

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

Masahiro Kanai^{1,2,3,4,5,*}, Roy Elzur^{1,2,3}, Wei Zhou^{1,2,3}, Global Biobank Meta-analysis Initiative, Mark J Daly^{1,2,3,6}, Hilary K Finucane^{1,2,3,*}

¹Analytic and Translational Genetics Unit, Massachusetts General Hospital, Boston, MA, USA, ²Program in Medical and Population Genetics, Broad Institute of MIT and Harvard, Cambridge, MA, USA, ³Stanley Center for Psychiatric Research, Broad Institute of MIT and Harvard, Cambridge, MA, USA, ⁴Department of Biomedical Informatics, Harvard Medical School, Boston, MA, USA, ⁵Department of Statistical Genetics, Osaka University Graduate School of Medicine, Suita, Japan, ⁶Institute for Molecular Medicine Finland (FIMM), University of Helsinki, Helsinki, Finland.

* Corresponding authors: Masahiro Kanai (mkanai@broadinstitute.org) and Hilary K Finucane (finucane@broadinstitute.org)

Abstract

Meta-analysis is pervasively used to combine multiple genome-wide association studies (GWAS) into a more powerful whole. To resolve causal variants, meta-analysis studies typically apply summary statistics-based fine-mapping methods as they are applied to single-cohort studies. However, it is unclear whether heterogeneous characteristics of each cohort (*e.g.*, ancestry, sample size, phenotyping, genotyping, or imputation) affect fine-mapping calibration and recall. Here, we first demonstrate that meta-analysis fine-mapping is substantially miscalibrated in simulations when different genotyping arrays or imputation panels are included. To mitigate these issues, we propose a summary statistics-based QC method, SLALOM, that identifies suspicious loci for meta-analysis fine-mapping by detecting outliers in association statistics based on ancestry-matched local LD structure. Having validated SLALOM performance in simulations and the GWAS Catalog, we applied it to 14 disease endpoints from the Global Biobank Meta-analysis Initiative and found that 68% of loci showed suspicious patterns that call into question fine-mapping accuracy. These predicted suspicious loci were significantly depleted for having likely causal variants, such as nonsynonymous variants, as a lead variant (2.8x; Fisher's exact $P = 6.2 \times 10^{-4}$). Compared to fine-mapping results in individual biobanks, we found limited evidence of fine-mapping improvement in the GBMI meta-analyses. Although a full solution requires complete synchronization across cohorts, our approach identifies likely spurious results in meta-analysis fine-mapping. We urge extreme caution when interpreting fine-mapping results from meta-analysis.

1 Introduction

2 Meta-analysis is pervasively used to combine multiple genome-wide association studies (GWAS)
3 from different cohorts¹. Previous GWAS meta-analyses have identified thousands of loci
4 associated with complex diseases and traits, such as type 2 diabetes^{2,3}, schizophrenia^{4,5},
5 rheumatoid arthritis^{6,7}, body mass index⁸, and lipid levels⁹. These meta-analyses are typically
6 conducted in large-scale consortia (e.g., the Psychiatric Genomics Consortium [PGC], the Global
7 Lipids Genetics Consortium [GLGC], and the Genetic Investigation of Anthropometric Traits
8 [GIANT] consortium) to increase sample size while harmonizing analysis plans across
9 participating cohorts in every possible aspect (e.g., phenotype definition, quality-control [QC]
10 criteria, statistical model, and analytical software) by sharing summary statistics as opposed to
11 individual-level data, thereby avoiding data protection issues and variable legal frameworks
12 governing individual genome and medical data around the world. The Global Biobank Meta-
13 analysis Initiative (GBMI)¹⁰ is one such large-scale, international effort, which aims to establish a
14 collaborative network spanning 19 biobanks from four continents (total $n = 2.1$ million) for
15 coordinated GWAS meta-analyses, while addressing the many benefits and challenges in meta-
16 analysis and subsequent downstream analyses.

17
18 One such challenging downstream analysis is statistical fine-mapping^{11–13}. Despite the great
19 success of past GWAS meta-analyses in locus discovery, individual causal variants in associated
20 loci are largely unresolved. Identifying causal variants from GWAS associations (i.e., fine-
21 mapping) is challenging due to extensive linkage disequilibrium (LD, the correlation among
22 genetic variants), the presence of multiple causal variants, and limited sample sizes, but is rapidly
23 becoming achievable with high confidence in individual cohorts^{14–17} owing to the recent
24 development of large-scale biobanks^{18–20} and scalable fine-mapping methods^{21–23} that enable
25 well-powered, accurate fine-mapping using in-sample LD from large-scale individual-level data.

26
27 After conducting GWAS meta-analysis, previous studies^{2,7,9,24–30} have applied existing summary
28 statistics-based fine-mapping methods (e.g., approximate Bayes factor [ABF]^{31,32}, CAVIAR³³,
29 PAINTOR^{34,35}, FINEMAP^{21,22}, and SuSiE²³) just as they are applied to single-cohort studies,
30 without considering or accounting for the unavoidable heterogeneity among cohorts (e.g.
31 differences in sample size, phenotyping, genotyping, or imputation). Such heterogeneity could
32 lead to false positives and miscalibration in meta-analysis fine-mapping (**Fig. 1**). For example,
33 case-control studies enriched with more severe cases or ascertained with different phenotyping
34 criteria may disproportionately contribute to genetic discovery, even when true causal effects for
35 genetic liability are exactly the same between these studies and less severe or unascertained
36 ones. Quantitative traits like biomarkers could have phenotypic heterogeneity arising from
37 different measurement protocols and errors across studies. There might be genuine biological
38 mechanisms too, such as gene–gene (GxG) and gene–environment (GxE) interactions and
39 (population-specific) dominance variation (e.g., rs671 and alcohol dependence³⁶), that introduce
40 additional heterogeneity across studies^{37,38}. In addition to phenotyping, differences in genotyping
41 and imputation could dramatically undermine fine-mapping calibration and recall at single-variant
42 resolution, because differential patterns of missingness and imputation quality across constituent
43 cohorts of different sample sizes can disproportionately diminish association statistics of
44 potentially causal variants. Finally, although more easily harmonized than phenotyping and
45 genotyping data, subtle differences in QC criteria and analytical software may further exacerbate
46 the effect of heterogeneity on fine-mapping.

47
48 An illustrative example of such issues can be observed in the *TYK2* locus (19p13.2) in the recent
49 meta-analysis from the COVID-19 Host Genetics Initiative (COVID-19 HGI; **Supplementary Fig.**

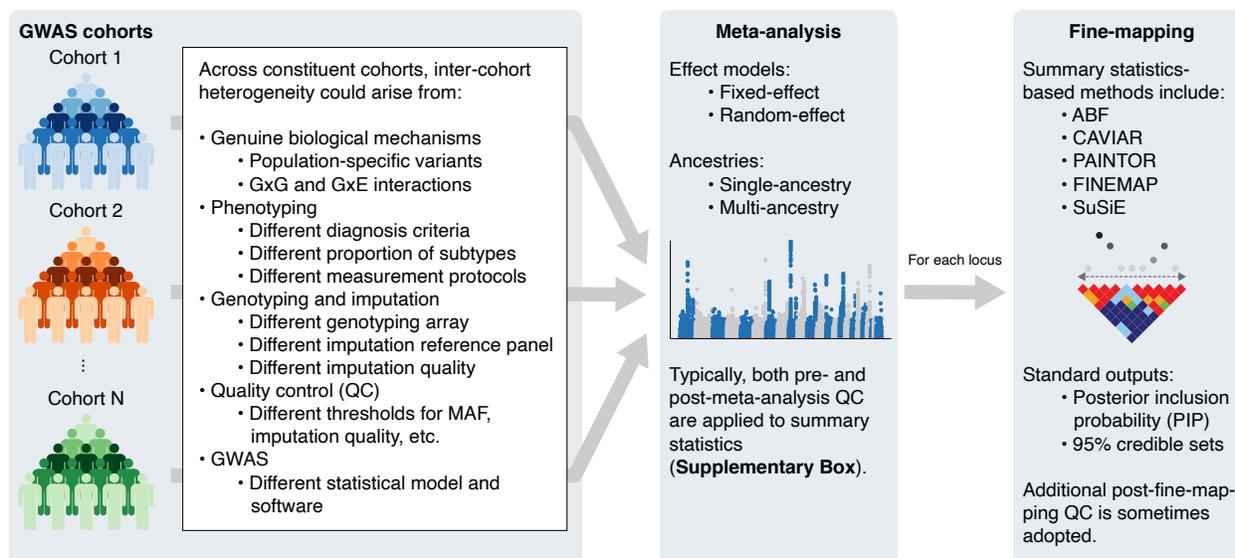
50 1)³⁹. This locus is known for protective associations against autoimmune diseases^{6,24}, while a
51 complete *TYK2* loss of function results in a primary immunodeficiency⁴⁰. Despite strong LD ($r^2 =$
52 0.82) with a lead variant in a locus (rs74956615; $P = 9.7 \times 10^{-12}$), a known functional missense
53 variant rs34536443 (p.Pro1104Ala) that reduces *TYK2* function^{41,42} did not achieve genome-wide
54 significance ($P = 7.5 \times 10^{-7}$), primarily due to its missingness in two more cohorts than
55 rs74956615. This serves as just one example of the major difficulties with meta-analysis fine-
56 mapping at single-variant resolution. Indeed, the COVID-19 HGI cautiously avoided an *in-silico*
57 fine-mapping in the flagship to prevent spurious results³⁹.

58
59 Only a few studies have carefully addressed these concerns in their downstream analyses. The
60 Schizophrenia Working Group of PGC, for example, recently updated their largest meta-analysis
61 of schizophrenia⁵ (69,369 cases and 236,642 controls), followed by a downstream fine-mapping
62 analysis using FINEMAP²¹. Unlike many other GWAS consortia, since PGC has access to
63 individual-level genotypes for a majority of samples, they were able to apply standardized sample
64 and variant QC criteria and impute variants using the same reference panel, all uniformly
65 processed using the RICOPILI pipeline⁴³. This harmonized procedure was crucial for properly
66 controlling inter-cohort heterogeneity and thus allowing more robust meta-analysis fine-mapping
67 at single-variant resolution. Furthermore, PGC's direct access to individual-level data enabled
68 them to compute in-sample LD matrices for multiple causal variant fine-mapping, which prevents
69 the significant miscalibration that results from using an external LD reference¹⁴⁻¹⁶. A 2017 fine-
70 mapping study of inflammatory bowel disease also benefited from access to individual-level
71 genotypes and careful pre- and post-fine-mapping QC⁴⁴. For a typical meta-analysis consortium,
72 however, many of these steps are infeasible as full genotype data from all cohorts is not available.
73 For such studies, a new approach to meta-analysis fine-mapping in the presence of the many
74 types of heterogeneity is needed. Until such a method is developed, QC of meta-analysis fine-
75 mapping results deserves increased attention.

76
77 While existing variant-level QC procedures are effective for limiting spurious associations in
78 GWAS (**Supplementary Box**)⁴⁵, they do not suffice for ensuring high-quality fine-mapping results.
79 In some cases, they even hurt fine-mapping quality, because they can i) cause or exacerbate
80 differential patterns of missing variants across cohorts, and ii) remove true causal variants as well
81 as suspicious variants. Thus, additional QC procedures that retain consistent variants across
82 cohorts for consideration but limit poor-quality fine-mapping results are needed. A recently
83 proposed method called DENTIST⁴⁶, for example, performs summary statistics QC to improve
84 GWAS downstream analyses, such as conditional and joint analysis (GCTA-COJO⁴⁷), by
85 removing variants based on estimated heterogeneity between summary statistics and reference
86 LD. Although DENTIST was also applied prior to fine-mapping (FINEMAP²¹), simulations only
87 demonstrated that it could improve power for detecting the correct number of causal variants in a
88 locus, not true causal variants. This motivated us to develop a new fine-mapping QC method for
89 better calibration and recall at single-variant resolution and to demonstrate its performance in
90 large-scale meta-analysis.

91
92 Here, we first demonstrate the effect of inter-cohort heterogeneity in meta-analysis fine-mapping
93 via realistic simulations with multiple heterogeneous cohorts, each with different combinations of
94 genotyping platforms, imputation reference panels, and genetic ancestries. We propose a
95 summary statistics-based QC method, SLALOM (ssuspicious loci analysis of meta-analysis
96 summary statistics), that identifies suspicious loci for meta-analysis fine-mapping by detecting
97 association statistics outliers based on local LD structure, building on the DENTIST method.
98 Applying SLALOM to 14 disease endpoints from the Global Biobank Meta-analysis Initiative¹⁰ as
99 well as 467 meta-analysis summary statistics from the GWAS Catalog⁴⁸, we demonstrate that

100 suspicious loci for fine-mapping are widespread in meta-analysis and urge extreme caution when
101 interpreting fine-mapping results from meta-analysis.
102
103



104
105 **Fig. 1 | Schematic overview of meta-analysis fine-mapping.**
106

107 Results

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

109 Existing fine-mapping methods^{21,23,31} assume that all association statistics are derived from a
110 single-cohort study, and thus do not model the per-variant heterogeneity in effect sizes and
111 sample sizes that arise when meta-analyzing multiple cohorts (**Figure 1**). To evaluate how
112 different characteristics of constituent cohorts in a meta-analysis affect fine-mapping calibration
113 and recall, we conducted a series of large-scale GWAS meta-analysis and fine-mapping
114 simulations (**Supplementary Table 1–4; Methods**). Briefly, we simulated multiple GWAS cohorts
115 of different ancestries (10 European ancestry, one African ancestry and one East Asian ancestry
116 cohorts; $n = 10,000$ each) that were genotyped and imputed using different genotyping arrays
117 (Illumina Omni2.5, Multi-Ethnic Global Array [MEGA], and Global Screening Array [GSA]) and
118 imputation reference panels (the 1000 Genomes Project Phase 3 [1000GP3]⁴⁹, the Haplotype
119 Reference Consortium [HRC]⁵⁰, and the TOPMed⁵¹). For each combination of cohort, genotyping
120 array, and imputation panel, we conducted 300 GWAS with randomly simulated causal variants
121 that resemble the genetic architecture of a typical complex trait, including minor allele frequency
122 (MAF) dependent causal effect sizes⁵², total SNP heritability⁵³, functional consequences of causal
123 variants¹⁷, and levels of genetic correlation across cohorts (*i.e.*, true effect size heterogeneity; r_g
124 = 1, 0.9, and 0.5; see **Methods**). We then meta-analyzed the single-cohort GWAS results across
125 10 independent cohorts based on multiple *configurations* (different combinations of genotyping
126 arrays and imputation panels for each cohort) to resemble realistic meta-analysis of multiple
127 heterogeneous cohorts (**Supplementary Table 4**). We applied ABF fine-mapping to compute a
128 posterior inclusion probability (PIP) for each variant and to derive 95% and 99% credible sets
129 (CS) that contain the smallest set of variants covering 95% and 99% of probability of causality.
130 We evaluated the false discovery rate (FDR, defined as the proportion of variants with PIP > 0.9
131 that are non-causal) and compared against the expected proportion of non-causal variants if the

132 meta-analysis fine-mapping method were calibrated, based on PIP. More details of our simulation
133 pipeline are described in **Methods** and visually summarized in **Supplementary Fig. 2**.

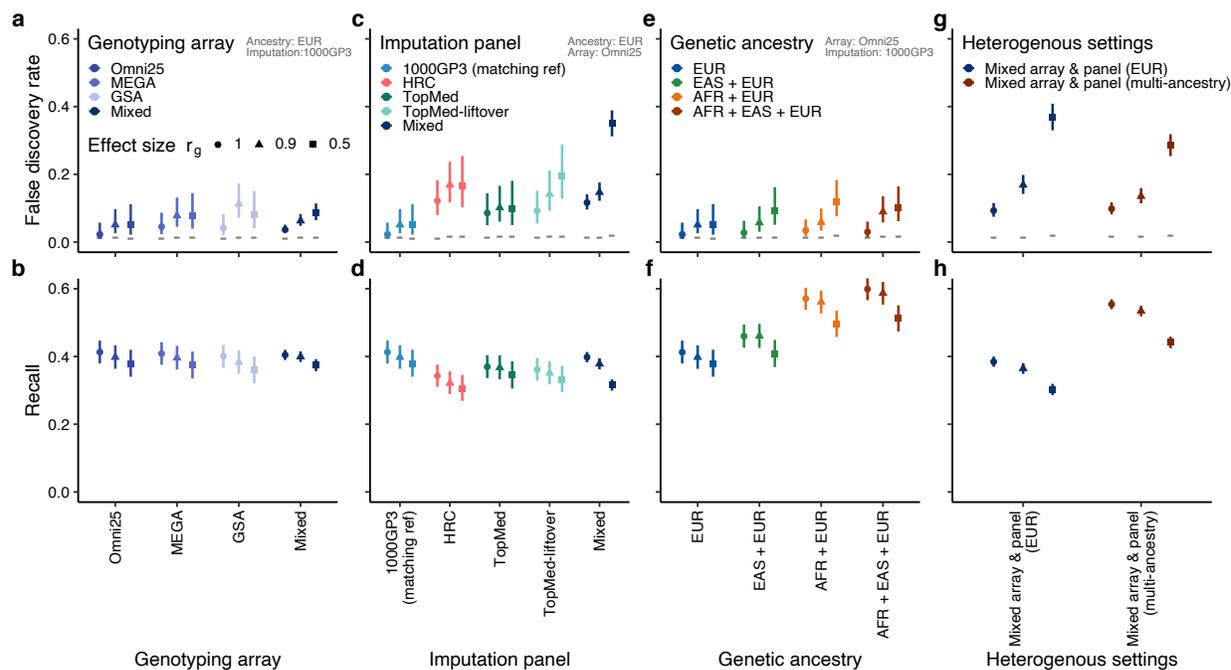
134
135 We found that FDR varied widely over the different configurations, reaching as high as 37% for
136 the most heterogeneous configurations (**Fig. 2**). We characterized the contributing factors to the
137 miscalibration. We first found that lower true effect size correlation r_g (*i.e.*, larger phenotypic
138 heterogeneity) always caused higher miscalibration and lower recall. Second, when using the
139 same imputation panel (1000GP3), use of less dense arrays (MEGA or GSA) led to moderately
140 inflated FDR (up to FDR = 11% vs. expected 1%), while use of multiple genotyping array did not
141 cause further FDR inflation (**Fig. 2a**). Third, when using the same genotyping array (Omni2.5),
142 use of imputation panels (HRC or TOPMed) that does not match our simulation reference
143 significantly affects miscalibration (up to FDR = 17% vs. expected 1%), and using multiple
144 imputation panels further increased miscalibration (up to FDR = 35% vs. expected 2%, **Fig. 2c**);
145 this setup is as bad as the most heterogeneous configuration using multiple genotyping arrays
146 and imputation panels (FDR = 37%). When TOPMed-imputed variants were lifted over from
147 GRCh38 to GRCh37, we observed FDR increases of up to 10%, likely due to genomic build
148 conversion failures (**Supplementary Note**)⁵⁴. Fourth, recall was not significantly affected by
149 heterogeneous genotyping arrays or imputation panels (**Fig. 2b,d**). Fifth, including multiple
150 genetic ancestries did not affect calibration when using the same genotyping array and imputation
151 panel (Omni 2.5 and 1000GP3; **Fig. 2e**) but significantly improved recall if African ancestry was
152 included (**Fig. 2f**). This is expected, given the shorter LD length in the African population
153 compared to other populations, which improves fine-mapping resolution⁵⁵. Finally, in the most
154 heterogeneous configurations where multiple genotyping arrays and imputation panels existed,
155 we observed a FDR of up to 37% and 28% for European and multi-ancestry meta-analyses,
156 respectively (vs. expected 2% for both), demonstrating that inter-cohort heterogeneity can
157 substantially undermine meta-analysis fine-mapping (**Fig. 2g,h**).

158
159 To further characterize observed miscalibration in meta-analysis fine-mapping, we investigated
160 the availability of GWAS variants in each combination of ancestry, genotyping array, and
161 imputation panel. Out of 3,285,617 variants on chromosome 3 that passed variant QC in at least
162 one combination (per-combination MAF > 0.001 and Rsq > 0.6; **Methods**), 574,261 variants
163 (17%) showed population-level gnomAD MAF > 0.001 in every ancestry that we simulated
164 (African, East Asian, and European). Because we used a variety of imputation panels, we
165 retrieved population-level MAF from gnomAD. Of these 574,261 variants, 389,219 variants (68%)
166 were available in every combination (**Supplementary Fig. 3a**). This fraction increased from 68%
167 to 73%, 74%, and 76% as we increased gnomAD MAF thresholds to > 0.005, 0.01, and 0.05,
168 respectively, but never reached 100% (**Supplementary Fig. 4**). Notably, we observed a
169 substantial number of variants that are unique to a certain genotyping array and an imputation
170 panel, even when we restricted to 344,497 common variants (gnomAD MAF > 0.05) in every
171 ancestry (**Supplementary Fig. 3b**). For example, there are 34,317 variants (10%) that were
172 imputed in the 1000GP3 and TOPMed reference but not in the HRC. Likewise, we observed
173 33,106 variants (10%) that were specific to the 1000GP3 reference and even 3,066 variants (1%)
174 that were imputed in every combination except for East Asian ancestry with the GSA array and
175 the TOPMed reference. When using different combinations of gnomAD MAF thresholds (> 0.001,
176 0.005, 0.01, or 0.05 in every ancestry) and Rsq thresholds (> 0.2, 0.4, 0.6, or 0.8), we observed
177 the largest fraction of shared variants (78%) was achieved with gnomAD MAF > 0.01 and Rsq >
178 0.2 while the largest number of the shared variants (427,494 variants) was achieved with gnomAD
179 MAF > 0.001 and Rsq > 0.2, leaving it unclear which thresholds would be preferable in the context
180 of fine-mapping (**Supplementary Fig. 4**).

181

182 The remaining 2,711,356 QC-passing variants in our simulations (gnomAD MAF ≤ 0.001 in at
 183 least one ancestry) further exacerbate variable coverage of the available variants
 184 (**Supplementary Fig. 3c**). Of these, the largest proportion of variants (39%) were only available
 185 in African ancestry, followed by African and European (but not in East Asian) available variants
 186 (7%), European-specific variants (6%), and East Asian-specific variants (5%). Furthermore,
 187 similar to the aforementioned common variants, we found a substantial number of variants that
 188 are unique to a certain combination. Altogether, we observed that only 393,471 variants (12%)
 189 out of all the QC-passing 3,285,617 variants were available in every combination
 190 (**Supplementary Fig. 3d**). These observations recapitulate that different combinations of genetic
 191 ancestry, genotyping array, imputation panels, and QC thresholds substantially affect the
 192 availability of common, well-imputed variants for association testing⁵³.

193
 194 Thus, the different combinations of genotyping and imputation cause each cohort in a meta-
 195 analysis to have a different set of variants, and consequently variants can have very different
 196 overall sample sizes. In our simulations with the most heterogeneous configurations, we found
 197 that 66% of the false positive loci (where a non-causal [false positive] variant was assigned PIP
 198 > 0.9) had different sample sizes for true causal and false positive variants (median
 199 maximum/minimum sample size ratio = 1.4; **Supplementary Fig. 5**). Analytically, we found that
 200 at common meta-analysis sample sizes and genome-wide significant effect size regimes, when
 201 two variants have similar marginal effects, the one with the larger sample size will usually achieve
 202 a higher ABF PIP (**Supplementary Note**). This elucidates the mechanism by which sample size
 203 imbalance can lead to miscalibration.
 204



205
 206 **Fig. 2 | Evaluation of false discovery rate (FDR) and recall in meta-analysis fine-mapping simulations.** We
 207 evaluated FDR and recall in meta-analysis fine-mapping using different genotyping arrays (a,b), imputation
 208 reference panels (c,d), genetic ancestries (e, f), and more heterogeneous settings by combining these (g, h). As shown in top-
 209 right gray labels, the EUR ancestry, the Omni2.5 genotyping array and/or the 1000GP3 reference panel were used
 210 unless otherwise stated. FDR is defined as the proportion of non-causal variants with PIP > 0.9 . Horizontal gray lines
 211 represent $1 - \text{mean PIP}$, i.e. expected FDR were the method calibrated. Recall is defined as the proportion of true
 212 causal variants in the top 1% PIP bin. Shapes correspond to the true effect size correlation r_g across cohorts which
 213 represent a phenotypic heterogeneity parameter (the lower r_g , the higher phenotypic heterogeneity).

214 Overview of the SLALOM method

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

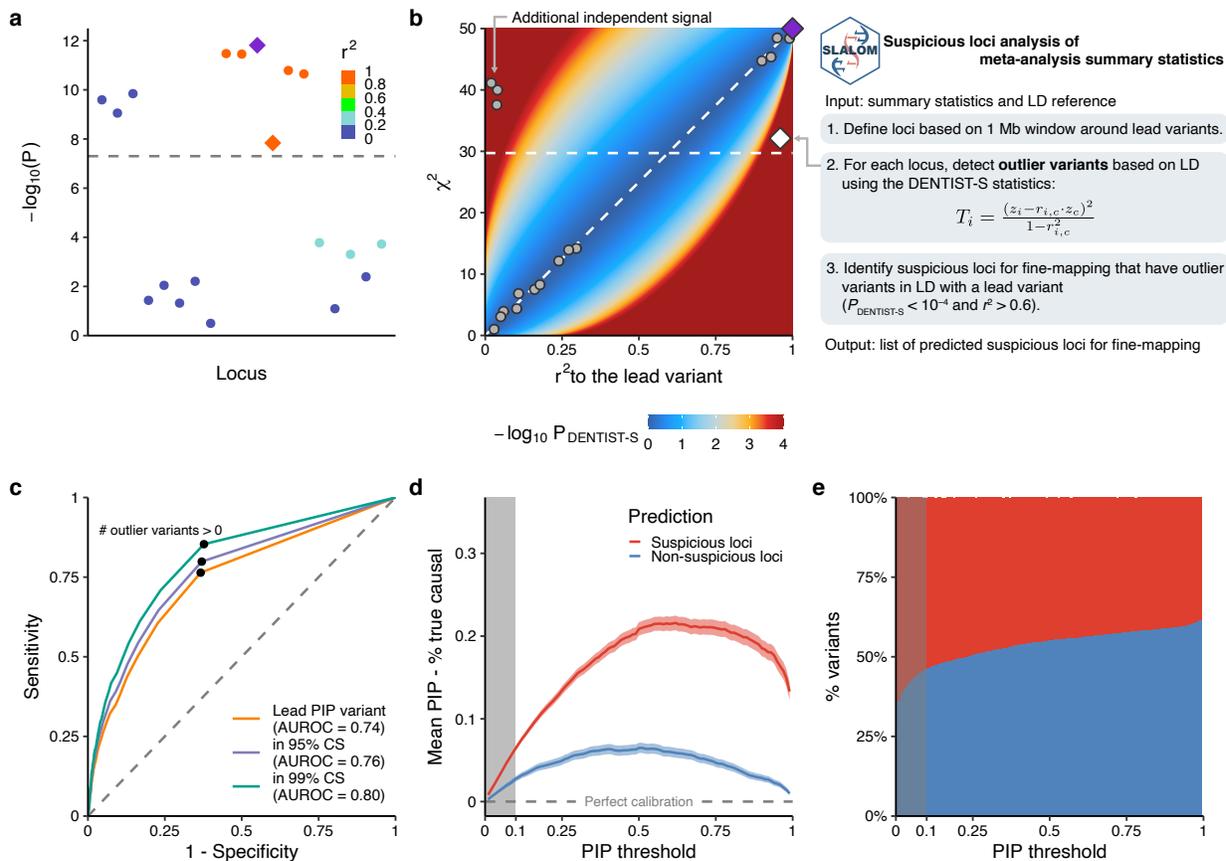
222
223 To detect outlier variants, we first assume a single causal variant per associated locus. Then the
224 marginal z-score z_i for a variant i should be approximately equal to $r_{i,c} \cdot z_c$ where z_c is the z-score
225 of the causal variant c , and $r_{i,c}$ is a correlation between variants i and c . For each variant in meta-
226 analysis summary statistics, we first test this relationship using a simplified version of the
227 DENTIST statistics⁴⁶, DENTIST-S, based on the assumption of a single causal variant. The
228 DENTIST-S statistics for a given variant i is written as

$$229 \quad T_i = \frac{(z_i - r_{i,c} \cdot z_c)^2}{1 - r_{i,c}^2} \quad (1)$$

230 which approximately follows a χ^2 distribution with 1 degree of freedom⁴⁶. Since the true causal
231 variant and LD structure are unknown in real data, we approximate the causal variant as the lead
232 PIP variant in the locus (the variant with the highest PIP) and use a large-scale external LD
233 reference from gnomAD⁵⁷, either an ancestry-matched LD for a single-ancestry meta-analysis or
234 a sample-size-weighted LD by ancestries for a multi-ancestry meta-analysis (**Methods**).

235
236 SLALOM then evaluates whether each locus is “suspicious”—that is, has a pattern of meta-
237 analysis statistics and LD that appear inconsistent and therefore call into question the fine-
238 mapping accuracy. By training on loci with maximum PIP > 0.9 in the simulations, we determined
239 that the best-performing criterion for classifying loci as true or false positives is whether a locus
240 has a variant with $r^2 > 0.6$ to the lead and DENTIST-S P -value < 1.0×10^{-4} (**Methods**). Using this
241 criterion we achieved an area under the receiver operating characteristic curve (AUROC) of 0.74,
242 0.76, and 0.80 for identifying whether a true causal variant is a lead PIP variant, in 95% credible
243 set (CS), and in 99% CS, respectively (**Fig. 3c**). We further validated the performance of SLALOM
244 using all the loci in the simulations and observed significantly higher miscalibration in predicted
245 suspicious loci than in non-suspicious loci (up to 16% difference in FDR at PIP > 0.9; **Fig. 3d**).
246 Given the relatively lower miscalibration and specificity at low PIP thresholds (**Fig. 3d,e**), in
247 subsequent real data analysis we restricted the application of SLALOM to loci with maximum PIP
248 > 0.1 (**Methods**).

249



250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269

Fig. 3 | Overview of the SLALOM method. **a,b.** An illustrative example of the SLALOM application. **a.** In an example locus, two independent association signals are depicted: i) 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 ii) 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 vs. marginal χ^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 χ^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) were 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 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 (AUROC = 0.74), 2) in 95% CS (AUROC = 0.76), and 3) in 99% CS (AUROC = 0.80). 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 – 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,e** represent a PIP ≤ 0.1 region as we excluded loci with maximum PIP ≤ 0.1 in the actual SLALOM analysis based on these panels.

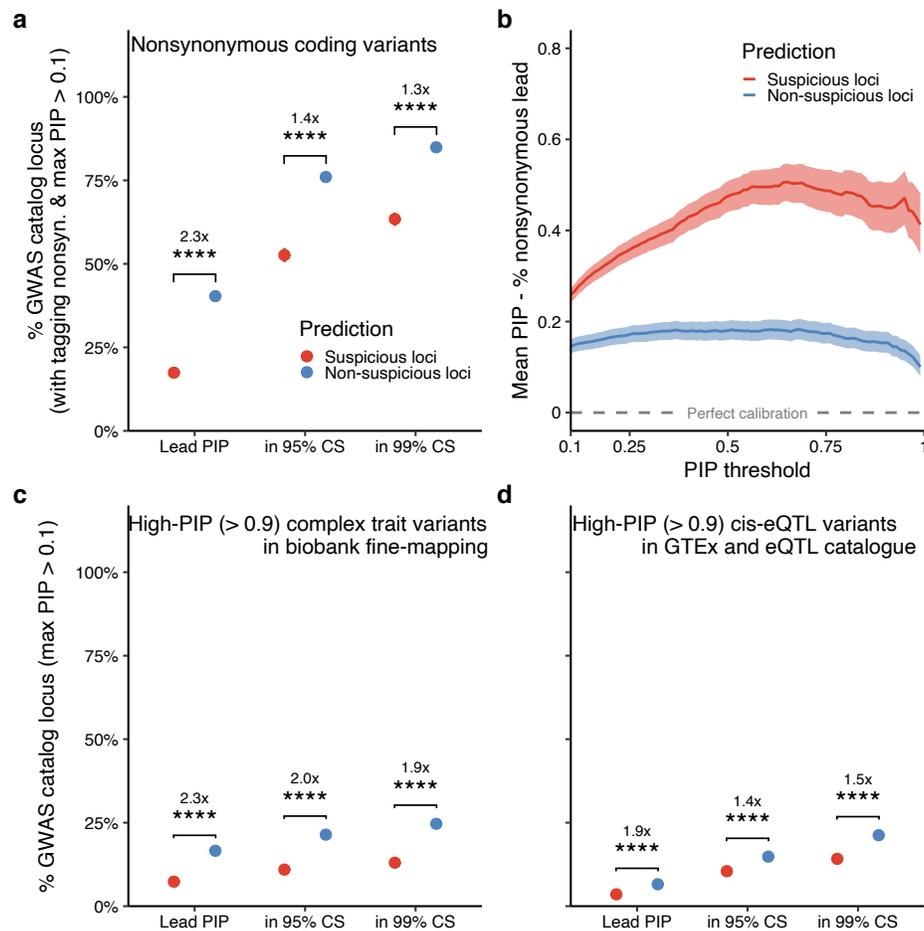
270 **Widespread suspicious loci for fine-mapping in existing meta-analysis summary** 271 **statistics**

272 Having assessed the performance of SLALOM in simulations, we applied SLALOM to 467 meta-
273 analysis summary statistics in the GWAS Catalog⁴⁸ that are publicly available with a sufficient
274 discovery sample size ($N > 10,000$; **Supplementary Table 5; Methods**) to quantify the
275 prevalence of suspicious loci in existing studies. These summary statistics were mostly European
276 ancestry-only meta-analysis (63%), followed by multi-ancestry (31%), East Asian ancestry-only
277 (3%), and African ancestry-only (2%) meta-analyses. Across 467 summary statistics from 96
278 publications, we identified 28,925 loci with maximum PIP > 0.1 (out of 35,864 genome-wide
279 significant loci defined based on 1 Mb window around lead variants; **Methods**) for SLALOM
280 analysis, of which 8,137 loci (28%) were predicted suspicious (**Supplementary Table 6**).

281
282 To validate SLALOM performance in real data, we restricted our analysis to 6,065 loci that have
283 maximum PIP > 0.1 and that contain nonsynonymous coding variants (predicted loss-of-function
284 [pLoF] and missense) in LD with the lead variant ($r^2 > 0.6$). Given prior evidence^{16,17,44} that such
285 nonsynonymous variants are highly enriched for being causal, we tested the validity of our method
286 by whether they achieve the highest PIP in the locus (*i.e.*, successful fine-mapping) in suspicious
287 vs. non-suspicious loci (**Methods**). While 40% (1,557 / 3,860) of non-suspicious loci successfully
288 fine-mapped nonsynonymous variants, only 17% (384 / 2,205) of suspicious loci did,
289 demonstrating a significant depletion (2.3x) of successfully fine-mapped nonsynonymous variants
290 in suspicious loci (Fisher's exact $P = 3.6 \times 10^{-79}$; **Fig. 4a**). We also tested whether
291 nonsynonymous variants belonged to 95% and 99% CS and again observed significant depletion
292 (1.4x and 1.3x, respectively; Fisher's exact $P < 4.6 \times 10^{-100}$). In addition, when we used a more
293 stringent r^2 threshold (> 0.8) for selecting loci that contain nonsynonymous variants, we also
294 confirmed significant enrichment (Fisher's exact $P < 6.1 \times 10^{-65}$; **Supplementary Figure 6**). To
295 quantify potential fine-mapping miscalibration in the GWAS Catalog, we investigated the
296 difference between mean PIP for lead variants and fraction of lead variants that are
297 nonsynonymous; assuming that nonsynonymous variants in these loci are truly causal, this
298 difference equals the difference between the true and reported fraction of lead PIP variants that
299 are causal. We observed differences between 26–51% and 10–18% under different PIP
300 thresholds in suspicious and non-suspicious loci, respectively (**Fig. 4b**), marking 45% and 15%
301 for high-PIP (> 0.9) variants.

302
303 We further assessed SLALOM performance in the GWAS Catalog meta-analyses by leveraging
304 high-PIP (> 0.9) complex trait and *cis*-eQTL variants that were rigorously fine-mapped^{16,17} in large-
305 scale biobanks (Biobank Japan [BBJ]⁵⁸, FinnGen²⁰, and UK Biobank [UKBB]¹⁹) and eQTL
306 resources (GTEx⁵⁹ v8 and eQTL Catalogue⁶⁰). Among the 27,713 loci analyzed by SLALOM
307 (maximum PIP > 0.1) that contain a lead variant that was included in biobank fine-mapping, 17%
308 (3,266 / 19,692) of the non-suspicious loci successfully fine-mapped one of the high-PIP GWAS
309 variants in biobank fine-mapping, whereas 7% (589 / 8,021) of suspicious loci did, showing a
310 significant depletion (2.3x) of the high-PIP complex trait variants in suspicious loci (Fisher's exact
311 $P = 4.6 \times 10^{-100}$; **Fig. 4c**). Similarly, among 26,901 loci analyzed by SLALOM that contain a lead
312 variant that was included in *cis*-eQTL fine-mapping, we found a significant depletion (1.9x) of the
313 high-PIP *cis*-eQTL variants in suspicious loci, where 7% (1,247 / 18,976) of non-suspicious loci
314 vs. 4% (281 / 7,925) of suspicious loci successfully fine-mapped one of the high-PIP *cis*-eQTL
315 variants (Fisher's exact $P = 2.6 \times 10^{-24}$; **Fig. 4d**). We observed the same significant depletions of
316 the high-PIP complex trait and *cis*-eQTL variants in suspicious loci that belonged to 95% and 99%
317 CS set (**Fig. 4c,d**).

318
319

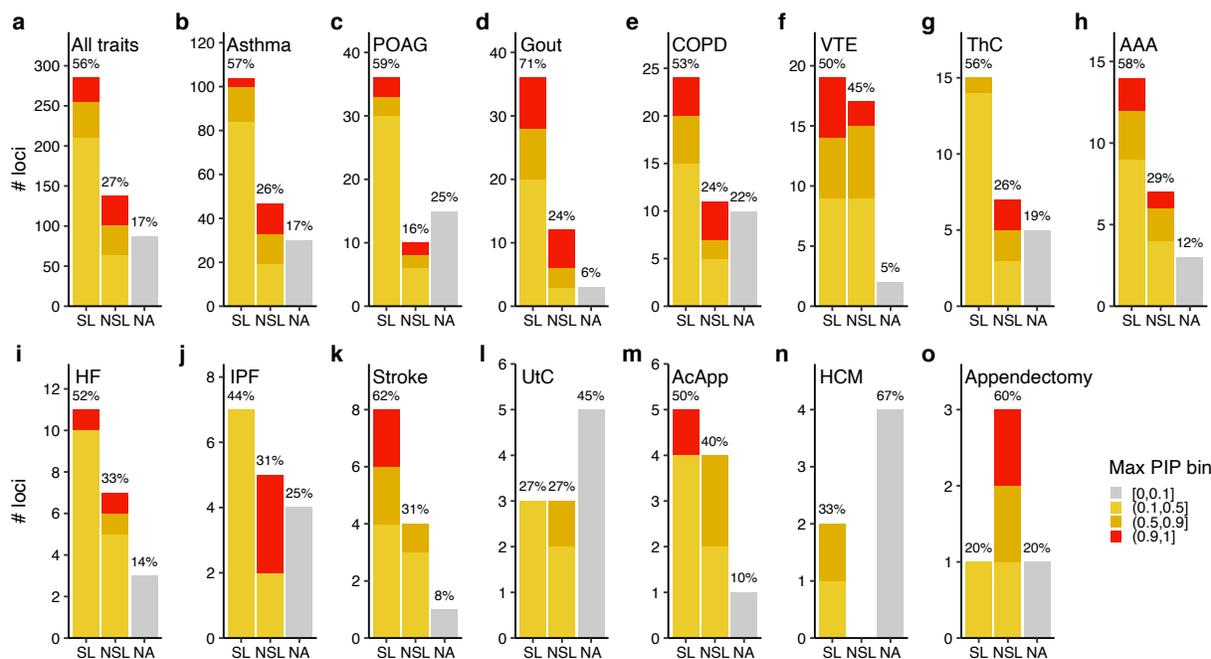


320
 321 **Fig. 4 | Evaluation of SLALOM performance in the GWAS Catalog summary statistics.** a,c,d. Depletion of likely
 322 causal variants in predicted suspicious loci. We evaluated whether (a) nonsynonymous coding variants (pLoF and
 323 missense), (c) high-PIP (> 0.9) complex trait variants in biobank fine-mapping, and (d) high-PIP (> 0.9) *cis*-eQTL
 324 variants in GTEx v8 and eQTL Catalogue were lead PIP variants, in 95% CS, or in 99% CS in suspicious vs. non-
 325 suspicious loci. Depletion was calculated by relative risk (*i.e.* a ratio of proportions; **Methods**). Error bars correspond
 326 to 95% confidence intervals using bootstrapping. Significance represents a Fisher's exact test *P*-value (*, *P* < 0.05; **,
 327 < 0.01; ***, < 0.001; ****, < 10⁻⁴). b. Plot of the estimated difference between true and reported proportion of causal
 328 variants in the loci tagging nonsynonymous variants (*r*² > 0.6 with the lead variants) in the GWAS Catalog under different
 329 PIP thresholds. Analogous to **Fig. 3b**, assuming nonsynonymous variants in these loci are truly causal, the mean PIP
 330 for lead variants minus the fraction of lead variants that are nonsynonymous above the threshold is equal to the
 331 difference between true and reported proportion of causal variants..
 332

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

334 Next, we applied SLALOM to meta-analysis summary statistics of 14 disease endpoints from the
 335 GBMI¹⁰. These summary statistics were generated from a meta-analysis of 2.1 million individuals
 336 in total across 19 biobanks, representing six different genetic ancestry groups of approximately
 337 33,000 African, 18,000 Admixed American, 31,000 Central and South Asian, 341,000 East Asian,
 338 1.8 million European, and 1,600 Middle Eastern individuals (**Supplementary Table 7**). Among
 339 509 genome-wide significant loci across the 14 traits, we found that 87 loci (17%) showed
 340 maximum PIP < 0.1, thus not being further considered by SLALOM. Of the remaining 422 loci
 341 with maximum PIP > 0.1, SLALOM identified that 285 loci (68%) were suspicious loci for fine-
 342 mapping (**Fig. 5a**; **Supplementary Table 8**). The fraction of suspicious loci and their maximum
 343 PIP varied by trait, reflecting different levels of statistical power (e.g., sample sizes, heritability,
 344 and local LD structure) as well as inter-cohort heterogeneity (**Fig. 5b–o**).

346 While the fraction of suspicious loci (68%) in the GBMI meta-analyses is higher than in the GWAS
 347 Catalog (28%), there might be multiple reasons for this discrepancy, including association
 348 significance, sample size, ancestral diversity, and study-specific QC criteria. For example, the
 349 GBMI summary statistics were generated from multi-ancestry, large-scale meta-analyses of
 350 median sample size of 1.4 million individuals across six ancestries, while 63% of the 467 summary
 351 statistics from the GWAS Catalog were only in European-ancestry studies and 83% had less than
 352 0.5 million discovery samples. Nonetheless, predicted suspicious loci for fine-mapping were
 353 prevalent in both the GWAS Catalog and the GBMI.



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

363
364

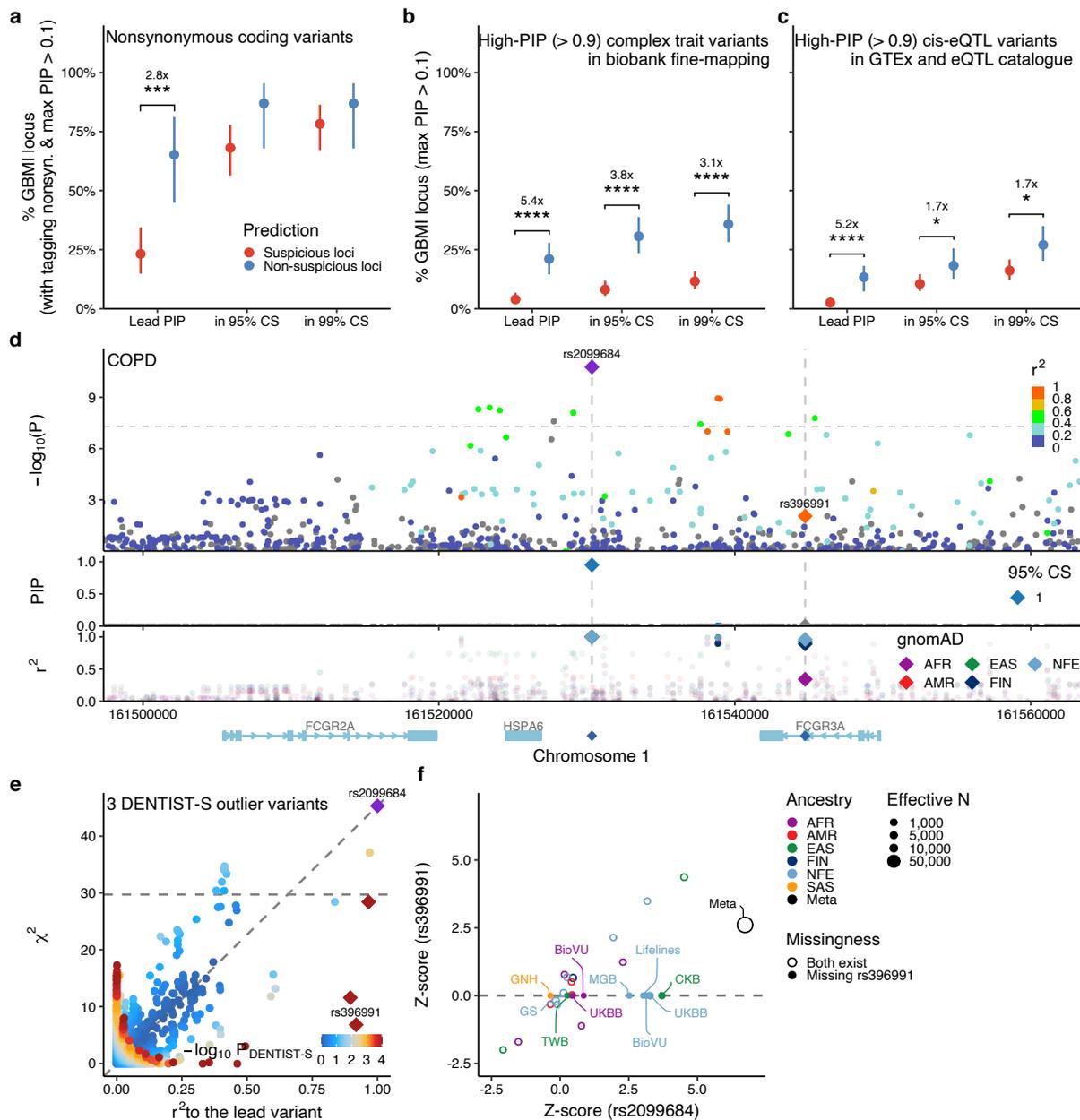
365 Using nonsynonymous (pLoF and missense) and high-PIP (> 0.9) complex trait and *cis*-eQTL
366 variants, we recapitulated a significant depletion of these likely causal variants in predicted
367 suspicious loci (2.8x, 5.4x, and 5.2x for nonsynonymous, high-PIP complex trait, and high-PIP
368 *cis*-eQTL variants being a lead PIP variant, respectively; Fisher's exact $P < 6.2 \times 10^{-4}$), confirming
369 our observation in the GWAS Catalog analysis (**Fig. 6a–c**).

370
371 In 15/23 non-suspicious loci harboring a nonsynonymous variant, the nonsynonymous variant
372 had the highest PIP. These included known missense variants such as rs116483731 (p.Arg20Gln)
373 in *SPDL1* for idiopathic pulmonary fibrosis (IPF)^{61,62} and rs28929474 (p.Glu366Lys) in *SERPINA1*
374 for chronic obstructive pulmonary disease (COPD)^{63,64}. In addition, we observed successful fine-
375 mapping in 2 novel loci for asthma, i) rs41286560 (p.Pro558Thr) in *RTL1*, a missense variant
376 known for decreasing height^{65,66} and ii) rs34187696 (p.Gly337Val) in *ZSCAN5A*, a known
377 missense variant for increasing monocyte count³⁰.

378
379 To characterize fine-mapping failures in suspicious loci, we examined suspicious loci in which a
380 nonsynonymous variant did not achieve the highest PIP. For example, the *FCGR2A/FCGR3A*
381 (1q23.3) locus for COPD contained a genome-wide significant lead intergenic variant rs2099684
382 ($P = 1.7 \times 10^{-11}$) which is in LD ($r^2 = 0.92$) with a missense variant rs396991 (p.Phe176Val) of
383 *FCGR3A* (**Fig. 6d**). This locus was not previously reported for COPD, but is known for
384 associations with autoimmune diseases (e.g., inflammatory bowel disease⁴⁴, rheumatoid
385 arthritis⁷, and systemic lupus erythematosus⁶⁷) and encodes the low-affinity human FC-gamma
386 receptors that bind to the Fc region of IgG and activate immune responses⁶⁸. Notably, this locus
387 contains copy number variations that contribute to the disease associations in addition to single-
388 nucleotide variants, which makes genotyping challenging^{68,69}. Despite strong LD with the lead
389 variant, rs396991 did not achieve genome-wide significance ($P = 9.1 \times 10^{-3}$), showing a significant
390 deviation from the expected association ($P_{\text{DENTIST-S}} = 5.3 \times 10^{-41}$; **Fig. 6e**). This is primarily due to
391 missingness of rs396991 in 8 biobanks out of 17 ($N_{\text{eff}} = 76,790$ and 36,781 for rs2099684 and
392 rs396991, respectively; **Fig. 6f**), which is caused by its absence from major imputation reference
393 panels (e.g., 1000GP⁴⁹, HRC⁵⁰, and UK10K⁷⁰) despite having a high MAF in every population
394 (MAF = 0.24–0.34 in African, admixed American, East Asian, European, and South Asian
395 populations of gnomAD⁵⁷).

396
397 Sample size imbalance across variants was pervasive in the GBMI meta-analyses⁷¹, and was
398 especially enriched in predicted suspicious loci—84% of suspicious loci vs. 24% of non-
399 suspicious loci showed a maximum/minimum effective sample size ratio > 2 among variants in
400 LD ($r^2 > 0.6$) with lead variants (a median ratio = 4.2 and 1.2 in suspicious and non-suspicious
401 loci, respectively; **Supplementary Fig. 7**). These observations are consistent with our
402 simulations, recapitulating that sample size imbalance results in miscalibration for meta-analysis
403 fine-mapping. Notably, we observed a similar issue in other GBMI downstream analyses (e.g.,
404 polygenic risk score [PRS]⁷¹ and drug discovery⁷²), where predictive performance improved
405 significantly after filtering out variants with maximum $N_{\text{eff}} < 50\%$. Although fine-mapping methods
406 cannot simply take this approach because it inevitably reduces calibration and recall by removing
407 true causal variants, other meta-analysis downstream analyses that primarily rely on polygenic
408 signals rather than individual variants should consider this filtering as an extra QC step.

409
410



411
 412 **Fig. 6 | Evaluation of SLALOM performance in the GBMI summary statistics.** a–c. Similar to Fig. 4, we evaluated
 413 whether (a) nonsynonymous coding variants (pLoF and missense), (b) high-PIP (> 0.9) complex trait variants in biobank
 414 fine-mapping, and (c) high-PIP (> 0.9) cis-eQTL variants in GTEx v8 and eQTL Catalogue were lead PIP variants, in
 415 95% CS, or in 99% CS in suspicious vs. non-suspicious loci. Depletion was calculated by relative risk (i.e. a ratio of
 416 proportions; **Methods**). Error bars correspond to 95% confidence intervals using bootstrapping. Significance represents
 417 a Fisher's 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
 418 COPD. The top panel shows a Manhattan plot, where the lead variant rs2099684 (purple diamond) and a missense
 419 variant rs396991 (orange diamond) are highlighted. Color represents r^2 values to the lead variant. Horizontal line
 420 represents a genome-wide significance threshold ($P = 5.0 \times 10^{-8}$). The middle panel shows PIP from ABF fine-mapping.
 421 Color represents whether variants belong to a 95% CS. The bottom panel shows r^2 values with the lead variant in
 422 gnomAD populations. e. A diagnosis plot showing r^2 values to the lead variant vs. marginal χ^2 . Color represents $-\log_{10}$
 423 $P_{\text{DENTIST-S}}$ values. Outlier marginal variants with $P_{\text{DENTIST-S}} < 10^{-4}$ are depicted in red with a diamond shape. Diagonal line
 424 represents an expected marginal association. Horizontal line represents a genome-wide significance threshold. f. Z-
 425 scores of the lead variant (rs2099684) vs. the missense variant (rs396991) in the constituent cohorts of the meta-
 426 analysis. Open and closed circles represent whether both variants exist in a cohort or rs396991 is missing. Circle size
 427 corresponds to an effective sample size. Color represents genetic ancestry.

428 **Comparison of fine-mapping results between the GBMI meta-analyses and individual** 429 **biobanks**

430 Motivated by successful validation of SLALOM performance, we investigated whether fine-
431 mapping confidence and resolution were improved in the GBMI meta-analyses over individual
432 biobanks. To this end, we used our fine-mapping results^{16,17} of nine disease endpoints (asthma⁶⁴,
433 COPD⁶⁴, gout, heart failure⁷³, IPF⁶², primary open angle glaucoma⁷⁴, thyroid cancer, stroke⁷⁵, and
434 venous thromboembolism⁷⁶) in BBJ⁵⁸, FinnGen²⁰, and UKBB¹⁹ Europeans that also contributed to
435 the GBMI meta-analyses for the same traits.

436
437 To perform an unbiased comparison of PIP between the GBMI meta-analysis and individual
438 biobanks, we investigated functional enrichment of fine-mapped variants based on top PIP
439 rankings in the GBMI and individual biobanks (top 0.5%, 0.1%, and 0.05% PIP variants in the
440 GBMI vs. maximum PIP across BBJ, FinnGen, and UKBB; **Methods**). Previous studies have
441 shown that high-PIP (> 0.9) complex trait variants are significantly enriched for well-known
442 functional categories, such as coding (pLoF, missense, and synonymou), 5'/3' UTR, promoter,
443 and *cis*-regulatory element (CRE) regions (DNase I hypersensitive sites [DHS] and H3K27ac)^{16,17}.
444 Using these functional categories, we found no significant enrichment of variants in the top PIP
445 rankings in the GBMI over individual biobanks (Fisher's exact $P > 0.05$; **Fig. 7a**) except for variants
446 in the promoter region (1.8x; Fisher's exact $P = 3.1 \times 10^{-4}$ for the top 0.1% PIP variants). We
447 observed similar trends regardless of whether variants were in suspicious or non-suspicious loci
448 (**Fig. 7b,c**). To examine patterns of increased and decreased PIP for individual variants, we also
449 calculated PIP difference between the GBMI and individual biobanks, defined as $\Delta\text{PIP} = \text{PIP}$
450 (GBMI) – maximum PIP across BBJ, FinnGen, and UKBB (**Supplementary Fig. 8,9**). We
451 investigated functional enrichment based on ΔPIP bins and observed inconsistent enrichment
452 results using different ΔPIP thresholds (**Supplementary Fig. 10**). Finally, to test whether fine-
453 mapping resolution was improved in the GBMI over individual biobanks, we compared the size of
454 95% CS after restricting them to cases where a GBMI CS overlapped with an individual biobank
455 CS from BBJ, FinnGen, or UKBB (**Methods**). We observed the median 95% CS size of 2.5 and
456 2.5 in non-suspicious loci for the GBMI and individual biobanks, respectively, and 5 and 15 in
457 suspicious loci, respectively (**Supplementary Fig. 11**). The smaller credible set size in suspicious
458 loci in GBMI could be due to improved resolution or to increased miscalibration. These results
459 provide limited evidence of overall fine-mapping improvement in the GBMI meta-analyses over
460 what is achievable by taking the best result from individual biobanks.

461
462 Individual examples, however, provide insights into the types of fine-mapping differences that can
463 occur. To characterize the observed differences in fine-mapping confidence and resolution, we
464 further examined non-suspicious loci with $\Delta\text{PIP} > 0.5$ in asthma. In some cases, the increased
465 power and/or ancestral diversity of GBMI led to improved fine-mapping: for example, an intergenic
466 variant rs1888909 (~18 kb upstream of *IL33*) showed $\Delta\text{PIP} = 0.99$ (PIP = 1.0 and 0.008 in GBMI
467 and FinnGen, respectively; **Fig. 7d**), which was primarily owing to increased association
468 significance in a meta-analysis ($P = 3.0 \times 10^{-86}$, 7.4×10^{-2} , 3.6×10^{-16} , and 1.9×10^{-53} in GBMI,
469 BBJ, FinnGen, and UKBB Europeans, respectively) as well as a shorter LD length in the African
470 population than in the European population (LD length = 4 kb vs. 41 kb for variants with $r^2 > 0.6$
471 with rs1888909 in the African and European populations, respectively; $N_{\text{eff}} = 4,270$ for Africans in
472 the GBMI asthma meta-analysis; **Supplementary Fig. 12**). This variant was also fine-mapped for
473 eosinophil count in UKBB Europeans (PIP = 1.0; $P = 1.3 \times 10^{-314}$)¹⁶ and was previously reported
474 to regulate *IL33* gene expression in human airway epithelial cells via allele-specific transcription
475 factor binding of OCT-1 (POU2F1)⁷⁷. Likewise, we observed a missense variant rs16903574
476 (p.Phe319Leu) in *OTULINL* showed $\Delta\text{PIP} = 0.79$ (PIP = 1.0 and 0.21 in GBMI and UKBB

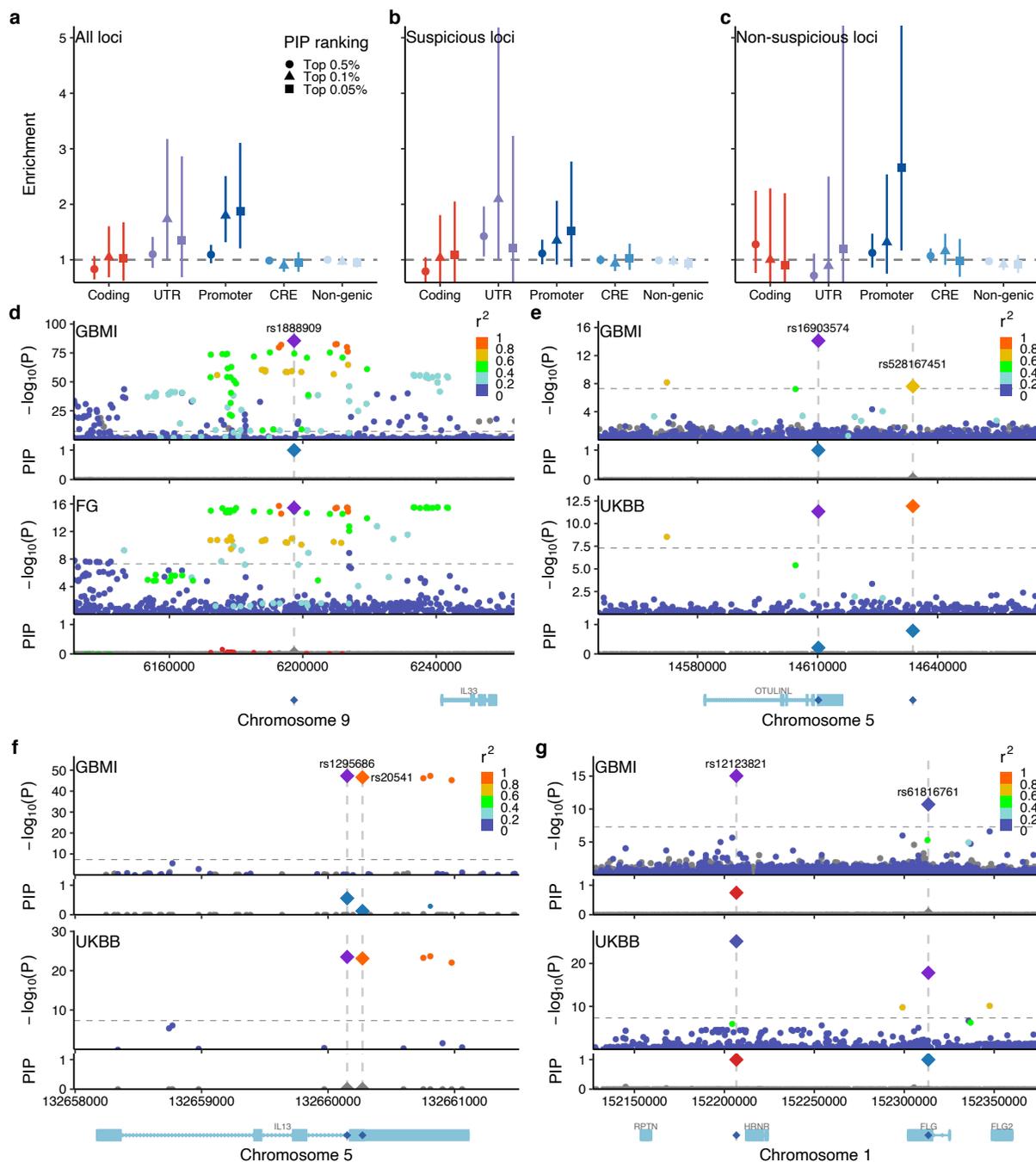
477 Europeans, respectively; **Fig. 7e**) owing to improved association significance ($P = 7.7 \times 10^{-15}$ and
478 4.7×10^{-12} in GBMI and UKBB Europeans, respectively).

479

480 However, we also observed very high Δ PIP for variants that are not likely causal. For example,
481 we observed that an intronic variant rs1295686 in *IL13* showed Δ PIP = 0.56 (PIP = 0.56 and
482 0.0002 in GBMI and UKBB Europeans, respectively; **Fig. 7f**), despite having strong LD with a
483 nearby missense variant rs20541 (p.Gln144Arg; $r^2 = 0.96$ with rs1295686) which only showed
484 Δ PIP = 0.13 (PIP = 0.13 and 0.0001 in GBMI and UKBB Europeans, respectively). The missense
485 variant rs20541 showed PIP = 0.23 and 0.15 for a related allergic disease, atopic dermatitis, in
486 BBJ and FinnGen, respectively¹⁷, and was previously shown to induce STAT6 phosphorylation
487 and up-regulate CD23 expression in monocytes, promoting IgE synthesis⁷⁸. Although the GBMI
488 meta-analysis contributed to prioritizing these two variants (sum of PIP = 0.69 vs. 0.0003 in GBMI
489 and UKBB Europeans, respectively), the observed Δ PIP was higher for rs1295686 than for
490 rs20541.

491

492 While increasing sample size in meta-analysis improves association significance, we also found
493 negative Δ PIP due to losing the ability to model multiple causal variants. A stop-gained variant
494 rs61816761 (p.Arg501Ter) in *FLG* showed Δ PIP = -1.0 (PIP = 6.4×10^{-5} and 1.0 in GBMI and
495 UKBB Europeans, respectively; **Fig. 7g**), which was primarily owing to a nearby lead variant
496 rs12123821 (~17 kb downstream of *HRNR*; $r^2 = 0.0$ with rs61816761). This lead variant
497 rs12123821 showed greater significance than rs61816761 in GBMI ($P = 9.3 \times 10^{-16}$ and $2.0 \times 10^{-$
498 11 for rs12123821 and rs61816761, respectively) as well as in UKBB Europeans ($P = 7.1 \times 10^{-26}$
499 and 1.5×10^{-18}). While our biobank fine-mapping^{16,17} assigned PIP = 1.0 for both variants based
500 on multiple causal variant fine-mapping (*i.e.*, FINEMAP²¹ and SuSiE²³), our ABF fine-mapping in
501 the GBMI meta-analysis was only able to assign PIP = 0.74 for the lead variant rs12123821 due
502 to a single causal variant assumption. This recapitulates the importance of multiple causal variant
503 fine-mapping in complex trait fine-mapping^{16,17}—however, we note that multiple causal variant
504 fine-mapping with an external LD reference is extremely error-prone as previously reported^{14–16}.



505
 506 **Fig. 7 | Fine-mapping improvement and retrogression in the GBMI meta-analyses over individual biobanks.** a–
 507 c. Functional enrichment of variants in each functional category based on top PIP rankings in the GBMI and individual
 508 biobanks (maximum PIP of BBJ, FinnGen, and UKBB). Shape corresponds to top PIP ranking (top 0.5%, 0.1%, and
 509 0.05%). Enrichment was calculated by a relative risk (*i.e.* a ratio of proportions; **Methods**). Error bars correspond to
 510 95% confidence intervals using bootstrapping. **d–g**. Locuszoom plots for the same non-suspicious locus of asthma in
 511 the GBMI meta-analysis and an individual biobank (BBJ, FinnGen, or UKBB Europeans) that showed the highest PIP
 512 in our biobank fine-mapping. Colors in the Manhattan panels represent r^2 values to the lead variant. In the PIP panels,
 513 only fine-mapped variants in the 95% CS are colored, where the same colors are applied between the GBMI meta-
 514 analysis and an individual biobank based on merged CS as previously described. Horizontal line represents a genome-
 515 wide significance threshold ($P = 5.0 \times 10^{-8}$). **d**. rs1888909 for asthma in the GBMI and FinnGen. **e**. rs16903574 for
 516 asthma in the GBMI and UKBB Europeans. Nearby rs528167451 was also highlighted, which was in strong LD ($r^2 =$
 517 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

518 GBMI and UKBB Europeans. A nearby missense, rs20541, showed lower PIP than rs1295686 despite having strong
519 LD ($r^2 = 0.96$). **g.** rs12123821 for asthma in the GBMI and UKBB Europeans. Nearby stop-gained rs61816761 was
520 independent of rs12123821 ($r^2 = 0.0$) and not fine-mapped in the GBMI due to a single causal variant assumption in
521 the ABF fine-mapping.

522 Discussion

523 In this study, we first demonstrated in simulations that meta-analysis fine-mapping is substantially
524 miscalibrated when constituent cohorts are heterogeneous in phenotyping and imputation. To
525 mitigate this issue, we developed SLALOM, a summary statistics-based QC method for identifying
526 suspicious loci in meta-analysis fine-mapping. Applying SLALOM to 14 disease endpoints from
527 the GBMI meta-analyses¹⁰ as well as 467 summary statistics from the GWAS Catalog⁴⁸, we
528 observed widespread suspicious loci in meta-analysis summary statistics, suggesting that meta-
529 analysis fine-mapping is often miscalibrated in real data too. Indeed, we demonstrated that the
530 predicted suspicious loci were significantly depleted for having likely causal variants as a lead PIP
531 variant, such as nonsynonymous variants, high-PIP (> 0.9) GWAS and cis-eQTL fine-mapped
532 variants from our previous fine-mapping studies^{16,17}. Our method provides better calibration in
533 non-suspicious loci for meta-analysis fine-mapping, generating a more reliable set of variants for
534 further functional characterization.

535
536 We have found limited evidence of improved fine-mapping in the GBMI meta-analyses over
537 individual biobanks. A few empirical examples in this study as well as other previous
538 studies^{7,9,26,27,30} suggested that multi-ancestry, large-scale meta-analysis could have potential to
539 improve fine-mapping confidence and resolution owing to increased statistical power in
540 associations and differential LD pattern across ancestries. However, we have highlighted that the
541 observed improvement in PIP could be due to sample size imbalance in a locus, miscalibration,
542 and technical confoundings too, which further emphasizes the importance of careful investigation
543 of fine-mapped variants identified through meta-analysis fine-mapping.

544
545 As high-confidence fine-mapping results in large-scale biobanks and molecular QTLs continue to
546 become available^{16,17,60}, we propose alternative approaches for prioritizing candidate causal
547 variants in a meta-analysis. First, these high-confidence fine-mapped variants have been a
548 valuable resource to conduct a “PheWAS”¹⁶ to match with associated variants in a meta-analysis,
549 which provides a narrower list of candidate variants assuming they would equally be functional
550 and causal in related complex traits or tissues/cell-types. Second, a traditional approach based
551 on tagging variants (e.g., $r^2 > 0.6$ with lead variants, or PICS⁷⁹ fine-mapping approach that only
552 relies on a lead variant and LD) can be still highly effective, especially for known functional
553 variants such as nonsynonymous coding variants. As we highlighted in this and previous³⁹
554 studies, potentially causal variants in strong LD with lead variants might not achieve genome-
555 wide significance because of missingness and heterogeneity.

556
557 While using an external LD reference for fine-mapping has been shown to be extremely error-
558 prone^{14–16}, we find here that it can be useful for flagging suspicious loci, even when it does not
559 perfectly represent the in-sample LD structure of the meta-analyzed individuals. However, our
560 use of external LD reference comes with several limitations. For example, due to the finite sample
561 size of external LD reference, rare or low-frequency variants have larger uncertainties around r^2
562 than common variants. Moreover, our r^2 values in a multi-ancestry meta-analysis are currently
563 approximated based on a sample-size-weighted average of r^2 across ancestries as previously
564 suggested⁸⁰, but this can be different from actual r^2 . These uncertainties around r^2 affect SLALOM
565 prediction performance and should be modeled appropriately for further method development. On
566 the other hand, we find it challenging to use a LD reference when true causal variants are located

567 within a complex region (e.g., major histocompatibility complex [MHC]), or are entirely missing
568 from standard LD or imputation reference panels, especially for structural variants. These
569 limitations are not specific to meta-analysis fine-mapping, and separate fine-mapping methods
570 based on bespoke imputation references have been developed (e.g., HLA⁸¹, KIR⁸², and variable
571 numbers of tandem repeats [VNTR]⁸³).

572
573 In addition, there are several methodological limitations of SLALOM. First, our simulations only
574 include one causal variant per locus. Although additional independent causal variants would not
575 affect SLALOM precision (but decrease recall), multiple *correlated* causal variants in a locus
576 would violate SLALOM assumptions and could lead to some DENTIST-S outliers that are not due
577 to heterogeneity or missingness but rather simply a product of tagging multiple causal variants in
578 LD. In fact, our previous studies have illustrated infrequent but non-zero presence of such
579 correlated causal variants in complex traits^{16,17}. Second, SLALOM prediction is not perfect.
580 Although fine-mapping calibration is certainly better in non-suspicious vs. suspicious loci,
581 SLALOM has low precision, and we still observe some miscalibration in non-suspicious loci.
582 Finally, SLALOM is a per-locus QC method and does not calibrate per-variant PIPs. Further
583 methodological development that properly models heterogeneity, missingness, multiple causal
584 variants, and LD uncertainty across multiple cohorts and ancestries is needed to refine per-variant
585 calibration and recall in meta-analysis fine-mapping.

586
587 We have found evidence in our simulations and real data of severe miscalibration of fine-mapping
588 results from GWAS meta-analysis; for example, we estimate that the difference between true and
589 reported proportion of causal variants is 20% and 45% for high-PIP (> 0.9) variants in suspicious
590 loci from the simulations and the GWAS Catalog, respectively. Our SLALOM method helps to
591 exclude spurious results from meta-analysis fine-mapping; however, even fine-mapping results in
592 SLALOM-predicted “non-suspicious” loci remain somewhat miscalibrated, showing estimated
593 differences between true and reported proportion of causal variants of 4% and 15% for high-PIP
594 variants in the simulations and the GWAS Catalog, respectively. We thus urge extreme caution
595 when interpreting PIPs computed from meta-analyses until improved methods are available. We
596 recommend that researchers looking to identify likely causal variants employ complete
597 synchronization of study design, case/control ascertainment, genomic profiling, and analytical
598 pipeline, or rely more heavily on functional annotations, biobank fine-mapping, or molecular QTLs.
599

600 **Data availability**

601 The GBMI summary statistics for the 14 endpoints are available at
602 <https://www.globalbiobankmeta.org/resources> and are browserble at the GBMI PheWeb⁸⁴
603 website (<http://results.globalbiobankmeta.org/>).

604 **Code availability**

605 The SLALOM software is available at <https://github.com/mkanai/slalom>. Custom scripts to
606 perform all the analyses and generate all the figures are available at
607 <https://github.com/mkanai/slalom-paper>.

608 **Acknowledgements**

609 We acknowledge all the participants of the 19 biobanks that have contributed to the GBMI.
610 Biobank-specific acknowledgements are included in the **Supplementary Note**. We thank H.
611 Huang, A.R. Martin, B.M. Neale, Y. Okada, K. Tsoo, J.C. Ulirsch, Y. Wang, and all the members
612 of Finucane and Daly labs for their helpful feedback. M.K. was supported by a Nakajima
613 Foundation Fellowship and the Masason Foundation. H.K.F. was funded by NIH grant DP5
614 OD024582.

615 **Competing interests**

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

617 **Methods**

618 **Meta-analysis fine-mapping simulation**

619 To benchmark fine-mapping performance in meta-analysis, we simulated a large-scale, realistic
620 GWAS meta-analysis and performed fine-mapping under different scenarios. An overview of our
621 simulation pipeline is summarized in **Supplementary Fig. 2**.

622 *Simulated true genotype*

623 Using HAPGEN2⁸⁵ with the 1000 Genomes Project Phase 3 reference⁴⁹, we simulated “true”
624 genotypes of chromosome 3 for multiple independent cohorts from African, East Asian, and
625 European ancestries. For each independent cohort from a given ancestry, we simulated 10,000
626 individuals each using the default parameters, with an ancestry-specific effective population size
627 set to 17,469, 14,269, and 11,418 for Africans, East Asians, and Europeans, respectively, as
628 recommended⁸⁵. To mimic sample size imbalance of different ancestries in the current meta-
629 analyses, we simulated 10 independent European cohorts, 1 African cohort, and 1 East Asian
630 cohort.

631
632 To restrict our analysis to unrelated samples, we computed sample relatedness based on KING
633 kinship coefficients⁸⁶ using PLINK 2.0 (ref. ⁸⁷) and removed monozygotic twins, duplicated
634 individuals, or first-degree relatives with the coefficient threshold of 0.177. The detailed sample
635 sizes of unrelated individuals for each cohort is summarized in **Supplementary Table 1**.

636 *Genotyping and imputation*

637 To simulate realistic genotyping and imputation procedures, we first virtually genotyped each
638 cohort by restricting variants to those that are available on different genotyping arrays. We
639 selected three major genotyping arrays from Illumina, Inc. (Omni2.5, Multi-Ethnic Global Array
640 [MEGA], and Global Screening Array [GSA]) that have different densities of genotyping probes
641 (**Supplementary Table 2**). For each cohort, we created three virtually genotyped datasets by
642 retaining variants that are genotyped on each array. For the sake of simplicity, we assumed no
643 genotyping errors occurred between true genotypes and virtually genotyped data—however, in
644 practice, genotyping error is one of the major sources of unexpected confounding (e.g., see recent
645 discussions here^{88,89}) and should be treated carefully.

646
647 For each pair of cohort and genotyping array, we then imputed missing variants using different
648 imputation reference panels. We used the Michigan Imputation Server
649 (<https://imputationserver.sph.umich.edu/>)⁹⁰ and the TOPMed Imputation Server
650 (<https://imputation.biodatacatalyst.nhlbi.nih.gov/>)⁵¹ with the default parameters, using three
651 publicly available reference panels: the 1000 Genomes Project Phase 3 (version 5; $n = 2,504$;
652 1000GP3)⁴⁹, the Haplotype Reference Consortium (version r1.1; $n = 32,470$; HRC)⁵⁰, and the
653 TOPMed (version R2; $n = 97,256$)⁵¹. Briefly, for each input, the imputation server created chunks
654 of 20 Mb, applied the standard QC, pre-phased each chunk with Eagle2 (ref. ⁹¹), and imputed
655 non-genotyped variants using a specified reference panel with Minimac4
656 (<https://genome.sph.umich.edu/wiki/Minimac4>). The detailed documentation of the imputation
657 pipeline is available on the Michigan and TOPMed websites and has been described elsewhere⁹⁰.

658
659 We applied post-imputation QC by only keeping variants with $MAF > 0.001$ and imputation $Rsq >$
660 0.6 . Because the TOPMed panel is based on GRCh38 while the 1000GP3 and the HRC panels
661 are on GRCh37, we lifted over TOPMed variants from GRCh38 to GRCh37 to meta-analyze with
662 other cohorts. We excluded any variants which were lifted over to different chromosomes or for
663 which the conversion failed. The number of virtually genotyped and imputed variants for each
664 combination of cohort, genotyping array, and imputation panel is summarized in **Supplementary**
665 **Table 3**.

666 *True phenotype*

667 We simulated 300 true phenotypes that resemble observed complex trait genetic architecture and
668 phenotypic heterogeneity across cohorts. Based on previous literature, we set parameters as
669 follows: 1) 50% of 1 Mb loci contain a true causal variant⁹²; 2) probability of being causal is
670 proportional to functional enrichments of variant consequences (pLoF, missense, synonymous,
671 UTR5/3, promoter, cis-regulatory region, and non-genic) for fine-mapped variants as estimated in
672 a previous complex trait fine-mapping study¹⁷; 3) per-allele causal effect sizes have a variance
673 proportional to $[2p(1-p)]^\alpha$ where p represents a maximum MAF across the three ancestries
674 (AFR, EAS, and EUR) and α is set to be -0.38 (ref. ⁵²); and 4) total SNP-heritability h_g^2 for
675 chromosome 3 equals 0.03 (ref. ⁵³). For the sake of simplicity, we randomly draw a single true
676 causal variant per locus because ABF assumes a single causal variant^{31,32}. We draw true causal
677 variants from 1,150,893 non-ambiguous single-nucleotide variants in 1000GP3 that showed MAF
678 > 0.01 in at least one of the three ancestries (AFR, EAS, or EUR) and were not located within
679 conversion-unstable positions (CUP)⁵⁴ between the human genome builds GRCh37 and
680 GRCh38. To mimic phenotypic heterogeneity across cohorts in real-world meta-analysis (due to
681 e.g., different ascertainment, measurement error, or true effect size heterogeneity), we introduced
682 cross-cohort genetic correlation of true effect sizes r_g which is set to be one of 1, 0.9, or 0.5. For
683 a true causal variant j , true causal effect sizes β_j across cohorts were randomly drawn from $\beta_j \sim$

684 $MVN(0, \Sigma)$ where diagonal elements of Σ were set to be $\sigma_g^2 \cdot [2p(1-p)]^\alpha$ and off-diagonal
685 elements of Σ were set to be $r_g \cdot \sigma_g^2 \cdot [2p(1-p)]^\alpha$. σ_g^2 was determined by $\sigma_g^2 = h_g^2 /$
686 $\sum_j [2p(1-p)]^{1+\alpha}$. For each cohort, true phenotype y was computed via $y = X\beta + \varepsilon$ where X
687 is the above true genotype matrix from HAPGEN2 and $\varepsilon_i \sim N(0, 1 - \sigma_g^2)$ i.i.d. We simulated 100
688 true phenotypes for each of $r_g = 1, 0.9, \text{ and } 0.5$, respectively.

689 GWAS

690 For each combination of phenotype, cohort, genotyping chip, and imputation panel, we conducted
691 GWAS via a standard linear regression as implemented in PLINK 2.0 using imputed dosages. For
692 covariates, we included top 10 principal components that were calculated based on true
693 genotypes after restricting to unrelated samples. We only used LD-pruned variants with MAF >
694 0.01 for PCA.

695 Meta-analysis

696 To simulate meta-analyses that resemble real-world settings, we generated multiple
697 configurations of the above GWAS results to meta-analyze across 10 independent cohorts.
698 Briefly, we chose configurations based on the following settings: 1) 10 EUR cohorts are genotyped
699 and imputed using the same genotyping array (one of GSA, MEGA, or Omni2.5) and the same
700 imputation panel (one of 1000GP3, HRC, TOPMed, or TOPMed-liftover); 2) 10 cohorts consisting
701 of multiple ancestries (9 EUR + 1 AFR/EAS cohorts or 8 EUR + 1 AFR + 1 EAS cohorts), with all
702 cohorts genotyped and imputed using the same array (Omni2.5) and the same panel (1000GP3);
703 3) 10 EUR or multi-ancestry cohorts are genotyped using the same array (Omni2.5) but imputed
704 using different panels across cohorts; 4) 10 EUR or multi-ancestry cohorts are imputed using the
705 same panel (1000GP3) but genotyped using different arrays across cohorts; 5) 10 EUR or multi-
706 ancestry cohorts are genotyped and imputed using different arrays and panels across cohorts.
707 For settings 3–5, we randomly draw a combination of a genotyping array and an imputation panel
708 for each cohort five times each for 10 EUR and multi-ancestry cohorts. In total, we generated 45
709 configurations as summarized in **Supplementary Table 4**.

710
711 For each configuration, we conducted a fixed-effect meta-analysis based on inverse-variance
712 weighted betas and standard errors using a modified version of PLINK 1.9
713 (https://github.com/mkanai/plink-ng/tree/add_se_meta).

714 Fine-mapping

715 For each meta-analysis, we defined fine-mapping regions based on a 1 Mb window around each
716 genome-wide significant lead variant and applied $ABF^{31,32}$ using prior effect size variance of $\sigma_0^2 =$
717 0.04. We set a prior variance of effect size to be 0.04 which was taken from Wakefield et al³¹ and
718 is commonly used in meta-analysis fine-mapping studies^{2,7}. We computed posterior inclusion
719 probability (PIP) and 95% credible set (CS) for each locus and evaluated whether true causal
720 variants were correctly fine-mapped.

721 The SLALOM method

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

724 *Locus definition*

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

729 *DENTIST-S outlier detection*

730 For each variant in a locus, we computed DENTIST-S statistics using equation (1) based on the
731 assumption of a single causal variant. DENTIST-S P-values ($P_{\text{DENTIST-S}}$) were computed using the
732 χ^2 distribution with 1 degree of freedom. We applied ABF^{31,32} using prior effect size variance of
733 $\sigma_0^2 = 0.04$ and used the lead PIP variant (the variant with the highest PIP) as an approximation of
734 the causal variant in the locus. To retrieve correlation r among the variants, we used publicly
735 available LD matrices from gnomAD⁵⁷ v2 as external LD reference for African, Admixed American,
736 East Asian, Finnish, and non-Finnish European populations. When multiple populations exist, we
737 computed a sample-size-weighted average of r^2 using per-variant sample sizes for each
738 population as previously suggested⁸⁰. We excluded variants without r^2 available in gnomAD from
739 the analysis. Since gnomAD v2 LD matrices are based on the human genome assembly GRCh37,
740 variants were lifted over to GRCh38 if the input summary statistics were based on GRCh38.

741
742 We determined DENTIST-S outlier variants using two thresholds: 1) $r^2 > \rho$ to the lead and 2)
743 $P_{\text{DENTIST-S}} < \tau$. The thresholds ρ and τ were set to $\rho = 0.6$ and $\tau = 1.0 \times 10^{-4}$ based on the training
744 in simulations as described below.

745 *Suspicious loci prediction*

746 We predicted whether a locus is suspicious or non-suspicious for fine-mapping based on the
747 number of DENTIST-S outlier variants in the locus $> \kappa$. To determine the best-performing
748 thresholds (ρ , τ , and κ), we used loci with maximum PIP > 0.9 in the simulations for training.
749 Positive conditions were defined as whether a true causal variant in a locus is 1) a lead PIP
750 variant, 2) in 95% CS, and 3) in 99% CS. We computed AUROC across different thresholds ($\rho =$
751 $0, 0.1, 0.2, \dots, 0.9$; $-\log_{10} \tau = 0, 0.5, 1, \dots, 10$; and $\kappa = 0, 1, 2, \dots$) and chose $\rho = 0.6$, $\tau = 1.0 \times 10^{-4}$,
752 and $\kappa = 0$ that showed the highest AUROC for all the aforementioned positive conditions. Using
753 all the loci in the simulations, we then evaluated fine-mapping miscalibration (defined as mean
754 PIP – fraction of true causal variants) at different PIP thresholds in suspicious and non-suspicious
755 loci and decided to only apply SLALOM to loci with maximum PIP > 0.1 owing to relatively lower
756 miscalibration and specificity of SLALOM at lower PIP thresholds.

757 **GWAS Catalog analysis**

758 We retrieved full GWAS summary statistics publicly available on the GWAS Catalog⁴⁸. Out of
759 33,052 studies from 5,553 publications registered at the GWAS Catalog (as of January 12, 2022),
760 we selected 467 studies from 96 publications that have 1) full harmonized summary statistics
761 preprocessed by the GWAS Catalog with non-missing variant ID, marginal beta, and standard
762 error columns, 2) a discovery sample size of more than 10,000 individuals, 3) African (including
763 African American, Afro-Caribbean, and Sub-Saharan African), admixed American (Hispanic and
764 Latin American), East Asian, or European samples based on their broad ancestral category
765 metadata, 4) at least one genome-wide significant association ($P < 5.0 \times 10^{-8}$), and 5) our manual
766 annotation as a meta-analysis rather than a single-cohort study (**Supplementary Table 5**). We
767 applied SLALOM to the 467 summary statistics and identified 35,864 genome-wide significant loci
768 (based on 1 Mb window around lead variants), of which 28,925 loci with maximum PIP > 0.1 were

769 further classified into suspicious and non-suspicious loci. Since per-variant sample sizes were not
770 available, we used overall sample sizes of each ancestry (African, Admixed American, East Asian,
771 and European) to calculate the weighted-average of r^2 . All the variants were harmonized into the
772 human genome assembly GRCh38 by the GWAS Catalog.

773 **GBMI analysis**

774 We used meta-analysis summary statistics of 14 disease endpoints from the GBMI
775 (**Supplementary Table 7**). These meta-analyses were conducted using up to 1.8 million
776 individuals across 19 biobanks, representing six different genetic ancestry groups (approximately
777 33,000 African, 18,000 Admixed American, 31,000 Central and South Asian, 341,000 East Asian,
778 1.8 million European, and 1,600 Middle Eastern individuals). Detailed procedures of the GBMI
779 meta-analyses were described in the GBMI flagship manuscript¹⁰.

780
781 Across the 14 summary statistics, we defined 503 genome-wide significant loci ($P < 5.0 \times 10^{-8}$)
782 based on a 1 Mb window around each lead variant and merged them if they overlapped. We
783 applied SLALOM to 422 loci with maximum PIP > 0.1 based on the ABF fine-mapping and
784 predicted whether they were suspicious or non-suspicious for fine-mapping. We used per-variant
785 sample sizes of each ancestry (African, Admixed American, East Asian, Finnish, and non-Finnish
786 European) to calculate the weighted-average of r^2 . Since gnomAD LD matrices were not available
787 for Central and South Asian and Middle Eastern, we did not use their sample sizes for the
788 calculation. All the variants were processed on the human genome assembly GRCh38.

789 **Fine-mapping results of complex traits and *cis*-eQTL**

790 We retrieved our previous fine-mapping results for 1) complex traits in large-scale biobanks
791 (BBJ⁵⁸, FinnGen²⁰, and UKBB¹⁹ Europeans)^{16,17} and 2) *cis*-eQTLs in GTEx⁵⁹ v8 and eQTL
792 Catalogue⁶⁰. Briefly, we conducted multiple-causal-variant fine-mapping (FINEMAP^{21,22} and
793 SuSiE²³) of complex trait GWAS (# unique traits = 148) and *cis*-eQTL gene expression (# unique
794 tissues/cell-types = 69) using summary statistics and in-sample LD. Detailed fine-mapping
795 methods are described elsewhere^{16,17}.

796
797 In this study, we collected 1) high-PIP GWAS variants that achieved PIP > 0.9 for any traits in any
798 biobank and 2) high-PIP *cis*-eQTL variants that achieved PIP > 0.9 for any gene expression in
799 any tissues/cell-types. All the variants were originally processed on the human genome assembly
800 GRCh37 and lifted over to the GRCh38 for comparison.

801 *Additional fine-mapping results*

802 To compare with the GBMI meta-analyses, we additionally conducted multi-causal-variant fine-
803 mapping of four additional endpoints (gout, heart failure, thyroid cancer, and venous
804 thromboembolism) that were not fine-mapped in our previous study^{16,17}. We used exactly the
805 same fine-mapping pipeline (FINEMAP^{21,22} and SuSiE²³) as described previously^{16,17}. For UKBB
806 Europeans, to use the exact same samples that contributed to the GBMI, we used individuals of
807 European ancestry ($n = 420,531$) as defined in the Pan-UKBB project
808 (<https://pan.ukbb.broadinstitute.org>), instead of those of “white British ancestry” ($n = 361,194$)
809 used in our previous study^{16,17}.

810 **Enrichment analysis of likely causal variants**

811 To validate SLALOM performance, we asked whether suspicious and non-suspicious loci were
812 enriched for having likely causal variants as a lead PIP variant, and for containing them in the
813 95% and 99% CS. We defined likely causal variants using 1) nonsynonymous coding variants,
814 *i.e.*, pLoF and missense variants annotated⁹³ by the Ensembl Variant Effect Predictor (VEP) v101
815 (using GRCh38 and GENCODE v35), 2) the high-PIP (> 0.9) complex trait fine-mapped variants,
816 and 3) the high-PIP (> 0.9) *cis*-eQTL fine-mapped variants from our previous studies as described
817 above.

818
819 We estimated enrichment for suspicious and non-suspicious loci as a relative risk (*i.e.*, a ratio of
820 proportion of variants) between being in suspicious/non-suspicious loci and having the annotated
821 likely causal variants as a lead PIP variant (or containing them in the 95% or 99% CS). That is, a
822 relative risk = (proportion of non-suspicious loci having the annotated variants as a lead PIP
823 variant) / (proportion of suspicious loci having the annotated variants as a lead PIP variant). We
824 computed 95% confidence intervals using bootstrapping.

825 **Comparison of fine-mapping results between the GBMI and individual biobanks**

826 To directly compare with fine-mapping results from the GBMI meta-analyses, we used our fine-
827 mapping results of nine disease endpoints (asthma⁶⁴, COPD⁶⁴, gout, heart failure⁷³, IPF⁶², primary
828 open angle glaucoma⁷⁴, thyroid cancer, stroke⁷⁵, and venous thromboembolism⁷⁶) in BBJ⁵⁸,
829 FinnGen²⁰, and UKBB¹⁹ Europeans that were also part of the GBMI meta-analyses for the same
830 traits. For comparison, we computed the maximum PIP for each variant and the minimum size of
831 95% CS across BBJ, FinnGen, and UKBB. We restricted the 95% CS in biobanks to those that
832 contain the lead variants from the GBMI. We defined the PIP difference between the GBMI and
833 individual biobanks as $\Delta\text{PIP} = \text{PIP}(\text{GBMI}) - \text{the maximum PIP across the biobanks}$.

834
835 We conducted functional enrichment analysis to compare between the GBMI meta-analysis and
836 individual biobanks because unbiased comparison of PIP requires conditioning on likely causal
837 variants independent of the fine-mapping results, and functional annotations have been shown to
838 be enriched for causal variants. Using functional categories (coding [pLoF, missense, and
839 synonymous], 5'/3' UTR, promoter, and CRE) from our previous study^{16,17}, we estimated
840 functional enrichments of variants in each functional category based on 1) top PIP rankings and
841 2) ΔPIP bins. Since fine-mapping PIP in the GBMI meta-analysis can be miscalibrated, we
842 performed a comparison based on top PIP rankings to assess whether the ordering given by
843 GBMI PIPs is more informative than the ordering given by the biobanks. For the top PIP rankings,
844 we took the top 0.5%, 0.1%, and 0.05% variants based on the PIP rankings in the GBMI and
845 individual biobanks. We computed enrichment as a relative risk = (proportion of top X% PIP
846 variants in the GBMI that are in the annotation) / (proportion of top X% PIP variants in the
847 individual biobanks that are in the annotation). For ΔPIP bins, we defined three bins using different
848 thresholds ($\theta = 0.01, 0.05, \text{ and } 0.1$): 1) decreased PIP bin, $\Delta\text{PIP} < -\theta$, 2) null bin, $-\theta \leq \Delta\text{PIP} \leq \theta$,
849 and 3) increased PIP bin, $\theta < \Delta\text{PIP}$. We computed enrichment as a relative risk = (proportion of
850 variants in the decreased/increased PIP bin that are in the annotation) / (proportion of variants in
851 the null PIP bin). We combined coding, UTR, and promoter categories for this analysis due to the
852 limited number of variants for each bin.

853

854 References

- 855 1. Evangelou, E. & Ioannidis, J. P. a. Meta-analysis methods for genome-wide association
856 studies and beyond. *Nat. Rev. Genet.* **14**, 379–389 (2013).
- 857 2. Mahajan, A. *et al.* Fine-mapping type 2 diabetes loci to single-variant resolution using high-
858 density imputation and islet-specific epigenome maps. *Nat. Genet.* **50**, 1505–1513 (2018).
- 859 3. Spracklen, C. N. *et al.* Identification of type 2 diabetes loci in 433,540 East Asian
860 individuals. *Nature* **582**, 240–245 (2020).
- 861 4. Ripke, S. *et al.* Biological insights from 108 schizophrenia-associated genetic loci. *Nature*
862 **511**, 421–427 (2014).
- 863 5. The Schizophrenia Working Group of the Psychiatric Genomics Consortium, Ripke, S.,
864 Walters, J. T. R. & O'Donovan, M. C. Mapping genomic loci prioritises genes and implicates
865 synaptic biology in schizophrenia. *bioRxiv* (2020) doi:10.1101/2020.09.12.20192922.
- 866 6. Okada, Y. *et al.* Genetics of rheumatoid arthritis contributes to biology and drug discovery.
867 *Nature* **506**, 376–381 (2014).
- 868 7. Ishigaki, K. *et al.* Trans-ancestry genome-wide association study identifies novel genetic
869 mechanisms in rheumatoid arthritis. *bioRxiv* (2021) doi:10.1101/2021.12.01.21267132.
- 870 8. Locke, A. E. *et al.* Genetic studies of body mass index yield new insights for obesity
871 biology. *Nature* **518**, 197–206 (2015).
- 872 9. Graham, S. E. *et al.* The power of genetic diversity in genome-wide association studies of
873 lipids. *Nature* **600**, 675–679 (2021).
- 874 10. Zhou, W. *et al.* Global Biobank Meta-analysis Initiative: powering genetic discovery across
875 human diseases. *bioRxiv* (2021) doi:10.1101/2021.11.19.21266436.
- 876 11. Visscher, P. M. *et al.* 10 Years of GWAS Discovery: Biology, Function, and Translation.
877 *Am. J. Hum. Genet.* **101**, 5–22 (2017).
- 878 12. Shendure, J., Findlay, G. M. & Snyder, M. W. Genomic Medicine-Progress, Pitfalls, and
879 Promise. *Cell* **177**, 45–57 (2019).
- 880 13. Schaid, D. J., Chen, W. & Larson, N. B. From genome-wide associations to candidate
881 causal variants by statistical fine-mapping. *Nat. Rev. Genet.* **19**, 491–504 (2018).
- 882 14. Ulirsch, J. C. *et al.* Interrogation of human hematopoiesis at single-cell and single-variant
883 resolution. *Nat. Genet.* **51**, 683–693 (2019).
- 884 15. Weissbrod, O. *et al.* Functionally informed fine-mapping and polygenic localization of
885 complex trait heritability. *Nat. Genet.* **52**, 1355–1363 (2020).
- 886 16. Ulirsch, J. C. & Kanai, M. An annotated atlas of causal variants underlying complex traits
887 and gene expression. *Under review*.
- 888 17. Kanai, M. *et al.* Insights from complex trait fine-mapping across diverse populations.
889 *medRxiv* (2021) doi:10.1101/2021.09.03.21262975.
- 890 18. Nagai, A. *et al.* Overview of the BioBank Japan Project: Study design and profile. *Journal of*
891 *Epidemiology* **27**, S2–S8 (2017).
- 892 19. Bycroft, C. *et al.* The UK Biobank resource with deep phenotyping and genomic data.
893 *Nature* **562**, 203–209 (2018).
- 894 20. Kurki, M. I. *et al.* FinnGen: Unique genetic insights from combining isolated population and
895 national health register data. *bioRxiv* (2022) doi:10.1101/2022.03.03.22271360.
- 896 21. Benner, C. *et al.* FINEMAP: Efficient variable selection using summary data from genome-
897 wide association studies. *Bioinformatics* **32**, 1493–1501 (2016).
- 898 22. Benner, C., Havulinna, A. S., Salomaa, V., Ripatti, S. & Pirinen, M. Refining fine-mapping:
899 effect sizes and regional heritability. *bioRxiv* 318618 (2018) doi:10.1101/318618.
- 900 23. Wang, G., Sarkar, A., Carbonetto, P. & Stephens, M. A simple new approach to variable
901 selection in regression, with application to genetic fine mapping. *J. R. Stat. Soc. Series B*
902 *Stat. Methodol.* **25**, 1 (2020).

- 903 24. Onengut-Gumuscu, S. *et al.* Fine mapping of type 1 diabetes susceptibility loci and
904 evidence for colocalization of causal variants with lymphoid gene enhancers. *Nat. Genet.*
905 **47**, 381–386 (2015).
- 906 25. Levey, D. F. *et al.* Bi-ancestral depression GWAS in the Million Veteran Program and meta-
907 analysis in >1.2 million individuals highlight new therapeutic directions. *Nat. Neurosci.* **24**,
908 954–963 (2021).
- 909 26. Gharakhani, P. *et al.* Genome-wide meta-analysis identifies 127 open-angle glaucoma loci
910 with consistent effect across ancestries. *Nat. Commun.* **12**, 1258 (2021).
- 911 27. Chen, J. *et al.* The trans-ancestral genomic architecture of glycemic traits. *Nat. Genet.* **53**,
912 840–860 (2021).
- 913 28. Zhou, W. *et al.* GWAS of thyroid stimulating hormone highlights pleiotropic effects and
914 inverse association with thyroid cancer. *Nat. Commun.* **11**, 1–13 (2020).
- 915 29. Wightman, D. P. *et al.* A genome-wide association study with 1,126,563 individuals
916 identifies new risk loci for Alzheimer’s disease. *Nat. Genet.* **53**, 1276–1282 (2021).
- 917 30. Chen, M.-H. *et al.* Trans-ethnic and Ancestry-Specific Blood-Cell Genetics in 746,667
918 Individuals from 5 Global Populations. *Cell* **182**, 1198–1213.e14 (2020).
- 919 31. Wakefield, J. A Bayesian measure of the probability of false discovery in genetic
920 epidemiology studies. *Am. J. Hum. Genet.* **81**, 208–227 (2007).
- 921 32. Wakefield, J. Bayes factors for genome-wide association studies: comparison with P-
922 values. *Genet. Epidemiol.* **33**, 79–86 (2009).
- 923 33. Hormozdvari, F., Kostem, E., Kang, E. Y., Pasaniuc, B. & Eskin, E. Identifying Causal
924 Variants at Loci with Multiple Signals of Association. *Genetics* **198**, 497–508 (2014).
- 925 34. Kichaev, G. *et al.* Integrating functional data to prioritize causal variants in statistical fine-
926 mapping studies. *PLoS Genet.* **10**, e1004722 (2014).
- 927 35. Kichaev, G. & Pasaniuc, B. Leveraging Functional-Annotation Data in Trans-ethnic Fine-
928 Mapping Studies. *Am. J. Hum. Genet.* **97**, 260–271 (2015).
- 929 36. Li, D., Zhao, H. & Gelernter, J. Strong protective effect of the aldehyde dehydrogenase
930 gene (ALDH2) 504Iys (*2) allele against alcoholism and alcohol-induced medical diseases
931 in Asians. *Hum. Genet.* **131**, 725–737 (2012).
- 932 37. Brown, B. C., Ye, C. J., Price, A. L. & Zaitlen, N. Transethnic Genetic-Correlation Estimates
933 from Summary Statistics. *Am. J. Hum. Genet.* **99**, 76–88 (2016).
- 934 38. Shi, H. *et al.* Population-specific causal disease effect sizes in functionally important
935 regions impacted by selection. *bioRxiv* 803452 (2020) doi:10.1101/803452.
- 936 39. COVID-19 Host Genetics Initiative. Mapping the human genetic architecture of COVID-19.
937 *Nature* **600**, 472–477 (2021).
- 938 40. Dendrou, C. A. *et al.* Resolving *TYK2* locus genotype-to-phenotype differences in
939 autoimmunity. *Sci. Transl. Med.* **8**, 363ra149 (2016).
- 940 41. Couturier, N. *et al.* Tyrosine kinase 2 variant influences T lymphocyte polarization and
941 multiple sclerosis susceptibility. *Brain* **134**, 693–703 (2011).
- 942 42. Li, Z. *et al.* Two rare disease-associated *Tyk2* variants are catalytically impaired but
943 signaling competent. *J. Immunol.* **190**, 2335–2344 (2013).
- 944 43. Lam, M. *et al.* RICOPILI: Rapid Imputation for COnsortias PIpeLine. *Bioinformatics* **36**,
945 930–933 (2020).
- 946 44. Huang, H. *et al.* Fine-mapping inflammatory bowel disease loci to single-variant resolution.
947 *Nature* **547**, 173–178 (2017).
- 948 45. Winkler, T. W. *et al.* Quality control and conduct of genome-wide association meta-
949 analyses. *Nat. Protoc.* **9**, 1192–1212 (2014).
- 950 46. Chen, W. *et al.* Improved analyses of GWAS summary statistics by reducing data
951 heterogeneity and errors. *Nat. Commun.* **12**, 7117 (2021).

- 952 47. Yang, J. *et al.* Conditional and joint multiple-SNP analysis of GWAS summary statistics
953 identifies additional variants influencing complex traits. *Nat. Genet.* **44**, 369–75, S1–3
954 (2012).
- 955 48. Buniello, A. *et al.* The NHGRI-EBI GWAS Catalog of published genome-wide association
956 studies, targeted arrays and summary statistics 2019. *Nucleic Acids Res.* **47**, D1005–
957 D1012 (2019).
- 958 49. The 1000 Genomes Project Consortium. A global reference for human genetic variation.
959 *Nature* **526**, 68–74 (2015).
- 960 50. McCarthy, S. *et al.* A reference panel of 64,976 haplotypes for genotype imputation. *Nat.*
961 *Genet.* **48**, 1279–1283 (2016).
- 962 51. Taliun, D. *et al.* Sequencing of 53,831 diverse genomes from the NHLBI TOPMed Program.
963 *Nature* **590**, 290–299 (2021).
- 964 52. Schoech, A. P. *et al.* Quantification of frequency-dependent genetic architectures in 25 UK
965 Biobank traits reveals action of negative selection. *Nat. Commun.* **10**, 790 (2019).
- 966 53. Yang, J. *et al.* Genome partitioning of genetic variation for complex traits using common
967 SNPs. *Nat. Genet.* **43**, 519–525 (2011).
- 968 54. Ormond, C., Ryan, N. M., Corvin, A. & Heron, E. A. Converting single nucleotide variants
969 between genome builds: from cautionary tale to solution. *Brief. Bioinform.* **22**, (2021).
- 970 55. Asimit, J. L., Hatzikotoulas, K., McCarthy, M., Morris, A. P. & Zeggini, E. Trans-ethnic study
971 design approaches for fine-mapping. *Eur. J. Hum. Genet.* **24**, 1330–1336 (2016).
- 972 56. Marchini, J. & Howie, B. Genotype imputation for genome-wide association studies. *Nat.*
973 *Rev. Genet.* **11**, 499–511 (2010).
- 974 57. Karczewski, K. J. *et al.* The mutational constraint spectrum quantified from variation in
975 141,456 humans. *Nature* **581**, 434–443 (2020).
- 976 58. Sakaue, S. *et al.* A cross-population atlas of genetic associations for 220 human
977 phenotypes. *Nat. Genet.* **53**, 1415–1424 (2021).
- 978 59. The GTEx Consortium. The GTEx Consortium atlas of genetic regulatory effects across
979 human tissues. *Science* **369**, 1318–1330 (2020).
- 980 60. Kerimov, N. *et al.* A compendium of uniformly processed human gene expression and
981 splicing quantitative trait loci. *Nat. Genet.* **53**, 1290–1299 (2021).
- 982 61. Koskela, J. T. *et al.* Genetic variant in SPDL1 reveals novel mechanism linking pulmonary
983 fibrosis risk and cancer protection. *bioRxiv* (2021) doi:10.1101/2021.05.07.21255988.
- 984 62. Partanen, J. J. *et al.* Leveraging global multi-ancestry meta-analysis in the study of
985 Idiopathic Pulmonary Fibrosis genetics. *bioRxiv* (2021) doi:10.1101/2021.12.29.21268310.
- 986 63. Foreman, M. G. *et al.* Alpha-1 Antitrypsin PiMZ Genotype Is Associated with Chronic
987 Obstructive Pulmonary Disease in Two Racial Groups. *Ann. Am. Thorac. Soc.* **14**, 1280–
988 1287 (2017).
- 989 64. Tsuo, K. *et al.* Multi-ancestry meta-analysis of asthma identifies novel associations and
990 highlights the value of increased power and diversity. *bioRxiv* (2021)
991 doi:10.1101/2021.11.30.21267108.
- 992 65. Benonisdottir, S. *et al.* Epigenetic and genetic components of height regulation. *Nat.*
993 *Commun.* **7**, 13490 (2016).
- 994 66. Marouli, E. *et al.* Rare and low-frequency coding variants alter human adult height. *Nature*
995 **542**, 186–190 (2017).
- 996 67. Langefeld, C. D. *et al.* Transancestral mapping and genetic load in systemic lupus
997 erythematosus. *Nat. Commun.* **8**, 16021 (2017).
- 998 68. Hargreaves, C. E. *et al.* Fcγ receptors: genetic variation, function, and disease. *Immunol.*
999 *Rev.* **268**, 6–24 (2015).
- 1000 69. Franke, L. *et al.* Association analysis of copy numbers of FC-gamma receptor genes for
1001 rheumatoid arthritis and other immune-mediated phenotypes. *Eur. J. Hum. Genet.* **24**, 263–
1002 270 (2016).

- 1003 70. UK10K Consortium *et al.* The UK10K project identifies rare variants in health and disease.
1004 *Nature* **526**, 82–90 (2015).
- 1005 71. Wang, Y. *et al.* Global biobank analyses provide lessons for computing polygenic risk
1006 scores across diverse cohorts. *medRxiv* 2021.11.18.21266545 (2021).
- 1007 72. Namba, S. *et al.* A practical guideline of genomics-driven drug discovery in the era of global
1008 biobank meta-analysis. *bioRxiv* (2021) doi:10.1101/2021.12.03.21267280.
- 1009 73. Wu, K.-H. H. *et al.* Polygenic risk score from a multi-ancestry GWAS uncovers susceptibility
1010 of heart failure. *bioRxiv* (2021) doi:10.1101/2021.12.06.21267389.
- 1011 74. Lo Faro, V. *et al.* Genome-wide association meta-analysis identifies novel ancestry-specific
1012 primary open-angle glaucoma loci and shared biology with vascular mechanisms and cell
1013 proliferation. *bioRxiv* (2021) doi:10.1101/2021.12.16.21267891.
- 1014 75. Surakka, I. *et al.* Multi-ancestry meta-analysis identifies 2 novel loci associated with
1015 ischemic stroke and reveals heterogeneity of effects between sexes and ancestries.
1016 *bioRxiv* (2022) doi:10.1101/2022.02.28.22271647.
- 1017 76. Wolford, B. *et al.* Multi-ancestry GWAS for venous thromboembolism identifies novel loci
1018 followed by experimental validation. *In preparation*.
- 1019 77. Aneas, I. *et al.* Asthma-associated genetic variants induce IL33 differential expression
1020 through an enhancer-blocking regulatory region. *Nat. Commun.* **12**, 6115 (2021).
- 1021 78. Vladich, F. D. *et al.* IL-13 R130Q, a common variant associated with allergy and asthma,
1022 enhances effector mechanisms essential for human allergic inflammation. *J. Clin. Invest.*
1023 **115**, 747–754 (2005).
- 1024 79. Farh, K. K.-H. *et al.* Genetic and epigenetic fine mapping of causal autoimmune disease
1025 variants. *Nature* **518**, 337–343 (2015).
- 1026 80. Wojcik, G. L. *et al.* Genetic analyses of diverse populations improves discovery for complex
1027 traits. *Nature* **570**, 514–518 (2019).
- 1028 81. Luo, Y. *et al.* A high-resolution HLA reference panel capturing global population diversity
1029 enables multi-ancestry fine-mapping in HIV host response. *Nat. Genet.* **53**, 1504–1516
1030 (2021).
- 1031 82. Sakaue, S. *et al.* Decoding the diversity of killer immunoglobulin-like receptors by deep
1032 sequencing and a high-resolution imputation method. *Cell Genomics* **2**, (2022).
- 1033 83. Mukamel, R. E. *et al.* Protein-coding repeat polymorphisms strongly shape diverse human
1034 phenotypes. *Science* **373**, 1499–1505 (2021).
- 1035 84. Gagliano Taliun, S. A. *et al.* Exploring and visualizing large-scale genetic associations by
1036 using PheWeb. *Nat. Genet.* **52**, 550–552 (2020).
- 1037 85. Su, Z., Marchini, J. & Donnelly, P. HAPGEN2: Simulation of multiple disease SNPs.
1038 *Bioinformatics* **27**, 2304–2305 (2011).
- 1039 86. Manichaikul, A. *et al.* Robust relationship inference in genome-wide association studies.
1040 *Bioinformatics* **26**, 2867–2873 (2010).
- 1041 87. Chang, C. C. *et al.* Second-generation PLINK: rising to the challenge of larger and richer
1042 datasets. *Gigascience* **4**, 1–16 (2015).
- 1043 88. Wei, X. & Nielsen, R. CCR5- Δ 32 is deleterious in the homozygous state in humans. *Nat.*
1044 *Med.* **25**, 909–910 (2019).
- 1045 89. Maier, R. *et al.* No statistical evidence for an effect of CCR5- Δ 32 on lifespan in the UK
1046 Biobank cohort. *Nat. Med.* **26**, 178–180 (2020).
- 1047 90. Das, S. *et al.* Next-generation genotype imputation service and methods. *Nat. Genet.* **48**,
1048 1284–1287 (2016).
- 1049 91. Loh, P.-R. *et al.* Reference-based phasing using the Haplotype Reference Consortium
1050 panel. *Nat. Genet.* **48**, 1443–1448 (2016).
- 1051 92. Loh, P.-R. *et al.* Contrasting genetic architectures of schizophrenia and other complex
1052 diseases using fast variance-components analysis. *Nat. Genet.* **47**, 1385–1392 (2015).
- 1053 93. McLaren, W. *et al.* The Ensembl Variant Effect Predictor. *Genome Biol.* **17**, 122 (2016).

1054 **Global Biobank Meta-analysis Initiative**

1055 Wei Zhou^{1,2,3}, Masahiro Kanai^{1,2,3,4,5}, Kuan-Han H Wu⁶, Rasheed Humaira^{7,8,9}, Kristin Tsuo^{1,2,3},
1056 Jibril B Hirbo^{10,11}, Ying Wang^{1,2,3}, Arjun Bhattacharya¹², Huiling Zhao¹³, Shinichi Namba⁵, Ida
1057 Surakka¹⁴, Brooke N Wolford⁶, Valeria Lo Faro^{15,16,17}, Esteban A Lopera-Maya¹⁸, Kristi Läll¹⁹,
1058 Marie-Julie Favé²⁰, Sinéad B Chapman^{2,3}, Juha Karjalainen^{1,2,3,21}, Mitja Kurki^{1,2,3,21}, Maasha
1059 Mutaamba^{1,2,3,21}, Ben M Brumpton²², Sameer Chavan²³, Tzu-Ting Chen²⁴, Michelle Daya²³, Yi
1060 Ding^{12,25}, Yen-Chen A Feng²⁶, Christopher R Gignoux²³, Sarah E Graham¹⁴, Whitney E Hornsby¹⁴,
1061 Nathan Ingold²⁷, Ruth Johnson^{12,28}, Triin Laisk¹⁹, Kuang Lin²⁹, Jun Lv³⁰, Iona Y Millwood^{29,31}, Priit
1062 Palta^{19,21}, Anita Pandit³², Michael Preuss³³, Unnur Thorsteinsdottir³⁴, Jasmina Uzunovic²⁰,
1063 Matthew Zawistowski³², Xue Zhong^{10,35}, Archie Campbell³⁶, Kristy Crooks²³, Geertruida h De
1064 Bock³⁷, Nicholas J Douville^{38,39}, Sarah Finer⁴⁰, Lars G Fritsche³², Christopher J Griffiths⁴⁰, Yu
1065 Guo⁴¹, Karen A Hunt⁴², Takahiro Konuma^{5,43}, Riccardo E Marioni³⁶, Jansonius Nomdo¹⁵, Snehal
1066 Patil³², Nicholas Rafaels²³, Anne Richmond⁴⁴, Jonathan A Shortt²³, Peter Straub^{10,35}, Ran Tao^{35,45},
1067 Brett Vanderwerff³², Kathleen C Barnes²³, Marike Boezen³⁷, Zhengming Chen^{29,31}, Chia-Yen
1068 Chen⁴⁶, Judy Cho³³, George Davey Smith^{13,47}, Hilary K Finucane^{1,2,3}, Lude Franke¹⁸, Eric
1069 Gamazon^{35,48}, Andrea Ganna^{1,2,21}, Tom R Gaunt¹³, Tian Ge^{49,50}, Hailiang Huang^{1,2}, Jennifer
1070 Huffman⁵¹, Clara Lajonchere^{52,53}, Matthew H Law²⁷, Liming Li³⁰, Cecilia M Lindgren⁵⁴, Ruth JF
1071 Loos³³, Stuart MacGregor²⁷, Koichi Matsuda⁵⁵, Catherine M Olsen²⁷, David J Porteous³⁶, Jordan
1072 A Shavit⁵⁶, Harold Snieder³⁷, Richard C Trembath⁵⁷, Judith M Vonk³⁷, David Whiteman²⁷, Stephen
1073 J Wicks²³, Cisca Wijmenga¹⁸, John Wright⁵⁸, Jie Zheng¹³, Xiang Zhou³², Philip Awadalla^{20,59},
1074 Michael Boehnke³², Nancy J Cox^{10,60}, Daniel H Geschwind^{52,61,62}, Caroline Hayward⁴⁴, Kristian
1075 Hveem²², Eimear E Kenny⁶³, Yen-Feng Lin^{24,64,65}, Reedik Mägi¹⁹, Hilary C Martin⁶⁶, Sarah E
1076 Medland²⁷, Yukinori Okada^{5,67,68,69,70}, Aarno V Palotie^{1,2,21}, Bogdan Pasaniuc^{12,25,52,61,71}, Serena
1077 Sanna^{18,72}, Jordan W Smoller⁷³, Kari Stefansson³⁴, David A van Heel⁴², Robin G Walters^{29,31},
1078 Sebastian Zoellner³², **Biobank Japan, BioMe, BioVU, Canadian Partnership for Tomorrow,**
1079 **China Kadoorie Biobank Collaborative Group, Colorado Center for Personalized Medicine,**
1080 **deCODE Genetics, Estonian Biobank, FinnGen, Generation Scotland, Genes & Health,**
1081 **LifeLines, Mass General Brigham Biobank, Michigan Genomics Initiative, QIMR Berghofer**
1082 **Biobank, Taiwan Biobank, The HUNT Study, UCLA ATLAS Community Health Initiative, UK**
1083 **Biobank, Alicia R Martin^{1,2,3}, Cristen J Willer^{6,14,74*}, Mark J Daly^{1,2,3,21*}, Benjamin M Neale^{1,2,3*}**

1084
1085 *These authors jointly supervised the initiative
1086

1087 ¹Analytic and Translational Genetics Unit, Department of Medicine, Massachusetts General
1088 Hospital, Boston, MA, USA, ²Stanley Center for Psychiatric Research, Broad Institute of MIT and
1089 Harvard, Cambridge, MA, USA, ³Program in Medical and Population Genetics, Broad Institute of
1090 MIT and Harvard, Cambridge, MA, USA, ⁴Department of Biomedical Informatics, Harvard Medical
1091 School, Boston, MA, USA, ⁵Department of Statistical Genetics, Osaka University Graduate
1092 School of Medicine, Suita 565-0871, Japan, ⁶Department of Computational Medicine and
1093 Bioinformatics, University of Michigan, Ann Arbor, MI, USA, ⁷K. G. Jebsen Center for Genetic
1094 Epidemiology, Department of Public Health and Nursing, Faculty of Medicine and Health, NTNU,
1095 Norwegian University of Science and Technology, Trondheim, Norway, ⁸Division of Medicine and
1096 Laboratory Sciences, University of Oslo, Norway, ⁹MRC Integrative Epidemiology Unit, Population
1097 Health Sciences, Bristol Medical School, University of Bristol, Bristol, UK, ¹⁰Department of
1098 Medicine, Division of Genetic Medicine, Vanderbilt University Medical Center, Nashville, TN, USA,
1099 ¹¹Vanderbilt Genetic Institute, Vanderbilt University Medical Center, Nashville, TN, USA,
1100 ¹²Department of Pathology and Laboratory Medicine, David Geffen School of Medicine, University
1101 of California, Los Angeles, Los Angeles, CA, USA, ¹³MRC Integrative Epidemiology Unit (IEU),
1102 Bristol Medical School, University of Bristol, Oakfield House, Oakfield Grove, Bristol, BS8 2BN,
1103 UK, ¹⁴Department of Internal Medicine, University of Michigan, Ann Arbor, MI, USA, ¹⁵University

1104 of Groningen, UMCG, Department of Ophthalmology, Groningen, the Netherlands, ¹⁶Department
1105 of Clinical Genetics, Amsterdam University Medical Center (AMC), Amsterdam, the Netherlands,
1106 ¹⁷Department of Immunology, Genetics and Pathology, Science for Life Laboratory, Uppsala
1107 University, Uppsala, Sweden, ¹⁸University of Groningen, UMCG, Department of Genetics,
1108 Groningen, the Netherlands, ¹⁹Estonian Genome Centre, Institute of Genomics, University of
1109 Tartu, Tartu, Estonia, ²⁰Ontario Institute for Cancer Research, Toronto, ON, Canada, ²¹Institute
1110 for Molecular Medicine Finland, University of Helsinki, Helsinki, Finland, ²²K.G. Jebsen Center for
1111 Genetic Epidemiology, Department of Public Health and Nursing, NTNU, Norwegian University of
1112 Science and Technology, Trondheim, 7030, Norway, ²³University of Colorado - Anschutz Medical
1113 Campus, Aurora, CO, USA, ²⁴Center for Neuropsychiatric Research, National Health Research
1114 Institutes, Miaoli, Taiwan, ²⁵Bioinformatics Interdepartmental Program, University of California,
1115 Los Angeles, Los Angeles, CA, USA, ²⁶Division of Biostatistics, Institute of Epidemiology and
1116 Preventive Medicine, College of Public Health, National Taiwan University, Taiwan, ²⁷QIMR
1117 Berghofer Medical Research Institute, Brisbane, Australia, ²⁸Department of Computer Science,
1118 University of California, Los Angeles, Los Angeles, CA, USA, ²⁹Nuffield Department of Population
1119 Health, University of Oxford, Oxford, UK, ³⁰Department of Epidemiology and Biostatistics, School
1120 of Public Health, Peking University Health Science Center, Beijing, China, ³¹MRC Population
1121 Health Research Unit, University of Oxford, Oxford, UK, ³²Department of Biostatistics and Center
1122 for Statistical Genetics, University of Michigan, Ann Arbor, MI, USA, ³³The Charles Bronfman
1123 Institute for Personalized Medicine, Icahn School of Medicine at Mount Sinai, New York, NY, USA,
1124 ³⁴deCODE Genetics/Amgen inc., 101, Reykjavik, Iceland, ³⁵Vanderbilt Genetics Institute,
1125 Vanderbilt University Medical Center, Nashville, TN, USA, ³⁶Centre for Genomic and
1126 Experimental Medicine, Institute of Genetics and Cancer, University of Edinburgh, Edinburgh, UK,
1127 ³⁷University of Groningen, UMCG, Department of Epidemiology, Groningen, the Netherlands,
1128 ³⁸Department of Anesthesiology, Michigan Medicine, Ann Arbor, MI, USA, ³⁹Institute of Healthcare
1129 Policy & Innovation, University of Michigan, Ann Arbor, MI, USA, ⁴⁰Wolfson Institute of Population
1130 Health, Queen Mary University of London, London, UK, ⁴¹Chinese Academy of Medical Sciences,
1131 Beijing, China, ⁴²Blizard Institute, Queen Mary University of London, London, UK, ⁴³Central
1132 Pharmaceutical Research Institute, JAPAN TOBACCO INC., Takatsuki 569-1125, Japan,
1133 ⁴⁴Medical Research Council Human Genetics Unit, Institute of Genetics and Cancer, University
1134 of Edinburgh, Edinburgh, UK, ⁴⁵Department of Biostatistics, Vanderbilt University Medical Center,
1135 Nashville, TN, USA, ⁴⁶Biogen, Cambridge, MA, USA, ⁴⁷NIHR Bristol Biomedical Research Centre,
1136 Bristol, UK, ⁴⁸MRC Epidemiology Unit, University of Cambridge, Cambridge, UK, ⁴⁹Psychiatric
1137 and Neurodevelopmental Genetics Unit, Center for Genomic Medicine, Massachusetts General
1138 Hospital, Boston, MA, USA, ⁵⁰Center for Precision Psychiatry, Massachusetts General Hospital,
1139 Boston, MA, USA, ⁵¹Centre for Population Genomics, VA Boston Healthcare System, Boston,
1140 MA, USA, ⁵²Institute of Precision Health, University of California, Los Angeles, Los Angeles, CA,
1141 USA, ⁵³Program in Neurogenetics, Department of Neurology, David Geffen School of Medicine,
1142 University of California, Los Angeles, Los Angeles, CA, USA, ⁵⁴Big Data Institute, Li Ka Shing
1143 Centre for Health Information and Discovery, University of Oxford, Oxford, UK, ⁵⁵Department of
1144 Computational Biology and Medical Sciences, Graduate school of Frontier Sciences, The
1145 University of Tokyo, Tokyo, Japan, ⁵⁶University of Michigan, Department of Pediatrics, Ann Arbor
1146 MI 48109, ⁵⁷School of Basic and Medical Biosciences, Faculty of Life Sciences and Medicine,
1147 King's College London, London, UK, ⁵⁸Bradford Institute for Health Research, Bradford Teaching
1148 Hospitals National Health Service (NHS) Foundation Trust, Bradford, UK, ⁵⁹Department of
1149 Molecular Genetics, University of Toronto, Toronto, ON, Canada, ⁶⁰Vanderbilt Genetics Institute,
1150 Vanderbilt University Medical Center, Nashville, TN, USA, ⁶¹Department of Human Genetics,
1151 David Geffen School of Medicine, University of California, Los Angeles, Los Angeles, CA, USA,
1152 ⁶²Department of Neurology, David Geffen School of Medicine, University of California, Los
1153 Angeles, Los Angeles, CA, USA, ⁶³Institute for Genomic Health, Icahn School of Medicine at
1154 Mount Sinai, New York, NY, USA, ⁶⁴Department of Public Health & Medical Humanities, School

1155 of Medicine, National Yang Ming Chiao Tung University, Taipei, Taiwan, ⁶⁵Institute of Behavioral
1156 Medicine, College of Medicine, National Cheng Kung University, Tainan, Taiwan, ⁶⁶Medical and
1157 Population Genomics, Wellcome Sanger Institute, Hinxton, UK, ⁶⁷Center for Infectious Disease
1158 Education and Research (CiDER), Osaka University, Suita 565-0871, Japan, ⁶⁸Laboratory of
1159 Statistical Immunology, Immunology Frontier Research Center (WPI-IFReC), Osaka University,
1160 Suita 565-0871, Japan, ⁶⁹Laboratory for Systems Genetics, RIKEN Center for Integrative Medical
1161 Sciences, Yokohama, Japan, ⁷⁰Integrated Frontier Research for Medical Science Division,
1162 Institute for Open and Transdisciplinary Research Initiatives, Osaka University, Suita 565-0871,
1163 Japan, ⁷¹Department of Computational Medicine, David Geffen School of Medicine, University of
1164 California, Los Angeles, Los Angeles, CA, USA, ⁷²Institute for Genetics and Biomedical Research,
1165 National Research Council, Cagliari 09100, Italy, ⁷³Psychiatric and Neurodevelopmental Genetics
1166 Unit, Center for Genomic Medicine, Massachusetts General Hospital, Boston, MA, ⁷⁴Department
1167 of Human Genetics, University of Michigan, Ann Arbor, MI, USA.