

1
2
3
4
5

RiVIERA-beta: Joint Bayesian inference of risk variants and tissue-specific epigenomic enrichments across multiple complex human diseases

6
7
8
9
10

Yue Li^{1,2,†} and Manolis Kellis^{1,2,†}

¹Computer Science and Artificial Intelligence Lab, Massachusetts Institute
of Technology, 32 Vassar St, Cambridge, Massachusetts 02139, USA

²The Broad Institute of Harvard and MIT, 415 Main Street, Cambridge,
Massachusetts 02142, USA

[†]Correspondence to liyue@mit.edu or manoli@mit.edu

11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31

Genome wide association studies (GWAS) provide a powerful approach for uncovering disease-associated variants in human, but fine-mapping the causal variants remains a challenge. This is partly remedied by prioritization of disease-associated variants that overlap GWAS-enriched epigenomic annotations. Here, we introduce a new Bayesian model RiVIERA-beta (Risk Variant Inferece using Epigenomic Reference Annotations) for inference of driver variants by modelling summary statistics p-values in Beta density function across multiple traits using hundreds of epigenomic annotations. In simulation, RiVIERA-beta promising power in detecting causal variants and causal annotations, the multi-trait joint inference further improved the detection power. We applied RiVIERA-beta to model the existing GWAS summary statistics of 9 autoimmune diseases and Schizophrenia by jointly harnessing the potential causal enrichments among 848 tissue-specific epigenomics annotations from ENCODE/Roadmap consortium covering 127 cell/tissue types and 8 major epigenomic marks. RiVIERA-beta identified meaningful tissue-specific enrichments for enhancer regions defined by H3K4me1 and H3K27ac for Blood T-Cell specifically in the 9 autoimmune diseases and Brain-specific enhancer activities exclusively in Schizophrenia. Moreover, the variants from the 95% credible sets exhibited high conservation and enrichments for GTEx whole-blood eQTLs located within transcription-factor-binding-sites and DNA-hypersensitive-sites. Furthermore, joint modeling the nine immune traits by simultaneously inferring and exploiting the underlying epigenomic correlation between traits further improved the functional enrichments compared to single-trait models.

32 1 Introduction

33 Genome wide association studies (GWAS) can help gain numerous insights on the genetic
34 basis of complex diseases, and ultimately contribute to personalized risk prediction and pre-
35 cision medicine [1–4]. However, fine-mapping the exact causal variants is challenging due to
36 linkage disequilibrium (LD) and the lack of ability to interpret the function of noncoding
37 variants, which contribute to about 90% of the current GWAS catalog (40.7% intergenic
38 and 48.6% intronic; [5]). On the other hand, several lines of evidence have been proposed
39 to help interpret non-coding genetic signals, in order to gain insights into potential regula-
40 tory functions. In particular, epigenomic annotations can pinpoint locations of biochemical
41 activity indicative of cis-regulatory functions [6, 7]. Indeed, comparison with genome-wide
42 annotations of putative regulatory elements has shown enrichment of GWAS variants in
43 enhancer-associated histone modifications, regions of open chromatin, and conserved non-
44 coding elements [3, 6, 8–12], indicating they may play gene-regulatory roles. These enrich-
45 ments have been used to predict relevant cell types and non-coding annotations for specific
46 traits [6, 9, 13]. Furthermore, many complex traits potentially share causal mechanisms such
47 as autoimmune diseases [14, 15] and psychiatric disorders [16, 17]. Thus, methods that jointly
48 model the intrinsic comorbidity implicated in the GWAS summary statistics of the related
49 traits may confer higher statistical power of causal variants detection. Recently, several
50 methods were developed to utilize the wealth of genome-wide annotations primarily pro-
51 vided by ENCODE consortium to predict causal variants and novel risk variants that are
52 weakly associated in complex traits. Pickrell (2014) developed a statistical approach called
53 fgwas that models association statistics of a given trait and used regularized logistic func-
54 tion to simultaneously learn the relevant annotations. To account for LD, fgwas assumes at
55 most one causal variants per locus via a softmax function. Kichaev et al. (2014) recently
56 developed a multivariate Gaussian framework called PAINTOR, which allows for more than
57 one causal SNP but at most three to be located within a single locus by considering all of
58 the combinatorial settings [18]. Chung et al. (2014) used a maximum likelihood framework
59 called GPA to infer driver variants shared among multiple traits by modeling the corre-
60 sponding GWAS p-values as Beta distributions with an option of using one or more sets
61 of annotations to improve the power detecting causal variants [19]. Although useful, these
62 methods are often designed to simultaneously operate on a small number of independent
63 annotations due to some computational constraints. Moreover, most methods only operate
64 on one trait at a time whereas exploiting the correlation between traits at the epigenomic
65 annotation level may prove useful for shared causal mechanisms that go beyond the level of
66 individual variants.

67 In this article, we describe a novel Bayesian framework called RiVIERA-beta (Risk
68 Variant Interference using Epigenomic Reference Annotations to model Beta likelihood of
69 GWAS summary statistics p-values). The main novelty of RiVIERA-beta is the ability to
70 perform efficient Bayesian inference of the intrinsic causal signals across multiple traits while

71 simultaneously inferring and exploiting enrichment signals and their correlation between
72 traits over hundreds of tissue-specific epigenomic annotations. We achieve this efficiently
73 via stochastic sampling of loci and powerful Hamiltonian Monte Carlo sampling of model
74 parameters [20]. We first use simulation to demonstrate the utility of RiVIERA-beta in
75 prioritizing driver variants and detecting functional epigenomic annotations. We then apply
76 RiVIERA-beta to some of the most well-powered GWAS datasets, consisting of 9 immunolog-
77 ical disorders from ImmunoBase [15] and Schizophrenia 2014 data from Psychiatric Genomic
78 Consortium [21]. To infer tissue-specific epigenomic enrichments, we utilize the largest com-
79 pendium of epigenomic annotations to date from ENCODE/Roadmap Consortia, consisting
80 of 848 annotations including 8 major epigenomic marks across 127 distinct cell types [7].
81 This allows us to revisit the GWAS of these 10 common complex disorders by inferring their
82 underlying regulatory variants implicated at the tissue-specific epigenomic contexts.

83 2 MATERIALS AND METHODS

84 GWAS summary statistics

85 The GWAS summary statistics for the nine immune diseases were obtained from ImmunoBase
86 (March 17, 2015) [15]. The nine diseases are: Autoimmune Thyroid Disease (ATD), Celiac
87 Disease (CEL), Juvenile Idiopathic Arthritis (JIA), Multiple Sclerosis (MS), Narcolepsy
88 (NAR), Primary Biliary Cirrhosis (PBC), Psoriasis (PSO), Rheumatoid Arthritis (RA),
89 Type 1 Diabetes (T1D). We imputed the p-values of un-genotyped SNPs using FAPI and
90 1000 Genome European data (Phase 1 version 3) [22]. We then obtained the p-values of
91 SNPs that fall within the pre-defined risk loci available from ImmunoBase for each of the
92 9 immune traits. For all analyses, we filtered out risk loci or variants in the MHC regions
93 or sex chromosomes X and Y. The Schizophrenia 2014 (SCZ2) summary data containing
94 642846 observed and imputed SNPs were obtained from Psychiatric Genomic Consortium
95 (PGC) [21]. Among these, 54132 SNPs fall within the 105 SCZ-associated loci of the au-
96 toosomes (chr 1-22) defined by PGC (we filtered out the 3 loci on chromosome X). **Table 1**
97 summarizes the total number of SNPs and risk loci for each individual GWAS that were
98 subject to the proposed fine-mapping analyses.

99 Roadmap epigenome data

100 Roadmap epigenome data were obtained from Roadmap epigenomic web portal (March,
101 2015). Peaks were defined if their p-values were below 0.01 (i.e., following the definition of
102 “Narrow Peaks” [7]). In total, there are 848 epigenome tracks, including 8 epigenomic marks
103 namely H3K4me1, H3K4me3, H3K36me3, H3K27me3, H3K9me3, H3K27ac, H3K9ac and
104 DNase I in 127 cell or tissue types, which were grouped into 19 categories [7]. To associate
105 each SNP with the annotations, we overlapped their genomic coordinates with each bigWig
106 epigenome track making use of the R packages *rtracklayer* and *GenomicAlignments*. SNPs
107 that fall within a peak of an annotations will have value 1 otherwise 0 for that annotation.
108 The resulting matrix is a $V_d \times K$ input matrix containing the epigenomic values across
109 $K = 848$ marks for each of the V_d SNPs in disease d .

110 Running existing fine-mapping software on simulated data

111 fgwas

112 The software `fgwas` [23] (version 0.3.4) were downloaded from GitHub. We prepared the
113 input for `fgwas` (1) the Z scores calculated as the t-statistics of the linear coefficients of
114 the genotype of each variant fitted separately by least square regression on the simulated
115 continuous phenotypes (**Materials and methods**) and (2) 100 discretized epigenomic an-
116 notations at $p < 0.01$. To enable fine-mapping, we issued `-fine` flag and specify the region
117 numbers for each SNP in the input file as required by the software. As part of the outputs
118 from `fgwas`, we obtained ‘PPA’ and ‘estimate’ for the causal variants and influences of each
119 epigenomic annotations, respectively.

120 GPA

121 GPA (0.9-3) [19] was downloaded from GitHub and run with default settings. Same as above,
122 we set the annotations to one at p-value < 0.01 and 0 otherwise. To test for trait-relevant
123 annotations, we followed the package vignette. Briefly, we fit two GPA models with and
124 without the annotation and compared the two models by `aTest` function from GPA, which
125 performs likelihood-ratio (LR) test via χ^2 approximation, and obtained the enrichment scores
126 as the $-\log_{10}$ p-value.

127 PAINTOR

128 PAINTOR (version 2.1) was downloaded from GitHub [18]. As suggested in the docu-
129 mentation, we prepared a list of input files for every locus including summary statistics as
130 t-statistics, LD matrices, and binary epigenomic annotations. We ran the software with
131 default setting with assumption of at most two causal variants per locus. We then extracted
132 the ‘Posterior.Prob’ and ‘Enrichment.Values’ as the model predictions for causal variants
133 and causal annotations, respectively.

134 Details of RiVIERA-beta Bayesian model

135 Inference of empirical prior π_{vd}

136 We first define the empirical prior function of a variant v being associated with disease d as
137 a logistic function:

$$\pi_{vd} = [1 + \exp(-[\sum_k w_{kd}e_{vk} + w_{0d}])]^{-1} \quad (1)$$

138 where $w_{kd} \in \mathbf{w}_d$ denotes the linear coefficient or the influence of the k^{th} epigenomic mark
139 affecting disease d and w_{0d} is the linear bias.

We assume that epigenomic causal effect w_{kd} follows a multivariate Gaussian distribution with zero mean and unknown covariance:

$$w_{kd} \sim \mathcal{N}(0, \Lambda_w^{-1}) \quad (2)$$

$$\Lambda_w \sim \mathcal{W}(\Lambda_0, \nu_0) \quad (3)$$

140 where Λ_w is a $D \times D$ inverse covariance matrix $\Lambda_w = \Sigma_w^{-1}$ to model the pairwise epigenomic
 141 correlation among D diseases. It follows a Wishart distribution with identity matrix as prior
 142 (i.e., by default, we assume *a priori* no correlation between the target traits) and $\nu_0 = 0$ (i.e.,
 143 by default, we did observe any samples *a priori* that are indicative of the correlation between
 144 any two diseases being modeled). The hyperparameters can be easily modified to incorporate
 145 prior belief on the correlation between any two diseases of interests.

Additionally, the bias w_{0d} follows a Gaussian distribution with unknown variance and mean determined based on our prior belief of the causal fraction π_0 :

$$w_{0d} \sim \mathcal{N}(\text{logit}(\pi_0), \lambda_{0d}^{-1}) \quad (4)$$

$$\lambda_{w_{0d}} \sim \Gamma(\alpha_0, \beta_0) \quad (5)$$

146 where $\text{logit}(\pi_0) = \log \frac{\pi_0}{1-\pi_0}$. By default, we set π_0 to 0.01, implying that 1% of the SNPs in
 147 the risk loci are expected to be causal when no functional enrichment. We set $\alpha = 0.01$ and
 148 $\beta = 0.0001$ to enable a broad hyperprior for w_{0d} .

149 Notably, w_{kd} can be interpreted as enrichment coefficient for annotation k in disease d ,
 150 where a positive w_{kd} will increase the causal prior π_{vd} when $e_{vk} = 1$. During the training,
 151 however, w_{kd} may become negative, which makes the interpretation difficult. Thus, we
 152 constrain w_{kd} to be non-negative values, which involves imposing infinitely high potential
 153 energy for negative w_{kd} . More details are described in **Supplementary Text 1**.

154 Inference of variant causality c_{vd} given prior π_{vd} and model parameters μ_d, ϕ_d

Because the target association variable a_{vd} for variant v in disease d represents p-values, which are continuous and restricted to the interval (0, 1), we assume that it follows a Beta distribution with unknown mean μ_d and unknown precision ϕ_d :

$$a_{vd} \sim \mathcal{B}(\mu_d, \phi_d) \quad (6)$$

Note that we re-parameterize Beta density function from the traditional “rate” p and “shape” q parameters, and instead use mean $\mu = p/(p + q)$ and precision $\phi = p + q$, as per [24, 25]. Specifically, the density function of association variable a_{vd} is defined as follows:

$$f(a_{vd}; \mu_d, \phi_d) = \frac{\Gamma(\phi_d)}{\Gamma(\mu_d \phi_d) \Gamma((1 - \mu_d) \phi_d)} a_{vd}^{(\mu_d \phi_d - 1)} (1 - a_{vd})^{(1 - \mu_d) \phi_d - 1} \quad (7)$$

Further, we let the mean μ_d and precision ϕ_d follow Beta and uniform prior, respectively:

$$\mu_d \sim \mathcal{B}(\mu_0, \phi_0) \quad (8)$$

$$\phi_d \sim \mathcal{U}(0, \phi_{\max}) \quad (9)$$

155 where the hyperparameters (μ_0, ϕ_0) reflect *a priori* belief on the p-value signal of a causal
 156 variant. By default, we set $\mu_0 = 0.1$ and $\phi_0 = 2$. If $\phi_{\max} = \infty$, ϕ follows an improper prior.
 157 Because it is unlikely to have a very large ϕ , by default, we set ϕ_{\max} to 1000. Notably, as
 158 long as ϕ_{\max} is large, the inference results remain the same with different ϕ_{\max} values.

With the prior $p(c_{vd}|\mathbf{w}_d, \mathbf{e}_v) \equiv \pi_{vd}$ and likelihood $p(a_{vd}|\mu_d, \phi_d) \equiv f(a_{vd}; \mu_d, \phi_d)$ established, the posterior probability of association (PPA) [26] of variant v being causal for disease d then follows:

$$p(c_{vd}|a_{vd}, \pi_{vd}) = \frac{p(a_{vd}|c_{vd})p(c_{vd}|\mathbf{w}_d, \mathbf{e}_v)}{\sum_{v' \in \mathcal{V}_b} p(a_{v'd}|c_{v'd})p(c_{v'd}|\mathbf{w}_d, \mathbf{e}_{v'})} \quad (10)$$

159 where \mathcal{V}_b represent all variants within locus b . The 95% credible set \mathcal{C}_{bd} for each locus b is
 160 the minimal number of SNPs $v' \in \mathcal{C}_{bd}$ in the locus such that $\sum_{v' \in \mathcal{C}_{bd}} p(c_{v'd}|a_{v'd}, \pi_{v'd}) \geq 0.95$.

161 Joint posterior distribution

The complete likelihood density function treating c_{vd} as missing values is defined as:

$$\begin{aligned} \mathcal{L} &= \prod_{v,d} f(a_{vd}, \pi_{vd}, c_{vd} | \mathbf{e}_v, \mu_d, \phi_d) \\ &= \prod_{v,d} [\pi_{vd} p(a_{vd} | \mu_d, \phi_d)]^{c_{vd}} (1 - \pi_{vd})^{(1-c_{vd})} \end{aligned} \quad (11)$$

The logarithmic joint posterior density function is then:

$$\begin{aligned} \log p(\Theta|\mathcal{D}) &= \log f(\mu, \phi, \mathbf{W}, \Lambda_w, \lambda_{0d} | \mathbf{E}, \mathbf{c}_d, \pi_d, \mathbf{a}_d) \\ &\propto \log f(\Lambda_w | \Lambda_0, \nu_0) + \sum_d \log f(\lambda_{0d} | \alpha_0, \beta_0) \\ &\quad + \log f(\mathbf{W} | \Lambda_w) + \sum_d \log f(w_{0d} | \mu_{w_0}, \lambda_{0d}) \\ &\quad + \log f(\mu_d | \mu_0, \phi_0) + \log f(\phi_d) \\ &\quad + \sum_{v,d} \log f(a_{vd}, \pi_{vd}, c_{vd} | \mathbf{e}_v, \mu_d, \phi_d) \end{aligned} \quad (12)$$

In principle, causality is inferred by integrating out all nuisance parameters:

$$p(c_{vd} | \mathbf{a}_d, \mathbf{e}_{vd}) = \int f(c_{vd} | \mathbf{a}_d, \mathbf{e}_v, \mu_d, \phi_d, \mathbf{w}_d, \Lambda_w) f(\mu_d, \phi_d, \mathbf{w}_d, \Lambda_w | \mathbf{a}_d, \mathbf{e}_v) d\mu_d, \phi_d, \mathbf{w}_d, \Lambda_w \quad (13)$$

162 which is not tractable. We employ Markov Chain Monte Carlo (MCMC) to sample from the
 163 joint posterior.

164 Markov Chain Monte Carlo

165 We use Gibbs sampling [27] to sample the precision matrix Λ_w of epigenomic effects from
 166 the posterior distribution. Specifically, Gibbs sampling requires a closed form posterior
 167 distribution. Due to the conjugacy of the Wishart prior of epigenomic precision Λ_w to the
 168 multivariate normal distribution of epigenomic effect \mathbf{W} , the posterior of the epigenomic
 169 precision matrix Λ_w also follows Wishart distribution [28]:

$$\Lambda_w | \mathbf{W} \sim \mathcal{W}((\Lambda_0^{-1} + \mathbf{S})^{-1}, \nu_0 + K) \quad (14)$$

170 where \mathbf{S} is the sample variance of \mathbf{W} , i.e., $S = \mathbf{W}^T \mathbf{W}$.

171 Similarly, we sample λ_{0d} from Gamma posterior distribution:

$$\lambda_{0d}|w_{0d} \sim \Gamma(\alpha_0 + 0.5, (\beta_0 + \frac{(w_{0d} - \mu_{w_0})^2}{2})^{-1}) \quad (15)$$

172 To sample epigenomic effects \mathbf{w}_d , prior bias w_{0d} , causal mean μ_d , causal precision ϕ_d for
 173 disease $d = 1, \dots, D$, we employ a more powerful gradient-based sampling scheme namely
 174 Hamiltonian Monte Carlo (also known as hybrid Monte Carlo) (HMC) [20, 29], exploiting
 175 the fact that the joint posterior of our model is differentiable with respect to the model
 176 parameters $\mu_d, \phi_d, w_{kd}, w_{0d}$ (**Supplementary Text S1**). Finally, after discarding $t\%$ models
 177 accepted before the burn-in period (default: $t=20\%$), we obtain the Bayesian estimates of
 178 PPA by averaging the corresponding values computed over the T' individual models accepted
 179 throughout the T MCMC runs.

180 Bayesian fold-enrichment tests for epigenomic annotations

Due to co-linearity among the epigenomic annotations, directly using w_{kd} to assess the epige-
 nomic enrichment for annotation k may be misleading. We propose an heuristic approach
 to assess the log fold-enrichment of the full prior model over the alternative prior with the
 effect of annotation k for disease d removed (i.e., $\mathbf{w}_{d \setminus k}, w_{kd} = 0$):

$$f_{kd} = \log \int p(\mathbf{w}_d) \frac{p(\mathbf{c}_d | \mathbf{w}_d, \mathbf{e}_v)}{p(\mathbf{c}_d | \mathbf{w}_{d \setminus k}, w_{kd} = 0, \mathbf{e}_v)} d\mathbf{w}_d \quad (16)$$

$$\approx \frac{1}{T'} \sum_{t=1}^{T'} \log \frac{1}{|\mathcal{C}_d|} \sum_{v \in \mathcal{C}_d} \frac{p(c_{vd} | \mathbf{w}_d^{(t)})}{p(c_{vd} | \mathbf{w}_{d \setminus k}^{(t)}, w_{kd}^{(t)} = 0)} \quad (17)$$

181 where $p(c_{vd} | \mathbf{w}_d^{(t)}, \mathbf{e}_v)$ is the logistic prior based on Eq 1, \mathcal{C}_d is the union of all the 95% credible
 182 sets across loci for disease d : $\mathcal{C}_d = \bigcup_b \mathcal{C}_{bd}$. Notably, under the optional constraint that
 183 $w_{kd} \geq 0$, f_{kd} is always positive, which directly translates to fold-enrichment of annotation k
 184 conditioned on all the other annotations $k' \neq k$. The 95% Bayesian credible interval for f_{kd}
 185 are obtained from the T' MCMC runs. The significance of each annotation k is determined
 186 based on the ranking of its lower bound f_{kd} (i.e., the 2.5% quantile of f_{kd}).

187 Alternatively, we can estimate the fold-enrichment for each annotation simply based on
 188 the ratio of estimated fraction of causal variants in an annotation e_{vk} over the fraction of
 189 all of the variants in that annotation $\frac{\sum_v c_v e_{vk} / \sum_v c_v}{\sum_v e_{vk} / V}$, where c_v is the PPA for SNP v . This
 190 is more efficient and accurate when the underlying causal variants were randomly sampled
 191 from the annotations as done in the simulation.

192 Stochastic gradient updates per locus

193 Directly updating model parameters based on the gradients of all GWAS loci at each MCMC
 194 iteration is inefficient and results in poor HMC acceptance rate. Instead, at each MCMC
 195 update, we randomly sample one locus and update the model parameters (which are shared
 196 across loci) based on that locus. We find this approach quite efficient in capturing meaningful

197 causal properties such as causal signals and relevant epigenomes that are shared across all
198 risk loci. Together, we outline the overall algorithm of the proposed Bayesian model in
199 Algorithm S2 (**Supplementary Text S1**).

200 **GWAS simulation**

201 To assess the power of the proposed fine-mapping model in identifying causal variants and
202 compare it with existing methods, we implemented a simulation pipeline adapted from [18].
203 Briefly, the simulation can be divided into three stages (1) simulate genotypes based on
204 the haplotypes from 1000 Genome European data (phase 1 version 3) using HapGen2 [30]
205 (**Supplementary Fig. S1**); (2) simulate epigenomic enrichments and subsequently sam-
206 ple causal variants accordingly using 100 Roadmap annotations selected from each of the
207 19 categories of primary tissue/cell types (**Supplementary Fig. S1**); (3) simulate liability
208 phenotype plus the random noise to obtain the desired heritability (fixed at 0.25) and sub-
209 sequently the GWAS summary statistics in terms p-values and z-scores via ordinary least
210 square regression. Details are described in **Supplementary Text**.

211 **Gene ontology enrichment analysis**

212 We obtained the latest gene annotations from Ensembl database (version 80) programmat-
213 ically via biomaRt package [31], which resulted in 10,801 gene ontology (GO) terms in
214 biological processes (BP). To assign SNPs to genes, we performed lift-over to map the SNPs
215 from hg19 to hg38 using rtracklayer [32] and assigned each SNP to a gene if it is located
216 within 35 kb up and 10 kb downstream of that gene. The resulting Ensembl gene identifiers
217 were matched with those genes in each GO-BP category. We then performed hypergeomet-
218 ric tests on each GO-BP term for all of potential *in-cis* target genes of the SNPs in each
219 trait and adjusted for multiple testings using Benjamini-Hocherg family-wise Type I error
220 correction method [33]. For the 9 immune traits, the enrichment signals are strong so we set
221 the cutoff at $FDR < 0.005$; for Schizophrenia, we set $FDR < 0.2$.

222 **RiVIERA-beta software**

223 RiVIERA-beta is available as an open-source R package with documented functions and
224 walk-through examples described in the vignette. Most functions were implemented in C++
225 by integrating *Rcpp* and *RcppArmadillo* libraries [34]. These libraries enabled us to apply
226 RiVIERA-beta to large matrices very efficiently with compiled code and having much lesser
227 memory overhead than a naïve R implementation. RiVIERA-beta is available at Github
228 (<https://github.mit.edu/liyue/rivieraBeta>).

229 **3 RESULTS**

230 **RiVIERA model overview**

231 The fundamental hypothesis of our model is that non-coding disease associations are driven
232 by disruption of regulatory elements of common activity patterns (e.g., motifs of sequence-

specific regulators), thus leading to gene expression changes and ultimately phenotypic changes at the cellular or organism level between case and control individuals. Our RiVIERA-beta Bayesian model aims to infer the probability that a given variant v is a driver for disease d by modeling the corresponding GWAS association statistic for that variant using a vector of genome-wide epigenomic annotations (\mathbf{e}_v). Given a set of B risk loci, the inputs to RiVIERA-beta are GWAS summary statistics in terms of p-values and a set of discrete or continuous epigenomic annotations (**Fig. 1a**). In this study, we used binary signals to ease interpretation of the functional enrichments. We train RiVIERA-beta by repeatedly sampling one locus at each iteration to efficiently learn the intrinsic (i.e., locus-independent) causal signals. **Fig. 1b** depicts RiVIERA-beta as probabilistic graphical model [35]. The observed variable of our model is the GWAS association values (in terms of p-values) a_{vd} for each variant v in each disease d . We assume that a_{vd} follows a Beta distribution with unknown mean and dispersion parameters. The effect of each annotation on each trait is learned as global annotation-by-disease weight matrix \mathbf{w} , which follows a D -dimensional multivariate normal distribution with zero mean and $D \times D$ disease-disease covariance Λ_w . The prior probability π_{vd} that a variant v is causal in disease d is essentially a linear combination of the weighted genomic annotations \mathbf{e}_v , which reflects the disease-associated active histone marks and DNA accessibility in the 127 cell types (**Materials and methods**). The outputs of the model (**Fig. 1c**) are (a) posterior probability of association (PPA) c_{vd} that variant v is causal in disease d ; (b) the Bayesian fold-enrichment estimates f_{kd} based on the ratio between the full prior model with all annotations over the null prior model with all annotations except for annotation k .

Method comparison using GWAS simulation

The goal of the simulation is to evaluate the model's power to predict (1) causal variants in each locus; (2) the relevant annotations that determine which variants are causal. To this end, we simulated GWAS summary statistics based on 1000 Genome European data (Phase 1 release 3) (**Supplementary Fig. S1**) and 100 representative epigenomic annotations (**Supplementary Fig. S1**) (**Materials and methods**). We performed a series of power analyses over 500 simulation runs.

First, we examined how well the posterior probabilities were calibrated by taking the credible SNPs that contribute to 95% posterior mass inferred by each method (**Supplementary Fig. S2**). As expected, when our model assumption of single-causal variant per locus holds, we observe that our model is well calibrated (**Fig. S2**), where the 95% credible SNPs indeed correspond to approximately 95% of the causal variants. When there are more than one causal variants per locus, the 95% credible SNPs include on average 50% the true causal SNPs (**Supplementary Fig. S2**).

Because the number of variants within the credible set differs depending on the concentration of the posterior probabilities inferred by each method, we sought to control that bias by evaluating the proportion of identified causal variants as a function of the absolute number of selected variants. When the assumption of one-causal-variant-per-locus holds, we observed comparable or better performance of RiVIERA-beta compared to existing methods (**Fig. 2**). As expected, when the assumption is violated, our current model is second to PAINTOR, which is able to infer multiple causal variants per locus (**Supplementary**

276 **Fig. S3**). We also examined the correlation between the functional enrichments estimated
277 by each method and the underlying epigenomic enrichments that were used to simulate the
278 causal variants. The performance of the four methods are comparable with the proposed
279 model achieving a slightly better correlation (**Supplementary Fig. S4**).

280 **Applications to immune and psychiatric disorders**

281 To demonstrate RiVIERA-beta in a real-world application, we used it to investigate 10
282 complex diseases including 9 immune diseases with summary statistics obtained from Im-
283 munoBase [15] and Schizophrenia from Psychiatric Genomic Consortium (PGC) [21] (**Table 1**).
284 We used 848 epigenomic annotations from ENCODE/Roadmap consortium (**Materials and**
285 **methods**) to build a functional prior for each trait to aid fine-mapping and conduct cell-type
286 specific epigenomic enrichment analyses [7]. We first applied RiVIERA-beta to the 10 traits
287 separately to examine individual causal signals and then demonstrated RiVIERA-beta's ca-
288 pability to operate on the 9 immune traits and the improved detection power compared to
289 the single-trait model.

290 **RiVIERA-beta detected meaningful tissue-specific enhancers in test** 291 **GWAS traits**

292 We first sought to confirm the validity of the model through its ability to identify mean-
293 ingful cell-types or tissues for each trait. To this end, we selected the top 5% (i.e., the top
294 43) of the 848 annotations for each disease based on the corresponding Bayesian estimates
295 of the lower bounds of the 95% credible interval (**Supplementary Table S1; Materials**
296 **and methods**). We then performed hypergeometric tests on enrichments of each of the 19
297 categories grouped by Roadmap consortium based on the cell types and tissues [7]. Indeed,
298 we observed a significant enrichment for Blood & T-cell for all 9 immune disorders but not
299 for Schizophrenia, which exhibits exclusive epigenomic enrichments in the Brain category
300 (Hypergeometric adjusted p-values < 0.05) (**Fig. 3a**). Additionally, we also observed mod-
301 est enrichments for B-cell and Thymus tissue in the 9 immune traits. We then examined the
302 enrichment status for the 8 epigenomic marks. Indeed, enhancer marks namely H3me4me1
303 and/or H3K27ac are most significantly enriched among all 8 marks (q-values < 0.05). In
304 addition, H3K4me3 associated with promoter is also enriched in most immune traits. In-
305 terestingly, we also observed a modest enrichment of H3K9me3 in Schizophrenia but not in
306 the immune traits. We further ascertained the enrichment results by re-running RiVIERA-
307 beta on the permuted data matrix and observed diminishment of the meaningful enrichment
308 observed above (**Supplementary Fig. S5**).

309 **SNPs in the credible set exhibit promising regulatory potentials**

310 The variants in the credible set are more enriched for functional elements. Inspired by
311 the promising tissue-specific enhancer enrichment results obtained above, we refined our
312 RiVIERA-beta model by re-training it on the top 5% (or 43) annotations on each trait
313 using the same GWAS data. For each locus in each trait, we then constructed 95% credible
314 set (**Supplementary Table S2; Materials and methods**). On average, we were able

315 to construct a rather small credible set ranging from 4 to 25 SNPs per locus for the 10
316 traits (**Table 1**). As a comparison, we extracted the same number of SNPs with the most
317 significant GWAS p-values from each locus. For ease of reference, we named our SNPs in the
318 credible set as “credible SNP” and the GWAS counterpart as “GWAS SNP”. Compared to
319 GWAS SNPs, the credible SNPs exhibit substantially higher averaged placental conservation
320 scores (phastCons46way obtained from UCSC database) across most traits (**Fig. 4 CONS**).

321 Moreover, the credible SNPs were significantly enriched for expression quantitative trait
322 loci (eQTL) that are in the regulatory regions. Specifically, we obtained in total 806,847
323 GTE_x whole-blood eQTL-SNPs (version 6) [36] and retained 122,549 and 23,973 eQTL-
324 SNPs that overlap with transcription factor binding sites derived from 1,772 TF recognition
325 motifs [37] and digital genomic footprints (DGF) at 6-bp resolution derived from DNaseI data
326 in CD cells using method described in [38], respectively as well as 6,743 eQTL-SNPs that
327 overlapped with both the TFBS and DGF regions. We then performed hypergeometric tests
328 to assess the significance of overlap between the credible/GWAS SNPs and the regulatory-
329 eQTL SNPs. Indeed, our credible SNPs exhibit much higher enrichments for those eQTL-
330 SNPs, suggesting their regulatory potentials elucidated based on the enhancer activities by
331 our proposed RiVIERA-beta model (**Fig. 4; Supplementary Table S3**).

332 **Gene-centric analysis revealed enrichment for meaningful biological** 333 **processes**

334 Genes adjacent to the SNPs in credible sets are significantly enriched for disease-specific
335 biological processes. In particular, we observed significant enrichments of many immune-
336 related processes for the *in-cis* genes for which the SNPs in the credible set are within
337 35 kb upstream or 10 kb downstream (**Fig. 5; Supplementary Table S4; Materials**
338 **and methods**). For instance, regulation of T cell homeostatic proliferation, regulation of
339 interferon-gamma-mediated signaling pathway, and regulation of type I interferon-mediated
340 signal pathways are among the most significantly enriched GO terms in 5 or 6 out of the 9
341 immune traits. In contrast, the enrichments for Schizophrenia are dominated by GO terms
342 involving synaptic processes and neuronal differentiation/development. The enrichment re-
343 sults are mostly consistent between the credible genes and the genes derived from the same
344 number of SNPs chosen based on the GWAS p-values (GWAS-genes).

345 Intriguingly, we observed a highly significant enrichment for keratinization (GO:0031424)
346 and epidermis (e.g., skin) development (GO:0008544) exclusively for Psoriasis. In particular,
347 17 genes among the 241 credible genes belong to keratinization and epidermis development,
348 which contain in total 49 and 121 genes, respectively ($q < 9 \times 10^{-18}$, $q < 2 \times 10^{-10}$).
349 Indeed, Psoriasis is mainly characterized as a chronic skin disease with epidermal hyper-
350 proliferation [39,40]. In contrast, there are only 6 out of 157 GWAS-genes are defined in
351 each of two GO categories ($q < 0.001$).

352 To further ascertain the RiVIERA-beta fine-mapping results, we created a visualization
353 scheme for each of the 469 risk loci across 10 traits examined (**Supplementary Fig. S6**).
354 **Fig. 6** displays two example loci for Type 1 diabetes (chr17: 37383069-38239012) and
355 Schizophrenia (chr7: 104598324-105062839). The upper panel displays the RiVIERA-beta model
356 prior, the genetic signals from GWAS -log p-values, and RiVIERA-beta PPA. Red colored

357 and diamond shape points are GTEx whole-blood eQTL SNP and top SNPs included into
358 60% credible set (we used 60% to not clutter the plot with the remaining SNPs in the 95%
359 credible set that exhibit low PPA). Intuitively, SNPs with high PPA exhibit both prominent
360 genetic and epigenetic signals. Thus, to infer causal variants, RiVIERA-beta efficiently took
361 into account not only the GWAS signals derived from the genetic data but also the prior
362 signals mainly driven by the weighted epigenomic profiles. The middle panel illustrates the
363 cumulative density for each epigenomic profiles weighted by the tissue-specific enrichment
364 estimates.

365 Consistent with the overall enrichment results (**Fig. 3**), we observed prominent enrich-
366 ments for the enhancer regions predominantly in blood T cells for all of the 9 immune traits
367 and brain tissue for Schizophrenia. The bottom tracks display the external functional infor-
368 mation (i.e., not in the training data) including conservation score, genes, transcription factor
369 binding sites based on motif matches that may further aid variant selection for downstream
370 experimental validation (please refer to **Supplementary Table S2** for detailed information).
371 We also visualized the signals within the of Psoriasis-associated risk region ch1:152536784-
372 152785170, which harbors genes involved in keratinization and epidermis development as
373 mentioned above. Interestingly, as an exception of most other immune-susceptible loci, the
374 underlying epigenomic profiles exhibit prominent signals not only in blood T cell but also
375 in epithelia enhancer regions (**Supplementary Fig. S6**). However, the associated SNPs
376 exhibit rather weak genetic signal perhaps due to lower allele frequencies.

377 **Multi-trait causal inference improved functional enrichments in** 378 **most immune traits**

379 Exploiting epigenomic correlation between highly related immune diseases improved func-
380 tional enrichments in several traits. We performed multi-trait causal inference over all of the
381 9 autoimmune traits by jointly applying our RiVIERA-beta to 364 risk loci concatenated
382 together from the 9 immune traits using 174 epigenomic annotations which was a union of
383 unique annotations from the top 43 annotations for each individual trait. We focused only on
384 the 9 immune GWAS (i.e., leaving out Schizophrenia) because they commonly utilized the
385 same genotyping platform namely ImmunoChip. The multi-trait GWAS summary statistics
386 triggered RiVIERA-beta to infer the disease covariance matrix and sample disease-specific
387 epigenomic weights from the joint posterior with modified zero-mean multivariate normal
388 prior that takes into account the sampled disease covariance (**Materials and methods**).
389 As a results, RiVIERA-beta sampled correlated epigenomic weights between traits more
390 frequently compared to the single-trait model.

391 We constructed the 95% credible sets for each trait using the disease-specific PPA derived
392 from the joint model and assessed the functional enrichments as above (**Supplementary**
393 **Table S6**). Notably, the cross-trait model exploited 174 annotations as apposed to 43
394 annotations used by the single-trait model. To examine whether the improved enrichments
395 were due to the increased number of annotations, we re-ran a single-trait model for each
396 of the 9 traits separately each using the 174 annotations. Compared to the 95% credible
397 set constructed based on the single-trait causal inference using the top 43 annotations, we
398 observed smaller 95% credible sets for 8 out of the 9 immune traits (**Table 1**), suggesting

399 that the multi-trait joint inference further improved the model confidence in some of the
400 highly related traits.

401 More importantly, we observed a much more improved enrichment for the GTEx whole-
402 blood eQTL SNPs located within open chromatin regions or digital genomic footprints in
403 most of the immune traits (**Fig. 7; Supplementary Table S5**). Thus, the joint inference
404 further improved the regulatory potential through following the Hamiltonian trajectory that
405 is more consistent with the epigenomic correlation pattern between the related immune traits.
406 We also repeated the GO enrichment analysis on the 95% credible set and found that the
407 enriched GO terms were mostly immune-specific biological processes and consistent with the
408 above single-trait analyses (**Supplementary Fig. S7; Supplementary Table S7**).

409 4 DISCUSSION

410 Understanding the genetic basis of complex traits hinges upon powerful integrative meth-
411 ods to map genotypes to phenotypes [41]. Fine-mapping causal variants has been a highly
412 active and fruitful research in the past several years [9, 18, 42–44]. However, most existing
413 methods typically operate solely on genetic data by estimating each SNP of being causal
414 conditioned on the lead index SNPs in the same LD block, which are typically approximated
415 by the 1000 Genome data [9, 15, 45, 46]. With the recent availability of large-scale functional
416 genomic data provided by ENCODE/Roadmap consortia, there is an urgent need to incor-
417 porate these valuable information in a principled way as a form of Bayesian prior. In this
418 article, we describe a novel Bayesian fine-mapping method RiVIERA-beta to re-prioritize
419 GWAS summary statistics based on their epigenomic contexts. The main contribution of
420 RiVIERA-beta is the ability to systematically construct from a targeted set of susceptible
421 loci a Bayesian credible set of SNPs, which exhibit plausible tissue-specific regulatory prop-
422 erties implicated in the large epigenomic data compendium either in a single trait or across
423 multiple traits.

424 One benefit of using the raw epigenomic annotations rather than using the inferred signals
425 such as ChromHMM [7] or GenoSkyline [47] states derived from the annotations is that
426 it eases the interpretation of the actual relevant epigenomic marks in the relevant tissue
427 types and facilitates downstream experimental efforts to assay the specific marks in the
428 disease-specific cell types. However, the correlation of the epigenomic marks will make
429 difficult estimating the underlying functional enrichments. Moreover, we choose to model
430 the summary statistics rather than genotypes because it is not always possible to obtain
431 individual-level phenotype-genotype data particularly for large-scale meta-analysis. Thus,
432 effective methods based on summary statistics may benefit wider research communities than
433 methods that solely operate on individual-level genotype data [18, 19, 23]. Moreover, our
434 model *requires only p-values* because it uses Beta distribution to model the likelihood. In
435 contrast, fgwas requires both the z-scores and the standard error from the linear regression
436 used in the GWAS to estimate the Wakefield approximate Bayes factors. While some recent
437 GWAS summary statistics include those information, there are many do not have z-scores
438 and/or standard error of the linear model but only p-values (e.g., the ImmunoChip data
439 we used in our studies for the 9 immune traits). When the standard error is not available
440 in a given GWAS summary statistics, fgwas needs to estimate it from the minor allele

441 frequency of a reference panel such as 1000 Genome, which may not be accurate depending
442 the study cohorts. Additionally, modeling p-values via Beta density function only has more
443 relaxing model assumption than modeling z-scores via normal density although both methods
444 are highly effective in practice.

445 Overall, SNPs included into the credible set exhibit both significant GWAS signal and
446 high prior. In some cases, however, SNPs that were added to the credible set in each locus
447 do not exhibit significant GWAS p-values (**Supplementary Table S2,S6**). This generally
448 occurs when the genetic signals in those loci are weak relative to the SNPs in other loci
449 for the same trait, and the model functional prior eventually dominates the SNP prioritiza-
450 tion. Thus, we recommend considering these variants cautiously when designing downstream
451 experiments.

452 One important assumption of our model is that there is one causal variant per locus,
453 which is reflected by the normalization of variants within each locus so that they sum to
454 1 in order to obtain PPA and construct 95% credible sets [23]. When this assumption
455 holds, the posterior probabilities are well calibrated (**Supplementary Fig. S2**). However,
456 as demonstrated in our simulation, when this assumption is violated, the PPA is not well
457 calibrated (**Supplementary Fig. S2,S3**). Other existing method such as PAINTOR [18]
458 and CAVIAR [48] employ multivariate normal distribution to model all of the variants within
459 a locus using LD reference panel estimated from 1000 Genome data as the covariance matrix,
460 which allows inferring more than one causal variants per locus. While CAVIER used only
461 summary statistics, PAINTOR is able to employ functional annotations to aid fine-mapping.
462 Both methods require computing the likelihood density across a combinatorial set of causal
463 configurations and therefore still needs to assume at most an arbitrarily small number of
464 causal variants, typically below 10 causal SNPs per locus.

465 As future works, we will explore potential ways to enable efficient inference of more
466 than one causal variants per locus. Furthermore, we will also explore the potential gain
467 of incorporating trans-ethnic data, which was effectively demonstrated by the trans-ethnic
468 version of the PAINTOR model [49]. Moreover, in addition to modeling the epigenomic
469 correlation between traits, variant prioritization may further benefit by jointly inferring the
470 comorbidity at the individual SNP level [19], gene level [50], and/or pathway level [17].
471 Together, we believe that RiVIERA-beta will serve as a valuable tool complementary to the
472 existing methods in identifying novel risk variants through tissue-specific epigenome-aware
473 fine-mapping as well as aiding the selection of the appropriate cell types and epigenomic
474 marks for more focused investigations of the disruptions of chromatin states by the disease-
475 specific causal variants.

476 5 ACKNOWLEDGEMENTS

477 We thank Yongjin Park, Gerald Quon, Abhishek Sarkar, and Zhizhuo Zhang for the helpful
478 discussions.

479 **Supplementary Data**

480 Supplementary Data are available at NAR Online: Supplementary Text S1,S2; Supplemen-
481 tary Tables S1-S5; Supplementary Figures S1-S6.

482 **6 FUNDING**

483 National Institutes of Health (NIH) [R01-HG004037, RC1-HG005334, R01-HG008155]. Fund-
484 ing for open access charge: NIH [R01 HG004037].

485 **6.0.1 Conflict of interest statement.**

486 None declared.

487 **References**

- 488 1. Burton, P. R. *et al.* Genome-wide association study of 14,000 cases of seven common
489 diseases and 3,000 shared controls. *Nature* **447**, 661–678 (2007).
- 490 2. Wray, N. R., Goddard, M. E. & Visscher, P. M. Prediction of individual genetic risk
491 of complex disease. *Current Opinion in Genetics & Development* **18**, 257–263 (2008).
- 492 3. Hindorff, L. A. *et al.* Potential etiologic and functional implications of genome-wide
493 association loci for human diseases and traits. *Proceedings of the National Academy
494 of Sciences* **106**, 9362–9367 (2009).
- 495 4. Visscher, P. M., Brown, M. A., McCarthy, M. I. & Yang, J. Five years of gwas
496 discovery. *The American Journal of Human Genetics* **90**, 7–24 (2012).
- 497 5. Welter, D. *et al.* The nhgri gwas catalog, a curated resource of snp-trait associations.
498 *Nucleic acids research* **42**, D1001–D1006 (2014).
- 499 6. Ernst, J. *et al.* Mapping and analysis of chromatin state dynamics in nine human cell
500 types. *Nature* **473**, 43–49 (2011).
- 501 7. Consortium, R. E. *et al.* Integrative analysis of 111 reference human epigenomes.
502 *Nature* **518**, 317–330 (2015).
- 503 8. Ward, L. D. & Kellis, M. Interpreting noncoding genetic variation in complex traits
504 and human disease. *Nature Biotechnology* **30**, 1095–1106 (2012).
- 505 9. Trynka, G. *et al.* Chromatin marks identify critical cell types for fine mapping complex
506 trait variants. *Nature genetics* **45**, 124–130 (2013).
- 507 10. Maurano, M. T. *et al.* Systematic localization of common disease-associated variation
508 in regulatory dna. *Science* **337**, 1190–1195 (2012).

- 509 11. Hnisz, D. *et al.* Super-enhancers in the control of cell identity and disease. *Cell* **155**,
510 934–947 (2013).
- 511 12. Lindblad-Toh, K. *et al.* A high-resolution map of human evolutionary constraint using
512 29 mammals. *Nature* **478**, 476–482 (2011).
- 513 13. Schaub, M. A., Boyle, A. P., Kundaje, A., Batzoglou, S. & Snyder, M. Linking disease
514 associations with regulatory information in the human genome. *Genome research* **22**,
515 1748–1759 (2012).
- 516 14. Burton, P. R. *et al.* Association scan of 14,500 nonsynonymous SNPs in four diseases
517 identifies autoimmunity variants. *Nature Genetics* **39**, 1329–1337 (2007).
- 518 15. Onengut-Gumuscu, S. *et al.* Fine mapping of type 1 diabetes susceptibility loci and
519 evidence for colocalization of causal variants with lymphoid gene enhancers. *Nature*
520 *Genetics* **47**, 381–386 (2015).
- 521 16. of the Psychiatric Genomics Consortium, C.-D. G. *et al.* Identification of risk loci
522 with shared effects on five major psychiatric disorders: a genome-wide analysis. *The*
523 *Lancet* **381**, 1371–1379 (2013).
- 524 17. O’Dushlaine, C. *et al.* Psychiatric genome-wide association study analyses implicate
525 neuronal, immune and histone pathways. *Nature Neuroscience* **18**, 199–209 (2015).
- 526 18. Kichaev, G. *et al.* Integrating functional data to prioritize causal variants in statistical
527 fine-mapping studies. *PLoS genetics* **10**, e1004722 (2014).
- 528 19. Chung, D., Yang, C., Li, C., Gelernter, J. & Zhao, H. GPA: A Statistical Approach to
529 Prioritizing GWAS Results by Integrating Pleiotropy and Annotation. *PLoS Genetics*
530 **10**, e1004787 (2014).
- 531 20. Duane, S., Kennedy, A. D., Pendleton, B. J. & Roweth, D. Hybrid monte carlo.
532 *Physics letters B* **195**, 216–222 (1987).
- 533 21. Ripke, S. *et al.* Biological insights from 108 schizophrenia-associated genetic loci.
534 *Nature* **511**, 421–427 (2014).
- 535 22. Kwan, J. S., Li, M.-X., Deng, J.-E. & Sham, P. C. Fapi: Fast and accurate p-value
536 imputation for genome-wide association study. *European Journal of Human Genetics*
537 **NA** (2015).
- 538 23. Pickrell, J. K. Joint Analysis of Functional Genomic Data and Genome-wide Associ-
539 ation Studies of 18 Human Traits. *The American Journal of Human Genetics* **94**,
540 559–573 (2014).
- 541 24. Ferrari, S. & Cribari-Neto, F. Beta regression for modelling rates and proportions.
542 *Journal of Applied Statistics* **31**, 799–815 (2004).

- 543 25. Bayes, C. L., Bazán, J. L. & García, C. A New Robust Regression Model for Pro-
544 portions. *Bayesian Analysis* **7**, 841–866 (2012). URL [http://projecteuclid.org/
545 euclid.ba/1354024464](http://projecteuclid.org/euclid.ba/1354024464).
- 546 26. Stephens, M. & Balding, D. J. Bayesian statistical methods for genetic association
547 studies. *Nature reviews Genetics* **10**, 681–690 (2009).
- 548 27. Geman, S. & Geman, D. Stochastic relaxation, gibbs distributions, and the bayesian
549 restoration of images. *Pattern Analysis and Machine Intelligence, IEEE Transactions
550 on PAMI-6*, 721–741 (1984).
- 551 28. Bernardo, J. M. & Smith, A. F. *Bayesian theory*, vol. 405 (John Wiley & Sons, 2009).
- 552 29. Neal, R. M. Mcmc using hamiltonian dynamics. *Handbook of Markov Chain Monte
553 Carlo* **2** (2011).
- 554 30. Su, Z., Marchini, J. & Donnelly, P. HAPGEN2: simulation of multiple disease SNPs.
555 *Bioinformatics (Oxford, England)* **27**, 2304–2305 (2011).
- 556 31. Durinck, S., Spellman, P. T., Birney, E. & Huber, W. Mapping identifiers for the
557 integration of genomic datasets with the r/bioconductor package biomart. *Nature
558 protocols* **4**, 1184–1191 (2009).
- 559 32. Lawrence, M., Gentleman, R. & Carey, V. rtracklayer: an r package for interfacing
560 with genome browsers. *Bioinformatics* **25**, 1841–1842 (2009).
- 561 33. Benjamini, Y. & Hochberg, Y. Controlling the false discovery rate: a practical and
562 powerful approach to multiple testing. *Journal of the Royal Statistical Society. Series
563 B (Methodological)* **57**, 289–300 (1995).
- 564 34. Eddelbuettel, D. *et al.* Rcpp: Seamless r and c++ integration. *Journal of Statistical
565 Software* **40**, 1–18 (2011).
- 566 35. Koller, D. & Friedman, N. *Probabilistic graphical models: principles and techniques*
567 (MIT press, 2009).
- 568 36. Ardlie, K. G. *et al.* The Genotype-Tissue Expression (GTEx) pilot analysis: Mul-
569 titissue gene regulation in humans. *Science* **348**, 648–660 (2015). URL [http:
570 //www.sciencemag.org/content/348/6235/648.full](http://www.sciencemag.org/content/348/6235/648.full).
- 571 37. Kheradpour, P. & Kellis, M. Systematic discovery and characterization of regulatory
572 motifs in encode tf binding experiments. *Nucleic acids research* **42**, 2976–2987 (2014).
- 573 38. Neph, S. *et al.* An expansive human regulatory lexicon encoded in transcription factor
574 footprints. *Nature* **489**, 83–90 (2012).
- 575 39. of Psoriasis Consortium & the Wellcome Trust Case Control Consortium 2, G. A.
576 *et al.* A genome-wide association study identifies new psoriasis susceptibility loci and
577 an interaction between hla-c and erap1. *Nature genetics* **42**, 985–990 (2010).

- 578 40. Tsoi, L. C. *et al.* Identification of 15 new psoriasis susceptibility loci highlights the
579 role of innate immunity. *Nature genetics* **44**, 1341–1348 (2012).
- 580 41. Ritchie, M. D., Holzinger, E. R., Li, R., Pendergrass, S. A. & Kim, D. Methods of
581 integrating data to uncover genotype-phenotype interactions. *Nature Reviews Genetics*
582 **16**, 85–97 (2015).
- 583 42. Maller, J. B. *et al.* Bayesian refinement of association signals for 14 loci in 3 common
584 diseases. *Nature genetics* **44**, 1294–1301 (2012).
- 585 43. Wen, X., Luca, F. & Pique-Regi, R. Cross-population joint analysis of eQTLs: fine
586 mapping and functional annotation. *PLoS genetics* **11**, e1005176 (2015).
- 587 44. Wallace, C. *et al.* Dissection of a complex disease susceptibility region using a bayesian
588 stochastic search approach to fine mapping. *PLoS Genet* **11**, e1005272 (2015).
- 589 45. The 1000 Genomes Project Consortium. A map of human genome variation from
590 population-scale sequencing. *Nature* **467**, 1061–1073 (2010).
- 591 46. Farh, K. K.-H. *et al.* Genetic and epigenetic fine mapping of causal autoimmune
592 disease variants. *Nature* **518**, 337–343 (2015).
- 593 47. Lu, Q., Powles, R. L., Wang, Q., He, B. J. & Zhao, H. Integrative Tissue-Specific
594 Functional Annotations in the Human Genome Provide Novel Insights on Many Com-
595 plex Traits and Improve Signal Prioritization in Genome Wide Association Studies.
596 *PLoS Genetics* **12**, e1005947 (2016).
- 597 48. Chen, W. *et al.* Fine mapping causal variants with an approximate bayesian method
598 using marginal test statistics. *Genetics* **NA**, genetics–115 (2015).
- 599 49. Kichaev, G. & Pasaniuc, B. Leveraging Functional-Annotation Data in Trans-ethnic
600 Fine-Mapping Studies. *American journal of human genetics* **97**, 260–271 (2015).
- 601 50. Hormozdiari, F., Kichaev, G., Yang, W.-Y., Pasaniuc, B. & Eskin, E. Identification
602 of causal genes for complex traits. *Bioinformatics* **31**, i206–i213 (2015).

603 Figure Legends

Figure 1: RiVIERA-beta model overview. **(a)** Inputs to RiVIERA-beta are GWAS summary statistics and epigenomic annotations for B risk loci. At a given iteration, the model samples one locus and tries to learn the intrinsic causal signals implicated in the corresponding GWAS summary data and epigenomic profiles. Highlighted variant is the causal variant based on the simulated data. **(b)** The probabilistic graphical model representation of RiVIERA-beta [35]. Variables for which distribution is defined are in circle. Epigenomic profiles are treated as observed values with no circle. The variable in shaded circle are observed (i.e., GWAS association a_{vd} and variables in unshaded circle are unobserved. The variables in red are observed and variables in blue are the variables of interest (i.e., causal indicator). The two colored plates represent K annotations (red) and V variants (blue). We model the GWAS association a_{vd} of variant v in terms of p-value sampled from Beta distribution with unknown precision ϕ_d and mean μ_d , which respectively follow an uninformative prior and a Beta distribution with hyperparameters μ_0, ϕ_0 . The latent variable c_{vd} indicates whether variant v is causal in disease d . On top of it, we dedicate an empirical prior as a linear combination of the epigenomic profile e_{vk} weighted by the epigenomic influence w_{kd} , which follows multivariate normal with zero mean and a $D \times D$ inverse covariance or precision matrix Λ_w^{-1} , where D is the number of traits that are being modeled. The linear bias w_{0d} expresses the prior belief of the causal fraction π_0 (default: 0.01). **(c)** Outputs from the model are posterior probabilities of association (PPA) for each variant in each locus, the 95% credible set containing the minimal number of SNPs whose PPA sum to 0.95 or greater, and the Bayesian estimates of the fold-enrichment of each annotation.

Figure 2: Model performance on simulated datasets. Proportion of causal variants were identified by each method as a function of increasing number of top variants selected.

Figure 3: Predicted tissue-specific epigenomic enrichments in the 10 GWAS traits. **(a)**. Hypergeometric enrichment for each of the 19 primary tissue categories using the top 5% or 43 annotations out of the 848 annotations in total for each trait based on the lower bound of the 95% credible interval of the Bayesian fold-enrichment estimates by our RiVIERA-beta; model; **(b)** enrichments for the 8 epigenomic marks among the top 43 annotations for each trait. Y-axis is the logarithmic q-values, which are the corrected p-values from the hypergeometric tests for multiple testing across traits and tissue groups or marks by Benjamini-Hochberg method [33]. On both plots, horizontal dashed bars indicate standard statistical threshold of $FDR < 0.05$.

Figure 4: Functional enrichments of credible SNPs. The top left panel displays the averaged phastCons46way conservation scores for variants in the 95% credible set (`cred_snp`) and the same number of SNPs chosen based on GWAS p-values (`gwas_snp`). The three other panels illustrate hypergeometric enrichments in terms of the $-\log_{10}$ q-values corrected for multiple testing over the 10 traits of the selected variants for GTEx whole blood eQTL located within transcription factor binding sites based on sequence motif (TFBS) (eQTL+TFBS) and genomic digital footprint (DGF) (eQTL+DGF), and eQTL in both TFBS and DGF (eQTL+TFBS+DGF).

Figure 5: Gene ontology enrichments across the 10 traits. Rows are the GO biological processes and columns are the 10 traits. Color intensities in each cell reflect the significance level in terms of $-\log_{10}$ p-value. Asterisks indicate q-values above significant cutoff after correcting for multiple testings ($FDR < 0.2$). GO names that match the pattern ‘synap|neuro|nerve’ are colored blue to highlight their exclusive association with ‘Schizophrenia’ (also in blue). Notably, GO terms ‘keratinization’ and ‘epidermis development’ (highlighted in the red box) are exclusively enriched for Psoriasis. Diseases were ordered based on hierarchical clustering based on the Pearson correlation of their GO enrichment scores.

Figure 6: Visualization of fine-mapping results. Top track: the upper panel display the RiVIERA-beta prior, genetic signals of GWAS $-\log_{10}$ p-values, and RiVIERA-beta PPA; the middle track illustrates the cumulative density of weighted epigenomic profiles colored based on the epigenomic group; the bottom tracks shows the conservation, gene annotations (Gencode 19), transcription factor binding sites (TFBS), and SNP positions. The red colored and bigger diamond plots indicate whole-blood GTEx eQTL SNPs and SNPs included into the 60% credible set, respectively. For illustration purpose, only one risk locus for Type 1 diabetes and one for Schizophrenia are shown above. The full display of 469 risk loci were in **Supplementary Fig. S6**.

Figure 7: Enrichments for eQTL using credible SNPs constructed from multi-trait joint inference. Credible SNPs for each trait were constructed based on PPA inferred by the joint RiVIERA-beta model over the 9 immune traits using 174 annotations, which are the union of the top 43 annotations detected from each trait individually. We then assessed the hypergeometric enrichments of the 95% credible sets for the GTEx whole-blood eQTL that are within DNA hypersensitive sites as defined by the genomic digital footprint data [38]. We compared these enrichment scores derived from the multi-trait model (`cred_snp_mt`) to the enrichments derived from the single-trait models either running on 43 annotations (`cred_snp_st43`) or on the 174 annotations (`cred_snp_st174`). The latter was included to control for the improvements due to the increased number of annotations (from 43 to 174).

604 Tables

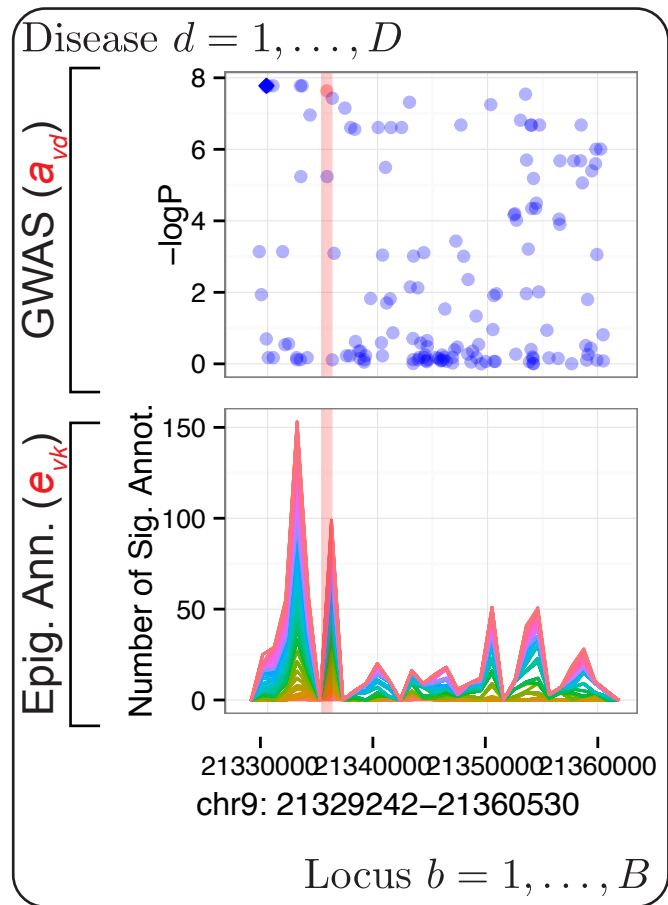
Table 1: GWAS data summary

Abbrev	Trait	Total	Loci	gwSNPs	cSNP_st	cSNPs_mt
ATD	Autoimmune Thyroid Disease	4206	8	630	38	49
CEL	Celiac Disease	29784	39	2592	344	211
JIA	Juvenile Idiopathic Arthritis	13427	22	3	440	223
MS	Multiple Sclerosis	61360	104	2096	884	339
NAR	Narcolepsy	1316	3	62	22	16
PBC	Primary Biliary Cirrhosis	14573	19	2498	172	111
PSO	Psoriasis	24832	34	457	305	171
RA	Rheumatoid Arthritis	38207	78	470	1978	719
SCZ2	Schizophrenia	54132	105	5217	2481	NA
T1D	Type 1 Diabetes	41945	57	2832	826	327

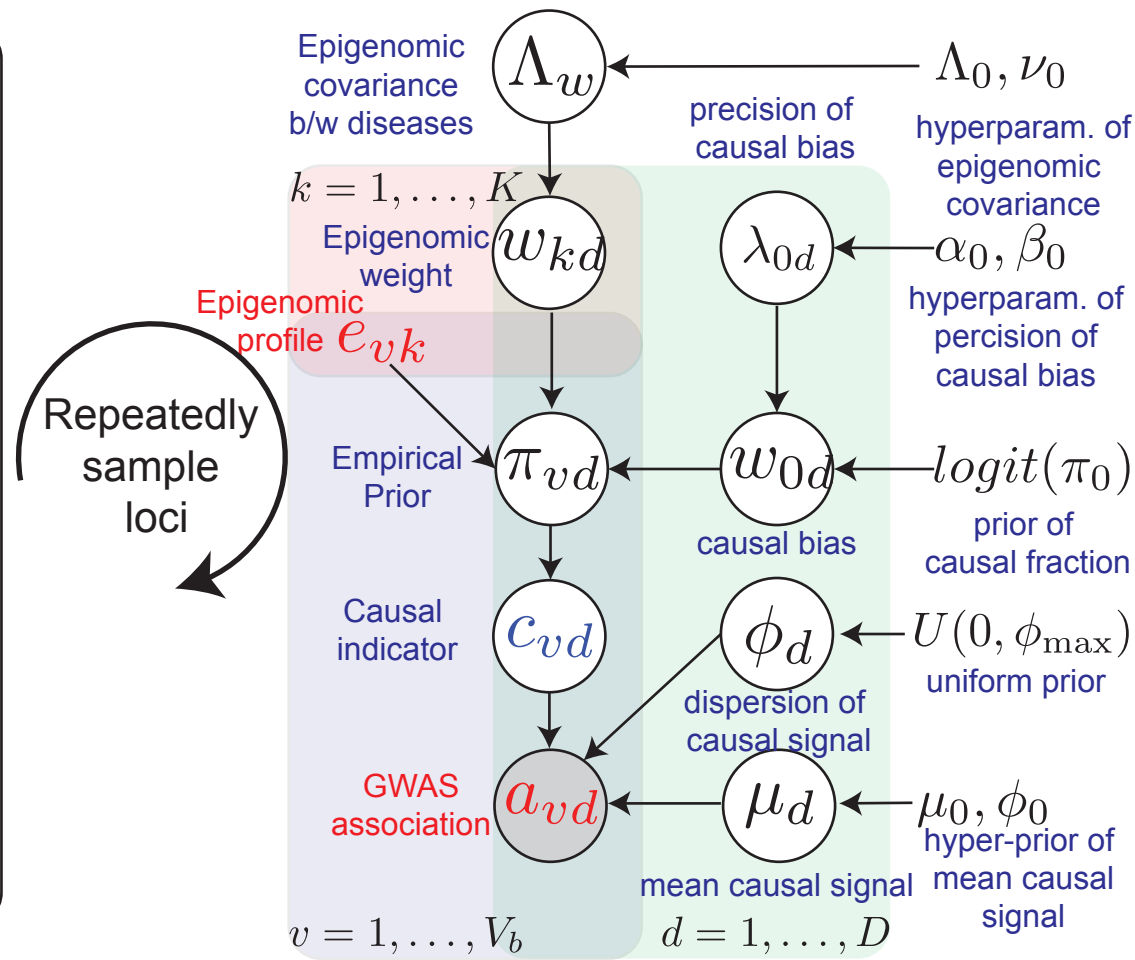
We investigated 10 GWAS traits as listed above. Abbrev: abbreviation of the trait names; Total: total number of SNPs in the risk loci with imputed and observed summary statistics; Loci: total number of risk loci for each trait; gwSNPs: SNPs that pass GWAS cutoff $p < 5e-8$; cSNP_st: total number of SNPs that are included into the 95% credible set based on single-trait risk inference using RiVIERA-beta; cSNP_mt: SNPs in 95% credible set constructed based on multi-trait joint risk inference using RiVIERA-beta across the 9 immune traits (without SCZ2).

Figure 1

a. Input



b. Model



c. Output

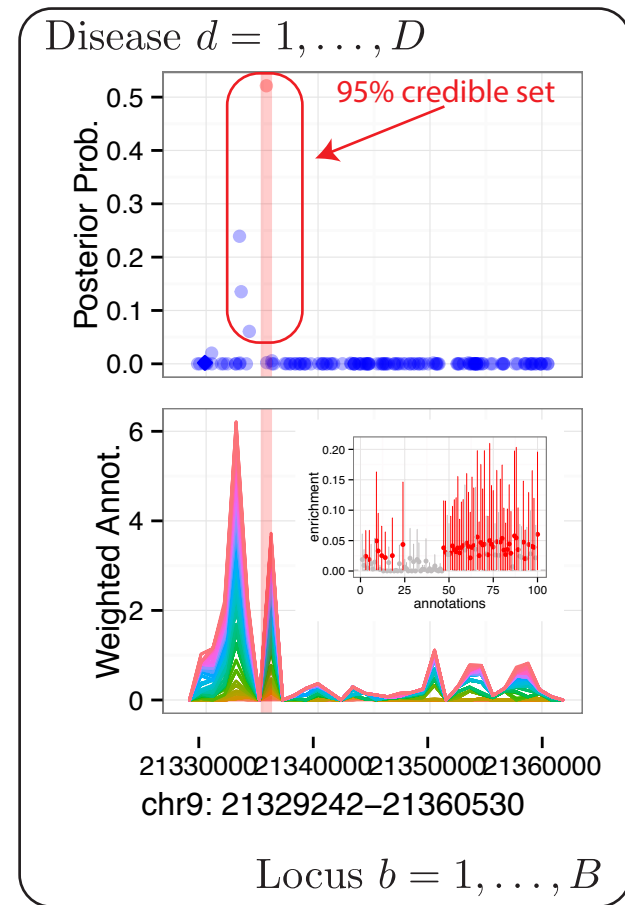


Figure 2

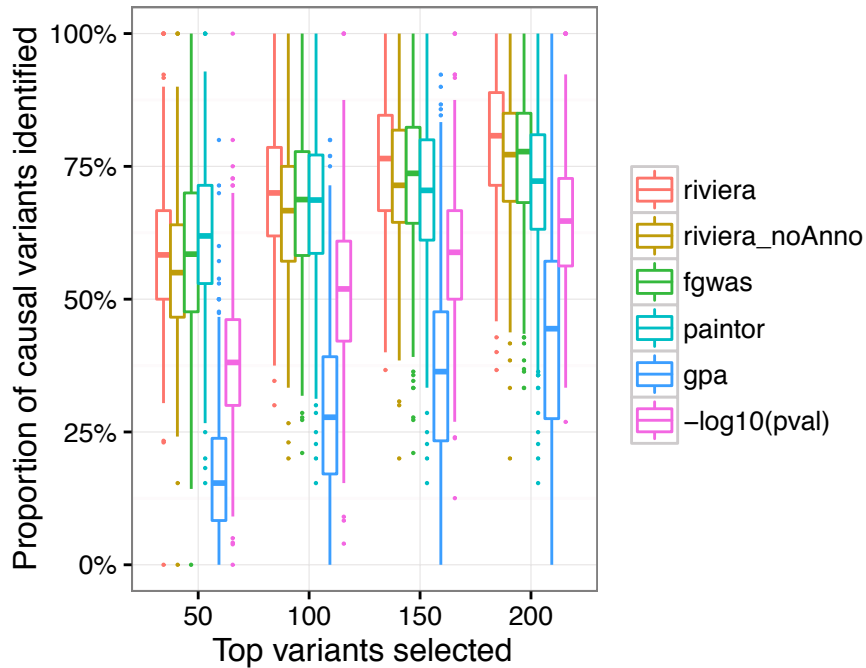


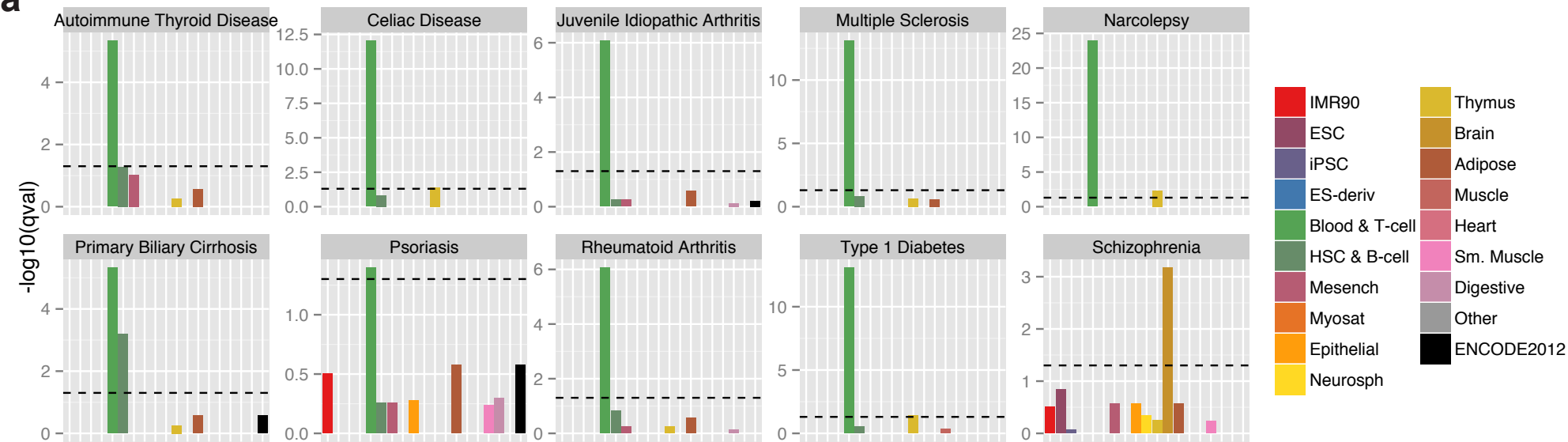
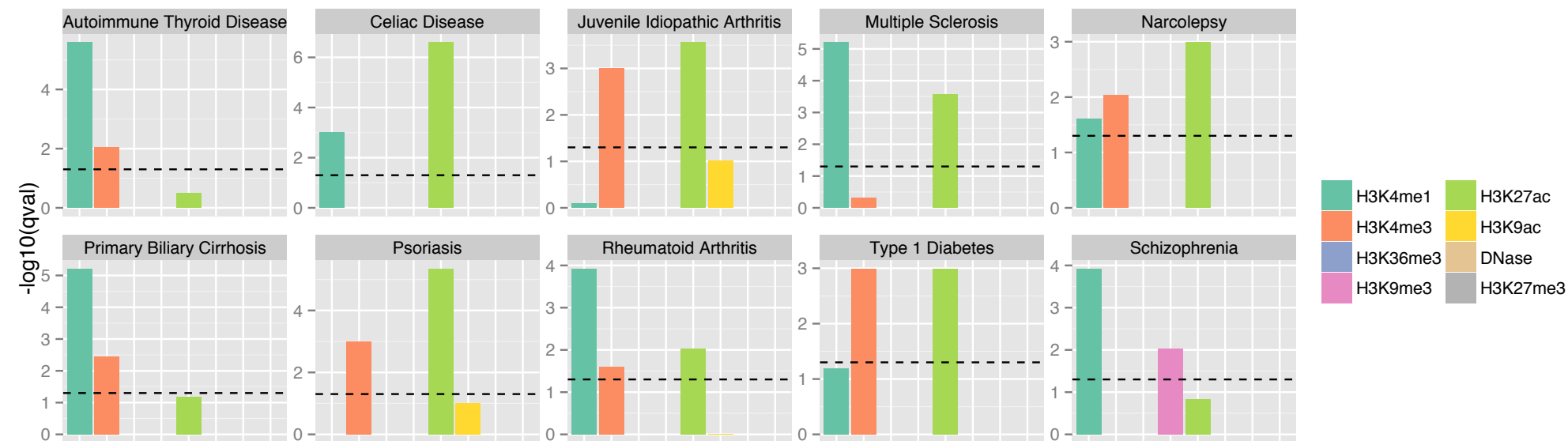
Figure 3**a****b**

Figure 4

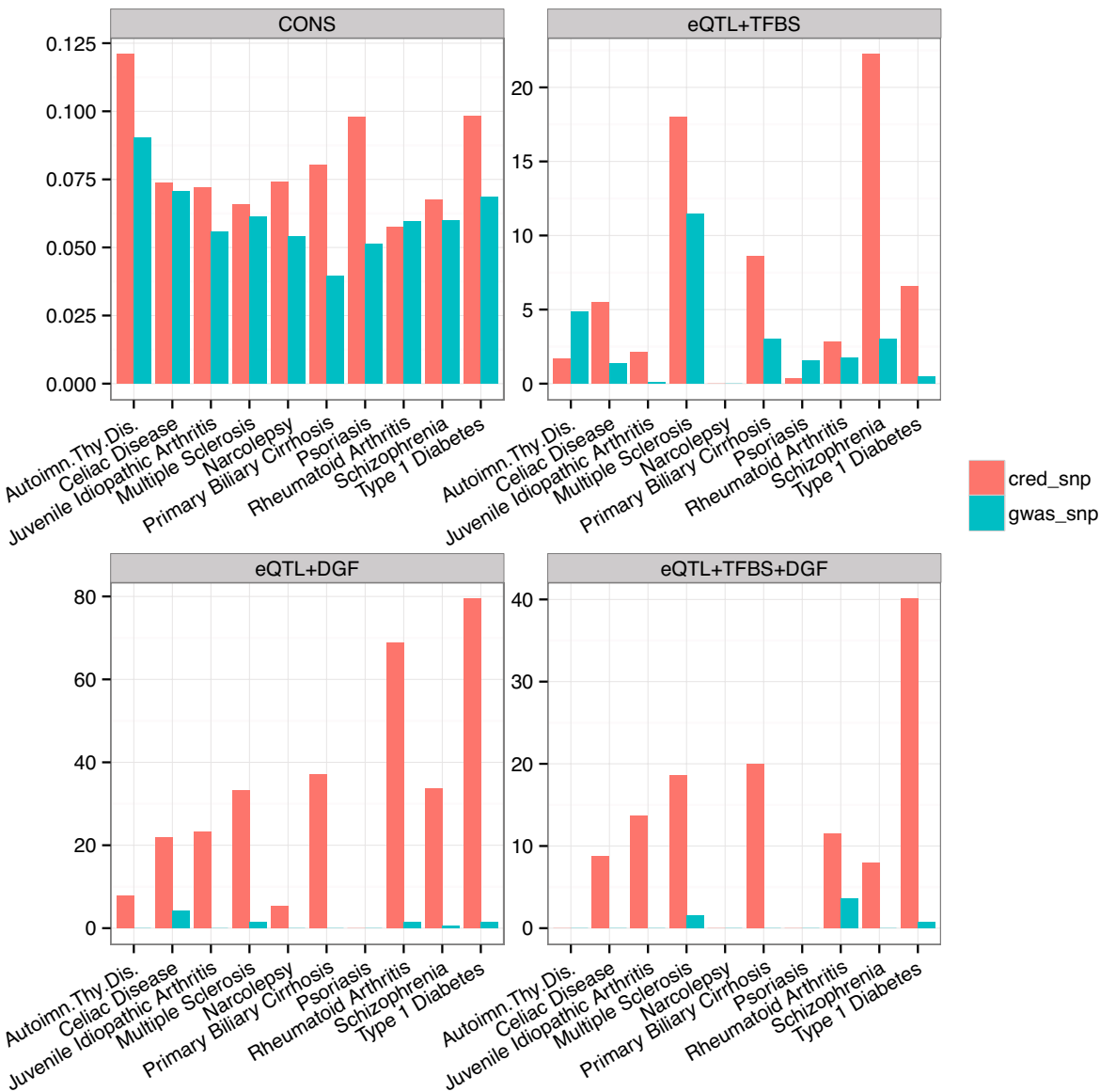


Figure 5

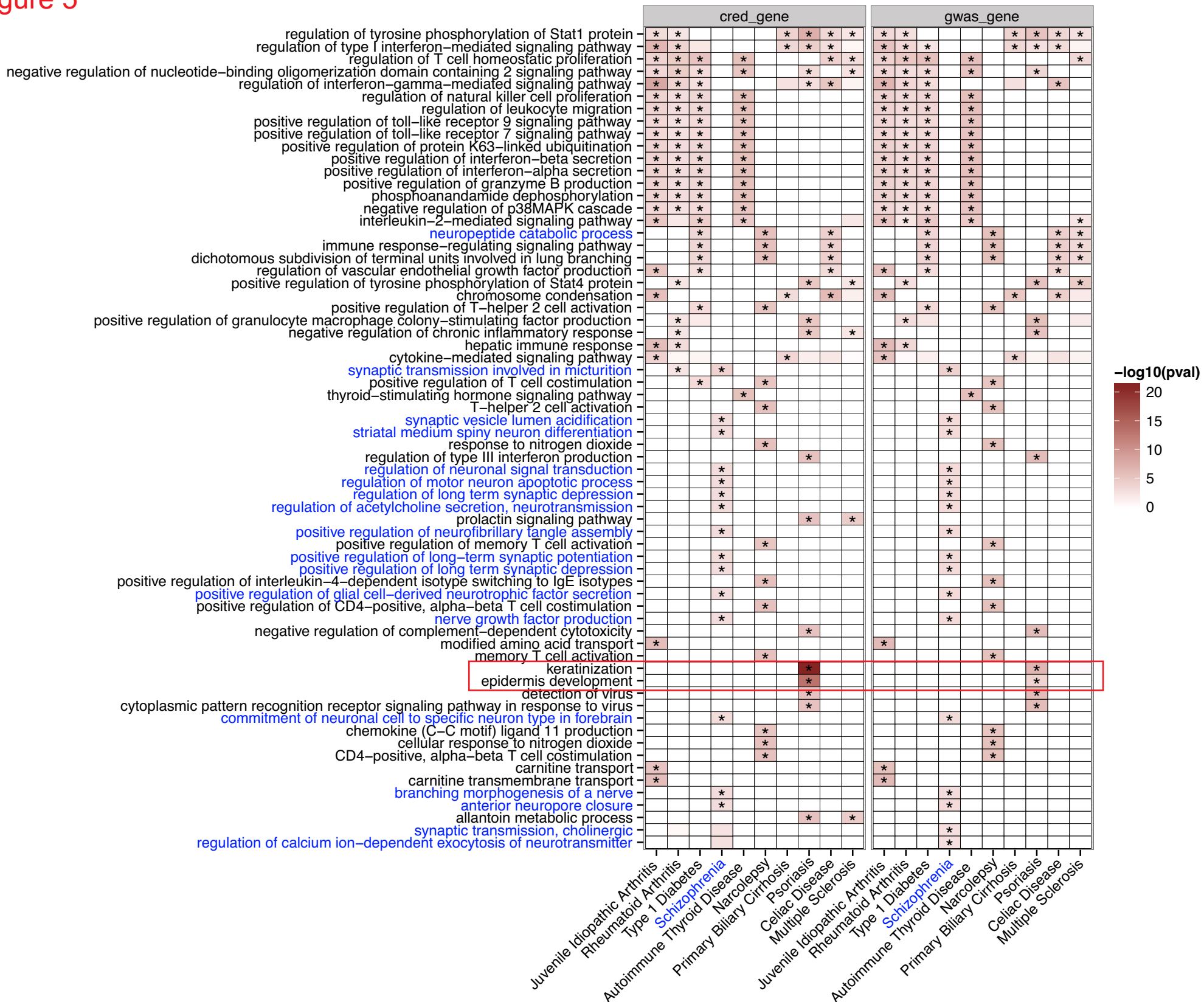
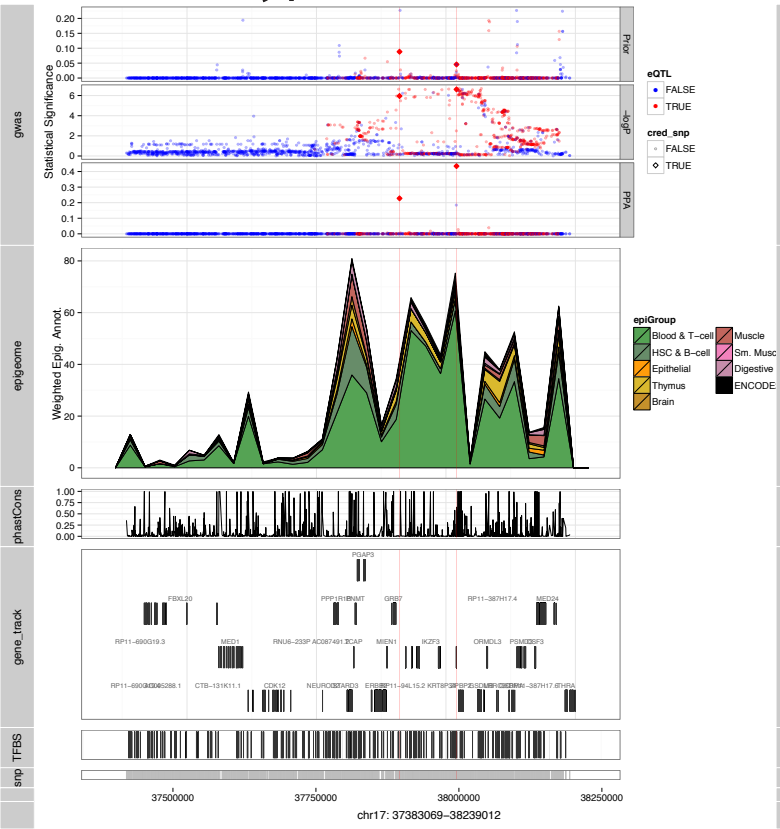


Figure 6

Type 1 Diabetes



Schizophrenia

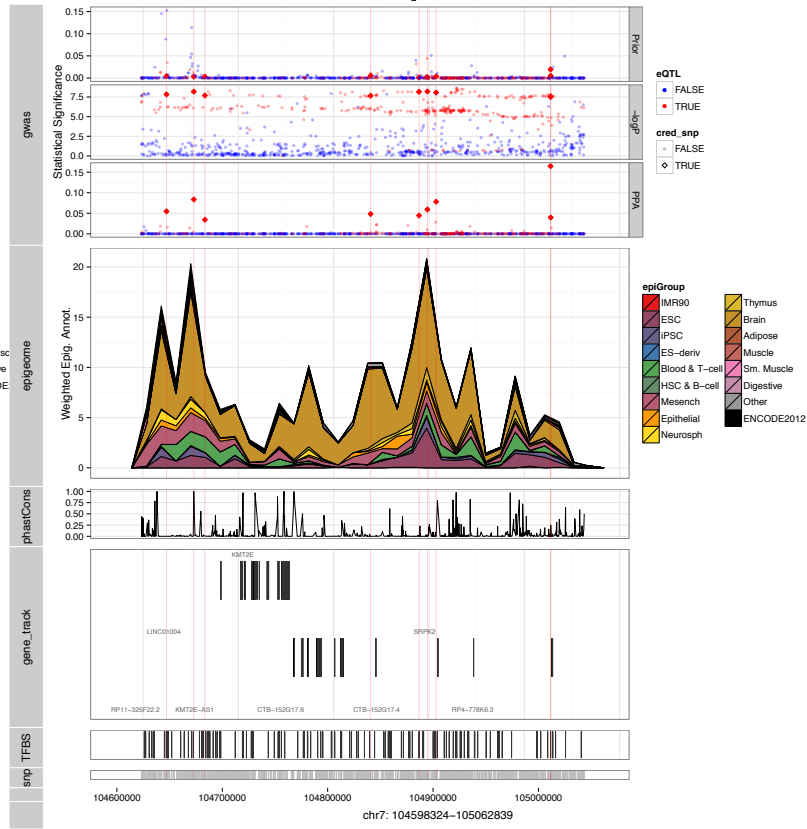


Figure 7

