1

2

3

4

# A scalable Bayesian method for integrating functional information in genome-wide association studies

7

Jingjing Yang[1], Lars G. Fritsche[1,2], Xiang Zhou[1,*], Gonçalo Abecasis[1,*], on behalf of the International Age-related Macular Degeneration Genomics Consortium (IAMDGC)

10

11

[1]Center for Statistical Genetics, Department of Biostatistics, University of Michigan School of Public Health, 1415 Washington Heights, Ann Arbor, MI 48109, USA.

[2]K.G. Jebsen Center for Genetic Epidemiology, Department of Public Health, NTNU, Norwegian University of Science and Technology, Trondheim, Norway.

16

17

*Correspondence authors

E-mail: xzhousph@umich.edu (XZ) and goncalo@umich.edu (GA)

20

1

# 1 Abstract

2 Although genome-wide association studies (GWASs) have identified many risk loci for

3 complex traits and common diseases, most of the identified associations reside in

4 noncoding regions and have unknown biological functions. Recent genomic sequencing

5 studies have produced a rich resource of annotations that help characterize the function of

6 genetic variants. Integrative analysis that incorporates these functional annotations into

7 GWAS can help elucidate the biological mechanisms underlying the identified associations

8 and help prioritize causal-variants. Here, we develop a novel, flexible Bayesian variable

9 selection model with efficient computational techniques for such integrative analysis.

10 Different from previous approaches, our method models the effect-size distribution and

11 probability of causality for variants with different annotations and jointly models genome-

12 wide variants to account for linkage disequilibrium (LD), thus prioritizing associations based

13 on the quantification of the annotations and allowing for multiple causal-variants per locus.

14 Our efficient computational algorithm dramatically improves both computational speed and

15 posterior sampling convergence by taking advantage of the block-wise LD structures of

16 human genomes. With simulations, we show that our method accurately quantifies the

17 functional enrichment and performs more powerful for identifying true causal-variants than

18 several competing methods. The power gain brought up by our method is especially

19 apparent in cases when multiple causal-variants in LD reside in the same locus. We also

20 apply our method for an in-depth GWAS of age-related macular degeneration with 33,976

21 individuals and 9,857,286 variants. We find the strongest enrichment for causality among

22 non-synonymous variants (54x more likely to be causal, 1.4x larger effect-sizes) and

23 variants in active promoter (7.8x more likely, 1.4x larger effect-sizes), as well as identify 5

24 potentially novel loci in addition to the 32 known AMD risk loci. In conclusion, our method is

25  shown to efficiently integrate functional information in GWASs, helping identify causal-

26  variants and underlying biology.

27

## Author summary

29  We propose a novel Bayesian hierarchical model to account for linkage disequilibrium (LD)

30  and multiple functional annotations in GWAS, paired with an expectation-maximization

31  Markov chain Monte Carlo (EM-MCMC) computational algorithm to jointly analyze genome-

32  wide variants. Our method improves the MCMC convergence property to ensure accurate

33  Bayesian inference of the quantifications of the functional enrichment pattern and fine-

34  mapped association results. By applying our method to the real GWAS of age-related

35  macular degeneration (AMD) with various functional annotations (i.e., gene-based,

36  regulatory, and chromatin states), we find that the variants of non-synonymous, coding, and

37  active promoter annotations have the highest causal probability and the largest effect-sizes.

38  In addition, our method produces fine-mapped association results in the identified risk loci,

39  two of which are shown as examples (*C2/CFB/SKIV2L* and *C3*) with justifications by

40  haplotype analysis, model comparison, and conditional analysis. Therefore, we believe our

41  integrative method will be useful for quantifying the enrichment pattern of functional

42  annotations in GWAS, and then prioritizing associations with respect to the learned

43  functional enrichment pattern.

44

45  **Keywords**: functional information; genome-wide association study (GWAS); Bayesian

46  variable selection regression (BVSR); expectation-maximization (EM); Markov chain Monte

47  Carlo (MCMC); age-related macular degeneration (AMD).

48

## Introduction

50    Genome-wide association studies (GWASs) have identified thousands of genetic loci

51    for complex traits and diseases, providing new insights into the underlying genetic

52    architecture [1-5]. Each associated locus typically contains hundreds of variants in linkage

53    disequilibrium (LD) [6,7], most of which are of unknown function and located outside

54    protein-coding regions. Unsurprisingly, the biological mechanisms underlying the identified

55    associations are often unclear [8] and pinpointing causal variants is difficult [9].

56    Recent functional genomic studies help understand and pinpoint causal variants and

57    mechanisms [10-12]. Genetic variants can be annotated based on the genomic location

58    (e.g., coding, intronic, and intergenic), role in determining protein structure and function

59    (e.g., Sorting Intolerant From Tolerant (SIFT) [13] and Polymorphism Phenotyping

60    (PolyPhen) [14] scores), ability to regulate gene expression (e.g., expression quantitative

61    trait loci (eQTL) and allelic specific expression (ASE) evidence [15,16]), biochemical

62    function (e.g., DNase I hypersensitive sites (DHS), metabolomic QTL (mQTL) evidence

63    [17], and chromatin states [18-20]), evolutionary significance (e.g., Genomic Evolutionary

64    Rate Profiling (GERP) annotations [21]), and a combination of different types of annotation

65    (e.g., CADD [22]). Many statistical methods, including stratified LD score regression [23]

66    and MINQUE [24], can now evaluate the role of functional annotations in GWASs through

67    heritability analysis. Preliminary studies also show higher proportions of associated variants

68    in protein-coding exons, regulatory regions, and cell-type-specific DHSs [25-27].

69    Integrating functional information into GWASs is expected to help identify and

70    prioritize true causal associations. However, accomplishing this goal in practice requires

71    methods to account for both LD and computational cost. Consider two recent methods,

72    Fgwas [26] and PAINTOR [27], as examples: Fgwas assumes that variants are

73    independent and there is at most one causal variant per locus, modeling no LD, which

74    dramatically improves computational speed and allows Fgwas to be applied at genome-

4

75 wide scale; PAINTOR accounts for LD, assuming the possibility of multiple association

76 signals per locus, but is computationally slow and can only be used to fine-map small

77 regions.

78     Here, we pair a flexible Bayesian method with an efficient computational algorithm.

79 Together the two represent an attractive means to incorporate functional information into

80 association mapping. Our model accounts for genotype correlation due to LD, allows for

81 multiple causal variants per locus and, importantly, shares information genome-wide to

82 increase association-mapping power. Our algorithm takes advantage of the local LD

83 structure in the human genome [28-30] and refines previous Markov chain Monte Carlo

84 (MCMC) algorithms to greatly improve mixing, which is key when searching for causal

85 variants among many associated variants in LD (but less important in other applications

86 such as modeling total genomic heritability). We refer to our method as the Bayesian

87 functional GWAS (bfGWAS). Below, we illustrate the benefits of our method with extensive

88 simulations and real data analyses of a large-scale GWAS on age-related macular

89 degeneration (AMD) [31] with 33,976 individuals and 9,857,286 genotyped or imputed

90 variants. The software for bfGWAS is freely available at https://github.com/yjingj/bfGWAS.

91

92 ## Results

93 **Method overview**

94     Our method is based on the standard Bayesian variable selection regression (BVSR)

95 model (Online Methods and Text S1; Fig S1(a)), allowing for annotations that classify

96 variants into $K$ non-overlapping categories. We assume that variants in annotation category

97 $q$ share a "spike-and-slab" prior [32,33] for effect-sizes, $\beta_i \sim \pi_q N\left(0, \tau^{-1}\sigma_q^2\right) + \left(1 - \right.$

98 $\left.\pi_q\right)\delta_0(\beta_i)$. This model implies effect sizes are normally distributed as $\beta_i \sim N\left(0, \tau^{-1}\sigma_q^2\right)$ with

99 probability $\pi_q$, or set to zero with probability $(1 - \pi_q)$, with $\delta_0(\beta_i)$ denoting the point-mass

100    function at 0. Here, $\pi_q$ represents the (unknown) causal probability for variants in the $q$th

101    category and $\sigma_q^2$ represents the (unknown) corresponding effect-size variance. An

102    enhancement to previous Bayesian models [33-35] is that we model both the proportion of

103    associated variants and their effect-size distribution in each annotation category.

104      Our goal is to simultaneously make inference on category specific parameters

105    $(\pi_q, \sigma_q^2)$ that represent the importance of each functional category, and on the variant

106    specific parameters —— effect-size $\beta_i$ and the probability of $\beta_i \neq 0$ (referred as posterior

107    inclusion probability ( $PP_i$ ), representing association evidence). Our model shares

108    information among genome-wide to estimate category specific parameters, which then

109    inform the variant specific parameters. As a result, variant associations will be prioritized

110    based on the inferred importance of functional categories.

111      Because standard MCMC algorithms suffer from heavy computational burden and

112    poor mixing of posterior samples for large GWASs, we develop a novel scalable

113    expectation-maximization MCMC (or EM-MCMC) algorithm. Our algorithm is based on the

114    observation that LD decays exponentially with distance and displays local block-wise

115    structure along the human genome [28-30,36,37]. This observation allows us to decompose

116    the complex joint likelihood of our model into a product of block-wise likelihoods (Online

117    Methods and Text S1). Intuitively, conditional on a common set of category specific

118    parameters $(\pi_q, \sigma_q^2)$, we can infer $(\beta_i, PP_i)$ by running the MCMC algorithm per genome-

119    block. A diagram of this EM-MCMC algorithm is shown in Fig S1(b).

120      Running MCMC per genome-block facilitates parallel computing and reduces the

121    search space. Unlike previous MCMC algorithms for GWAS that use proposal distributions

122    based only on marginal association evidence (such as implemented in GEMMA [38]), our

123    MCMC algorithm uses a proposal distribution that favors variants near the "causal" variants

124    being considered in each iteration, and prioritizes among these neighboring variants based

125  on their conditional association evidence (see Text S1). Our strategy dramatically improves

126  the MCMC mixing property, encouraging our method to explore different combinations of

127  potentially causal variants in each locus (Fig S2). In addition, we implemented memory

128  reduction techniques that reduce memory usage up to 97%, effectively reducing the

129  required physical memory from 120 GB (usage by GEMMA [38]) to 3.6 GB for a GWAS with

130  ~33K individuals and ~400K genotyped variants (Online Methods and Text S1).

131  In practice, we segment the whole genome into blocks of 5,000 ~ 10,000 variants,

132  based on marginal association evidence, genomic distance, and LD. We always ensure

133  variants in LD ($R^2$ >0.1) with significant signals (P-values $<5 \times 10^{-8}$) are in the same block

134  (Online Methods). We first initialize category specific parameters $(\pi_q, \sigma_q^2)$, then run the

135  MCMC algorithm per block (E-step), summarize the MCMC posterior estimates of $(\beta_i, PP_i)$

136  across all blocks to update $(\pi_q, \sigma_q^2)$ (M-step), and repeat the block-wise EM-MCMC steps

137  until $(\pi_q, \sigma_q^2)$ estimates converge (Fig S1(b)).

138  In addition, we calculate the regional posterior inclusion probability (regional-PP) per

139  block that is the proportion of MCMC iterations with at least one "causal" variant (see Text

140  S1). Because Bayesian PP might be split among multiple variants in high LD, the threshold

141  of regional-PP >0.95 (conservatively analogous to false discovery rate 0.05) is used for

142  identifying loci.

143

144  **Simulation**

145  We simulated phenotypes with the genotype data (chromosomes 20-22) from the

146  AMD GWAS [31], including 33,976 individuals and 241,500 variants with minor allele

147  frequency (MAF) >0.1. We segmented this small genome into 50 x 2.5Mb blocks, each with

148  ~5,000 variants. Within each block, we marked a 25KB continuous region (starting 37.5Kb

149  from the beginning of a block) as the causal locus and randomly selected two causal single

7

150    nucleotide polymorphisms (SNPs) per locus. We simulated two complementary annotations

151    to classify variants into "coding" and "noncoding" groups, where the coding variants account

152    for ~1% overall variants but ~10% variants within the causal loci (matching the pattern in

153    the real AMD data). We simulated two scenarios: (i) coding variants ~44x enriched among

154    causal variants (30 coding vs. 70 noncoding); (ii) no enrichment (1 coding vs. 99

155    noncoding). A total of 15% of phenotypic variance was divided equally among causal

156    variants. We compared bfGWAS with single variant likelihood-ratio test, conditional analysis

157    (CA), and Fgwas. The single variant test P-value (also referred to as P-value), conditioned

158    P-value, Fgwas posterior association probability (PP, see Online Methods), and our

159    Bayesian PP were used as criteria to identify associations.

160        We first compared power of different methods using average ROC curves[27,33]

161    across 100 simulation replicates. Fgwas was more powerful than P-value at low false-

162    positive rates (FPR), presumably because Fgwas incorporates annotation information (Fig

163    1(a)). However, with high false-positive rates, Fgwas underperformed P-value, presumably

164    because Fgwas incorrectly assumes one variant per locus. In contrast, bfGWAS (modeling

165    LD and allowing multiple causal variants per locus) outperformed both Fgwas and P-value

166    for false-positive rates in (0, 0.01). Importantly, the advantage of bfGWAS became more

167    pronounced with increasing sample size (Fig S3). Specifically, the power (based on

168    FPR=0.5%) of bfGWAS increased from 48% to 64% as the sample size increased from 20K

169    to 33K, while the power of Fgwas only increased from 52% to 56% and the power of P-

170    values only increased from 47% to 52%. In addition, with sample size 33K and the

171    threshold of regional-PP >0.95, bfGWAS has power 92.3% to identify associated loci,

172    versus Fgwas with 88.6% power. The advantage of bfGWAS with large sample size

173    suggests that bfGWAS can better extract the richer information available as sample size

174    increases.

175    In a typical GWAS, researchers identify a series of associated loci and then examine

176    associated variants within each locus independently. We examined the ability of each

177    method to prioritize the true causal variants in each locus. Since we simulated two causal

178    SNPs per locus (SNP1 and SNP2), we examine the power for identifying each of these

179    separately (Fig 1(b)). All methods have the same median rank for causal SNP1 (typically,

180    ranked 3rd rank among 150 SNPs in the locus by P-value, Fgwas and bfGWAS),

181    suggesting that the strongest signal in a locus can often be identified without incorporating

182    functional information. The median rank for the second causal SNP2 was the 7th by

183    bfGWAS, 12th by Fgwas, 17th by P-value, and 18th by conditional analysis — suggesting

184    that incorporating functional information improves power to identify multiple signals in a

185    locus. Stratified results based on the LD between two causal variants further demonstrate

186    that bfGWAS has the highest power for identifying the weaker signal, especially when both

187    SNPs are in high LD (Fig S4).

188    Both bfGWAS and Fgwas correctly identified enrichment in scenario (i) and properly

189    controlled for the type I error of enrichment in scenario (ii), despite some numerical issues

190    for Fgwas (Supplement Fig 5). Moreover, bfGWAS estimated the effect-size variance per

191    annotation. For all 100 simulation replicates under both scenarios, the 95% confidence

192    intervals of the log-ratio of estimated effect-size variances between coding and noncoding

193    overlapped with 0 (Fig S6), suggesting effect-size variances were similar between two

194    annotations (matching the simulated truth).

195    In summary, our simulation studies show that, in comparisons with competing

196    methods, bfGWAS has higher power, especially in loci with multiple associated variants and

197    when the sample size is large. Further, bfGWAS produces enrichment parameter estimates

198    that can help with interpretation of association results.

199

**GWAS of AMD**

Next, we applied our method to a GWAS of age-related macular degeneration (AMD) with 16,144 advanced cases and 17,832 controls, for a total of 33,976 unrelated European individuals. A total of 439,350 variants were genotyped on a customized Exome-Chip, and then imputed up to 12,023,830 variants in 1000 Genomes Project Phase 1 [39,40]. We analyzed 9,866,744 (~10M) low-frequency and common variants (MAF >0.5%) with three types of genomic annotations: gene-based functional annotations by SeattleSeq, summarized regulatory annotations [41], and the chromatin states profiled in nine human cell types from chromHMM [42,43].

**Coding variation and AMD.**

We used SeattleSeq to classify variants according to their impact on coding sequences (Table S1) and then applied our method bfGWAS and Fgwas. bfGWAS identified 37 loci out of 1,063 considered genome-blocks with regional-PP >0.95 (Tables S2, S3, and S5), including 32 among the 34 known AMD loci [31] and 5 potentially novel loci. Using the threshold of Bayesian PP >0.1068 (roughly equivalent to the P-value $5 \times 10^{-8}$ based on permutations of AMD data; Fig S7), we identified 150 associated variants (Fig S9(a); Table S3), with 47 distributed among 42,005 non-synonymous variants, 4 among 67,165 synonymous coding variants, 54 among 3,679,235 intronic variants, 18 among 5,512,423 intergenic variants (including non-annotated variants), and 27 among 565,916 "other-genomic" variants (UTR, non-coding exons, upstream and downstream of genes). Very roughly, this corresponds to fraction of associated variants of ~1:1,000 among non-synonymous variants, 1:15,000 among synonymous variants, 1:100,000 among intronic variants, 1:300,000 among intergenic variants and 1:20,000 among "other-genomic" variants.

225     Similarly, Fgwas identified 46 loci by regional-PP >0.95, including all 34 known loci

226     and 12 potentially novel loci (Tables S2, S4, and S6; Fig S9(b)). Since Fgwas analyzed the

227     whole genome as 4,934 segments (each with 2,000 variants) and, thus, partitioned the

228     genome somewhat differently than our method. Fgwas identified 178 associated variants

229     with Fgwas PP >0.1068, including 24 non-synonymous, 13 coding-synonymous, 42 intronic,

230     40 intergenic, and 59 other-genomic signals. Compared with bfGWAS, the proportion of loci

231     that contain at least one non-synonymous variant with PP >0.1068 is significantly smaller

232     (11 out of 46 by Fgwas vs. 18 of 37 by bfGWAS; P-value = 0.017). Similarly, the proportion

233     of non-synonymous variants prioritized by Fgwas is also significantly smaller (24 out of 178

234     by Fgwas vs. 47 of 150 by bfGWAS; P-value $=7.7 \times 10^{-5}$), indicating that bfGWAS places

235     greater weight on coding variants —— which, as a group, appears to have both a higher

236     prior probability of association and larger effect sizes when associated.

237     Besides replicating the association results within known AMD loci[31], bfGWAS

238     identified five novel loci (Table S5): missense *rs7562391/PPIL3*, *rs61751507/CPN1*,

239     *rs2232613/LBP*, downstream *rs114348558/ZNRD1-AS1*, and splice *rs6496562/ABHD2*.

240     These loci were also identified by Fgwas (Table S6) with different top association variants

241     for *CPN1* (coding-synonymous *rs61733667*) and *ZNRD1-AS1* (downstream *rs116112857*).

242     Interestingly, there are several connections between these potentially novel loci and known

243     AMD loci. For example, the protein encoded by *LBP* is part of the lipid transfer protein

244     family (which also includes *CETP* among the known AMD risk loci) that promotes the

245     exchange of neutral lipids and phospholipids between plasma lipoproteins [44]. Similarly,

246     *ZNRD1-AS1* has been associated with lipid metabolisms [45] and *ABHD2* has been

247     associated with coronary artery disease [46], two other traits where the AMD loci encoding

248     *CETP*, *APOE*, and *LIPC* are also involved. The gene *CPN1* has been associated with age-

249     related disease (specifically, hearing impairment [47]).

250

**Multiple signals in a single locus.** We use two examples to illustrate the importance of studying multiple signals in a single locus. Our first example focuses on a 1Mb region around locus *C2/CFB/SKIV2L* on chromosome 6 where 1,862 variants have P-values < $5 \times 10^{-8}$. There are an estimated 4 independent signals in the region by conditional analysis [31], 21 variants with Fgwas PP >0.1068, 11 with Bayesian PP >0.1068 by the standard Bayesian variable selection regression (BVSR) method that models no functional information, and 12 with Bayesian PP >0.1068 by bfGWAS. Interestingly, the alternative methods (P-value, Fgwas, and BVSR) identified intronic SNP *rs116503776/SKIV2L/NELFE* as the top candidates (P-value = $2.1 \times 10^{-114}$; Fgwas PP = 0.912; BVSR PP = 1.0), while bfGWAS identified two missense SNPs *rs4151667/C2/CFB* (P-value = $1.4 \times 10^{-44}$; bfGWAS PP = 0.917) and *rs115270436/SKIV2L/NELFE* (P-value = $2.8 \times 10^{-99}$; SBA PP = 0.633) as the top functional candidates (Fig 2; Tables S2-S4).

A haplotype analysis describing the odds ratios (ORs) for all possible haplotypes for SNPs *rs116503776, rs4151667*, and *rs115270436*, helps clarify the region. Intronic SNP *rs116503776* with the smallest P-value appears to be associated with the phenotype by tagging the other two missense SNPs (Table S15). In particular, haplotypes with *rs116503776* can either increase or decrease risk, depending on alleles at the other two SNPs. To further confirm the importance of the missense SNPs *rs4151667* and *rs115270436*, we compared the AIC/BIC/loglikelihood between two models: one model with top two independent signals (*rs116503776* and *rs114254831*) identified by single-variant conditional analysis[31], versus the other model with top two signals (*rs4151667* and *rs115270436*) identified by bfGWAS. As expected, the second model has smaller AIC/BIC and larger loglikelihood than the first one (Table S16). Thus, we can see that while alternative methods (P-value, Fgwas, and BVSR) focus on the SNP with the smallest P-

275    value, our bfGWAS method finds an alternative pairing of missense signals that better

276    accounts for all data.

277        Our second example focuses on a 1Mb region around gene *C3* on chromosome 19

278    (Fig S10) with 112 genome-wide significant variants with P-value $<5 \times 10^{-8}$. Fgwas only

279    discovered a single missense signal, *rs2230199* with the most significant P-value=$1.7 \times$

280    $10^{-77}$ (top blue triangle in Fig S10(a, c)). However, both BVSR and bfGWAS identified 2

281    missense variants with PPs = 1.0, and 5 intronic variants with 0.11< PPs <0.18. The top two

282    missense signals *rs2230199* and *rs147859257* (241 base pairs apart) were confirmed by

283    conditional analysis [31], where the second signal *rs147859257* has conditioned P-

284    value=$6.0 \times 10^{-33}$ (the purple triangle in Fig S10(b, d), overlapping with *rs2230199*). These

285    two missense signals match the interpretation of previous studies [48-50]. Because other 5

286    intronic variants (*rs11569479, rs11569470, rs201063729, rs10408682, rs11569466*) are in

287    high LD with between variant $R^2$ >0.98, we believe this is the third independent signal

288    whose Bayesian PP was split among 5 variants in high LD by bfGWAS.

289

290    **Enrichment analysis.** bfGWAS estimated that non-synonymous variants are 10-100 times

291    more likely to be causal than variants in other categories and that they also have larger

292    effect-sizes (Fig 3(a, b)). To better compare enrichment among multiple categories, we

293    define two new sets of parameters (Text S1). The first set of parameters, $(\pi_q/\pi_{avg})$, is

294    defined to contrast the posterior association probability estimate $(\pi_q)$ for each category to

295    the genome-wide average $(\pi_{avg})$. The second set of parameters $(\sigma_q^2/\sigma_{avg}^2)$ is similarly

296    defined to contrast the effect-size variance from each category to the genome-wide

297    average. Moreover, the square root of the effect-size variance reflects the effect-size

298    magnitude because of the prior assumption for the effect-size in our model.

299    Compared to the genome-wide average probability of causality $\pi_{avg}$=$4.3 \times 10^{-06}$ (Fig

300    S12(a)), we found that non-synonymous category were 54x more likely to be causal (P-

301    value= $7.24 \times 10^{-84}$); that coding-synonymous and other variants were 4.3x and 2.2x more

302    likely (P-values = 0.005, 0.003); and that intergenic 0.7x less likely (P-value=$4.9 \times 10^{-6}$);

303    while the intronic variants matched the genome-wide average (P-value=0.659). In addition,

304    compared to the genome-wide average effect-size variance ($\sigma_{avg}^2 = 0.02$; Fig S12(b)), we

305    found that the effect size variance of was 1.9x larger for non-synonymous variants (P-

306    value=0.014; i.e., 1.4x larger effect-size); and 0.4x smaller for variants in the intronic

307    category (P-value=$4.5 \times 10^{-06}$); remaining categories were not significantly different (P-

308    values >0.2). The estimated enrichment parameters by Fgwas show a similar pattern,

309    although the contrast of the estimated enrichment for non-synonymous versus other

310    annotations is not as pronounced as by bfGWAS (Fig S11(a)).

311

312    **Analysis with regulatory annotations**

313    Second, we analyzed the GWAS data of AMD with the summarized regulatory

314    annotations[41]: coding, UTR, promoter (defined as within 2KB of a transcription starting

315    site), DHS in any of 217 cell types, intronic, intergenic, and "others" (not annotated as any

316    of the previous six categories). Overall GWAS results were similar as the ones described in

317    previous context (Tables S7-S10). Compared to the genome-wide average association

318    probability ($\pi_{avg}$=$4.03 \times 10^{-6}$; Fig S12(c)), we found that the association probability of the

319    coding category was 28x higher (P-value $< 2.2 \times 10^{-16}$); the promoter was 2.6x (P-

320    value=0.028) higher; the intergenic and "others" were 0.5x and 0.9x less (P-values =

321    $5.3 \times 10^{-4}$, 0.033); while the DHS and intronic were not significantly different (P-values

322    >0.1). In addition, compared to the genome-wide average effect-size variance ($\sigma_{avg}^2 =$

323    $0.024$), we found that the effect-size variance of the coding category was 1.9x larger (P-

14

324 value=0.019; i.e., 1.4x larger effect-size); the DHS and intronic were 0.5x less (P-values =

325 0.011, 0.007); while the promoter, intergenic, and "others" were not significantly different (P-

326 values >0.1; Fig S12(d)). Here, Fgwas identified a slightly different enrichment pattern (Fig

327 S11(b)), where UTR was identified as the second most enriched category. This is

328 presumably because Fgwas assumes one causal variant per locus and tends to prioritize

329 the variant with the smallest P-value in each locus, e.g., UTR variants

330 *rs1142/KMT2E/SPRK2* and *rs10422209/CNN2* have the highest Fgwas PP and the

331 smallest P-value in their respective locus (Tables S2 and S8).

332

333 **Analysis with chromatin states**

334 Last, we considered the annotations of seven chromatin states obtained with

335 ChromHMM in nine human cell types[43]: active promoter (APromoter), poised promoter

336 (PPromoter), strong enhancer (SEnhancer), weak enhancer (WEnhancer), insulator,

337 transcription elongation (TxnElong), repetitive/copy number variation (CNV). Nine human

338 cell types include: embryonic stem cells (H1-hESC), erythrocytic leukaemia cells (K562), B-

339 lymphoblastoid cells (GM12878), hepatocellular carcinoma cells (HepG2), umbilical vein

340 endothelial cells (HUVEC), skeletal muscle myoblasts (HSMM), normal lung fibroblasts

341 (NHLF), normal epidermal keratinocytes (NHEK) and mammary epithelial cells (HMEC).

342 With each set of chromatin states profiled in one cell type, we applied bfGWAS on

343 the GWAS data of AMD, and then examined the list of variants that contribute 95%

344 posterior probabilities in the identified loci with regional-PP >95%. We found that the results

345 by accounting for the chromatin states profiled in the erythrocytic leukaemia cells (K562)

346 gave the shortest list (average 14 variants per locus; Table S17), and the enrichment

347 analysis results of other cell types were slightly different (Figs S13-S15).

348        Here, we present the results of accounting for the chromatin states profiled in the

349      K562 cell type (Fig 3(e, f); Tables S11-S14). Compared to the genome-wide average

350      association probability ($\pi_{avg} = 4.0 \times 10^{-6}$; Fig S12(e)), the association probability was 7.8x

351      higher for the active promoter category (P-value = $7.4 \times 10^{-10}$), 3x higher for the strong

352      enhancer category (P-value=0.013), 2.6x higher for the weak enhancer category (P-value =

353      0.002), 1.8x higher for the transcription elongation category (P-value = 0.002), 0.4x less for

354      the CNV category (P-values = 0.004). In addition, the effect-size variances of associated

355      variants in active promoter and strong enhancer were found 2x larger than the genome-

356      wide average ($\sigma^2_{avg} = 0.022$; P-values = 0.048, 0.073), while the effect-size variances of

357      weak enhancer, transcription elongation, and CNV categories were not significantly

358      different (P-values >0.1; Fig S12(f)).

359        Note that the Bayesian enrichment estimates of the poised promoter and insulator

360      categories are the same as their priors (not plotted in Fig 3(e, f)), suggesting that bfGWAS

361      identified no associations in these two categories. Again, Fgwas identified a similar

362      enrichment pattern (Fig S11(c)).

363

## Discussion

365        Here, we describe a scalable Bayesian hierarchical method, bfGWAS, for integrating

366      functional information in GWASs to help prioritize functional associations and understand

367      underlying genetic architecture. bfGWAS models both association probability and effect-

368      size distribution as a function of annotation categories for improving fine-mapping

369      resolution. Unlike previous methods [26,27], bfGWAS accounts for LD and allows for the

370      possibility of multiple association signals per locus while remaining capable of genome-wide

371      inference. Further, bfGWAS employs an improved MCMC sampling strategy to greatly

372 improve the mixing of MCMC samples, which ensures the capability of identifying a list of

373 association candidates.

374 By simulation studies, we demonstrated that bfGWAS had higher power than Fgwas

375 and conditioned P-value for identifying multiple signals in a single locus by accounting for

376 both functional information and LD. We also showed that bfGWAS accurately estimated the

377 enrichment patterns under scenarios with or without enrichment for one annotation in

378 simulations. In the real analysis using the AMD GWAS data and three different types of

379 annotations, by bfGWAS, we obtained posterior association probabilities and effect-size

380 variances for variants of considered annotation categories, as well as an improved list of

381 fine-mapped association signals. In addition, we replicated the findings of 32 out of 34

382 known AMD risk loci, as well as identified 5 potentially novel loci by bfGWAS. Further, we

383 gave two fine-mapped AMD loci *C2/CFB/SKIV2L* and *C3* by bfGWAS as examples with

384 justifications by haplotype analysis, model comparison, and previous findings. Thus, we

385 believe our method is useful for understanding the underlying genetic architecture of

386 complex traits and diseases, for efficiently integrating functional information into GWASs.

387

388 Our flexible framework allows for many further extensions. For example, it can be

389 extended to deal with overlapping or quantitative annotations (Text S1). These extensions

390 will allow us to investigate the importance of a broader class of annotations (e.g. Combined

391 Annotation Dependent Depletion (CADD) scores, MAF, and eQTL evidence). Importantly,

392 as the development of new genomic assays and computational tools enables new variant

393 annotations, simultaneous modeling of available annotations will be critical to identify the

394 set of annotations that are important for a specific trait. Then extending bfGWAS to select

395 relevant annotations would be useful.

396    bfGWAS makes a key assumption that the variant correlation matrix has a block-

397    wise structure, which allows us to segment the genome into approximately independent

398    blocks, analyze variants per block by MCMC, and summarize genome-wide information by

399    an EM algorithm. In parallel to our study, many recent studies have also explored the

400    benefits of dividing the human genome into approximately independent LD blocks to

401    facilitate genome-wide analyses [26,51]. Although the standard segmentation methods

402    (e.g., based on genomic location [51] as we adopted here, or the number of variants per

403    block [26]) are often sufficient in practice, we expect that a better segmentation method [30]

404    based on LD blocks will likely further increase the association mapping power.

405    The biggest limitation of bfGWAS is probably computational cost, as we perform

406    MCMC using the complete genotype data. Specifically, bfGWAS took 5,000 CPU hours (~5

407    hours with parallel computations on 1,000 CPUs for the 1,063 genome-blocks) to analyze

408    the AMD GWAS data with 33,976 individuals and 9,857,286 variants. Implementing

409    bfGWAS with summary statistics is expected to reduce the computation cost significantly,

410    which is part of our continuing project. In addition, the variational approximation [52,53] and

411    other approximations [54,55] of MCMC may provide an efficient alternative for posterior

412    inference in large GWAS.

## Materials and Methods

### Bayesian variable selection regression model

Our method is based on the standard Bayesian variable selection regression (BVSR) model

$$\boldsymbol{y}_{n \times 1} = \boldsymbol{X}_{n \times p} \boldsymbol{\beta}_{p \times 1} + \boldsymbol{\epsilon}_{n \times 1}, \quad \beta_i \sim \pi_i N(0, \tau^{-1}\sigma_i^2) + (1 - \pi_i)\delta_0(\beta_i), \qquad \epsilon_i \sim N(0, \tau^{-1}),$$

where $n$ denotes the number of individuals and $p$ denotes the number of genetic variants; $\boldsymbol{y}_{n \times 1}$ is the phenotype vector; $\boldsymbol{X}_{n \times p}$ is the genotype matrix; $\boldsymbol{\beta}_{p \times 1}$ is a vector of genetic effect-sizes where each element $\beta_i$ follows a spike-and-slab prior (known as the point-normal distribution) ---- that is, $\beta_i$ follows a normal distribution $N(0, \tau^{-1}\sigma_i^2)$ with probability $\pi_i$, or $\beta_i$ is set as 0 with probability $(1 - \pi_i)$ and a point mass density function $\delta_0(\beta_i)$ at 0 ($\delta_0(\beta_i) = 1$ if $\beta_i = 0$, $\delta_0(\beta_i) = 0$ otherwise) [32,33]; and $\epsilon_i$ is the residual error that independently and identically follows a normal distribution $N(0, \tau^{-1})$. We assume that both the phenotype vector $\boldsymbol{y}_{n \times 1}$ and columns of the genotype matrix $\boldsymbol{X}_{n \times p}$ are centered, thus dropping the intercept. Although this model is developed for quantitative traits, we can treat binary phenotypes (e.g., cases and controls) as quantitative following previous approaches[33,35].

### Bayesian hierarchical model accounting for functional information

For integrating functional information into the above BVSR model, we classify all variants into disjoint categories by assuming one annotation per variant. We further assume that variants in the same functional category have the same spike-and-slab prior for the effect-sizes, i.e., $\pi_i = \pi_q, \sigma_i^2 = \sigma_q^2$ for the $q$ th category. Consequently, $\pi_q$ denotes the category specific causal probability and $\sigma_q^2$ denotes the category specific effect-size variance (the square root of $\sigma_q^2$ reflects the magnitude of effect size). Although we focus on discrete non-overlapping annotations in this paper, our method can be extended to overlapping and continuous annotations (Text S1).

19

436    We assume a Bayesian hierarchical framework[34] of BVSR with the following

437    independent hyper priors:

$$\pi_q \sim Beta(a_q, b_q), \qquad \sigma_q^2 \sim IG(k_1, k_2), \qquad \pi_q \perp \sigma_q^2,$$

438    where $\pi_q$ follows a Beta distribution with positive shape parameters $a_q$ and $b_q$, $\sigma_q^2$ follows

439    an Inverse-Gamma distribution with shape parameter $k_1$ and scale parameter $k_2$. In order

440    to adjust for the unbalanced distribution of functional annotations among all variants and

441    enforce a sparse model in our analysis, we choose values for $a_q$ and $b_q$ such that the Beta

442    distribution has mean $\frac{a_q}{a_q+b_q} = 10^{-6}$ with $(a_q + b_q)$ equal to the number of variants in

443    category $q$. We set $k_1 = k_2 = 0.1$ in our analysis to induce non-informative prior for $\sigma_q^2$. Note

444    that $\tau$ is fixed at the phenotype variance value in our Bayesian inferences (Text S1).

445    **Bayesian references**

446    We introduce a latent indicator vector $\boldsymbol{\gamma_{p \times 1}}$ to facilitate computation, where each

447    binary element $\gamma_i$ indicates whether $\beta_i = 0$ by $\gamma_i = 0$, or $\beta_i \sim N(0, \tau^{-1}\sigma_i^2)$ by $\gamma_i = 1$.

448    Equivalently,

$$\gamma_i \sim \text{Bernoulli}(\pi_i), \qquad \boldsymbol{\beta_{-\gamma}} \sim \boldsymbol{\delta_0}, \qquad \boldsymbol{\beta_\gamma} \sim MVN_{|\gamma|}(\boldsymbol{0}, \tau^{-1}\boldsymbol{V_\gamma}),$$

449    where $|\boldsymbol{\gamma}|$ denotes the number of 1's in $\boldsymbol{\gamma}$; $\boldsymbol{\beta_{-\gamma}}$ denotes the sub-vector of $\boldsymbol{\beta_{p \times 1}}$

450    corresponding to variants with $\gamma_i = 0$; $\boldsymbol{\beta_\gamma}$ denotes the sub-vector of $\boldsymbol{\beta_{p \times 1}}$ corresponding to

451    variants with $(\gamma_j = 1; j = 1, \dots, |\boldsymbol{\gamma}|)$; and $\boldsymbol{V_\gamma}$ denotes the sub-matrix of the diagonal matrix

452    $\boldsymbol{V_{p \times p}}$ whose $ith$ diagonal element is $V_{ii} = \sigma_i^2$. Consequently, the expectation of $\gamma_i$ is an

453    estimate of the posterior inclusion probability (PP) for the $i$th variant, $E[\gamma_i] = Prob(\gamma_i = 1) =$

454    $PP_i$.

455        For the described Bayesian hierarchical model above, the posterior joint distribution

456    is proportional to

$$P(\boldsymbol{\beta}, \boldsymbol{\gamma}, \boldsymbol{\pi}, \boldsymbol{\sigma}^2, \tau \mid \boldsymbol{y}, \boldsymbol{X}, \boldsymbol{A}) \propto P(\boldsymbol{y}|\boldsymbol{X}, \boldsymbol{\beta}, \boldsymbol{\gamma}, \tau)P(\boldsymbol{\beta}, |\boldsymbol{A}, \boldsymbol{\pi}, \boldsymbol{\sigma}^2, \boldsymbol{\gamma}, \tau)P(\boldsymbol{\gamma}|\boldsymbol{\pi})P(\boldsymbol{\pi})P(\boldsymbol{\sigma}^2)P(\tau),$$

457    where $\boldsymbol{\pi} = \left(\pi_1, \ldots, \pi_Q\right)^T$, $\boldsymbol{\sigma}^2 = \left(\sigma_1^2, \ldots, \sigma_Q^2\right)^T$, $\boldsymbol{A}$ is the $p \times Q$ matrix of binary annotations, and

458    $Q$ is the total number of annotations. The goal is to estimate the category specific

459    parameters $(\boldsymbol{\pi}, \boldsymbol{\sigma}^2)$ and the variant specific parameters $(\boldsymbol{\beta}, E[\boldsymbol{\gamma}])$ from their posterior

460    distributions, conditioning on the data $(\boldsymbol{y}, \boldsymbol{X}, \boldsymbol{A})$. Here, the category specific parameters

461    denote the shared characteristics among all variants with the same annotation, which are

462    also called enrichment parameters.

463    **EM-MCMC algorithm**

464        The basic idea of the EM-MCMC algorithm is to segment the whole genome into

465    approximately independent blocks each with 5,000 ~ 10,000 variants; run MCMC algorithm

466    per block with fixed category specific parameter values $(\boldsymbol{\pi}, \boldsymbol{\sigma}^2)$ to obtain posterior estimates

467    of $(\boldsymbol{\beta}, E[\boldsymbol{\gamma}])$ (E-step); then summarize the genome-wide posterior estimates of $(\boldsymbol{\beta}, E[\boldsymbol{\gamma}])$ and

468    update values of $(\boldsymbol{\pi}, \boldsymbol{\sigma}^2)$ by maximizing their posterior likelihoods (M-step). Repeat such

469    EM-MCMC iterations for a few times until the estimates of $(\boldsymbol{\pi}, \boldsymbol{\sigma}^2)$ (maximum a posteriori

470    estimates, i.e., MAPs) converge (Fig S1).

471        We derive the log-posterior-likelihood functions for $(\boldsymbol{\pi}, \boldsymbol{\sigma}^2)$ and the analytical

472    formulas for their MAPs. In addition, we construct their confidence intervals using Fisher

473    information, whose analytical forms are derived for our Bayesian hierarchical model

474    (Supplement Information). In our practical analyses, we find that, in general, with about 5

475    EM iterations, the estimates for $(\boldsymbol{\pi}, \boldsymbol{\sigma}^2)$ would achieve convergence. Our method of

476    conducting GWAS with functional information by using the above Bayesian hierarchical

477   model and EM-MCMC algorithm is referred as "Scalable Functional Bayesian Association"

478   (bfGWAS).

**Convergence diagnosis**

480   Here, the MCMC algorithm is essentially a random walk over all possible linear

481   regression models with combinations of variants, which can start with either a model

482   containing multiple significant variants by sequential conditional analysis or the most

483   significant variant by P-value. In each MCMC iteration, a new model is proposed by

484   including an additional variant, or deleting one variant from the current model, or switching

485   one variant within the current model with one outside; and then up to acceptation or

486   rejection by the Metropolis-Hastings algorithm (Text S1). Importantly, we refine the

487   standard proposal strategy for the switching step, by prioritizing variants in the

488   neighborhood of the switch candidate according to their conditional association evidence

489   (e.g., P-values conditioning on variants, except the switch candidate, in the current model).

490   As a result, this MCMC algorithm encourages our method to explore different combinations

491   of potentially causal variants in each locus, and significantly improves the mixing property.

492   We used the potential scale reduction factor (PSRF) [56] to quantitatively diagnose

493   MCMC mixing property. PSRF is essentially a ratio between the average within-chain

494   variance of the posterior samples and the overall-chain variance with multiple MCMC

495   chains. From the example plots of the PSRFs of Bayesian PPs (Fig S2), for 58 top

496   marginally significant SNPs (with P-values $< 5 \times 10^{-8}$) in the WTCCC GWAS data of

497   Crohn's disease [1], we can see that about half of the PSRF values by the standard MCMC

498   algorithm (used in GEMMA [35]) exceed 1.2, suggesting the standard MCMC algorithm has

499   poor mixing property. In contrast, the PSRF values by our MCMC algorithm are within the

500   range of (0.9, 1.2), suggesting that our MCMC algorithm has greatly improved mixing

501   property.

**Computational technics**

502

503        We employ two computational technics to save memory in the bfGWAS software.

504 One is to save all genotype data as unsigned characters in memory, because unsigned

505 characters are equivalent to unsigned integers in (0, 256) that can be easily converted to

506 genotype values within the range of (0.0, 2.0) by multiplying with 0.01. This technic saves

507 up to 90% memory comparing with saving genotypes in double type. Second, with an

508 option of in-memory compression, bfGWAS will further save additional 70% memory. As a

509 result, we can decrease the memory usage from ~120 GB (usage by GEMMA[35]) to

510 ~3.6GB, for a typical GWAS dataset with ~33K individuals and ~500K variants.

511        The bfGWAS software wraps a C++ executable file for the E-step (MCMC algorithm)

512 and an R script for the M-step together by a Makefile, which is generated by a Perl script

513 and enables parallel computation through submitting jobs. Generally, 50K MCMC iterations

514 with ~5K variants and ~33K individuals take about 300MB memory and 1hr CPU time on a

515 1.6GHz core, where the computation cost is of order $O(nm^2)$ with the sample size ($n$) and

516 number of variants ($m$) considered in the linear models during MCMC iterations (usually

517 $m < 10$). The computation cost for M-step is almost negligible because of analytical

518 formulas for the MAPs.

**Fgwas**

519

520        In this paper, the Fgwas results were generated by using summary statistics from

521 single variant likelihood-ratio tests and the same annotation information used by bfGWAS.

522 Fgwas[26] produces variant-specific posterior association probabilities (PPs), segment-

523 specific PPs, and enrichment estimates for all annotations. To avoid the issue of failing

524 convergence, we used segment size of 2,000 variants for Fgwas in both simulations and

525 real data analyses. As a result, the final Fgwas PP is given by the product of the variant-

526   specific PP and the corresponding segment–specific PP, and the Fgwas regional-PP is

527   given by the highest segment-specific PP in a region or genome block.

528   **Simulation data**

529   We used genotype data on Chromosome 20-22 from the AMD GWAS (33,976

530   individuals and 241,500 variants with MAF>0.1) to simulate quantitative phenotypes from

531   the standard linear regression model $y_i = X_i^T \beta + \epsilon_i, \ i = 1, \dots, 33976$, where $X_i$ is the

532   genotype vector of the $ith$ individual and $\epsilon_i$ is the noise term generated from $N(0, \sigma_\epsilon^2)$. We

533   segmented the genotype data into 50x2.5Mb blocks each with ~5,000 variants. Within each

534   block, we marked a ~25Kb continuous region (starting 37.5Kb from the beginning of a

535   block) as the causal locus and randomly selected two causal SNPs per locus. Two

536   complementary annotations ("coding" vs. "noncoding") were simulated, where the coding

537   variants account for ~1% overall variants but ~10% variants within the causal loci (matching

538   the pattern in the real AMD analysis). We selected positive effect-size vector $\beta$ and noise

539   variance $\sigma_\epsilon^2$ such that a total of 15% phenotypic variance was equally explained by causal

540   SNPs. We controlled the enrichment-fold of coding variants by varying the number of

541   coding variants among these 100 causal SNPs.

542   We compared bfGWAS with P-value, conditioned P-value, and Fgwas. In the

543   simulation studies, P-values were obtained from a series of likelihood-ratio tests based on

544   the standard linear regression model. P-values conditioning on the top significant variant

545   per locus were used to identify the second signal by conditional analysis. Fgwas was

546   implemented with summary statistics from single variant tests and the segment size of

547   2,000 variants (selected to avoid convergence issues). We failed to include PAINTOR in the

548   comparison, because PAINTOR cannot complete the analysis for one block in >1,000 CPU

24

549  hours (on a 2.5GHz, 64-bit CPU) and is thus expected to require >1 million CPU hours for a

550  genome-wide analysis.

551  **GWAS data of AMD**

552  In the GWAS data of AMD, the advanced AMD cases – including wet cases with

553  choroidal neovascularization (CNV, when accompanied by angiogenesis) and dry cases

554  with geographic atrophy (GA, when angiogenesis is absent) – and control subjects were

555  gathered across 26 studies, with DNA samples collected and genotyped centrally [39]. All

556  genotypes were generated by a customized chip that contains (i) the usual genome-wide

557  variant content, (ii) exome content comparable to the Exome chip (protein-altering variants

558  across all exons), (iii) variants in known AMD risk loci (protein-altering variants and

559  previously associated variants), and (iv) previously observed and predicted variation in

560  *TIMP3* and *ABCA4* (two genes implicated in monogenic retinal dystrophies). The genotyped

561  variants (439,350) were then imputed to the 1000 Genomes reference panel (Phase I) [40],

562  resulting a total of 12,023,830 variants.

563  The software bfGWAS used dosage genotype data and standardized phenotypes.

564  Phenotypes were first coded quantitatively with 1's for cases and 0's for controls; second

565  corrected for the first and second principle components, age, gender, and source of DNA

566  samples; and then standardized to have mean 0 and standard deviation 1. In order to make

567  the Bayesian inferences scalable to the AMD GWAS data (33,976 individuals, 9,866,744

568  variants with MAF >0.5%), we segmented the whole genome into 1,063 non-overlapped

569  blocks, such that each block has length ~2.5Mb (containing ~10,000 variants) and all

570  previously identified loci along with variants in LD ($R^2$ >0.1) were not split. Then we applied

571  the EM-MCMC algorithm with 5 EM steps and 50,000 MCMC iterations (including 50,000

572  extra burn-ins).

25

573        For comparison, P-values were obtained by a series of likelihood-ratio tests, using

574    the same "quantitative" phenotype vector as used by bfGWAS; Fgwas was implemented

575    with the summary statistics from single variant tests and the segment size of 2,000 variants

576    (resulting 4,934 segments); and a standard Bayesian variable selection regression (BVSR)

577    method that models no functional information was also applied.

578        Three types of genomic annotations were considered for analyzing the AMD data:

579    gene-based functional annotations of SNPs and small indels from SeattleSeq

580    (http://snp.gs.washington.edu/SeattleSeqAnnotation138/index.jsp), summarized regulatory

581    annotations [41], and the chromatin states profiled respectively in nine human cell types

582    from chromHMM [19,42,43]. For variants annotated with multiple functions, we used the

583    most severe function in the analysis: non-synonymous > coding-synonymous > other-

584    genomic > intronic > intergenic for the gene-based annotations; coding > UTR > promoter >

585    DHS > intronic > intergenic > "others" for the summarized regulatory annotations; active

586    promoter > poised promoter > strong enhancer > weak enhancer > insulator > transcription

587    elongation > CNV for the chromatin states.

588 **Software**

589        Our software bfGWAS is freely available on Github

590    (https://github.com/yjingj/bfGWAS).

591

## Figure captions

593 Fig1. Power comparison among bfGWAS, Fgwas, P-value, and conditional analysis (CA),
594 with 100 simulation replicates and the complete sample size 33,976. (a) Average ROC
595 curves. (b) Boxplot of the ranks of the true causal SNP1 (with smaller P-value) and SNP2.

596 Fig2. ZoomLocus plots around *rs4151667* in the locus *C2/CFB/SKIV2L* using the
597 association criteria by P-value, standard BVSR, Fgwas, and bfGWAS. (a) –log10(P-values)
598 by single variant tests. (b) Bayesian PPs by BVSR. (c) Fgwas PPs. (d) Bayesian PPs by

599 bfGWAS. The top cyan squares in panels (a, b, c) denote the intronic variant *rs116503776*;
600 the purple triangle in (d) denotes the non-synonymous variant *rs4151667*; shapes denote
601 different annotations (triangle point up Δ for non-syn, circle o for coding-syn, square □ for
602 intronic, diamond ◊ for intergenic, and triangle point down ∇ for other-genomic).

603 Fig3. Category specific (enrichment) parameter estimates with 95% error bars by bfGWAS,
604 with gene-based annotations, regulatory annotations, and chromatin-states profiled in the
605 K562 cell line. (a, c, e) Causal probabilities. (b, d, f) Effect-size variances. The estimates
606 that are the same as their priors are not plotted, e.g., estimates of UTR in (c, d), estimates
607 of the active/poised promoter in (e, f). The estimate of the effect-size variance for the
608 "Others" category in (d) is also close to the prior because of low region-association
609 evidence, hence it has a wide 95% error bar.

610

611

## Supporting information

613 Table S1-S17. Tables S1-S17.

614 Text S1. Supplementary text containing more details about the bfGWAS method.

615 Fig S1. Flowchart of the bfGWAS. (a) Bayesian hierarchical model. (b) Diagram of the EM-
616 MCMC algorithm.

617 Fig S2. Plots of the PSRF values for the Bayesian PPs of 58 top marginally significant
618 SNPs (using the WTCCC GWAS data of Crohn's disease) with 3, 8, 15, and 20 MCMC
619 chains. (a) Standard MCMC algorithm as in GEMMA. (b) Our MCMC algorithm in bfGWAS.
620 PSRF within (0.9, 1.2) suggests good mixing property.

621 Fig S3. Average ROC curves of 100 simulation studies by P-value (single variant test),
622 Fgwas, and bfGWAS with various sample sizes.

623 Fig S4. Prioritized ranks of the true causal SNP1 (pink) and SNP2 (cyan) by P-value (single
624 variant test), conditional analysis (CA), Fgwas, and bfGWAS, stratified by the LD between
625 SNP1 and SNP2. Our bfGWAS method outperforms the other three methods for identifying
626 SNP2 in all scenarios, while the relative performance of the conditional analysis and Fgwas
627 for identifying SNP2 depends on LD. In particular, when SNP1 and SNP2 are in low LD, the
628 conditional analysis outperforms Fgwas. This presumably is due to the wrong Fgwas
629 assumption of one causal per locus, and the reduced association strength for SNP2 by
630 conditioning on the top significant variant (SNP1 or a proxy for SNP1) that is in high LD with
631 SNP2.

632 Fig S5. Estimates of Fgwas enrichment parameters and log-relative-risk $\ln(\pi_0/\pi_1)$ by
633 bfGWAS, along with 95% confidence intervals. (a) Simulation Scenario I with ~44x
634 enrichment in coding. (b) Scenario II with no enrichment. No enrichment is estimated when
635 the 95% confidence interval covers 0, while enrichment for coding is estimated with the
636 95% confidence interval above 0.

637 Fig S6. Estimates of the log-ratio of effect-size variances $\ln(\sigma_0^2/\sigma_1^2)$ by bfGWAS, along with
638 95% confidence intervals. (a) Scenario I with ~44x enrichment in coding. (b) Scenario II with

639  no enrichment. The effect-sizes of both groups in Scenarios I and II were simulated from
640  the same normal distribution, thus the 95% confidence intervals covering 0 suggest that
641  bfGWAS estimates similar effect-size variances between two categories.

642  Fig S7. Sorted top Bayesian PP by bfGWAS versus the sorted top -log10(P-values) of
643  single variant tests for 100 GWASs with AMD genotype data and permuted phenotypes.
644  Here, the P-value $5 \times 10^{-8}$ roughly corresponds to Bayesian PP 0.1068.

645  Fig S8. Manhattan plot with -log10(P-values) by single variant tests for the GWAS of AMD,
646  where variants with Bayesian PP >0.1068 by standard Bayesian variable selection
647  regression (BVSR, modeling no functional information) are colored.

648  Fig S9. Manhattan plots with -log10(P-values) by single variant tests for the GWAS of AMD
649  with gene-based annotations, where different shapes denote different functional
650  annotations. (a) Variants with Bayesian PP >0.1068 bfGWAS are colored. (b) Variants with
651  Fgwas PP >0.1068 are colored.

652  Fig S10. ZoomLocus plots around *rs147859257* in the locus *C3* using the association
653  criteria by P-value, standard BVSR, Fgwas, and bfGWAS. (a) –log10(P-values) by single
654  variant tests. (b) Bayesian PPs by BVSR. (c) Fgwas PPs. (d) Bayesian PPs by bfGWAS.
655  The purple triangle in (b, d) denotes variant *rs147859257*; the blue triangle in (a, c) denotes
656  the top significant variant (*rs2230199*) by P-value.

657  Fig S11. Fgwas enrichment estimates with 95% error bars, using various functional
658  annotations. (a) Gene-based annotations. (b) Regulatory annotations. (c) Chromatin states
659  profiled in the K562 cell line.

660  Fig S12. Ratios of enrich parameters versus the respective genome-wide averages, along
661  with 95% confidence intervals, using various functional annotations. (a, c, e) Causal
662  probability ratios $\pi_q/\pi_{avg}$. (b, d, f) Effect-size variance ratios $\sigma_q^2/\sigma_{avg}^2$.

663  Fig S13. Enrichment parameter estimates with 95% error bars by bfGWAS and Fgwas, with
664  chromatin states profiled in different cell lines. (a, b, c) H1-hESC cell line. (d, e, f) GM12878
665  cell line. (g, h, i) HepG2 cell line. Missing bfGWAS estimates are due to none enrichment
666  (estimates are the same as their priors), while missing Fgwas estimates are due to
667  convergence issues.

668  Fig S14. Enrichment parameter estimates with 95% error bars by bfGWAS and Fgwas, with
669  chromatin states profiled in different cell lines. (a, b, c) HUVEC cell line. (d, e, f) HSMM cell
670  line. (d, h, i) NHLF cell line.

671  Fig S15. Enrichment parameter estimates with 95% error bars by bfGWAS and Fgwas, with
672  chromatin states profiled in different cell lines. (a, b, c) NHEK cell line. (d, e, f) NHEK cell
673  line.
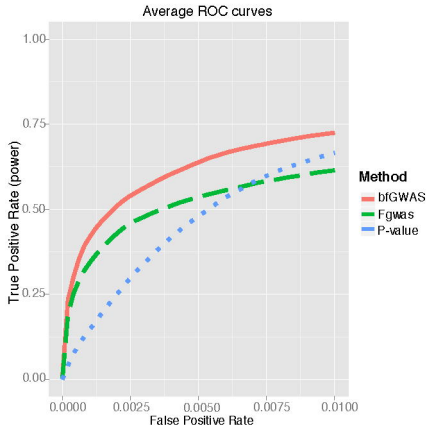
674

675

676  **References**

677    1. Wellcome Trust Case Control C (2007) Genome-wide association study of 14,000 cases of seven
678         common diseases and 3,000 shared controls. Nature 447: 661-678.
679    2. McCarthy MI, Abecasis GR, Cardon LR, Goldstein DB, Little J, et al. (2008) Genome-wide
680         association studies for complex traits: consensus, uncertainty and challenges. Nat Rev
681         Genet 9: 356-369.
682    3. Voight BF, Scott LJ, Steinthorsdottir V, Morris AP, Dina C, et al. (2010) Twelve type 2 diabetes
683         susceptibility loci identified through large-scale association analysis. Nat Genet 42: 579-
684         589.
685    4. Visscher PM, Brown MA, McCarthy MI, Yang J (2012) Five years of GWAS discovery. Am J Hum
686         Genet 90: 7-24.
687    5. Global Lipids Genetics C, Willer CJ, Schmidt EM, Sengupta S, Peloso GM, et al. (2013) Discovery
688         and refinement of loci associated with lipid levels. Nat Genet 45: 1274-1283.
689    6. Hirschhorn JN, Daly MJ (2005) Genome-wide association studies for common diseases and
690         complex traits. Nat Rev Genet 6: 95-108.
691    7. Yu J, Pressoir G, Briggs WH, Vroh Bi I, Yamasaki M, et al. (2006) A unified mixed-model method
692         for association mapping that accounts for multiple levels of relatedness. Nat Genet 38:
693         203-208.
694    8. Hindorff LA, Sethupathy P, Junkins HA, Ramos EM, Mehta JP, et al. (2009) Potential etiologic
695         and functional implications of genome-wide association loci for human diseases and
696         traits. Proc Natl Acad Sci U S A 106: 9362-9367.
697    9. Yang J, Ferreira T, Morris AP, Medland SE, Genetic Investigation of ATC, et al. (2012)
698         Conditional and joint multiple-SNP analysis of GWAS summary statistics identifies
699         additional variants influencing complex traits. Nat Genet 44: 369-375, S361-363.
700    10. Carithers LJ, Moore HM (2015) The Genotype-Tissue Expression (GTEx) Project. Biopreserv
701         Biobank 13: 307-308.
702    11. Dixon JR, Jung I, Selvaraj S, Shen Y, Antosiewicz-Bourget JE, et al. (2015) Chromatin
703         architecture reorganization during stem cell differentiation. Nature 518: 331-336.
704    12. Kellis M, Wold B, Snyder MP, Bernstein BE, Kundaje A, et al. (2014) Defining functional DNA
705         elements in the human genome. Proc Natl Acad Sci U S A 111: 6131-6138.
706    13. Kumar P, Henikoff S, Ng PC (2009) Predicting the effects of coding non-synonymous variants
707         on protein function using the SIFT algorithm. Nat Protoc 4: 1073-1081.
708    14. Adzhubei I, Jordan DM, Sunyaev SR (2013) Predicting functional effect of human missense
709         mutations using PolyPhen-2. Curr Protoc Hum Genet Chapter 7: Unit7 20.
710    15. Pickrell JK, Marioni JC, Pai AA, Degner JF, Engelhardt BE, et al. (2010) Understanding
711         mechanisms underlying human gene expression variation with RNA sequencing. Nature
712         464: 768-772.
713    16. Tung J, Zhou X, Alberts SC, Stephens M, Gilad Y (2015) The genetic architecture of gene
714         expression levels in wild baboons. Elife 4.
715    17. Lea AJ, Tung J, Zhou X (2015) A Flexible, Efficient Binomial Mixed Model for Identifying
716         Differential DNA Methylation in Bisulfite Sequencing Data. PLoS Genet 11: e1005650.
717    18. Pique-Regi R, Degner JF, Pai AA, Gaffney DJ, Gilad Y, et al. (2011) Accurate inference of
718         transcription factor binding from DNA sequence and chromatin accessibility data. Genome
719         Res 21: 447-455.
720    19. Ernst J, Kellis M (2012) ChromHMM: automating chromatin-state discovery and
721         characterization. Nat Methods 9: 215-216.
722    20. McVicker G, van de Geijn B, Degner JF, Cain CE, Banovich NE, et al. (2013) Identification of
723         Genetic Variants That Affect Histone Modifications in Human Cells. Science 342: 747-749.

29

724   21. Cooper GM, Stone EA, Asimenos G, Program NCS, Green ED, et al. (2005) Distribution and
725         intensity of constraint in mammalian genomic sequence. Genome Res 15: 901-913.
726   22. Kircher M, Witten DM, Jain P, O'Roak BJ, Cooper GM, et al. (2014) A general framework for
727         estimating the relative pathogenicity of human genetic variants. Nat Genet 46: 310-315.
728   23. Finucane HKaB-S, Brendan and Gusev, Alexander and Trynka, Gosia and Reshef, Yakir and
729         Loh, Po-Ru and Anttila, Verneri and Xu, Han and Zang, Chongzhi and Farh, Kyle and Ripke,
730         Stephan and Day, Felix R and ReproGen Consortium and Schizophrenia Working Group of
731         the Psychiatric Genomics Consortium and The RACI Consortium and Purcell, Shaun and
732         Stahl, Eli and Lindstrom, Sara and Perry, John R B and Okada, Yukinori and Raychaudhuri,
733         Soumya and Daly, Mark J and Patterson, Nick and Neale, Benjamin M and Price, Alkes L
734         (2015) Partitioning heritability by functional annotation using genome-wide association
735         summary statistics. Nat Genet 47: 1228--1235.
736   24. Zhou X (2016) A Unified Framework for Variance Component Estimation with Summary
737         Statistics in Genome-wide Association Studies. bioRxiv.
738   25. Schork AJ, Thompson WK, Pham P, Torkamani A, Roddey JC, et al. (2013) All SNPs are not
739         created equal: genome-wide association studies reveal a consistent pattern of enrichment
740         among functionally annotated SNPs. PLoS Genet 9: e1003449.
741   26. Pickrell JK (2014) Joint analysis of functional genomic data and genome-wide association
742         studies of 18 human traits. Am J Hum Genet 94: 559-573.
743   27. Kichaev G, Yang WY, Lindstrom S, Hormozdiari F, Eskin E, et al. (2014) Integrating functional
744         data to prioritize causal variants in statistical fine-mapping studies. PLoS Genet 10:
745         e1004722.
746   28. Gabriel SB, Schaffner SF, Nguyen H, Moore JM, Roy J, et al. (2002) The structure of haplotype
747         blocks in the human genome. Science 296: 2225-2229.
748   29. Wall JD, Pritchard JK (2003) Haplotype blocks and linkage disequilibrium in the human
749         genome. Nat Rev Genet 4: 587-597.
750   30. Berisa T, Pickrell JK (2016) Approximately independent linkage disequilibrium blocks in
751         human populations. Bioinformatics 32: 283-285.
752   31. Fritsche LG, Igl W, Bailey JN, Grassmann F, Sengupta S, et al. (2015) A large genome-wide
753         association study of age-related macular degeneration highlights contributions of rare and
754         common variants. Nat Genet.
755   32. Chipman H, George EI, McCulloch RE (2001) The Practical Implementation of Bayesian Model
756         Selection. In: Lahiri P, editor. Model selection. Beachwood, OH: Institute of Mathematical
757         Statistics. pp. 65-116.
758   33. Guan Y, Stephens M (2011) Bayesian variable selection regression for genome-wide
759         association studies and other large-scale problems. 1780-1815.
760   34. Carbonetto P, Stephens M (2013) Integrated enrichment analysis of variants and pathways in
761         genome-wide association studies indicates central role for IL-2 signaling genes in type 1
762         diabetes, and cytokine signaling genes in Crohn's disease. PLoS Genet 9: e1003770.
763   35. Zhou X, Carbonetto P, Stephens M (2013) Polygenic modeling with bayesian sparse linear
764         mixed models. PLoS Genet 9: e1003264.
765   36. Marchini J, Howie B, Myers S, McVean G, Donnelly P (2007) A new multipoint method for
766         genome-wide association studies by imputation of genotypes. Nat Genet 39: 906-913.
767   37. Wen X, Stephens M (2014) Bayesian Methods for Genetic Association Analysis with
768         Heterogeneous Subgroups: From Meta-Analyses to Gene-Environment Interactions. Ann
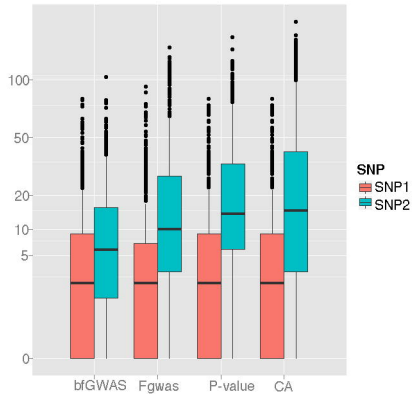769         Appl Stat 8: 176-203.

770  38. Zhou X, Stephens M (2012) Genome-wide efficient mixed-model analysis for association
771      studies. Nat Genet 44: 821-824.
772  39. Fritsche LG, Igl W, Cooke Bailey JN, Grassman F, Sengupta S, et al. (in press) Insights into Rare
773      and Common Genetic Variation From a Large Study of Age-Related Macular Degeneration.
774      Nature genetics.
775  40. Genomes Project C, Auton A, Brooks LD, Durbin RM, Garrison EP, et al. (2015) A global
776      reference for human genetic variation. Nature 526: 68-74.
777  41. Gusev A, Lee SH, Trynka G, Finucane H, Vilhjalmsson BJ, et al. (2014) Partitioning heritability
778      of regulatory and cell-type-specific variants across 11 common diseases. Am J Hum Genet
779      95: 535-552.
780  42. Ernst J, Kellis M (2010) Discovery and characterization of chromatin states for systematic
781      annotation of the human genome. Nat Biotechnol 28: 817-825.
782  43. Ernst J, Kheradpour P, Mikkelsen TS, Shoresh N, Ward LD, et al. (2011) Mapping and analysis
783      of chromatin state dynamics in nine human cell types. Nature 473: 43-49.
784  44. Masson D, Jiang XC, Lagrost L, Tall AR (2009) The role of plasma lipid transfer proteins in
785      lipoprotein metabolism and atherogenesis. J Lipid Res 50 Suppl: S201-206.
786  45. Kettunen J, Tukiainen T, Sarin AP, Ortega-Alonso A, Tikkanen E, et al. (2012) Genome-wide
787      association study identifies multiple loci influencing human serum metabolite levels. Nat
788      Genet 44: 269-276.
789  46. Nikpay M, Goel A, Won HH, Hall LM, Willenborg C, et al. (2015) A comprehensive 1,000
790      Genomes-based genome-wide association meta-analysis of coronary artery disease. Nat
791      Genet 47: 1121-1130.
792  47. Fransen E, Bonneux S, Corneveaux JJ, Schrauwen I, Di Berardino F, et al. (2015) Genome-wide
793      association analysis demonstrates the highly polygenic character of age-related hearing
794      impairment. Eur J Hum Genet 23: 110-115.
795  48. Helgason H, Sulem P, Duvvari MR, Luo H, Thorleifsson G, et al. (2013) A rare nonsynonymous
796      sequence variant in C3 is associated with high risk of age-related macular degeneration.
797      Nat Genet 45: 1371-1374.
798  49. Seddon JM, Yu Y, Miller EC, Reynolds R, Tan PL, et al. (2013) Rare variants in CFI, C3 and C9 are
799      associated with high risk of advanced age-related macular degeneration. Nat Genet 45:
800      1366-1370.
801  50. Zhan X, Larson DE, Wang C, Koboldt DC, Sergeev YV, et al. (2013) Identification of a rare
802      coding variant in complement 3 associated with age-related macular degeneration. Nat
803      Genet 45: 1375-1379.
804  51. Loh PR, Bhatia G, Gusev A, Finucane HK, Bulik-Sullivan BK, et al. (2015) Contrasting genetic
805      architectures of schizophrenia and other complex diseases using fast variance-
806      components analysis. Nat Genet 47: 1385-1392.
807  52. Jordan MI, Ghahramani Z, Jaakkola TS, Saul LK (1999) An Introduction to Variational Methods
808      for Graphical Models. Machine Learning 37: 183-233.
809  53. Carbonetto P, Stephens M (2012) Scalable Variational Inference for Bayesian Variable
810      Selection in Regression, and Its Accuracy in Genetic Association Studies. 73-108.
811  54. Rue H, Martino S, Chopin N (2009) Approximate Bayesian inference for latent Gaussian
812      models by using integrated nested Laplace approximations. Journal of the Royal Statistical
813      Society: Series B (Statistical Methodology) 71: 319-392.
814  55. Singh SaW, Michael and McCallum, Andrew (2012) Monte Carlo MCMC: efficient inference by
815      approximate sampling: Association for Computational Linguistics. 1104-1113 p.

816     **56. Gelman A, Rubin DB (1992) Inference from Iterative Simulation Using Multiple Sequences.**
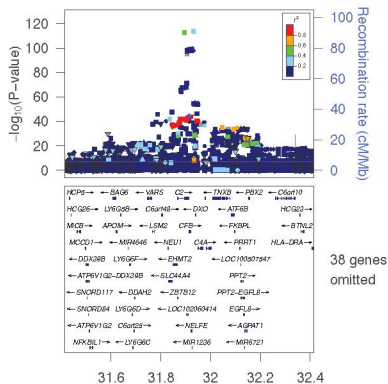817          **Statistical Science 7: 457-472.**
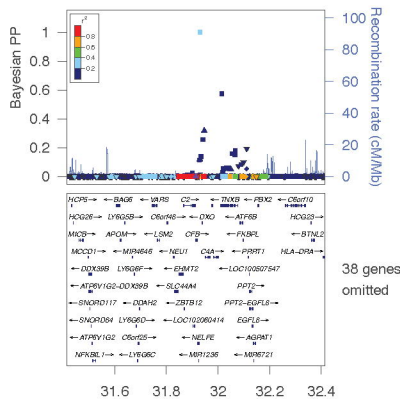
(a)

Average ROC curves
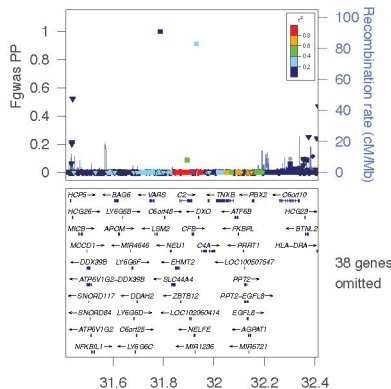
(b)

(a)

Locus Around rs4151667
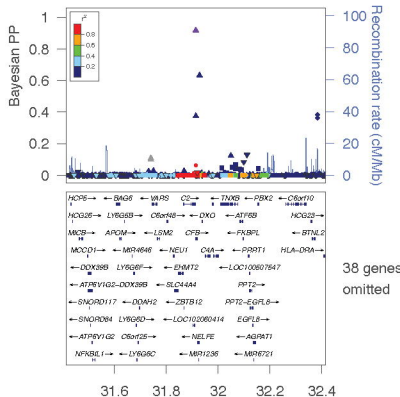
(b)

Locus Around rs4151667
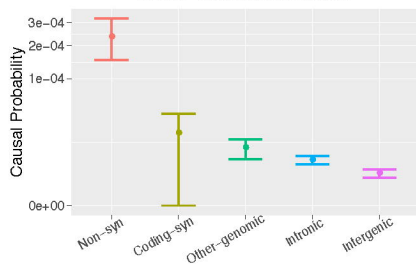
(c)

Locus Around rs4151667
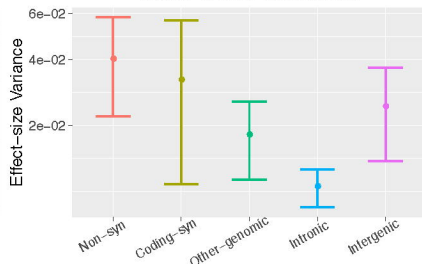
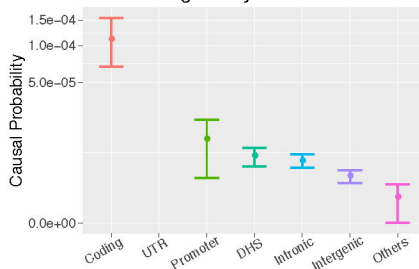(d)

Locus Around rs4151667
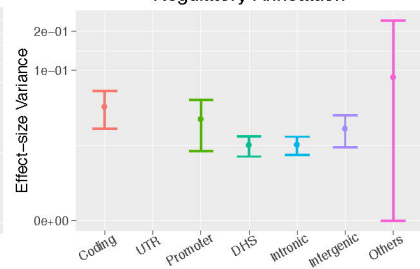
(a) Gene-based Annotation
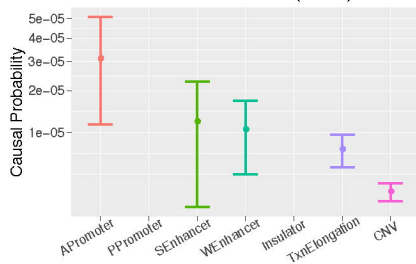
(b) Gene-based Annotation

(c) Regulatory Annotation

(d) Regulatory Annotation

(e) Chromatin States (K562)

(f) Chromatin States (K562)