1    **BayesFM: a software program to fine-map multiple causative variants in**

2    **GWAS identified risk loci.**

3    *Ming Fang[1,2] & Michel Georges[1]*

4    [1]Unit of Animal Genomics, GIGA-R & Faculty of Veterinary Medicine, University of

5    Liège (B34), 1 Avenue de l'Hôpital, 4000-Liège

6    [2]Life Science College, Heilongjiang Bayi Agricultural University, 163319-Daqing

7    Correspondence: michel.georges@ulg.ac.be

8

9    **Absract**

10    We herein describe a new method to fine-map GWAS-identified risk loci based

11    on the Bayesian Least Absolute Shrinkage Selection Operator (LASSO) combined

12    with a Monte Carlo Markov Chain (MCMC) approach, and corresponding

13    software package (BayesFM). We characterize the performances of BayesFM

14    using simulated data, showing that it outperforms standard forward selection

15    both in terms of sensitivity and specificity. We apply the method to the *NOD2*

16    locus, a well-established risk locus for Crohn's disease, in which we identify 13

17    putative independent signals.

18    **Introduction**

19    Thousands of risk loci have been identified by Genome Wide Association Studies

20    (GWAS) affecting nearly all studied common complex diseases in humans

21    (Welter et al., 2014, and http://www.ebi.ac.uk/gwas/). However, for the vast

22    majority of risk loci the causative variants and genes remain unknown. This

23    knowledge is essential to fully capitalize on the investments in GWAS, including

24    for the development of improved diagnostic and therapeutic applications.

25    A number of issues complicate the identification of the causative variants by

26    association analysis. The first is that the utilized case-control cohorts are usually

27    genotyped for only a subset of the variants segregating in the population. Ideally,

28    fine-mapping would require sequencing of all cases and controls in the

29    chromosome regions of interest, if not the entire genome. This will remain

30    difficult to achieve, at least in the short term. At present, the best alternative is

31    genotype imputation using for instance the data from the 1,000 Genomes Project

32    as reference set. However, the reliability of the imputed genotypes is not perfect,

33    particularly for low frequency and rare variants. Thus the information content

34    varies between variants. In other words, the effective number of genotypes may

35    vary between variants, precluding fair comparison of the strength of their

36    association.

37    A second issue is the difficulty, when using the most commonly applied single-

38    marker analyses (i.e. testing for disease association one marker at a time), to

39    distinguish the association patterns of the causative variants from that of

40    "passenger" variants that are merely in linkage disequilibrium (LD) with

41    causative variants. In the case of allelic homogeneity (one causative variant

42    only), one "asymptotically" expects the causative variant to show the strongest

43    single-marker association (highest –log(p) value) of all variants.  But in the real

44    world the causative variant may be overshadowed by passenger variants that by

45    chance (or as a result of unaccounted confounding effects), and given the limited

46    size of the case-control cohorts, appear more strongly associated with the

47    disease.  The situation becomes even trickier in the case of allelic heterogeneity

48    (i.e. the segregation of multiple causative risk variants that may or may not be in

49    LD), a scenario that is likely to be very common.  In this case, the lead SNP in

50    single-marker analysis may be a "ghost" variant that is in LD with two or more

51    causative variants, hence being "asymptotically" more strongly associated with

52    the disease than either of them.   Also, in the case of allelic heterogeneity,

53    causative variants are by definition bound to exist amongst the "non lead"

54    variants in single-marker analysis.

55    Improving the mapping resolution by analysis of association requires the

56    development of statistical models that allow inclusion of confounding factors,

57    estimation of the effects of individual variants conditional on the other ones (i.e.

58    multi-marker analysis to distinguish causative from passenger variants), and

59    identification of the best amongst the large number of possible models (i.e. what

60    combination of variants, assumed to be causative, explains the data best).

61    We herein introduce a software package (BayesFM) that uses Bayesian Least

62    Absolute Shrinkage Selection Operator (LASSO) combined with a Monte Carlo

63    Markov Chain (MCMC) to achieve that goal.  After describing the underlying

64    algorithms, we test BayesFM on simulated data and compare its performances

65    with that of a more standard "forward selection" approach.  We then describe the

66    results obtained when applying BayesFM to the *NOD2* locus, a well established

67    risk locus for Inflammatory Bowel Disease (f.i. Jostins et al., 2015; Huang et al.,

68    2016).

69

70    **Materials & Methods**

71    **BayesFM algorithm**

72    ***Assumptions.*** We assume that GWAS studies have identified one or more risk

73    loci for a disease of interest in an available case-control cohort.   We further

74    assume that – within the identified risk loci - the genotypes of array-interrogated

75    SNPs have been augmented in cases and controls with genotypes of as many

76    variants as possible either by imputation or by sequencing.  We herein propose

77    an approach that will model disease outcome as a function of the genotypes at

78    one or more variants in a given risk locus with the aim to fine map that locus, i.e.

79    to identify causative variants within that locus.  Fine-mapping is conducted one

80    risk locus at the time.  Risk loci defined by GWAS typically span ~250 Kb, contain

81    ~5 genes (range: 0 to >50) and encompass thousands of common and low

82    frequency variants.

83    ***Model.*** The proposed model is based on the standard assumption of an

84    underlying, normally distributed liability *y* with threshold *t*, such that individuals

85    for which *y > t* are affected and individuals for which *y ≤* t are healthy. We model

86    the liability of individual *i* ($y_i$) as

$$y_i = \mu + \sum_{k}^{s} \alpha_k z_{ik} + \sum_{j=1}^{m} \beta_j x_{ij} + \varepsilon_i$$

87    where $\mu$ is the population mean, $\alpha_k$ is the effect of principal component (PC) $k$ of

88    $s$, $z_{ik}$ is the value of PC $k$ for individual $i$, $\beta_j$ is the effect of variant $j$ of $m$, $x_{ij}$ is the

89    dosage of the alternative allele of variant $j$ for individual $i$, and $\varepsilon_i$ is the residual

90    error term for individual $i$. PCs were included to correct for population

91    stratification, and values of $z_{ik}$ were computed using standard procedures. Other

92    fixed effects could be added to the model in exactly the same way as PCs. $m$, i.e.

93    the number of causative variants in the risk locus, was arbitrarily set at 20,

94    meaning that we did not anticipate more than 20 independent effects per risk

95    region. In other words, we consider that there can be multiple causative variants

96    for each risk region, but that this number cannot exceed 20. The challenge is to

97    find the "at most 20" causative variants amongst the thousands of genotyped or

98    imputed variants in each locus. $\varepsilon_i$ is assumed to be normally distributed with

99    mean 0 and variance $\sigma_E^2$.

100   ***Prior distributions.*** Following Sorensen and Gianola (2002), the values of $t$ and

101   $\sigma_E^2$ are fixed at 0 and 1, respectively. The individual liabilities, $y_i$, are assumed to

102   be normally distributed

$$\pi(y_i) = N\left(y_i \,\middle|\, \mu + \sum_{k=1}^{s} \alpha_k p_{ik} + \sum_{j=1}^{m} \beta_j x_{ij}, 1\right)$$

103   The population mean, $\mu$, and effects of the PCs capturing population stratification

104   are assumed to follow uniform distributions ($\pi(\mu) \propto 1; \pi(\alpha_k) \propto 1$). Following

105   Fang et al. (2012), the prior distribution of the effect of variant $j$ from $m$, $\beta_j$, is

106   assumed to follow a double-exponential distribution:

$$\pi(\beta_j) = \frac{\lambda_j}{2} e^{-\lambda_j |\beta_j|}$$

107    which is factorized in three sub-priors: (i) normally distributed $\pi\left(\beta_j\middle|\tau_j^2\right) =$

108    $N\left(\beta_j\middle|0,\tau_j^2\right)$; (ii) exponentially distributed $\pi\left(\tau_j^2\middle|\lambda_j\right) = \frac{\lambda_j}{2}e^{-\lambda_j^2\tau_j^2/2}$; (iii) gamma

109    distributed $\pi\left(\frac{\lambda_j^2}{2}\right) = \Gamma(\xi,\xi), \xi \rightarrow 0$.

110    ***Posterior distributions for Gibbs sampling.*** Effects $\beta_j$ are sampled from normal

111    distributions with mean

$$\overline{\beta_j} = \left(\sum\nolimits_{i=1}^{n} x_{ij}^2 + \frac{1}{\tau_j^2}\right)^{-1} \sum\nolimits_{i=1}^{n} x_{ij}\left(y_i - \mu - \sum\nolimits_{k=1}^{s}\alpha_k p_{ik} - \sum\nolimits_{l\neq j}^{m}\beta_l x_{il}\right)$$

112    and variance

$$\sigma_{\beta_j}^2 = \left(\sum\nolimits_{i=1}^{n} x_{ij}^2 + \frac{1}{\tau_j^2}\right)^{-1},$$

113    in which $n$ is the total number of analyzed individuals (cases + controls). $1/\tau_j^2$

114    are sampled from inverse Gaussian distributions

$$\pi\left(\tau_j^2\middle|y,\dots\right) = \text{InvGauss}\left(\sqrt{\frac{\lambda_j^2}{\beta_j^2}}, \lambda_j^2\right), j = 1,\dots,m$$

115    The hyper-parameters $\lambda_j^2$ are sampled from gamma distributions

$$\pi\left(\lambda_j^2\middle|y,\dots\right) = \Gamma\left(1,\frac{\tau_j^2}{2}\right), j = 1,\dots,m$$

116    PC effects, $\alpha_k$, are sampled from normal distributions with mean

$$\overline{\alpha_k} = \left(\sum\nolimits_{i=1}^{n} p_{ik}^2\right)^{-1} \sum\nolimits_{i=1}^{n}\left(y_i - \mu - \sum\nolimits_{l\neq k}^{s}\alpha_l p_{il} - \sum\nolimits_{j=1}^{m}\beta_j x_{ij}\right)$$

117    and variance

$$\sigma_{\alpha_k}^2 = \left( \sum_{i=1}^n p_{ik}^2 \right)^{-1}$$

118      The population mean, $\mu$, is sampled from a normal distribution with mean

$$\bar{\mu} = \frac{1}{n} \sum_{i=1}^n \left( y_i - \sum_{k=1}^s \alpha_k p_{ik} - \sum_{j=1}^m \beta_j x_{ij} \right)$$

119      and variance $1/n$.

120      For affected individuals ($\gamma_i = 1$), the liabilities, $y_i$, are sampled from the

121      truncated normal distributions (such that $y_i > t$) with density

$$\pi(y_i | \gamma_i = 1, \dots) = \frac{N\left(y_i \big| \mu + \sum_{k=1}^s \alpha_k p_{ik} + \sum_{j=1}^m \beta_j x_{ij}, 1\right)}{1 - \Phi_t\left(\mu + \sum_{k=1}^s \alpha_k p_{ik} + \sum_{j=1}^m \beta_j x_{ij}, 1\right)}$$

122      For unaffected individuals ($\gamma_1 = 0$), the liabilities, $y_i$, are sampled from the

123      truncated normal distributions (such that $y_i \leq t$) with density

$$(y_i | \gamma_i = 0, \dots) = \frac{N\left(y_i \big| \mu + \sum_{k=1}^s \alpha_k p_{ik} + \sum_{j=1}^m \beta_j x_{ij}, 1\right)}{\Phi_t\left(\mu + \sum_{k=1}^s \alpha_k p_{ik} + \sum_{j=1}^m \beta_j x_{ij}, 1\right)}$$

124      In these, $\Phi_t\left(\mu + \sum_{k=1}^s \alpha_k p_{ik} + \sum_{j=1}^m \beta_j x_{ij}, 1\right)$ corresponds to the cumulative

125      density from $-\infty$ to $t$.

126      ***Variant sampling using the Metropolis-Hastings algorithm.*** We first

127      hierarchically cluster variants that are in high LD using $(1-r^2)$ as distance

128      measure and the "single linkage" approach implemented with the R "hclust"

129      package. By doing so variants in distinct clusters will never have $r^2 > C$, yet

130      variants within clusters may have $r^2 < C$. We tested C values of 0.9 and 0.5. The

131      $m$ (=20) variants to include in the model are sampled such that each cluster can

132      only be represented by one variant. At each round of the MCMC chain, we

133    sequentially attempt to swap each of the *m* variants in the model with a better

134    one, to ultimately find the best overall combination of variants. The substitute

135    variants are selected 50% of the time from variants from the same cluster, and

136    50% of the time from variants of unrepresented clusters. In other words, this

137    means that the MCMC chain spends halve of its time searching for the best

138    possible variants within clusters, and halve of its time for the best possible

139    clusters. When a substitute variant is selected, the probability to "accept" it is

$$\alpha = min\left(1, \frac{\pi\left(y\middle|x_{\cdot j}^{new}, \ldots\right)}{\pi\left(y\middle|x_{\cdot j}^{old}, \ldots\right)}\right)$$

140    where

$$\pi(y) = \frac{1}{\sqrt{2\pi}}\exp\left(\frac{1}{2}\sum\left(y_i - \mu - \sum\alpha_k p_{ik} - \sum\beta_j x_{ij}\right)^2\right)$$

141    and $x_{\cdot j}^{new}$ and $x_{\cdot j}^{old}$ correspond to the genotype indicator variables for the "new"

142    and "old" positions, respectively.

143    ***Implementation of the MCMC chain.*** We initiate the chain by assigning values

144    randomly to all variables (within their legal boundaries). We then sequentially

145    update the position of the *m* variants by either choosing a variant in another,

146    unrepresented cluster (50%) or in the same cluster (50%), using the M-H

147    algorithm described above. The likelihood of the new proposition is computed

148    with the parameter values of the previous cycle. Corresponding $\beta_j$, $\tau_j^{-2}$ and $\lambda_j$

149    are updated by 50 rounds of Gibbs sampling. After each round of update of the *m*

150    variants, we update $\mu$, the $\alpha_k$'s and $y_i$'s by one round of Gibbs sampling. The

151    complete process was repeated 500,000 (simulated data) or 1 million times (real

152    data).  The first 100,000 (simulated data) or 500,000 cycles (real data) were

153    used as burn-in and ignored when compiling the summary statistics.

154    ***Summarizing the results.***   We computed posterior probabilities (PP) for

155    clusters as well as individual variants from the proportion of MCMC cycles in

156    which they were included in the model.   Within clusters, we defined "credible

157    sets" of markers by ranking them on PP and considering the minimal set of

158    markers that would jointly account for 95% of the PP of the cluster.

159    Clusters with posterior probability ≥ 0.50 were retained for further validation by

160    fitting their corresponding lead SNPs jointly in a logistic regression model.

161    Clusters exceeding the set significance threshold were considered to be positive.

162    We used thresholds of $10^{-4}$, $10^{-6}$ and $10^{-8}$ that might be considered as locus-

163    specific, multi-locus (~100 loci; cfr. Huang et al., 2016) and genome-wide

164    thresholds.

165

166    **Forward Selection**

167    The performance of BayesFM was compared with that of a standard forward

168    selection approach implemented by logistic regression in R.   Significance

169    thresholds were the same as defined above ($10^{-4}$, $10^{-6}$ and $10^{-8}$).   We built

170    credible sets associated with selected "lead" (*l*) variants by computing the PP for

171    all *n* variants in high LD ($r^2 > C$, as defined above) with *l*.  The PP probability of

172    variant *j* of *n* was computed as:

$$PP_j = L_j \bigg/ \sum_{i=1}^{n} L_i$$

173     where $L_j$ is the maximum likelihood of the data considering variant $j$ in the

174     model. Likelihoods were computed using the glm() function (binomial family,

175     logit link function, and logLik) in R. Credible sets (associated with a given lead

176     variant) corresponded to the smallest set of variants that would jointly account

177     for 95% of the PP. Variants included in a credible set were ignored when

178     pursuing the forward selection.

179

180     **Datasets**

181     ***Simulated data.*** We took advantage of the Immunochip dataset of the

182     International IBD Genetics Consortium (IIBDGC) and Multiple Sclerosis Genetics

183     Consortium (IMSGC), consisting of 18,967 Crohn's disease cases, 14,628

184     ulcerative colitis cases, and 34,257 controls, all of European ancestry. We

185     randomly selected a genomic region corresponding to a GWAS-identified risk

186     locus for Inflammatory Bowel Disease (chr5: 40,286,967-40.818,088) with 2,978

187     markers either interrogated by the Immunochip (936), or imputed from the

188     1,000 Genomes project with quality score > 0.4 (2042) (Huang et al., 2016).

189     Within this region, we randomly selected one (model I), three (model II) or five

190     (model III) variants with MAF ≥ 0.01, to act as causative variants. The variance

191     explained by the locus ($\sigma_L^2$) was set at 2% (see hereafter). In the cases with

192     multiple causative variants, the proportion of the variance explained by the

193     distinct causative variants was set at 4/7, 2/7 and 1/7 (three variants, model II)

194     or 16/31, 8/31, 4/31, 2/31, 1/31 (five variants, model III).

195    The effect of a causative variant $j$ ($\beta_j$) was assumed to be additive and to have

196    numerical value $\beta_j = \sqrt{\sigma_j^2/\sigma^2(x_{ij})}$, as $\sigma_j^2 = \sigma^2(\beta_j x_{ij})$. The value of $\sigma_j^2$, the

197    variance due to variant $j$, was set as described above, i.e. 2% (model I) or the

198    corresponding fraction of 2% (models II & III). The value of $\sigma^2(x_{ij})$, where $x_{ij}$ is

199    – as before - the genotype dosage if individual $i$ for variant $j$, was computed from

200    the Immunochip data. The procedure described thus far does not account for

201    the LD that may exist between the multiple causative variants in models II&III,

202    which may cause $\sigma_L^2$ to deviate from 2%. Indeed,

$$\sigma_L^2 = \sigma^2\left(\sum_{j=1}^{1,3\ or\ 5} \beta_j x_{ij}\right) = \sum_{j=1}^{1,3\ or\ 5} \sigma^2(\beta_j x_{ij}) + 2\sum_{j<j'}^{1,3\ or\ 5} \beta_j \beta_{j'} covar(x_{ij} x_{ij'})$$

$$\geq \sum_{j=1}^{1,3\ or\ 5} \sigma^2(\beta_j x_{ij})$$

203    Thus, we rescaled the effects $\beta_j$ to $\beta_j^* = \beta_j\sqrt{2\%/\sigma_L^2}$. Substituting $\beta_j^*$ for $\beta_j$ in the

204    previous equation indeed gives $\alpha_L^2 = 2\%$.

205    To generate a simulated case-control cohort we sampled 1 million individuals

206    from the Immunochip dataset with replacement and without discrimination of

207    real case-control status. For each of these, we generated a liability, $y_i$, as (i) the

208    sum of the genotype effects at the 1, 3 or 5 causative variants: $\sum_{j=1}^{1,3or5} \beta_j\, x_{ij}$ , plus

209    (ii) a residual effect, $\varepsilon_i$, drawn from a normal distribution with mean 0 and

210    variance of 1. Thus the variance explained by the locus as a fraction of the total

211    liability variance is in fact $0.02/(1 + 0.02) \approx 0.02$. Assuming an incidence of the

212    disease of 1/400, we kept the 2,500 individuals with highest liability as case

213    cohort. We randomly sampled 2,500 individuals from the remaining 1,000,000-

214    2,500 = 997,500 individuals to serve as controls.  We generated 100 such

215    simulated case-control datasets to compare the performance of BayesFM with

216    that of a more standard forward selection procedure.  When analyzing the

217    simulated datasets, we did not fit PC in the statistical models.

218    ***Real data.***  For the analysis of real data, we likewise took advantage of the

219    Immunochip dataset of the International IBD Genetics Consortium (IIBDGC) and

220    Multiple Sclerosis Genetics Consortium (IMSGC), consisting of 18,967 Crohn's

221    disease cases, and 34,257 controls, all of European ancestry.  We selected the

222    genomic  region  corresponding  to  the  GWAS-identified  risk  locus  for

223    Inflammatory Bowel Disease encompassing the *NOD2* gene (chr16: 50,692,364-

224    50,847,022) with 1,048 markers either interrogated by the Immunochip (283),

225    or imputed from the 1,000 Genomes project with quality score > 0.4 (765)

226    (Huang et al., 2016).  The analysis was restricted to Crohn's disease.

227

228    **Results**

229    ***Simulated data.*** As expected, the True Positive Rate (TPR, i.e. the proportion of

230    true positive signals amongst the total number of true signals (true positives plus

231    false  negatives))  was  decreasing  for  both  methods  (BayesFM  and  FS)  with

232    increasing significance threshold and decreasing variance accounted for by the

233    considered variant (Table 1 and Supplemental Table 1).   Mapping resolution

234    (defined as the size of the credible sets) was comparable between BayesFM and

235    FS and only very mildly affected by significance threshold and variance explained.

236    As expected, it decreased with LD threshold used to define clusters/credible sets

237    (average number of variants per cluster/credible set of ~25 ($r^2$=0.9) versus ~30

238    ($r^2$=0.5)).  Otherwise, LD threshold ($r^2$=0.5 or 0.9) had only very modest effects

239    on TPR .

240    For given thresholds and variance explained, the TPR tended to be slightly better

241    for BayesFM than for FS, especially at the higher significance thresholds, but the

242    differences were modest (Table 1 and Supplemental Table 1).  However, when

243    considering models II and III, characterized by multiple causative variants, the

244    False Discovery Rate (FDR, i.e. the proportion of false positives amongst the total

245    number of positive signals (true positives and false positives)) was considerably

246    higher for FS than for BayesFM (Table 1 and Supplemental Table 1).  Thus, while

247    BayesFM and FS appear to have comparable sensitivity, BayesFM outperforms FS

248    in generating a smaller number of false positives in situations of allelic

249    heterogeneity.

250    To examine whether BayesFM and FS might be complementary and might best

251    be used in combination as done in Huang et al. (2016), we measured the TPR and

252    FDR for a specific scenario (model III, log(1/p) threshold 8) considering (i) both

253    approaches separately, (ii) overlapping findings, and (iii) approach-specific

254    findings.  As expected the TPR was highest when considering both approaches

255    individually.   For the examined scenario, the sensitivity measured by the TPR

256    was 31% for BayesFM and 28% for FS.  The corresponding FDRs were 16% for

257    BayesFM and 24% for FS.    Hence and as mentioned before, BayesFM

258    outperformed FS especially with regards to specificity in this scenario of allelic

259    heterogeneity.  When only considering positive results found by both methods

260    (overlapping findings), the sensitivity dropped only very slightly when compared

261 to FS alone (TPR = 26%), while the specificity increased considerably especially

262 when compared to FS (FDR = 12.5%).  Considering BayesFM-specific findings in

263 addition to the overlapping findings (BayesFM and FS) increased the yield of true

264 positives by nearly 20%, with a still reasonable FDR of 28%.   Considering FS-

265 specific findings in addition to the overlapping finings would only increase the

266 yield of true positives by 9%, with an abysmal FDR of 68% (Figure 1).

267 There are at least two scenarios in which BayesFM is expected to beat Forward

268 Selection.  The first is a situation in which a passenger variant is in LD with two

269 or more causative variants and single-handedly accounts for a higher proportion

270 of the variance than any of the causative variants alone.   Standard Forward

271 Selection will then irreversibly include it in the model, which may subsequently

272 preclude the actual causative variants from entering it.   An example of such a

273 scenario, previously referred to as "ghost" effect, is illustrated in Figure 2A, and

274 Table 2.   It generates both false positives and false negatives with Forward

275 Selection. A second scenario where BayesFM is expected to outperform Forward

276 Selection is when two causative variants are in LD, and the risk alleles are in

277 repulsion hence neutralizing each other effects.   An example of such a situation

278 is shown in Figure 2B and Table 2.

279 ***Real data.*** We then examined the results obtained with BayesFM on the *NOD2*

280 locus (chr16: 50,692,364-50,847,022), a well-established risk locus for Crohn's

281 disease.  We analyzed the dataset of the IIBDGC described in Huang et al. (2016)

282 and comprising 18,967 Crohn's disease cases and 34,257 matched controls.

283 Table 3 summarizes the results that were obtained using either $r^2 > 0.9$ or $r^2 >$

284 0.5 as LD threshold to define clusters of variants (see M&M).  For both analyses,

285    we report the clusters with PP > 0.50. Thirteen such signals were obtained with

286    $r^2 > 0.9$, and fourteen with $r^2 > 0.5$. The average number of variants per signal

287    was 1.15 with $r^2 > 0.9$ and 3.14 with $r^2 > 0.5$. Single variant resolution was

288    obtained for 11/13 signals with $r^2 > 0.9$ and 8/14 signals with $r^2 > 0.5$,

289    highlighting the remarkable resolving power that can be achieved for at least

290    some loci with BayesFM. The log(1/p) value obtained by fitting the lead variant

291    (with highest PP) of each signal in a multivariate logistic regression, exceeded 6

292    for 10/13 signals with $r^2 > 0.9$ and 10/14 signals with $r^2 > 0.9$. The lowest

293    log(1/p) value was 3.33 for the signal that was detected using $r^2 > 0.5$ only. It

294    was > 4 for all others.

295    Using very stringent criteria, nine independent signals were retained for the

296    same locus and reported in Huang et al. (2016)(Table 3). All but one of these

297    were detected by BayesFM, whether using $r^2 > 0.9$ or 0.5. These included four

298    signals corresponding to single variant resolution of non-synonymous (NS)

299    variants in *NOD2* (fs1007insC, R702W, G908R, N289S, all with PP~1), and one

300    signal with two-variant resolution corresponding to two non-synonymous

301    variants each (V793M and S431L; resolved to V793M by BayesFM when setting

302    $r^2 > 0.9$). The remarkable enrichment in NS variants testifies of the specificity of

303    the fine-mapping methods utilized in Huang et al. (2016) including BayesFM.

304    The signal that remained undetected by BayesFM was characterized by a PP of

305    0.24 with $r^2 > 0.9$ and 0.20 with $r^2 > 0.5$. It is worth noting that the corresponding

306    lead variant (*rs104895467*) has a Phastcons conservation score of 1, increasing

307    the likelihood of it being a genuine causative variant.

308    Amongst the six additional putative signals detected by BayesFM when applying

309    more lenient thresholds, three were characterized by a NS *NOD2* lead variants

310    (A585T, R676C and A891D).    This remarkable enrichment in NS variants

311    strongly suggests that at least some of these signals are true.

312    Two of the detected signals are characterized by very common risk alleles (31%

313    and 61% in cases, respectively).  Both of these signals are characterized by fairly

314    large clusters when using either FS or BayesFM ($r^2>0.5$), indicating that many

315    variants are in high LD with the corresponding causative variants.   Four of the

316    signals are characterized by low frequency risk alleles (1% < frequency < 5% in

317    controls) and include the well-known *fs1007incC*, *R702W*, *G908R* and *rs72796367*

318    variants.  The remaining signals correspond to rare risk alleles with frequencies

319    below 1% in controls.

320

321    **Discussion**

322    Identifying causative variants in GWAS-defined risk loci is important in order to

323    gain a better understanding of the molecular mechanisms underlying inherited

324    disease predisposition, including the identification of the causative genes.   A

325    number of fine-mapping strategies have been explored to achieve this goal using

326    association information in case-control cohorts.   These include Bayesian

327    approaches to define credible sets that are likely to contain the causative

328    variants (f.i. Wellcome Trust Case Control Consortium et al., 2012; van de Bunt et

329    al., 2015).  However, the corresponding methods make the unlikely assumption

330    of allelic homogeneity at the considered loci, i.e. the occurrence of single

331  causative variants only.  Searching for multiple independent causative variants

332  in a given locus – a more realistic scenario - is most often conducted using

333  variations of stepwise forward selection approaches.  In these approaches - if

334  deemed significant - the strongest signal is sequentially added as covariate to the

335  model.  Two of the methods utilized in Huang et al. (2016) are advanced

336  Bayesian versions of this approach.  It is relatively easy to imagine scenarios in

337  which these forward selection approaches may either miss true signals or

338  incorporate non-causative variants into the model (see results section on

339  simulated data).  Alternative Bayesian variable selection approaches, combined

340  with Monte Carlo techniques, have been devised to overcome these limitations.

341  In the field of fine-mapping, these include BimBam (Servin & Stephens, 2007),

342  GUESSFM (Wallace et al., 2015), and BayesFM presented in this manuscript.

343  Simulations (including in this study) indicate that they are generally superior to

344  the other approaches.  In our simulations, BayesFM appeared to have improved

345  sensitivity and specificity when compared to a standard implementation of

346  forward selection.  Our analyses confirm the benefits of selecting signals that are

347  detected by forward selection *and* BayesFM, the strategy followed in part by

348  Huang et al. (2016).  It considerably improves specificity with limited impact on

349  sensitivity.  We nevertheless note that the signals detected by BayesFM only

350  (and not FS) appear to be characterized by an acceptable specificity, while FS-

351  only signals were in essence not trustworthy  (Figure 1).

352  Ideally, fine-mapping should be done with complete sequence information in the

353  utilized case-controls.  While this may become possible in the near future as

354  sequencing costs continue to diminish, present studies typically augment

355   genotyping data from SNP arrays with imputation using the for instance the

356   1,000 Genomes Project data as reference (f.i. Huang et al., 2016). The imputation

357   accuracy varies between variants and is typically inferior for low frequency

358   variants. This is very likely affecting the precision of fine-mapping and very

359   difficult to overcome by. An accurate evaluation of the impact of imputation

360   accuracy on the outcome of fine-mapping, including with Bayesian model

361   selection approaches such as BayesFM is needed.

362   In its present version, BayesFM assumes that the identified risk variants operate

363   "additively", i.e. we ignore the possibility of dominance within variants and

364   epistatic interaction between variants. BimBam allows modeling of dominance

365   effect within variants (Servin & Stephens, 2007). The impact of this

366   simplification on power and accuracy remains to be determined. Modeling these

367   higher order effects implies the estimation of additional parameters. As a

368   consequence, impact on power and accuracy is likely to be a function of sample

369   size.

370   The results of our simulations (data simulated under a simplified additive

371   model) indicate that fine-mapping results need to be considered with caution.

372   FDR was > 10% even when considering signals detected both by FS and BayesFM,

373   and this is likely to be an underestimate. Nevertheless, applying FS-type

374   methods in conjunction with BayesFM on a large case-control cohort for

375   Inflammatory Bowel Disease lead to the fine-mapping of 42 signals with

376   resolution < 5 variants. The causality of the corresponding variants, especially

377   when non-coding, can now be tested directly by CRISPR-CAS9 technology in cell

378   culture systems.

379

## Acknowledgments

385

**Software availability.** BayesFM can be downloaded from
387　https://sourceforge.net/projects/bayesfm-mcmc-v1-0

388

## References

390　Fang M, Jiang D, Yang R, Fu W, Pu L, Gao H, Wang G, Yu L. 2012. Improved LASSO

391　priors for shrinkage quantitative trait loci mapping, Theor Appl Genet 124: 1315-

392　1324.

393　Huang H, Fang, Jostins L, et al. 2016. Association mapping of inflammatory bowel

394　disease loci to single variant resolution. Nature, under revision.

395　Jostins L, Ripke S , Weersma R, Duerr R, McGovern D, et al. 2012. Host–microbe

396　interactions have shaped the genetic architecture of inflammatory bowel disease,

397　Nature 491: 119–124.

398　Servin B, Stephens M. 2007. Imputation-based analysis of association studies:

399　candidate regions and quantitative traits. PLoS Genet 3: e114.

400  Sorensen D & Gianola D. 2002. Likelihood, Bayesian and MCMC Methods in

401  Quantitative Genetics. *Springer*.

402  Van de Bunt M, Cortes A, IGAS Consortium, Brown MA, Morris AP, McCarthy MI.

403  2015. Evaluating the performance of fine-mapping strategies at common variant
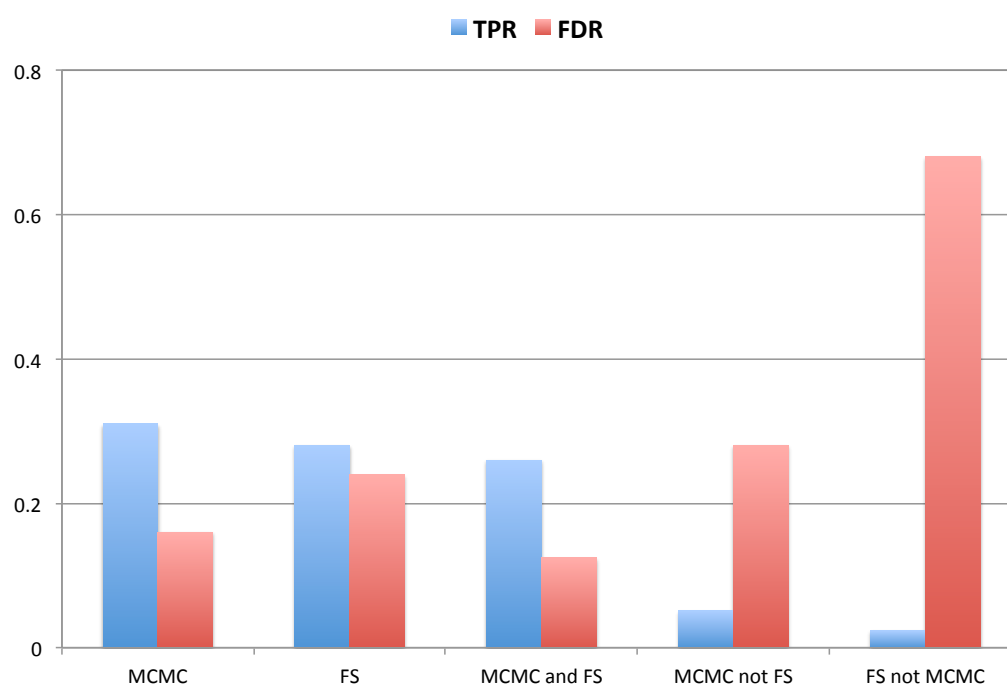
404  GWAS loci. PLoS Genet 11: e1005535.

405  Wallace C, Cutler AJ, Pontikos N, Pekalski ML, Burren OS, Cooper JD, Garcia AR,

406  Ferreira RC, Guo H, Walker NM, Smyth DJ, Rich SS, Onengut-Gumuscu S, Sawcer

407  SJ, Ban M, Richardson S, Toff JA, Wicker LS. 2015. Dissection of a complex disease

408  susceptibility region using a Bayesian stochastic search approach to fine

409  mapping. PLoS Genet 11: e1005272.

410  Wellcome Trust Case Control Consortium, Maller JB, McVean G, Bymes J,

411  Vukcevic D, et al. 2012. Bayesian refinement of association signals for 14 loci in 3

412  common diseases. Nat Genet 44 :1294-1301.

413  Welter D, MacArthur J, Morales J, Burdett T, Hall P, Junkins H, Klemm A, Flicek P,

414  Manolio T, Hindorff L, and Parkinson H. 2014. The NHGRI GWAS Catalog, a

415  curated resource of SNP-trait associations. Nucleic Acids Research 42 (Database

416  issue): D1001-D1006.

417

418

**Figure 1:** Simulated data. True Positive Rates (TPR) and False Discovery rates (FDR) obtained when (i) considering BayesFM (MCMC) and Forward Selection (FS) separately, (ii) when considering overlapping results (MCMC and FS), (iii) when considering method-specific results (MCMC not FS, and FS not MCMC). The scenario under consideration was model III (five causative variants) and a genome-wide significance threshold of log(1/p) = 8.
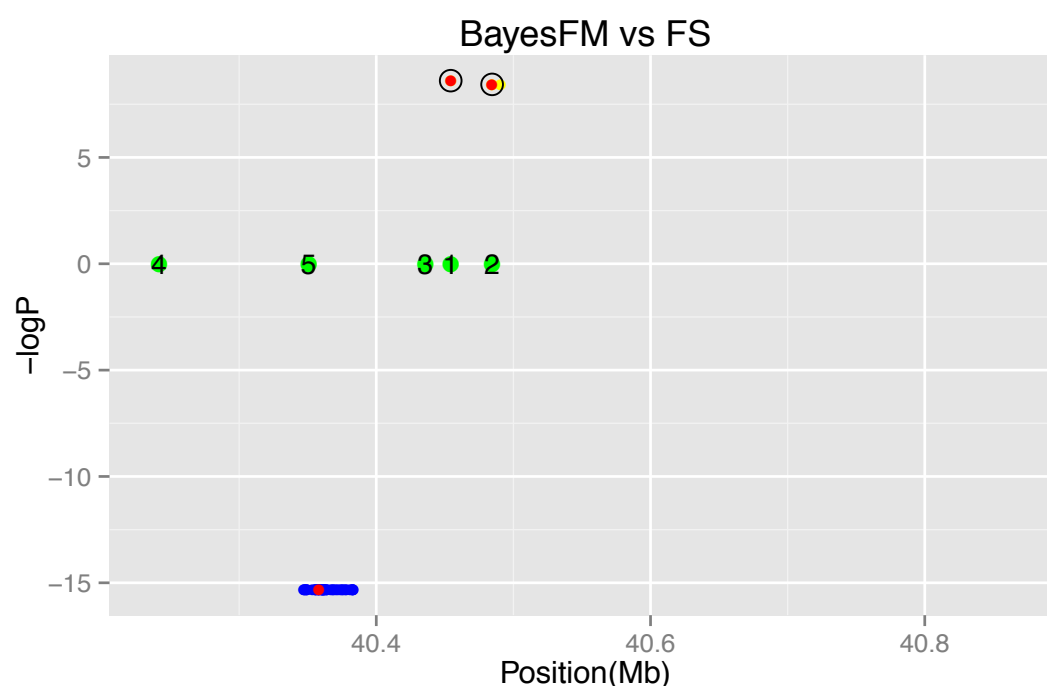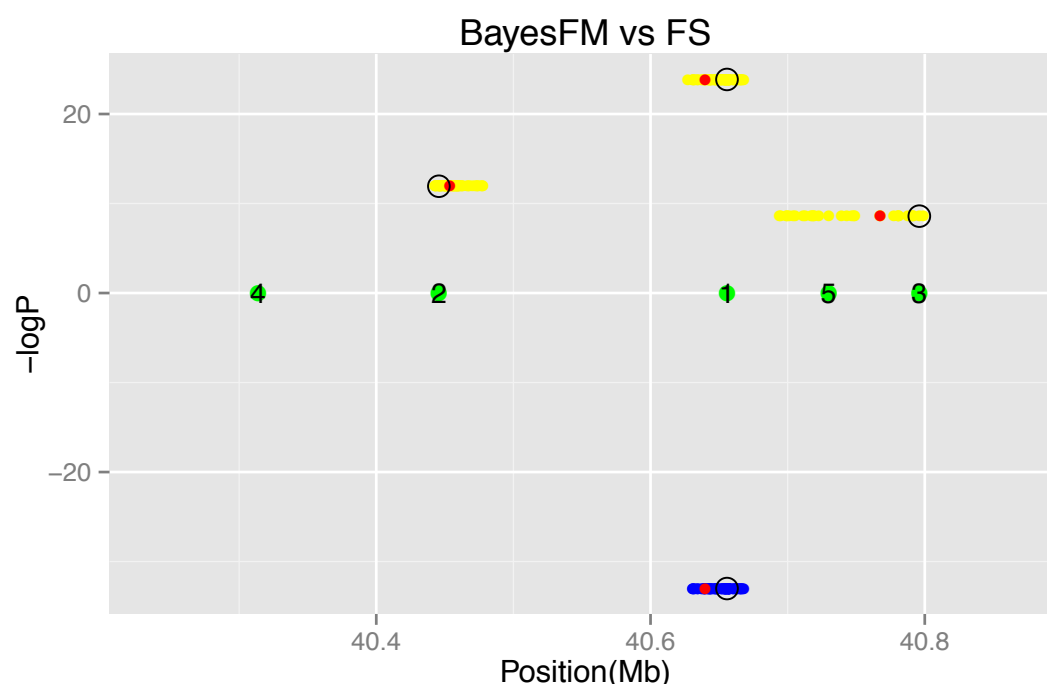
427



428

429 **Figure 2:** Simulated data. Statistical significance of disease association for
430 credible sets of variants selected (PP ≥ 50%) using BayesFM (log(1/p), positive
431 values) or Forward Selection (-log(1/p), negative values), estimated by
432 multivariate logistic regression. The positions of the five simulated causative
433 variants (model III) are shown by the numbered green dots. The yellow and blue
434 dots mark the position of the variants in the credible sets identified by BayesFM
435 and Forward Selection, respectively. The red dots mark the positions of the lead
436 variants in the corresponding credible sets. Causative variants within credible
437 sets are circled. **A.** Example of a "ghost QTL" effect. Forward Selection

438   erroneously identifies a cluster of passenger variants that is in LD with multiple
439   causative variants thereby single-handedly achieving a higher significance than
440   any of the causative variants.   It is therefore erroneously and irreversibly
441   introduced into the forward selection model.   BayesFM avoids this trap and
442   correctly identifies at least causative variants 1 and 2. **B.** Example of two
443   causative variants in LD with an excess of haplotypes with risk alleles in
444   repulsion.  By modeling them simultaneously, BayesFM uncovers risk alleles 2
445   and 3. By modeling them sequentially, Forward Selection misses both 2 and 3 as
446   they neutralize each other's effects.