

SparsePro: an efficient genome-wide fine-mapping method integrating summary statistics and functional annotations

Wenmin Zhang¹, Hamed Najafabadi^{1,2,3}, and Yue Li^{1,4,*}

¹Quantitative Life Sciences, McGill University, Montreal, Canada

²Department of Human Genetics, McGill University, Montreal, Canada

³McGill Genome Centre, Montreal, Canada

⁴School of Computer Science, McGill University, Montreal, Canada;

*Correspondence to yueli@cs.mcgill.ca

Abstract

Identifying causal variants from genome-wide association studies (GWASs) is challenging due to widespread linkage disequilibrium (LD). Functional annotations of the genome may help prioritize variants that are biologically relevant and thus improve fine-mapping of GWAS results. However, classical fine-mapping methods have a high computational cost, particularly when the underlying genetic architecture and LD patterns are complex. Here, we propose a novel approach, SparsePro, to efficiently conduct functionally informed statistical fine-mapping. Our method enjoys two major innovations: First, by creating a sparse low-dimensional projection of the high-dimensional genotype data, we enable a linear search of causal variants instead of a combinatorial search of causal configurations used in most

20 existing methods; Second, we adopt a probabilistic framework with a highly efficient varia-
21 tional expectation-maximization algorithm to integrate statistical associations and functional
22 priors. We evaluate SparsePro through extensive simulations using resources from the UK
23 Biobank. Compared to state-of-the-art methods, SparsePro achieved more accurate and
24 well-calibrated posterior inference with greatly reduced computation time. We demonstrate
25 the utility of SparsePro by investigating the genetic architecture of five functional biomarkers
26 of vital organs. We show that, compared to other methods, the causal variants identified by
27 SparsePro are highly enriched for expression quantitative trait loci and explain a larger pro-
28 portion of trait heritability. We also identify potential causal variants contributing to the ge-
29 netically encoded coordination mechanisms between vital organs, and pinpoint target genes
30 with potential pleiotropic effects. In summary, we have developed an efficient genome-wide
31 fine-mapping method with the ability to integrate functional annotations. Our method may
32 have wide utility in understanding the genetics of complex traits as well as in increasing the
33 yield of functional follow-up studies of GWASs. SparsePro software is available on GitHub at
34 <https://github.com/zhwm/SparsePro>.

35 **1 Introduction**

36 Establishment of large biobanks and advances in genotyping and sequencing technologies
37 have enabled large-scale genome-wide association studies (GWASs) [1–3]. Although GWASs
38 have revealed extensive associations between genetic variants and traits of interest, under-
39 standing the genetic architecture underlying these genetic associations remains challenging
40 [4–6], mainly because GWASs typically rely on univariate regression models, which are not
41 able to distinguish the causal variants from other variants in linkage disequilibrium (LD) [5, 7,
42 8].

43 Several statistical fine-mapping approaches have been proposed for identifying causal vari-
44 ants in GWASs while considering the underlying LD patterns. For instance, BIMBAM [9], CAVIAR
45 [10] and CAVIARBF [11] estimate the posterior inclusion probabilities (PIPs) in a pre-defined
46 locus by evaluating multivariate Gaussian likelihood enumerating all possible configurations.

47 FINEMAP [12] accelerates the inference with a shotgun stochastic search focusing on the most
48 likely subset of causal configurations. However, the number of causal configurations required
49 to evaluate can grow combinatorially as the number of causal variants increases, thus tremen-
50 dously increasing the computational cost if multiple causal variants exist. SuSiE [13] introduces
51 an iterative Bayesian stepwise selection algorithm for variable selection, which can also be ap-
52 plied to statistical fine-mapping with greatly improved computational efficiency.

53 Furthermore, it has been recognized that functional annotations of the genome may help pri-
54 oritize variants that are biologically relevant, thus improving fine-mapping of GWAS results [8].
55 For example, PAINTOR [14] and RiVIERA [15] empirically estimate the impacts of functional
56 annotations from statistical evidence, which improves the accuracy of fine-mapping but has a
57 high computational cost, especially when multiple causal SNPs exist in the same locus. Poly-
58 Fun [16] adopts stratified LD score regression [17] to effectively partition total trait heritabil-
59 ity into annotation-dependent heritability estimates, and uses these estimates of annotation-
60 tagged heritability to specify functional priors for fine-mapping methods.

61 In this work, we present a unified probabilistic framework called *Sparse Projections to Causal*
62 *Effects* (SparsePro) for statistical fine-mapping with the capacity to incorporate functional anno-
63 tations. Accompanied with an efficient variational expectation-maximization inference algorithm
64 [18], SparsePro achieves superior accuracy in identifying causal variants as well as computa-
65 tional efficiency compared to the state-of-the-art approaches in both simulation studies and real
66 data analyses. We further demonstrate the utility of SparsePro in genome-wide fine-mapping of
67 functional biomarkers for five vital organs in human.

68 **2 Materials and Methods**

69 **2.1 SparsePro method overview**

70 To fine-map causal SNPs, our method takes two lines of evidence (**Figure 1**). First, from esti-
71 mated marginal associations between genetic variants and a complex trait of interest, accom-
72 panied by matched LD information, we can group correlated genetic variants together and as-

73 sess their effects jointly. Then, we infer the contribution of each SNP towards each group of
74 causal effect separately to obtain posterior inclusion probabilities (PIPs). Second, optionally,
75 if we have knowledge about any functional annotations which may be enriched for the causal
76 SNPs, we can estimate the relative enrichment of these annotations, and re-prioritize SNPs
77 according to the enrichments of these annotations. As outputs, our model yields functionally
78 informed PIP for each SNP and the enrichment estimates of candidate functional annotations.

79 **2.2 Our contributions in the context of the existing methods**

80 Our work is related to two existing methods, SuSiE [13] and PolyFun [16]. Inspired by the “sum
81 of single effects” model in SuSiE, we introduce a sparse projection of the genotype in our model
82 specification so that the identification of causal variants and estimation of causal effect sizes
83 are separated. This sparse projection avoids exhaustively evaluating the combinatorial num-
84 ber of causal configurations. For statistical inference, SuSiE adopts an iterative Bayesian step-
85 wise selection algorithm that operates on the Bayes Factors (BFs) [13]. Here, we use a paired
86 mean field variational inference algorithm [18] to jointly update the variational parameters for
87 the causal effects and causal indicators of each SNP. Moreover, we have adapted our algorithm
88 to directly work with GWAS summary statistics and provided appropriate estimates for the hy-
89 perparameters including trait variance and heritability estimates. To enable functionally informed
90 fine-mapping, PolyFun uses genome-wide heritability estimates from LD score regression to
91 set the functional priors for fine-mapping methods [16]. In contrast, we aggregate the genome-
92 wide statistical fine-mapping evidence by maximizing the evidence lower bound of SparsePro
93 to prioritize relevant annotations and robustly derive genome-wide functional priors.

94 2.3 SparsePro model specification

We assume the following data generative process (**Figure 1**) for a continuous polygenic trait.

First, the prior probability $\tilde{\pi}_g$ for the g -th SNP ($g \in \{1, \dots, G\}$) being causal is defined as:

$$\tilde{\pi}_g = \text{softmax}(\mathbf{A}_g \mathbf{w}) = \frac{\exp(\mathbf{A}_g \mathbf{w})}{\sum_{g'=1}^G \exp(\mathbf{A}_{g'} \mathbf{w})}$$

95 where \mathbf{A}_g is the $1 \times M$ annotation row vector of M candidate annotations for the g -th SNP; and

96 \mathbf{w} is the $M \times 1$ vector of logarithm of relative enrichment. Here, we use the *softmax* function

97 to ensure the prior probabilities sum up to 1. If no functional information is provided, the prior

98 probability of being causal is considered equal for all SNPs, i.e. $\tilde{\pi}_g = \frac{1}{G}$.

We assume that there exist K independent causal effects, and that

$$\mathbf{s}_k \sim \text{Multinomial}(1, \tilde{\boldsymbol{\pi}})$$

99 where $\tilde{\boldsymbol{\pi}} = (\tilde{\pi}_1, \dots, \tilde{\pi}_G)$ and \mathbf{s}_k is a binary indicator vector of length G indicating which SNP is the

100 causal SNP under the k -th ($k \in \{1, \dots, K\}$) causal effect.

Then, the causal effect sizes are sampled from a normal distribution, i.e.

$$\beta_k \sim \mathcal{N}(0, \tau_{\beta_k}^{-1})$$

Finally, the continuous trait $\mathbf{y}_{N \times 1}$ over N individuals is generated as follows:

$$\mathbf{y} = \mathbf{X} \sum_k \mathbf{s}_k \beta_k + \boldsymbol{\epsilon}$$

or in matrix form:

$$\mathbf{y} = \mathbf{X} \mathbf{S} \boldsymbol{\beta} + \boldsymbol{\epsilon}$$

101 where $\mathbf{X}_{N \times G}$ is the full genotype matrix, $\mathbf{S}_{G \times K}$ is the sparse projection matrix, $\boldsymbol{\beta}_{K \times 1}$ is the causal

102 effect vector, and $\boldsymbol{\epsilon}_{N \times 1} \sim \mathcal{N}(0, \tau_y^{-1} \mathbf{I}_N)$ denotes the variance not attributable to the modelled ge-

103 netic effects.

104 2.4 A variational inference algorithm for Bayesian fine-mapping

With this model specification (**Figure 1**), finding the causal variants is equivalent to inferring the sparse projections s_k and the effect sizes β_k given y and \mathbf{X} for $k \in \{1, \dots, K\}$:

$$p(\mathbf{S}, \boldsymbol{\beta} | y, \mathbf{X}, \tilde{\boldsymbol{\pi}}, \tau_\beta, \tau_y) = \frac{p(y, \mathbf{S}, \boldsymbol{\beta} | \mathbf{X}, \tilde{\boldsymbol{\pi}}, \tau_\beta, \tau_y)}{p(y | \mathbf{X}, \tilde{\boldsymbol{\pi}}, \tau_\beta, \tau_y)}$$

As the number of possible causal configurations grows combinatorial with G , the exact posterior solution is intractable because of the marginal likelihood in the denominator. Unlike most existing fine-mapping approaches using sampling-based methods to search through a subset of possible causal configurations [12, 19], we adopt a paired mean field factorization of variational family to approximate the posterior [18]:

$$q(\mathbf{S}, \boldsymbol{\beta}) = \prod_k q(s_k, \beta_k) = \prod_k q(s_k)q(\beta_k | s_k)$$

105 This variational distribution preserves the dependency between s_k and β_k . It has been shown
106 that the paired mean field variational family has similar mode and shape as the desired poste-
107 rior distribution, and that such inference can achieve high accuracy with substantially improved
108 computational efficiency [18].

To find the best approximation, we minimize the Kullback-Leibler (KL) divergence between the posterior distribution and the proposed variational distribution, which is equivalent to maximizing the evidence lower bound (ELBO) [20]:

$$ELBO = E_q[\log p(y, \mathbf{S}, \boldsymbol{\beta} | \mathbf{X}, \tilde{\boldsymbol{\pi}}, \tau_\beta, \tau_y)] - E_q[\log q(\mathbf{S}, \boldsymbol{\beta})]$$

109 Based on the mean field assumptions [18], this optimization can be conducted iteratively for
110 the k^{th} causal effect and the g^{th} SNP with the following closed-form updates until convergence
111 (derivation details are available in **Supplementary Notes**).

112 We update posterior effect size for the g -th SNPs in the k -th causal effect:

$$\mu_{kg}^* = \frac{\tau_y}{\tau_{kg}^*} (\mathbf{X}_g^\top \mathbf{y} - \mathbf{X}_g^\top \mathbf{X} \sum_{k' \neq k} \gamma_{k'}^* \circ \boldsymbol{\mu}_{k'}^*) \quad (1)$$

113 with

$$\tau_{kg}^* = \mathbf{X}_g^\top \mathbf{X}_g \tau_y + \tau_{\beta_k} \quad (2)$$

114 where \circ represents element-wise multiplication of vectors.

115 We then update the posterior probability of the g -th SNP being causal in the k -th causal ef-
116 fect:

$$\gamma_{kg}^* = \text{softmax}(\log \tilde{\pi}_g - \frac{1}{2} \log \frac{\tau_{kg}^*}{2\pi} + \frac{\tau_{kg}^* \mu_{kg}^{*2}}{2}) \quad (3)$$

117 For fine-mapping, we take the maximum of these K probabilities as the PIP for SNP g : $\gamma_g^* =$
118 $\max(\gamma_{1g}^*, \dots, \gamma_{Kg}^*)$.

119 2.5 Adaptation to GWAS summary statistics

120 The above variational inference algorithm requires access to large datasets containing both
121 individual-level genotype \mathbf{X} and phenotype data \mathbf{y} . Since a growing number of GWASs have
122 released publicly available summary statistics (i.e., marginal effect size estimate $\hat{\beta}_g$ and its
123 standard error se_g for the g -th SNP), we adapt SparsePro to directly operate on these sum-
124 mary statistics with additional information from an LD reference panel (i.e. estimates of pair-
125 wise SNP-SNP Pearson correlation).

126 Specifically, if we have reasonable surrogates for $\mathbf{X}_g^\top \mathbf{X}_g$, $\mathbf{X}^\top \mathbf{X}$ and $\mathbf{X}_g^\top \mathbf{y}$, we can plug them
127 into Equations (1), (2), and (3). We include two forms of reformulation depending on whether
128 the genotypes are standardized to have zero mean and unit variance in the GWAS.

1. If the genotypes are standardized, we have

$$\mathbf{X}_g^\top \mathbf{X}_g = N$$

$$\mathbf{X}^\top \mathbf{X} = N * LD$$

$$\mathbf{X}_g^\top \mathbf{y} = N \hat{\beta}_g$$

129 where N is the sample size.

2. If the genotypes are not standardized, we have

$$\hat{\beta}_g = (\mathbf{X}_g^\top \mathbf{X}_g)^{-1} \mathbf{X}_g^\top \mathbf{y}$$

$$se_g = \sqrt{\text{var}(\mathbf{y})(\mathbf{X}_g^\top \mathbf{X}_g)^{-1}}$$

Therefore,

$$\mathbf{X}_g^\top \mathbf{X}_g = \frac{\text{var}(\mathbf{y})}{(se_g^2)}$$

$$\mathbf{X}^\top \mathbf{X} = LD * (\mathbf{se}^\top \mathbf{se})$$

$$\mathbf{X}_g^\top \mathbf{y} = \mathbf{X}_g^\top \mathbf{X}_g * \hat{\beta}_g$$

130 Notably, if \mathbf{y} has been standardized to have unit variance prior to a GWAS, we naturally
 131 supply $\text{var}(\mathbf{y}) = 1$. Otherwise, it can be estimated as $\text{var}(\mathbf{y}) = 2Np(1 - p)se^2$ where N
 132 (the study sample size), p (minor allele frequencies), and se (standard errors of effect size
 133 estimates) are usually available in GWAS summary statistics.

134 **2.6 Variational expectation-maximization for integrating functional anno-** 135 **tations**

136 To estimate the relative enrichment of functional annotations and further prioritize variants,
 137 we adopt a variational expectation-maximization scheme to maximize ELBO with respect to
 138 the logarithm of relative enrichment (\mathbf{w}) first and then use the estimate $\hat{\mathbf{w}}$ to calculate $\tilde{\pi}_g$ (prior
 139 probability of being causal) for each SNP.

Suppose we have M candidate annotations and A_{gm} ($m \in \{1, \dots, M\}$) is a 0/1 indicator denoting whether the g -th SNP has the m -th annotation. By setting the derivative of ELBO with respect to w_m to 0 and solving for w_m , we have the following estimate for the logarithm of rele-

vant enrichment (detailed in **Supplementary Notes**),

$$w_m = \log\left(\frac{r_1/r_0}{k_1/k_0}\right)$$

where

$$k_1 = \sum_g [A_{gm} = 1] \text{softmax}\left(\sum_{m' \neq m} A_{gm'} w_{m'}\right)$$

$$k_0 = \sum_g [A_{gm} = 0] \text{softmax}\left(\sum_{m' \neq m} A_{gm'} w_{m'}\right)$$

$$r_1 = \sum_{k,g} [A_{gm} = 1] \gamma_{kg}^*$$

$$r_0 = \sum_{k,g} [A_{gm} = 0] \gamma_{kg}^*$$

We note that this metric is equivalent to the logarithm of a relative risk, thus its standard error can be calculated as

$$se(w_m) = \sqrt{\frac{1}{r_1} + \frac{1}{r_0} - \frac{1}{k_1} - \frac{1}{k_0}}$$

We evaluate the significance of annotation enrichment with the log likelihood ratio test (G-test) [21]. Only annotations which demonstrate statistical significance are included in our model to update the prior probability of being causal for each SNP. Specifically,

$$\tilde{\pi}_g = \text{softmax}\left(\sum_m A_{gm} \hat{w}_m\right)$$

140 This functionally informed prior helps prioritize causal SNPs in addition to statistical evidence.

141 2.7 Hyperparameter settings

142 We have three hyperparameters: number of causal effect K , inverse of the unexplained vari-
 143 ance τ_y and inverse variance of causal effect sizes τ_{β_k} in our model. As we show in **Supple-**
 144 **mentary Notes**, our model is not sensitive to the setting of K as long as K is larger than the

145 actual number of independent effects, except that increasing K marginally increases the com-
146 putation time.

147 We set τ_y as

$$\tau_y = \frac{1}{\text{var}(y) * (1 - h^2)}$$

148 where h^2 is the local SNP heritability that can be estimated by a modified Heritability Estimation
149 from Summary Statistics (HESS) [22] based on GWAS summary statistics (**Supplementary**
150 **Notes**)

151 We set τ_β as

$$\tau_\beta = \frac{k}{\text{var}(y) * h^2}$$

152 for each of the independent causal effects. We use $k \in \{1, \dots, K\}$ to account for different effect
153 sizes and to improve model identifiability.

154 2.8 Simulation studies

155 We conducted simulations to showcase the efficiency and utility of our method. We leveraged
156 resources from the UK Biobank [1]. Specifically, we first retained 353,606 White British an-
157 cestry participants by excluding one individual from each pair of closely related individuals
158 (who had a 3rd degree or closer relationship). We then retrieved the genotypes of these indi-
159 viduals based on 271,699 SNPs which had a minor allele frequency ≥ 0.001 and an imputa-
160 tion quality score ≥ 0.6 on chromosome 22. Next, we sampled 50 causal SNPs with a two-
161 fold relative enrichment amongst SNPs that were annotated as “conserved sequences” [23],
162 “DNase I hypersensitive sites” (DHS) [24], “non-synonymous” [25], or that overlapped with his-
163 tone marks H3K27ac [26] or H3K4me3 [24]. We used the GCTA GWAS simulation pipeline [27]
164 to simulate a continuous trait with a per-chromosome heritability of 0.01. We tested the associ-
165 ation between each SNP and this simulated trait, and obtained GWAS summary statistics using
166 the fastGWA software [28]. This process was replicated 22 times to imitate a GWAS. We ob-

167 tained LD information calculated using the UK Biobank participants from https://alkesgroup.broadinstitute.org/UKBB_LD/ [16]. These LD matrices were generated for genome-wide
168 SNPs binned into sliding windows of 3 Mb where two neighboring windows had a 2-Mb over-
169 lap.
170

171 We applied SparsePro to the GWAS summary statistics with the above LD information, and
172 iterated over all sliding windows, first without any functional annotation information. We de-
173 noted the fine-mapping results as “SparsePro-”. Next, we aggregated the results from all 22
174 replications to estimate the relative enrichment for ten binary functional annotations. In addition
175 to the five annotations simulated to be enriched of causal SNPs, we also included five anno-
176 tations without enrichment: “actively transcribed regions” [29], “transcription start sites” [29],
177 “promoter regions” [30], “5'-untranslated regions” [25], and “3'-untranslated regions” [25].

178 Annotations with a G-test p-value $< 1 \times 10^{-6}$ were selected to conduct functionally informed
179 fine-mapping, and the results were denoted as “SparsePro+”. τ_{β} and τ_{γ} were set according to
180 aforementioned empirical estimates. PIPs for SNPs in the 1-Mb centre of each 3-Mb sliding
181 window were extracted.

182 **2.9 Method comparisons using simulated data**

183 We also performed fine-mapping with some of the state-of-the-art methods. To perform fine-
184 mapping with conditional and joint (COJO) analyses [31] and FINEMAP [12], we first selected
185 COJO lead SNPs based on GWAS summary statistics by performing stepwise model selec-
186 tion implemented in the GCTA-COJO software [27]. We then applied FINEMAP with shotgun
187 stochastic search to SNPs in a 1-MB window centered at each COJO-identified lead SNP. We
188 wrote an in-house script using the “susie_rss” function to perform genome-wide fine-mapping
189 with SuSiE in the same sliding windows as SparsePro. We aggregated summary statistics from
190 22 replications and used PolyFun with the “baselineLF2.2.UKB” model [16] to calculate func-
191 tional priors. The “baselineLF2.2.UKB” model contained all annotations used in SparsePro as
192 well as additional pre-computed LD-related annotations for optimal performance of PolyFun
193 [16]. The estimated priors were provided to SuSiE via “prior_weights” and to FINEMAP via the

194 --prior-snps option, respectively. The maximal number of causal SNPs in each locus was set
195 to 5 for all methods.

196 We compared the performance of these methods in terms of precision (1 - false discovery
197 rate), recall, calibration of PIPs, as well as computation time, all evaluated on a 2.1 GHz CPU
198 node on Compute Canada.

199 **2.10 Fine-mapping genetic determinants of functional biomarkers for vi-** 200 **tal organs**

201 To investigate the genetic coordination mechanisms of vital organs, we performed GWAS in
202 the UK Biobank [1] for five functional biomarkers: forced expiratory volume in one second to
203 forced vital capacity (FEV1-FVC) ratio for lung function, estimated glomerular filtration rate for
204 kidney function, pulse rate for heart function, total protein for liver function and blood glucose
205 level for pancreatic islet function. For each trait, we first regressed out the effects of age, age²,
206 sex, genotyping array, recruitment centre, and the first 20 genetic principal components before
207 inverse normal transforming the residuals to z-scores that had zero mean and unit variance.
208 We then performed GWAS analysis on the resulting z-scores with the fastGWA software [27,
209 28] to obtain summary statistics.

210 Using the summary statistics and the matched LD information [16], we performed genome-
211 wide fine-mapping with SparsePro-, SparsePro+, SuSiE and PolyFun-informed SuSiE as de-
212 scribed in the simulation analyses (**Section 2.9**), except that the number of causal effects was
213 set to 9 for each LD region to account for potentially more causal variants.

214 To evaluate the biological relevance of SNPs fine-mapped by different methods, we assessed
215 their relative enrichment in tissue-specific expression quantitative loci (eQTL). Tissue-specific
216 eQTL identified in the most recent release of the Genotype-Tissue Expression (GTEx) project
217 [32, 33] were obtained from <https://gtexportal.org/home/datasets>. The eQTL information
218 was not used by any functionally informed fine-mapping methods.

219 Additionally, we calculated trait heritability conferred by fine-mapped SNPs with SparsePro-
220 and SparsePro+, respectively, at several commonly used PIP thresholds for determining causal

221 variants: 0.50, 0.80, 0.90, 0.95, and 0.99. The adjusted R^2 obtained from multivariate linear re-
222 gression of the z-scores (i.e. inverse normal transformed trait residuals after regressing out co-
223 variate effects) against all fine-mapped SNPs was used as a surrogate of the SNP heritability.
224 We compared these results to heritability captured by the same number of SNPs fine-mapped
225 by SuSiE and PolyFun-informed SuSiE, separately at each PIP threshold. For instance, if SparsePro-
226 identified J SNPs with a PIP > 0.5 , we would select J SNPs with the highest PIP determined
227 by SuSiE and compare the adjusted R^2 . Notably, this analysis evaluates predictive associations
228 instead of actual causality, hence the adjusted R^2 is not a direct indicator of the validity of the
229 fine-mapping results.

230 We selected SNPs with a PIP > 0.8 to explore possible pleiotropic effects using phenogram
231 [34]. Loci with potential pleiotropic effects were visualized using LocusZoom [35].

3 Results

3.1 SparsePro demonstrates superior performance in simulation

We performed simulations based on real genotype data from UK Biobank (**Materials and Methods**). We observed that SparsePro consistently demonstrated superior accuracy in identifying true causal variants. SparsePro without annotation (SparsePro-) achieved an area under the precision-recall curve (AUPRC) of 0.3699, higher than the AUPRC of 0.2677 by FINEMAP and the AUPRC of 0.3573 by SuSiE (**Figure 2A**). Notably, SparsePro had a substantially higher precision at the same recall rates (**Figure 2A**). For example, at the recall rate of 25%, SparsePro achieved greater than 95% precision, which is highly desirable in fine-mapping because only a small number of the prioritized SNPs will be experimentally validated *in vivo* or *in vitro* in practice (**Figure 2A**).

Moreover, SparsePro can incorporate functional priors (**Supplementary Table S1**) with improved fine-mapping power. SparsePro+ achieved an AUPRC of 0.4636, outperforming both functionally informed FINEMAP (AUPRC = 0.3088) and functionally informed SuSiE (AUPRC = 0.4042) with functional priors derived by PolyFun. As expected, we also found that the performance of SparsePro was not sensitive to the pre-specified number of independent causal effects (**Supplementary Table S2** and **Supplementary Notes**).

Compared to FINEMAP and SuSiE, the PIPs yielded by SparsePro appeared to be much more calibrated. It has been shown that for a well-calibrated fine-mapping method, the mean PIP of all SNPs with a PIP above a certain threshold should be equal to the precision if these SNPs were to be considered causal variants [16]. Here, we found that the mean PIP of all SNPs considered to be causal variants by SparsePro was almost identical to the desired precision at any threshold (**Figure 2B**). In contrast, the PIPs generated by FINEMAP and SuSiE appeared to be inflated (**Figure 2B**).

For instance, if SNPs with a $PIP > 0.8$ were to be considered as causal variants, SparsePro- and SparsePro+ would both have a median precision across simulations of 95% and 100% respectively (**Figure 2C**). The selected SNPs by FINEMAP (median precision = 50% only)

259 and SuSiE (median precision = 71% only) included an excessive proportion of false positives,
260 even with functional priors (median precision = 77% for FINEMAP and 79% for SuSiE; **Fig-**
261 **ure 2C**). The high precision by SparsePro was consistent for all frequently used PIP thresholds
262 (**Figure 2C**) although FINEMAP and SuSiE sometimes have a slightly higher recall.

263 Furthermore, SparsePro conferred not only higher fine-mapping precision, but also higher
264 computational efficiency. In our simulation, it took only an hour to fine-map chromosome 22,
265 which was 6.5 times faster than FINEMAP and 3 times faster than SuSiE (**Figure 2D** and **Sup-**
266 **plementary Table S3**).

267 **3.2 Fine-mapped SNPs by SparsePro are more enriched in tissue-specific** 268 **eQTL and confer higher trait heritability**

269 We performed GWAS in the UK Biobank [1] for five functional biomarkers: FEV1-FVC ratio
270 (lung function), estimated glomerular filtration rate (kidney function), pulse rate (heart function),
271 total protein (liver function) and blood glucose level (pancreatic islet function). Genome-wide
272 fine-mapping of five functional biomarkers based on the UK Biobank population using Sparse-
273 Pro identified multiple potentially causal variants (**Supplementary Table S4**). To assess bio-
274 logical relevance of the fine-mapping results, we estimated the relative enrichment of causal
275 signals in tissue-specific eQTL for each trait (**Materials and Methods**). We found that the fine-
276 mapped SNPs were significantly enriched in tissue-specific eQTL for all five biomarkers, while
277 results based on SparsePro-/+ showed the strongest enrichment (**Figure 3A**). For example,
278 for total protein, the fine-mapped SNPs determined by SparsePro- were 4.00-fold (95% CI:
279 3.25-4.92) more likely to be liver-specific eQTL than non-fine-mapped SNPs, compared to a
280 1.54-fold (95% CI: 1.35-1.75) enrichment based on fine-mapped SNPs by SuSiE. While SuSiE
281 was substantially improved by functional priors derived from PolyFun with a 2.20-fold (95% CI:
282 1.97-2.45) enrichment, the fine-mapped SNPs by SparsePro+ exhibited the highest biological
283 relevance, being 4.06-fold (95% CI: 3.31-4.97) more likely to be liver-specific eQTL.

284 Moreover, at most PIP thresholds, the SNPs fine-mapped by SparsePro- explained a higher
285 proportion of phenotypic variance based on all UK Biobank subjects (**Methods**) compared to

286 the same number of the most likely causal SNPs identified by SuSiE (**Figure 3B** and **Sup-**
287 **plementary Table S5**). With the functional annotations (**Supplementary Table S5**), the fine-
288 mapped SNPs by SparsePro+ consistently achieved a higher SNP heritability for estimated
289 glomerular filtration rate, FEV1-FVC ratio, as well as total protein compared to the PolyFun-
290 informed SuSiE; although for glucose and pulse rate, PolyFun-informed SuSiE was able to
291 identify SNPs with a slightly higher predictive performance at certain PIP thresholds (**Figure 3C**
292 and **Supplementary Table S6**).

293 **3.3 Pleiotropic effects of SNPs rs1260326 and rs5742915 on the func-** 294 **tions of multiple vital organs**

295 Overall, we observed considerable polygenicity for the five biomarkers of the vital organs (**Figure 4A**).
296 Interestingly, at the PIP threshold of 0.80, we found two potentially causal variants for three of
297 the five biomarkers. Specifically, SNP rs1260326 (**Figure 4B**), a missense variant (Leu446Pro)
298 in gene *GCKR*, was fine-mapped for glomerular filtration rate (PIP = 1.000), blood glucose level
299 (PIP = 0.998), pulse rate (PIP = 0.823) and total protein level (PIP = 1.000). Notably, this spe-
300 cific variant has been found to be significantly associated with a wide variety of glycemic traits
301 [36] and other quantitative traits for metabolic syndromes and comorbidities [37, 38], and has
302 been implicated in the functions of liver and other vital organs [39–41].

303 Another SNP, rs5742915 (**Figure 4C**), a missense variant (Phe645Leu) in gene *PML* was
304 fine-mapped for FEV1-FVC ratio (PIP = 0.858), pulse rate (PIP = 1.000) and total protein level
305 (PIP = 0.987). This variant has also been associated with other quantitative biomarkers of poly-
306 genic traits featuring development and metabolism, including birth weight [42], height [43], ap-
307 pendicular lean mass [44], and age at menarche [45]. These findings, along with other SNPs
308 exhibiting pleiotropic effects at somewhat lower PIP thresholds (**Supplementary Table S4**) pre-
309 sented promising genetic targets for experimental validations in a larger effort towards under-
310 standing the mechanisms of genetic coordination among vital organs.

311 4 Discussion

312 Accurately identifying trait-determining and disease-causing variants is fundamental in genet-
313 ics and particularly important for appropriately interpreting GWAS results [5, 8]. In this work,
314 we developed SparsePro, an efficient fine-mapping method to help prioritize causal variants for
315 complex traits, possibly with prior functional information. Through genome-wide simulations, we
316 showed that SparsePro was highly accurate and computationally efficient compared to existing
317 methods. By fine-mapping genetic associations with five biomarkers for vital organ functions,
318 we demonstrated that SparsePro identified candidate variants that were biologically relevant,
319 including two variants with pleiotropic effects, which might indicate genetically encoded coordi-
320 nation among vital organs.

321 Compared to the existing methods, SparsePro has three important features. First of all, we
322 use an efficient variational inference algorithm to approximate the posterior distribution of the
323 causal variant indicators instead of exhaustively searching through all possible causal config-
324 urations or performing stepwise regression. As a result, SparsePro can be significantly faster
325 than the existing fine-mapping methods, such as FINEMAP [12], and is more than twice as
326 fast as SuSiE [13], which is a similar variable selection framework but implements an itera-
327 tive Bayesian stepwise selection procedure. The substantially improved computational effi-
328 ciency enables statistical fine-mapping of large chunks of the genome instead of analyzing
329 genetic associations on a per-locus basis as in most existing follow-up studies of GWASs. In
330 our simulation studies, compared to locus-wise fine-mapping based on COJO-identified lead
331 SNPs, such a genome-wide fine-mapping requires neither a pre-specified p-value threshold
332 (e.g. $p < 5 \times 10^{-8}$) for determining candidate loci nor an arbitrary number of causal effects per
333 locus. If functional annotations are available, the estimation of functional enrichment may also
334 be more robust by including more variants with little additional computational overhead.

335 Second, we utilize a paired mean field variational family, where the causal effect and the ca-
336 sual indicator are coupled in the variational distribution. This ensures that our approximation
337 matches closely with the true posterior distribution of the causal variant indicators [18]. As a
338 result, SparsePro yielded better-calibrated PIPs compared to existing fine-mapping methods.

339 Third, given GWAS summary statistics, we provide estimates for hyperparameters including
340 τ_y and τ_β that are reasonable in the context of polygenic trait genetics. Consequently, at several
341 commonly used PIP thresholds for defining causal variants, SparsePro showed improved control
342 of false positives, demonstrated higher precision in identifying causal variants in simulation
343 and obtained stronger enrichment for tissue-specific eQTL in real data application.

344 Last, we propose and implement a probabilistic model that coherently integrates statistical
345 evidence and functional prior information. The key difference between SparsePro+ and other
346 methods that leverage functional priors, such as PolyFun [16] and PAINTOR [14], is that each
347 annotation is tested for its relevance with the trait of interest before being used to derive the
348 priors in our model. Therefore, functional annotations serve as complementary evidence when
349 statistical evidence is not sufficient to discern causal variants. Based on our results, it seems
350 that this approach distills better prior information from the functional annotations compared to
351 the aforementioned methods.

352 We note that SparsePro can be further improved with the following future directions. First,
353 SparsePro generally requires that the supplied LD reference panel matches well with that of the
354 GWAS study population to guarantee proper convergence. While we advocate the public availability
355 of the in-sample LD information along with the GWAS summary statistics, a more robust
356 model is needed to account for mismatched LD information. Second, SparsePro currently supports
357 only binary annotations while compatibility with continuous annotations is also desirable.
358 Last, the current variational expectation-maximization scheme might not accurately estimate
359 the joint enrichment of highly correlated annotations. Performing variable selection beforehand
360 or effectively aggregating enrichment estimates may enable the inclusion of multiple correlated
361 informative annotations, such as cell type-specific annotations to further improve the utility of
362 SparsePro.

363 In summary, SparsePro is an efficient genome-wide fine-mapping method with the ability of
364 integrate functional annotations. We envision its wide utility in understanding the genetic architecture
365 of complex traits, identifying target genes, and increasing the yield of functional follow-up
366 studies of GWASs.

367 5 Figure Legends

368 **Figure 1.** SparsePro overview. The data generating process of SparsePro is depicted in a
369 plate model with shaded nodes represent observed variables and unshaded nodes represent
370 latent variables. The trait y is generated from K causal effects, where the k -th causal effect
371 size $\beta_k \sim \mathcal{N}(0, \tau_{\beta_k})$. We use a sparse projection $\mathbf{s}_k \sim \text{Multinomial}(1, \hat{\boldsymbol{\pi}})$ of genotype to indicate
372 causal variant for the k -th effect. Given the causal effect sizes and sparse indicators of causal
373 variants, the target trait y_i for individual i follows a normal distribution $y_i \sim \mathcal{N}(\mathbf{X}_i \sum_k \mathbf{s}_k \beta_k, \tau_y^{-1})$.
374 To help prioritize variants with functional annotations, we assume the prior probability of being
375 causal $\hat{\pi}_g$ for the g -th variant as $\hat{\pi}_g = \text{softmax}(\mathbf{A}_g \mathbf{w})$ where \mathbf{A}_g is a $M \times 1$ functional annotation
376 vector and \mathbf{w} is the $M \times 1$ vector of annotation enrichment coefficients. We adopt an efficient
377 variational inference algorithm to jointly estimate both causal effect sizes and sparse indicators
378 and an expectation-maximization scheme for estimating annotation enrichment coefficients \mathbf{w}
379 as detailed in Section 2.

380 **Figure 2.** SparsePro demonstrated improved accuracy and computational efficiency in genome-
381 wide simulation results. (A) Precision-Recall curves. The inset shows the area under the preci-
382 sion recall curve (AUPRC) for each method. (B) Calibration of posterior inclusion probabilities
383 (PIPs). The y-axis is the mean PIPs for all SNPs considered as causal variants, correspond-
384 ing to the expected precision at different PIP cutoffs. The x-axis represents the actual preci-
385 sion at different PIP cutoffs. The black dashed line indicates an optimal calibration, where the
386 expected precision perfectly matches the observed precision. (C) Precision and recall rates
387 obtained at five frequently used PIP thresholds. Error bars indicate inter-quartile ranges. (D)
388 Comparison of computational time. Boxes denote inter-quartile ranges and the line inside each
389 box indicates the median running time. The color legends are displayed at the bottom of the
390 figure.

391 **Figure 3.** Biological relevance of fine-mapped SNPs for five biomarkers, each for a distinct
392 vital organ. (A) Relative enrichment of causal signals in tissue-specific eQTL. Target traits and

393 the corresponding organs are indicated. Estimates of relative enrichment with 95% confidence
394 intervals are plotted on a logarithmic scale. (B) Comparison of the proportion of total trait vari-
395 ance explained by fine-mapped SNPs between SparsePro- and SuSiE.(C) Comparison of the
396 proportion of total variance explained by fine-mapped SNPs between SparsePro+ and PolyFun
397 informed SuSiE. Fine-mapped SNPs were identified at five PIP thresholds. As a surrogate of
398 SNP heritability, the proportion of trait variance explained was obtained from multivariate linear
399 regression adjusted R^2 . In this multivariate regression, we regress the inverse normal trans-
400 formed trait residuals against all fine-mapped SNPs after adjusting for covariate effects. We se-
401 lected the same number of top-ranked SNPs for each method separately at each PIP threshold
402 **(Materials and Methods)**.

403 **Figure 4.** Fine-mapping genetic associations for five functional biomarkers of vital organs. (A)
404 Illustration of genome-wide distribution of fine-mapped SNPs on 22 chromosomes. SNPs with
405 a posterior inclusion probability > 0.80 were indicated as colored solid circles. Two loci with
406 potential pleiotropic effects on four and three vital organ biomarkers respectively were high-
407 lighted by red dashed rectangles. Locus zoom plots were presented for these two loci: (B) lo-
408 cus with fine-mapped SNP rs1260326. and (C) locus with fine-mapped SNP rs5742915. SNPs
409 in a ± 500 kb window are included, colored by r^2 with the corresponding fine-mapped SNP.

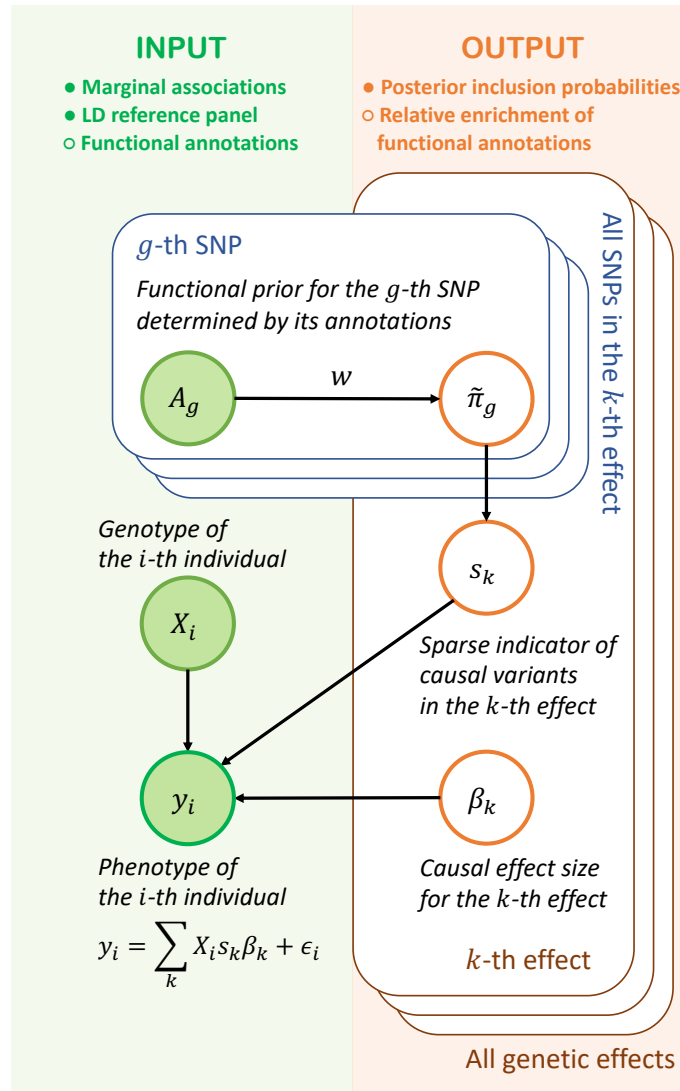


Figure 1: SparsePro overview. The data generating process of SparsePro is depicted in a plate model with shaded nodes represent observed variables and unshaded nodes represent latent variables. The trait y is generated from K causal effects, where the k -th causal effect size $\beta_k \sim \mathcal{N}(0, \tau_{\beta_k})$. We use a sparse projection $s_k \sim \text{Multinomial}(1, \hat{\pi})$ of genotype to indicate causal variant for the k -th effect. Given the causal effect sizes and sparse indicators of causal variants, the target trait y_i for individual i follows a normal distribution $y_i \sim \mathcal{N}(\mathbf{X}_i \sum_k s_k \beta_k, \tau_y^{-1})$. To help prioritize variants with functional annotations, we assume the prior probability of being causal $\hat{\pi}_g$ for the g -th variant as $\hat{\pi}_g = \text{softmax}(\mathbf{A}_g \mathbf{w})$ where \mathbf{A}_g is a $M \times 1$ functional annotation vector and \mathbf{w} is the $M \times 1$ vector of annotation enrichment coefficients. We adopt an efficient variational inference algorithm to jointly estimate both causal effect sizes and sparse indicators and an expectation-maximization scheme for estimating annotation enrichment coefficients \mathbf{w} as detailed in Section 2.

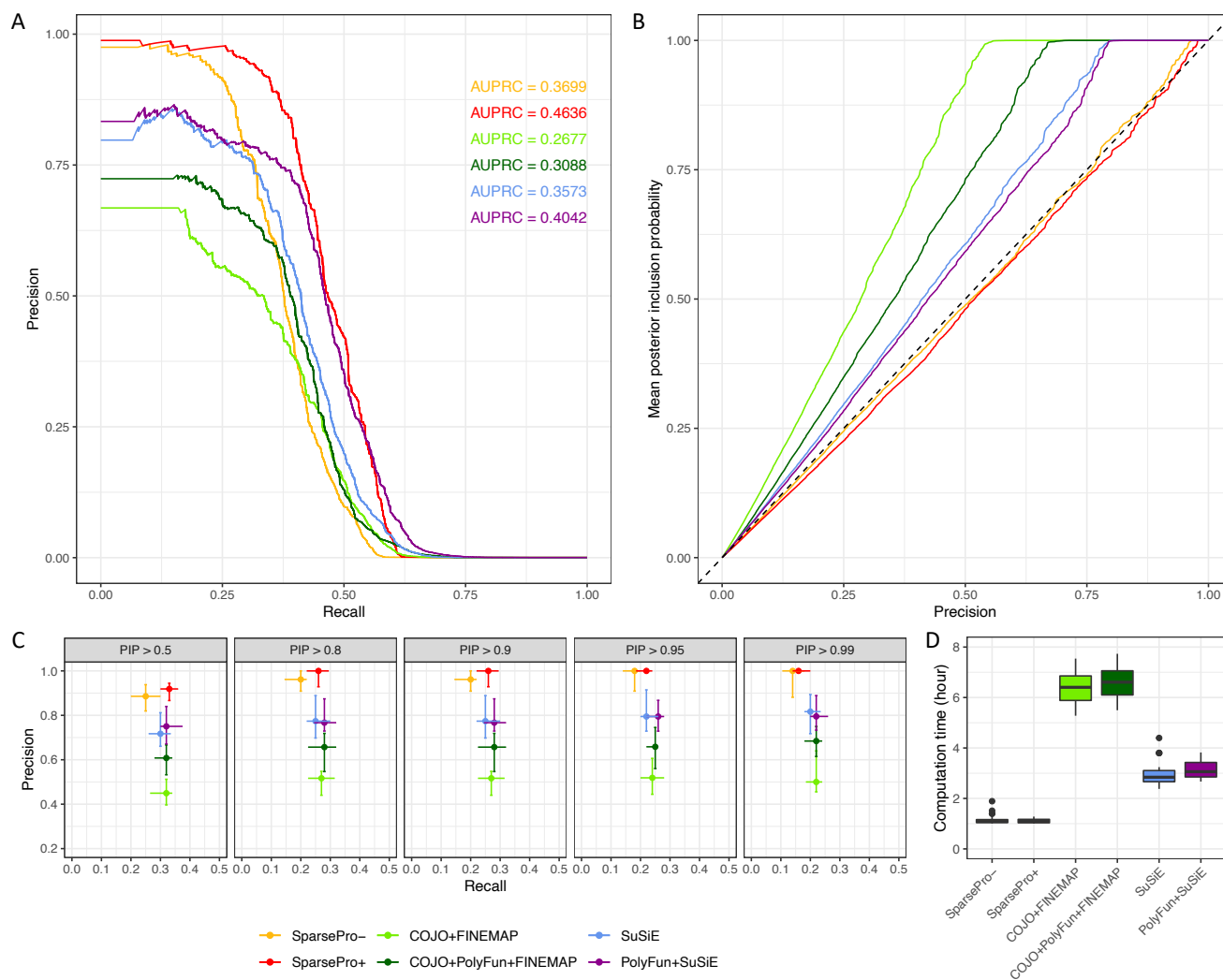


Figure 2: SparsePro demonstrated improved accuracy and computational efficiency in genome-wide simulation results. (A) Precision-Recall curves. The inset shows the area under the precision recall curve (AUPRC) for each method. (B) Calibration of posterior inclusion probabilities (PIPs). The y-axis is the mean PIPs for all SNPs considered as causal variants, corresponding to the expected precision at different PIP cutoffs. The x-axis represents the actual precision at different PIP cutoffs. The black dashed line indicates an optimal calibration, where the expected precision perfectly matches the observed precision. (C) Precision and recall rates obtained at five frequently used PIP thresholds. Error bars indicate inter-quartile ranges. (D) Comparison of computational time. Boxes denote inter-quartile ranges and the line inside each box indicates the median running time. The color legends are displayed at the bottom of the figure.

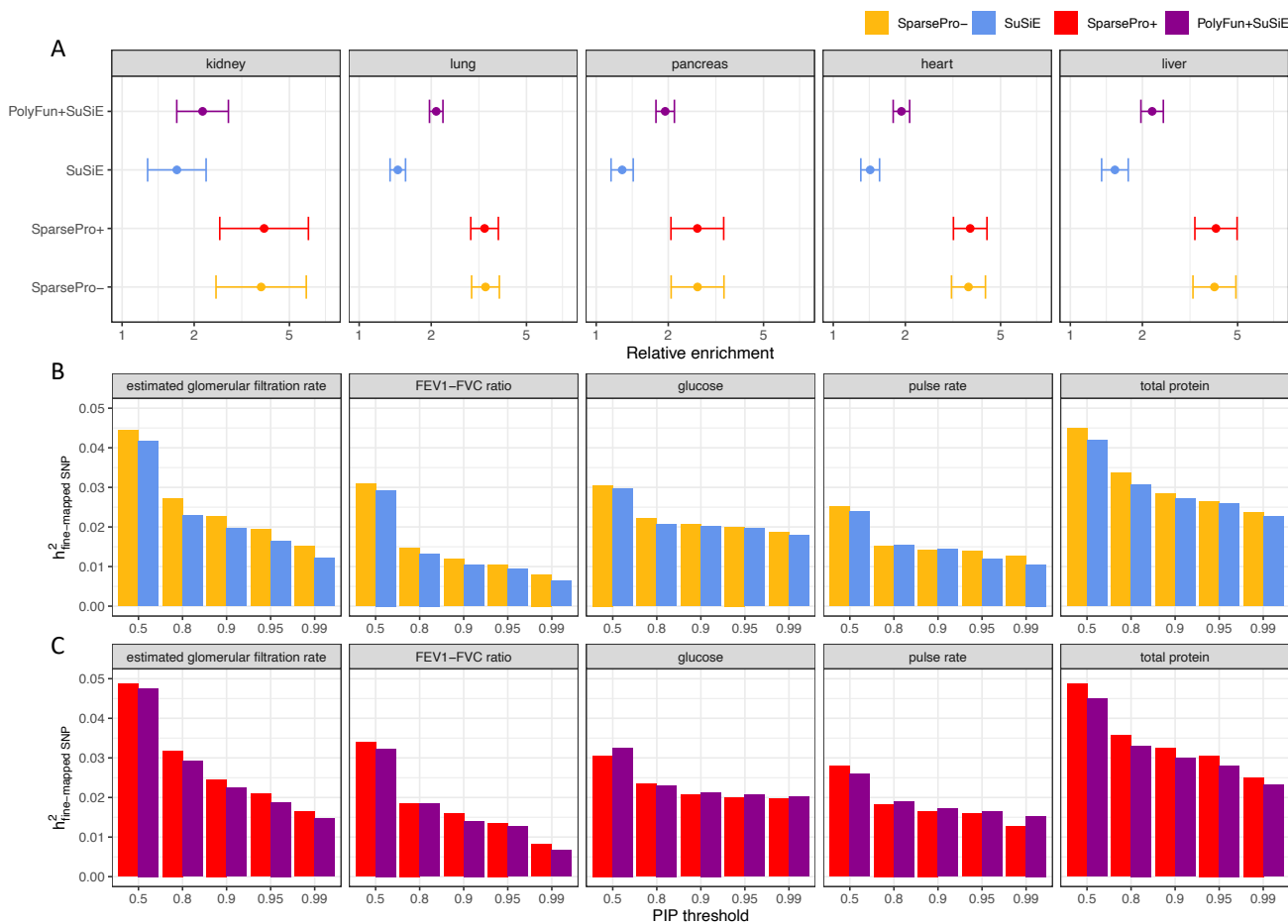


Figure 3: Biological relevance of fine-mapped SNPs for five biomarkers, each for a distinct vital organ. (A) Relative enrichment of causal signals in tissue-specific eQTL. Target traits and the corresponding organs are indicated. Estimates of relative enrichment with 95% confidence intervals are plotted on a logarithmic scale. (B) Comparison of the proportion of total trait variance explained by fine-mapped SNPs between SparsePro- and SuSiE. (C) Comparison of the proportion of total variance explained by fine-mapped SNPs between SparsePro+ and PolyFun informed SuSiE. Fine-mapped SNPs were identified at five PIP thresholds. As a surrogate of SNP heritability, the proportion of trait variance explained was obtained from multivariate linear regression adjusted R^2 . In this multivariate regression, we regress the inverse normal transformed trait residuals against all fine-mapped SNPs after adjusting for covariate effects. We selected the same number of top-ranked SNPs for each method separately at each PIP threshold (**Materials and Methods**).

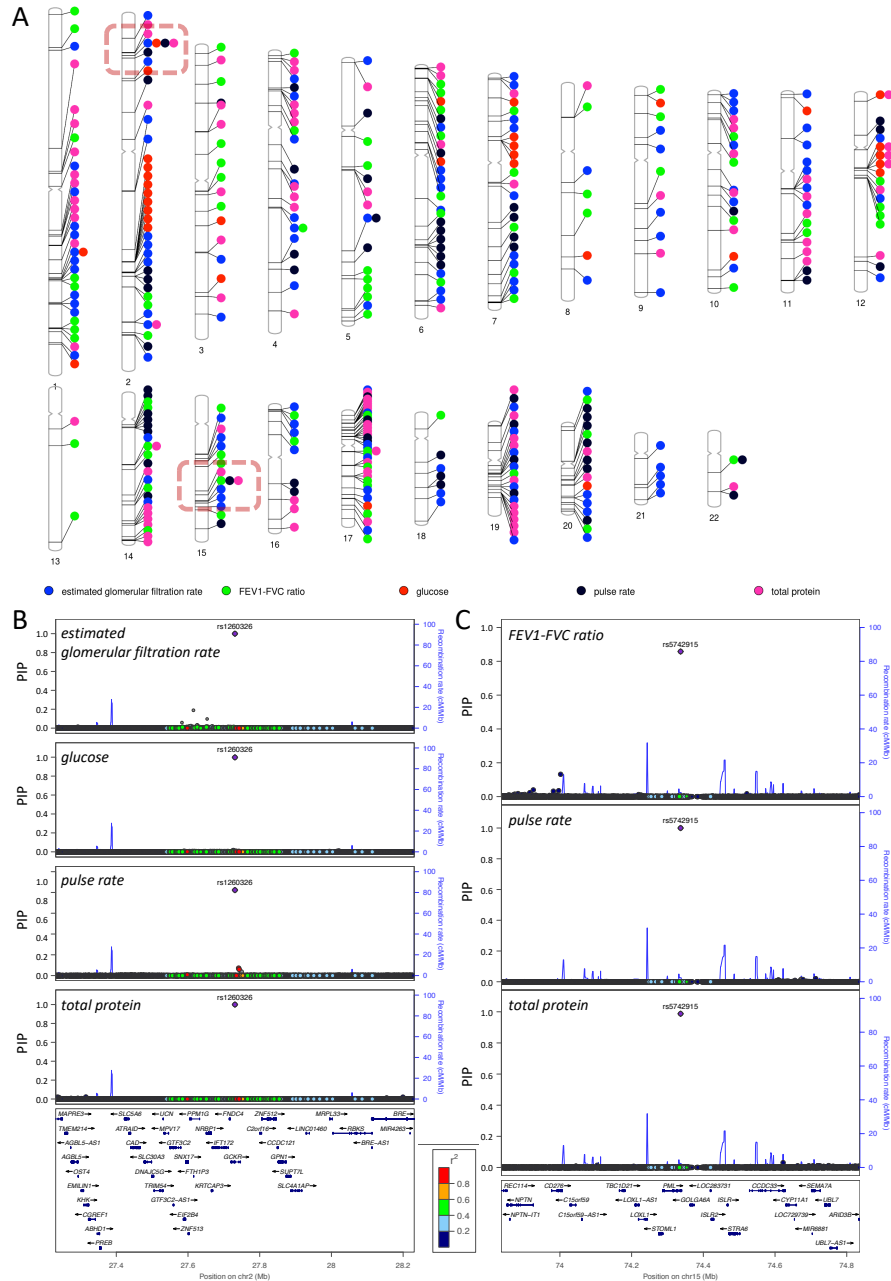


Figure 4: Fine-mapping genetic associations for five functional biomarkers of vital organs. (A) Illustration of genome-wide distribution of fine-mapped SNPs on 22 chromosomes. SNPs with a posterior inclusion probability > 0.80 were indicated as colored solid circles. Two loci with potential pleiotropic effects on four and three vital organ biomarkers respectively were highlighted by red dashed rectangles. Locus zoom plots were presented for these two loci: (B) locus with fine-mapped SNP rs1260326. and (C) locus with fine-mapped SNP rs5742915. SNPs in a ± 500 kb window are included, colored by r^2 with the corresponding fine-mapped SNP.

410 **6 Acknowledgements**

411 YL is supported by Natural Sciences and Engineering Research Council (NSERC) Discovery
412 Grant (RGPIN-2019-0621), Fonds de recherche Nature et technologies (FRQNT) New Ca-
413 reer (NC-268592), and Canada First Research Excellence Fund Healthy Brains for Healthy Life
414 (HBHL) initiative New Investigator start-up award (G249591). This study has been conducted
415 using UK Biobank Resources under Application Number 45551. This study was enabled, in
416 part, by support from Calcul Québec and Compute Canada. WZ has been supported by a doc-
417 toral training fellowship from the Healthy Brains, Healthy Lives Program, funded by the Canada
418 First Research Excellence Fund (CFREF), Quebec's Ministère de l'Économie et de l'Innovation
419 (MEI), and the Fonds de recherche du Québec (FRQS, FRQSC and FRQNT). H.S.N. holds a
420 Canada Research Chair funded by the Canadian Institutes of Health Research.

421 **7 Author contributions**

422 W.Z and Y.L have conceived the study and developed the methodology. W.Z created the com-
423 putational software and ran the analyses. All authors interpreted the results. W.Z. drafted the
424 initial manuscript. H.S.N and Y.L supervised this study and revised the manuscript critically.

425 **8 Disclosures**

426 The authors declare no conflict of interest.

427 **9 Data and Software Availability**

428 SparsePro is an open-access software and publicly available at <https://github.com/zhwm/>
429 SparsePro. All simulation and plotting scripts to reproduce this study are publicly available at
430 https://github.com/zhwm/SparsePro_Paper. Individual-level phenotype and genotype data
431 from the UK Biobank are available upon successful application to its research committee. GCTA

432 were downloaded from https://cnsgenomics.com/software/gcta/bin/gcta_1.93.2beta.zip.
433 FINEAMP were downloaded from http://www.christianbenner.com/finemap_v1.4_x86_64.tgz.
434 SuSiE (version 0.11.42) were installed from CRAN. PolyFun were installed from <https://github.com/omerwe/polyfun>.
435 UK Biobank LD information was downloaded from https://alkesgroup.broadinstitute.org/UKBB_LD/.
436 Tissue-specific eQTL were obtained from https://storage.googleapis.com/gtex_analysis_v8/single_tissue_qtl_data/GTEX_Analysis_v8_eQTL_EUR.tar.

439 References

- 440 1. Bycroft, C. *et al.* The UK Biobank resource with deep phenotyping and genomic data. *Nature* **562**, 203–209 (2018).
441
- 442 2. Canela-Xandri, O., Rawlik, K. & Tenesa, A. An atlas of genetic associations in UK Biobank. *Nature genetics* **50**, 1593–1599 (2018).
443
- 444 3. Loh, P.-R., Kichaev, G., Gazal, S., Schoech, A. P. & Price, A. L. Mixed-model association
445 for biobank-scale datasets. *Nature genetics* **50**, 906–908 (2018).
- 446 4. Visscher, P. M. *et al.* 10 years of GWAS discovery: biology, function, and translation. *The American Journal of Human Genetics* **101**, 5–22 (2017).
447
- 448 5. Schaid, D. J., Chen, W. & Larson, N. B. From genome-wide associations to candidate
449 causal variants by statistical fine-mapping. *Nature Reviews Genetics* **19**, 491–504 (2018).
- 450 6. Stranger, B. E., Stahl, E. A. & Raj, T. Progress and promise of genome-wide association
451 studies for human complex trait genetics. *Genetics* **187**, 367–383 (2011).
- 452 7. Benner, C. *et al.* Prospects of fine-mapping trait-associated genomic regions by using
453 summary statistics from genome-wide association studies. *The American Journal of Human Genetics* **101**, 539–551 (2017).
454
- 455 8. Spain, S. L. & Barrett, J. C. Strategies for fine-mapping complex traits. *Human molecular genetics* **24**, R111–R119 (2015).
456

- 457 9. Servin, B. & Stephens, M. Imputation-based analysis of association studies: candidate
458 regions and quantitative traits. *PLoS genetics* **3**, e114 (2007).
- 459 10. Hormozdiari, F., Kostem, E., Kang, E. Y., Pasaniuc, B. & Eskin, E. Identifying causal vari-
460 ants at loci with multiple signals of association. *Genetics* **198**, 497–508 (2014).
- 461 11. Chen, W. *et al.* Fine mapping causal variants with an approximate Bayesian method using
462 marginal test statistics. *Genetics* **200**, 719–736 (2015).
- 463 12. Benner, C. *et al.* FINEMAP: efficient variable selection using summary data from genome-
464 wide association studies. *Bioinformatics* **32**, 1493–1501 (2016).
- 465 13. Wang, G., Sarkar, A., Carbonetto, P. & Stephens, M. A simple new approach to variable
466 selection in regression, with application to genetic fine mapping. *Journal of the Royal Sta-*
467 *tistical Society: Series B (Statistical Methodology)* **82**, 1273–1300 (2020).
- 468 14. Kichaev, G. *et al.* Integrating functional data to prioritize causal variants in statistical fine-
469 mapping studies. *PLoS genetics* **10**, e1004722 (2014).
- 470 15. Li, Y. & Kellis, M. Joint Bayesian inference of risk variants and tissue-specific epigenomic
471 enrichments across multiple complex human diseases. *Nucleic acids research* **44**, e144–
472 e144 (2016).
- 473 16. Weissbrod, O. *et al.* Functionally informed fine-mapping and polygenic localization of
474 complex trait heritability. *Nature Genetics* **52**, 1355–1363 (2020).
- 475 17. Finucane, H. K. *et al.* Heritability enrichment of specifically expressed genes identifies
476 disease-relevant tissues and cell types. *Nature genetics* **50**, 621–629 (2018).
- 477 18. Titsias, M. & Lázaro-Gredilla, M. Spike and slab variational inference for multi-task and
478 multiple kernel learning. *Advances in neural information processing systems* **24**, 2339–
479 2347 (2011).
- 480 19. Kichaev, G. *et al.* Improved methods for multi-trait fine mapping of pleiotropic risk loci.
481 *Bioinformatics* **33**, 248–255 (2017).
- 482 20. Blei, D. M., Kucukelbir, A. & McAuliffe, J. D. Variational inference: A review for statisticians.
483 *Journal of the American statistical Association* **112**, 859–877 (2017).

- 484 21. Woolf, B. The log likelihood ratio test (the G-test). *Annals of human genetics* **21**, 397–409
485 (1957).
- 486 22. Shi, H., Kichaev, G. & Pasaniuc, B. Contrasting the genetic architecture of 30 complex
487 traits from summary association data. *The American Journal of Human Genetics* **99**, 139–
488 153 (2016).
- 489 23. Lindblad-Toh, K. *et al.* A high-resolution map of human evolutionary constraint using 29
490 mammals. *Nature* **478**, 476–482 (2011).
- 491 24. Trynka, G. *et al.* Chromatin marks identify critical cell types for fine mapping complex trait
492 variants. *Nature genetics* **45**, 124–130 (2013).
- 493 25. Wang, K., Li, M. & Hakonarson, H. ANNOVAR: functional annotation of genetic variants
494 from high-throughput sequencing data. *Nucleic acids research* **38**, e164–e164 (2010).
- 495 26. Hnisz, D. *et al.* Super-enhancers in the control of cell identity and disease. *Cell* **155**, 934–
496 947 (2013).
- 497 27. Yang, J., Lee, S. H., Goddard, M. E. & Visscher, P. M. GCTA: a tool for genome-wide com-
498 plex trait analysis. *The American Journal of Human Genetics* **88**, 76–82 (2011).
- 499 28. Jiang, L. *et al.* A resource-efficient tool for mixed model association analysis of large-
500 scale data. *Nature genetics* **51**, 1749–1755 (2019).
- 501 29. Hoffman, M. M. *et al.* Integrative annotation of chromatin elements from ENCODE data.
502 *Nucleic acids research* **41**, 827–841 (2013).
- 503 30. Villar, D. *et al.* Enhancer evolution across 20 mammalian species. *Cell* **160**, 554–566
504 (2015).
- 505 31. Yang, J. *et al.* Conditional and joint multiple-SNP analysis of GWAS summary statistics
506 identifies additional variants influencing complex traits. *Nature genetics* **44**, 369–375 (2012).
- 507 32. Lonsdale, J. *et al.* The genotype-tissue expression (GTEx) project. *Nature genetics* **45**,
508 580–585 (2013).
- 509 33. Consortium, G. *et al.* The GTEx Consortium atlas of genetic regulatory effects across hu-
510 man tissues. *Science* **369**, 1318–1330 (2020).

- 511 34. Wolfe, D., Dudek, S., Ritchie, M. D. & Pendergrass, S. A. Visualizing genomic information
512 across chromosomes with PhenoGram. *BioData mining* **6**, 1–12 (2013).
- 513 35. Pruim, R. J. *et al.* LocusZoom: regional visualization of genome-wide association scan
514 results. *Bioinformatics* **26**, 2336–2337 (2010).
- 515 36. Chen, J. *et al.* The trans-ancestral genomic architecture of glycemic traits. *Nature genet-*
516 *ics* **53**, 840–860 (2021).
- 517 37. Huang, L. O. *et al.* Genome-wide discovery of genetic loci that uncouple excess adiposity
518 from its comorbidities. *Nature Metabolism* **3**, 228–243 (2021).
- 519 38. Vuckovic, D. *et al.* The polygenic and monogenic basis of blood traits and diseases. *Cell*
520 **182**, 1214–1231 (2020).
- 521 39. Chen, V. L. *et al.* Genome-wide association study of serum liver enzymes implicates di-
522 verse metabolic and liver pathology. *Nature communications* **12**, 1–13 (2021).
- 523 40. Pazoki, R. *et al.* Genetic analysis in European ancestry individuals identifies 517 loci as-
524 sociated with liver enzymes. *Nature communications* **12**, 1–12 (2021).
- 525 41. Bell, S. *et al.* A genome-wide meta-analysis yields 46 new loci associating with biomark-
526 ers of iron homeostasis. *Communications biology* **4**, 1–14 (2021).
- 527 42. Warrington, N. M. *et al.* Maternal and fetal genetic effects on birth weight and their rele-
528 vance to cardio-metabolic risk factors. *Nature genetics* **51**, 804–814 (2019).
- 529 43. Wood, A. R. *et al.* Defining the role of common variation in the genomic and biological ar-
530 chitecture of adult human height. *Nature genetics* **46**, 1173–1186 (2014).
- 531 44. Pei, Y.-F. *et al.* The genetic architecture of appendicular lean mass characterized by asso-
532 ciation analysis in the UK Biobank study. *Communications biology* **3**, 1–13 (2020).
- 533 45. Kichaev, G. *et al.* Leveraging polygenic functional enrichment to improve GWAS power.
534 *The American Journal of Human Genetics* **104**, 65–75 (2019).

SparsePro Supplementary Information

Wenmin Zhang¹, Hamed Najafabadi^{1,2,3}, Yue Li^{1,4,*}

¹Quantitative Life Sciences, McGill University, Montreal, Canada

²Department of Human Genetics, McGill University, Montreal, Canada

³McGill Genome Centre, Montreal, Canada

⁴School of Computer Science, McGill University

*Correspondence: yueli@cs.mcgill.ca

1 Supplementary Notes

1.1 SparsePro is not sensitive to hyperparameter K

The number of causal effects K is an important hyperparameter in statistical fine-mapping. In methods that exhaustively search through causal configurations, the computation time increases combinatorially with K since the number of candidate causal configurations also grows combinatorially. In contrast, in SparsePro, the computation time increases linearly with K . In practice, most of the computation time is spent on loading LD information, thus the computation time varies only slightly with $K \in \{5, 7, 9, 11\}$. The output of SparsePro is not sensitive to the choice of K as long as K is greater than or equal to the actual number of causal effects. In our simulation studies, we found that with $K = 7, 9$, or 11 , the resulting PIPs were extremely highly correlated with those based on $K = 5$, and the overall AUPRC metrics were also highly consistent (**Supplementary Table S2**).

547 **1.2 Modified HESS estimates for hyperparameters τ_y and τ_β**

548 Local heritability estimates are useful in setting hyperparameters for SparsePro. Shi et al. [22]
549 provided an unbiased estimator for local heritability estimation based on summary statistics:

$$\hat{h}_g = \frac{N\hat{\boldsymbol{\beta}}^T \mathbf{R}^{-1} \hat{\boldsymbol{\beta}} - P}{N - P}$$

550 where \mathbf{R} is the LD matrix, $\hat{\boldsymbol{\beta}}$ is GWAS summary effect size, N is the sample size in the GWAS
551 and P is the number of SNPs considered in a locus. However, this estimate requires that when
552 generating summary statistics, both genotypes and phenotypes should be standardized to have
553 zero mean and unit variance. Since summary statistics generated by some GWAS pipelines do
554 not specifically standardize the genotypes and phenotypes, we modified the HESS estimator to
555 account for the non-unit variance:

$$\hat{h}_g = \frac{(\hat{\boldsymbol{\beta}} \circ \mathbf{v})^T (\mathbf{X}^T \mathbf{X})^{-1} (\hat{\boldsymbol{\beta}} \circ \mathbf{v}) - \text{var}(\mathbf{y})P}{\text{var}(\mathbf{y})(N - P)}$$

556 where \circ represents element-wise multiplication and \mathbf{v} is a $P \times 1$ vector: $v_p = \mathbf{X}_p^T \mathbf{X}_p$ for the p -th
557 SNP with genotype vector \mathbf{X}_p . This estimate can be adapted to directly operate on summary
558 statistics as explained in **Materials and Methods**.

559 **1.3 Full derivation of variational EM algorithm:**

As has been described in **Materials and Methods**, based on the data generative process, for the k -th causal effect, we have:

$$\mathbf{s}_k \sim \text{Multinomial}(1, \tilde{\boldsymbol{\pi}})$$

$$\beta_k \sim \mathcal{N}(0, \tau_{\beta_k}^{-1})$$

$$\mathbf{y} = X \sum_k \mathbf{s}_k \beta_k + \boldsymbol{\epsilon}$$

560 with $\epsilon_i \sim N(0, \tau_y^{-1})$. Therefore, we have the joint probability:

$$p(\mathbf{y}, \mathbf{S}, \boldsymbol{\beta} | \mathbf{X}, \tilde{\boldsymbol{\pi}}, \tau_\beta, \tau_y) = p(\mathbf{y} | \mathbf{X}, \mathbf{S}, \boldsymbol{\beta}, \tau_y) \prod_k p(\beta_k | \tau_{\beta_k}) \prod_k p(\mathbf{s}_k | \tilde{\boldsymbol{\pi}}) \quad (4)$$

The goal of fine-mapping is to infer the posterior probability, and in particular, of the sparse projection \mathbf{S} (from here we make the dependency on hyperparameters implicit for the ease of notation):

$$p(\mathbf{S}, \boldsymbol{\beta} | \mathbf{y}, \mathbf{X}) = \frac{p(\mathbf{y}, \mathbf{S}, \boldsymbol{\beta} | \mathbf{X})}{p(\mathbf{y} | \mathbf{X})}$$

561 We use a paired mean field factorized [18] variational family $q(\mathbf{S}, \boldsymbol{\beta})$ to approximate the pos-
562 terior:

$$q(\mathbf{S}, \boldsymbol{\beta}) = \prod_k q(\mathbf{s}_k, \beta_k) = \prod_k q(\mathbf{s}_k) q(\beta_k | \mathbf{s}_k)$$

563 Note that in this variational family, we do not specify the form of the distribution; rather, we
564 only specify the dependency of β_k on \mathbf{s}_k and that all K causal effects are independent of each
565 other. Also, the form of the variational family does not depend on any observed data.

566 To better approximate the posterior distribution with members of the variational family, we
567 aim to minimize the KL divergence between the posterior distribution and the proposed varia-
568 tional distribution, which is equivalent to maximizing the ELBO [20]:

$$ELBO = E_{q(\mathbf{S}, \boldsymbol{\beta})}[\log p(\mathbf{y}, \mathbf{S}, \boldsymbol{\beta} | \mathbf{X})] - E_{q(\mathbf{S}, \boldsymbol{\beta})}[\log q(\mathbf{S}, \boldsymbol{\beta})]$$

To maximize the above ELBO, the following requirement should be satisfied for each k :

$$\log q(\mathbf{s}_k, \beta_k) = E_{q(\mathbf{S}_{\setminus k}, \boldsymbol{\beta}_{\setminus k})}[(\mathbf{y}, \mathbf{S}, \boldsymbol{\beta} | \mathbf{X})]$$

where $E_{q(\mathbf{S}_{\setminus k}, \boldsymbol{\beta}_{\setminus k})}$ is the expectation with respect to the variational distribution excluding the k -th

component. With the joint probability provided in Equation (4) we have

$$\begin{aligned} \log p(\mathbf{y}, \mathbf{S}, \boldsymbol{\beta} | \mathbf{X}) &= \log p(\mathbf{y} | \mathbf{X}, \mathbf{S}, \boldsymbol{\beta}) + \sum_k \log p(\beta_k | \tau_{\beta_k}) + \sum_k (\mathbf{s}_k | \tilde{\boldsymbol{\pi}}) \\ &= \frac{N}{2} \log \frac{\tau_y}{2\pi} - \frac{\tau_y}{2} (\mathbf{y} - \mathbf{X}(\sum_k \mathbf{s}_k \beta_k))^\top (\mathbf{y} - \mathbf{X}(\sum_k \mathbf{s}_k \beta_k)) \\ &\quad + \sum_k \left(\frac{1}{2} \log \frac{\tau_{\beta_k}}{2\pi} - \frac{\tau_{\beta_k}}{2} \beta_k^2 \right) + \sum_k \sum_g s_{kg} \log \tilde{\pi}_g \end{aligned}$$

569 Taking expectation with respect to the variational distribution excluding the k -th component
570 and plugging in $s_{kg} = 1$ and $\mathbf{s}_{k \setminus g} = \mathbf{0}$ for all SNPs excluding the g -th SNP, we can obtain the joint
571 distribution of the k -th effect as:

$$\log q(s_{kg} = 1, \mathbf{s}_{k \setminus g} = \mathbf{0}, \beta_k) = \text{const} - \frac{\tau_{\beta_k}}{2} \beta_k^2 - \frac{\tau_y}{2} \mathbf{X}_g^\top \mathbf{X}_g \beta_k^2 + \tau_y \beta_k \mathbf{X}_g^\top (\mathbf{y} - \mathbf{X} \tilde{\boldsymbol{\beta}}_{\setminus k}) + \log \tilde{\pi}_g \quad (5)$$

572 where

$$\tilde{\boldsymbol{\beta}}_{\setminus k} = E_{q(\mathbf{s}_{\setminus k}, \boldsymbol{\beta}_{\setminus k})} \left[\sum_{k' \neq k} \mathbf{s}_{k'} \beta_{k'} \right] = \sum_{k' \neq k} \gamma_{k'}^* \circ \mu_{k'}^*$$

573 We recognize that

$$q(\beta_k | s_{kg}=1, \mathbf{s}_{k \setminus g} = \mathbf{0}) \sim \mathcal{N}(\mu_{kg}^*, \tau_{kg}^*)$$

By matching sufficient statistics for this normal distribution, we can obtain the following variational parameters for updates:

$$\begin{aligned} \tau_{kg}^* &= \tau_y \mathbf{X}_g^\top \mathbf{X}_g + \tau_{\beta_k} \\ \mu_{kg}^* &= \frac{\tau_y}{\tau_{kg}^*} \mathbf{X}_g^\top (\mathbf{y} - \mathbf{X} \tilde{\boldsymbol{\beta}}_{\setminus k}) \end{aligned}$$

By integrating out β_k in Equation (5), we obtain that

$$\log q(s_{kg} = 1, \mathbf{s}_{k \setminus g} = \mathbf{0}) = \log \tilde{\pi}_g - \frac{1}{2} \log \frac{\tau_{kg}^*}{2\pi} + \frac{1}{2} \tau_{kg}^* \mu_{kg}^{*2} + const$$

Therefore, the posterior probability of the g -th SNP being causal in the k -th effect can be estimated as:

$$\gamma_{kg}^* := q(s_{kg} = 1, \mathbf{s}_{k \setminus g} = \mathbf{0}) = softmax(\log \tilde{\pi}_g - \frac{1}{2} \log \frac{\tau_{kg}^*}{2\pi} + \frac{1}{2} \tau_{kg}^* \mu_{kg}^{*2})$$

This completes the variational expectation step in our inference algorithm. When functional annotations are available, we use the following maximization step to integrate relevant annotations. After the expectation step, we have that

$$\begin{aligned} ELBO &= const + \sum_{k,g} \gamma_{k,g}^* \log \tilde{\pi}_g \\ &= const + \sum_{k,g} \gamma_{k,g}^* \log \frac{\exp(\mathbf{A}_g \mathbf{w})}{\sum_g \exp(\mathbf{A}_g \mathbf{w})} \\ &= const + \sum_{k,g} \gamma_{k,g}^* [\mathbf{A}_g \mathbf{w} - \log(\sum_g \exp(\mathbf{A}_g \mathbf{w}))] \end{aligned}$$

To maximize ELBO with respect to the relative enrichment of the m -th candidate annotation,

we take partial derivatives of ELBO with respect to w_m and set it to 0 to solve for w_m :

$$\begin{aligned}
 \frac{\partial ELBO}{\partial w_m} &= \sum_{k,g} \gamma_{k,g}^* \left[A_{gm} - \frac{\sum_g A_{gm} \exp(\mathbf{A}_g \mathbf{w})}{\sum_g \exp(\mathbf{A}_g \mathbf{w})} \right] \\
 &= \sum_{k,g} \gamma_{k,g}^* \left[A_{gm} - \frac{\sum_g A_{gm} \exp(A_{gm} w_m) \exp(\sum_{m' \neq m} A_{gm'} w_{m'})}{\sum_g \exp(A_{gm} w_m) \exp(\sum_{m' \neq m} A_{gm'} w_{m'})} \right] \\
 &= \sum_{k,g} \gamma_{k,g}^* \left[A_{gm} - \frac{\sum_g A_{gm} \exp(A_{gm} w_m) \text{softmax}(\sum_{m' \neq m} A_{gm'} w_{m'})}{\sum_g \exp(A_{gm} w_m) \text{softmax}(\sum_{m' \neq m} A_{gm'} w_{m'})} \right] \\
 &= \sum_{k,g} [A_{gm} = 1] \gamma_{kg}^* \\
 &\quad - \sum_{k,g} \gamma_{kg}^* \frac{e^{w_m} \sum_g [A_{gm} = 1] \text{softmax}(\sum_{m' \neq m} A_{gm'} w_{m'})}{e^{w_m} \sum_g [A_{gm} = 1] \text{softmax}(\sum_{m' \neq m} A_{gm'} w_{m'}) + \sum_g [A_{gm} = 0] \text{softmax}(\sum_{m' \neq m} A_{gm'} w_{m'})} \\
 &= r_1 - (r_1 + r_0) \frac{k_1 e^{w_m}}{k_1 e^{w_m} + k_0} \\
 &= 0
 \end{aligned}$$

where

$$\begin{aligned}
 k_1 &= \sum_g [A_{gm} = 1] \text{softmax}(\sum_{m' \neq m} A_{gm'} w_{m'}) \\
 k_0 &= \sum_g [A_{gm} = 0] \text{softmax}(\sum_{m' \neq m} A_{gm'} w_{m'}) \\
 r_1 &= \sum_{k,g} [A_{gm} = 1] \gamma_{kg}^* \\
 r_0 &= \sum_{k,g} [A_{gm} = 0] \gamma_{kg}^*
 \end{aligned}$$

We then obtain:

$$\frac{k_1 e^{w_m}}{k_1 e^{w_m} + k_0} = \frac{r_1}{r_1 + r_0}$$

574

and solve for:

$$w_m = \log \left(\frac{r_1/r_0}{k_1/k_0} \right)$$

575 Notably, this estimate is analogous to a relative risk estimate in a 2×2 contingency table. Sup-
576 pose we consider one annotation, then k_1 corresponds to the number of variants with this spe-
577 cific annotation while k_0 corresponds to the number of variants without the annotation. Mean-
578 while, r_0 corresponds to the sum of posterior probability for variants with the annotation while r_1
579 corresponds to the sum of posterior probability for variants without the annotation.

Similarly, the standard error of this estimate can be calculated based on the standard error of a relative risk:

$$se(\hat{w}_m) = \sqrt{\frac{1}{r_1} + \frac{1}{r_0} - \frac{1}{k_1} - \frac{1}{k_0}}$$

580 Finally, we can evaluate the statistical significance of enrichment with the log likelihood ratio
581 test (G-test) [21].

582 2 Supplementary Table Legends

583 **Supplementary Table S1** Relative enrichment of functional priors in simulation studies.

584 **Supplementary Table S2** Comparison of fine-mapping results based on different hyperpa-
585 rameter settings of the number of causal effects K .

586 **Supplementary Table S3** Details of computation time by each method.

587 **Supplementary Table S4** Fine-mapping results for five functional biomarkers based on the
588 UK Biobank, including genetic variants with a PIP > 0.1 .

589 **Supplementary Table S5** Relative enrichment of functional annotations for five functional
590 biomarkers.

591 **Supplementary Table S6** Comparison of fine-mapped SNP heritability for five functional
592 biomarkers.