

Supplemental information

Global Biobank Meta-analysis Initiative:

Powering genetic discovery across human disease

Wei Zhou, Masahiro Kanai, Kuan-Han H. Wu, Humaira Rasheed, Kristin Tsuo, Jibril B. Hirbo, Ying Wang, Arjun Bhattacharya, Huiling Zhao, Shinichi Namba, Ida Surakka, Brooke N. Wolford, Valeria Lo Faro, Esteban A. Lopera-Maya, Kristi Läll, Marie-Julie Favé, Juulia J. Partanen, Sinéad B. Chapman, Juha Karjalainen, Mitja Kurki, Mutaamba Maasha, Ben M. Brumpton, Sameer Chavan, Tzu-Ting Chen, Michelle Daya, Yi Ding, Yen-Chen A. Feng, Lindsay A. Guare, Christopher R. Gignoux, Sarah E. Graham, Whitney E. Hornsby, Nathan Ingold, Said I. Ismail, Ruth Johnson, Triin Laisk, Kuang Lin, Jun Lv, Iona Y. Millwood, Sonia Moreno-Grau, Kisung Nam, Priit Palta, Anita Pandit, Michael H. Preuss, Chadi Saad, Shefali Setia-Verma, Unnur Thorsteinsdottir, Jasmina Uzunovic, Anurag Verma, Matthew Zawistowski, Xue Zhong, Nahla Afifi, Kawthar M. Al-Dabhani, Asma Al Thani, Yuki Bradford, Archie Campbell, Kristy Crooks, Geertruida H. de Bock, Scott M. Damrauer, Nicholas J. Douville, Sarah Finer, Lars G. Fritsche, Eleni Fthenou, Gilberto Gonzalez-Arroyo, Christopher J. Griffiths, Yu Guo, Karen A. Hunt, Alexander Ioannidis, Nomdo M. Jansonius, Takahiro Konuma, Ming Ta Michael Lee, Arturo Lopez-Pineda, Yuta Matsuda, Riccardo E. Marioni, Babak Moatamed, Marco A. Nava-Aguilar, Kensuke Numakura, Snehal Patil, Nicholas Rafaels, Anne Richmond, Agustin Rojas-Muñoz, Jonathan A. Shortt, Peter Straub, Ran Tao, Brett Vanderwerff, Manvi Vernekar, Yogasudha Veturi, Kathleen C. Barnes, Marike Boezen, Zhengming Chen, Chia-Yen Chen, Judy Cho, George Davey Smith, Hilary K. Finucane, Lude Franke, Eric R. Gamazon, Andrea Ganna, Tom R. Gaunt, Tian Ge, Hailiang Huang, Jennifer Huffman, Nicholas Katsanis, Jukka T. Koskela, Clara Lajonchere, Matthew H. Law, Liming Li, Cecilia M. Lindgren, Ruth J.F. Loos, Stuart MacGregor, Koichi Matsuda, Catherine M. Olsen, David J. Porteous, Jordan A. Shavit, Harold Snieder, Tomohiro Takano, Richard C. Trembath, Judith M. Vonk, David C. Whiteman, Stephen J. Wicks, Cisca Wijmenga, John Wright, Jie Zheng, Xiang Zhou, Philip Awadalla, Michael Boehnke, Carlos D. Bustamante, Nancy J. Cox, Segun Fatumo, Daniel H. Geschwind, Caroline Hayward, Kristian Hveem, Eimear E. Kenny, Seunggeun Lee, Yen-Feng Lin, Hamdi Mbarek, Reedik Mägi, Hilary C. Martin, Sarah E. Medland, Yukinori Okada, Aarno V. Palotie, Bogdan Pasaniuc, Daniel J. Rader, Marylyn D. Ritchie, Serena Sanna, Jordan W. Smoller, Kari Stefansson, David A. van Heel, Robin G. Walters, Sebastian Zöllner, Biobank of the Americas, Biobank Japan Project, BioMe, BioVU, CanPath - Ontario Health Study, China Kadoorie Biobank Collaborative Group, Colorado Center for Personalized Medicine, deCODE Genetics, Estonian Biobank, FinnGen, Generation Scotland, Genes & Health Research Team, LifeLines, Mass General Brigham Biobank, Michigan Genomics Initiative, National Biobank of Korea, Penn Medicine BioBank, Qatar Biobank, The QSkin Sun and Health Study, Taiwan Biobank, The HUNT Study, UCLA ATLAS

Figure S1. Disease prevalence varies across biobanks. Relates to Figure 1 and Table S2.

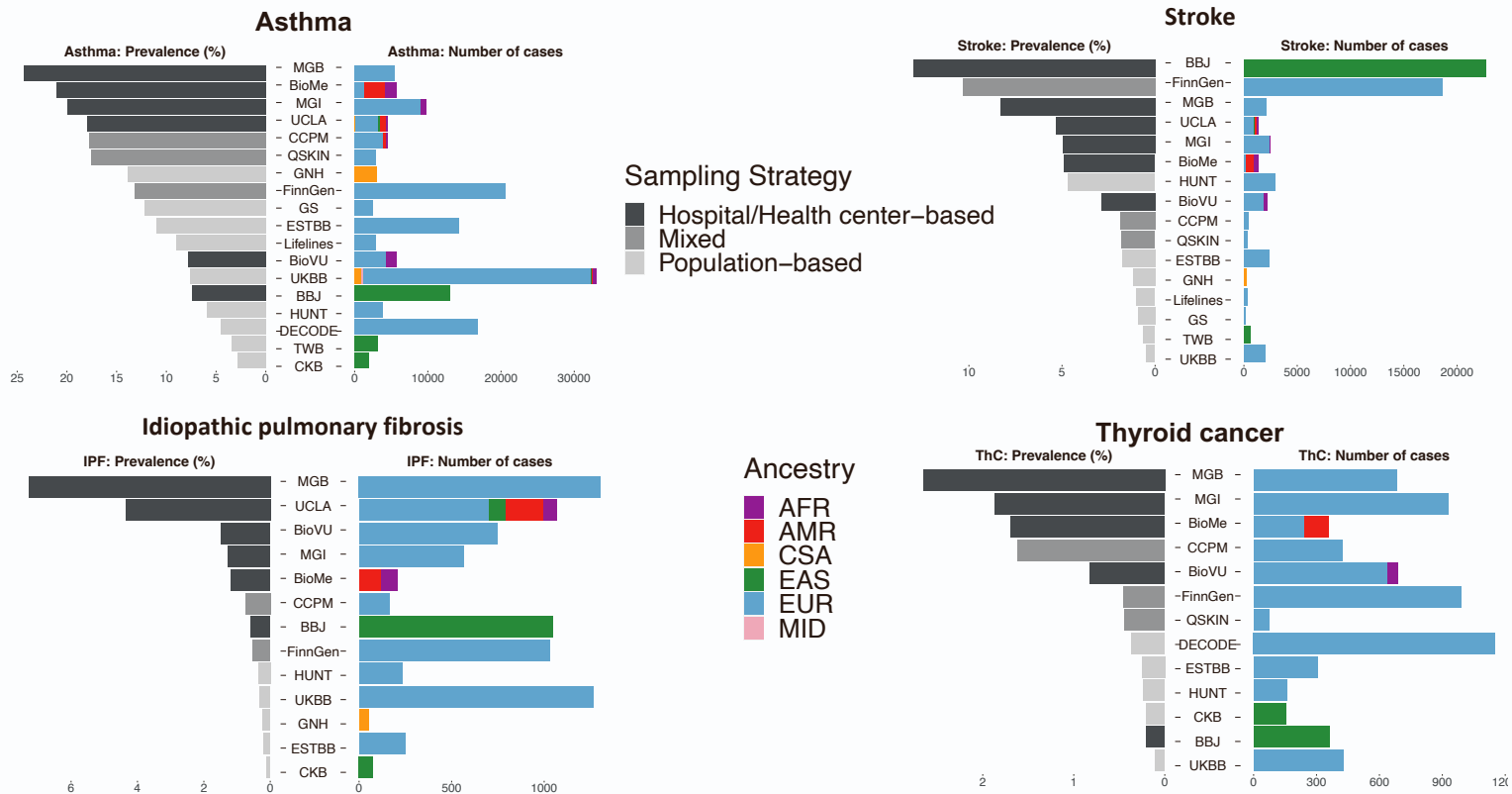


Figure S2. Disease prevalence varies by different sample recruiting strategies. A. Box plots for prevalence by three sampling strategies. B. Box plots to compare prevalence between population-based and hospital/health center-based biobanks. (**, $P < 0.01$, *, $P < 0.05$, unpaired Wilcoxon test). The center lines in box plots represent median and box limits are upper and lower quartiles. Relates to Figure 1 and Table S2.

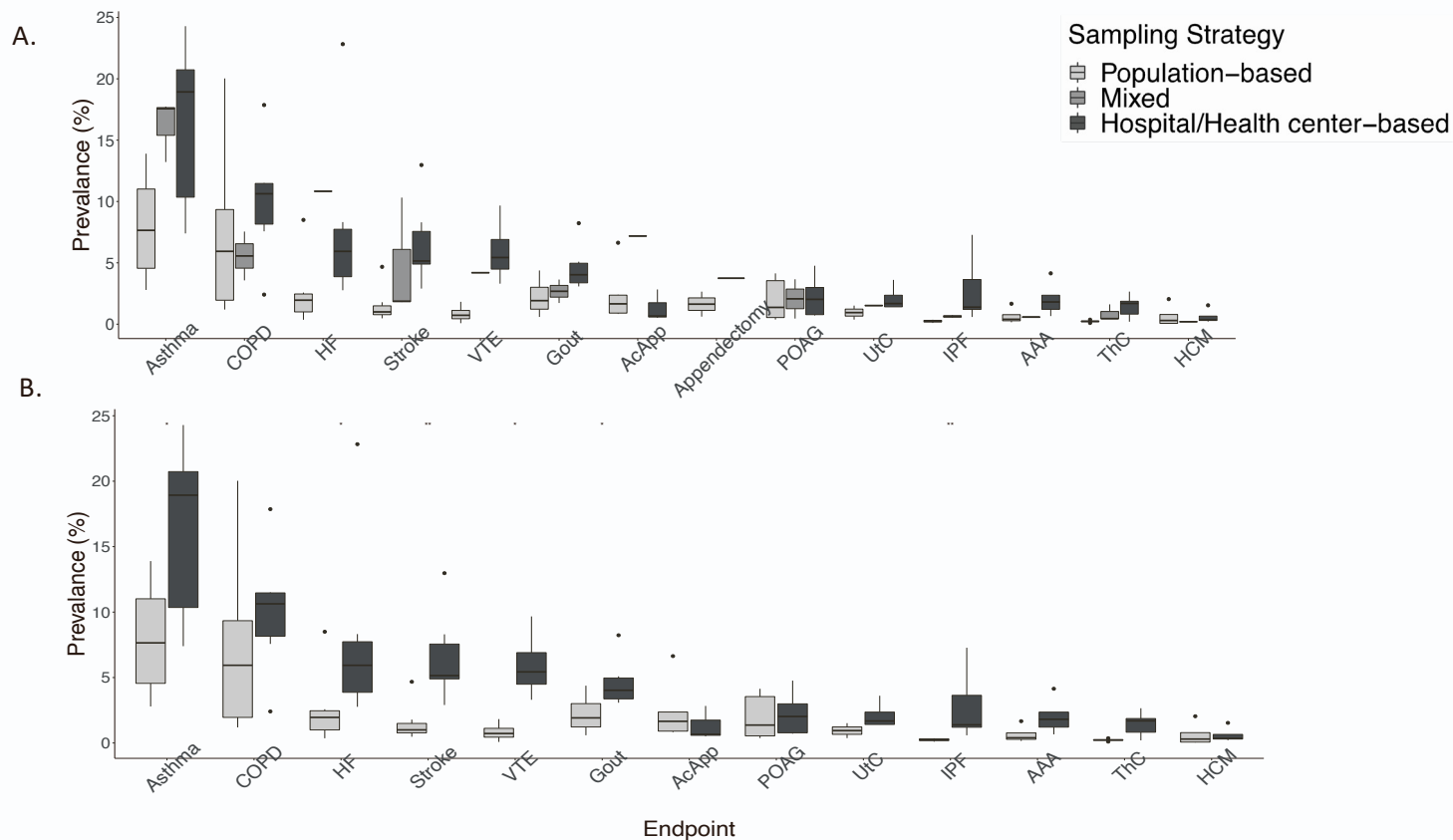
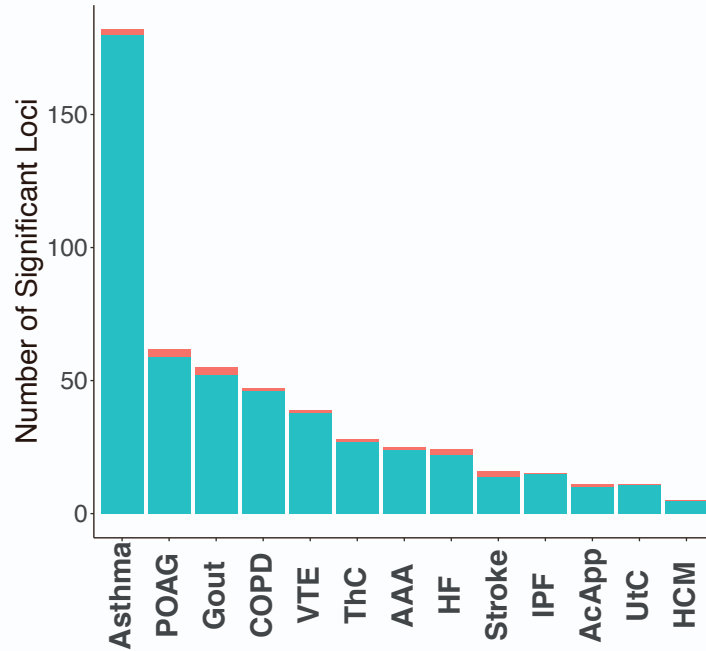


Figure S3. Additional significant loci identified by the meta-regression approach implemented in MR-MEGA¹¹ to account for effect size heterogeneity across different data sets in meta-analyses compared to the fixed-effect meta-analyses. Relates to Figure 5 and Table S10.



- Additional loci by accounting for heterogeneity of variant effects in different ancestries
- Loci identified by fixed-effect meta-analysis (assuming same variant effects in all ancestries)

Figure S4. A. Additional genetic variants analyzed due to incorporating non-European samples. EUR: genetic variants observed in samples with European ancestry. non EUR: genetic variants only observed in samples with non-European ancestry. The highest minor allele frequency (MAF) among non EUR ancestry was used in the plot. B. Distribution of the number of biobanks in which the genetic variants were tested. Relates to Figure 3, Table S3, and Table S6.

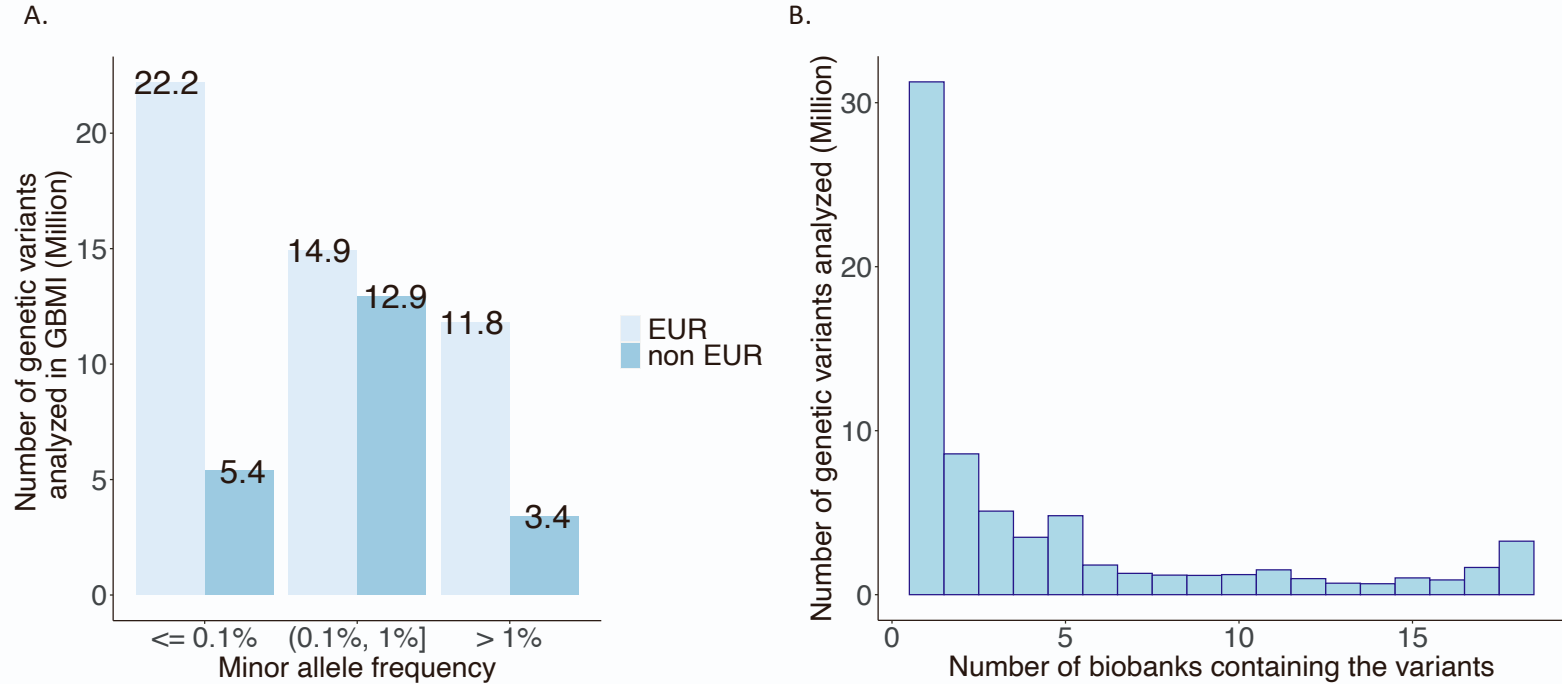
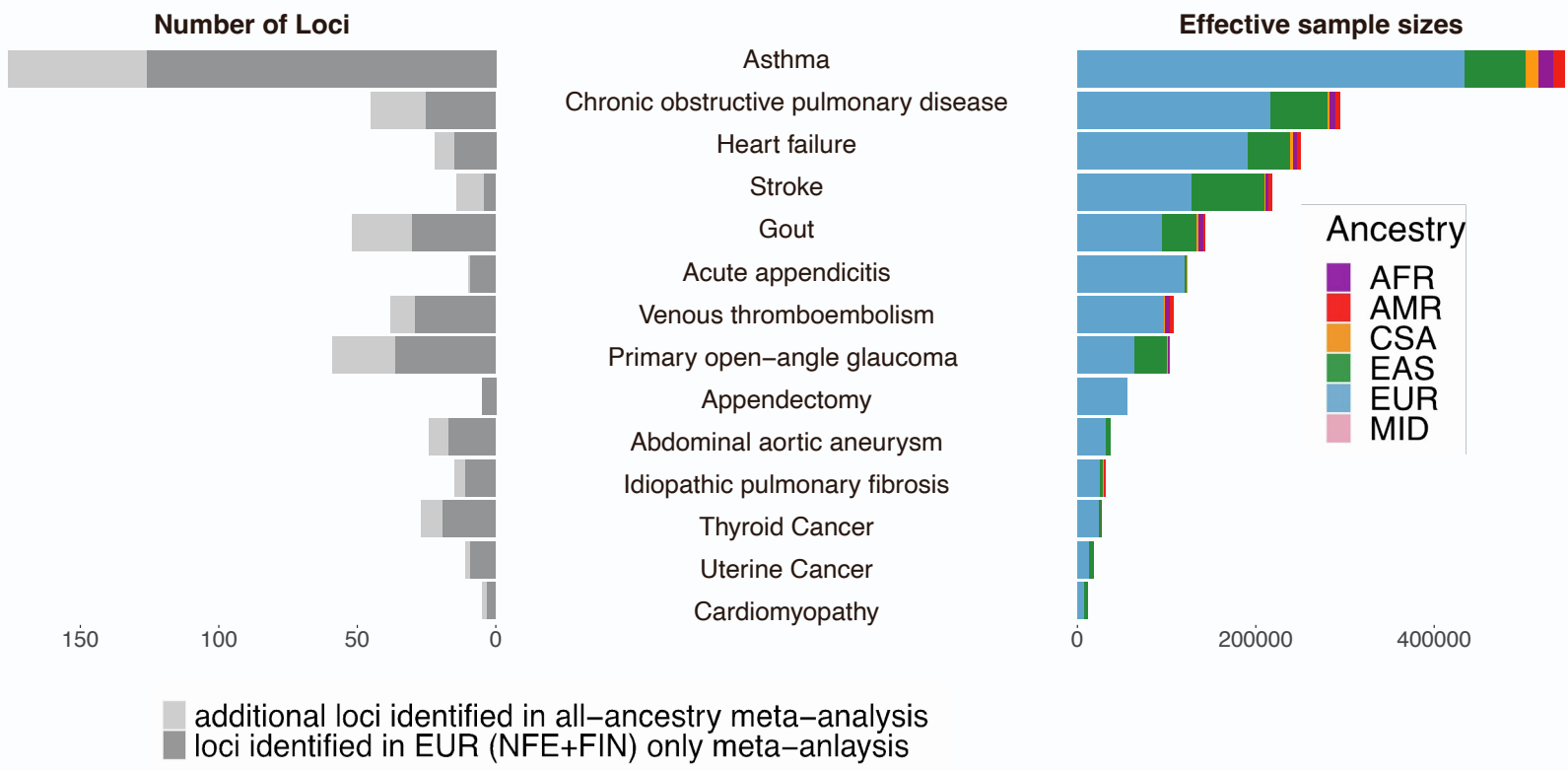


Figure S5. Additional significant loci identified when non-European samples were included in the meta-analysis. Relates to Figure 3 and Table S12.

A. Summary of loci by endpoints



B. Forest plots for additional loci identified when non-European samples were included in the meta-analysis. Error bars represent 95% confidence intervals of effect size estimates (log of odds ratios).

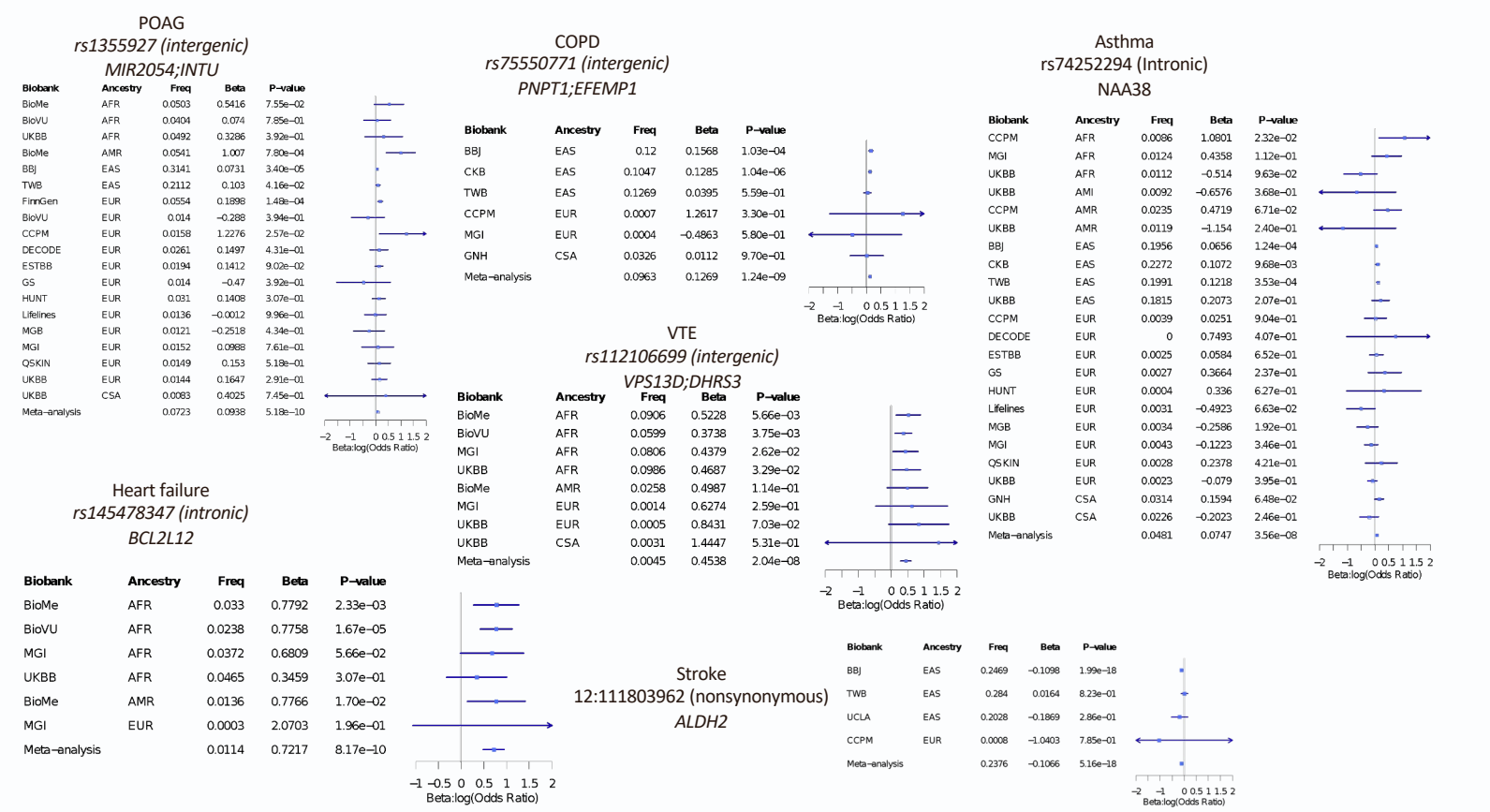


Figure S6. All-biobank meta-analysis results stratified by sex were filtered to identify regions with different effect sizes in men and women (P-value for Cochran's Q test < 0.002). Error bars represent 95% confidence intervals of effect size estimates (log of odds ratios). Relates to Figure 4 and Table S13.

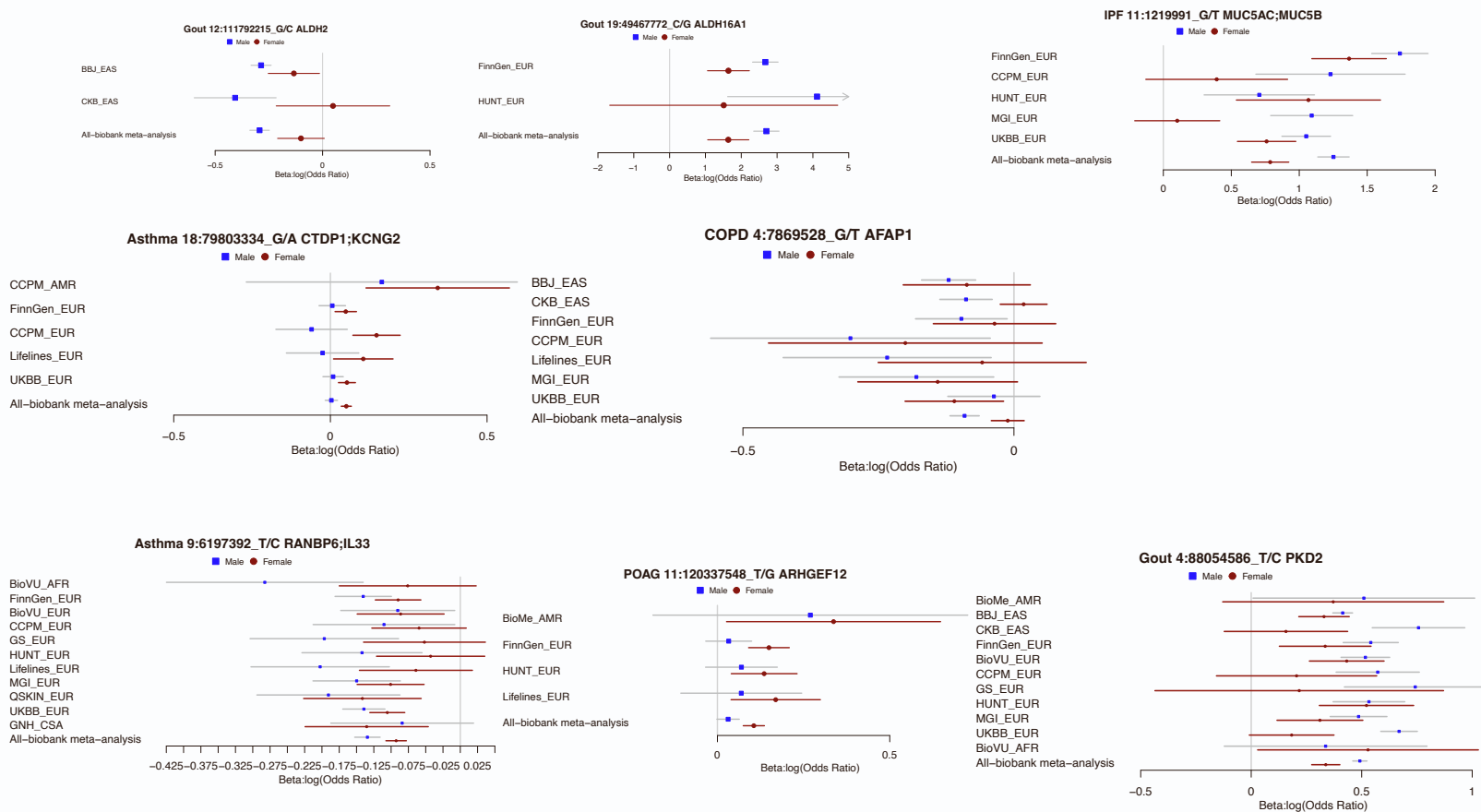


Figure S7. Plots of A. prevalence ratios and B. effective sample sizes ratios in females and in males by ancestry across endpoints. The prevalence is calculated as number of cases / (number of cases + number of controls). The effective sample sizes is calculated as $4 / (1 / \text{number of cases} + 1 / \text{number of controls})$. The center lines in box plots represent median and box limits are upper and lower quartiles. Relates to Figure 3 and Table S2.



Figure S8. The slopes of Deming regression for effect sizes for index variants in each biobank and leave-one-biobank-out meta-analysis (LOBO) pair are plotted against the effective sample sizes. Index variants with association p-values < 1×10^{-10} in the all-biobank meta-analysis were used for the regression. Biobanks, in which at least three index variants passed the cutoff, are plotted. Biobanks are annotated by phenotype source, sampling strategy and sample ancestry. The dotted line indicates $y=1$. A positive slope indicates that effect size estimates of the top hits are higher in the leave-one-biobank-out (LOBO) meta-analysis than in the individual biobank and a negative slope suggests lower effect size estimates in LOBO meta-analysis than in the individual biobank. The effective sample sizes is calculated as $4/(1/\text{case number} + 1/\text{control number})$. Error bars represent 95% confidence intervals of the slope estimates. Relates to Figure 4 and Methods.

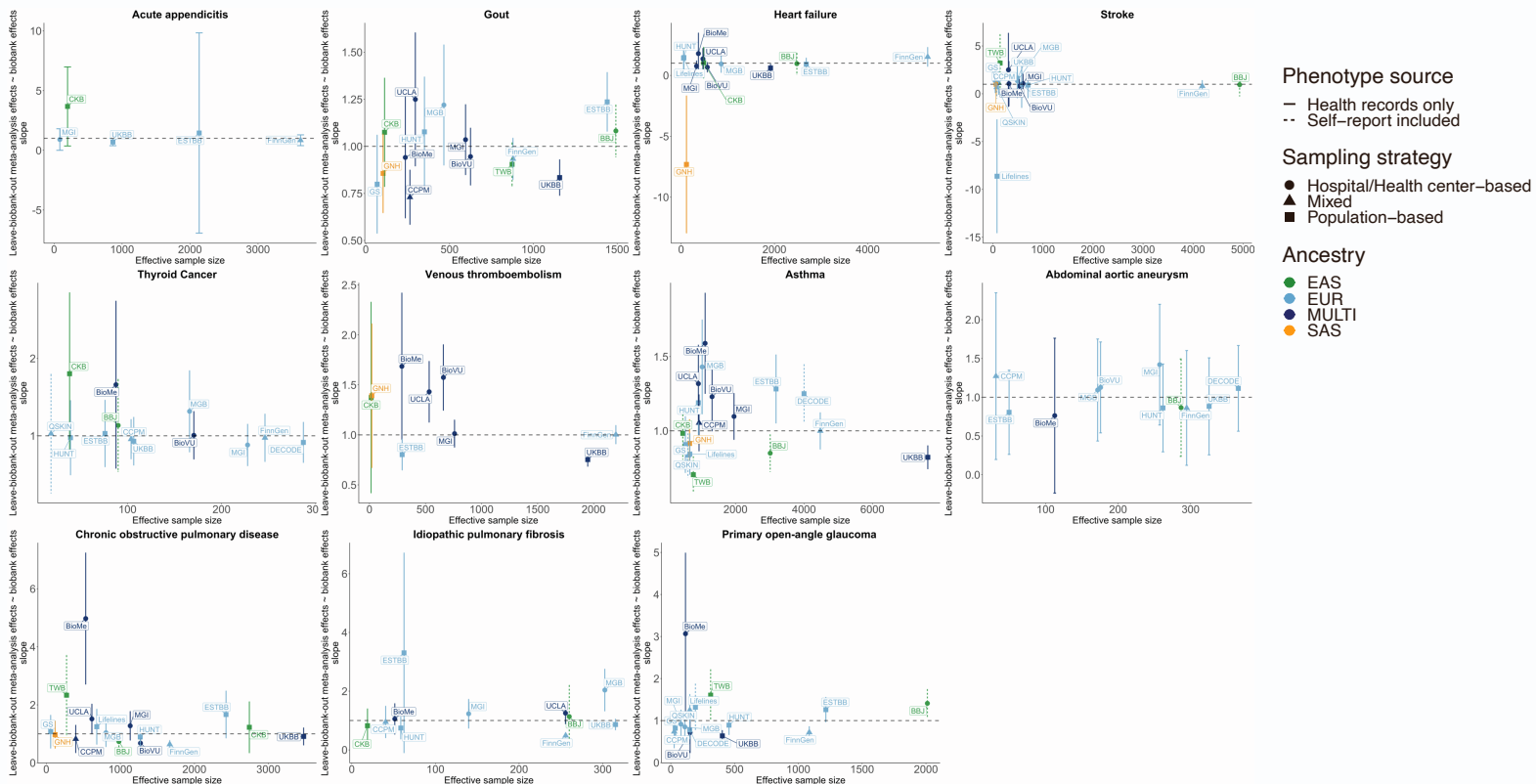


Figure S9. Genetic correlation between each biobank and leave-on-biobank meta-analysis in GBMI. Error bars represent 95% confidence intervals of genetic correlation estimates. Relates to Figure 4, Table S9, and Methods.

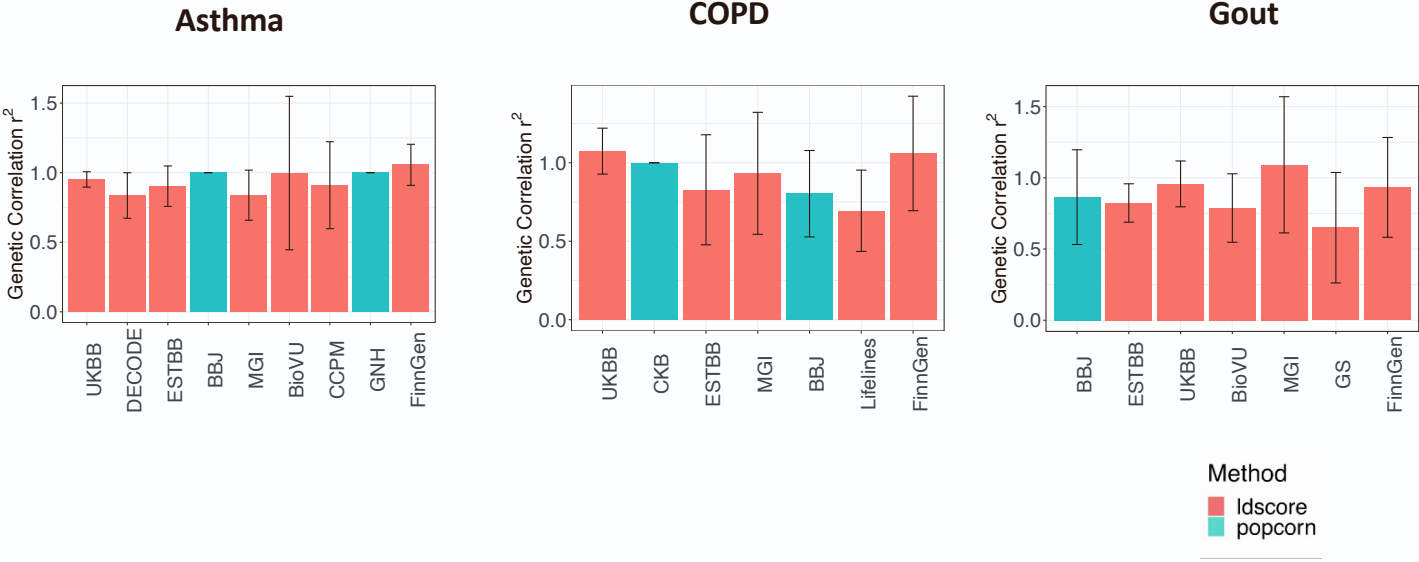


Figure S10. Scatter plots of the effect size estimates in population-based biobanks and hospital-/healthcare-based biobanks with the Deming regression lines and the slope estimates (intercepts were fixed to 0). Loci that were genome-wide significant in all-biobank meta-analyses and have p -value $< 1 \times 10^{-6}$ in both meta-analyses of population-based biobanks and hospital-/healthcare-based biobanks, respectively, were included in the analyses. Endpoints that have more than 5 loci included in the analyses were plotted. Error bars represent 95% confidence intervals of the effect size estimates. Relates to Figure 1.

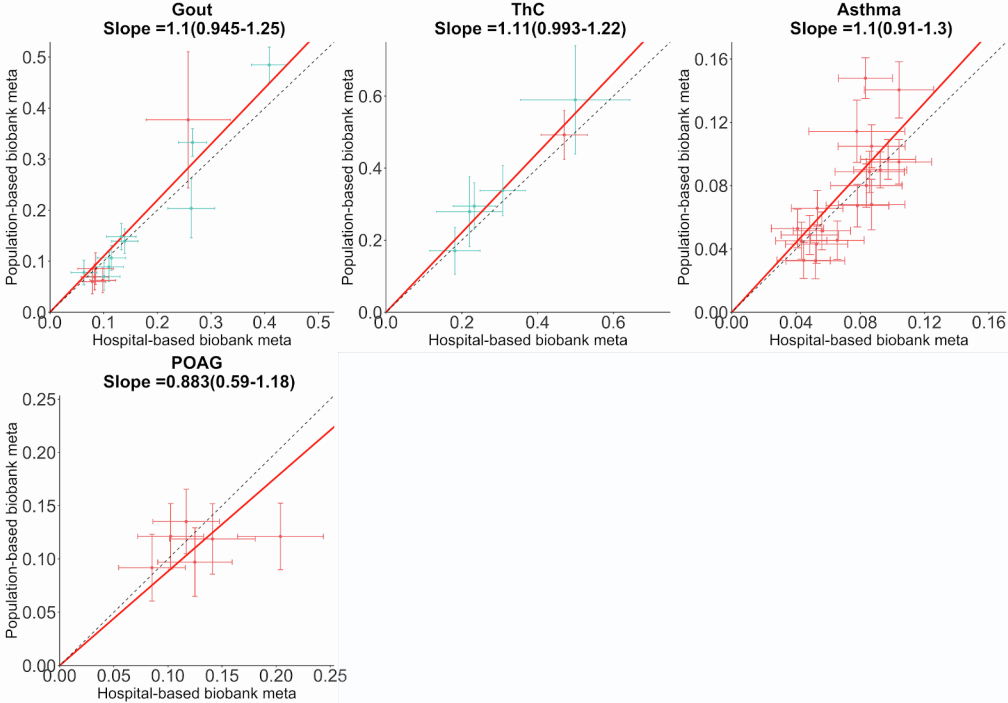


Figure S11. All 18 loci identified by previous GWAS for Asthma¹⁵ have more significant p-values in all-biobank meta-analysis. Relates to Table S15 and Methods.

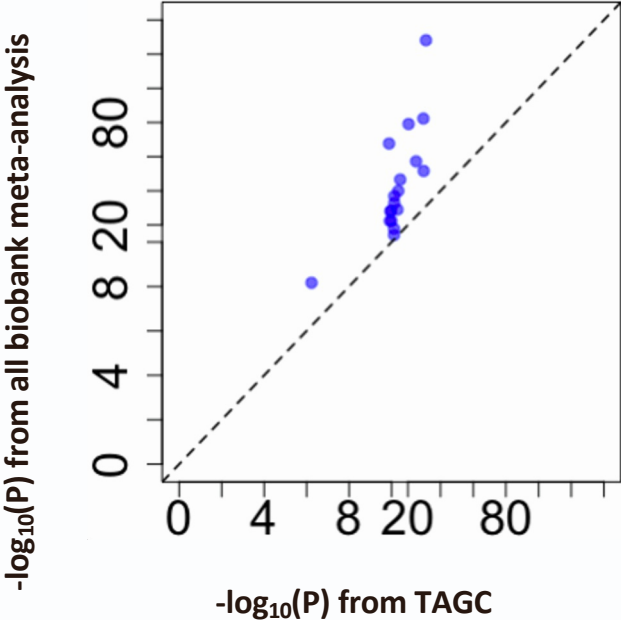


Figure S12. Number of genes prioritized by different methods: PoPs (top 1%), DEPICT (FDR < 0.05), TWAS (P < 2.5 x 10⁻⁶), PWMR (P < 0.001, Colocalization probability > 0.7), nearest genes around the top hits (Nearest gene, for intergenic variants, the nearest gene on each side will be included if both are located within 50kb from the top hit). Relates to Table S22 and Methods.

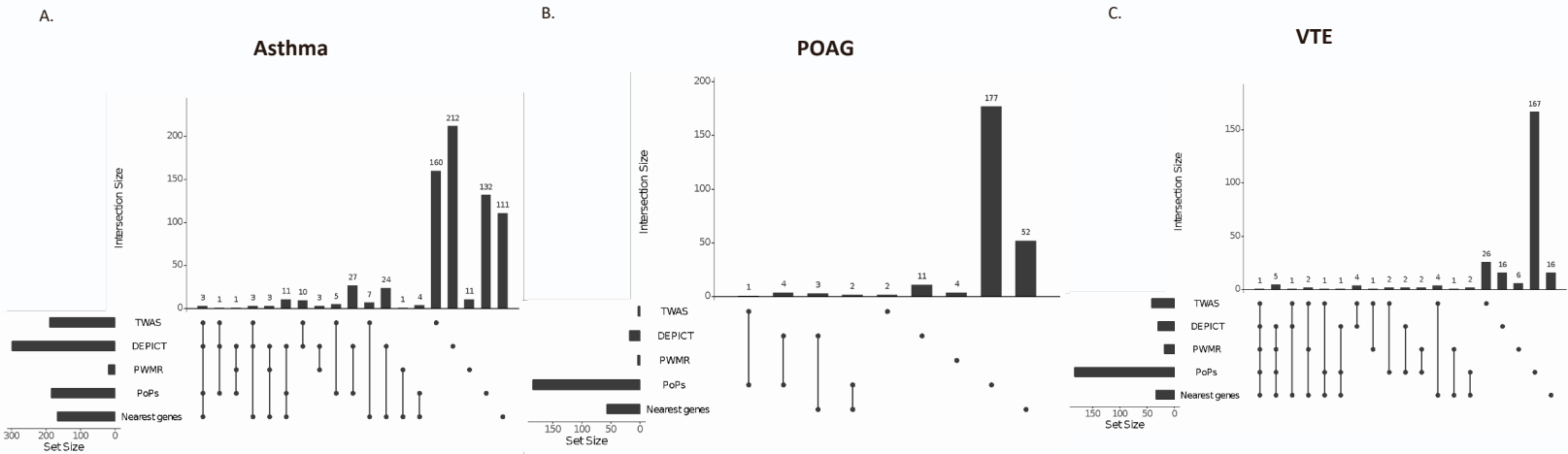


Figure S13. For VTE, A. number of genes prioritized by different methods: PoPs (top 1% and top 0.1%), DEPICT (FDR < 0.05), TWAS ($P < 2.5 \times 10^{-6}$), PWMR ($P < 0.001$, Colocalization probability > 0.7), nearest genes around the top hits (Nearest gene, for intergenic variants, the nearest gene on each side will be included if both are located within 50kb from the top hit). B. precision and recall of the different gene prioritization methods based on a gold standard set of 41 VTE genes that was curated prior to the meta-analysis by medical and molecular genetics experts in VTE²⁵. C. . precision and recall of the different gene prioritization methods based on 13 genes (*ADAMTS13,F10,F2,F5,F7,FGA,FGB,FGG,PROC,PROS1,PROZ,THBD,VWF*) in the gold standard sets that fall within 1Mb around VTE top hits in GBMI meta-analysis. Relates to Table S24 and Methods.

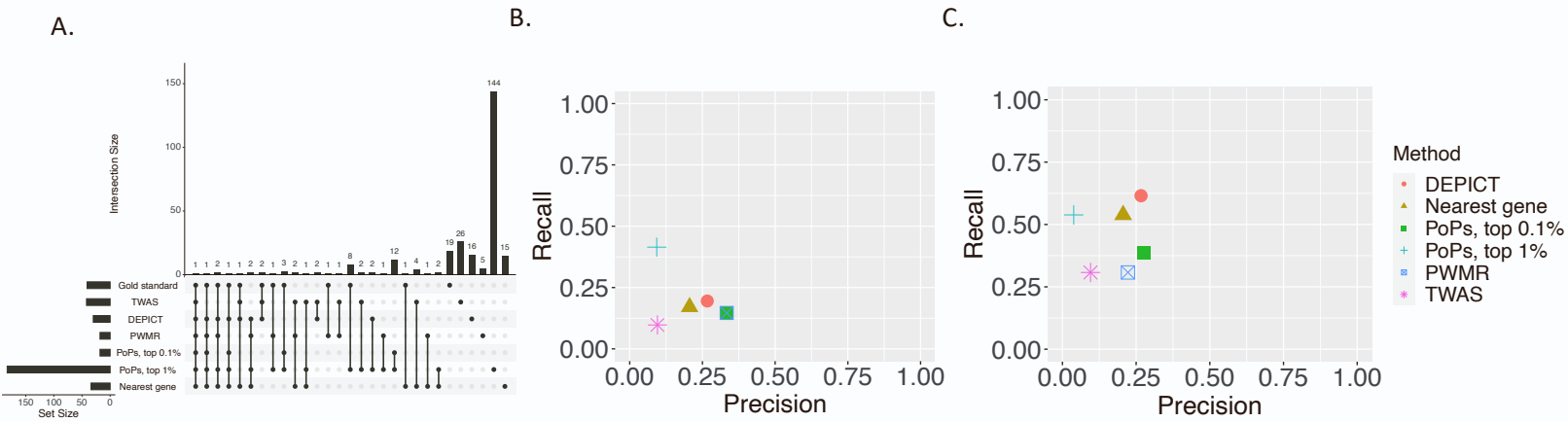


Figure S14. Improved polygenic risk scores (PRS) prediction accuracy using GBMI meta-analysis results compared to TAGC summary statistics. Error bars represent 95% confidence intervals. Relates to Methods.

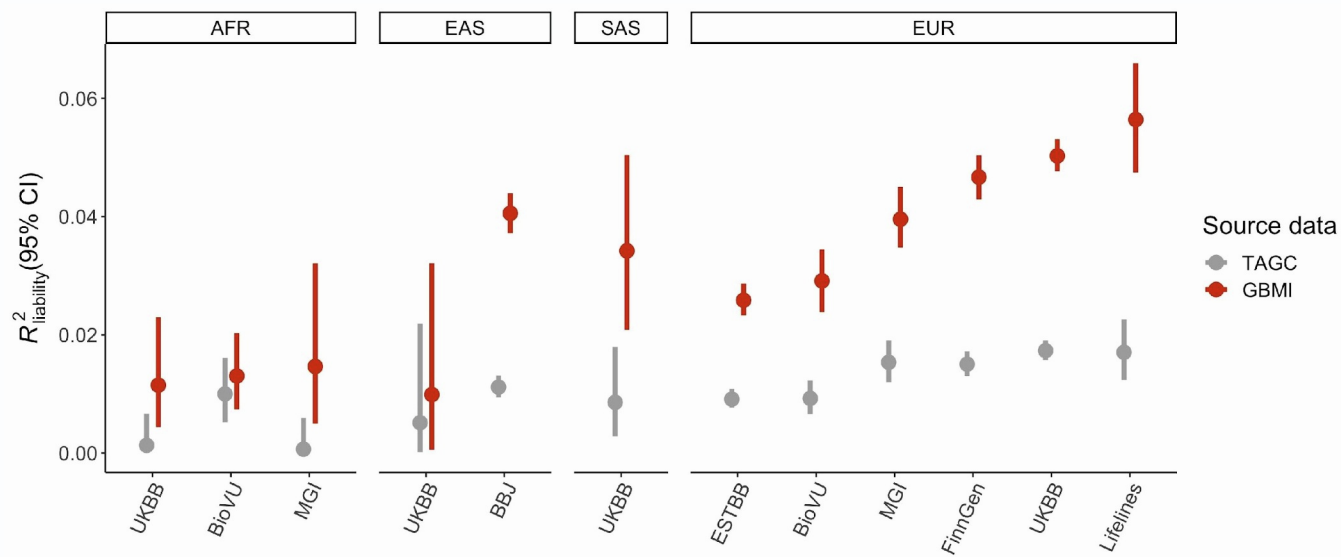


Figure S15. Palindromic SNPs with potential strand flip and genetic variants with different allele frequencies compared to gnomAD were flagged when included in the meta-analyses. Relates to Figure 4 and Methods.

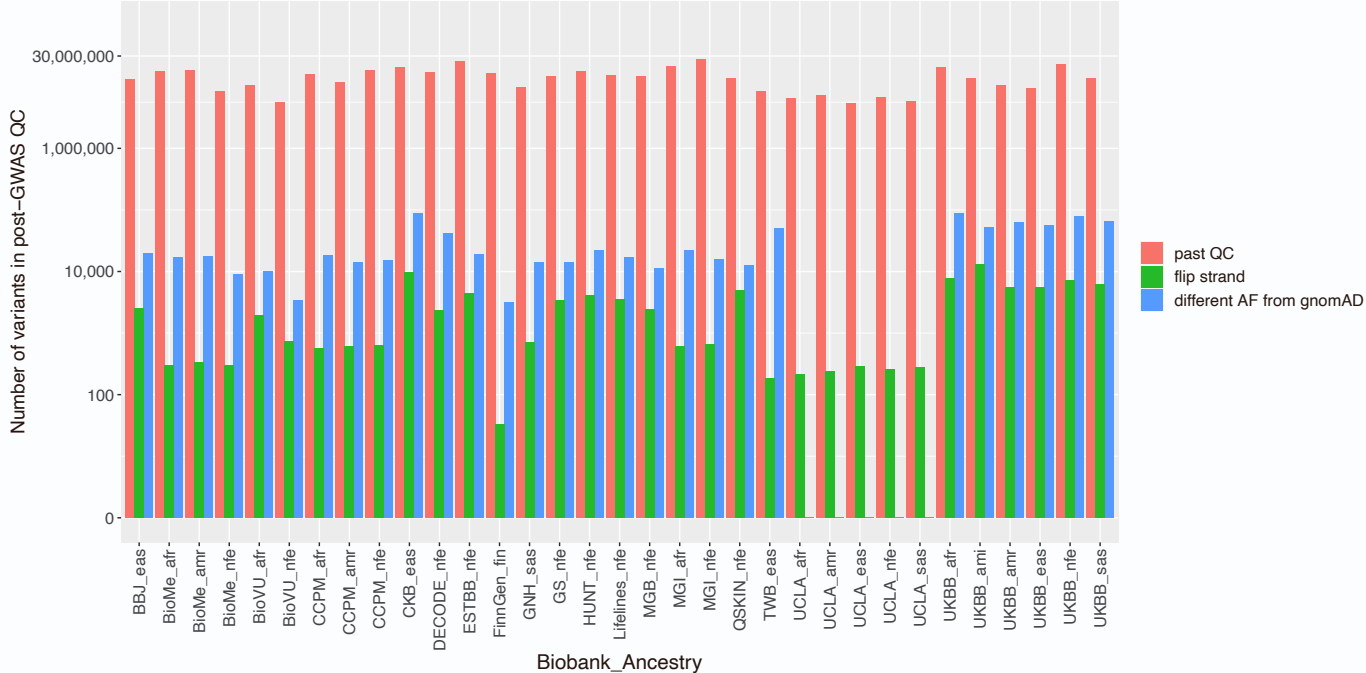
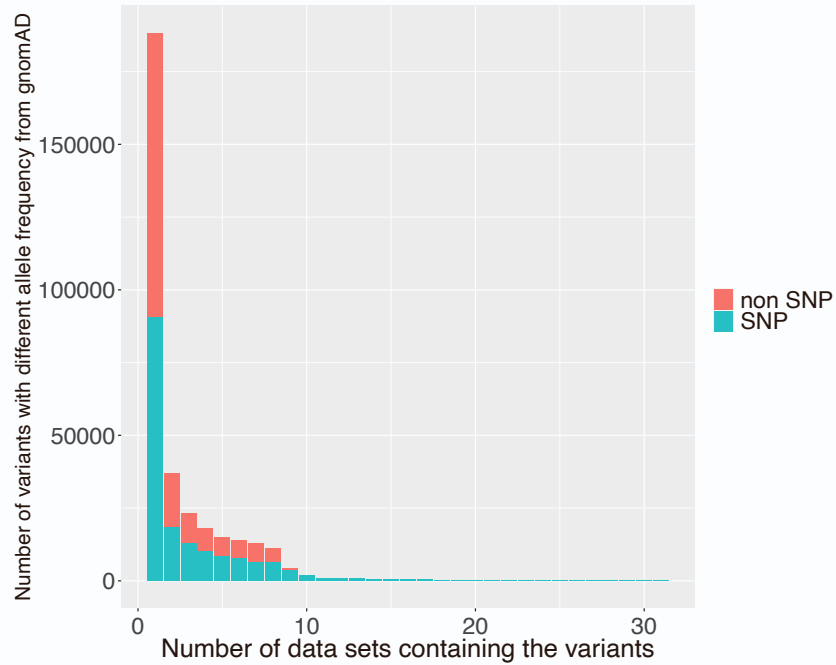


Figure S16. The distribution of number of biobanks that contain the genetic variants with different allele frequencies compared to gnomAD. Relates to Figure 4 and Methods.



Supplemental Reference list

1. Nagai A, Hirata M, Kamatani Y, et al. Overview of the BioBank Japan Project: Study design and profile. *J Epidemiol.* 2017;27(3S):S2-S8. doi:10.1016/j.je.2016.12.005
2. Abul-Husn NS, Soper ER, Braganza GT, et al. Implementing genomic screening in diverse populations. *Genome Med.* 2021;13(1):17. doi:10.1186/s13073-021-00832-y
3. Bowton EA, Collier SP, Wang X, et al. Phenotype-Driven Plasma Biobanking Strategies and Methods. *J Pers Med.* 2015;5(2):140-152. doi:10.3390/jpm5020140
4. Chen Z, Chen J, Collins R, et al. China Kadoorie Biobank of 0.5 million people: survey methods, baseline characteristics and long-term follow-up. *Int J Epidemiol.* 2011;40(6):1652-1666. doi:10.1093/ije/dyr120
5. Aquilante CL, Kao DP, Trinkley KE, et al. Clinical implementation of pharmacogenomics via a health system-wide research biobank: the University of Colorado experience. *Pharmacogenomics.* 2020;21(6):375-386. doi:10.2217/pgs-2020-0007
6. Gudbjartsson DF, Helgason H, Gudjonsson SA, et al. Large-scale whole-genome sequencing of the Icelandic population. *Nat Genet.* 2015;47(5):435-444. doi:10.1038/ng.3247
7. Leitsalu L, Haller T, Esko T, et al. Cohort Profile: Estonian Biobank of the Estonian Genome Center, University of Tartu. *Int J Epidemiol.* 2015;44(4):1137-1147. doi:10.1093/ije/dyt268
8. Kurki MI, Karjalainen J, Palta P, et al. FinnGen: Unique genetic insights from combining isolated population and national health register data. *bioRxiv.* Published online March 6, 2022. doi:10.1101/2022.03.03.22271360
9. Smith BH, Campbell A, Linksted P, et al. Cohort Profile: Generation Scotland: Scottish Family Health Study (GS:SFHS). The study, its participants and their potential for genetic research on health and illness. *Int J Epidemiol.* 2013;42(3):689-700. doi:10.1093/ije/dys084
10. Finer S, Martin HC, Khan A, et al. Cohort Profile: East London Genes & Health (ELGH), a community-based population genomics and health study in British Bangladeshi and British Pakistani people. *Int J Epidemiol.* 2020;49(1):20-21i. doi:10.1093/ije/dyz174
11. Scholtens S, Smidt N, Swertz MA, et al. Cohort Profile: Lifelines, a three-generation cohort study and biobank. *Int J Epidemiol.* 2015;44(4):1172-1180. doi:10.1093/ije/dyu229
12. Zawistowski, Fritsche, Pandit, et al. The Michigan Genomics Initiative: a biobank linking genotypes and electronic clinical records in Michigan Medicine patients. *In preparation.* Published online December 16, 2021. doi: 10.1101/2021.12.15.21267864

13. Karlson EW, Boutin NT, Hoffnagle AG, Allen NL. Building the Partners HealthCare Biobank at Partners Personalized Medicine: Informed Consent, Return of Research Results, Recruitment Lessons and Operational Considerations. *J Pers Med*. 2016;6(1). doi:10.3390/jpm6010002
14. Olsen CM, Green AC, Neale RE, et al. Cohort profile: the QSkin Sun and Health Study. *Int J Epidemiol*. 2012;41(4):929-929i. doi:10.1093/ije/dys107
15. Feng YCA, Chen CY, Chen TT, et al. Taiwan Biobank: a rich biomedical research database of the Taiwanese population. *medRxiv*. Published online 2021.
16. Brumpton BM, Graham S, Surakka I, et al. The HUNT Study: a population-based cohort for genetic research. *bioRxiv*. Published online December 25, 2021. doi:10.1101/2021.12.23.21268305
17. Johnson R, Ding Y, Venkateswaran V, et al. Leveraging genomic diversity for discovery in an EHR-linked biobank: the UCLA ATLAS Community Health Initiative. *bioRxiv*. Published online September 23, 2021. doi:10.1101/2021.09.22.21263987
18. Sudlow C, Gallacher J, Allen N, et al. UK biobank: an open access resource for identifying the causes of a wide range of complex diseases of middle and old age. *PLoS Med*. 2015;12(3):e1001779. doi:10.1371/journal.pmed.1001779
19. Dummer TJB, Awadalla P, Boileau C, et al. The Canadian Partnership for Tomorrow Project: a pan-Canadian platform for research on chronic disease prevention. *CMAJ*. 2018;190(23):E710-E717. doi:10.1503/cmaj.170292
20. Kim Y, Han BG, KoGES group. Cohort profile: The Korean genome and epidemiology study (KoGES) consortium. *Int J Epidemiol*. 2017;46(2):e20. doi:10.1093/ije/dyv316
21. Al Thani A, Fthenou E, Paparrodopoulos S, et al. Qatar Biobank cohort study: Study design and first results. *Am J Epidemiol*. 2019;188(8):1420-1433. doi:10.1093/aje/kwz084
22. Mbarek H, Devadoss Gandhi G, Selvaraj S, et al. Qatar genome: Insights on genomics from the Middle East. *Hum Mutat*. 2022;43(4):499-510. doi:10.1002/humu.24336
23. Lopez-Pineda A, Vernekar M, Grau SM, et al. Validating and automating learning of cardiometabolic polygenic risk scores from direct-to-consumer genetic and phenotypic data: implications for scaling precision health research. *bioRxiv*. Published online March 3, 2022. doi:10.1101/2022.03.01.22271722
24. Deming WE. Statistical adjustment of data. 1943;261. <https://psycnet.apa.org/fulltext/1944-00642-000.pdf>

25. Demenais F, Margaritte-Jeannin P, Barnes KC, et al. Multiancestry association study identifies new asthma risk loci that colocalize with immune-cell enhancer marks. *Nat Genet.* 2017;50(1):42-53. doi:10.1038/s41588-017-0014-7
26. Klarin D, Verma SS, Judy R, et al. Genetic Architecture of Abdominal Aortic Aneurysm in the Million Veteran Program. *Circulation.* 2020;142(17):1633-1646. doi:10.1161/CIRCULATIONAHA.120.047544
27. Tin A, Marten J, Halperin Kuhns VL, et al. Target genes, variants, tissues and transcriptional pathways influencing human serum urate levels. *Nat Genet.* 2019;51(10):1459-1474. doi:10.1038/s41588-019-0504-x
28. Partanen JJ, Häppölä P, Zhou W, et al. Leveraging global multi-ancestry meta-analysis in the study of Idiopathic Pulmonary Fibrosis genetics. *medRxiv.* Published online December 31, 2021:2021.12.29.21268310. doi:10.1101/2021.12.29.21268310
29. Tsuo K, Zhou W, Wang Y, et al. Multi-ancestry meta-analysis of asthma identifies novel associations and highlights the value of increased power and diversity. *medRxiv.* Published online December 7, 2021:2021.11.30.21267108. doi:10.1101/2021.11.30.21267108
30. Allen RJ, Guillen-Guio B, Oldham JM, et al. Genome-Wide Association Study of Susceptibility to Idiopathic Pulmonary Fibrosis. *Am J Respir Crit Care Med.* 2020;201(5):564-574. doi:10.1164/rccm.201905-1017OC
31. Wolford BN, Zhao Y, Surakka I, et al. Multi-ancestry GWAS for venous thromboembolism identifies novel loci followed by experimental validation in zebrafish. Published online June 27, 2022. doi:10.1101/2022.06.21.22276721
32. Thibord F, Klarin D, Brody JA, et al. Cross-Ancestry Investigation of Venous Thromboembolism Genomic Predictors. *medRxiv.* Published online March 8, 2022:2022.03.04.22271003. doi:10.1101/2022.03.04.22271003
33. Andersson RE, Olaison G, Tysk C, Ekbohm A. Appendectomy and protection against ulcerative colitis. *N Engl J Med.* 2001;344(11):808-814. doi:10.1056/NEJM200103153441104
34. Gill D, Cameron AC, Burgess S, et al. Urate, Blood Pressure, and Cardiovascular Disease: Evidence From Mendelian Randomization and Meta-Analysis of Clinical Trials. *Hypertension.* 2021;77(2):383-392. doi:10.1161/HYPERTENSIONAHA.120.16547
35. Sinnott-Armstrong N, Tanigawa Y, Amar D, et al. Genetics of 35 blood and urine biomarkers in the UK Biobank. *Nat Genet.* 2021;53(2):185-194. doi:10.1038/s41588-020-00757-z
36. Huffman JE, Albrecht E, Teumer A, et al. Modulation of genetic associations with serum urate levels by body-mass-index in humans. *PLoS One.* 2015;10(3):e0119752. doi:10.1371/journal.pone.0119752

37. Xu L, Shi Y, Zhuang S, Liu N. Recent advances on uric acid transporters. *Oncotarget*. 2017;8(59):100852-100862. doi:10.18632/oncotarget.20135
38. Terkeltaub R, Curtiss LK, Tenner AJ, Ginsberg MH. Lipoproteins containing apoprotein B are a major regulator of neutrophil responses to monosodium urate crystals. *J Clin Invest*. 1984;73(6):1719-1730. doi:10.1172/JCI111380
39. So AK, Martinon F. Inflammation in gout: mechanisms and therapeutic targets. *Nat Rev Rheumatol*. 2017;13(11):639-647. doi:10.1038/nrrheum.2017.155
40. Pers TH, Karjalainen JM, Chan Y, et al. Biological interpretation of genome-wide association studies using predicted gene functions. *Nat Commun*. 2015;6:5890. doi:10.1038/ncomms6890
41. Liu S, Xie Z, Daugherty A, et al. Mineralocorticoid receptor agonists induce mouse aortic aneurysm formation and rupture in the presence of high salt. *Arterioscler Thromb Vasc Biol*. 2013;33(7):1568-1579. doi:10.1161/ATVBAHA.112.300820
42. Lo Faro V, Bhattacharya A, Zhou W, et al. Genome-wide association meta-analysis identifies novel ancestry-specific primary open-angle glaucoma loci and shared biology with vascular mechanisms and cell proliferation. *medRxiv*. Published online January 17, 2022:2021.12.16.21267891. doi:10.1101/2021.12.16.21267891
43. Weeks EM, Ulirsch JC, Cheng NY, et al. Leveraging polygenic enrichments of gene features to predict genes underlying complex traits and diseases. *bioRxiv*. Published online September 10, 2020. doi:10.1101/2020.09.08.20190561
44. Bhattacharya A, Hirbo JB, Zhou D, et al. Best practices for multi-ancestry, meta-analytic transcriptome-wide association studies: lessons from the Global Biobank Meta-analysis Initiative. *medRxiv*. Published online November 29, 2021:2021.11.24.21266825. doi:10.1101/2021.11.24.21266825
45. Zhao H, Rasheed H, Nøst TH, et al. Proteome-wide Mendelian randomization in global biobank meta-analysis reveals multi-ancestry drug targets for common diseases. *medRxiv*. Published online January 11, 2022:2022.01.09.21268473. doi:10.1101/2022.01.09.21268473
46. Rabe KF, Brusselle G, Castro M, et al. Dupilumab shows rapid and sustained suppression of inflammatory biomarkers in corticosteroid (CS)-dependent severe asthma patients in LIBERTY ASTHMA VENTURE. In: *Allergy and Immunology*. European Respiratory Society; 2018. doi:10.1183/13993003.congress-2018.pa5003
47. Castro M, Corren J, Pavord ID, et al. Dupilumab Efficacy and Safety in Moderate-to-Severe Uncontrolled Asthma. *N Engl J Med*. 2018;378(26):2486-2496. doi:10.1056/NEJMoa1804092

48. DeVries A, Wlasiuk G, Miller SJ, et al. Epigenome-wide analysis links SMAD3 methylation at birth to asthma in children of asthmatic mothers. *J Allergy Clin Immunol.* 2017;140(2):534-542. doi:10.1016/j.jaci.2016.10.041