.fi/biopankki/), THL Biobank (https://thl.fi/fi/web/thl-biopank), Helsinki Biobank (https://www.terveyskyla.fi/helsinginbiopankki/), Northern Finland Biobank Borealis (https://www.ppshp.fi/Tutkimus-ja-opetus/Biopankki), Finnish Clinical Biobank Tampere (https://www.tays.fi/biopankki), Biobank of Eastern Finland (https:// ita-suomenbiopankki.fi), Central Finland Biobank (https://www.ksshp.fi/fi-FI/Potilaalle/ Biopankki), Finnish Red Cross Blood Service Biobank (https://www.bloodservice.fi/ Research%20Projects/biobanking) and Terveystalo Biobank Finland (https://www. terveystalo.com/fi/Yritystietoa/Terveystalo-Biopankki/Biopankki/). S.S. was in part supported by the Mochida Memorial Foundation for Medical and Pharmaceutical Research, Kanae Foundation for the Promotion of Medical Science, Astellas Foundation for Research on Metabolic Disorders and the JCR Grant for Promoting Basic Rheumatology. M. Kanai was supported by a Nakajima Foundation Fellowship and the Masason Foundation. Y. Tanigawa is in part supported by a Funai Overseas Scholarship from the Funai Foundation for Information Technology and the Stanford University School of Medicine. M.A.R. is in part supported by the National Human Genome Research Institute (NHGRI) of the National Institutes of Health (NIH) under award no. R01HG010140, and an NIH Center for Multi- and Cross-population Mapping of Mendelian and Complex Diseases grant (no. 5U01 HG009080). The content is solely the responsibility of the authors and does not necessarily represent the official views of the NIH. Y.O. was supported by the Japan Society for the Promotion of Science (JSPS) KAKENHI (grant nos. 19H01021, 20K21834) and AMED (grant nos. JP21km0405211, JP21ek0109413, JP21ek0410075, JP21gm4010006 and JP21km0405217), JST Moonshot R&D Grant Number JPMJMS2021, Takeda Science Foundation and the Bioinformatics Initiative of Osaka University Graduate School of Medicine, Osaka University.

#### Author contributions

S.S., M. Kanai and Y.O. conceived the study. S.S., M. Kanai, Y. Tanigawa., M.A.R. and Y.O. wrote the manuscript. S.S., M. Kanai, J.K., M. Kurki, T. Konuma, Kenichi Yamamoto, M.A., K. Ishigaki, Kazuhiko Yamamoto, Y. Kamatani, A.P., M.J.D. and Y.O. conducted GWAS data analysis. S.S., Y. Tanigawa and M.A.R. conducted statistical decomposition

### NATURE GENETICS

analysis. S.S., S.K., A.N., G.T. and Y.O. conducted metabolome analysis. A.S., K.S., W.O., K. Yamaji, K.T., S.A., Y. Takahashi, T.S., N.S., H.Y., S. Minami, S. Murayama, K. Yoshimori, S.N., D.O., M.H., A.M., Y. Koretsune, K. Ito, C.T., T.Y., I.K., T. Kadowaki, M.Y., Y.N., M. Kubo, Y.M., Kazuhiko Yamamoto and K.M. collected and managed samples and data. A.P. and M.J.D. coordinated collaboration with FinnGen.

#### **Competing interests**

M.A.R. is on the SAB of 54Gene and the Computational Advisory Board for Goldfinch Bio and has advised BioMarin, Third Rock Ventures, MazeTx and Related Sciences. The funders had no role in study design, data collection and analysis, decision to publish or preparation of the manuscript. The remaining authors declare no competing interests.

#### **Additional information**

 $\label{eq:constraint} \textbf{Extended data} is available for this paper at https://doi.org/10.1038/s41588-021-00931-x.$ 

**Supplementary information** The online version contains supplementary material available at https://doi.org/10.1038/s41588-021-00931-x.

**Correspondence and requests for materials** should be addressed to Saori Sakaue, Koichi Matsuda or Yukinori Okada.

**Peer review information** *Nature Genetics* thanks Caroline Hayward, Marylyn Ritchie, and the other, anonymous, reviewer(s) for their contribution to the peer review of this work.

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

### **NATURE GENETICS**

### ARTICLES



**Extended Data Fig. 1 Overview of this study.** We performed 220 deep-phenotype GWASs in BioBank Japan, including 108 novel GWASs ever conducted in East Asian population. We performed trans-biobank meta-analyses with UK Biobank and FinnGen ( $n_{total} = 628,000$ ), resulting in discovery of 5,343 novel loci. All summary statistics are openly shared through pheweb.jp web portal. As downstream analyses, we performed (i) cross-population comparison of pleiotropy and genetic correlation, (ii) comprehensive HLA fine-mapping, and (iii) statistical decomposition of a matrix of summary statistics to gain insights into biology underlying current disease classifications, by incorporating functional genomics, metabolomics, and biomarker data.



Extended Data Fig. 2 | See next page for caption.



d





### **NATURE GENETICS**

# ARTICLES

**Extended Data Fig. 2 | Locus plots for representative loci. (a)** Regional association plots for Pulmonary Tuberculosis (PTB) in BBJ are shown. The lead variant (rs140780894) is colored in pink, and colors of other dots indicate linkage disequilibrium measure  $r^2$  with the lead variant. (b) Regional association plots for cholelithiasis in BBJ are shown. The lead variant (rs715) is colored in pink, and colors of other dots indicate linkage disequilibrium measure  $r^2$  with the lead variant. (c) Regional association plots for gastric diseases in BBJ at the *PSCA* locus in gastric ulcer, gastric cancer, and gastric polyp are shown. Rs2976397, which was a lead variant in gastric ulcer, is colored in pink, and colors of other dots indicate linkage disequilibrium measure  $r^2$  with the lead variant. (d) Regional association plots at the *FUT3* locus in gall bladder polyp and cholelithiasis in BBJ are shown. Rs28362459, which was a lead variant in gall bladder polyp and cholelithiasis in BBJ are shown. Rs28362459, which was a lead variant in gall bladder polyp and cholelithiasis in BBJ are shown. The lead variant (rs56043070) is colored in pink, and colors of other dots indicate linkage disequilibrium measure  $r^2$  with the lead variant. (f) Regional association plots for salicylic acids prescription in BBJ are shown. The lead variant (rs151193009) is colored in pink, and colors of other dots indicate linkage disequilibrium measure  $r^2$  with the lead variant (rs151193009) is colored in pink, and colors of other dots indicate linkage disequilibrium measure  $r^2$  with the lead variant (rs151193009) is colored in pink, and colors of other dots indicate linkage disequilibrium measure  $r^2$  with the lead variant (rs151193009) is colored in pink, and colors of other dots indicate linkage disequilibrium measure  $r^2$  with the lead variant.



**Extended Data Fig. 3 | The effect size correlation between BBJ GWAS and European GWAS.** The marginal effect sizes of genome-wide significant variants across traits in BBJ are compared with those in European GWAS. Each plot represents a variant, and is colored based on the significance in European GWAS as shown in the left top legend. Pearson's correlation *r* and *P* value (two-sided) between BBJ GWAS and European GWAS are also shown in the legend.

### **NATURE GENETICS**

a Phenotypic correlation matrix

### ARTICLES





**Extended Data Fig. 4 | Phenotypic correlation across 220 phenotypes in BBJ. a.** Heatmap of pair-wise phenotypic correlation matrix. The color of the cells indicates the value of correlation *r* as shown in a color scale at the bottom. The traits (rows and columns) were hierarchically clustered by hclust library in R. **b.** Silhouette score for clustering of closely related phenotypes with different number of clusters (Supplementary Notes).



**Extended Data Fig. 5 | The degree of pleiotropy in BBJ after accounting for phenotypic or genetic correlations.** The Manhattan-like plots show the number of significant associations ( $P < 5.0 \times 10^{-8}$ ) at each tested genetic variant in Japanese. **a**. For all traits ( $n_{trait} = 220$ ; as shown in Fig. 2a). **b**. After accounting for phenotypic correlations. **c**. After accounting for genetic correlations.

### a Japanese



Extended Data Fig. 6 | Genetic correlation matrices across populations. The matrices describe pairwise genetic correlation r<sub>g</sub> in Japanese GWAS (**a**; n = 5,565) and in European GWAS (**b**; n = 10,878), which was estimated by bivariate LD score regression. A color of the cells indicates the value of  $r_g$ as shown in a color scale at the bottom. The traits (rows and columns) were hierarchically clustered by hclust library in R, and trait domains are displayed as colored boxes (see Methods).

0



**Extended Data Fig. 7 | Network representation of the TSVD analysis.** Two-dimensional illustration of interconnection among 159 diseases and 40 latent components. Plots in blue indicate each trait's statistics, and plots in pink indicate the latent components derived by TSVD. White lines represent the contribution of each phenotype in each component. The width of the lines indicates the strength of the contribution based on the squared cosine score.



**Extended Data Fig. 8 | Heatmap representation of squared cosine scores of diseases to components.** The components (rows) are shown from 1 (top) to 40 (bottom), and the diseases (columns) are sorted based on the contribution of each component to the disease based on the squared cosine score (from component 1 to 40). Each cell is colored based on the squared cosine score of a given trait to a given component, as shown in a color scale at the bottom right.



Extended Data Fig. 9 | Enrichment analyses of genes explaining each component with tissue specificity. A heatmap representation of the enrichment analyses of genes explaining each component with tissue-specific genes defined by GTEx expression profile (a) and regulatory vocabulary from ENCODE3 data (b). Each cell is colored based on P<sub>enrichment</sub> from Fisher's exact tests to assess the enrichment of the genes comprising each component within each tissue-specific gene set as shown in a color scale at the bottom right.

### **b** ENCODE3



**Extended Data Fig. 10 | Genetic variants analyzed in the three cohorts.** The Venn diagram showing the number of genetic variants analyzed in this study in each of the three cohorts (BBJ, UKB, and FinnGen) and overlapping variants across the cohorts.

# nature research

Saori Sakaue, Koichi Matsuda, and Yukinori Corresponding author(s): Okada

Last updated by author(s): May 25, 2021

# **Reporting Summary**

Nature Research wishes to improve the reproducibility of the work that we publish. This form provides structure for consistency and transparency in reporting. For further information on Nature Research policies, see our <u>Editorial Policies</u> and the <u>Editorial Policy Checklist</u>.

### **Statistics**

For	all st	atistical analyses, confirm that the following items are present in the figure legend, table legend, main text, or Methods section.				
n/a	Cor	firmed				
		The exact sample size (n) for each experimental group/condition, given as a discrete number and unit of measurement				
	$\square$	A statement on whether measurements were taken from distinct samples or whether the same sample was measured repeatedly				
		The statistical test(s) used AND whether they are one- or two-sided Only common tests should be described solely by name; describe more complex techniques in the Methods section.				
		A description of all covariates tested				
		A description of any assumptions or corrections, such as tests of normality and adjustment for multiple comparisons				
		A full description of the statistical parameters including central tendency (e.g. means) or other basic estimates (e.g. regression coefficient) AND variation (e.g. standard deviation) or associated estimates of uncertainty (e.g. confidence intervals)				
		For null hypothesis testing, the test statistic (e.g. F, t, r) with confidence intervals, effect sizes, degrees of freedom and P value noted Give P values as exact values whenever suitable.				
$\boxtimes$		For Bayesian analysis, information on the choice of priors and Markov chain Monte Carlo settings				
		For hierarchical and complex designs, identification of the appropriate level for tests and full reporting of outcomes				
		Estimates of effect sizes (e.g. Cohen's d, Pearson's r), indicating how they were calculated				
	I	Our web collection on <u>statistics for biologists</u> contains articles on many of the points above.				

### Software and code

Policy information about <u>availability of computer code</u>								
Data collection	No software was used in data collection.							
Data analysis	We used publicly available software for the data analysis (SAIGE0.37, BOLT-LMM2.3.4, LDSC (1.0.1), R 3.6.3, plink 1.9 and 2.0, python3, Hail 0.2, Eagle2, Minimac3, IMPUTE4, METAL (released on 2011-03-25), ANNOVAR, Popcorn (version 1.0), SNP2HLA, FINEMAP version1.3.1, susieR version 0.8.1.0521, GREP, ssimp version 0.5.5, SHAPEIT2 (r837).							

For manuscripts utilizing custom algorithms or software that are central to the research but not yet described in published literature, software must be made available to editors and reviewers. We strongly encourage code deposition in a community repository (e.g. GitHub). See the Nature Research guidelines for submitting code & software for further information.

### Data

Policy information about availability of data

All manuscripts must include a data availability statement. This statement should provide the following information, where applicable:

- Accession codes, unique identifiers, or web links for publicly available datasets
- A list of figures that have associated raw data
- A description of any restrictions on data availability

The genotype data of BBJ used in this study are available from the Japanese Genotype-phenotype Archive (JGA) with accession code JGAS000114/JGAD000123 and JGAS000114/JGAD000220 which can be accessed through application at https://humandbs.biosciencedbc.jp/en/hum0014-latest. The UKB analysis was conducted via application number 47821. The genotype and phenotype data can be accessed through application at https://www.ukbiobank.ac.uk.This study used the FinnGen release 3 data. Summary results can be accessed through application at https://www.finngen.fi/en/access\_results. All summary statistics of 220 GWASs (BioBank Japan, European, and cross-population meta-analyses) are deposited at the National Bioscience Database Center (NBDC) Human Database with the accession code hum0197. We also provide an interactive visualization of Manhattan, Locus Zoom, and PheWAS plots with downloadable GWAS summary statistics at our

PheWeb.jp website [https://pheweb.jp/]. The summary statistics of metabolite GWASs in the Japanese population (Tohoku Medical Megabank Organization) are being prepared as a different project and in manuscript preparation (Koshiba et al.). The previous version of the partial statistics is publicly available at https:// jmorp.megabank.tohoku.ac.jp/202008/gwas/TGA000003.

# Field-specific reporting

Please select the one below that is the best fit for your research. If you are not sure, read the appropriate sections before making your selection.

Life sciences

Behavioural & social sciences Ecologi

Ecological, evolutionary & environmental sciences

For a reference copy of the document with all sections, see <u>nature.com/documents/nr-reporting-summary-flat.pdf</u>

# Life sciences study design

All studies must disclose on these points even when the disclosure is negative.

Sample size	All sample size in GWASs in this study is summarized in Supplementary Table 2 and 5. We did not perform sample size calculation but included the maximum number of individuals in each cohort who passed the QC threshold. This strategy maximizes the statistical power in each cohort and we also performed the cross-population meta-analysis to further increase the power.
Data exclusions	All samples were selected based on quality-control criteria in each cohort, which is summarized in Method section.
Replication	We compared all signal identified in BBJ GWASs with corresponding but independent GWASs in UK Biobank and FinnGen. We confirmed high replicability (directional concordance of effects= 94.2%, P<1E-325 in sign test).
Randomization	We did not need to use randomization in this study because this is a genotype-phenotype association study. All the samples with available accessibility to genotype and phenotype data were included in the analysis.
Blinding	We did not apply blinding of the samples because this is a genotype-phenotype association study and no intervention was conducted in our study.

# Reporting for specific materials, systems and methods

We require information from authors about some types of materials, experimental systems and methods used in many studies. Here, indicate whether each material, system or method listed is relevant to your study. If you are not sure if a list item applies to your research, read the appropriate section before selecting a response.

#### Materials & experimental systems

#### Methods

	, ,		
n/a	Involved in the study	n/a	Involved in the study
$\boxtimes$	Antibodies	$\boxtimes$	ChIP-seq
$\boxtimes$	Eukaryotic cell lines	$\boxtimes$	Flow cytometry
$\boxtimes$	Palaeontology and archaeology	$\boxtimes$	MRI-based neuroimaging
$\boxtimes$	Animals and other organisms		
	Human research participants		
$\boxtimes$	Clinical data		
$\boxtimes$	Dual use research of concern		

### Human research participants

Policy information about studies involving human research participants

Population characteristics	BBJ is a prospective biobank that collaboratively collected DNA and serum samples from 12 medical institutions in Japan and recruited approximately 200,000 participants, mainly of Japanese ancestry (Supplementary Note). Mean age of participants at recruitment was 63.0 years old, and 46.3% were female. All study participants had been diagnosed with one or more of 47 target diseases by physicians at the cooperating hospitals. We previously conducted GWASs of 42 out of the 47 target diseases. The UK Biobank project is a population-based prospective cohort that recruited approximately 500,000 people across the United Kingdom. Mean age of participants at recruitment was 56.8 years old, and 53.8% were female. FinnGen is a public–private partnership project combining genotype data from Finnish biobanks and digital health record data from Finnish health registries. Mean age of participants at DNA sample collection was 51.8 years old, and 56.3% were female.			
Recruitment	All study participants in BBJ had been diagnosed with one or more of 47 target diseases by physicians at the cooperating hospitals. Participants were registered to the cohort from June 2003 to March 2008, and their clinical information was collected annually via interviews and medical record reviews until 2013. The UK Biobank project recruited approximately 500,000 people aged between 40–69 years from 2006 to 2010 from across the United Kingdom. FinnGen is a public–private partnership project. Six regional and three country-wide Finnish biobanks participate in FinnGen. Additionally, data from previously established population and disease-based cohorts are utilized. Participants' health outcomes are followed up by			

linking to the national health registries (1969–2016), which collect information from birth to death. Each biobank has specific population context, and we also note that the BBJ is a hospital-based cohort, UKB is a populationbased cohort, and FinnGen is a mixture of them. The coherent results across three biobanks mitigated concerns over potential biases.

Ethics oversight

All the participants provided written informed consent approved from ethics committees of the Institute of Medical Sciences, the University of Tokyo and RIKEN Center for Integrative Medical Sciences.

Note that full information on the approval of the study protocol must also be provided in the manuscript.