# Fine-Mapping and Credible Set Construction using a Multi-population Joint Analysis of Marginal Summary Statistics from Genome-wide Association Studies

**Jiayi Shen[1], Lai Jiang[1], Kan Wang[1], Anqi Wang[2], Fei Chen[2], Paul J. Newcombe[3], Christopher A. Haiman[2,4], David V. Conti[1, 2, 4, *]**

## Abstract

Recent advancement in Genome-wide Association Studies (GWAS) comes from not only increasingly larger sample sizes but also the shifted focus towards underrepresented populations. Multi-population GWAS may increase power to detect novel risk variants and improve fine-mapping resolution by leveraging evidence from diverse populations and accounting for the difference in linkage disequilibrium (LD) across ethnic groups. Here, we expand upon our previous approach for single-population fine-mapping through Joint Analysis of Marginal SNP Effects (JAM) to a multi-population analysis (mJAM). Under the assumption that true causal variants are

[1] Division of Biostatistics, Department of Population and Public Health Sciences, Keck School of Medicine, University of Southern California, Los Angeles, California, 90032, USA

[2] Center for Genetic Epidemiology, Keck School of Medicine, University of Southern California, Los Angeles, California, 90033, USA

[3] MRC Biostatistics Unit, University of Cambridge, Cambridge, CB2 0SR, United Kingdom

[4] Norris Comprehensive Cancer Center, University of Southern California, Los Angeles, California, 90033, USA

*Correspondence: dconti@med.usc.edu

15  common across studies, we implement a novel version of JAM that conditions on

16  multiple SNPs while explicitly incorporating the different LD structures across

17  populations. The mJAM framework can be used to first select index variants using the

18  mJAM likelihood with any feature selection approach. In addition, we present a novel

19  approach leveraging the ideas of mediation to construct credible sets for these index

20  variants. Construction of such credible sets can be performed given any existing index

21  variants. We illustrate the implementation of the mJAM likelihood through two

22  implementations: mJAM-SuSiE (a Bayesian approach) and mJAM-Forward selection.

23  Through simulation studies based on realistic effect sizes and levels of LD, we

24  demonstrated that mJAM performs better than other existing multi-ethnic methods for

25  constructing concise credible sets that include the underlying causal variants. In real

26  data examples taken from the most recent multi-population prostate cancer GWAS, we

27  showed several practical advantages of mJAM over other existing methods.

# Introduction

The development of high-throughput genotyping and genotype imputation has boosted the application of genome-wide association studies (GWAS) which is now a standard approach to identify susceptibility loci or genomic regions for many complex diseases and traits[1,2]. However, the linkage disequilibrium (LD) of single-nucleotide polymorphisms (SNPs) makes it challenging to determine the true causal variant(s) within a region or to further prioritize genetic variants for functional studies[2,3].

Fine-mapping is a post-GWAS approach which seeks to specify the underlying causal variant and quantify the strength of effect given existing evidence that a certain region is likely to contain at least one causal signal. Many methods for fine-mapping often start with a lead SNP – the SNP with the smallest $p$-value within one region – and then they examine additional highly correlated neighboring SNPs in the region using different strategies such as setting a threshold on pairwise correlation ($r^2$) with the lead SNP[2]. These approaches are intuitive but do not jointly analyze all the SNPs within a region. In addition, they often do not generalize easily to the investigation in multiple populations.

More recent and advanced fine-mapping approaches attempt to jointly or conditionally analyze all SNPs within a region, and include stepwise regression[2], penalized regression[4-7], and Bayesian methods[8-11]. Conditional step-wise selection has been used to discover multiple signals at a locus with individual level data[12,13]. However, stepwise selection can be very unstable with a large amount of highly correlated SNPs and the $P$-values of the signals in the final selected model tend to be conservative[2].

50  Alternative selection approaches with individual level data are penalized regression

51  models, such as lasso[4] and elastic net[5], and Bayesian methods, such as CAVIAR[10] and

52  Sum of Single Effect models (SuSiE)[11]. In contrast to step-wise selection, penalized

53  regression techniques are potentially more stable because the penalty term encourages

54  shrinkage of effect estimates towards zero resulting in sparsity and robust estimation.

55  However, penalized often do not perform well with highly correlated SNPs and they do

56  not represent the uncertainty in effect estimation and model selection[2]. In contrast, fully

57  Bayesian methods compute posterior probabilities for models within the model space to

58  infer the probability of causality for each SNP and often result in credible sets to

59  measure the fine-mapping resolution using these probabilities. Ideally, exact inference

60  is possible by enumerating all possible models or combinations of SNPs but the model

61  space increases so rapidly that exhaustive searches become impractical as the number

62  of SNPs increases. Stochastic search algorithms are often used to perform inference on

63  posterior distributions. For example, piMASS[8] and BVS[14,15] both use Markov chain

64  Monte Carlo (MCMC) algorithm to search through the model space, while the later can

65  also incorporate external annotations as prior information in the Bayesian model

66  selection to further prioritize causal SNPs.

67      In addition to analyses performed on individual-level data, methods for fine-

68  mapping using only summary statistics from GWAS are becoming more widely applied

69  [16-19]. In general, these methods use reference samples to estimate the correlations

70  between SNPs and then integrate the correlation structure with modified marginal SNP

71  summary statistics in a multivariable regression framework to approximate the

72  corresponding individual-level analysis. Differences between methods are due to

73     variations in the assumptions for residual error and algorithms for model selection[16-19].

74     For example, FINEMAP[16] places a Gaussian prior for causal effect estimates and

75     adopts a shotgun stochastic search algorithm to prioritize the search to a set of most

76     likely important causal configurations. The original implementation of the Joint Analysis

77     of Marginal SNP Effects (JAM)[19] invokes a Cholesky transformation on the linear

78     regression likelihood, adopts a $g$-prior for effect estimates, and then implements a

79     computationally efficient reversible jump MCMC stochastic search algorithm.

80        Leveraging the information across multiple ethnic groups or ancestry populations

81     can enhance the power of fine-mapping[20-22]. Different ancestry groups may have

82     distinct LD structures due to different evolutionary and migration histories[23,24]. For

83     example, compared to non-African Americans, African Americans have smaller LD

84     blocks with weaker correlations as the number of recombination events for each region

85     is expected to be higher[25]. If a true causal variant exists across populations, its

86     corresponding estimated association across populations should be more consistent

87     than the estimated association for proxy SNPs with different LD across populations[26-28].

88     Therefore, integrating the difference in the LD structures across populations can

89     potentially narrow the credible set that a causal variant resides in and improve the

90     resolution of the fine-mapping[29,30].

91        Here, we present an extension of the single-population fine-mapping through

92     JAM to a multi-population setting by fitting a multi-SNP joint model, "mJAM". mJAM

93     assumes that the true causal variant(s) share the same effect across ancestry groups

94     and it explicitly accounts for different LD structures across ancestry groups in the joint

95     model. The mJAM likelihood allows for different feature selection procedures to be

96    performed on summary statistics obtained from multiple populations. This includes

97    Bayesian variable selection approaches that also yield credible sets or more

98    conventional approaches for only selecting certain SNPs. When combined with

99    approaches that only select specific SNPs, mJAM conditional models can further be

100   used in a mediation type framework to construct credible sets for these index variants in

101   a multi-population analysis. We illustrate this flexibility with two computationally efficient

102   implementations of mJAM: "mJAM-SuSiE" for Bayesian variable selection with native

103   SuSiE credible sets, and "mJAM-Forward" for frequentist forward selection of index

104   SNPs and subsequent credible set construction. Through simulation studies with

105   realistic effect size and various patterns of LD, we compare mJAM-SuSiE and mJAM-

106   Forward with other multi-population approaches, including the most commonly used

107   fixed-effect meta-analysis, COJO[17] with pooled LD structure and meta-analyzed

108   summary statistics, and MsCAVIAR[31], a Bayesian fine-mapping approach that allows for

109   an arbitrary number of causal variants in a region. We then applied these methods to

110   three known regions for prostate cancer to demonstrate the practical advantages of

111   mJAM.

# Material and methods

### Multi-population JAM

114        To simplify notation and without loss of generality, we consider the scenario with

115   three populations. Within each population and for a given set of $p$ SNPs within each

116   region, a linear phenotypic model is used.

$$y^{(i)} = G^{(i)}\beta_{global} + \epsilon, \text{ for } i = 1, 2, 3 \tag{1}$$

117    where $y^{(i)}$ is a $N^{(i)} \times 1$ vector of mean-centered phenotypic trait values for the $i^{th}$

118    population, with $N^{(i)}$ being the sample size of the $i^{th}$ population; $G^{(i)}$ is a $N^{(i)} \times p$ matrix

119    of individual-level genotype data for the $i^{th}$ group, where each SNP has been centered

120    to its mean; $\beta_{global} \in \mathbb{R}^P$ denotes the joint effect of the given set of $p$ SNPs. $\epsilon \sim N(0, \sigma^2)$

121    where $\sigma^2$ is the residual variance. It is assumed that all three populations share the

122    same joint effect size, i.e., $\beta_{global}$, and the same residual variance.

123        Akin to a meta-regression, a second-stage model describes the relationship

124    between the population joint effect estimates and the underlying true effect:

$$\begin{pmatrix} \widehat{\beta}^{(1)} \\ \widehat{\beta}^{(2)} \\ \widehat{\beta}^{(3)} \end{pmatrix} = \begin{pmatrix} I_P \\ I_P \\ I_P \end{pmatrix} \beta_{global} + \delta \tag{2}$$

125    where $\delta \sim N(0, \tau^2)$ and $\widehat{\beta}^{(i)} \in \mathbb{R}^P$ is the vector of estimated joint SNP effects for the $i^{th}$

126    population.

127        Equation ( 1 ) and ( 2 ) together form a two-stage model when individual-level

128    data are available. The first stage is three separate linear phenotypic models whereas

129    the second stage fits a fixed-effect meta-analysis model that combines all populations

130    together. By replacing the $\beta_{global}'s$ in ( 1 ) with ( 2 ), we have the following linear fixed-

131    effect model that incorporates the individual-level data of all populations:

$$\begin{pmatrix} y^{(1)} \\ y^{(2)} \\ y^{(3)} \end{pmatrix} = \begin{pmatrix} G^{(1)} & 0 & 0 \\ 0 & G^{(2)} & 0 \\ 0 & 0 & G^{(3)} \end{pmatrix} \begin{pmatrix} I_P \\ I_P \\ I_P \end{pmatrix} \beta_{global} + \epsilon' \tag{3}$$

132

133    With summary data in which only the marginal effect sizes and their standard errors are

134    available, it is also possible to estimate the joint effect size, $\boldsymbol{\beta}_{global}$, with an additional

135    reference sample that estimates the LD between the SNPs[2,32]. Thus, Equation (3) can

136    be used with only GWAS summary statistics with a modified mJAM likelihood after

137    linear transformation:

$$\boldsymbol{G}_c\boldsymbol{I}_c\boldsymbol{y}_c \sim MVN\left(\left((\boldsymbol{G}_c\ \boldsymbol{I}_c)'\ \boldsymbol{G}_c\ \boldsymbol{I}_c\right)\boldsymbol{\beta}_{global},\ \sigma^2\left((\boldsymbol{G}_c\ \boldsymbol{I}_c)'\ \boldsymbol{G}_c\ \boldsymbol{I}_c\right)\right) \qquad (4)$$

138    where $\boldsymbol{G}_c, \boldsymbol{y}_c, \boldsymbol{I}_c$ denotes $\begin{pmatrix} \boldsymbol{G}^{(1)} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \boldsymbol{G}^{(2)} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \boldsymbol{G}^{(3)} \end{pmatrix}, \begin{pmatrix} \boldsymbol{y}^{(1)} \\ \boldsymbol{y}^{(2)} \\ \boldsymbol{y}^{(3)} \end{pmatrix}$ and $\begin{pmatrix} \boldsymbol{I}_P \\ \boldsymbol{I}_P \\ \boldsymbol{I}_P \end{pmatrix}$ in Equation ( 3 )

139    respectively. By expanding each matrix, we have $\boldsymbol{G}_c\boldsymbol{I}_c\boldsymbol{y}_c = \sum_{i=1}^{3}\boldsymbol{G}^{(i)}{}'\ \boldsymbol{y}^{(i)}$ and

140    $(\boldsymbol{G}_c\ \boldsymbol{I}_c)'\ \boldsymbol{G}_c\ \boldsymbol{I}_c = \sum_{i=1}^{3}\boldsymbol{G}^{(i)}{}'\ \boldsymbol{G}^{(i)}$ where $\boldsymbol{G}^{(i)}{}'\ \boldsymbol{G}^{(i)}$ and $\boldsymbol{G}^{(i)}{}'\ \boldsymbol{y}^{(i)}$ are population-specific

141    statistics and can be estimated by population-specific GWAS summary statistics and a

142    reference genotype matrix or LD matrix. Detailed derivation can be found in

143    Supplemental Methods.

144    ***Index SNP Selection and Credible Set Construction for Fine Mapping***

145    mJAM establishes a multi-SNP model within each population with corresponding

146    population-specific LD, while jointly estimating a fixed-effects summary estimate of

147    effect. The mJAM likelihood presented in Equation ( 4 ) can be used in a wide variety of

148    existing feature selection approaches which are applicable to the mJAM statistics

149    shown in Equation ( 4 ). Possible approaches for index SNP selection in mJAM includes

150    stepwise selection[2], Ridge regression[7], and Bayesian approaches such as SuSiE[11].

151      We adopt a forward selection approach based on conditional *P*-value for index

152    SNP selection because of its computational efficiency and straightforward interpretation.

153    We define our implementation of "mJAM-Forward" as a two-step approach in which a

154    first step relies on a conventional stepwise forward selection to select an additional

155    index SNP based on its corresponding *P*-value from a mJAM model conditional on any

156    previous index SNP(s). We incorporate a *g*-prior to stabilize effect estimates[33]. To avoid

157    fitting models with highly correlated SNPs we include a pruning process within Algorithm

158    1.

159      The second step for mJAM-Froward is to define a multi-population credible set

160    for each index SNP. Here, we fit two mJAM models for each candidate credible set

161    SNP, *W,* located within a region of an index SNP, *X.* These models demonstrate that

162    the candidate credible set SNP is: 1) associated with disease marginally, and 2) that the

163    index SNP mediates the effect of the candidate SNP on the disease. The first model

164    includes *W* by itself to yield a probability that *W* is associated with the trait. This model

165    also provides a posterior distribution for the marginal effect for the candidate credible

166    set SNP.

$$Pr(M_W \mid Data) = \frac{p(M_W)BF[M_W:M_{Null}]}{\sum_w p(M_W)BF[M_W:M_{Null}]} \qquad (5)$$

167    where $p(M_W)$ is the prior density of one-SNP model that includes *W* and $BF[M_W:M_{Null}]$

168    is the Bayes factor of one-SNP model with *W* to the null model. See Supplemental

169    Methods for detailed expression of $BF[M_W:M_{Null}]$ with the incorporation of a *g*-prior of

170    the effect estimates. The second model conditions on the index SNP, *X,* to obtain a

171    posterior estimate for an adjusted effect estimate for the credible set SNP. Borrowing

172    from a mediation framework[34], we then calculate the probability that the index SNP

173    mediates the candidate credible set SNP effect using the two models (Figure 1).

$$Pr(Mediation|Data) = Pr(|\tau_W - \tau'_W| > 0|Data)$$

$$(6)$$

174    where $\tau_W$ is the total effect of the candidate credible set SNP on the outcome and $\tau'_W$ is

175    the direct effect. A strong mediation effect indicates that the observed marginal effect of

176    the candidate credible set SNP on the outcome is mainly due to its indirect effect

177    through its strong correlation with the index SNP, and not due to a direct effect on the

178    outcome. These two model probabilities are then combined to calculate the probability

179    that a candidate SNP is a credible set SNP, Posterior Credible Set Probability (PCSP).

$$PCSP_W := Pr(M_W \mid Data) \cdot Pr(|\tau_W - \tau'_W| > 0|Data)$$

$$(7)$$

180    PCSP are then scaled over all SNPs in the region and used to define a 95% credible set

181    of cross-population SNPs.

182    **Algorithm 1 Pseudo algorithm for fitting mJAM-Forward and credible set**

183    **construction in a region**

---

Input data: $\widehat{\boldsymbol{\beta}}^{(i)}$, $se(\widehat{\boldsymbol{\beta}}^{(i)})$, sample size of GWAS $N_{GWAS}$, Effect Allele Frequencies $\mathbf{EAF}^{(i)}$, $\boldsymbol{G}_R^{(i)}$ for each study indexed by $i$

Input arguments: LD threshold for index SNP selection, conditional *P*-value threshold

1. Compute mJAM statistics $(\boldsymbol{G}_c\,\boldsymbol{I}_c)'\boldsymbol{G}_c\,\boldsymbol{I}_c$, $(\boldsymbol{G}_c\,\boldsymbol{I}_c)'\boldsymbol{y}_c$, and $\boldsymbol{y}_c'\boldsymbol{y}_c$
2. Compute marginal mJAM *P*-values under $g$ prior specification for all testing SNPs in the region
3. While the smallest conditional *P*-value (or marginal *P*-value in the first round only) is smaller than threshold:
4.     Identify the testing SNP with the smallest conditional *P*-value as the next index SNP
5.     Construct credible set of the new index SNP
6.     Prune out SNPs in LD with the new index SNP based on LD threshold

> 7.    Compute the conditional mJAM *P*-value for all remaining SNPs in the region
> 8.  Stop until no SNP in the region has conditional p-value smaller than threshold
>
> Return index SNP(s), corresponding conditional *P*-value(s) and credible set(s).

184    We also integrate the mJAM likelihood and summary statistics into a Bayesian

185    selection method that indicates index SNPs and simultaneously estimates credible set

186    SNPs, "mJAM-SuSiE" (See Supplemental Methods for the pseudo-algorithm of fitting

187    mJAM-SuSiE) [35].

### *Incorporating missing variants in mJAM*

189    In genetic association studies with more than one cohort or study, it is common

190    that a particular SNP might be available in some studies but missing in the others[36]. A

191    notable practical feature of the mJAM framework is that it allows for these SNPs with

192    missing information to be analyzed without being filtered or removed. This is

193    accomplished with a simple modification by substituting a value of zero in the identity

194    matrix in Equation ( 3 ) and ( 4 ). Such modification then allows for observed statistics

195    from other populations to be used but removes the contribution from the population in

196    which it is missing but does not alter the algorithm nor the fitting process. This

197    modification is applicable either when the SNP is missing in the reference panel or

198    when the population-specific GWAS summary statistics are not available for the SNP

199    (See Supplemental Methods for more details).

### *Simulation Study on Structured LD*

201    We conducted a simulation study to compare the performance of the two mJAM

202    implementations (mJAM-SuSiE and mJAM-Forward) with three commonly used

203    alternative approaches: fixed-effect meta-analysis, COJO stepwise selection and

204    MsCAVIAR. Fixed-effect meta-analysis takes an inverse-variance weighted average of

205    the marginal estimates from individual studies or populations. COJO approximates the

206    conditional and joint effect from summary statistics and single reference LD and then

207    implements a stepwise selection based on conditional *P*-values. Additionally, for use of

208    COJO on multiple populations, the summary-level statistics come from the fixed-effect

209    meta-analysis across all populations and the reference LD can be obtained from either

210    the pooled individual-level genotype data or a subset of meta-analysis sample. We used

211    the former as the reference LD for COJO in our simulations. MsCAVIAR is built upon a

212    Bayesian multivariate normal framework first described as CAVIAR[10] to account for

213    between-study or between-population heterogeneity using a random-effects model. To

214    compare the performance of index SNP selection across these multi-population

215    approaches, we use three metrics: number of selected index SNPs, sensitivity/power,

216    and positive predictive value (PPV). In addition, since MsCAVIAR, mJAM-SuSiE, and

217    mJAM-Forward output credible set(s) within each region, we compare credible set

218    performance using the number of credible set(s), size of each credible set,

219    sensitivity/power, PPV and empirical coverage. For the two non-Bayesian methods, FE

220    and COJO, we consider the group of SNPs with meta-analyzed or conditional *P*-values

221    less than a Bonferroni-corrected significance level as a single credible set for the

222    purpose of performance comparison.

223        We performed two sets of scenarios: 1) simulated correlation structures with the

224    same block LD structures across populations; and 2) simulated correlation based on

225    real genetic correlation structures observed in the study cohort from Elucidating Loci

226    Involved in Prostate Cancer Susceptibility (ELLIPSE) OncoArray Consortium[21]. For the

227    first set of scenarios, we first simulated a baseline scenario where each population has

228    3 individual association studies with N = 5,000 each to closely represent the real-life

229    situation where there are multiple association studies performed for each ethnic group

230    (total sample size = 5,000 × 3 studies/population × 3 population = 45,000). A total of 50

231    SNPs are simulated in 5 blocks of 10 SNPs. Within each block of 10 SNPs, the pair-

232    wise correlations are uniformly set to a constant value $r^2$ across ancestries for simplicity.

233    $r^2$ varies from 0, $0.6^2$ and $0.9^2$ to represent independent, moderate LD and high LD

234    scenarios. Corresponding LD heatmaps are shown in Figure S1. We then selected a

235    single causal SNP with an effect size of 0.03 for a standard normal outcome. The

236    baseline scenario was extended by varying parameters, including the ratio of sample

237    sizes between each population, levels of LD, the total number of causal SNPs and

238    corresponding effect sizes.

239    **_Simulation Study with Real Data_**

240          To better capture realistic LD patterns, we performed simulations based on real

241    correlation within three ancestry groups (Europeans, African Americans, and East

242    Asians) from the ELLIPSE OncoArray Consortium[21]. The available sample sizes for

243    these three populations are 93,749 Europeans, 9,531 African Americans, and 2,075

244    Asians. We simulated 120 SNPs within a 1334 kb region from chromosome 2 using a

245    multivariate normal model with an estimated correlation structure from individual-level

246    genotypes. The heatmap of this region for each ethnicity is shown in Figure S2. In each

247    simulation, we randomly chose one SNP out of a selected LD block to be the causal

248    SNP with effect size being 0.04, resulting in an empirical average -log10(_P_-value) of the

249    most significance variant of 7.75 ($P$-value $\approx 1.8 \times 10^{-8}$) averaged across 500

250    simulations.

251    ***Applied examples***

252        To illustrate mJAM on real data, we applied the methods on three regions using

253    summary statistics from the latest cross-ancestry prostate cancer association study[37]

254    across four ancestry groups, including 122,188 prostate cancer cases and 604,640

255    controls of European ancestry, 19,391 cases and 61,608 controls of African ancestry,

256    10,809 cases and 95,790 controls of East Asian ancestry, and 3,931 cases and 26,405

257    controls from Hispanic populations. Within each region, we applied mJAM-Forward to

258    select index SNP(s) using population-specific summary statistics and reference dosage

259    for each population. Then we constructed mJAM credible set(s) by including top SNPs

260    ranked by their mJAM posterior probabilities until those SNPs included in the credible

261    set reached a cumulative posterior probability of 95%. Reference dosage were obtained

262    from the Prostate Cancer Association Group to Investigate Cancer-Associated

263    Alterations in the Genome and Collaborative Oncological Gene-Environment Study

264    Consortium [PRACTICAL iCOGS], the Elucidating Loci Involved in Prostate Cancer

265    Susceptibility OncoArray Consortium [ELLIPSE OncoArray], the African Ancestry

266    Prostate Cancer Consortium [AAPC GWAS], GWAS of prostate cancer in Latinos

267    [LAPC GWAS] and Japanese [JAPC GWAS][21]. Results from mJAM-Forward are

268    compared with those from mJAM-SuSiE, COJO and MsCAVIAR.

# Results

### *Simulation Study on Artificial LD*

Under the baseline scenario (50 SNPs in total, 1 causal SNP with an effect size

of 0.03, 3 studies per population, and balanced sample size across populations), the

95% credible sets from mJAM-Forward, mJAM-SuSiE, and MsCAVIAR were well

calibrated to the specified coverage level (Figure S3). Both mJAM-Forward and mJAM-

SuSiE preserved relatively high sensitivity in terms of including the true causal SNP in

its credible set (sensitivity = 0.86 and 0.64 respectively, Figure 2A). Although

MsCAVIAR had the highest sensitivity (0.99) under the baseline scenario, its average

credible set size was much larger (9.47 for MsCAVIAR; 2.12 for mJAM-Forward and

0.78 for mJAM-SuSiE, Figure 2C), thus leading to a much lower PPV (0.39, Figure 2B).

mJAM-SuSiE had the highest PPV (0.89) among the methods we compared, meaning

that it had the highest proportion of true causal over the total number of credible set

SNPs on average, followed by mJAM-Forward (0.58) (Figure 2B). In terms of credible

set sensitivity, PPV and average CS size the methods had similar patterns of

performance for scenarios expanded beyond the baseline to various LD structures,

imbalanced sample size across populations, and 3 causal SNPs (Figure S4).

In terms of identifying the true causal variant as an index SNP (i.e. sensitivity),

mJAM-Forward and MsCAVIAR had the best performance under moderate LD

scenarios (Figure 3A) with a sensitivity was 0.73, and 0.72 respectively. However, for

these two methods, mJAM-Forward had a better PPV was 0.81, compared to

MsCAVIAR (0.72). In comparison, mJAM-SuSiE had poor sensitivity (0.62) but a higher

291 PPV (0.88). COJO had a similar performance with mJAM-Forward under independent

292 LD scenarios but its sensitivity and PPV worsened compared to mJAM-Forward as the

293 level of LD increases. Though COJO performs a similar stepwise selection as mJAM-

294 Forward, unlike mJAM-Forward that specifically accounted for population-specific LD,

295 COJO uses meta-analyzed marginal summary statistics and pooled LD panel which

296 makes it difficult to identify the true common variants through disentangling the

297 population-specific LD structure. All methods selected on average 1 index SNP among

298 500 simulations, close to the true number of causals (Figure 3C).  For MsCAVIAR pre-

299 specification is required so we set the value to 1 for all scenarios. Notably for practical

300 implementation, for a small number of scenarios (40%), mJAM-SuSiE did not select any

301 index SNP under independent or moderate LD scenarios, leading to relatively low

302 sensitivity compared other methods when averaged over replicates (Figure 3A).

303 When the pairwise correlation within each LD block increased, the average

304 credible set sizes for all methods increased correspondingly (Figure 2C). As a result,

305 under high LD scenarios, the PPV of identifying the true causal(s) in a credible set

306 decreased to a noticeable extent for MsCAVIAR, mJAM-Forward, and FE (Figure 2B).

307 Though mJAM-Forward's PPV dropped due to the increase in credible set sizes on

308 average, mJAM-Forward was still able to retain a sensitivity of 0.88 under the high LD

309 scenario. mJAM-SuSiE achieved the highest PPV (0.81, Figure 2B) among all methods

310 under high LD scenarios while retaining relatively high sensitivity and small credible set

311 size. However, mJAM-SuSiE's sensitivity was relatively low compared to mJAM-

312 Forward and MsCAVIAR under independent or moderate LD scenarios (Figure 2A).

313    Despite of mJAM-SuSiE's outstanding performance under high LD scenarios with

314    moderate causal effect size, we noticed that its results were very sensitive to the

315    marginal significance of the true causal SNPs. To represent a real-life situation where a

316    lead variant within a region has an extremely significant marginal $P$-value, we expanded

317    the baseline scenario with 1 true causal SNP to additional scenarios with increasing

318    significance of the true causal SNP, where the average -log10($P$-value) of the true

319    causal ranges from 5 to 263 (mimicking significance often found in applied GWAS).

320    Under increasingly high-power scenarios, mJAM-Forward consistently selected 1

321    credible set regardless of the significance of the true causal whereas the average

322    number of credible sets by mJAM-SuSiE increased as the statistical significance (i.e.

323    effective power) increased (Figure 4B). As a result, mJAM-SuSiE selected more false

324    positive SNPs within the credible sets when the true causal SNP has high observed

325    marginal significance. In addition, the empirical coverage of mJAM-SuSiE's credible

326    sets dropped below the expected level quickly after the true causal SNP became more

327    significant (Figure 4A). In contrast, mJAM-Forward's credible sets remained well-

328    calibrated.

329    To explore the impact of two types of missingness on the performance of mJAM-

330    Forward, we modified our simulation studies with artificial LD structure to include a

331    missing SNP in LD with the causal SNP, or with the missing SNP as the causal SNP

332    itself. The flexibility of mJAM likelihood (Equation 2) allows us to incorporate SNPs with

333    missing information in some studies or populations in the analysis. We found that when

334    the missing SNP is in LD with the causal SNP, mJAM-Forward has stable performance

335    in comparison to when there is no missingness (Figure S6).  When the causal SNP is

336 missing, mJAM-Forward still preserves the power both to select the causal SNP as the

337 index SNP and to include the causal SNP in its credible set.

### Simulation Study on Real LD

339 When applied to the simulated data on the 120-SNP region on chromosome 2,

340 mJAM-Forward, mJAM-SuSiE and MsCAVIAR selected on average around 1 index

341 SNP whereas COJO selected 1.5 index SNPs, indicating a slight increase in false

342 positive signals. mJAM-Forward had highest sensitivity and PPV of identifying the true

343 causal from a complicated LD structure as an index SNP (Table 1). In terms of credible

344 set performance, MsCAVIAR demonstrated high empirical coverage of its credible set

345 as well as high sensitivity compared to the other two mJAM methods. However, such

346 high sensitivity and PPV was achieved at the cost of a much larger size for the credible

347 sets. The average size of the 95% CS of MsCAVIAR is 56.52, even larger than the

348 number of SNPs that reached marginal genome-wide significant ($5 \times 10^{-8}$) in a fixed-

349 effect meta-analysis (48.88). On the other hand, the average credible set size for

350 mJAM-Forward and mJAM-SuSiE was 19.70 and 18.37 respectively. Meanwhile, both

351 approaches preserved reasonably high sensitivity and empirical coverage.

### Applied example 1: a single-hit region on chromosome 12

353 The first applied example is a 1013 kb region on chromosome 12 which consists

354 of 276 SNPs with a marginal meta-analyzed *P*-value $< 10^{-3}$ and minor allele frequency

355 (MAF) > 2%. Figure S7 shows the LD structure for the four ancestry groups in this

356 analysis. None of the SNPs in this region reached genome-wide significance in any

357 population-specific analyses (Figure 5B) but after multi-population meta-analysis 48

358    SNPs are genome-wide significant (Figure 5A).  By setting a conditional *P*-value

359    threshold at $5 \times 10^{-8}$, mJAM-Forward identified one index SNP at 12:109994870:A:T

360    (meta P-value = $3.5 \times 10^{-10}$) with a corresponding 95% credible set of 41 SNPs. The

361    median $r^2$ between the credible set SNPs with the index SNP is 0.998 for European LD,

362    0.979 for African, 0.990 for Hispanic and 0.996 for East Asian. COJO identified the

363    same index SNP, 12:109994870:A:T. MsCAVIAR reported a slightly larger 95% credible

364    set than mJAM-Forward, consisting of 45 SNPs (Figure S8). The index SNP of

365    MsCAVIAR's credible set is 12:109998097:A:G (meta *P*-value = $3.7 \times 10^{-10}$) whose $r^2$

366    with 12:109994870:A:T is greater than 0.99 in all four ancestry groups. This index SNP,

367    12:109998097:A:G, is included in a mJAM-Forward credible set only when coverage is

368    increased to 99%; whereas the index SNP for mJAM-Forward, 12:109994870:A:T, is

369    included in the 95% MsCAVIAR credible set. mJAM-SuSiE estimates a single 95%

370    credible set with 28 total SNPs and a unique single index SNP, 12:109996343:A:C

371    (meta *P*-value = $2.2 \times 10^{-9}$) which is also included in both credible sets of mJAM-

372    Forward and MsCAVIAR. The median $r^2$ within a credible set is also greater than 0.99

373    for all ancestry groups (Table S1). The index SNP from mJAM-Forward was also

374    included in its credible set (Figure S8B).

375    ***Applied example 2: Asian-driven signals on chromosome 10***

376         As a second example, we conducted an analysis on a chromosome 10 region

377    which consists of 412 SNPs after QC and spans around 1571 kb. Figure S9 shows the

378    LD structure in this region separately for European, African, East Asian, and Hispanic

379    populations. This region contains two clear signals with meta-analyzed *P*-value < $10^{-15}$,

380    which are mainly driven by the results from East Asian and African populations (Figure

381    6). In this example, mJAM-Forward identified two index SNPs, 10:80835998:C:T (meta

382    *P*-value = $9 \times 10^{-21}$ and 10:80238015:C:T (meta *P*-value = $1 \times 10^{-19}$) (Figure 6A). The

383    95% mJAM-Forward credible set for the first index SNP, 10:80835998:C:T, contains 3

384    SNPs in total and there are 45 SNPs in the credible set for the second index SNP. The

385    minimum r$^2$ between the mJAM-Forward credible set SNPs with its own index SNP is no

386    less than 0.95 in European, East Asian and Hispanic populations, and no less than 0.81

387    in African ancestry populations (Table S2). COJO identified two index SNPs,

388    10:80835998:C:T and 10:80240493:A:G. 10:80835998:C:T is the same as one of the

389    index SNPs selected by mJAM-Forward and 10:80240493:A:G is included in the mJAM-

390    Forward 95% credible set of 10:80238015:C:T. Since MsCAVIAR does not support

391    reporting more than one distinctive credible set, we split this region into two adjacent

392    regions and applied MaCAVIAR on these two subregions separately. MsCAVIAR

393    selected the same 3-SNP 95% credible set (Figure S10) with index SNP being

394    10:80835998:C:T, and another 45-SNP credible set with index SNP being

395    10:80238015:C:T where 42 of them are replicated in the mJAM-Forward credible set.

396    mJAM-SuSiE also identified the same 3-SNP credible set (95%) with the same index

397    SNP 10:80835998:C:T but did not identify any credible set around 10:80238015:C:T.

398    Instead, it reported two additional credible sets at 10:80260938:V1 (meta *P*-value =

399    $2 \times 10^{-10}$) and 10:80476778:V1 (meta *P*-value = $4 \times 10^{-4}$) (Figure S10), and the

400    credible set size is 2 and 5 respectively.

401    ***Applied example 3: Secondary signal within 40kb region of a leading SNP***

402           The third applied example illustrates a scenario where there is a secondary

403    signal within close proximity of the leading SNP in a chromosome 11 region. This region

404    spans 335.5 kb and consists of 191 SNPs that passed QC. The population-specific LD

405    structure and Manhattan plot of multi-population meta-analysis results are shown in

406    Figure S11 and Figure 7. The lead variant, 11:102401661:C:T, has a multi-population

407    meta-analyzed $P$-value of $1.5 \times 10^{-38}$ and mJAM-Forward identified a secondary index

408    SNP, 11:102440927:A:G, only 39 kb away which has a meta $P$-value of $4.9 \times 10^{-11}$.

409    The $r^2$ between these two index SNPs is less than 0.01 in all four ancestry groups

410    (Figure S11), suggesting statistical independence between these two SNPs. COJO

411    selected the same primary index SNP, 11:102401661:C:T, and a different secondary

412    index, 11:102433309:A:G, which has a meta $P$-value of $1.3 \times 10^{-7}$ and is highly

413    correlated with 11:102440927:A:G ($r^2$ = 0.79 in EUR; 0.55 in AA; 0.87 in LA and 0.99 in

414    ASN). mJAM-SuSiE also selected two credible sets in this region: the first set has 2

415    SNPs which are both replicated in mJAM-Forward's first credible set; the second set

416    has 26 SNPs where 24 of them are found in mJAM-Forward's second set. However, the

417    index SNP of the second set in mJAM-SuSiE is one with lower marginal significance

418    (meta $P$-value = $6.3 \times 10^{-5}$) compared to mJAM-Forward.

419        Both mJAM-SuSiE and mJAM-Forward are able to identify multiple sets within

420    one region without any pre-defined number of causal variants. On the other hand, the

421    implementation of MsCAVIAR requires users to specify the maximum number of causal

422    variants in a region to enumerate all possible causal configurations. Gauging the

423    possible number of causal variants can be difficult when secondary signals are located

424    close to the lead variant. In this example, the secondary signal is located only 39 kb

425    away from the leading variant, and visual inspection of the Manhattan plot (Figure 7)

426    suggests only one peak. Even if we specify the number of causal variants to be two

427     when applying MsCAVIAR to this region, MsCAVIAR reports only one credible set such

428     that the posterior probability of this set containing 2 causal variants is at least 0.95.

429     Thus, it becomes difficult to separate the selected credible set SNPs into two distinctive

430     groups. When the number of causal variants is set to two, MsCAVIAR selected 24

431     SNPs among which the 2 SNPs with highest posterior probability are 11:102401661:C:T

432     and 11:102396607:C:T (Figure S12). However, these two SNPs are in high LD and thus

433     are likely linked to a single underlying causal signal and not indicative of multiple

434     independent signals.

## Discussion

436        As integrating studies from ancestrally diverse populations may increase power

437     to detect novel variant and improve fine-mapping resolution[22,38,39], we extend our

438     previous single-population fine-mapping through JAM to a multi-population approach,

439     mJAM. mJAM requires only population-specific summary statistics and population-

440     specific reference LD panels, which are more accessible than individual-level data to

441     many researchers. mJAM explicitly incorporates the different LD structures across

442     populations to yield conditional estimates of SNP effects from a single joint model. The

443     mJAM framework can be used to first select index SNPs using existing feature selection

444     approaches, such as forward stepwise selection[2], Bayesian model selection[8,9,11], or

445     regularized regression[6,7]. To demonstrate this flexibility, we have implemented mJAM

446     through two implementations of feature selection: mJAM-SuSiE (a Bayesian approach)

447     and mJAM-forward selection. We also combine the forward selection implementation

448     with a second step to identify credible set SNPs. This step works given any set of index

449   SNPs within a region by estimating a posterior credible set probability (PCSP) for a SNP

450   defined as a combination of two component probabilities: one models the marginal

451   association between the candidate SNP and the outcome; the other models the

452   mediation effect of the index SNP on the candidate SNP, borrowing from a mediation

453   framework. These PCSPs are then used to construct credible sets. The closed-formed

454   expression for PCSP allows computational efficient construction of credible sets,

455   compared to other Bayesian approaches that often use computationally intensive

456   algorithms to obtain posterior distributions. It also allows credible set construction from

457   any index SNP list allowing researchers to apply other feature selection methods or use

458   existing lists or knowledge to determine index SNP.

459   The two-stage model framework utilized in mJAM builds upon previous work

460   highlighting the use of hierarchical JAM (hJAM)[40], an approach for the joint analysis of

461   marginal summary statistics that incorporates a prior information matrix. This matrix

462   characterizes the SNPs and can include information such as SNP effects on gene

463   expression analogous to TWAS or on intermediates biomarkers analogous to Mendelian

464   randomization. mJAM is an extension to hJAM in that it replaces the prior information

465   matrix in hJAM with a stacked identity matrix, $\begin{pmatrix} I_P \\ I_P \\ I_P \end{pmatrix}$, as described in Methods section.

466   The stacked identify matrix can be interpreted as our prior believe on the joint SNP

467   effect estimates that all populations share the same true effect sizes.

468   In a set of realistic simulation settings, both mJAM implementations

469   demonstrated the ability to infer the number of independent signals within a region, to

470   differentiate signals from noise, and to achieve a sufficient level of sensitivity while

471  preserving high fine-mapping resolution through small-sized credible sets. We also

472  investigated the impact of imbalanced sample size across populations on model

473  performance and demonstrated that all methods showed a similar decrease in terms of

474  sensitivity and PPV when the sample size is imbalanced but the total sample size

475  remains constant (Figure S4).  mJAM is described using three populations in simulation

476  studies and we apply mJAM to real data with four distinct populations. In practice,

477  mJAM can be used to analyze a large number of studies or population-specific

478  summary statistics facilitating flexibility in application. Thus, analyses do not need to be

479  limited to aggregating continental ancestry populations, but can include numerous, more

480  specific ancestry appropriate reference panels to aggregate data across many studies

481  (Figure S5). However, as with all summary statistic approaches that rely on reference

482  panels, the ability to disentangle highly correlated SNPs will be driven by the sample

483  sizes[41] and LD within and between the reference panels used[42]. In addition, another

484  practical limitation to many summary statistics-based approaches is the requirement for

485  complete summary statistics and refence data for all SNPs across all studies and

486  populations analyzed[36]. Missingness can be due to the difference in genotyping arrays

487  used by different studies, or rare variants not being captured due to limited sample size

488  in certain studies. Filtering too many variants might be dangerous because as less

489  information is used to disentangle the LD structure within each region and potentially

490  missing the causal variant. An important feature of mJAM is that it will work even in the

491  presence of differential missingness across studies or populations utilizing all

492  information that is available.

493    In the simulation study with artificial LD structures, mJAM-SuSiE resulted in

494    outstanding performance under high LD scenarios, achieving both high sensitivity and

495    high PPV. However, as the significance of the causal variant(s) within a region

496    increases, mJAM-SuSiE tends to break down selecting more false positive signals with

497    each in separate credible sets. This results in a substantial decrease in the empirical

498    coverage of mJAM-SuSiE credible sets. In practice, we recommend limiting the

499    application of mJAM-SuSiE to only regions with SNPs with modest marginal statistically

500    significance or to screen for any potential false positive credible sets before interpreting

501    mJAM-SuSiE's credible sets after estimation.

502    We also carried out a case study of prostate cancer where mJAM is applied to

503    several prostate cancer susceptible regions. Through three different regions with

504    different characteristics in number of estimated independent signals and underlying LD

505    within and between populations, we demonstrated the practical advantages of mJAM-

506    Forward, including allowing more than one causal variant within a region, outputting

507    individual credible sets corresponding to each index, and easily interpretable index

508    variants with conditional estimates. In addition to the three applied examples shown

509    here, mJAM has been applied to perform index variants selection across all regions in

510    the latest multi-population prostate cancer GWAS[37] which is currently under review.

511    For all approaches that use marginal summary statistics and reference data,

512    careful consideration and construction of the correlation matrices is important. This

513    includes using a reference panel with ancestry and LD that matches the population in

514    which the original marginal summary statistics were estimated[41,43]. The methods also

515    require that the correlation matrix used is full rank and positive-definite which is often

516  driven by the sample size of the data and the frequency of the SNPs. For mJAM such

517  consideration must be considered across all populations used in the analysis. Firstly, for

518  rare variants, mJAM estimates of multi-population effect and standard errors that can be

519  different from the marginal meta-analyzed estimates which use inverse-variance

520  weighting. mJAM estimation from summary statistics assume Hardy-Weinberg

521  equilibrium which some variants, especially rare variants, might not satisfy. In addition,

522  many rare variants will also have large effect sizes and large standard errors from the

523  population-specific summary statistics thus resulting in more uncertainty in multi-

524  population analysis compared to variants that are common across all populations.

525  Secondly, in regions with extremely significant lead variants from a well-powered

526  GWAS, even small degrees of LD can pull the marginal and conditional effect estimates

527  of other variants away from the null. Thus, false positive signals might be selected if we

528  apply the same threshold for index SNP selection and LD pruning. For such regions,

529  researchers may consider setting a higher significance threshold for secondary signal

530  selection and a more stringent LD threshold for pruning out correlated signals.

531  In conclusion, mJAM offers a flexible and efficient modeling framework for multi-

532  population fine-mapping that first selects index variants and then constructs credible

533  sets. One key assumption in mJAM is that causal variants and their effect sizes are

534  similar across all populations and there exists evidence suggesting that common causal

535  variants tend to have consistent effect sizes across populations[26-28]. In future research,

536  we plan to relax the current mJAM assumption to allow different true effect sizes across

537  populations. Other potential future directions include follow-up functional analyses

538  based on mJAM credible sets and polygenic risk score models based on mJAM fine-

539    mapped results. mJAM is currently available as a R package for fine-mapping of

540    specific regions and can easily be adapted for genome-wide applications.

541

542  **Tables**

543  **Table 1 Comparison of model performance on data simulated from real LD**

544  **structure.**

| | | Method | | | | |
|---|---|---|---|---|---|---|
| | | mJAM-Forward | mJAM-SuSiE | FE | COJO | MsCAVIAR |
| Credible Set Performance | Sensitivity[a] | 0.930 | 0.910 | 0.972 | - | 0.994 |
| | PPV[b] | 0.064 | 0.069 | 0.024 | - | 0.022 |
| | CS size[c] | 19.70 | 18.37 | 48.88 | - | 56.62 |
| | CS coverage[d] | 0.934 | 0.940 | - | - | 1.000 |
| Index SNP Performance | Sensitivity[e] | 0.218 | 0.174 | - | 0.186 | 0.134 |
| | PPV[f] | 0.219 | 0.021 | - | 0.144 | 0.134 |
| | Number of selected index | 1.00 | 0.97 | - | 1.51 | 1.00 |

545  Abbreviations: FE, fixed-effect meta-analysis; CS, credible set, PPV, positive predictive value.

546  [a] proportion of true causal SNPs being selected in a credible set, averaged over 500 simulations

547  [b] proportion of true causal SNPs over the total number of selected credible set SNPs, averaged over
548  500 simulations.

549  [c] total number of SNPs included in a credible set, averaged over all 95% credible sets in 500
550  simulations.

551  [d] proportion of 95% credible sets in 500 simulations that included at least one true causal SNP.

552  [e] proportion of true causal SNPs being selected as an index SNP, averaged over 500 simulations.

553  [f] proportion of true causal SNPs over the total number of selected index SNPs, averaged over 500
554  simulations.

555

556

557  **Declaration of interests**

558    The authors declare no competing interests.

## Acknowledgements

## Data and code availability

563    Both mJAM-Forward and mJAM-SuSiE are available as an R package at

564    https://github.com/USCbiostats/hJAM/R. The codes used for simulations and real data

565    examples are available at https://github.com/USCbiostats/hJAM/manuscript_codes.

## References

567    1.    Buniello, A., MacArthur, J.A.L., Cerezo, M., Harris, L.W., Hayhurst, J., Malangone, C., McMahon,
568          A., Morales, J., Mountjoy, E., and Sollis, E. (2019). The NHGRI-EBI GWAS Catalog of published
569          genome-wide association studies, targeted arrays and summary statistics 2019. Nucleic acids
570          research *47*, D1005-D1012.
571    2.    Schaid, D.J., Chen, W., and Larson, N.B. (2018). From genome-wide associations to candidate
572          causal variants by statistical fine-mapping. Nature Reviews Genetics *19*, 491-504.
573    3.    Visscher, P.M., Brown, M.A., McCarthy, M.I., and Yang, J. (2012). Five years of GWAS discovery.
574          The American Journal of Human Genetics *90*, 7-24.
575    4.    Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. Journal of the Royal
576          Statistical Society: Series B (Methodological) *58*, 267-288.
577    5.    Cho, S., Kim, H., Oh, S., Kim, K., and Park, T. (2009). Elastic-net regularization approaches for
578          genome-wide association studies of rheumatoid arthritis. BMC Proceedings *3*, S25.
579          10.1186/1753-6561-3-s7-s25.
580    6.    Ayers, K.L., and Cordell, H.J. (2010). SNP Selection in genome-wide and candidate gene studies
581          via penalized logistic regression. Genetic Epidemiology *34*, 879-891. 10.1002/gepi.20543.
582    7.    Hastie, T. (2020). Ridge Regularization: An Essential Concept in Data Science. Technometrics *62*,
583          426-433. 10.1080/00401706.2020.1791959.
584    8.    Guan, Y., and Stephens, M. (2011). Bayesian variable selection regression for genome-wide
585          association studies and other large-scale problems. The Annals of Applied Statistics, 1780-1815.
586    9.    Chen, W., Larrabee, B.R., Ovsyannikova, I.G., Kennedy, R.B., Haralambieva, I.H., Poland, G.A.,
587          and Schaid, D.J. (2015). Fine Mapping Causal Variants with an Approximate Bayesian Method
588          Using Marginal Test Statistics. Genetics *200*, 719-736. 10.1534/genetics.115.176107.
589    10.   Hormozdiari, F., Kostem, E., Kang, E.Y., Pasaniuc, B., and Eskin, E. (2014). Identifying causal
590          variants at loci with multiple signals of association. Genetics *198*, 497-508.
591          papers3://publication/doi/10.1534/genetics.114.167908.

592    11.    Zou, Y., Carbonetto, P., Wang, G., Stephens, M.V.O.P., and Stephens, V.O.P.M. (2022). Fine-
593           mapping from summary data with the "Sum of Single Effects" model. bioRxiv.
594           https://doi.org/10.1101/2021.11.03.467167.

595    12.    Galarneau, G., Palmer, C.D., Sankaran, V.G., Orkin, S.H., Hirschhorn, J.N., and Lettre, G. (2010).
596           Fine-mapping at three loci known to affect fetal hemoglobin levels explains additional genetic
597           variation. Nature genetics *42*, 1049.

598    13.    Trynka, G., Hunt, K.A., Bockett, N.A., Romanos, J., Mistry, V., Szperl, A., Bakker, S.F., Bardella,
599           M.T., Bhaw-Rosun, L., and Castillejo, G. (2011). Dense genotyping identifies and localizes
600           multiple common and rare variant association signals in celiac disease. Nature genetics *43*,
601           1193-1201.

602    14.    Quintana, M.A., Schumacher, F.R., Casey, G., Bernstein, J.L., Li, L., and Conti, D.V. (2012).
603           Incorporating prior biologic information for high-dimensional rare variant association studies.
604           Human heredity *74*, 184-195.

605    15.    Quintana, M., and Conti, D. (2013). Integrative variable selection via Bayesian model
606           uncertainty. Statistics in medicine *32*, 4938-4953.

607    16.    Benner, C., Spencer, C.C., Havulinna, A.S., Salomaa, V., Ripatti, S., and Pirinen, M. (2016).
608           FINEMAP: efficient variable selection using summary data from genome-wide association
609           studies. Bioinformatics *32*, 1493-1501.

610    17.    Yang, J., Ferreira, T., Morris, A.P., Medland, S.E., Madden, P.A., Heath, A.C., Martin, N.G.,
611           Montgomery, G.W., Weedon, M.N., Loos, R.J., et al. (2012). Conditional and joint multiple-SNP
612           analysis of GWAS summary statistics identifies additional variants influencing complex traits. Nat
613           Genet *44*, 369-375, S361-363. 10.1038/ng.2213.

614    18.    Verzilli, C., Shah, T., Casas, J.P., Chapman, J., Sandhu, M., Debenham, S.L., Boekholdt, M.S.,
615           Khaw, K.T., Wareham, N.J., and Judson, R. (2008). Bayesian meta-analysis of genetic association
616           studies with different sets of markers. The American Journal of Human Genetics *82*, 859-872.

617    19.    Newcombe, P.J., Conti, D.V., and Richardson, S. (2016). JAM: A Scalable Bayesian Framework for
618           Joint Analysis of Marginal SNP Effects. Genet Epidemiol *40*, 188-201. 10.1002/gepi.21953.

619    20.    Mahajan, A., Go, M.J., Zhang, W., Below, J.E., Gaulton, K.J., Ferreira, T., Horikoshi, M., Johnson,
620           A.D., Ng, M.C., and Prokopenko, I. (2014). Genome-wide trans-ancestry meta-analysis provides
621           insight into the genetic architecture of type 2 diabetes susceptibility. Nature genetics *46*, 234-
622           244.

623    21.    Conti, D.V., Darst, B.F., Moss, L.C., Saunders, E.J., Sheng, X., Chou, A., Schumacher, F.R., Al
624           Olama, A.A., Benlloch, S., and Dadaev, T. (2021). Trans-ancestry genome-wide association meta-
625           analysis of prostate cancer identifies new susceptibility loci and informs genetic risk prediction.
626           Nature genetics *53*, 65-75.

627    22.    Asimit, J.L., Hatzikotoulas, K., McCarthy, M., Morris, A.P., and Zeggini, E. (2016). Trans-ethnic
628           study design approaches for fine-mapping. European journal of human genetics *24*, 1330-1336.

629    23.    Kichaev, G., and Pasaniuc, B. (2015). Leveraging functional-annotation data in trans-ethnic fine-
630           mapping studies. The American Journal of Human Genetics *97*, 260-271.

631    24.    Teo, Y.Y., Ong, R.T., Sim, X., Tai, E.S., and Chia, K.S. (2010). Identifying candidate causal variants
632           via trans‐population fine‐mapping. Genetic epidemiology *34*, 653-664.

633    25.    Campbell, M.C., and Tishkoff, S.A. (2008). African genetic diversity: implications for human
634           demographic history, modern human origins, and complex disease mapping. Annu. Rev.
635           Genomics Hum. Genet. *9*, 403-433.

636    26.    Shi, H., Burch, K.S., Johnson, R., Freund, M.K., Kichaev, G., Mancuso, N., Manuel, A.M., Dong, N.,
637           and Pasaniuc, B. (2020). Localizing Components of Shared Transethnic Genetic Architecture of
638           Complex Traits from GWAS Summary Data. The American Journal of Human Genetics *106*, 805-
639           817. 10.1016/j.ajhg.2020.04.012.

640    27.    Zanetti, D., and Weale, M.E. (2018). Transethnic differences in GWAS signals: A simulation study.
641           Annals of Human Genetics *82*, 280-286. 10.1111/ahg.12251.
642    28.    Marigorta, U.M., and Navarro, A. (2013). High Trans-ethnic Replicability of GWAS Results Implies
643           Common Causal Variants. PLoS Genetics *9*, e1003566. 10.1371/journal.pgen.1003566.
644    29.    Hu, Y., Tanaka, T., Zhu, J., Guan, W., Wu, J.H., Psaty, B.M., McKnight, B., King, I.B., Sun, Q., and
645           Richard, M. (2017). Discovery and fine-mapping of loci associated with MUFAs through trans-
646           ethnic meta-analysis in Chinese and European populations. Journal of lipid research *58*, 974-981.
647    30.    Morris, A.P. (2011). Transethnic meta‐analysis of genomewide association studies. Genetic
648           epidemiology *35*, 809-822.
649    31.    Lapierre, N., Taraszka, K., Huang, H., He, R., Hormozdiari, F., and Eskin, E. (2021). Identifying
650           causal variants by fine mapping across multiple studies. PLOS Genetics *17*, e1009733.
651           10.1371/journal.pgen.1009733.
652    32.    Yang, J., Ferreira, T., Morris, A.P., Medland, S.E., Madden, P.A.F., Heath, A.C., Martin, N.G.,
653           Montgomery, G.W., Weedon, M.N., Loos, R.J., et al. (2012). Conditional and joint multiple-SNP
654           analysis of GWAS summary statistics identifies additional variants influencing complex traits.
655           Nature Genetics *44*, 369-375. 10.1038/ng.2213.
656    33.    Goel, P., and Zellner, A. (1986). On assessing prior distributions and Bayesian regression analysis
657           with g-prior distributions. In Bayesian inference and decision techniques: Essays in Honor of
658           Bruno De Finetti, pp. 233-243.
659    34.    Robins, J.M., and Greenland, S. (1992). Identifiability and exchangeability for direct and indirect
660           effects. Epidemiology *3*, 143-155. 10.1097/00001648-199203000-00013.
661    35.    Wang, G., Sarkar, A., Carbonetto, P., and Stephens, M. (2020). A simple new approach to
662           variable selection in regression, with application to genetic fine mapping. Journal of the Royal
663           Statistical Society Series B-Statistical Methodology *82*, 1273-1300. 10.1111/rssb.12388.
664    36.    Jiang, Y., Chen, S., Mcguire, D., Chen, F., Liu, M., Iacono, W.G., Hewitt, J.K., Hokanson, J.E.,
665           Krauter, K., Laakso, M., et al. (2018). Proper conditional analysis in the presence of missing data:
666           Application to large scale meta-analysis of tobacco use phenotypes. PLOS Genetics *14*,
667           e1007452. 10.1371/journal.pgen.1007452.
668    37.    Wang, A., Shen, J., Rodriguez, A., Saunders, E., Chen, F., Darst, B., Sheng, X., Xu, Y., Chou, A.,
669           Benlloch, S., et al. (2022). Improving prostate cancer risk prediction through multi-ancestry
670           genome-wide discovery of 187 novel risk variants. [Manuscript submitted for publication].
671    38.    Mahajan, A., Rodan, R., Aylin, Le, H., Thu, Gaulton, J., Kyle, Haessler, J., Stilp, M., Adrienne,
672           Kamatani, Y., Zhu, G., Sofer, T., Puri, S., et al. (2016). Trans-ethnic Fine Mapping Highlights
673           Kidney-Function Genes Linked to Salt Sensitivity. The American Journal of Human Genetics *99*,
674           636-646. 10.1016/j.ajhg.2016.07.012.
675    39.    Mahajan, A., Spracklen, C.N., Zhang, W., Ng, M.C.Y., Petty, L.E., Kitajima, H., Yu, G.Z., Rüeger, S.,
676           Speidel, L., Kim, Y.J., et al. (2022). Multi-ancestry genetic study of type 2 diabetes highlights the
677           power of diverse populations for discovery and translation. Nature Genetics *54*, 560-572.
678           10.1038/s41588-022-01058-3.
679    40.    Jiang, L., Xu, S., Mancuso, N., Newcombe, P.J., and Conti, D.V. (2021). A Hierarchical Approach
680           Using Marginal Summary Statistics for Multiple Intermediates in a Mendelian Randomization or
681           Transcriptome Analysis. Am J Epidemiol *190*, 1148-1158. 10.1093/aje/kwaa287.
682    41.    Benner, C., Havulinna, A.S., Järvelin, M.-R., Salomaa, V., Ripatti, S., and Pirinen, M. (2017).
683           Prospects of Fine-Mapping Trait-Associated Genomic Regions by Using Summary Statistics from
684           Genome-wide Association Studies. The American Journal of Human Genetics *101*, 539-551.
685           10.1016/j.ajhg.2017.08.012.
686    42.    Peterson, R.E., Kuchenbaecker, K., Walters, R.K., Chen, C.-Y., Popejoy, A.B., Periyasamy, S., Lam,
687           M., Iyegbe, C., Strawbridge, R.J., Brick, L., et al. (2019). Genome-wide Association Studies in
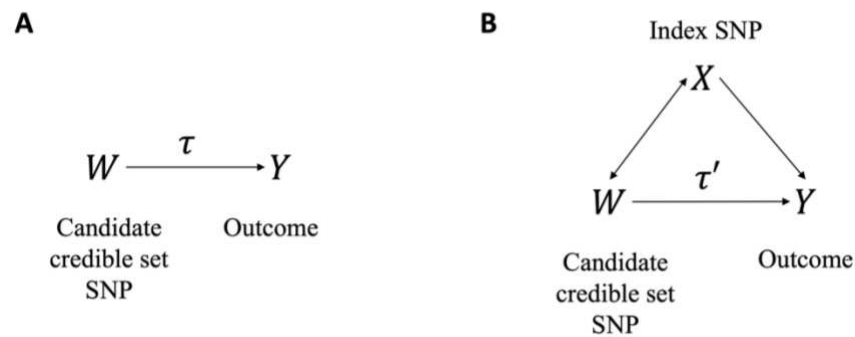
688    Ancestrally Diverse Populations: Opportunities, Methods, Pitfalls, and Recommendations. Cell
689    *179*, 589-603. 10.1016/j.cell.2019.08.051.
690 43.  Li, B., and Ritchie, M. (2021). From GWAS to Gene: Transcriptome-Wide Association Studies and
691    Other Methods to Functionally Understand GWAS Discoveries. Frontiers in Genetics *12*.
692    10.3389/fgene.2021.713230.

693

**A**

$$W \xrightarrow{\ \ \tau\ \ } Y$$

Candidate      Outcome
credible set
SNP

**B**        Index SNP

$$X$$

$$W \xrightarrow{\ \ \tau'\ \ } Y$$

Candidate      Outcome
credible set
SNP

**Figure 1 The direct acyclic graphs (DAG) for the probability that the index SNP mediates the candidate credible set SNP effect.**

*(A) Model with the candidate credible set SNP,W, by itself. $\tau$ is the total effect of W on Y. (B) Model with W and X, the index SNP. $\tau'$ is the direct effect of W on Y.*

**Figure 2 Credible set performance in simulation studies with artificial LD structure.**

*(A) Sensitivity, i.e. the proportion of 500 simulations where the true causal SNP was selected in a credible set. (B) Positive Predictive Value (PPV), i.e., the proportion of true causal SNP over the credible set size, averaged over 500 iterations. (C) Average CS size.*

**Figure 3 Performance of index SNP(s) selection in simulation studies with artificial LD structure.**

*(A) Sensitivity, i.e. the proportion of 500 simulations where the true causal SNP was selected in an index SNP. (B) Positive Predictive Value (PPV), i.e., the proportion of causal SNP selected as an index over all selected indices, averaged over 500 iterations. (C) Number of index SNP(s) selected, averaged over 500 iterations.*

**Figure 4 Credible set behaviour of mJAM-SuSiE and mJAM-Forward as causal SNP significance increases.**

*Simulations were conducted under baseline scenario setting (1 causal SNP out of 50 SNPs in total which are divided into 5 LD blocks) with varying effect sizes. The average empirical -log10(P-value) of the causal SNP ranged from 5 to 263, covering most situations seen in practice. Red dashed line indicates requested coverage which is set to be 0.95 for both methods. (A) Empirical credible set coverage; (B) Average number of credible sets selected among 500 simulations.*

**Figure 5 Manhattan plot for mJAM-Forward credible sets at chromosome 12 position 109194870 to 110794870.**

*(A) y-axis is meta-analyzed -log10(P-value) from multi-ethnic analysis; (B) y-axis is -log10(P-value) from ethnic-specific analysis.*

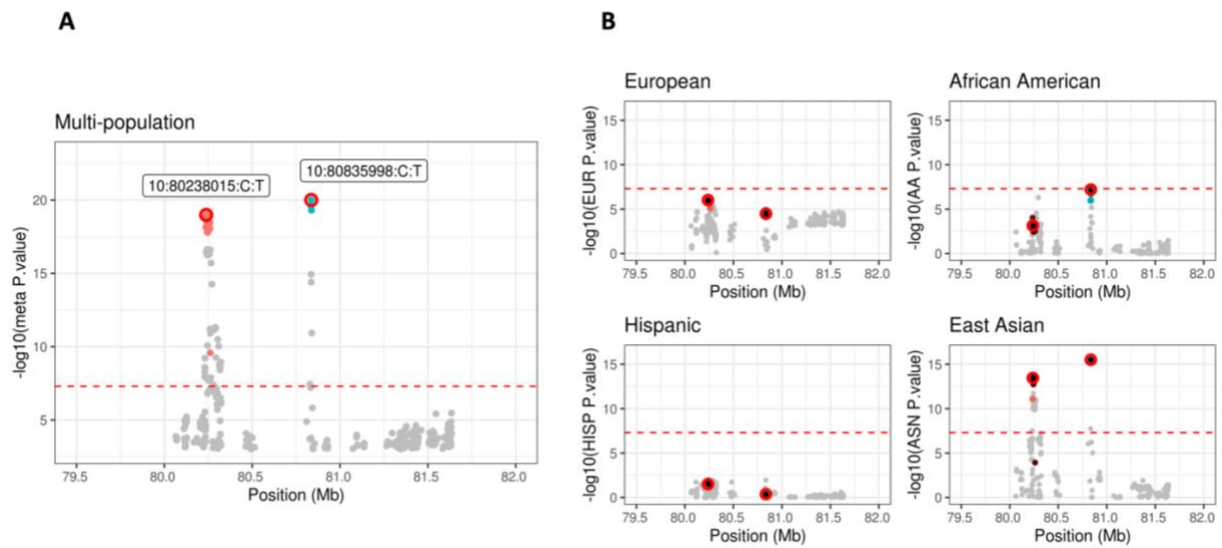**Figure 6 Manhattan plot for mJAM-Forward credible sets at chromosome 10 position 79436999 to 81635998.**

*(A) y-axis is meta-analyzed -log10(P-value) from multi-ethnic analysis; (B) y-axis is -log10(P-value) from ethnic-specific analysis.*
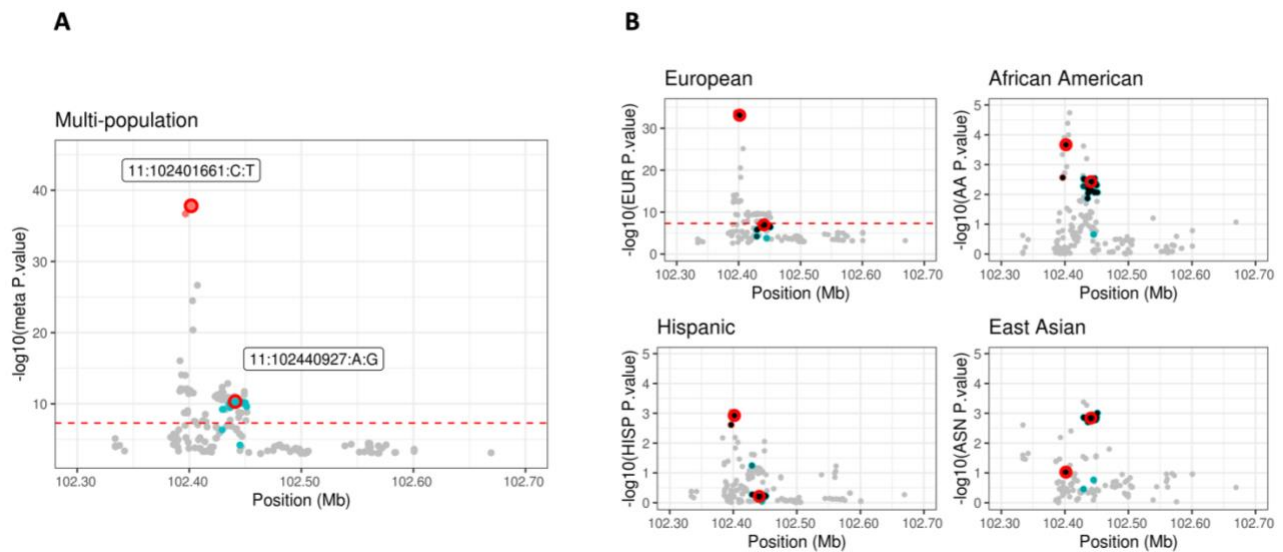
**Figure 7 Manhattan plot for mJAM-Forward credible sets SNPs at chromosome 11 position 101601661 to 103201661.**

*(A) y-axis is meta-analyzed -log10(P-value) from multi-ethnic analysis; (B) y-axis is -log10(P-value) from ethnic-specific analysis.*