1 A global atlas of genetic associations of 220 deep phenotypes

2

Saori Sakaue^{1-5, 46, *}, Masahiro Kanai^{1, 5-9, 46}, Yosuke Tanigawa¹⁰, Juha Karjalainen^{5-7,9}. Mitja 3 Kurki^{5-7,9}, Seizo Koshiba^{11,12}, Akira Narita¹¹, Takahiro Konuma¹, Kenichi Yamamoto^{1,13}, 4 Masato Akiyama^{2,14}, Kazuyoshi Ishigaki²⁻⁵, Akari Suzuki¹⁵, Ken Suzuki¹, Wataru Obara¹⁶, 5 Ken Yamaji¹⁷, Kazuhisa Takahashi¹⁸, Satoshi Asai^{19,20}, Yasuo Takahashi²¹, Takao 6 Suzuki²², Nobuaki Shinozaki²², Hiroki Yamaguchi²³, Shiro Minami²⁴, Shigeo Murayama²⁵, 7 Kozo Yoshimori²⁶, Satoshi Nagayama²⁷, Daisuke Obata²⁸, Masahiko Higashiyama²⁹, 8 Akihide Masumoto³⁰, Yukihiro Koretsune³¹, FinnGen, Kaoru Ito³², Chikashi Terao², 9 Toshimasa Yamauchi³³, Issei Komuro³⁴, Takashi Kadowaki³³, Gen Tamiya^{11,12,35,36}, 10 Masayuki Yamamoto^{11,12,35}, Yusuke Nakamura^{37,38}, Michiaki Kubo³⁹, Yoshinori Murakami⁴⁰, 11 Kazuhiko Yamamoto¹⁵, Yoichiro Kamatani^{2,41}, Aarno Palotie^{5,9,42}, Manuel A. Rivas¹⁰, Mark J. 12 Dalv^{5-7,9}. Koichi Matsuda^{43, *}. Yukinori Okada^{1,2,41,44,45, *} 13 14

Department of Statistical Genetics, Osaka University Graduate School of Medicine, Suita,
 Japan.

17 2. Laboratory for Statistical and Translational Genetics, RIKEN Center for Integrative

- 18 Medical Sciences, Yokohama, Japan
- 19 3. Center for Data Sciences, Harvard Medical School, Boston, MA, USA

20 4. Divisions of Genetics and Rheumatology, Department of Medicine, Brigham and
21 Women's Hospital, Harvard Medical School, Boston, MA, USA

- 5. Program in Medical and Population Genetics, Broad Institute of Harvard and MIT,Cambridge, MA, USA
- 6. Analytic and Translational Genetics Unit, Massachusetts General Hospital, Boston, MA,USA

26 7. Stanley Center for Psychiatric Research, Broad Institute of Harvard and MIT, Cambridge,

- 27 MA, USA
- 28 8. Department of Biomedical Informatics, Harvard Medical School, Boston, MA, USA
- 29 9. Institute for Molecular Medicine Finland (FIMM), University of Helsinki, Helsinki, Finland
- 30 10. Department of Biomedical Data Science, School of Medicine, Stanford University,
- 31 Stanford, CA, USA
- 32 11. Tohoku Medical Megabank Organization, Tohoku University, Sendai, Japan
- 33 12. The Advanced Research Center for Innovations in Next-Generation Medicine (INGEM),
- 34 Sendai, Japan

- 13. Department of Pediatrics, Osaka University Graduate School of Medicine, Suita, Japan.
- 36 14. Department of Ophthalmology, Graduate School of Medical Sciences, Kyushu University,
- 37 Fukuoka, Japan
- 38 15. Laboratory for Autoimmune Diseases, RIKEN Center for Integrative Medical Sciences,
- 39 Yokohama, Japan
- 40 16. Department of Urology, Iwate Medical University, Iwate, Japan
- 41 17. Department of Internal Medicine and Rheumatology, Juntendo University Graduate
- 42 School of Medicine, Tokyo, Japan
- 43 18. Department of Respiratory Medicine, Juntendo University Graduate School of Medicine,
- 44 Tokyo, Japan
- 45 19. Division of Pharmacology, Department of Biomedical Science, Nihon University School
- 46 of Medicine, Tokyo, Japan
- 47 20. Division of Genomic Epidemiology and Clinical Trials, Clinical Trials Research Center,
- 48 Nihon University School of Medicine, Tokyo, Japan
- 49 21. Division of Genomic Epidemiology and Clinical Trials, Clinical Trials Research Center,
- 50 Nihon University School of Medicine, Tokyo, Japan
- 51 22. Tokushukai Group, Tokyo, Japan
- 52 23. Department of Hematology, Nippon Medical School, Tokyo, Japan
- 53 24. Department of Bioregulation, Nippon Medical School, Kawasaki, Japan
- 54 25. Tokyo Metropolitan Geriatric Hospital and Institute of Gerontology, Tokyo, Japan
- 55 26. Fukujuji Hospital, Japan Anti-Tuberculosis Association, Tokyo, Japan
- 56 27. The Cancer Institute Hospital of the Japanese Foundation for Cancer Research, Tokyo,57 Japan
- 58 28. Center for Clinical Research and Advanced Medicine, Shiga University of Medical
- 59 Science, Otsu, Japan
- 60 29. Department of General Thoracic Surgery, Osaka International Cancer Institute, Osaka,
- 61 Japan
- 62 30. Aso lizuka Hospital, Fukuoka, Japan
- 63 31. National Hospital Organization Osaka National Hospital, Osaka, Japan
- 64 32. Laboratory for Cardiovascular Genomics and Informatics, RIKEN Center for Integrative
- 65 Medical Sciences, Yokohama, Japan
- 66 33. Department of Diabetes and Metabolic Diseases, Graduate School of Medicine, The
- 67 University of Tokyo, Tokyo, Japan
- 68 34. Department of Cardiovascular Medicine, Graduate School of Medicine, The University of
- 69 Tokyo, Tokyo, Japan

- 70 35. Graduate School of Medicine, Tohoku University, Sendai, Japan
- 71 36. Center for Advanced Intelligence Project, RIKEN, Tokyo, Japan
- 72 37. Human Genome Center, Institute of Medical Science, The University of Tokyo, Tokyo,
- 73 Japan
- 38. Cancer Precision Medicine Center, Japanese Foundation for Cancer Research, Tokyo,
- 75 Japan
- 76 39. RIKEN Center for Integrative Medical Sciences, Yokohama, Japan
- 40. Division of Molecular Pathology, Institute of Medical Science, The University of Tokyo,
- 78 Tokyo, Japan
- 79 41. Laboratory of Complex Trait Genomics, Department of Computational Biology and
- 80 Medical Sciences, Graduate School of Frontier Sciences, The University of Tokyo, Tokyo,
- 81 Japan
- 42. Psychiatric & Neurodevelopmental Genetics Unit, Department of Psychiatry, Analytic
- 83 and Translational Genetics Unit, Department of Medicine, and the Department of Neurology,
- 84 Massachusetts General Hospital, Boston, MA, USA
- 43. Department of Computational Biology and Medical Sciences, Graduate school of
- 86 Frontier Sciences, the University of Tokyo, Tokyo, Japan
- 87 44. Laboratory of Statistical Immunology, Immunology Frontier Research Center
- 88 (WPI-IFReC), Osaka University, Suita, Japan
- 89 45. Integrated Frontier Research for Medical Science Division, Institute for Open and
- 90 Transdisciplinary Research Initiatives, Osaka University, Suita, Japan
- 91 46. These authors contributed equally: S Sakaue and M Kanai.
- 92 * Corresponding authors

Saori Sakaue, M.D., Ph.D. ssakaue@bwh.harvard.edu Center for Data Sciences, Harvard Medical School Koichi Matsuda, M.D., Ph.D. kmatsuda@edu.k.u-tokyo.ac.jp Department of Computational Biology and Medical Sciences, Graduate school of Frontier Sciences, The University of Tokyo Yukinori Okada, M.D., Ph.D. yokada@sg.med.osaka-u.ac.jp Department of Statistical Genetics, Osaka University Graduate School of Medicine

94 Abstract

95 The current genome-wide association studies (GWASs) do not yet capture sufficient 96 diversity in terms of populations and scope of phenotypes. To address an essential need to 97 expand an atlas of genetic associations in non-European populations, we conducted 220 98 deep-phenotype GWASs (disease endpoints, biomarkers, and medication usage) in 99 BioBank Japan (n = 179,000), by incorporating past medical history and text-mining results 100 of electronic medical records. Meta-analyses with the harmonized phenotypes in the UK 101 Biobank and FinnGen ($n_{\text{total}} = 628,000$) identified over 4,000 novel loci, which substantially 102 deepened the resolution of the genomic map of human traits, benefited from East Asian 103 endemic diseases and East Asian specific variants. This atlas elucidated the globally shared 104 landscape of pleiotropy as represented by the MHC locus, where we conducted 105 fine-mapping by HLA imputation. Finally, to intensify the value of deep-phenotype GWASs, 106 we performed statistical decomposition of matrices of phenome-wide summary statistics. 107 and identified the latent genetic components, which pinpointed the responsible variants and 108 shared biological mechanisms underlying current disease classifications across populations. 109 The decomposed components enabled genetically informed subtyping of similar diseases 110 (e.g., allergic diseases). Our study suggests a potential avenue for hypothesis-free 111 re-investigation of human disease classifications through genetics.

113 **Main**

114 Introduction

115 Medical diagnosis has been shaped through the description of organ dysfunctions and 116 extraction of shared key symptoms, which categorizes a group of individuals into a specific 117 disease to provide an optimal treatment. The earliest physicians in ancient Egypt empirically 118 made disease diagnoses based on clinical symptoms, palpitation, and auscultation (~2600 119 BC)¹. Since then, continuous efforts by physicians have sophisticated the disease 120 classifications through empirical categorization. An increased understanding of organ 121 functions and the availability of diagnostic tests including biomarkers and imaging techniques have further contributed to the current disease classifications, such as ICD10² 122 123 and phecode³.

124 In the past decades, genome-wide association studies (GWASs) have provided new 125 insights into the biological basis underlying disease diagnoses. The genetic underpinnings 126 enable us to re-interrogate the validity of historically- defined disease classifications. To this 127 end, a comprehensive catalog of disease genetics is warranted. However, current genetic 128 studies still lack the comprehensiveness in three ways; (i) population, in that the vast majority of GWASs have been predominated by European populations⁴, (ii) scope of 129 130 phenotypes, which have been limited to target diseases of a sampling cohort, and (iii) a 131 systematic method to interpret a plethora of summary results for understanding disease 132 pathogenesis and epidemiology. We thus need to promote equity in genetic studies by 133 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 in BioBank Japan project (BBJ), including 108 novel phenotypes in East Asian

136 populations. We then conducted GWASs for corresponding harmonized phenotypes in UK 137 Biobank (UKB) and FinnGen, and finally performed trans-ethnic meta-analyses ($n_{\text{total}} =$ 138 628,000). The association results elucidated trans-ethnically shared landscape of the 139 pleiotropy and genetic correlations across diseases. Furthermore, we applied DeGAs⁵ to 140 perform truncated singular-value decomposition (TSVD) on the matrix of GWAS summary 141 statistics of 159 diseases each in Japanese and European ancestries, and derived latent 142 components shared across the diseases. We interpreted the derived components by (i) 143 functional annotation of the genetic variants explaining the component, (ii) identification of 144 important cell types in which the genes contributing to the component are specifically 145 regulated, and (iii) projection of GWASs of biomarkers or metabolomes into the component 146 space. The latent components recapitulated the hierarchy of current disease classifications, 147 while different diseases sometimes converged on the same component which implicated the 148 shared biological pathway and relevant tissues. We also classified a group of similar 149 diseases (e.g., allergic diseases) into subgroups based on these components. Analogous to 150 the conventional hierarchical classification of diseases based on the shared symptoms, an 151 atlas of genetic studies resolved the shared latent structure behind human diseases, which 152 elucidated the genetic variants, genes, organs, and biological functions underlying human 153 diseases.

154

155 Results

156 GWAS of 220 traits in BBJ and trans-ethnic meta-analysis

157 Overview of this study is presented in **Extended Data Figure 1**. BBJ is a nationwide 158 biobank in Japan, and recruited participants based on the diagnosis of at least one of 47 target diseases (Supplementary Note)⁶. Along with the target disease status, deep 159 160 phenotype data, such as past medical history (PMH), drug prescription records (~ 7 million), 161 text data retrieved from electronic medical records (EMR), and biomarkers, have been 162 collected. Beyond the collection of case samples based on the pre-determined target 163 diseases, the PMH and EMR have provided broader insights into disease genetics, as shown in recently launched biobanks such as UKB⁷ and BioVU⁸. In this context, we curated 164 165 the PMH, performed text-mining of the EMR, and merged them with 47 target disease status.⁹ We created individual-level phenotype on 159 disease endpoints (38 target 166 167 diseases with median 1.25 times increase in case samples and 121 novel disease 168 endpoints) and 23 categories of medication usage. We then systematically mapped the 169 disease endpoints into phecode and ICD10, to enable harmonized GWASs in UKB and 170 FinnGen. We also analyzed a quantitative phenotype of 38 biomarkers in BBJ, of which 171 individual phenotype data are available in UKB¹⁰. Using genotypes imputed with the 1000 172 Genome Project phase 3 data (n = 2,504) and population-specific whole-genome 173 sequencing data (n = 1,037) as a reference panel¹¹, we conducted the GWASs of 159 binary 174 disease endpoints, 38 biomarkers, and 23 medication usages in ~179,000 individuals in BBJ 175 (Figure 1a-c, Supplementary Table 1 and 2 for phenotype summary). To maximize the statistical power, we used a linear mixed model implemented in SAIGE¹² (for binary traits) 176 and BOLT¹³ (for quantitative traits). By using linkage disequilibrium (LD)-score regression¹⁴, 177

178 we confirmed that the confounding biases were controlled in the GWASs (Supplementary 179 **Table 3**). In this expanded scope of GWASs in the Japanese population, we identified 396 180 genome-wide significant loci across 159 disease endpoints, 1891 across 38 biomarkers, of which 92 and 156 loci were novel, respectively ($P < 1.0 \times 10^{-8}$; see **Methods**, 181 182 **Supplementary Table 4**). We conducted the initial medication-usage GWASs in East Asian 183 populations, and detected 171 genome-wide significant loci across 23 traits (see Methods). 184 These signals underscore the value of (i) conducting GWASs in non-Europeans and (ii) 185 expanding scope of phenotypes by incorporating biobank resources such as PMH and EMR. 186 For example, we detected an East Asian-specific variant, rs140780894, at the MHC locus in pulmonary tuberculosis (PTB; Odds Ratio [OR] = 1.2, $P = 2.9 \times 10^{-23}$, Minor Allele 187 188 Frequency[MAF]_{EAS} = 0.24; **Extended Data Figure 2**), which was not present in European population (Minor Allele Count $[MAC]_{EUR} = 0$)¹⁵. PTB is a serious global health burden and 189 190 relatively endemic in Japan¹⁶ (annual incidence per 100,000 was 14 in Japan whereas 8 in 191 the United Kingdom and 3 in the United States in 2018 [World Health Organization, Global 192 Tuberculosis Report]). Because PTB, an infectious disease, can be treatable and remittable. 193 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 549⁹ to 7.800 case 194 195 individuals). Similar to this example, we identified a novel signal at rs190894416 at 7p14.2 (OR = 16, $P = 6.9 \times 10^{-9}$; **Extended Data Figure 3**) in dysentery, which is bloody diarrhea 196 197 caused by infection with Shigella bacillus that was once endemic in Japan when poor hygiene had been common¹⁷. We also identified novel signals in common diseases that 198 199 have not been target diseases but were included in the PMH record, such as rs715 at 3'UTR

of CPS1 in cholelithiasis (Extended Data Figure 4; OR = 0.87, $P = 9.6 \times 10^{-13}$) and 200 201 rs2976397 at the PSCA locus in gastric ulcer, gastric cancer, and gastric polyp (Extended **Data Figure 5**; OR = 0.86, $P = 6.1 \times 10^{-24}$). We detected pleiotropic functionally impactful 202 203 variants, such as a deleterious missense variant, rs28362459 (p.Leu20Arg), in FUT3 associated with gallbladder polyp (OR = 1.46, $P = 5.1 \times 10^{-11}$) and cholelithiasis (OR = 1.11, P 204 = 7.3×10^{-9} : **Extended Data Figure 6**), and a splice donor variant, rs56043070 (c.89+1G>A), 205 causing loss of function of GCSAML associated with urticaria (OR = 1.24, P = 6.9×10^{-12} ; 206 207 **Extended Data Figure 7**), which was previously reported to be associated with platelet and reticulocyte counts¹⁸. Medication-usage GWASs also provided interesting signals as an 208 209 alternative perspective for understanding disease genetics¹⁹. For example, individuals taking 210 HMG CoA reductase inhibitors (C10AA in Anatomical Therapeutic Chemical Classification 211 [ATC]) were likely to harbor variations at HMGCR (lead variant at rs4704210, OR = 1.11, P = 2.0×10⁻²⁷). Prescription of salicylic acids and derivatives (N02BA in ATC) were significantly 212 213 associated with a rare East Asian missense variant in PCSK9, rs151193009 (p.Arg93Cys; OR = 0.75, $P = 7.1 \times 10^{-11}$, MAF_{EAS}=0.0089, MAF_{EUR}=0.000; **Extended Data Figure 8**), which 214 215 might indicate a strong protective effect against the thromboembolic diseases in general. 216 To confirm that the signals identified in BBJ were replicable, we conducted GWASs of 217 corresponding phenotypes (i.e., disease endpoints and biomarkers) in UKB and FinnGen, 218 and collected summary statistics of medication usage GWAS recently conducted in UKB¹⁹ 219 (Supplementary Table 5). We then compared the effect sizes of the genome-wide 220 significant variants in BBJ with those in a European dataset across binary and quantitative

221 traits (see **Methods**). The loci identified in our GWASs were successfully replicated in the same effect direction (1,830 out of 1,929 [94.9%], $P < 10^{-325}$ in sign test) and with high 222 223

effect-size correlation (Extended Data Figure 9).

224 Motivated by the high replicability, we performed trans-ethnic meta-analyses of these 225 220 harmonized phenotypes across three biobanks (see Methods). We identified 1,362 226 disease-associated, 10,572 biomarker-associated, and 841 medication-associated loci in 227 total, of which 356, 3,576, and 236 were novel, respectively (Figure 1d, Supplementary 228 Table 6). All these summary results of GWASs are openly shared without any restrictions. 229 Together, we successfully expanded the genomic map of human complex traits in terms of 230 populations and scope of phenotypes through conducting deep-phenotype GWASs across 231 trans-ethnic nationwide biobanks.





Figure 1. Overview of the identified loci in the trans-ethnic meta-analyses of 220 deep phenotype GWASs.

- 236 (**a-c**) The pie charts describe the phenotypes analyzed in this study. The disease endpoints
- 237 (**a**; $n_{\text{trait}} = 159$) were categorized based on the ICD10 classifications (A to Z; **Supplementary**
- Table 1a), the biomarkers (b; $n_{\text{trait}} = 38$; Supplementary Table 1b) were classified into nine
- 239 categories, and medication usage was categorized based on the ATC system (A to S;
- 240 Supplementary Table 1c). (d) The genome-wide significant loci identified in the

- trans-ethnic meta-analyses and pleiotropic loci ($P < 1.0 \times 10^{-8}$). The traits (rows) are sorted as
- shown in the pie chart, and each dot represents significant loci in each trait. Pleiotropic loci
- are annotated by lines with a locus symbol.

245 The regional landscape of pleiotropy.

246 Because human traits are highly polygenic and the observed variations within the human 247 genome are finite in number, pleiotropy, where a single variant affects multiple traits, is 248 pervasive²⁰. While pleiotropy has been intensively studied in European populations by compiling previous GWASs^{20,21}, the landscape of pleiotropy in non-European populations 249 250 has remained elusive. By leveraging this opportunity for comparing the genetics of deep 251 phenotypes across populations, we sought to investigate the landscape of regional 252 pleiotropy in both Japanese and European populations. We defined the degree of pleiotropy as the number of significant associations per variant $(P < 1.0 \times 10^{-8})^{21}$. In the Japanese, 253 254 rs11066015 harbored the largest number of genome-wide significant associations (45 traits; 255 Figure 2a), which was in tight LD with a missense variant at the ALDH2 locus, rs671. 256 Following this, rs117326768 at the MHC locus (23 traits) and rs1260326 at the GCKR locus 257 (18 traits) were most pleiotropic. In Europeans, rs3132941 at the MHC locus harbored the 258 largest number of genome-wide significant associations (46 traits; Figure 2b), followed by 259 rs4766578 at the ATXN2/SH2B3 locus (38 traits) and rs4665972 at the GCKR locus (28 260 traits). Notably, the ALDH2 locus (pleiotropic in Japanese) and the MHC locus (pleiotropic in Japanese and Europeans) are known to be under recent positive selection^{22,23}. To 261 262 systematically assess whether pleiotropic regions in the genome were likely to be under 263 selection pressure in each of the populations, we investigated the enrichment of the 264 signatures of recent positive selection quantified by the metric singleton density score (SDS)²² values within the pleiotropic loci (see **Methods**). Intriguingly, when compared with 265 those under the null hypothesis, we observed significantly higher values of SDS χ^2 values 266 267 within the pleiotropic loci, and this fold change increased as the number of associations

- increased (i.e., more pleiotropic) in both Japanese and Europeans (Figure 2c and 2d). To
- summarize, the trans-ethnic atlas of genetic associations elucidated the broadly shared
- 270 landscape of pleiotropy, which implied a potential connection to natural selection signatures
- 271 affecting human populations.
- 272



274 Figure 2. Number of significant associations per variant.

(a, b) The Manhattan-like plots show the number of significant associations ($P < 1 \times 10^{-8}$) 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 χ^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.

283 Pleiotropic associations in HLA and ABO locus.

284 Given the strikingly high number of associations in both populations, we next sought to 285 fine-map the pleiotropic signals within the MHC locus. To this end, we imputed the classical 286 HLA alleles in BBJ and UKB, and performed association tests for 159 disease endpoints and 287 38 biomarkers (Figure 3a and 3b). After the fine-mapping and conditional analyses (see 288 Methods), we identified 94 and 153 independent association signals in BBJ and UKB, respectively (the regional threshold of significance was set to $P < 1.0 \times 10^{-6}$; Supplementary 289 290 Table 7). Overall, HLA-B in class I and HLA-DRB1 in class II harbored the largest number of 291 associations in both BBJ and UKB. For example, we successfully fine-mapped the strong signal associated with PTB to HLA-DR β 1 Ser57 (OR = 1.20, P = 7.1×10⁻¹⁹) in BBJ. This is 292 293 the third line of evidence showing the robust association of HLA with tuberculosis identified to date^{24,25}, and we initially fine-mapped the signal to *HLA-DRB1*. Interestingly, HLA-DRB1 294 295 at position 57 also showed pleiotropic associations with other autoimmune and 296 thyroid-related diseases, such as Grave's disease (GD), hyperthyroidism, Hashimoto's 297 disease, hypothyroidism, Sjogren's disease, chronic hepatitis B, and atopic dermatitis in BBJ. 298 Of note, the effect direction of the association of HLA-DR^{β1} Ser57 was the same between hyperthyroid status (OR = 1.29, $P = 2.6 \times 10^{-14}$ in GD and OR = 1.37, $P = 1.4 \times 10^{-8}$ in 299 hyperthyroidism) and hypothyroid status (OR = 1.50, $P = 9.0 \times 10^{-8}$ in Hashimoto's disease 300 and OR = 1.31, $P = 1.5 \times 10^{-7}$ in hypothyroidism), despite the opposite direction of thyroid 301 302 hormone abnormality. This association of HLA-DR^{β1} was also observed in Sjogren's 303 syndrome (OR = 2.04, $P = 7.9 \times 10^{-12}$), which might underlie the epidemiological

comorbidities of these diseases²⁶. Other novel associations in BBJ included HLA-DR^β1 304 Asn197 with sarcoidosis (OR = 2.07, $P = 3.7 \times 10^{-8}$), and four independent signals with 305 306 chronic sinusitis (i.e., HLA-DRA, HLA-B, HLA-A, and HLA-DQA1). 307 Another representative pleiotropic locus in the human genome is the ABO locus. We 308 performed ABO blood-type PheWAS in BBJ and UKB (Figure 3c and 3d). We estimated the 309 ABO blood type from three variants (rs8176747, rs8176746, and rs8176719 at 9q34.2)²⁷, 310 and associated them with the risk of diseases and quantitative traits for each blood group. A 311 variety of phenotypes, including common diseases such as myocardial infarction as well as 312 biomarkers such as blood cell traits and lipids, were strongly associated with the blood types 313 in both biobanks (Supplementary Table 8). We also 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 BBJ²⁸. 314 315



316

317 Figure 3. HLA and ABO association PheWAS.

318 (a,b) Significantly associated HLA genes identified by HLA PheWAS in BBJ (a) or in UKB (b) 319 are plotted. In addition to the top association signals of the phenotypes, independent 320 associations identified by conditional analysis are also plotted, and the primary association 321 signal is indicated by the plots with a gray border. The color of each plot indicates two-tailed 322 P values calculated with logistic regression (for binary traits) or linear regression (for 323 quantitative traits) as designated in the color bar at the bottom. The bars in green at the top 324 indicate the number of significant associations per gene in each of the populations. The 325 detailed allelic or amino acid position as well as statistics in the association are provided in 326 Supplementary Table 7.

- 327 (c,d) Significant associations identified by ABO blood-type PheWAS in BBJ (c) or in UKB (d)
- 328 are shown as boxes and colored based on the odds ratio. The size of each box indicates
- 329 two-tailed *P* values calculated with logistic regression (for binary traits) or linear regression
- 330 (for quantitative traits).
- 331

332 Genetic correlation elucidates the shared phenotypic domains across populations.

333 The interplay between polygenicity and pleiotropy suggests widespread genetic correlations among complex human traits²⁹. Genetic relationships among human diseases have 334 contributed to the refinement of disease classifications³⁰ and elucidation of the biology 335 underlying the epidemiological comorbidity²⁹. To obtain deeper insights into the 336 337 interconnections among human traits and compare them across populations, we computed 338 pairwise genetic correlations (r_q) across 106 traits (in Japanese) and 148 traits (in 339 Europeans) with Z-score for $h_{SNP}^2 > 2$, using bivariate LD score regression (see **Methods**). 340 We then defined the correlated trait domains by greedily searching for the phenotype blocks 341 with pairwise $r_{q} > 0.7$ within 70% of r_{q} values in the block on the hierarchically clustered 342 matrix of pairwise r_{a} values (**Extended Data Figure 10**). We detected domains of tightly 343 correlated phenotypes, such as (i) cardiovascular- acting medications, (ii) coronary artery 344 disease, (iii) type 2 diabetes- related phenotypes, (iv) allergy- related phenotypes, and (v) 345 blood-cell phenotypes in BBJ (**Extended Data Figure 10a**). These domains implicated the 346 shared genetic backgrounds on the similar diseases and their treatments (e.g., (ii) diseases 347 of the circulatory system in ICD10 and for coronary artery disease and their treatments) and 348 diagnostic biomarkers (e.g., (iii) glucose and HbA1c in type 2 diabetes). Intriguingly, the 349 corresponding trait domains were mostly identified in UKB as well (Extended Data Figure 350 **10b**). Thus, we confirmed that the current clinical boundaries for a spectrum of human 351 diseases broadly reflect the shared genetic etiology across populations, despite differences 352 in ethnicity and despite potential differences in diagnostic and prescription practices.

353

354 Deconvolution of a matrix of summary statistics of 159 diseases provides novel 355 insights into disease pathogenesis.

356 A major challenge in genetic correlation is that the r_{a} is a scholar value between two traits, 357 which summarizes the averaged correlation over the whole genome into just one metric³¹. 358 This approach is not straightforward in specifying a set of genetic variants driving the 359 observed correlation, which should pinpoint biological pathways and dysfunctional organs 360 explaining the shared pathogenesis. To address this, gathering of the genetic association 361 statistics of hundreds of different phenotypes can dissect genotype-phenotype association 362 patterns without a prior hypothesis, and identify latent structures underlying a spectrum of 363 complex human traits. In particular, matrix decomposition on the summary statistics is a promising approach^{5,32,33}, which derives orthogonal components that explain association 364 365 variance across multiple traits while accounting for linear genetic architectures in general. 366 This decomposition can address two challenges in current genetic correlation studies. First, 367 it informs us of genetic variants that explain the shared structure across multiple diseases, 368 thereby enabling functional interpretation of the component. Second, it can highlight 369 sub-significant associations and less powered studies, which are important in understanding 370 the contribution of common variants in rare disease genetics with a small number of case samples³² or in genetic studies in underrepresented populations where smaller statistical 371 372 power is inevitable.

Therefore, we applied DeGAs⁵ on a matrix of our disease GWAS summary statistics in Japanese and the meta-analyzed statistics in Europeans ($n_{disease} = 159$; **Figure 4a** and **4b**). To interpret the derived latent components, we annotated the genetic variants explaining each component (i) through GREAT genomic region ontology enrichment analysis³⁴, (ii)

through identification of relevant cell types implicated from tissue specific regulatory DNA (ENCODE3³⁵) and expression (GTEx³⁶) profiles, and (iii) by projecting biomarker GWASs and metabolome GWASs into the component space ($n_{biomarker}=38$, $n_{metabolite_EAS}=206$, $n_{metabolite_EUR}=248$; **Figure 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 in the input summary statistics

383 matrix (Extended Data Figure 11, 12).

384 Globally, hierarchically similar diseases as defined by the conventional ICD10 385 classification were explained by the same components, based on DeGAs trait squared cosine scores that quantifies component loadings⁵ (Figure 4c, d). This would be considered 386 387 as a hypothesis-free support of the historically defined disease classification. For example, 388 component 1 explained the genetic association patterns of diabetes (E10 and E11 in ICD10) 389 and component 2 explained those of cardiac and vascular diseases (100-183), in both 390 populations. Functional annotation enrichment of the genetic variants explaining these 391 components by GREAT showed that component 1 (diabetes component) was associated with abnormal pancreas size (binomial $P_{enrichment} = 7.7 \times 10^{-19}$) as a human phenotype, 392 393 whereas component 2 (cardiovascular disease component) was associated with xanthelasma (i.e., cholesterol accumulation on the eyelids; binomial $P_{\text{enrichment}} = 3.0 \times 10^{-10}$). 394 395 Further, the genes comprising component 1 were enriched in genes specifically expressed in the pancreas ($P_{enrichment} = 5.5 \times 10^{-4}$), and those comprising component 2 were enriched in 396 genes specifically expressed in the aorta ($P_{enrichment} = 1.9 \times 10^{-3}$; Extended Data Figure 13). 397 398 By projecting the biomarker and metabolite GWASs into this component space, we

399 observed that component 1 represented the genetics of glucose and HbA1c, and component 400 2 represented the genetics of blood pressure and lipids, all of which underscored the 401 biological relevance. Thus, this deconvolution-projection analysis elucidated the latent 402 genetic structure behind human diseases, which highlighted the underlying biological 403 functions, relevant tissues, and associated human phenotypes.

404 The latent components shared across diseases explained the common biology behind 405 etiologically similar diseases. For example, we identified that component 10 explained the 406 genetics of cholelithiasis (gall stone), cholecystitis (inflammation of gallbladder), and gall 407 bladder polyp (Figure 4e). The projection of European metabolite GWASs into the 408 component space identified that component 10 represented the metabolite GWAS in the 409 bilirubin metabolism pathway. Component 10 was composed of variants involved in 410 intestinal cholesterol absorption in the mouse phenotype (binomial $P_{\text{enrichment}} = 3.8 \times 10^{-10}$). 411 This is biologically relevant, since increased absorption of intestinal cholesterol is a major cause of cholelithiasis, which also causes cholecystitis³⁷. This projection analysis was also 412 413 applicable to the Japanese metabolites GWASs, which showed the connection between the 414 component 1 (diabetes component) and arginine and glucose levels, and between the 415 component 10 (gallbladder disease component) and glycine, which conjugates with bile acids³⁸. 416

Some components could be further utilized to boost understanding of the underpowered GWASs with the use of well-powered GWAS, and for identifying the contributor of shared genetics between different diseases. For example, we complemented underpowered varicose GWAS in BBJ ($n_{case} = 474$, genome-wide significant loci = 0) with higher-powered GWAS in Europeans ($n_{case} = 22,037$, genome-wide significant loci = 54), since both GWASs

422 were mostly represented by component 11, which was explained by variants related to abnormal vascular development (binomial $P_{\text{enrichment}} = 4.2 \times 10^{-7}$; Figure 4f). Another example 423 424 is component 27, which was shared with rheumatoid arthritis and systemic lupus 425 erythematosus, two distinct but representative autoimmune diseases. Component 27 was 426 explained by the variants associated with interleukin secretion and plasma cell number (binomial $P_{\text{enrichment}} = 6.1 \times 10^{-10}$ and 9.3×10^{-10} , respectively), and significantly enriched in the 427 428 DNase I hypersensitive site (DHS) signature of lymphoid tissue ($P_{\text{enrichment}} = 1.3 \times 10^{-4}$; Figure 429 4g). This might suggest the convergent etiology of the two autoimmune diseases, which 430 could not be elucidated by the genetic correlation alone.

431 Finally, we aimed at hypothesis-free categorization of diseases based on these 432 components. Historically, hypersensitivity reactions have been classified into four types (e.g., types I to IV)³⁹, but the clear sub-categorization of allergic diseases based on this 433 434 pathogenesis and whether the categorization can be achieved solely by genetics were 435 unknown. In our TSVD results, the allergic diseases (mostly J and L in ICD10) were 436 represented by the four components 3, 16, 26, and 34. By combining these components as 437 axis-1 (e.g., components 3 and 16) and axis-2 (e.g., components 26 and 34), and comparing 438 the cumulative variance explained by these axes, we defined axis-1 dominant allergic 439 diseases (e.g., asthma and allergic rhinitis) and axis-2 dominant allergic diseases (metal 440 allergy, contact dermatitis, and atopic dermatitis; Figure 4h). Intriguingly, the axis-1 441 dominant diseases corresponded etiologically well to type I allergy (i.e., immediate 442 hypersensitivity). The variants explaining axis-1 were biologically related to IgE secretion and Th₂ cells (binomial $P_{enrichment} = 9.9 \times 10^{-46}$ and 2.9×10^{-44} , respectively). Furthermore, 443

444 GWAS of eosinophil count was projected onto axis-1, which recapitulated the biology of type I allergy⁴⁰. In contrast, the axis-2 dominant diseases corresponded to type IV allergy (i.e., 445 446 cell-mediated delayed hypersensitivity). The variants explaining axis-2 were associated with IL-13 and interferon secretion (binomial $P_{\text{enrichment}} = 1.6 \times 10^{-10}$ and 5.2×10^{-9} , respectively), 447 and GWAS of C-reactive protein was projected onto axis-2, which was distinct from axis-141. 448 449 To summarize, our deconvolution approach (i) recapitulated the existing disease 450 classifications, (ii) clarified the underlying biological mechanisms and relevant tissues 451 shared among a spectrum of related diseases, and (iii) showed potential application for 452 genetics-driven categorization of human diseases.

453



455

Figure 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 BBJ: Biobank
Japan) was decomposed into latent components, which were interpreted by annotation of a

461 set of genetic variants driving each component and in the context of other GWASs through 462 projection. (b) A schematic representation of TSVD applied to decompose a summary 463 statistic matrix W to derive latent components. U, S, and V represent resulting matrices of 464 singular values (S) and singular vectors (U and V). (c) A heatmap representation of DeGAs 465 squared cosine scores of diseases (columns) to components (rows). The components are 466 shown from 1 (top) to 40 (bottom), and diseases are sorted based on the contribution of 467 each component to the disease measured by the squared cosine score (from component 1 468 to 40). Full results with disease and component labels are in Extended Data Figure 14. 469 (d) Results of TSVD of disease genetics matrix and the projection of biomarker genetics. 470 Diseases (left) and biomarkers (right) are colored based on the ICD10 classification and 471 functional categorization, respectively. The derived components (middle; from 1 to 40) are 472 colored alternately in blue or red. The squared cosine score of each disease to each 473 component and each biomarker to each component is shown as red and blue lines. The 474 width of the lines indicates the degree of contribution. The diseases with squared cosine 475 score > 0.3 in at least one component are displayed. Anth; anthropometry, BP; blood 476 pressure, Metab; metabolic, Prot; protein, Kidn; kidney-related, Ele; Electrolytes, Liver; 477 liver-related, Infl; Inflammatory, BC; blood cell. (e-h) Examples of disease-component 478 correspondence and the biological interpretation of the components by projection and 479 enrichment analysis using GREAT. A representative component explaining a group of 480 diseases based on the contribution score, along with responsible genes, functional 481 enrichment results GREAT, relevant tissues, and relevant biobarkers/metabolites is shown. 482 GB; gallbladder. RA; rheumatoid arthritis. SLE; systemic lupus erythematosus.

483

484 **Discussion**

485 Here, we performed 220 GWASs of human complex traits by incorporating the PMH and 486 EMR data in BBJ, substantially expanding the atlas of genotype-phenotype associations in 487 non-Europeans. We then systematically compared their genetic basis with GWASs of 488 corresponding phenotypes in Europeans. We confirmed the global replication of loci 489 identified in BBJ, and discovered 4,170 novel loci through trans-ethnic meta-analyses, 490 highlighting the value of conducting GWASs in diverse populations. The results are openly 491 shared through web resources, which will be a platform to accelerate further research such as functional follow-up studies and drug discovery⁴². Of note, leveraging these well-powered 492 493 GWASs, we observed that the genes associated with endocrine/metabolic, circulatory, and 494 respiratory diseases (E, I, and J by ICD10) were systematically enriched in targets of 495 approved medications treating those diseases (Extended Data Figure 15). This should 496 motivate us to use this expanded resource for genetics-driven novel drug discovery and 497 drug repositioning.

498 The landscape of regional pleiotropy was globally shared across populations, and 499 pleiotropic regions tended to have been under recent positive selection. Further elucidation 500 of pleiotropy in other populations is warranted to replicate our results. To highlight the utility 501 of deep phenotype GWASs, we finally decomposed the multi-ethnic genotype-phenotype 502 association patterns by TSVD. The latent components derived from TSVD pinpointed the 503 convergent biological mechanisms and relevant cell types across diseases, which can be 504 utilized for re-evaluation of existing disease classifications. The incorporation of biomarker 505 and metabolome GWAS summary statistics enabled further interpretation of the latent 506 components. Our approach suggested a potential avenue for restructuring of the medical

diagnoses through dissecting the shared genetic basis across a spectrum of diseases, as
analogous to the current disease diagnostics historically shaped through empirical
categorization of shared key symptoms across a spectrum of organ dysfunctions.

510 In conclusion, our study substantially expanded the atlas of genetic associations, 511 supported the historically-defined categories of human diseases, and should accelerate the 512 discovery of the biological basis contributing to complex human diseases.

513

514

515 Acknowledgments

516 We sincerely thank all the participants of BioBank Japan, UK Biobank, and FinnGen. This 517 research was supported by the Tailor-Made Medical Treatment program (the BioBank Japan 518 Project) of the Ministry of Education, Culture, Sports, Science, and Technology (MEXT), the 519 Japan Agency for Medical Research and Development (AMED). The FinnGen project is 520 funded by two grants from Business Finland (HUS 4685/31/2016 and UH 4386/31/2016) and 521 nine industry partners (AbbVie, AstraZeneca, Biogen, Celgene, Genentech, GSK, MSD, 522 Pfizer and Sanofi). Following biobanks are acknowledged for collecting the FinnGen project 523 samples: Auria Biobank (https://www.auria.fi/biopankki/), THL Biobank 524 (https://thl.fi/fi/web/thl-biopank), Helsinki Biobank 525 (https://www.terveyskyla.fi/helsinginbiopankki/), Northern Finland Biobank Borealis 526 (https://www.ppshp.fi/Tutkimus-ja-opetus/Biopankki), Finnish Clinical Biobank Tampere 527 (https://www.tays.fi/biopankki), Biobank of Eastern Finland (https://ita-suomenbiopankki.fi), 528 Central Finland Biobank (https://www.ksshp.fi/fi-FI/Potilaalle/Biopankki), Finnish Red Cross 529 Blood Service Biobank (https://www.bloodservice.fi/Research%20Projects/biobanking),

530	Terveystalo	Biobank	Finland
531	(https://www.terveystalo.com/fi/Yritystietoa/	Terveystalo-Biopankki/Biopankki/).	S.S. was in
532	part supported by The Mochida Memoria	al Foundation for Medical and Pha	armaceutical
533	Research . M.Kanai was supported by a Na	akajima Foundation Fellowship and t	he Masason
534	Foundation. Y.Tanigawa is in part suppor	ted by a Funai Overseas Scholarsl	nip from the
535	Funai Foundation for Information Techn	ology and the Stanford University	/ School of
536	Medicine. M.A.R. is in part supported by	y National Human Genome Resea	rch Institute
537	(NHGRI) of the National Institutes of Health	(NIH) under award R01HG010140 (I	M.A.R.), and
538	a National Institute of Health center for Mu	Ilti- and Trans-ethnic Mapping of Me	ndelian and
539	Complex Diseases grant (5U01 HG009080)). The content is solely the response	sibility of the
540	authors and does not necessarily represe	nt the official views of the National	Institutes of
541	Health. Y.O. was supported by the Japan	n Society for the Promotion of Scie	ence (JSPS)
542	KAKENHI (19H01021, 20K21834), and	d AMED (JP20km0405211, JP20)ek0109413,
543	JP20ek0410075, JP20gm4010006, and J	P20km0405217), Takeda Science	Foundation,
544	and Bioinformatics Initiative of Osaka U	niversity Graduate School of Medic	cine, Osaka
545	University.		

546

547 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 studies. S.S., Y. Tanigawa., and M.A.R. conducted statistical decomposition
analysis. S.S., S.T., A.N., G.T., and Y.O. conducted metabolome analysis. A.S., K.S., W.O.,

Ken Yamaji, K.T., S.A., Y.Takahashi, T.S., N.S., H.Y., S.Minami, S.Murayama, Kozo
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.

557

558 Competing Financial Interests

559 M.A.R. is on the SAB of 54Gene and Computational Advisory Board for Goldfinch Bio and 560 has advised BioMarin, Third Rock Ventures, MazeTx and Related Sciences. The funders 561 had no role in study design, data collection and analysis, decision to publish, or preparation 562 of the manuscript.

563

564 **Data availability**

565 The genotype data of BBJ used in this study are available from the Japanese 566 Genotype-phenotype Archive (JGA; http://trace.ddbj.nig.ac.jp/jga/index_e.html) with 567 accession code JGAD0000000123 and JGAS0000000114. The UKB analysis was 568 conducted via application number 47821. This study used the FinnGen release 3 data. 569 Summary statistics of BBJ GWAS and trans-ethnic meta-analysis will be publicly available 570 without any restrictions.

571

572 Code availability

573 We used publicly available software for the analyses. The used software is listed and 574 described in the **Method** section of our manuscript.

575

576 Methods

577 Genome-wide association study of 220 traits in BBJ

578 We conducted 220 deep phenotype GWASs in BBJ. BBJ is a prospective biobank that 579 collaboratively collected DNA and serum samples from 12 medical institutions in Japan and 580 recruited approximately 200,000 participants, mainly of Japanese ancestry (Supplementary 581 **Note**). All study participants had been diagnosed with one or more of 47 target diseases by 582 physicians at the cooperating hospitals. We previously conducted GWASs of 42 out of the 47 target diseases⁹. In this study, we newly curated the PMH records included in the clinical 583 584 data, and performed text-mining to retrieve disease records from the free-format EMR as 585 well. For disease phenotyping, we merged this information with the target disease status, 586 and defined the case status for 159 diseases with a case count > 50 (Supplementary Table 587 2). As controls, we used samples in the cohort without a given diagnosis or related 588 diagnoses, which was systematically defined by using the phecode framework³ 589 (Supplementary Table 1). For medication-usage phenotyping, we again retrieved 590 information by text-mining of 7,018,972 medication records. Then, we categorized each 591 medication trade name by using the ATC, World Health Organization. For biomarker 592 phenotyping, we used the same processing and quality control method as previously described (**Supplementary Table 2** for phenotype summary)^{10,43}. In brief, we excluded 593 594 measurements outside three times of interguartile range (IQR) of upper/lower quartile. For 595 individuals taking anti-hypertensive medications, we added 15 mmHg to systolic blood 596 pressure (SBP) and 10 mmHg to diastolic blood pressure (DBP). For individuals taking a 597 statin, we applied the following correction to the lipid measurements: i) Total cholesterol was 598 divided by 0.8; ii) measured LDL-cholesterol (LDLC) was adjusted as LDLC / 0.7; iii) derived

599 LDLC from the Friedewald was re-derived as (Total cholesterol / 0.8) - HDLC - (Triglyceride/600 5).

601 We genotyped participants with the Illumina HumanOmniExpressExome BeadChip or a 602 combination of the Illumina HumanOmniExpress and HumanExome BeadChips. Quality 603 control of participants and genotypes was performed as described elsewhere¹¹. In this 604 project, we analyzed 178,726 participants of Japanese ancestry as determined by the 605 principal component analysis (PCA)-based sample selection criteria. The genotype data 606 were further imputed with 1000 Genomes Project Phase 3 version 5 genotype (n = 2,504) 607 and Japanese whole-genome sequencing data (n = 1,037) using Minimac3 software. After 608 this imputation, we excluded variants with an imputation quality of Rsg < 0.7.

609 We conducted GWASs for binary traits (i.e., disease endpoints and medication usage) 610 by using a generalized linear mixed model implemented in SAIGE (version 0.37), which had 611 substantial advantages in terms of (i) maximizing the sample size by including genetically related participants, and (ii) controlling for case-control imbalance¹², which was the case in 612 many of the disease endpoints in this study. We included adjustments for age, age², sex, 613 614 agexsex, age²xsex, and top 20 principal components for as covariates used in step 1. For 615 sex-specific diseases, we alternatively adjusted for age, age², and the top 20 principal 616 components as covariates used in step 1, and we used only controls of the sex to which the 617 disease is specific. For the X chromosome, we conducted GWASs separately for males and females, and merged their results by inverse-variance fixed-effects meta-analysis⁴⁴. We 618 619 conducted GWASs for quantitative traits (i.e., biomarkers) by using a linear mixed model 620 implemented in BOLT-LMM (version 2.3.4). We included the same covariates as used in the 621 binary traits above.

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.

625

626 Harmonized genome-wide association study of 220 traits in UKB and FinnGen

627 We conducted the GWASs harmonized with BBJ in UKB and in FinnGen. The UK Biobank 628 project is a population-based prospective cohort that recruited approximately 500,000 629 people across the United Kingdom (Supplementary Note). We defined case and control 630 status of 159 disease endpoints, which were originally retrieved from the clinical information 631 in UKB and mapped to BBJ phenotypes via phecode (**Supplementary Table 1**). We also 632 analyzed 38 biomarker values provided by the UKB. The genotyping was performed using 633 either the Applied Biosystems UK BiLEVE Axiom Array or the Applied Biosystems UK 634 Biobank Axiom Array. The genotypes were further imputed using a combination of the 635 Haplotype Reference Consortium, UK10K, and 1000 Genomes Phase 3 reference panels by IMPUTE4 software⁷. In this study, we analyzed 361,194 individuals of white British genetic 636 637 ancestry as determined by the PCA-based sample selection criteria (see URLs). We 638 excluded the variants with (i) INFO score ≤ 0.8 , (ii) MAF ≤ 0.0001 (except for missense and protein-truncating variants annotated by VEP⁴⁵, which were excluded if MAF $\leq 1 \times 10^{-6}$), and 639 (iii) $P_{HWF} \le 1 \times 10^{-10}$. We conducted GWASs for 159 disease endpoints by using SAIGE with 640 641 the same covariates used in the BBJ GWAS. For biomarker GWASs, we used publicly 642 available summary statistics of UKB biomarker GWAS when available (see URLs), and 643 otherwise performed linear regression using PLINK software with the same covariates, excluding the genetically related individuals (the 1st, 2nd, or 3rd degree)⁷. For medication 644

645 usage GWASs, we used publicly available summary statistics of medication usage in UKB¹⁹, 646 which was organized by the ATC and thus could be harmonized with BBJ GWASs. 647 FinnGen is a public-private partnership project combining genotype data from Finnish 648 biobanks and digital health record data from Finnish health registries (Supplementary 649 **Notes**). For GWASs, we used the summary statistics of FinnGen release 3 data (see **URLs**). 650 The disease endpoints were mapped to BBJ phenotypes by using ICD10 code, and we 651 defined 129 out of 159 endpoints in BBJ. We did not conduct biomarker and 652 medication-related GWASs because the availability of these phenotypes was limited. 653 654 Meta-analysis, definition of significant loci, and annotation of the lead variants with 655 genome-wide significance 656 First, we performed intra-European meta-analysis when summary statistics of both UKB and 657 FinnGen were available, and then performed trans-ethnic meta-analysis across three or two 658 cohorts in 159 disease endpoints, 38 biomarker values, and 23 medication usage GWASs. 659 We conducted these meta-analyses by using the inverse-variance method and estimated

heterogeneity with Cochran's Q test with metal software⁴⁴. The summary statistics of 661 primary GWASs in BBJ and trans-ethnic meta-analysis GWASs are openly shared without 662 any restrictions.

660

We adopted the genome-wide significance threshold of $< 1.0 \times 10^{-8}$, as previously used 663 664 in similar projects in BBJ and UKB^{9,21}. We defined independent genome-wide significant loci 665 on the basis of genomic positions within ±500 kb from the lead variant. We considered a 666 trait-associated locus as novel when the locus within ±1 Mb from the lead variant did not 667 include any variants that were previously reported to be significantly associated with the

same disease. We basically searched for previous reports of known loci in the GWAS catalog¹⁸, but also referred to PubMed or preprints when the corresponding trait was not included in GWAS catalog or when the large-scale GWASs were released in the preprint server as of July 2020 (**Supplementary Table 9**).

672 We annotated the lead variants using ANNOVAR software, such as rsIDs in dbSNP

database (see **URLs**), the genomic region and closest genes, and functional consequences.

674 We also supplemented this with the gnomAD database¹⁵, and also looked for the allele

675 frequencies in global populations as an independent resource.

676

677 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 in the association test in BBJ and in European GWAS. We also performed the correlation tests with variants with $P_{EUR} < 0.05$ and to those with $P_{EUR} < 1.0 \times 10^{-8}$.

685

686 Evaluation of regional pleiotropy

We assessed the regional pleiotropy based on each tested genetic variant separately for BBJ GWASs and for European GWASs (i.e., intra- European 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

691 significant associations across 220 traits. We then annotated loci from the largest number of 692 associations ($n_{associations} > 9$ in BBJ and > 18 in Europeans) in **Figure 2a**, **b**.

693 Next, we assessed the recent natural selection signature within the pleiotropic loci 694 separately for Japanese and for Europeans. To do this, we first defined the pleiotropic loci 695 by identifying genetic variants that harbored a larger number of significant associations than 696 a given threshold. We varied this threshold from 1 to 40. Then, at each threshold, we calculated the sum of SDS χ^2 values within the pleiotropic loci, and compared this with the χ^2 697 698 distribution under the null hypothesis with a degree of freedom equal to the number of 699 variants in the loci. We thus estimated the SDS enrichment within the pleiotropic loci defined 700 by a given threshold as fold change and P value. The SDS values were obtained from the 701 web resource indicated in the original article on Europeans (see URLs) and provided by the authors on Japanese²³. The raw SDS values were normalized according to the derived allele 702 703 frequency as described previously.

704

705 Fine-mapping of HLA and ABO loci

706 We performed the fine-mapping of MHC associations in BBJ and UKB by HLA imputation⁴⁶. 707 In BBJ, we imputed classical HLA alleles and corresponding amino acid sequences using 708 the reference panel recently constructed from 1,120 individuals of Japanese ancestry by the 709 combination of SNP2HLA software, Eagle, and minimac3, as described previously⁴⁷. We 710 applied post-imputation quality control to keep the imputed variants with minor allele 711 frequency (MAF) \ge 0.5% and Rsg > 0.7. For each marker dosage that indicated the 712 presence or absence of an investigated HLA allele or an amino acid sequence, we 713 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)⁴⁸. We applied the same post-imputation quality control and performed the association tests as in BBJ.

719

720 Heritability and genetic correlation estimation

721 We performed LD score regression (see URLs) for GWASs of BBJ and Europeans to 722 estimate SNP-based heritability, potential bias, and pairwise genetic correlations. Variants in 723 the MHC region (chromosome 6:25–34 Mb) were excluded. We also excluded variants with $x^2 \supseteq > \Box 80$, as recommended previously⁴⁹. For heritability estimation, we used the baselineLD 724 725 model (version 2.2), which included 97 annotations that correct for bias in heritability 726 estimates⁵⁰. We note that we did not report liability-scale heritability, since population 727 prevalence of 159 diseases in each country was not always available, and the main 728 objective of this analysis was an assessment of bias in GWAS, rather than the accurate 729 estimation of heritability. We calculated the heritability Z-score to assess the reliability of 730 heritability estimation, and reported the LDSC results with Z-score for h^2_{SNP} is > 2 731 (Supplementary Table 3). For calculating pairwise genetic correlation, we again restricted 732 the target GWASs to those whose Z-score for h_{SNP}^2 is > 2, as recommended previously⁴⁹. In 733 total, we calculated genetic correlation for 106 GWASs in BBJ and 148 in European GWASs, 734 which resulted in 5,565 and 10,878 trait pairs, respectively.

To illustrate trait-by-trait genetic correlation, we hierarchically clustered the r_g values with hclust and colored them as a heatmap (**Extended Data Figure 10**). To adopt reliable

genetic correlations, we restricted the r_g values that had $P_{cor} < 0.05$. Otherwise, the r_g values were replaced with 0. We then defined the tightly clustered trait domains by greedily searching for the phenotype blocks with pairwise $r_g > 0.7$ within 70% of r_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 Figure 10**).

743

744 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 framework⁵. In this study, we first focused on 159 disease endpoint GWASs in BBJ and European GWAS (i.e., 318 in total) to derive latent components through TSVD. On constructing a Z-score matrix, we conducted variant-level QC. 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 were satisfied:

752 - P value of marginal association ≥ 0.001

753 - Standard error of beta value ≥ 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 either in BBJ or Europeans. We then performed LD pruning using PLINK software⁵¹ ("--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

760 populations) × 22,980 (M: variants). With a predetermined number of K, TSVD decomposed **W** into a product of three matrices: **U**, **S**, and **V**^T: **W** $\square = \square$ **USV**^T. **U** $\square = \square$ ($u_{i,k}$)_{i,k} is 761 762 an orthonormal matrix of size $N \square \times \square K$ whose columns are phenotype singular vectors, **S** is a diagonal matrix of size $K \square \times \square K$ whose elements are singular values, and $\mathbf{V} \square = \square (v_{i,k})_{i,k}$ is an 763 764 orthonormal matrix of size $M \supseteq \times \supseteq K$ whose columns are variant singular vectors. Here we set 765 K as 40, which together explained 36.7% of the total variance of the original matrix. This 766 value was determined by experimenting with different values from 20 to 100 and selecting 767 the informative and sufficient threshold. We used the TruncatedSVD module in the 768 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;

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

772 where

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

774

775 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;

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

778 For each component, we can thus rank the variants based on their contribution to the 779 component and calculate the cumulative contribution score. We defined a set of *contributing* 780 variants to a given component to include top-ranked variants that had high contribution 781 scores until the cumulative contribution score to the component exceeded 0.5. For these 782 variant sets contributing to the latent components, we performed the GREAT (version 4.0.4) binomial genomic region enrichment analysis³⁴ based on the size of the regulatory domain 783 784 of genes and guantified the significance of enrichment in terms of binomial fold enrichment 785 and binomial P value to biologically interpret these components. We used the human 786 phenotype and mouse genome informatics phenotype ontology, which contains manually 787 curated knowledge about the hierarchical structure of phenotypes and genotype-phenotype 788 mapping of human and mouse, respectively. The enriched annotation with a false discovery 789 rate (FDR) < 0.05 is considered significant and displayed in the figures.

790 For a gene set associated with the contributing variants with a given component (P<791 0.05), we sought to identify relevant cell types by integrating two datasets: (i) ENCODE3 792 DHS regulatory patterns across human tissues from non-negative matrix factorization and (ii) specifically expressed genes defined from GTEx data³⁶. In brief, a 793 $(NFM)^{35}$ 794 vocabulary (i.e., DHS patterns) for regulatory patterns was defined from the NFM of 3 million 795 DHSs x 733 human biosamples encompassing 438 cell and tissue types. Then, for each 796 regulatory vocabulary, GENCODE genes were assigned based on their overlying DHSs. The gene labeling result was downloaded from the journal website³⁵. We also defined genes 797 798 specifically expressed in 53 tissues from GTEx version 7 data, based on the top 5% of the tstatistics in each tissue as described elsewhere⁵². Then, for (i) each regulatory vocabulary 799

and (ii) 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.

802

803 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 (i.e., 76 rows) and 248 known metabolites of independent previous GWASs in the European population⁵³ × 22,980 variants (**Supplementary Table 10**). 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 software⁵⁴ (version 0.5.5 --ref 1KG/EUR --impute.maf 0.01), and otherwise we set the missing Z-scores to zero.

815 Projection of Metabolite GWASs in Japanese into the component space

To investigate whether the projection analysis is applicable to independent dataset, we conducted metabolite GWASs in Tohoku Medical Megabank Organization (ToMMo). ToMMo is a community-based biobank that combines medical and genome information from the participants in the Tohoku region of Japan⁵⁵. Detailed cohort description is presented in **Supplementary Notes**. In this study, we analyzed a total of 206 metabolites⁵⁶ measured by proton nuclear magnetic resonance (NMR) or liquid chromatography (LC)–MS

822 (Supplementary Table 11). For sample QC, we excluded samples meeting any of the 823 following criteria: (1) genotype call rate < 95%. (2) one individual from each pair of those in close genetic relation (PI HAT calculated by $PLINK^{51} \ge 0.1875$) based on call rate, and (3) 824 825 outliers from Japanese ancestry cluster based on the principal component analysis with 826 samples of 1KGP phase 3 data. For phenotype QC, we excluded (1) the measurements in 827 pregnant women, (2) those which took time from sampling to biobanking \geq 2 days, and (3) 828 phenotypic outlier defined as log-transformed measurements laying more than 4 SD from 829 the mean for each metabolite. The participants were genotyped with a custom SNP array for 830 the Japanese population (i.e., Japonica Array v2). For genotype QC, we excluded variants 831 meeting any of the following criteria: (1) call rate < 98%, (2) P value for Hardy–Weinberg equilibrium < $1.0 \square \times \square 10^{-6}$, and (3) MAF < 0.01. The QCed genotype data were pre-phased 832 833 by using SHAPEIT2 software (r837), and imputed by using IMPUTE4 software (r300.3) with 834 a combined reference panel of 1KGP phase3 (n = 2.504) and population specific WGS data 835 (i.e., 3.5KJPNv2; n = 3.552)⁵⁶. After imputation, we excluded variants with imputation INFO 836 < 0.7.

837 For GWASs, we obtained the residuals from a linear regression model of each of log-transformed metabolites adjusted for age, age², sex, time period from sampling to 838 839 biobanking, and top 20 genotype PCs. The residuals were then transformed by rank-based 840 inverse normalization. Association analysis of imputed genotype dosage with the normalized 841 residual of each metabolite was performed using PLINK2 software. We constructed the 842 Z-score matrices (**W**') of the Japanese metabolites GWASs (i.e., 206 rows) \times 22,980 843 variants, in which we applied the same QC to the Z-scores and set the missing Z-scores to 844 zero again. We then performed the projection as described above.

845

846 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)⁵⁷ 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.

854 References

855 1. Berger, D. A brief history of medical diagnosis and the birth of the clinical laboratory. 856 Part 1--Ancient times through the 19th century. MLO. Med. Lab. Obs. 31, (1999). 857 2. Organización Mundial de la Salud. International statistical classification of diseases 858 and related health problems, 10th revision (ICD-10). World Heal. Organ. (2016). 859 3. Denny, J. C. et al. Systematic comparison of phenome-wide association study of 860 electronic medical record data and genome-wide association study data. Nat. 861 *Biotechnol.* **31**, 1102–1110 (2013). 862 4. Martin, A. R. et al. Clinical use of current polygenic risk scores may exacerbate health 863 disparities. Nat. Genet. 51, 584-591 (2019). 864 5. Tanigawa, Y. et al. Components of genetic associations across 2,138 phenotypes in 865 the UK Biobank highlight adipocyte biology. Nat. Commun. 10, 1–14 (2019). 866 Nagai, A. et al. Overview of the BioBank Japan Project: Study design and profile. J. 6. 867 Epidemiol. 27, S2–S8 (2017). 868 7. Bycroft, C. et al. The UK Biobank resource with deep phenotyping and genomic data. 869 Nature 562, 203-209 (2018). 870 8. Ritchie, M. D. et al. Robust Replication of Genotype-Phenotype Associations across 871 Multiple Diseases in an Electronic Medical Record. Am. J. Hum. Genet. 86, 560–572 872 (2010). 873 9. Ishigaki, K. et al. Large-scale genome-wide association study in a Japanese 874 population identifies novel susceptibility loci across different diseases. Nat. Genet. 52, 875 669-679 (2020). 876 Kanai, M. et al. Genetic analysis of quantitative traits in the Japanese population links 10. 877 cell types to complex human diseases. Nat. Genet. 50, 390–400 (2018). 878 11. Akiyama, M. et al. Characterizing rare and low-frequency height-associated variants in 879 the Japanese population. Nat. Commun. 10, 4393 (2019). 880 12. Zhou, W. et al. Efficiently controlling for case-control imbalance and sample 881 relatedness in large-scale genetic association studies. Nat. Genet. 50, 1335–1341 882 (2018). 883 13. Loh, P. R. et al. Efficient Bayesian mixed-model analysis increases association power 884 in large cohorts. Nat. Genet. 47, 284–290 (2015). 885 14. Bulik-Sullivan, B. et al. LD score regression distinguishes confounding from 886 polygenicity in genome-wide association studies. Nat. Genet. 47, 291–295 (2015).

15. Karczewski, K. J. et al. The mutational constraint spectrum quantified from variation in

888 141,456 humans. *Nature* **581**, 434–443 (2020).

- 16. Hagiya, H. *et al.* Trends in incidence and mortality of tuberculosis in Japan: A
- population-based study, 1997-2016. *Epidemiology and Infection* vol. 147 (2019).
- 17. Kudoh, Y. & Sakai, S. Current Status of Bacterial Diarrheal Diseases in Japan. in
- 892 Bacterial Diarrheal Diseases 83–93 (Springer Netherlands, 1985).
- 893 doi:10.1007/978-94-009-4990-4_8.
- 894 18. Welter, D. *et al.* The NHGRI GWAS Catalog, a curated resource of SNP-trait
 895 associations. *Nucleic Acids Res.* 42, D1001-6 (2014).
- 896 19. Wu, Y. *et al.* Genome-wide association study of medication-use and associated
 897 disease in the UK Biobank. *Nat. Commun.* **10**, 1–10 (2019).
- 898 20. Watanabe, K. *et al.* A global overview of pleiotropy and genetic architecture in
 899 complex traits. *Nat. Genet.* 51, 1339–1348 (2019).
- 21. Canela-Xandri, O., Rawlik, K. & Tenesa, A. An atlas of genetic associations in UK
 Biobank. *Nat. Genet.* **50**, 1593–1599 (2018).
- 902 22. Field, Y. *et al.* Detection of human adaptation during the past 2000 years. *Science*903 (80-.). 354, 760–764 (2016).
- 904 23. Okada, Y. *et al.* Deep whole-genome sequencing reveals recent selection signatures
 905 linked to evolution and disease risk of Japanese. *Nat. Commun.* 9, 1631 (2018).
- 906 24. Qi, H. *et al.* Discovery of susceptibility loci associated with tuberculosis in Han
 907 Chinese. *Hum. Mol. Genet.* 26, 4752–4763 (2017).
- 908 25. Sveinbjornsson, G. *et al.* HLA class II sequence variants influence tuberculosis risk in
 909 populations of European ancestry. *Nat. Genet.* 48, 318–322 (2016).
- 910 26. Baldini, C., Ferro, F., Mosca, M., Fallahi, P. & Antonelli, A. The association of Sjögren
 911 syndrome and autoimmune thyroid disorders. *Frontiers in Endocrinology* vol. 9 121
 912 (2018).
- 913 27. Nakao, M. *et al.* ABO blood group alleles and the risk of pancreatic cancer in a
 914 Japanese population. *Cancer Sci.* **102**, 1076–1080 (2011).
- 915 28. Edgren, G. *et al.* Risk of gastric cancer and peptic ulcers in relation to ABO blood type:
 916 A cohort study. *Am. J. Epidemiol.* **172**, 1280–1285 (2010).
- 917 29. Bulik-Sullivan, B. *et al.* An atlas of genetic correlations across human diseases and
 918 traits. *Nat. Genet.* 47, 1236–1241 (2015).
- 919 30. Anttila, V. *et al.* Analysis of shared heritability in common disorders of the brain.
 920 Science (80-.). 360, (2018).

Shi, H., Mancuso, N., Spendlove, S. & Pasaniuc, B. Local Genetic Correlation Gives

921

31.

922 Insights into the Shared Genetic Architecture of Complex Traits. Am. J. Hum. Genet. 923 **101**, 737–751 (2017). 924 Burren, O. S. & Wallace, C. Informed dimension reduction of clinically-related 32. 925 genome-wide association. bioRxiv (2020). 926 33. Chasman, D. I., Giulianini, F., Demler, O. V. & Udler, M. S. Pleiotropy-Based 927 Decomposition of Genetic Risk Scores: Association and Interaction Analysis for Type 928 2 Diabetes and CAD. Am. J. Hum. Genet. 106, 646-658 (2020). 929 34. McLean, C. Y. et al. GREAT improves functional interpretation of cis-regulatory 930 regions. Nat. Biotechnol. 28, 495-501 (2010). 931 35. Meuleman, W. et al. Index and biological spectrum of human DNase I hypersensitive 932 sites. Nature 584, 244–251 (2020). 933 36. GTEx Consortium, F. et al. Genetic effects on gene expression across human tissues. 934 Nature 550, 204–213 (2017). 935 37. Portincasa, P. & Wang, D. Q. H. Intestinal absorption, hepatic synthesis, and biliary 936 secretion of cholesterol: Where are we for cholesterol gallstone formation? 937 Hepatology vol. 55 1313-1316 (2012). 938 38. Vessey, D. A. The biochemical basis for the conjugation of bile acids with either 939 glycine or taurine. *Biochem. J.* **174**, 621–626 (1978). 940 39. Coombs, R. R. A. & Gell, P. G. . The Classification of Allergic Reactions Underlying 941 Disease. in *Clinical Aspects of Immunology* 317–337 (1963). 942 40. Stone, K. D., Prussin, C. & Metcalfe, D. D. IgE, mast cells, basophils, and eosinophils. 943 J. Allergy Clin. Immunol. **125**, S73 (2010). 944 41. Kobayashi, K., Kaneda, K. & Kasama, T. Immunopathogenesis of delayed-type 945 hypersensitivity. Microsc. Res. Tech. 53, 241-245 (2001). 946 42. Okada, Y. et al. Genetics of rheumatoid arthritis contributes to biology and drug 947 discovery. Nature 506, 376-381 (2014). 948 Sakaue, S. et al. Trans-biobank analysis with 676,000 individuals elucidates the 43. 949 association of polygenic risk scores of complex traits with human lifespan. Nat. Med. 950 **26**, 542–548 (2020). 951 44. Willer, C. J., Li, Y. & Abecasis, G. R. METAL: Fast and efficient meta-analysis of 952 genomewide association scans. Bioinformatics 26, 2190-2191 (2010). 953 45. McLaren, W. et al. The Ensembl Variant Effect Predictor. Genome Biol. 17, 122 954 (2016). 47

855 46. Raychaudhuri, S. *et al.* Five amino acids in three HLA proteins explain most of the
856 association between MHC and seropositive rheumatoid arthritis. *Nat. Genet.* 44,
957 291–296 (2012).

- 47. Hirata, J. *et al.* Genetic and phenotypic landscape of the major histocompatibility
 complex region in the Japanese population. *Nat. Genet.* **51**, 470–480 (2019).
- 960 48. Jia, X. *et al.* Imputing Amino Acid Polymorphisms in Human Leukocyte Antigens.
- 961 *PLoS One* **8**, e64683 (2013).
- 49. 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).
- 965 50. Gazal, S. *et al.* Linkage disequilibrium-dependent architecture of human complex traits
 966 shows action of negative selection. *Nat. Genet.* 49, 1421–1427 (2017).
- 967 51. Purcell, S. et al. PLINK: A Tool Set for Whole-Genome Association and
- 968 Population-Based Linkage Analyses. Am. J. Hum. Genet. 81, 559–575 (2007).
- 969 52. Finucane, H. K. *et al.* Heritability enrichment of specifically expressed genes identifies
 970 disease-relevant tissues and cell types. *Nat. Genet.* **50**, 621–629 (2018).
- 971 53. Shin, S. Y. *et al.* An atlas of genetic influences on human blood metabolites. *Nat.*972 *Genet.* 46, 543–550 (2014).
- 973 54. Rüeger, S., McDaid, A. & Kutalik, Z. Evaluation and application of summary statistic
 974 imputation to discover new height-associated loci. *PLoS Genet.* 14, e1007371 (2018).
- 975 55. Kuriyama, S. et al. The Tohoku Medical Megabank Project: Design and mission. J.
- 976 *Epidemiol.* **26**, 493–511 (2016).
- 977 56. Tadaka, S. *et al.* JMorp: Japanese Multi Omics Reference Panel. *Nucleic Acids Res.*978 46, D551–D557 (2018).
- 979 57. Sakaue, S. & Okada, Y. GREP: Genome for REPositioning drugs. *Bioinformatics* 35, 3821–3823 (2019).
- 981
- 982

- 983 URLs
- 984 SDS values in UK10K provided by Pritchard's lab;
- 985 http://web.stanford.edu/group/pritchardlab/UK10K-SDS-values.zip
- 986 Summary statistics of biomarker GWASs in UKB by Neale's lab ;
- 987 http://www.nealelab.is/uk-biobank/ukbround2announcement
- 988 LDSC software; https://github.com/bulik/ldsc
- 989 FinnGen release 3 data; https://www.finngen.fi/en/access_results
- 990 dbSNP; https://www.ncbi.nlm.nih.gov/snp/
- 991 World Health Organization, Global Tuberculosis Report;
- 992 https://www.who.int/tb/publications/global_report/en/