



Published in final edited form as:

*Cell Genom.* 2022 December 14; 2(12): . doi:10.1016/j.xgen.2022.100210.

## Meta-analysis fine-mapping is often miscalibrated at single-variant resolution

Masahiro Kanai<sup>1,2,3,4,5,7,\*</sup>, Roy Elzur<sup>1,2,3</sup>, Wei Zhou<sup>1,2,3</sup>,

Global Biobank Meta-analysis Initiative,

Mark J. Daly<sup>1,2,3,6</sup>, Hilary K. Finucane<sup>1,2,3,\*</sup>

<sup>1</sup>Analytic and Translational Genetics Unit, Massachusetts General Hospital, Boston, MA 02142, USA

<sup>2</sup>Program in Medical and Population Genetics, Broad Institute of MIT and Harvard, Cambridge, MA 02142, USA

<sup>3</sup>Stanley Center for Psychiatric Research, Broad Institute of MIT and Harvard, Cambridge, MA 02142, USA

<sup>4</sup>Department of Biomedical Informatics, Harvard Medical School, Boston, MA 02115, USA

<sup>5</sup>Department of Statistical Genetics, Osaka University Graduate School of Medicine, Suita 565-0871, Japan

<sup>6</sup>Institute for Molecular Medicine Finland (FIMM), University of Helsinki, Helsinki, Finland

<sup>7</sup>Lead contact

### SUMMARY

Meta-analysis is pervasively used to combine multiple genome-wide association studies (GWASs). Fine-mapping of meta-analysis studies is typically performed as in a single-cohort study. Here, we first demonstrate that heterogeneity (e.g., of sample size, phenotyping, imputation) hurts calibration of meta-analysis fine-mapping. We propose a summary statistics-based quality-control (QC) method, suspicious loci analysis of meta-analysis summary statistics (SLALOM), that identifies suspicious loci for meta-analysis fine-mapping by detecting outliers in association statistics. We validate SLALOM in simulations and the GWAS Catalog. Applying SLALOM to 14 meta-analyses from the Global Biobank Meta-analysis Initiative (GBMI), we find that 67% of loci show suspicious patterns that call into question fine-mapping accuracy. These predicted suspicious loci are significantly depleted for having nonsynonymous variants as lead variant ( $2.7\times$ ; Fisher's exact  $p = 7.3 \times 10^{-4}$ ). We find limited evidence of fine-mapping improvement in the GBMI

---

This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>).

\*Correspondence: mkanai@broadinstitute.org (M.K.), finucane@broadinstitute.org (H.K.F.).

#### AUTHOR CONTRIBUTIONS

M.K., M.J.D., and H.K.F. designed the study. M.K., R.E., and W.Z. performed analyses. H.K.F. supervised this work. H.K.F. and M.K. obtained funding. M.K., R.E., M.J.D., and H.K.F. wrote the manuscript with input from all authors.

#### SUPPLEMENTAL INFORMATION

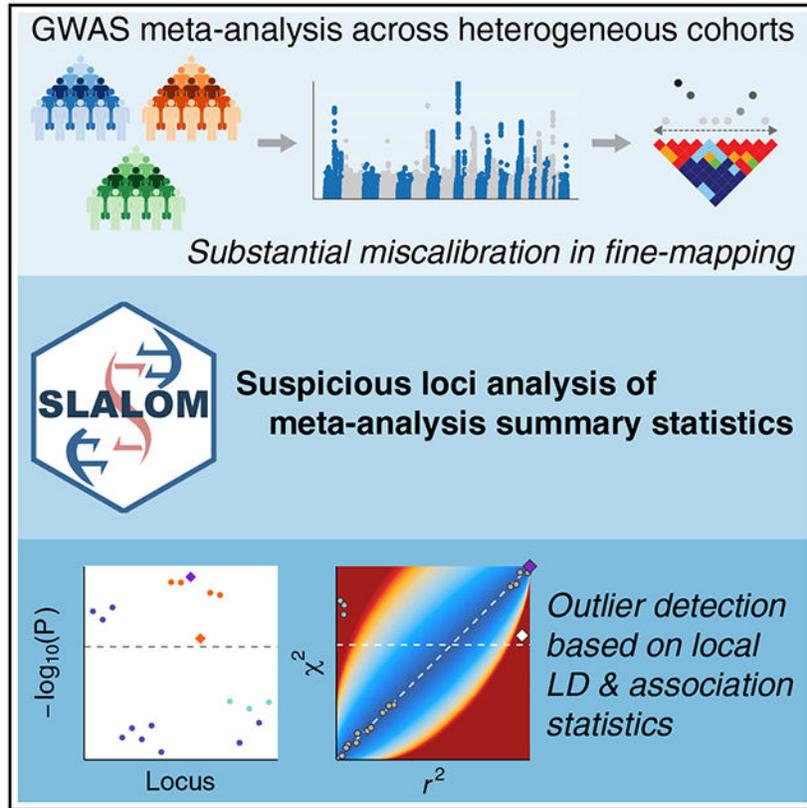
Supplemental information can be found online at <https://doi.org/10.1016/j.xgen.2022.100210>.

#### DECLARATION OF INTERESTS

M.J.D. is a founder of Maze Therapeutics. All other authors declare no competing interests.

meta-analyses compared with individual biobanks. We urge extreme caution when interpreting fine-mapping results from meta-analysis of heterogeneous cohorts.

**Graphical Abstract**



**In brief**

Genome-wide associations studies (GWASs), often performed as meta-analyses, have identified tens of thousands of disease-associated loci. Kanai et al. demonstrate via large-scale simulations and real data analysis that standard tools for pinpointing the causal variants underlying these associations can produce unreliable results when applied to GWAS meta-analyses.

**INTRODUCTION**

Meta-analysis is pervasively used to combine multiple genome-wide association studies (GWASs) from different cohorts.<sup>1</sup> Previous GWAS meta-analyses have identified thousands of loci associated with complex diseases and traits, such as type 2 diabetes,<sup>2,3</sup> schizophrenia,<sup>4,5</sup> rheumatoid arthritis,<sup>6,7</sup> body mass index,<sup>8</sup> and lipid levels.<sup>9</sup> These meta-analyses are typically conducted in large-scale consortia (e.g., the Psychiatric Genomics Consortium [PGC] and the Genetic Investigation of Anthropometric Traits [GIANT] consortium) to increase sample size while harmonizing analysis plans across participating cohorts in every possible aspect (e.g., phenotype definition, quality-control [QC] criteria, statistical model, and analytical software) by sharing summary statistics as opposed to

individual-level data, thereby avoiding data protection issues and variable legal frameworks governing individual genome and medical data around the world. The Global Biobank Meta-analysis Initiative (GBMI)<sup>10</sup> is one such large-scale, international effort, which aims to establish a collaborative network spanning 23 biobanks from four continents (total n = 2.2 million) for coordinated GWAS meta-analyses, while addressing the many benefits and challenges in meta-analysis and subsequent downstream analyses.

One such challenging downstream analysis is statistical fine-mapping.<sup>11-13</sup> Despite the great success of past GWAS meta-analyses in locus discovery, individual causal variants in associated loci are largely unresolved. Identifying causal variants from GWAS associations (i.e., fine-mapping) is challenging due to extensive linkage disequilibrium (LD, the correlation among genetic variants), the presence of multiple causal variants, and limited sample sizes, but is rapidly becoming achievable with high confidence in individual cohorts<sup>14-16</sup> owing to the recent development of large-scale biobanks<sup>17-19</sup> and scalable fine-mapping methods<sup>20-22</sup> that enable well-powered, accurate fine-mapping using in-sample LD from large-scale individual-level data.

After conducting GWAS meta-analysis, previous studies<sup>2,7,9,23-29</sup> have applied existing summary statistics-based fine-mapping methods (e.g., approximate Bayes factor [ABF],<sup>30,31</sup> CAVIAR,<sup>32</sup> PAINTOR,<sup>33,34</sup> FINEMAP,<sup>20,21</sup> and SuSiE<sup>22</sup>) just as they are applied to single-cohort studies, without considering or accounting for the unavoidable heterogeneity among cohorts (e.g. differences in sample size, phenotyping, genotyping, or imputation). Such heterogeneity could lead to false-positives and miscalibration in meta-analysis fine-mapping (Figure 1). For example, case-control studies enriched with more severe cases or ascertained with different phenotyping criteria may disproportionately contribute to genetic discovery, even when true causal effects for genetic liability are exactly the same between these studies and less severe or unascertained ones. Quantitative traits such as biomarkers could have phenotypic heterogeneity arising from different measurement protocols and errors across studies. There might be genuine biological mechanisms too, such as gene-gene (GxG) and gene-environment (GxE) interactions and (population-specific) dominance variation (e.g., rs671 and alcohol dependence<sup>35</sup>), that introduce additional heterogeneity across studies.<sup>36,37</sup> In addition to phenotyping, differences in genotyping and imputation could dramatically undermine fine-mapping calibration and recall at single-variant resolution, because differential patterns of missingness and imputation quality across constituent cohorts of different sample sizes can disproportionately diminish association statistics of potentially causal variants. Finally, although more easily harmonized than phenotyping and genotyping data, subtle differences in QC criteria and analytical software may further exacerbate the effect of heterogeneity on fine-mapping.

An illustrative example of such issues can be observed in the *TYK2* locus (19p13.2) in the recent meta-analysis from the COVID-19 Host Genetics Initiative (COVID-19 HGI; Figure S1).<sup>38</sup> This locus is known for protective associations against autoimmune diseases,<sup>6,23</sup> while a complete *TYK2* loss of function results in a primary immunodeficiency.<sup>39</sup> Despite strong LD ( $r^2 = 0.82$ ) with a lead variant in the locus (rs74956615;  $p = 9.7 \times 10^{-12}$ ), a known functional missense variant rs34536443 (p.Pro1104Ala) that reduces *TYK2* function<sup>40,41</sup> did not achieve genome-wide significance and was assigned a very low

posterior inclusion probability (PIP) in fine-mapping ( $p = 7.5 \times 10^{-7}$ ;  $\text{PIP} = 9.5 \times 10^{-4}$ ), primarily due to its missingness in two more cohorts than rs74956615. This serves as just one example of the major difficulties with meta-analysis fine-mapping at single-variant resolution. Indeed, the COVID-19 HGI cautiously avoided an *in silico* fine-mapping in the flagship to prevent spurious results.<sup>38</sup>

Only a few studies have carefully addressed these concerns in their downstream analyses. The Schizophrenia Working Group of PGC, for example, recently updated their largest meta-analysis of schizophrenia<sup>5</sup> (69,369 cases and 236,642 controls), followed by a downstream fine-mapping analysis using FINEMAP.<sup>20</sup> Unlike many other GWAS consortia, since PGC has access to individual-level genotypes for a majority of samples, they were able to apply standardized sample and variant QC criteria and impute variants using the same reference panel, all uniformly processed using the RICOPILI pipeline.<sup>42</sup> This harmonized procedure was crucial for properly controlling inter-cohort heterogeneity and thus allowing more robust meta-analysis fine-mapping at single-variant resolution. Furthermore, PGC's direct access to individual-level data enabled them to compute in-sample LD matrices for multiple-causal-variant fine-mapping, which prevents the significant miscalibration that results from using an external LD.<sup>14,15</sup> A 2017 fine-mapping study of inflammatory bowel disease also benefited from access to individual-level genotypes and careful pre- and post-fine-mapping QC.<sup>43</sup> For a typical meta-analysis consortium, however, many of these steps are infeasible as full genotype data from all cohorts are not available. For such studies, a new approach to meta-analysis fine-mapping in the presence of the many types of heterogeneity is needed. Until such a method is developed, QC of meta-analysis fine-mapping results deserves increased attention.

While existing variant-level QC procedures are effective for limiting spurious associations in GWAS (Data S1),<sup>44</sup> they do not suffice for ensuring high-quality fine-mapping results. In some cases, they even hurt fine-mapping quality, because they can (1) cause or exacerbate differential patterns of missing variants across cohorts, and (2) remove true causal variants as well as suspicious variants. Thus, additional QC procedures that retain consistent variants across cohorts for consideration but limit poor-quality fine-mapping results are needed. A recently proposed method called DENTIST,<sup>45</sup> for example, performs summary statistics QC to improve GWAS downstream analyses, such as conditional and joint analysis (GCTA-COJO<sup>46</sup>), by removing variants based on estimated heterogeneity between summary statistics and reference LD. Although DENTIST was also applied prior to fine-mapping (FINEMAP<sup>20</sup>), simulations only demonstrated that it could improve power for detecting the correct number of causal variants in a locus, not true causal variants. This motivated us to develop a new fine-mapping QC method for better calibration and recall at single-variant resolution and to demonstrate its performance in large-scale meta-analysis.

Here, we first demonstrate the effect of inter-cohort heterogeneity in meta-analysis fine-mapping via realistic simulations with multiple heterogeneous cohorts, each with different combinations of genotyping platforms, imputation reference panels, and genetic ancestries. We propose a summary statistics-based QC method, suspicious loci analysis of meta-analysis summary statistics (SLALOM), that identifies suspicious loci for meta-analysis fine-mapping by detecting association statistics outliers based on local LD structure,

building on the DENTIST method. Applying SLALOM to 14 disease endpoints from the GBMI<sup>10</sup> as well as 467 meta-analysis summary statistics from the GWAS Catalog,<sup>47</sup> we demonstrate that suspicious loci for fine-mapping are widespread in meta-analysis and urge extreme caution when interpreting fine-mapping results from meta-analysis.

## RESULTS

### Large-scale simulations demonstrate miscalibration in meta-analysis fine-mapping

Existing fine-mapping methods<sup>20,22,30</sup> assume that all association statistics are derived from a single-cohort study, and thus do not model the per-variant heterogeneity in effect sizes and sample sizes that arise when meta-analyzing multiple cohorts (Figure 1). To evaluate how different characteristics of constituent cohorts in a meta-analysis affect fine-mapping calibration and recall, we conducted a series of large-scale GWAS meta-analysis and fine-mapping simulations (Tables S1-S4; STAR Methods). Briefly, we simulated multiple GWAS cohorts of different ancestries (10 European ancestry, one African ancestry, and one East Asian ancestry cohorts;  $n = 10,000$  each) that were genotyped and imputed using different genotyping arrays (Illumina Omni2.5, Multi-Ethnic Global Array [MEGA], and Global Screening Array [GSA]) and imputation reference panels (the 1000 Genomes Project Phase 3 [1000GP3],<sup>48</sup> the Haplo-type Reference Consortium [HRC],<sup>49</sup> and the TOPMed<sup>50</sup>). For each combination of cohort, genotyping array, and imputation panel, we conducted 300 GWAS with randomly simulated causal variants that resemble the genetic architecture of a typical complex trait, including minor allele frequency (MAF) dependent causal effect sizes,<sup>51</sup> total SNP heritability,<sup>52</sup> functional consequences of causal variants,<sup>16</sup> and levels of genetic correlation across cohorts (i.e., true effect size heterogeneity;  $r_g = 1, 0.9, \text{ and } 0.5$ ; STAR Methods). We then meta-analyzed the single-cohort GWAS results across 10 independent cohorts based on multiple configurations (different combinations of genotyping arrays and imputation panels for each cohort) to resemble realistic meta-analysis of multiple heterogeneous cohorts (Table S4). We applied ABF fine-mapping to compute a PIP for each variant and to derive 95% and 99% credible sets (CSs) that contain the smallest set of variants covering 95% and 99% of probability of causality. We evaluated the false discovery rate (FDR, defined as the proportion of variants with  $\text{PIP} > 0.9$  that are non-causal) and compared against the expected proportion of non-causal variants if the meta-analysis fine-mapping method were calibrated, based on PIP. More details of our simulation pipeline are described in STAR Methods and visually summarized in Figure S2.

We found that FDR varied widely over the different configurations, reaching as high as 37% for the most heterogeneous configurations (Figure 2). We characterized the contributing factors to the miscalibration. We first found that lower true effect size correlation  $r_g$  (i.e., larger phenotypic heterogeneity) always caused higher miscalibration and lower recall. Second, when using the same imputation panel (1000GP3), use of less dense arrays (MEGA or GSA) led to moderately inflated FDR (up to  $\text{FDR} = 11\%$  versus expected 1%), while use of multiple genotyping array did not cause further FDR inflation (Figure 2C). Third, when using the same genotyping array (Omni2.5), use of imputation panels (HRC or TOPMed) that do not match our simulation reference significantly affects miscalibration (up to  $\text{FDR} = 17\%$  versus expected 1%), and using multiple imputation panels further increased



only 393,471 variants (12%) out of all the QC-passing 3,285,617 variants were available in every combination (Figure S4B). These observations recapitulate that different combinations of genetic ancestry, genotyping array, imputation panels, and QC thresholds substantially affect the availability of common, well-imputed variants for association testing.<sup>55</sup>

Thus, the different combinations of genotyping and imputation cause each cohort in a meta-analysis to have a different set of variants, and consequently variants can have very different overall sample sizes. In our simulations with the most heterogeneous configurations, we found that 66% of the false-positive loci (where a non-causal [false-positive] variant was assigned  $PIP > 0.9$ ) had different sample sizes for true causal and false-positive variants (median maximum/minimum sample size ratio = 1.4; Figure S6). Analytically, we found that at common meta-analysis sample sizes and genome-wide significant effect size regimes, when two variants have similar marginal effects, the one with the larger sample size will usually achieve a higher ABF PIP (Data S2; Figures S7-S9). This elucidates the mechanism by which sample size imbalance can lead to miscalibration.

### Overview of the SLALOM method

To address the challenges in meta-analysis fine-mapping discussed above, we developed SLALOM, a method that flags suspicious loci for meta-analysis fine-mapping by detecting outliers in association statistics based on deviations from expectation, estimated with local LD structure (STAR Methods). SLALOM consists of three steps: (1) defining loci and lead variants based on a 1 Mb window, (2) detecting outlier variants in each locus using meta-analysis summary statistics and an external LD reference panel, and (3) identifying suspicious loci for meta-analysis fine-mapping (Figures 3A and 3B).

To detect outlier variants, we first assume a single causal variant per associated locus. Then the marginal  $Z$  score  $z_i$  for a variant  $i$  should be approximately equal to  $r_{i,c} z_c$  where  $z_c$  is the  $Z$  score of the causal variant  $c$ , and  $r_{i,c}$  is a correlation between variants  $i$  and  $c$ . For each variant in meta-analysis summary statistics, we first test this relationship using a simplified version of the DENTIST statistics,<sup>45</sup> DENTIST-S, based on the assumption of a single causal variant. The DENTIST-S statistics for a given variant  $i$  is written as

$$T_i = \frac{(z_i - r_{i,c} z_c)^2}{1 - r_{i,c}^2} \quad (\text{Equation 1})$$

which approximately follows a distribution with 1 degree of freedom.<sup>45</sup> Since the true causal variant and LD structure are unknown in real data, we approximate the causal variant as the lead PIP variant in the locus (the variant with the highest PIP) and use a large-scale external LD reference from gnomAD,<sup>56</sup> either an ancestry-matched LD for a single-ancestry meta-analysis or a sample-size-weighted LD by ancestries for a multi-ancestry meta-analysis (STAR Methods). We note that the existence of multiple independent causal variants in a locus would not affect SLALOM precision but would decrease recall (see section “discussion”).

SLALOM then evaluates whether each locus is “suspicious”; that is, has a pattern of meta-analysis statistics and LD that appear inconsistent and therefore call into question the

fine-mapping accuracy. By training on loci with maximum PIP  $>0.9$  in the simulations, we determined that the best-performing criterion for classifying loci as true- or false-positives is whether a locus has a variant with  $r^2 > 0.6$  to the lead and DENTIST-S p-value  $< 1.0 \times 10^{-4}$  (STAR Methods). Using this criterion, we achieved an area under the receiver operating characteristic curve (AUROC) of 0.74, 0.76, and 0.80 for identifying whether a true causal variant is a lead PIP variant, in 95% CS, and in 99% CS, respectively (Figure 3C). Using different thresholds, we observed that the SLALOM performance is not very sensitive to thresholds near the threshold we chose (Figure S10). We further validated the performance of SLALOM using all the loci in the simulations and observed significantly higher miscalibration in predicted suspicious loci than in non-suspicious loci (up to 16% difference in FDR at PIP  $>0.9$ ; Figure 3D). We found that SLALOM-predicted suspicious loci tend to be from more heterogeneous configurations and the SLALOM sensitivity and specificity depend on the level of heterogeneity (Table S5). Given the lower miscalibration and specificity at low PIP thresholds (Figures 3D and 3E), in subsequent real data analysis we restricted the application of SLALOM to loci with maximum PIP  $>0.1$  (STAR Methods).

### Widespread suspicious loci for fine-mapping in existing meta-analysis summary statistics

Having assessed the performance of SLALOM in simulations, we applied SLALOM to 467 meta-analysis summary statistics in the GWAS Catalog<sup>47</sup> that are publicly available with a sufficient discovery sample size ( $N > 10,000$ ; Table S6; STAR Methods) to quantify the prevalence of suspicious loci in existing studies. These summary statistics were mostly European-ancestry-only meta-analyses (63%), followed by multi-ancestry (31%), East Asian ancestry-only (3%), and African ancestry-only (2%) meta-analyses. Across 467 summary statistics from 96 publications, we identified 28,925 loci with maximum PIP  $>0.1$  (out of 35,864 genome-wide significant loci defined based on 1-Mb window around lead variants; STAR Methods) for SLALOM analysis, of which 8,137 loci (28%) were predicted suspicious (Table S7).

To validate SLALOM performance in real data, we restricted our analysis to 6,065 loci that have maximum PIP  $>0.1$  and that contain nonsynonymous coding variants (predicted loss of function [pLoF] and missense) in LD with the lead variant ( $r^2 > 0.6$ ). Given prior evidence<sup>16,43,57</sup> that such nonsynonymous variants are highly enriched for being causal, we tested the validity of our method by whether they achieve the highest PIP in the locus (i.e., successful fine-mapping) in suspicious versus non-suspicious loci (STAR Methods). While 40% (1,557 out of 3,860) of non-suspicious loci successfully fine-mapped nonsynonymous variants, only 17% (384 out of 2,205) of suspicious loci did, demonstrating a significant depletion (2.3 $\times$ ) of successfully fine-mapped nonsynonymous variants in suspicious loci (Fisher's exact  $p = 3.6 \times 10^{-79}$ ; Figure 4A). We also tested whether nonsynonymous variants belonged to 95% and 99% CS and again observed significant depletion (1.4 $\times$  and 1.3 $\times$ , respectively; Fisher's exact  $p < 4.6 \times 10^{-100}$ ). In addition, when we used a more stringent  $r^2$  threshold ( $>0.8$ ) for selecting loci that contain nonsynonymous variants, we also confirmed significant enrichment (Fisher's exact  $p < 6.1 \times 10^{-65}$ ; Figure S11). To quantify potential fine-mapping miscalibration in the GWAS Catalog, we investigated the difference between mean PIP for lead variants and fraction of lead variants that are nonsynonymous; assuming that nonsynonymous variants in these loci are truly causal, this difference equals

the difference between the true and reported fraction of lead PIP variants that are causal. We observed differences between 26%–51% and 10%–18% under different PIP thresholds in suspicious and non-suspicious loci, respectively (Figure 4B), marking 45% and 15% for high-PIP (>0.9) variants.

We further assessed SLALOM performance in the GWAS Catalog meta-analyses by leveraging high-PIP (>0.9) complex trait and *cis*-eQTL variants that were rigorously fine-mapped<sup>16</sup> in large-scale biobanks (Biobank Japan [BBJ],<sup>58</sup> FinnGen,<sup>19</sup> and UK Biobank [UKBB]<sup>18</sup>) and eQTL resources (GTEx<sup>59</sup> v8 and eQTL Catalog<sup>60</sup>). Among the 27,713 loci analyzed by SLALOM (maximum PIP >0.1) that contain a lead variant that was included in biobank fine-mapping, 17% (3,266 out of 19,692) of the non-suspicious loci successfully fine-mapped one of the high-PIP GWAS variants in biobank fine-mapping, whereas 7% (589 out of 8,021) of suspicious loci did, showing a significant depletion (2.3×) of the high-PIP complex trait variants in suspicious loci (Fisher's exact  $p = 4.6 \times 10^{-100}$ ; Figure 4C). Similarly, among 26,901 loci analyzed by SLALOM that contain a lead variant that was included in *cis*-eQTL fine-mapping, we found a significant depletion (1.9×) of the high-PIP *cis*-eQTL variants in suspicious loci, where 7% (1,247 out of 18,976) of non-suspicious loci versus 4% (281 out of 7,925) of suspicious loci successfully fine-mapped one of the high-PIP *cis*-eQTL variants (Fisher's exact  $p = 2.6 \times 10^{-24}$ ; Figure 4D). We observed the same significant depletions of the high-PIP complex trait and *cis*-eQTL variants in suspicious loci that belonged to 95% and 99% CS set (Figures 4C and 4D).

### Suspicious loci for fine-mapping in the GBMI summary statistics

Next, we applied SLALOM to meta-analysis summary statistics of 14 disease endpoints from the GBMI.<sup>10</sup> These summary statistics were generated from a meta-analysis of up to 1.8 million individuals in total across 18 biobanks for discovery, representing six different genetic ancestry groups of approximately 33,000 African, 18,000 admixed American, 31,000 Central and South Asian, 341,000 East Asian, 1.4 million European, and 1,600 Middle Eastern individuals (Table S8). Among 489 genome-wide significant loci across the 14 traits (excluding the major histocompatibility complex [MHC] region; STAR Methods), we found that 82 loci (17%) showed maximum PIP <0.1, thus not being further considered by SLALOM. Of the remaining 407 loci with maximum PIP >0.1, SLALOM identified that 272 loci (67%) were suspicious loci for fine-mapping (Figure 5A; Table S9). The fraction of suspicious loci and their maximum PIP varied by trait, reflecting different levels of statistical power (e.g., sample sizes, heritability, and local LD structure) as well as inter-cohort heterogeneity (Figures 5B-5O).

While the fraction of suspicious loci (67%) in the GBMI meta-analyses is higher than in the GWAS Catalog (28%), there might be multiple reasons for this discrepancy, including association significance, sample size, ancestral diversity, and study-specific QC criteria. For example, the GBMI summary statistics were generated from multi-ancestry, large-scale meta-analyses of median sample size of 1.4 million individuals across six ancestries, while 63% of the 467 summary statistics from the GWAS Catalog were only in European-ancestry studies and 83% had less than 0.5 million discovery samples. Nonetheless, predicted suspicious loci for fine-mapping were prevalent in both the GWAS Catalog and the GBMI.

Using nonsynonymous (pLoF and missense) and high-PIP (>0.9) complex trait and *cis*-eQTL variants, we recapitulated a significant depletion of these likely causal variants in predicted suspicious loci (2.7×, 5.2×, and 5.1× for nonsynonymous, high-PIP complex trait, and high-PIP *cis*-eQTL variants being a lead PIP variant, respectively; Fisher's exact  $p < 7.3 \times 10^{-4}$ ), confirming our observation in the GWAS Catalog analysis (Figures 6A-6C).

In 15 out of 23 non-suspicious loci harboring a nonsynonymous variant, the nonsynonymous variant had the highest PIP. These included known missense variants such as rs116483731 (p.Arg20Gln) in *SPDL1* for idiopathic pulmonary fibrosis (IPF)<sup>61,62</sup> and rs28929474 (p.Glu366Lys) in *SERPINA1* for chronic obstructive pulmonary disease (COPD).<sup>63,64</sup> In addition, we observed successful fine-mapping in two novel loci for asthma: (1) rs41286560 (p.Pro558Thr) in *RTL1*, a missense variant known for decreasing height<sup>65,66</sup>, and (2) rs34187696 (p.Gly337Val) in *ZSCAN5A*, a known missense variant for increasing monocyte count.<sup>29</sup>

To characterize fine-mapping failures in suspicious loci, we examined suspicious loci in which a nonsynonymous variant did not achieve the highest PIP. For example, the *FCGR2A/FCGR3A* (1q23.3) locus for COPD contained a genome-wide significant lead intergenic variant rs2099684 ( $p = 1.7 \times 10^{-11}$ ), which is in LD ( $r^2 = 0.92$ ) with a missense variant rs396991 (p.Phe176Val) of *FCGR3A*; Figure 6D). This locus was not previously reported for COPD but is known for associations with autoimmune diseases (e.g., inflammatory bowel disease,<sup>43</sup> rheumatoid arthritis,<sup>7</sup> and systemic lupus erythematosus<sup>67</sup>) and encodes the low-affinity human FC-gamma receptors that bind to the Fc region of immunoglobulin (Ig) G and activate immune responses.<sup>68</sup> Notably, this locus contains copy number variations that contribute to the disease associations in addition to single-nucleotide variants, which makes genotyping challenging.<sup>68,69</sup> Despite strong LD with the lead variant, rs396991 did not achieve genome-wide significance ( $p = 9.1 \times 10^{-3}$ ), showing a significant deviation from the expected association ( $P_{\text{DENTIST-S}} = 5.3 \times 10^{-41}$ ; Figure 6E). This is primarily due to missingness of rs396991 in eight biobanks out of 17 ( $N_{\text{eff}} = 76,790$  and 36,781 for rs2099684 and rs396991, respectively; Figure 6F), which is caused by its absence from major imputation reference panels (e.g., 1000GP,<sup>48</sup> HRC,<sup>49</sup> and UK10K<sup>70</sup>) despite having a high MAF in every population (MAF = 0.24–0.34 in African, admixed American, East Asian, European, and South Asian populations of gnomAD<sup>56</sup>).

Sample size imbalance across variants was pervasive in the GBMI meta-analyses,<sup>71</sup> and was especially enriched in predicted suspicious loci: 84% of suspicious loci versus 24% of non-suspicious loci showed a maximum/minimum effective sample size ratio >2 among variants in LD ( $r^2 > 0.6$ ) with lead variants (a median ratio = 4.2 and 1.2 in suspicious and non-suspicious loci, respectively; Figure S12). These observations are consistent with our simulations, recapitulating that sample size imbalance results in miscalibration for meta-analysis fine-mapping. Notably, we observed a similar issue in other GBMI downstream analyses (e.g., polygenic risk score [PRS]<sup>71</sup> and drug discovery<sup>72</sup>), where predictive performance improved significantly after filtering out variants with maximum  $N_{\text{eff}} < 50\%$ . Although fine-mapping methods cannot simply take this approach because it inevitably reduces calibration and recall by removing true causal variants, other meta-

analysis downstream analyses that primarily rely on polygenic signals rather than individual variants should consider this filtering as an extra QC step.

### Comparison of fine-mapping results between the GBMI meta-analyses and individual biobanks

Motivated by successful validation of SLALOM performance, we investigated whether fine-mapping confidence and resolution were improved in the GBMI meta-analyses over individual biobanks. To this end, we used our fine-mapping results<sup>16</sup> of nine disease endpoints (asthma,<sup>64</sup> COPD,<sup>64</sup> gout, heart failure,<sup>73</sup> IPF,<sup>62</sup> primary open-angle glaucoma,<sup>74</sup> thyroid cancer, stroke,<sup>75</sup> and venous thromboembolism<sup>76</sup>) in BBJ,<sup>58</sup> FinnGen,<sup>19</sup> and UKBB<sup>18</sup> Europeans that also contributed to the GBMI meta-analyses for the same traits.

To perform an unbiased comparison of PIP between the GBMI meta-analysis and individual biobanks, we investigated functional enrichment of fine-mapped variants based on top PIP rankings in the GBMI and individual biobanks (top 0.5%, 0.1%, and 0.05% PIP variants in the GBMI versus maximum PIP across BBJ, FinnGen, and UKBB; STAR Methods). Previous studies have shown that high-PIP (>0.9) complex trait variants are significantly enriched for well-known functional categories, such as coding (pLoF, missense, and synonymous), 5'/3' UTR, promoter, and *cis*-regulatory element (CRE) regions (DNase I hypersensitive sites and H3K27ac).<sup>16</sup> Using these functional categories, we found no significant enrichment of variants in the top PIP rankings in the GBMI over individual biobanks (Fisher's exact  $p > 0.05$ ; Figure 7A) except for variants in the promoter region (1.8 $\times$ ; Fisher's exact  $p = 4.9 \times 10^{-4}$  for the top 0.1% PIP variants). We observed similar trends regardless of whether variants were in suspicious or non-suspicious loci (Figures 7B and 7C). To examine patterns of increased and decreased PIP for individual variants, we also calculated PIP difference between the GBMI and individual biobanks, defined as  $\Delta\text{PIP} = \text{PIP}(\text{GBMI}) - \text{maximum PIP across biobanks}$  (Figures S13 and S14). We investigated functional enrichment based on  $\Delta\text{PIP}$  bins and observed inconsistent enrichment results using different PIP thresholds (Figure S15). Finally, to test whether fine-mapping resolution was improved in the GBMI over individual biobanks, we compared the size of 95% CS after restricting them to cases where a GBMI CS overlapped with an individual biobank CS (STAR Methods). We observed the median 95% CS size of 2 and 2 in non-suspicious loci for the GBMI and individual biobanks, respectively, and 5 and 14 in suspicious loci, respectively (Figure S16). The smaller CS size in suspicious loci in GBMI could be due to improved resolution or to increased miscalibration. These results provide limited evidence of overall fine-mapping improvement in the GBMI meta-analyses over what is achievable by taking the best result from individual biobanks.

Individual examples, however, provide insights into the types of fine-mapping differences that can occur. To characterize the observed differences in fine-mapping confidence and resolution, we further examined non-suspicious loci with  $\Delta\text{PIP} > 0.5$  in asthma. In some cases, the increased power and/or ancestral diversity of GBMI led to improved fine-mapping: for example, an intergenic variant rs1888909 (~18 kb upstream of *IL33*) showed  $\Delta\text{PIP} = 0.99$  ( $\text{PIP} = 1.0$  and  $0.008$  in GBMI and FinnGen, respectively; Figure 7D), which was primarily owing to increased association significance in a meta-analysis

( $p = 3.0 \times 10^{-86}$ ,  $7.4 \times 10^{-2}$ ,  $3.6 \times 10^{-16}$ , and  $1.9 \times 10^{-53}$  in GBMI, BBJ, FinnGen, and UKBB Europeans, respectively) as well as a shorter LD length in the African population than in the European population (LD length = 4 versus 41 kb for variants with  $r^2 > 0.6$  with rs1888909 in the African and European populations, respectively;  $N_{\text{eff}} = 4,270$  for Africans in the GBMI asthma meta-analysis; Figure S17). This variant was also fine-mapped for eosinophil count in UKBB Europeans (PIP = 1.0;  $p = 1.3 \times 10^{-314}$ )<sup>16</sup> and was previously reported to regulate *IL33* gene expression in human airway epithelial cells via allele-specific transcription factor binding of OCT-1 (POU2F1).<sup>77</sup> Likewise, we observed a missense variant rs16903574 (p.Phe319Leu) in *OTULINL* showed PIP = 0.79 (PIP = 1.0 and 0.21 in GBMI and UKBB Europeans, respectively; Figure 7E) owing to improved association significance ( $p = 7.7 \times 10^{-15}$  and  $4.7 \times 10^{-12}$  in GBMI and UKBB Europeans, respectively).

However, we also observed very high PIP for variants that are not likely causal. For example, we observed that an intronic variant rs1295686 in *IL13* showed PIP = 0.56 (PIP = 0.56 and 0.0002 in GBMI and UKBB Europeans, respectively; Figure 7F), despite having strong LD with a nearby missense variant rs20541 (p.Gln144Arg;  $r^2 = 0.96$  with rs1295686), which only showed PIP = 0.13 (PIP = 0.13 and 0.0001 in GBMI and UKBB Europeans, respectively). The missense variant rs20541 showed PIP = 0.23 and 0.15 for a related allergic disease, atopic dermatitis, in BBJ and FinnGen, respectively,<sup>16</sup> and was previously shown to induce STAT6 phosphorylation and upregulate CD23 expression in monocytes, promoting IgE synthesis.<sup>78</sup> Although the GBMI meta-analysis contributed to prioritizing these two variants (sum of PIP = 0.69 versus 0.0003 in GBMI and UKBB Europeans, respectively), the observed PIP was higher for rs1295686 than for rs20541.

While increasing sample size in meta-analysis improves association significance, we also found negative PIP due to losing the ability to model multiple causal variants. A stop-gained variant rs61816761 (p.Arg501Ter) in *FLG* showed PIP = -1.0 (PIP =  $6.4 \times 10^{-5}$  and 1.0 in GBMI and UKBB Europeans, respectively; Figure 7G), which was primarily owing to a nearby lead variant rs12123821 (~17 kb downstream of *HRNR*;  $r^2 = 0.0$  with rs61816761). This lead variant rs12123821 showed greater significance than rs61816761 in GBMI ( $p = 9.3 \times 10^{-16}$  and  $2.0 \times 10^{-11}$  for rs12123821 and rs61816761, respectively) as well as in UKBB Europeans ( $p = 7.1 \times 10^{-26}$  and  $1.5 \times 10^{-18}$ ). While our biobank fine-mapping<sup>16</sup> assigned PIP = for both variants based on multiple-causal-variant fine-mapping (i.e., FINEMAP<sup>20</sup> and SuSiE<sup>22</sup>), our ABF fine-mapping in the GBMI meta-analysis was only able to assign PIP = 0.74 for the lead variant rs12123821 due to a single causal variant assumption. This recapitulates the importance of multiple-causal-variant fine-mapping in complex trait fine-mapping<sup>16</sup>; however, we note that multiple-causal-variant fine-mapping with an external LD reference is extremely error prone as previously reported.<sup>14,15</sup>

## DISCUSSION

In this study, we first demonstrated in simulations that meta-analysis fine-mapping is substantially miscalibrated when constituent cohorts are heterogeneous in phenotyping, genotyping, and imputation. To mitigate this issue, we developed SLALOM, a summary statistics-based QC method for identifying suspicious loci in meta-analysis fine-mapping. Applying SLALOM to 14 disease endpoints from the GBMI meta-analyses<sup>10</sup> as well as

467 summary statistics from the GWAS Catalog,<sup>47</sup> we observed widespread suspicious loci in meta-analysis summary statistics, suggesting that meta-analysis fine-mapping is often miscalibrated in real data too. Indeed, we demonstrated that the predicted suspicious loci were significantly depleted for having likely causal variants as a lead PIP variant, such as nonsynonymous variants, high-PIP (>0.9) GWAS, and *cis*-eQTL fine-mapped variants from our previous fine-mapping studies.<sup>16</sup> Our method provides better calibration in non-suspicious loci for meta-analysis fine-mapping, generating a more reliable set of variants for further functional characterization.

We have found limited evidence of improved fine-mapping in the GBMI meta-analyses over individual biobanks. A few empirical examples in this study as well as other previous studies<sup>7,9,25,26,29</sup> suggested that multi-ancestry, large-scale meta-analysis could have potential to improve fine-mapping confidence and resolution owing to increased statistical power in associations and differential LD pattern across ancestries. However, we have highlighted that the observed improvement in PIP could be due to sample size imbalance in a locus, miscalibration, and technical confoundings too, which further emphasizes the importance of careful investigation of fine-mapped variants identified through meta-analysis fine-mapping. Given practical challenges in data harmonization across different cohorts, a large-scale biobank with multiple ancestries (e.g., UK Biobank<sup>18</sup> and All of Us<sup>79</sup>) would likely benefit the most from meta-analysis fine-mapping across ancestries.

As high-confidence fine-mapping results in large-scale biobanks and molecular quantitative trait loci (QTLs) continue to become available,<sup>15,16,60</sup> we propose alternative approaches for prioritizing candidate causal variants in a meta-analysis. First, these high-confidence fine-mapped variants have been a valuable resource to conduct a phenome-wide association study (PheWAS) to match with associated variants in a meta-analysis, which provides a narrower list of candidate variants assuming they would equally be functional and causal in related complex traits or tissues/cell types. Second, a traditional approach based on tagging variants (e.g.,  $r^2 > 0.6$  with lead variants, or PICS<sup>57</sup> fine-mapping approach that only relies on a lead variant and LD) can still be highly effective, especially for known functional variants such as nonsynonymous coding variants. As we highlighted in this and previous<sup>38</sup> studies, potentially causal variants in strong LD with lead variants might not achieve genome-wide significance because of missingness and heterogeneity.

While using an external LD reference for fine-mapping has been shown to be extremely error prone,<sup>14,15</sup> we find here that it can be useful for flagging suspicious loci, even when it does not perfectly represent the in-sample LD structure of the meta-analyzed individuals. However, our use of external LD reference comes with several limitations. For example, due to the finite sample size of external LD reference, rare or low-frequency variants have larger uncertainties around  $r^2$  than common variants. Moreover, our  $r^2$  values in a multi-ancestry meta-analysis are currently approximated based on a sample-size-weighted average of  $r^2$  across ancestries as previously suggested,<sup>80</sup> but this can be different from actual  $r^2$ . These uncertainties around  $r^2$  affect SLALOM prediction performance and should be modeled appropriately for further method development. On the other hand, we find it challenging to use an LD reference when true causal variants are located within a complex region (e.g.,

MHC), or are entirely missing from standard LD or imputation reference panels, especially for structural variants. These limitations are not specific to meta-analysis fine-mapping, and separate fine-mapping methods based on bespoke imputation references have been developed (e.g., human leukocyte antigen [HLA],<sup>81</sup> killer cell immunoglobulin-like receptor [KIR],<sup>82</sup> and variable numbers of tandem repeats<sup>83</sup>).

We have found evidence in our simulations and real data of severe miscalibration of fine-mapping results from GWAS meta-analysis; for example, we estimate that the difference between true and reported proportion of causal variants is 20% and 45% for high-PIP (>0.9) variants in suspicious loci from the simulations and the GWAS Catalog, respectively. Our SLALOM method helps to exclude spurious results from meta-analysis fine-mapping; however, even fine-mapping results in SLALOM-predicted non-suspicious loci remain somewhat miscalibrated, showing estimated differences between true and reported proportion of causal variants of 4% and 15% for high-PIP variants in the simulations and the GWAS Catalog, respectively. We thus urge extreme caution when interpreting PIPs computed from meta-analyses until improved methods are available. We recommend that researchers looking to identify likely causal variants employ complete synchronization of study design, case/control ascertainment, genomic profiling, and analytical pipeline, or rely more heavily on functional annotations, biobank fine-mapping, or molecular QTLs.

### Limitations of the study

There are several methodological limitations of SLALOM. First, our simulations only include one causal variant per locus. Although additional independent causal variants would not affect SLALOM precision (but decrease recall), multiple *correlated* causal variants in a locus would violate SLALOM assumptions and could lead to some DENTIST-S outliers that are not due to heterogeneity or missingness but rather simply a product of tagging multiple causal variants in LD. In fact, our previous studies have illustrated infrequent but non-zero presence of such correlated causal variants in complex traits.<sup>16</sup> Second, SLALOM prediction is not perfect. Although fine-mapping calibration is certainly better in non-suspicious versus suspicious loci, SLALOM has low precision, and we still observe some miscalibration in non-suspicious loci. Optimal thresholds for SLALOM prediction might be different for other datasets. Third, SLALOM does not model effect size heterogeneity. Although SLALOM is able to detect suspicious loci due to effect size heterogeneity as the method is agnostic to the source of heterogeneity, methods that model effect size heterogeneity, such as MR-MEGA,<sup>84</sup> could improve SLALOM performance. Finally, SLALOM is a per-locus QC method and does not calibrate per-variant PIPs. Further methodological development that properly models heterogeneity, missingness, sample size imbalance, multiple causal variants, and LD uncertainty across multiple cohorts and ancestries is needed to refine per-variant calibration and recall in meta-analysis fine-mapping.

## STAR★METHODS

### RESOURCE AVAILABILITY

**Lead contact**—Further information and requests for resources and data should be directed to and will be fulfilled by the lead contact, Masahiro Kanai (mkanai@broadinstitute.org).

**Materials availability**—This study did not generate new unique reagents.

**Data and code availability**—The GBMI summary statistics for the 14 endpoints are publicly available and are browserble at the GBMI PheWeb website (<http://results.globalbiobankmeta.org/>). Example outputs from the meta-analysis fine-mapping simulation pipeline have been deposited at Harvard Dataverse. All original code has been deposited at Zenodo and is publicly available as of the date of publication. DOIs and links are listed in the key resources table. Any additional information required to reanalyze the data reported in this paper is available from the lead contact upon request.

### METHOD DETAILS

**Meta-analysis fine-mapping simulation**—To benchmark fine-mapping performance in meta-analysis, we simulated a large-scale, realistic GWAS meta-analysis and performed fine-mapping under different scenarios. An overview of our simulation pipeline is summarized in Figure S2.

**Simulated true genotype**—Using HAPGEN2<sup>85</sup> with the 1000 Genomes Project Phase 3 (ref. 48), we simulated “true” genotypes of chromosome 3 for multiple independent cohorts from African, East Asian, and European ancestries. For each independent cohort from a given ancestry, we simulated 10,000 individuals each using the default parameters, with an ancestry-specific effective population size set to 17,469, 14,269, and 11,418 for Africans, East Asians, and Europeans, respectively, as recommended.<sup>85</sup> To mimic sample size imbalance of different ancestries in the current meta-analyses, we simulated 10 independent European cohorts, 1 African cohort, and 1 East Asian cohort.

To restrict our analysis to unrelated samples, we computed sample relatedness based on KING kinship coefficients<sup>88</sup> using PLINK (ref. 86) and removed monozygotic twins, duplicated individuals, or first-degree relatives with the coefficient threshold of 0.177. The detailed sample sizes of unrelated individuals for each cohort is summarized in Table S1.

**Genotyping and imputation**—To simulate realistic genotyping and imputation procedures, we first virtually genotyped each cohort by restricting variants to those that are available on different genotyping arrays. We selected three major genotyping arrays from Illumina, Inc. (Omni2.5, Multi-Ethnic Global Array [MEGA], and Global Screening Array [GSA]) that have different densities of genotyping probes (Table S2). For each cohort, we created three virtually genotyped datasets by retaining variants that are genotyped on each array. For the sake of simplicity, we assumed no genotyping errors occurred between true genotypes and virtually genotyped data—however, in practice, genotyping error is one of the major sources of unexpected confounding (*e.g.*, see recent discussions here<sup>89,90</sup>) and should be treated carefully.

For each pair of cohort and genotyping array, we then imputed missing variants using different imputation reference panels. We used the Michigan Imputation Server (<https://imputationserver.sph.umich.edu/>)<sup>87</sup> and the TOPMed Imputation Server (<https://imputation.biodatacatalyst.nhlbi.nih.gov/>)<sup>50</sup> with the default parameters, using three publicly available reference panels: the 1000 Genomes Project Phase 3 (version 5;  $n = 2,504$ ; 1000GP3),<sup>48</sup> the Haplo-type Reference Consortium (version r1.1;  $n = 32,470$ ; HRC),<sup>49</sup> and the TOPMed (version R2;  $n = 97,256$ ).<sup>50</sup> Briefly, for each input, the imputation server created chunks of 20 Mb, applied the standard QC, pre-phased each chunk with Eagle2 (ref. 91), and imputed non-genotyped variants using a specified reference panel with Minimac4 (<https://genome.sph.umich.edu/wiki/Minimac4>). The detailed documentation of the imputation pipeline is available on the Michigan and TOPMed websites and has been described elsewhere.<sup>87</sup>

We applied post-imputation QC by only keeping variants with  $MAF > 0.001$  and imputation  $Rsq > 0.6$ . Because the TOPMed panel is based on GRCh38 while the 1000GP3 and the HRC panels are on GRCh37, we lifted over TOPMed variants from GRCh38 to GRCh37 to meta-analyze with other cohorts. We excluded any variants which were lifted over to different chromosomes or for which the conversion failed. The number of virtually genotyped and imputed variants for each combination of cohort, genotyping array, and imputation panel is summarized in Table S3.

**True phenotype**—We simulated 300 true phenotypes that resemble observed complex trait genetic architecture and phenotypic heterogeneity across cohorts. Based on previous literature, we set parameters as follows: 1) 50% of 1 Mb loci contain a true causal variant<sup>92</sup>; 2) probability of being causal is proportional to functional enrichments of variant consequences (pLoF, missense, synonymous, 5'/3' UTR, promoter, *cis*-regulatory region, and non-genic) for fine-mapped variants as estimated in a previous complex trait fine-mapping study<sup>16</sup>; 3) per-allele causal effect sizes have a variance proportional to where represents a maximum MAF across the three ancestries (AFR, EAS, and EUR) and is set to be  $-0.38$  (ref. 51); and 4) total SNP-heritability for chromosome 3 equals 0.03 (ref. 52). For the sake of simplicity, we randomly draw a single true causal variant per locus because ABF assumes a single causal variant.<sup>30,31</sup> We draw true causal variants from 1,150,893 non-ambiguous single-nucleotide variants in 1000GP3 that showed  $MAF > 0.01$  in at least one of the three ancestries (AFR, EAS, or EUR) and were not located within conversion-unstable positions (CUP)<sup>53</sup> between the human genome builds GRCh37 and GRCh38. To mimic phenotypic heterogeneity across cohorts in real-world meta-analysis (due to e.g., different ascertainment, measurement error, or true effect size heterogeneity), we introduced cross-cohort genetic correlation of true effect sizes  $r_g$  which is set to be one of 1, 0.9, or 0.5. For a true causal variant  $j$ , true causal effect sizes  $\beta_j$  across cohorts were randomly drawn from  $\beta \sim \text{MVN}(0, \Sigma)$  where diagonal elements of  $\Sigma$  were set to be  $\sigma_g^2 \bullet [2p(1-p)]^\alpha$  and off-diagonal elements of  $\Sigma$  were set to be  $r_g \bullet \sigma_g^2 [2p(1-p)]^\alpha$ .  $\sigma_g^2$  was determined by  $\sigma_g^2 = h_g^2 / \Sigma_j [2p(1-p)]^{1+\alpha}$ . For each cohort, true phenotype  $y$  was computed via  $y = X\beta + \epsilon$  where  $X$  is the above true genotype matrix from HAPGEN2 and  $\epsilon_i \sim N(0, 1 - \sigma_g^2)$  i.i.d. We simulated 100 true phenotypes for each of  $r_g = 1, 0.9, \text{ and } 0.5$ , respectively.

**GWAS**—For each combination of phenotype, cohort, genotyping chip, and imputation panel, we conducted GWAS via a standard linear regression as implemented in PLINK 2.0 using imputed dosages. For covariates, we included top 10 principal components that were calculated based on true genotypes after restricting to unrelated samples. We only used LD-pruned variants with MAF >0.01 for PCA.

**Meta-analysis**—To simulate meta-analyses that resemble real-world settings, we generated multiple *configurations* of the above GWAS results to meta-analyze across 10 independent cohorts. Briefly, we chose configurations based on the following settings: 1) 10 EUR cohorts are genotyped and imputed using the same genotyping array (one of GSA, MEGA, or Omni2.5) and the same imputation panel (one of 1000GP3, HRC, TOPMed, or TOPMed-liftover); 2) 10 cohorts consisting of multiple ancestries (9 EUR +1 AFR/EAS cohorts or 8 EUR +1 AFR +1 EAS cohorts), with all cohorts genotyped and imputed using the same array (Omni2.5) and the same panel (1000GP3); 3) 10 EUR or multi-ancestry cohorts are genotyped using the same array (Omni2.5) but imputed using different panels across cohorts; 4) 10 EUR or multi-ancestry cohorts are imputed using the same panel (1000GP3) but genotyped using different arrays across cohorts; 5) 10 EUR or multi-ancestry cohorts are genotyped and imputed using different arrays and panels across cohorts. For settings 3–5, we randomly draw a combination of a genotyping array and an imputation panel for each cohort five times each for 10 EUR and multi-ancestry cohorts. In total, we generated 45 configurations as summarized in Table S4.

For each configuration, we conducted a fixed-effect meta-analysis based on inverse-variance weighted betas and standard errors using a modified version of PLINK 1.9 ([https://github.com/mkanai/plink-ng/tree/add\\_se\\_meta](https://github.com/mkanai/plink-ng/tree/add_se_meta)).

**Fine-mapping**—For each meta-analysis, we defined fine-mapping regions based on a 1 Mb window around each genome-wide significant lead variant and applied ABF<sup>30,31</sup> using prior effect size variance of = 0.04. We set a prior variance of effect size to be 0.04 which was taken from Wakefield et al.<sup>30</sup> and is commonly used in meta-analysis fine-mapping studies.<sup>2,7</sup> We computed posterior inclusion probability (PIP) and 95% credible set (CS) for each locus and evaluated whether true causal variants were correctly fine-mapped.

**The SLALOM method**—SLALOM takes GWAS summary statistics and external LD reference as input and predicts whether a locus is suspicious for fine-mapping. SLALOM consists of the following three steps:

**Locus definition**—Consistent with common fine-mapping region definition, we defined loci based on a 1 Mb window around each genome-wide significant lead variant and merged them if they overlapped. We excluded the major histocompatibility complex (MHC) region (chr 6:25–36 Mb) from analysis due to extensive LD structure in the region.

**DENTIST-S outlier detection**—For each variant in a locus, we computed DENTIST-S statistics using Equation 1 based on the assumption of a single causal variant. DENTIST-S P-values ( $P_{\text{DENTIST-S}}$ ) were computed using the distribution with 1 degree of freedom. We applied ABF<sup>30,31</sup> using prior effect size variance of = 0.04 and used the lead PIP

variant (the variant with the highest PIP) as an approximation of the causal variant in the locus. To retrieve correlation  $r$  among the variants, we used publicly available LD matrices from gnomAD<sup>56</sup> v2 as external LD reference for African, Admixed American, East Asian, Finnish, and non-Finnish European populations. When multiple populations exist, we computed a sample-size-weighted average of  $r^2$  using per-variant sample sizes for each population as previously suggested.<sup>80</sup> We excluded variants without  $r^2$  available in gnomAD from the analysis. Since gnomAD v2 LD matrices are based on the human genome assembly GRCh37, variants were lifted over to GRCh38 if the input summary statistics were based on GRCh38.

We determined DENTIST-S outlier variants using two thresholds: 1)  $r^2 > \rho$  to the lead and 2)  $P_{\text{DENTIST-S}} < \tau$ . The thresholds  $\rho$  and  $\tau$  were set to  $\rho = 0.6$  and  $\tau = 1.0 \times 10^{-4}$  based on the training in simulations as described below.

**Suspicious loci prediction**—We predicted whether a locus is suspicious or non-suspicious for fine-mapping based on the number of DENTIST-S outlier variants in the locus  $> \kappa$ . To determine the best-performing thresholds ( $\rho$ ,  $\tau$ , and  $\kappa$ ), we used loci with maximum PIP  $> 0.9$  in the simulations for training. Positive conditions were defined as whether a true causal variant in a locus is 1) a lead PIP variant, 2) in 95% CS, and 3) in 99% CS. We computed AUROC across different thresholds ( $\rho = 0, 0.1, 0.2, \dots, 0.9$ ;  $-\log_{10} \tau = 0, 0.5, 1, \dots, 10$ ; and  $\kappa = 0, 1, 2, \dots$ ) and chose  $\rho = 0.6$ ,  $\tau = 1.0 \times 10^{-4}$ , and  $\kappa = 0$  that showed the highest AUROC for all the aforementioned positive conditions. Using all the loci in the simulations, we then evaluated fine-mapping miscalibration (defined as mean PIP – fraction of true causal variants) at different PIP thresholds in suspicious and non-suspicious loci and decided to only apply SLALOM to loci with maximum PIP  $> 0.1$  owing to relatively lower miscalibration and specificity of SLALOM at lower PIP thresholds.

**GWAS catalog analysis**—We retrieved full GWAS summary statistics publicly available on the GWAS Catalog.<sup>47</sup> Out of 33,052 studies from 5,553 publications registered at the GWAS Catalog (as of January 12, 2022), we selected 467 studies from 96 publications that have 1) full harmonized summary statistics preprocessed by the GWAS Catalog with non-missing variant ID, marginal beta, and SE columns, 2) a discovery sample size of more than 10,000 individuals, 3) African (including African American, Afro-Caribbean, and Sub-Saharan African), admixed American (Hispanic and Latin American), East Asian, or European samples based on their broad ancestral category metadata, 4) at least one genome-wide significant association ( $p < 5.0 \times 10^{-8}$ ), and 5) our manual annotation as a meta-analysis rather than a single-cohort study (Table S6). We applied SLALOM to the 467 summary statistics and identified 35,864 genome-wide significant loci (based on 1 Mb window around lead variants), of which 28,925 loci with maximum PIP  $> 0.1$  were further classified into suspicious and non-suspicious loci. Since per-variant sample sizes were not available, we used overall sample sizes of each ancestry (African, Admixed American, East Asian, and European) to calculate the weighted-average of  $r^2$ . All the variants were harmonized into the human genome assembly GRCh38 by the GWAS Catalog.

**GBMI analysis**—We used meta-analysis summary statistics of 14 disease endpoints from the GBMI (Table S8). These meta-analyses were conducted using up to 1.8 million

individuals across 18 biobanks for discovery, representing six different genetic ancestry groups (approximately African, 18,000 Admixed American, 31,000 Central and South Asian, 341,000 East Asian, 1.4 million European, and 1,600 Middle Eastern individuals). Detailed procedures of the GBMI meta-analyses were described in the GBMI flagship publication.<sup>10</sup>

Across the 14 summary statistics, we used 489 out of 500 genome-wide significant loci ( $p < 5.0 \times 10^{-8}$ ; 1 Mb window around each lead variant, as defined in the GBMI flagship publication<sup>10</sup>), excluding 11 loci that overlap with the MHC region. We applied SLALOM to 422 loci with maximum PIP  $>0.1$  based on the ABF fine-mapping and predicted whether they were suspicious or non-suspicious for fine-mapping. We used per-variant sample sizes of each ancestry (African, Admixed American, East Asian, Finnish, and non-Finnish European) to calculate the weighted-average of  $r^2$ . Since gnomAD LD matrices were not available for Central and South Asian and Middle Eastern, we did not use their sample sizes for the calculation. All the variants were processed on the human genome assembly GRCh38.

**Fine-mapping results of complex traits and *cis*-eQTL**—We retrieved our previous fine-mapping results for 1) complex traits in large-scale biobanks (BBJ,<sup>58</sup> FinnGen,<sup>19</sup> and UKBB<sup>18</sup> Europeans)<sup>16</sup> and 2) *cis*-eQTLs in GTEx<sup>59</sup> v8 and eQTL Catalogue<sup>60</sup>. Briefly, we conducted multiple-causal-variant fine-mapping (FINEMAP<sup>20,21</sup> and SuSiE<sup>22</sup>) of complex trait GWAS (# unique traits = 148) and *cis*-eQTL gene expression (# unique tissues/cell-types = 69) using summary statistics and in-sample LD. Detailed fine-mapping methods are described elsewhere.<sup>16</sup>

In this study, we collected 1) high-PIP GWAS variants that achieved PIP  $>0.9$  for any traits in any biobank and 2) high-PIP *cis*-eQTL variants that achieved PIP  $>0.9$  for any gene expression in any tissues/cell-types. All the variants were originally processed on the human genome assembly GRCh37 and lifted over to the GRCh38 for comparison.

**Additional fine-mapping results**—To compare with the GBMI meta-analyses, we additionally conducted multi-causal-variant fine-mapping of four additional endpoints (gout, heart failure, thyroid cancer, and venous thromboembolism) that were not fine-mapped in our previous study.<sup>16</sup> We used exactly the same fine-mapping pipeline (FINEMAP<sup>20,21</sup> and SuSiE<sup>22</sup>) as described previously.<sup>16</sup> For UKBB Europeans, to use the exact same samples that contributed to the GBMI, we used individuals of European ancestry ( $n = 420,531$ ) as defined in the Pan-UKBB project (<https://pan.ukbb.broadinstitute.org>), instead of those of “white British ancestry” ( $n = 361,194$ ) used in our previous study.<sup>16</sup>

**Enrichment analysis of likely causal variants**—To validate SLALOM performance, we asked whether suspicious and non-suspicious loci were enriched for having likely causal variants as a lead PIP variant, and for containing them in the 95 and 99% CS. We defined likely causal variants using 1) nonsynonymous coding variants, *i.e.*, pLoF and missense variants annotated<sup>93</sup> by the Ensembl Variant Effect Predictor (VEP) v101 (using GRCh38 and GENCODE v35), 2) the high-PIP ( $>0.9$ ) complex trait fine-mapped variants, and 3)

the high-PIP (>0.9) *cis*-eQTL fine-mapped variants from our previous studies as described above.

We estimated enrichment for suspicious and non-suspicious loci as a relative risk (*i.e.*, a ratio of proportion of variants) between being in suspicious/non-suspicious loci and having the annotated likely causal variants as a lead PIP variant (or containing them in the 95% or 99% CS). That is, a relative risk = (proportion of non-suspicious loci having the annotated variants as a lead PIP variant)/(proportion of suspicious loci having the annotated variants as a lead PIP variant). We computed 95% confidence intervals using bootstrapping.

**Comparison of fine-mapping results between the GBMI and individual biobanks**—To directly compare with fine-mapping results from the GBMI meta-analyses, we used our fine-mapping results of nine disease endpoints (asthma,<sup>64</sup> COPD,<sup>64</sup> gout, heart failure,<sup>73</sup> IPF,<sup>62</sup> primary open-angle glaucoma,<sup>74</sup> thyroid cancer, stroke,<sup>75</sup> and venous thromboembolism<sup>76</sup>) in BBJ,<sup>58</sup> FinnGen,<sup>19</sup> and UKBB<sup>18</sup> Europeans that were also part of the GBMI meta-analyses for the same traits. For comparison, we computed the maximum PIP for each variant and the minimum size of 95% CS across BBJ, FinnGen, and UKBB. We restricted the 95% CS in biobanks to those that contain the lead variants from the GBMI. We defined the PIP difference between the GBMI and individual biobanks as  $\text{PIP} = \text{PIP}(\text{GBMI}) - \text{the maximum PIP across the biobanks}$ .

We conducted functional enrichment analysis to compare between the GBMI meta-analysis and individual biobanks because unbiased comparison of PIP requires conditioning on likely causal variants independent of the fine-mapping results, and functional annotations have been shown to be enriched for causal variants. Using functional categories (coding [pLoF, missense, and synonymous], 5'/3' UTR, promoter, and CRE) from our previous study,<sup>16</sup> we estimated functional enrichments of variants in each functional category based on 1) top PIP rankings and 2) PIP bins. Since fine-mapping PIP in the GBMI meta-analysis can be miscalibrated, we performed a comparison based on top PIP rankings to assess whether the ordering given by GBMI PIPs is more informative than the ordering given by the biobanks. For the top PIP rankings, we took the top 0.5%, 0.1%, and 0.05% variants based on the PIP rankings in the GBMI and individual biobanks. We computed enrichment as a relative risk = (proportion of top X% PIP variants in the GBMI that are in the annotation)/(proportion of top X% PIP variants in the individual biobanks that are in the annotation). For PIP bins, we defined three bins using different thresholds ( $\theta = 0.01, 0.05, \text{ and } 0.1$ ): 1) decreased PIP bin,  $\text{PIP} < -\theta$ , 2) null bin,  $-\theta \leq \text{PIP} \leq \theta$ , and 3) increased PIP bin,  $\text{PIP} > \theta$ . We computed enrichment as a relative risk = (proportion of variants in the decreased/increased PIP bin that are in the annotation)/(proportion of variants in the null PIP bin). We combined coding, UTR, and promoter categories for this analysis due to the limited number of variants for each bin.

## QUANTIFICATION AND STATISTICAL ANALYSIS

All statistical analysis was performed using R 4.0.3, Hail 0.2, PLINK 1.9 and 2.0. All methodological details can be found in the method details, and all statistical tests are named as they are used.

## Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

## ACKNOWLEDGMENTS

We acknowledge all the participants and researchers of the 23 biobanks that have contributed to the GBMI. Biobank-specific acknowledgments are included in the Data S3. We thank H. Huang, A.R. Martin, B.M. Neale, Y. Okada, K. Tsuo, J.C. Ulirsch, Y. Wang, and all the members of Finucane and Daly labs for their helpful feedback. M.K. was supported by a Nakajima Foundation Fellowship and the Masason Foundation. H.K.F. was funded by NIH grant DP5 OD024582.

## CONSORTIA

GBMI: 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 Community Health Initiative, Uganda Genome Resource, UK Biobank, Alicia R. Martin, Cristen J. Willer, Mark J. Daly, Benjamin M. Neale. See the Supplemental PDF for consortium member affiliations.

## REFERENCES

1. Evangelou E, and Ioannidis J.P.a. (2013). Meta-analysis methods for genome-wide association studies and beyond. *Nat. Rev. Genet* 14, 379–389. [PubMed: 23657481]
2. Mahajan A, Taliun D, Thurner M, Robertson NR, Torres JM, Rayner NW, Payne AJ, Steinthorsdottir V, Scott RA, Grarup N, et al. (2018). Fine-mapping type 2 diabetes loci to single-variant resolution using high-density imputation and islet-specific epigenome maps. *Nat. Genet* 50, 1505–1513. [PubMed: 30297969]
3. Spracklen CN, Horikoshi M, Kim YJ, Lin K, Bragg F, Moon S, Suzuki K, Tam CHT, Tabara Y, Kwak S-H, et al. (2020). Identification of type 2 diabetes loci in 433, 540 East Asian individuals. *Nature* 582, 240–245. [PubMed: 32499647]
4. Schizophrenia Working Group of the Psychiatric Genomics Consortium; Neale BM, Corvin A, Walters JTR, Farh K-H, Holmans P.a., Lee P, Bulik-Sullivan B, Collier D.a., Huang H, et al. (2014). Biological insights from 108 schizophrenia-associated genetic loci. *Nature* 511, 421–427. [PubMed: 25056061]
5. Trubetskoy V, Pardiñas AF, Qi T, Panagiotaropoulou G, Awasthi S, Bigdeli TB, et al. (2022). Mapping genomic loci implicates genes and synaptic biology in schizophrenia. *Nature* 604, 502–508. [PubMed: 35396580]
6. Okada Y, Wu D, Trynka G, Raj T, Terao C, Ikari K, Kochi Y, Ohmura K, Suzuki A, Yoshida S, et al. (2014). Genetics of rheumatoid arthritis contributes to biology and drug discovery. *Nature* 506, 376–381. [PubMed: 24390342]
7. Ishigaki K, Sakaue S, Terao C, Luo Y, Sonehara K, Yamaguchi K, et al. (2021). Trans-ancestry genome-wide association study identifies novel genetic mechanisms in rheumatoid arthritis. Preprint at medRxiv. 10.1101/2021.12.01.21267132.
8. Locke AE, Kahali B, Berndt SI, Justice AE, Pers TH, Day FR, Powell C, Vedantam S, Buchkovich ML, Yang J, et al. (2015). Genetic studies of body mass index yield new insights for obesity biology. *Nature* 518, 197–206. [PubMed: 25673413]
9. Graham SE, Clarke SL, Wu K-HH, Kanoni S, Zajac GJM, Ramdas S, Surakka I, Ntalla I, Vedantam S, Winkler TW, et al. (2021). The power of genetic diversity in genome-wide association studies of lipids. *Nature* 600, 675–679. [PubMed: 34887591]
10. Zhou W, Kanai M, Wu K-HH, Rasheed H, Tsuo K, Hirbo JB, et al. (2022). Global Biobank Meta-analysis Initiative: Powering genetic discovery across human disease. *Cell Genomics* 2, 100192.
11. Visscher PM, Wray NR, Zhang Q, Sklar P, McCarthy MI, Brown MA, and Yang J (2017). 10 Years of GWAS discovery: biology, function, and translation. *Am. J. Hum. Genet* 101, 5–22. [PubMed: 28686856]
12. Shendure J, Findlay GM, and Snyder MW (2019). Genomic medicine-progress, pitfalls, and promise. *Cell* 177, 45–57. [PubMed: 30901547]
13. Schaid DJ, Chen W, and Larson NB (2018). From genome-wide associations to candidate causal variants by statistical fine-mapping. *Nat. Rev. Genet* 19, 491–504. [PubMed: 29844615]
14. Ulirsch JC, Lareau CA, Bao EL, Ludwig LS, Guo MH, Benner C, Satpathy AT, Kartha VK, Salem RM, Hirschhorn JN, et al. (2019). Interrogation of human hematopoiesis at single-cell and single-variant resolution. *Nat. Genet* 51, 683–693. [PubMed: 30858613]
15. Weissbrod O, Hormozdiari F, Benner C, Cui R, Ulirsch J, Gazal S, Schoech AP, van de Geijn B, Reshef Y, Márquez-Luna C, et al. (2020). Functionally informed fine-mapping and polygenic localization of complex trait heritability. *Nat. Genet* 52, 1355–1363. [PubMed: 33199916]

16. Kanai M, Ulirsch JC, Karjalainen J, Kurki M, Karczewski KJ, Fauman E, Wang QS, Jacobs H, Aguet F, Ardlie KG, et al. (2021). Insights from complex trait fine-mapping across diverse populations. Preprint at medRxiv. 10.1101/2021.09.03.21262975.
17. Nagai A, Hirata M, Kamatani Y, Muto K, Matsuda K, Kiyohara Y, Ninomiya T, Tamakoshi A, Yamagata Z, Mushirola T, et al. (2017). Overview of the BioBank Japan project: study design and profile. *J. Epidemiol* 27, S2–S8. [PubMed: 28189464]
18. Bycroft C, Freeman C, Petkova D, Band G, Elliott LT, Sharp K, Motyer A, Vukcevic D, Delaneau O, O'Connell J, et al. (2018). The UK Biobank resource with deep phenotyping and genomic data. *Nature* 562, 203–209. [PubMed: 30305743]
19. Kurki MI, Karjalainen J, Palta P, Sipilä TP, Kristiansson K, Donner K, et al. (2022). FinnGen: unique genetic insights from combining isolated population and national health register data. Preprint at medRxiv. 10.1101/2022.03.03.22271360.
20. Benner C, Spencer CCA, Havulinna AS, Salomaa V, Ripatti S, and Pirinen M (2016). FINEMAP: efficient variable selection using summary data from genome-wide association studies. *Bioinformatics* 32, 1493–1501. [PubMed: 26773131]
21. Benner C, Havulinna AS, Salomaa V, Ripatti S, and Pirinen M (2018). Refining fine-mapping: effect sizes and regional heritability. Preprint at bioRxiv. 10.1101/318618.
22. Wang G, Sarkar A, Carbonetto P, and Stephens M (2020). A simple new approach to variable selection in regression, with application to genetic fine mapping. *J. R. Stat. Soc. B* 82, 1273–1300.
23. Onengut-Gumuscu S, Chen W-M, Burren O, Cooper NJ, Quinlan AR, Mychaleckyj JC, Farber E, Bonnie JK, Szpak M, Schofield E, et al. (2015). Fine mapping of type 1 diabetes susceptibility loci and evidence for colocalization of causal variants with lymphoid gene enhancers. *Nat. Genet* 47, 381–386. [PubMed: 25751624]
24. Levey DF, Stein MB, Wendt FR, Pathak GA, Zhou H, Aslan M, Quaden R, Harrington KM, Nuñez YZ, Overstreet C, et al. (2021). Bi-ancestral depression GWAS in the Million Veteran Program and meta-analysis in >1.2 million individuals highlight new therapeutic directions. *Nat. Neurosci* 24, 954–963. [PubMed: 34045744]
25. Gharahkhani P, Jorgenson E, Hysi P, Khawaja AP, Pendergrass S, Han X, Ong JS, Hewitt AW, Segrè AV, Rouhana JM, et al. (2021). Genome-wide meta-analysis identifies 127 open-angle glaucoma loci with consistent effect across ancestries. *Nat. Commun* 12, 1258. [PubMed: 33627673]
26. Chen J, Spracklen CN, Marenne G, Varshney A, Corbin LJ, Luan J, Willems SM, Wu Y, Zhang X, Horikoshi M, et al. (2021). The trans-ancestral genomic architecture of glycemic traits. *Nat. Genet* 53, 840–860. [PubMed: 34059833]
27. Zhou W, Brumpton B, Kabil O, Gudmundsson J, Thorleifsson G, Weinstock J, Zawistowski M, Nielsen JB, Chaker L, Medici M, et al. (2020). GWAS of thyroid stimulating hormone highlights pleiotropic effects and inverse association with thyroid cancer. *Nat. Commun* 11, 3981–4013. [PubMed: 32769997]
28. Wightman DP, Jansen IE, Savage JE, Shadrin AA, Bahrami S, Holland D, Rongve A, Børte S, Winsvold BS, Drange OK, et al. (2021). A genome-wide association study with 1,126,563 individuals identifies new risk loci for Alzheimer's disease. *Nat. Genet* 53, 1276–1282. [PubMed: 34493870]
29. Chen M-H, Raffield LM, Mousas A, Sakaue S, Huffman JE, Moscati A, Trivedi B, Jiang T, Akbari P, Vuckovic D, et al. (2020). Trans-ethnic and ancestry-specific blood-cell genetics in 746,667 individuals from 5 global populations. *Cell* 182, 1198–1213.e14. [PubMed: 32888493]
30. Wakefield J (2007). A Bayesian measure of the probability of false discovery in genetic epidemiology studies. *Am. J. Hum. Genet* 81, 208–227. [PubMed: 17668372]
31. Wakefield J (2009). Bayes factors for genome-wide association studies: comparison with P-values. *Genet. Epidemiol* 33, 79–86. [PubMed: 18642345]
32. Hormozdiari F, Kostem E, Kang EY, Pasaniuc B, and Eskin E (2014). Identifying causal variants at loci with multiple signals of association. *Genetics* 198, 497–508. [PubMed: 25104515]
33. Kichaev G, Yang W-Y, Lindstrom S, Hormozdiari F, Eskin E, Price AL, Kraft P, and Pasaniuc B (2014). Integrating functional data to prioritize causal variants in statistical fine-mapping studies. *PLoS Genet.* 10, e1004722. [PubMed: 25357204]

34. Kichaev G, and Pasaniuc B (2015). Leveraging functional-annotation data in trans-ethnic fine-mapping studies. *Am. J. Hum. Genet* 97, 260–271. [PubMed: 26189819]
35. Li D, Zhao H, and Gelernter J (2012). Strong protective effect of the aldehyde dehydrogenase gene (ALDH2) 504Iys (\*2) allele against alcoholism and alcohol-induced medical diseases in Asians. *Hum. Genet* 131, 725–737. [PubMed: 22102315]
36. Brown BC, Asian Genetic Epidemiology Network Type 2 Diabetes Consortium; Ye CJ, Price AL, and Zaitlen N (2016). Transethnic genetic-correlation estimates from summary statistics. *Am. J. Hum. Genet* 99, 76–88. [PubMed: 27321947]
37. Shi H, Gazal S, Kanai M, Koch EM, Schoech AP, Siewert KM, et al. (2021). Population-specific causal disease effect sizes in functionally important regions impacted by selection. *Nat. Commun* 12, 1098. [PubMed: 33597505]
38. COVID-19 Host Genetics Initiative. (2021). Mapping the human genetic architecture of COVID-19. *Nature* 600, 472–477. [PubMed: 34237774]
39. Dendrou CA, Cortes A, Shipman L, Evans HG, Attfield KE, Jostins L, Barber T, Kaur G, Kuttikkatte SB, Leach OA, et al. (2016). Resolving *TYK2* locus genotype-to-phenotype differences in autoimmunity. *Sci. Transl. Med* 8, 363ra149.
40. Couturier N, Bucciarelli F, Nurtdinov RN, Debouverie M, Lebrun-Frenay C, Defer G, Moreau T, Confavreux C, Vukusic S, Cournu-Rebeix I, et al. (2011). Tyrosine kinase 2 variant influences T lymphocyte polarization and multiple sclerosis susceptibility. *Brain* 134, 693–703. [PubMed: 21354972]
41. Li Z, Gakovic M, Ragimbeau J, Eloranta M-L, Rönnblom L, Michel F, and Pellegrini S (2013). Two rare disease-associated Tyk2 variants are catalytically impaired but signaling competent. *J. Immunol* 190, 2335–2344. [PubMed: 23359498]
42. Lam M, Awasthi S, Watson HJ, Goldstein J, Panagiotaropoulou G, Trubetskoy V, Karlsson R, Frei O, Fan C-C, DeWitte W, et al. (2020). RICOPIIL: rapid imputation for COnsortias PipeLine. *Bioinformatics* 36, 930–933. [PubMed: 31393554]
43. Huang H, Fang M, Jostins L, Umi evi Mirkov M, Boucher G, Anderson CA, Andersen V, Cleyne I, Cortes A, Crins F, et al. (2017). Fine-mapping inflammatory bowel disease loci to single-variant resolution. *Nature* 547, 173–178. [PubMed: 28658209]
44. Winkler TW, Day FR, Croteau-Chonka DC, Wood AR, Locke AE, Mägi R, Ferreira T, Fall T, Graff M, Justice AE, et al. (2014). Quality control and conduct of genome-wide association meta-analyses. *Nat. Protoc* 9, 1192–1212. [PubMed: 24762786]
45. Chen W, Wu Y, Zheng Z, Qi T, Visscher PM, Zhu Z, and Yang J (2021). Improved analyses of GWAS summary statistics by reducing data heterogeneity and errors. *Nat. Commun* 12, 7117. [PubMed: 34880243]
46. Yang J, Ferreira T, Morris AP, Medland SE, Genetic Investigation of ANthropometric Traits GIANT Consortium; DIAbetes Genetics Replication And Meta-analysis DIAGRAM Consortium; Madden PAF, Heath AC, Martin NG, Montgomery GW, et al. (2012). Conditional and joint multiple-SNP analysis of GWAS summary statistics identifies additional variants influencing complex traits. *Nat. Genet* 44, 369–375. S1–S3. [PubMed: 22426310]
47. Buniello A, MacArthur JAL, Cerezo M, Harris LW, Hayhurst J, Malangone C, McMahon A, Morales J, Mountjoy E, Sollis E, et al. (2019). The NHGRI-EBI GWAS Catalog of published genome-wide association studies, targeted arrays and summary statistics 2019. *Nucleic Acids Res.* 47, D1005–D1012. [PubMed: 30445434]
48. 1000 Genomes Project Consortium; Auton A, Brooks LD, Durbin RM, Garrison EP, Kang HM, Korbel JO, Marchini JL, McCarthy S, McVean GA, and Abecasis GR (2015). A global reference for human genetic variation. *Nature* 526, 68–74. [PubMed: 26432245]
49. McCarthy S, Das S, Kretzschmar W, Delaneau O, Wood AR, Teumer A, Kang HM, Fuchsberger C, Danecek P, Sharp K, et al. (2016). A reference panel of 64, 976 haplotypes for genotype imputation. *Nat. Genet* 48, 1279–1283. [PubMed: 27548312]
50. Taliun D, Harris DN, Kessler MD, Carlson J, Szpiech ZA, Torres R, Taliun SAG, Corvelo A, Gogarten SM, Kang HM, et al. (2021). Sequencing of 53, 831 diverse genomes from the NHLBI TOPMed Program. *Nature* 590, 290–299. [PubMed: 33568819]

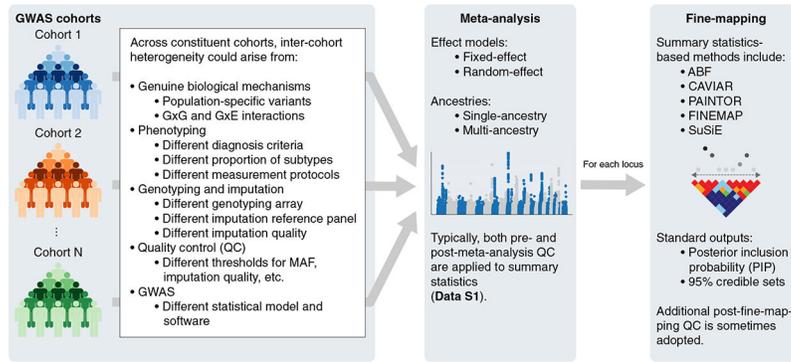
51. Schoech AP, Jordan DM, Loh P-R, Gazal S, O'Connor LJ, Balick J, Palamara PF, Finucane HK, Sunyaev SR, and Price AL (2019). Quantification of frequency-dependent genetic architectures in 25 UK Biobank traits reveals action of negative selection. *Nat. Commun* 10, 790. [PubMed: 30770844]
52. Yang J, Manolio T.a., Pasquale LR, Boerwinkle E, Caporaso N, Cunningham JM, de Andrade M, Feenstra B, Feingold E, Hayes MG, et al. (2011). Genome partitioning of genetic variation for complex traits using common SNPs. *Nat. Genet* 43, 519–525. [PubMed: 21552263]
53. Ormond C, Ryan NM, Corvin A, and Heron EA (2021). Converting single nucleotide variants between genome builds: from cautionary tale to solution. *Brief. Bioinform* 22, bbab069. 10.1093/bib/bbab069. [PubMed: 33822888]
54. Asimit JL, Hatzikotoulas K, McCarthy M, Morris AP, and Zeggini E (2016). Trans-ethnic study design approaches for fine-mapping. *Eur. J. Hum. Genet* 24, 1330–1336. [PubMed: 26839038]
55. Marchini J, and Howie B (2010). Genotype imputation for genome-wide association studies. *Nat. Rev. Genet* 11, 499–511. [PubMed: 20517342]
56. Karczewski KJ, Francioli LC, Tiao G, Cummings BB, Alföldi J, Wang Q, Collins RL, Laricchia KM, Ganna A, Birnbaum DP, et al. (2020). The mutational constraint spectrum quantified from variation in 141, 456 humans. *Nature* 581, 434–443. [PubMed: 32461654]
57. Farh KK-H, Marson A, Zhu J, Kleinewietfeld M, Housley WJ, Beik S, Shores N, Whitton H, Ryan RJH, Shishkin AA, et al. (2015). Genetic and epigenetic fine mapping of causal autoimmune disease variants. *Nature* 518, 337–343. [PubMed: 25363779]
58. Sakaue S, Kanai M, Tanigawa Y, Karjalainen J, Kurki M, Koshiba S, Narita A, Konuma T, Yamamoto K, Akiyama M, et al. (2021). A cross-population atlas of genetic associations for 220 human phenotypes. *Nat. Genet* 53, 1415–1424. [PubMed: 34594039]
59. GTEx Consortium (2020). The GTEx Consortium atlas of genetic regulatory effects across human tissues. *Science* 369, 1318–1330. [PubMed: 32913098]
60. Kerimov N, Hayhurst JD, Peikova K, Manning JR, Walter P, Kolberg L, Samovi a M, Sakthivel MP, Kuzmin I, Trevanion SJ, et al. (2021). A compendium of uniformly processed human gene expression and splicing quantitative trait loci. *Nat. Genet* 53, 1290–1299. [PubMed: 34493866]
61. Koskela JT, Häppölä P, Liu A, Partanen J, Genovese G, Artomov M, et al. (2021). Genetic variant in SPDL1 reveals novel mechanism linking pulmonary fibrosis risk and cancer protection. Preprint at medRxiv. 10.1101/2021.05.07.21255988.
62. Partanen JJ, Häppölä P, Zhou W, Lehisto AA, Ainola M, Sutinen D, et al. (2022). Leveraging global multi-ancestry meta-analysis in the study of idiopathic pulmonary fibrosis genetics. *Cell Genomics* 2, 100181.
63. Foreman MG, Wilson C, DeMeo DL, Hersh CP, Beaty TH, Cho MH, Ziniti J, Curran-Everett D, Criner G, Hokanson JE, et al. (2017). Alpha-1 Antitrypsin PiMZ genotype is associated with chronic obstructive pulmonary disease in two racial groups. *Ann. Am. Thorac. Soc* 14, 1280–1287. [PubMed: 28380308]
64. Tsuo K, Zhou W, Wang Y, Kanai M, Namba S, Gupta R, et al. (2021). Multi-ancestry meta-analysis of asthma identifies novel associations and highlights the value of increased power and diversity. Preprint at medRxiv. 10.1101/2021.11.30.21267108.
65. Benonisdottir S, Oddsson A, Helgason A, Kristjansson RP, Sveinbjornsson G, Oskarsdottir A, Thorleifsson G, Davidsson OB, Arnadottir GA, Sulem G, et al. (2016). Epigenetic and genetic components of height regulation. *Nat. Commun* 7, 13490. [PubMed: 27848971]
66. Marouli E, Graff M, Medina-Gomez C, Lo KS, Wood AR, Kjaer TR, Fine RS, Lu Y, Schurmann C, Highland HM, et al. (2017). Rare and low-frequency coding variants alter human adult height. *Nature* 542, 186–190. [PubMed: 28146470]
67. Langefeld CD, Ainsworth HC, Cunnigham Graham DS, Kelly JA, Comeau ME, Marion MC, Howard TD, Ramos PS, Croker JA, Morris DL, et al. (2017). Transancestral mapping and genetic load in systemic lupus erythematosus. *Nat. Commun* 8, 16021. [PubMed: 28714469]
68. Hargreaves CE, Rose-Zerilli MJ, Machado LR, Iriyama C, Hollox EJ, Cragg MS, and Strefford JC (2015). Fcγ receptors: genetic variation, function, and disease. *Immunol. Rev* 268, 6–24. [PubMed: 26497510]

69. Franke L, el Bannoudi H, Jansen DTSL, Kok K, Trynka G, Diogo D, Swertz M, Fransen K, Knevel R, Gutierrez-Achury J, et al. (2016). Association analysis of copy numbers of FC-gamma receptor genes for rheumatoid arthritis and other immune-mediated phenotypes. *Eur. J. Hum. Genet* 24, 263–270. [PubMed: 25966632]
70. UK10K Consortium; Walter K, Min JL, Huang J, Crooks L, Memari Y, McCarthy S, Perry JRB, Xu C, Futema M, et al. (2015). The UK10K project identifies rare variants in health and disease. *Nature* 526, 82–90. [PubMed: 26367797]
71. Wang Y, Namba S, Lopera-Maya EA, Kerminen S, Tsuo K, Lall K, Kanai M, Zhou W, Wu K-HH, Fave M-J, et al. (2021). Global biobank analyses provide lessons for computing polygenic risk scores across diverse cohorts. Preprint at medRxiv. 10.1101/2021.11.19.21266436.
72. Namba S, Konuma T, Wu K-H, Zhou W, and Okada Y; Global Biobank Meta-analysis Initiative (2022). A practical guideline of genomics-driven drug discovery in the era of global biobank meta-analysis. *Cell Genomics* 2, 100190.
73. Wu K-HH, Douville NJ, Konerman MC, Mathis MR, Hummel SL, Wolford BN, et al. (2021). Polygenic risk score from a multi-ancestry GWAS uncovers susceptibility of heart failure. Preprint at medRxiv. 10.1101/2021.12.06.21267389.
74. Faro VL, Bhattacharya A, Zhou W, Zhou D, Wang Y, Läll K, et al. (2021). Genome-wide association meta-analysis identifies novel ancestry-specific primary open-angle glaucoma loci and shared biology with vascular mechanisms and cell proliferation. Preprint at medRxiv. 10.1101/2021.12.16.21267891.
75. Surakka I, Wu K-H, Hornsby W, Wolford BN, Shen F, Zhou W, et al. (2022). Multi-ancestry meta-analysis identifies 2 novel loci associated with ischemic stroke and reveals heterogeneity of effects between sexes and ancestries. Preprint at medRxiv. 10.1101/2022.02.28.22271647.
76. Wolford BN, Zhao Y, Surakka I, Wu K-HH, Yu X, Richter CE, Bhatta L, Brumpton B, Desch K, Thibord F, et al. (2022). Multi-ancestry GWAS for venous thromboembolism identifies novel loci followed by experimental validation in zebrafish. Preprint at medRxiv. 10.1101/2022.06.21.22276721.
77. Aneas I, Decker DC, Howard CL, Sobreira DR, Sakabe NJ, Blaine KM, Stein MM, Hrusch CL, Montefiori LE, Tena J, et al. (2021). Asthma-associated genetic variants induce IL33 differential expression through an enhancer-blocking regulatory region. *Nat. Commun* 12, 6115. [PubMed: 34675193]
78. Vladich FD, Brazille SM, Stern D, Peck ML, Ghittoni R, and Vercelli D (2005). IL-13 R130Q, a common variant associated with allergy and asthma, enhances effector mechanisms essential for human allergic inflammation. *J. Clin. Invest* 115, 747–754. [PubMed: 15711639]
79. All of Us Research Program Investigators; Denny JC, Rutter JL, Goldstein DB, Philippakis A, Smoller JW, Jenkins G, and Dishman E (2019). The “all of us” Research program. *N. Eng. J. Med* 381, 668–676.
80. Wojcik GL, Graff M, Nishimura KK, Tao R, Haessler J, Gignoux CR, Highland HM, Patel YM, Sorokin EP, Avery CL, et al. (2019). Genetic analyses of diverse populations improves discovery for complex traits. *Nature* 570, 514–518. [PubMed: 31217584]
81. Luo Y, Kanai M, Choi W, Li X, Sakaue S, Yamamoto K, Ogawa K, Gutierrez-Arcelus M, Gregersen PK, Stuart PE, et al. (2021). A high-resolution HLA reference panel capturing global population diversity enables multi-ancestry fine-mapping in HIV host response. *Nat. Genet* 53, 1504–1516. [PubMed: 34611364]
82. Sakaue S, Hosomichi K, Hirata J, Nakaoka H, Yamazaki K, Yawata M, Yawata N, Naito T, Umeno J, Kawaguchi T, et al. (2022). Decoding the diversity of killer immunoglobulin-like receptors by deep sequencing and a high-resolution imputation method. *Cell Genomics* 2, 100101. 10.1016/j.xgen.2022.100101.
83. Mukamel RE, Handsaker RE, Sherman MA, Barton AR, Zheng Y, McCarroll SA, and Loh P-R (2021). Protein-coding repeat polymorphisms strongly shape diverse human phenotypes. *Science* 373, 1499–1505. [PubMed: 34554798]
84. Mägi R, Horikoshi M, Sofer T, Mahajan A, Kitajima H, Franceschini N, McCarthy MI, COGENT-Kidney Consortium T2D-GENES Consortium; Morris AP, and Morris AP (2017). Trans-ethnic meta-regression of genome-wide association studies accounting for ancestry increases power for

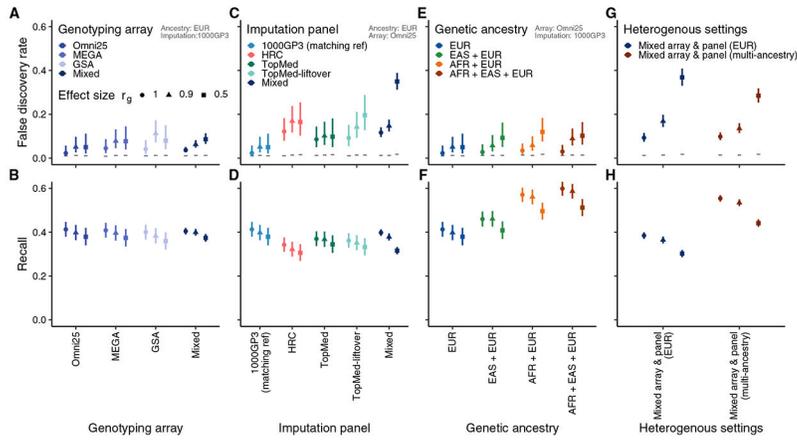
- discovery and improves fine-mapping resolution. *Hum. Mol. Genet* 26, 3639–3650. [PubMed: 28911207]
85. Su Z, Marchini J, and Donnelly P (2011). HAPGEN2: simulation of multiple disease SNPs. *Bioinformatics* 27, 2304–2305. [PubMed: 21653516]
86. Chang CC, Chow CC, Tellier LC, Vattikuti S, Purcell SM, and Lee JJ (2015). Second-generation PLINK: rising to the challenge of larger and richer datasets. *GigaScience* 4, 7. [PubMed: 25722852]
87. Das S, Forer L, Schönherr S, Sidore C, Locke AE, Kwong A, Vrieze SI, Chew EY, Levy S, McGue M, et al. (2016). Next-generation genotype imputation service and methods. *Nat. Genet* 48, 1284–1287. [PubMed: 27571263]
88. Manichaikul A, Mychaleckyj JC, Rich SS, Daly K, Sale M, and Chen W-M (2010). Robust relationship inference in genome-wide association studies. *Bioinformatics* 26, 2867–2873. [PubMed: 20926424]
89. Wei X, and Nielsen R (2019). CCR5- 32 is deleterious in the homozygous state in humans. *Nat. Med* 25, 909–910. [PubMed: 31160814]
90. Maier R, Akbari A, Wei X, Patterson N, Nielsen R, and Reich D (2020). No statistical evidence for an effect of CCR5- 32 on lifespan in the UK Biobank cohort. *Nat. Med* 26, 178–180. [PubMed: 31873311]
91. Loh P-R, Danecek P, Palamara PF, Fuchsberger C, A Reshef Y, K Finucane H, Schoenherr S, Forer L, McCarthy S, Abecasis GR, et al. (2016). Reference-based phasing using the Haplotype reference consortium panel. *Nat. Genet* 48, 1443–1448. [PubMed: 27694958]
92. Loh P-R, Bhatia G, Gusev A, Finucane HK, Bulik-Sullivan BK, Pollack SJ, Schizophrenia Working Group of Psychiatric Genomics Consortium; de Candia TR, Lee SH, Wray NR, et al. (2015). Contrasting genetic architectures of schizophrenia and other complex diseases using fast variance-components analysis. *Nat. Genet* 47, 1385–1392. [PubMed: 26523775]
93. McLaren W, Gil L, Hunt SE, Riat HS, Ritchie GRS, Thormann A, Flicek P, and Cunningham F (2016). The Ensembl variant effect predictor. *Genome Biol.* 17, 122. [PubMed: 27268795]

### Highlights

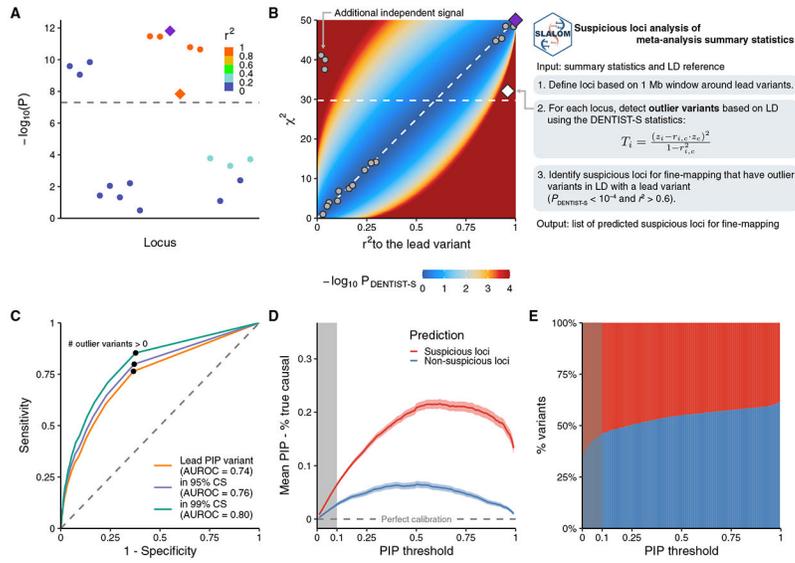
- Extensive simulation of meta-analyses to show substantial fine-mapping miscalibration
- SLALOM, a novel method that identifies suspicious loci for meta-analysis fine-mapping
- Significant depletion of likely causal variants in SLALOM-predicted suspicious loci
- Widespread suspicious loci for fine-mapping in current meta-analysis summary statistics



**Figure 1.** Schematic overview of meta-analysis fine-mapping



**Figure 2. Evaluation of FDR and recall in meta-analysis fine-mapping simulations**  
 We evaluated FDR and recall in meta-analysis fine-mapping using (A–H) different genotyping arrays (A and B), imputation reference panels (C and D), genetic ancestries (E and F), and more heterogeneous settings by combining these (G and H). As shown in top-right gray labels, the EUR ancestry, the Omni2.5 genotyping array, and/or the 1000GP3 reference were the method calibrated. FDR is defined as the proportion of non-causal variants with  $PIP > 0.9$ . Horizontal gray lines represent  $1 - \text{mean PIP}$ ; i.e., expected FDRs were the method calibrated. Recall is defined as the proportion of true causal variants in the top 1% PIP bin. Shapes correspond to the true effect size correlation  $r_g$  across cohorts that represent a phenotypic heterogeneity parameter (the lower  $r_g$ , the higher phenotypic heterogeneity). Error bars correspond to 95% confidence intervals.



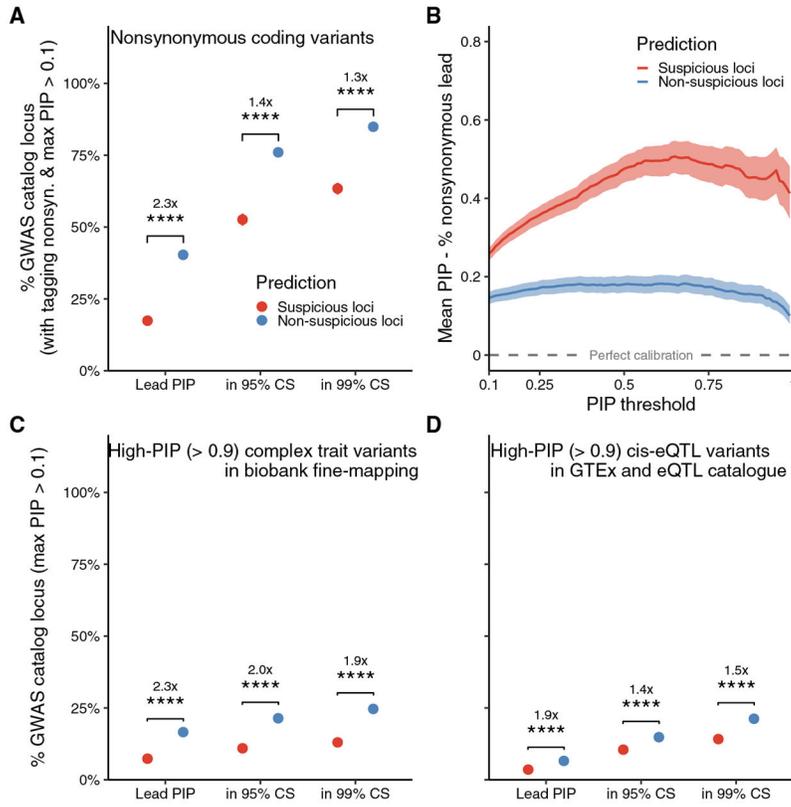
**Figure 3. Overview of the SLALOM method**

(A and B) An illustrative example of the SLALOM application. (A) In an example locus, two independent association signals are depicted: (1) the most significant signal that contains a lead variant (purple diamond) and five additional variants that are in strong LD ( $r^2 > 0.9$ ) with the lead variant, and (2) an additional independent signal ( $r^2 < 0.05$ ). There is one outlier variant (orange diamond) in the first signal that deviates from the expected association based on LD. (B) Step-by-step procedure of the SLALOM method. For outlier variant detection in a locus, a diagnosis plot of  $r^2$  values to the lead variant versus marginal  $\chi^2$  is shown to aid interpretation. Background color represents a theoretical distribution of  $-\log_{10} P_{\text{DENTIST-S}}$  values when a lead variant has a marginal  $\chi^2$  of 50, assuming no allele flipping. Points represent the variants depicted in the example locus (A), where the lead variant (purple diamond) and the outlier variant (white diamond) are highlighted. Diagonal line represents an expected marginal association. Horizontal dotted lines represent the genome-wide significance threshold ( $p < 5.0 \times 10^{-8}$ ).

(C). The receiver operating characteristic (ROC) curve of SLALOM prediction for identifying suspicious loci in the simulations. Positive conditions were defined as whether a true causal variant in a locus is (1) a lead PIP variant, (2) in 95% CS, and (3) in 99% CS. AUROC values are shown in the labels. Black points represent the performance of our adopted metric; i.e., whether a locus contains at least one outlier variant ( $P_{\text{DENTIST-S}} < 1.0 \times 10^{-4}$  and  $r^2 > 0.6$ ).

(D) Calibration plot in the simulations under different PIP thresholds. Calibration was measured as the mean PIP minus the fraction of true causal variants among variants above the threshold. Shadows around the lines represent 95% confidence intervals.

(E) The fraction of variants in predicted suspicious and non-suspicious loci under different PIP thresholds. Gray shadows in the panels (D and E) represent a PIP  $\leq 0.1$  region as we excluded loci with maximum PIP  $\leq 0.1$  in the actual SLALOM analysis based on these panels.

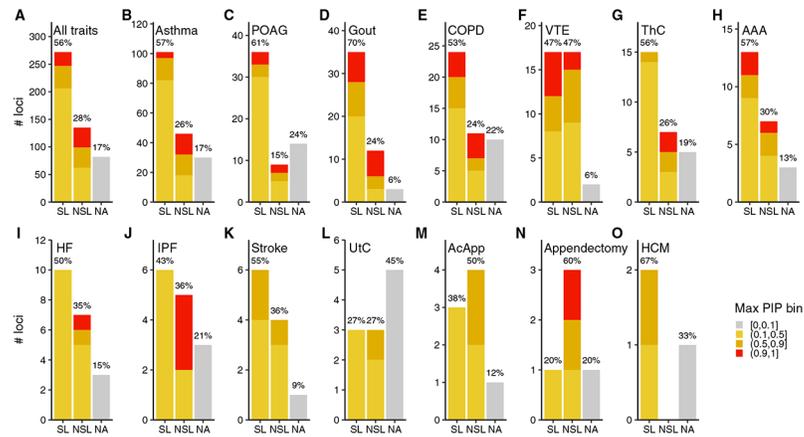


**Figure 4. Evaluation of SLALOM performance in the GWAS Catalog summary statistics**

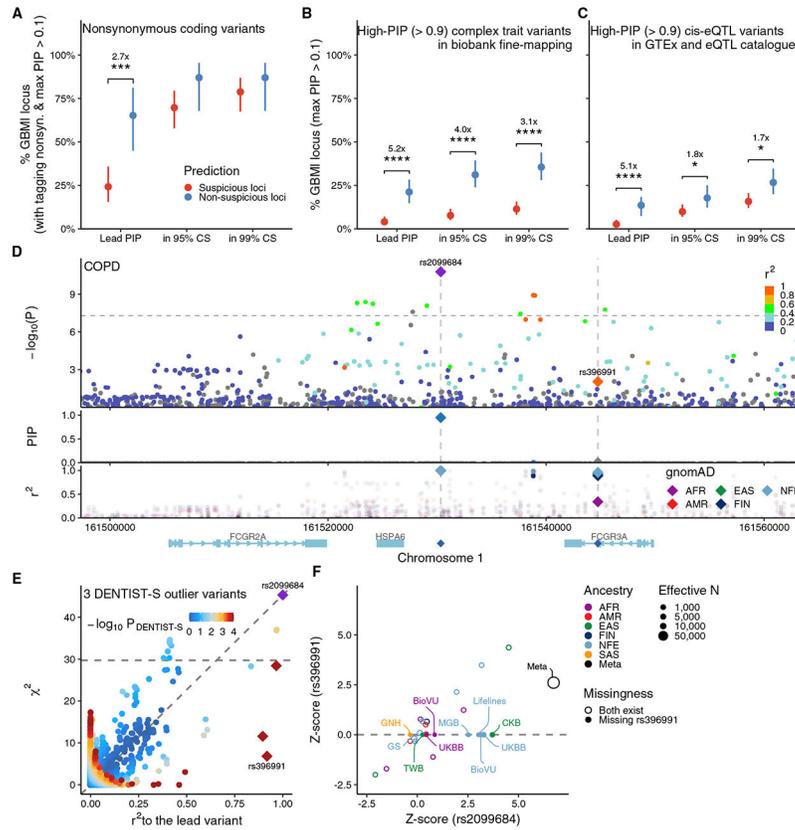
(A) Depletion of likely causal variants in predicted suspicious loci. We evaluated whether nonsynonymous coding variants (pLoF and missense) were lead PIP variants, in 95% CS, or in 99% CS in suspicious versus non-suspicious loci. Depletion was calculated by relative risk (i.e., a ratio of proportions; STAR Methods). Error bars, invisible due to their small size, correspond to 95% confidence intervals using bootstrapping. Significance represents a Fisher exact test p value (\* $p < 0.05$ , \*\* $p < 0.01$ , \*\*\* $p < 0.001$ , \*\*\*\* $p < 10^{-4}$ ).

(B) Plot of the estimated difference between true and reported proportion of causal variants in the loci tagging nonsynonymous variants ( $r^2 > 0.6$  with the lead variants) in the GWAS Catalog under different PIP thresholds. Analogous to Figure 3D, assuming nonsynonymous variants in these loci are truly causal, the mean PIP for lead variants minus the fraction of lead variants that are nonsynonymous above the threshold is equal to the difference between true and reported proportion of causal variants. Shadows around the lines represent 95% confidence intervals.

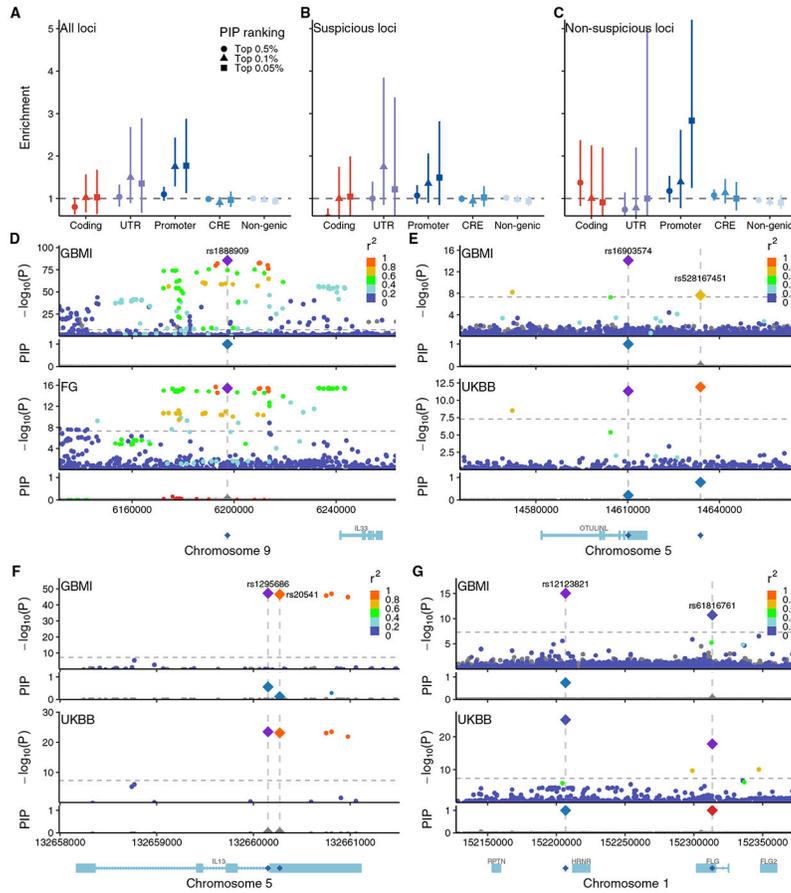
(C and D) Similar to (A), we evaluated whether (C) high-PIP (>0.9) complex trait variants in biobank fine-mapping and (D) high-PIP (>0.9) *cis*-eQTL variants in GTEx v8 and eQTL Catalog were lead PIP variants, in 95% CS, or in 99% CS in suspicious versus non-suspicious loci.



**Figure 5. SLALOM prediction results in the GBMI summary statistics**  
 (A–O) For (A) all 14 traits and (B–O) individual traits, a number of predicted suspicious (SL), non-suspicious (NSL), and non-applicable (NA; maximum PIP <0.1) loci were summarized. Individual traits are ordered by the total number of loci. Color represents the maximum PIP in a locus. Label represents the fraction of loci in each prediction category. AAA, abdominal aortic aneurysm; AcApp, acute appendicitis; COPD, chronic obstructive pulmonary disease; HCM, hypertrophic cardiomyopathy; HF, heart failure; IPF, idiopathic pulmonary fibrosis; POAG, primary open-angle glaucoma; ThC, thyroid cancer; UtC, uterine cancer; VTE, venous thromboembolism.



**Figure 6. Evaluation of SLALOM performance in the GBMI summary statistics**  
 (A–C) Similar to Figure 4, we evaluated whether (A) nonsynonymous coding variants (pLoF and missense), (B) high-PIP (>0.9) complex trait variants in biobank fine-mapping, and (C) high-PIP (>0.9) *cis*-eQTL variants in GTEx v8 and eQTL Catalog were lead PIP variants, in 95% CS, or in 99% CS in suspicious versus non-suspicious loci. Depletion was calculated by relative risk (i.e., a ratio of proportions; STAR Methods). Error bars correspond to 95% confidence intervals using bootstrapping. Significance represents a Fisher exact test p value (\* $p < 0.05$ , \*\* $p < 0.01$ , \*\*\* $p < 0.001$ , \*\*\*\* $p < 10^{-4}$ ).  
 (D) Locuszoom plot of the 1q23.3 locus for COPD. (Top) A Manhattan plot, where the lead variant rs2099684 (purple diamond) and a missense variant rs396991 (orange diamond) are highlighted. Color represents  $r^2$  values to the lead variant. Horizontal line represents a genome-wide significance threshold ( $p = 5.0 \times 10^{-8}$ ). (Middle) PIP from ABF fine-mapping. Color represents whether variants belong to a 95% CS. (Bottom)  $r^2$  values with the lead variant in gnomAD populations.  
 (E) A diagnosis plot showing  $r^2$  values to the lead variant versus marginal  $\chi^2$ . Color represents  $-\log_{10} P_{\text{DENTIST-S}}$  values. Outlier variants with  $P_{\text{DENTIST-S}} < 10^{-4}$  are depicted in red with a diamond shape. Diagonal line represents an expected marginal association. Horizontal line represents a genome-wide significance threshold.  
 (F) Z-scores of the lead variant (rs2099684) versus the missense variant (rs396991) in the constituent cohorts of the meta-analysis. Open and closed circles represent whether both variants exist in a cohort or rs396991 is missing. Circle size corresponds to an effective sample size. Color represents genetic ancestry.



**Figure 7. Fine-mapping improvement and retrogression in the GBMI meta-analyses over individual biobanks**

(A–C) Functional enrichment of variants in each functional category based on top PIP rankings in the GBMI and individual biobanks (maximum PIP of BBJ, FinnGen, and UKBB) using (A) all loci, (B) suspicious loci, or (C) non-suspicious loci. Shape corresponds to top PIP ranking (top 0.5%, 0.1%, and 0.05%). Enrichment was calculated by a relative risk (i.e., a ratio of proportions; STAR Methods). Error bars correspond to 95% confidence intervals using bootstrapping. (D and E) Locuszoom plots for the same non-suspicious locus of asthma in the GBMI meta-analysis and an individual biobank (BBJ, FinnGen, or UKBB Europeans) that showed the highest PIP in our biobank fine-mapping. Colors in the Manhattan panels represent  $r^2$  values to the lead variant. In the PIP panels, only fine-mapped variants in the 95% CS are colored, where the same colors are applied between the GBMI meta-analysis and an individual biobank based on merged CS as previously described. Horizontal line represents a genome-wide significance threshold ( $p = 5.0 \times 10^{-8}$ ).

(D) rs1888909 for asthma in the GBMI and FinnGen.

(E) rs16903574 for asthma in the GBMI and UKBB Europeans. Nearby rs528167451 was also highlighted, which was in strong LD ( $r^2 = 0.86$ ) and in the same 95% CS in UKBB Europeans, but not in the GBMI ( $r^2 = 0.67$ ).

(F) rs1295686 for asthma in the GBMI and UKBB Europeans. A nearby missense, rs20541, showed lower PIP than rs1295686 despite having strong LD ( $r^2 = 0.96$ ).

(G) rs12123821 for asthma in the GBMI and UKBB Europeans. Nearby stop-gained rs61816761 was independent of rs12123821 ( $r^2 = 0.0$ ) and not fine-mapped in the GBMI due to a single causal variant assumption in the ABF fine-mapping.

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript

## KEY RESOURCES TABLE

REAGENT or RESOURCE	SOURCE	IDENTIFIER
Deposited Data		
GBMI summary statistics	Zhou, W. et al., 2022 <sup>10</sup>	<a href="https://www.globalbiobankmeta.org/resources">https://www.globalbiobankmeta.org/resources</a>
BBJ fine-mapping results	Kanai, M. et al., 2021 <sup>16</sup>	<a href="https://humandbs.biosciencedbc.jp/en/hum0197-latest#hum0197.v5.gwas.v1">https://humandbs.biosciencedbc.jp/en/hum0197-latest#hum0197.v5.gwas.v1</a>
FinnGen fine-mapping results	Kanai, M. et al., 2021 <sup>16</sup>	<a href="https://www.finnngen.fi/en/access_results">https://www.finnngen.fi/en/access_results</a>
UKBB fine-mapping results	Kanai, M. et al., 2021 <sup>16</sup>	<a href="https://www.finucanelab.org/data">https://www.finucanelab.org/data</a>
GWAS Catalog	GWAS Catalog (as of January 12, 2022)	<a href="https://www.ebi.ac.uk/gwas/">https://www.ebi.ac.uk/gwas/</a>
Example outputs from the meta-analysis fine-mapping simulation pipeline	This study	<a href="https://doi.org/10.7910/DVN/M86OCQ">https://doi.org/10.7910/DVN/M86OCQ</a>
Software and Algorithms		
SLALOM	This study	<a href="https://github.com/mkanai/slalom">https://github.com/mkanai/slalom</a> , <a href="https://doi.org/10.5281/zenodo.6984388">https://doi.org/10.5281/zenodo.6984388</a>
Meta-analysis fine-mapping simulation pipeline	This study	<a href="https://github.com/mkanai/meta-finemapping-simulation">https://github.com/mkanai/meta-finemapping-simulation</a> , <a href="https://doi.org/10.5281/zenodo.6984391">https://doi.org/10.5281/zenodo.6984391</a>
Analysis code	This study	<a href="https://github.com/mkanai/slalom-paper">https://github.com/mkanai/slalom-paper</a> , <a href="https://doi.org/10.5281/zenodo.7010731">https://doi.org/10.5281/zenodo.7010731</a>
HAPGEN2	Su, Z. et al., 2011 <sup>85</sup>	<a href="https://mathgen.stats.ox.ac.uk/genetics_software/hapgen/hapgen2.html">https://mathgen.stats.ox.ac.uk/genetics_software/hapgen/hapgen2.html</a>
PLINK2.0	Chang, CC. et al., 2015 <sup>86</sup>	<a href="https://www.cog-genomics.org/plink/2.0/">https://www.cog-genomics.org/plink/2.0/</a>
Michigan Imputation Server	Das, S. et al., 2016 <sup>87</sup>	<a href="https://imputationserver.sph.umich.edu/">https://imputationserver.sph.umich.edu/</a>
TOPMed Imputation Server	Taliun, D. et al., 2021 <sup>50</sup>	<a href="https://imputation.biodatacatalyst.nhlbi.nih.gov/">https://imputation.biodatacatalyst.nhlbi.nih.gov/</a>
Hail	Hail team, 2022	<a href="https://hail.is/">https://hail.is/</a>