# B-LORE
## *Bayesian multiple logistic regression for case-control GWAS*

Saikat Banerjee[1], Lingyao Zeng[2], Heribert Schunkert[2] and Johannes Söding[1]

[1] Max Planck Institute for Biophysical Chemistry, 37077 Göttingen, Germany
[2] German Heart Centre, 80636 Munich, Germany

*saikat.banerjee@mpibpc.mpg.de, soeding@mpibpc.mpg.de*

## 1 MOTIVATION

In genome-wide association studies (GWAS), genetic variants are tested for disease association mostly using **SIMPLE REGRESSION**, one variant at a time. This is straightforward, fast and easy to interpret but ignores the complexity of the data. Improvement in GWAS methods have explored several directions:

**META-ANALYSIS** improve power by combining summary statistics from many studies. It is only used with simple regression.

**MULTIPLE REGRESSION** aggregate evidence from multiple nearby variants. It can distinguish disease-coupled variants from those which are merely correlated with a coupled variant. It requires full genotype data. Multiple logistic regression use inefficient sampling schemes.

**FUNCTIONAL GENOMICS** data from other sources improve finemapping *i.e.* pinpoint causal SNPs. Finemapping for a meta study use approximations for LD structure of each population.
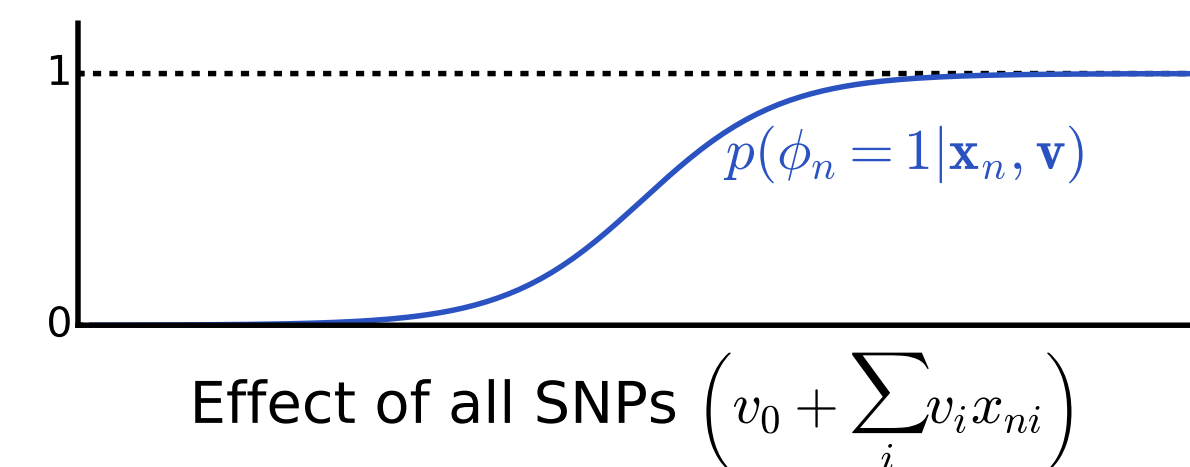
**Can we use multiple logistic regression in a meta-analysis of case control GWAS and prioritize variants with functional genomics data?**

We attempt to solve this in B-LORE, which uses a novel **quasi-Laplace approximation** to analytically integrate over variant effect sizes.

## 2 MODEL AND PRIORS

Probability of $n^{th}$ individual with genotype $\mathbf{x}_n$ to be diseased:

$$p(\phi_n = 1 \mid \mathbf{x}_n, \mathbf{v}) = \frac{\exp(\mathbf{v}^\intercal \mathbf{x}_n)}{1 + \exp(\mathbf{v}^\intercal \mathbf{x}_n)}$$



$p(\phi_n = 1 \mid \mathbf{x}_n, \mathbf{v})$

Effect of all SNPs $\left(v_0 + \sum_i v_i x_{ni}\right)$

Prior on effect sizes given hyperparameters $\boldsymbol{\theta}$ $(\pi, \mu, \sigma, \sigma_{\text{bg}})$, $p(v_i \mid \boldsymbol{\theta})$

$$= \underbrace{\pi_i \mathcal{N}(v_i \mid \mu, \sigma^2)}_{\text{Causal}} + \underbrace{(1-\pi_i)\mathcal{N}(v_i \mid 0, \sigma_{\text{bg}}^2)}_{\text{Non-causal}}$$

$$= \sum_{z_i=0,1} \pi_i^{z_i}(1-\pi_i)^{(1-z_i)} \times \mathcal{N}(v_i \mid \boldsymbol{\mu}_{\mathbf{z},i}, \text{diag}(\boldsymbol{\sigma}_{\mathbf{z},i}^2))$$

$$\mu_{\mathbf{z},i} = z_i \mu$$

$$\sigma_{\mathbf{z},i}^2 = \sigma_{\text{bg}}^2 + z_i\left(\sigma^2 - \sigma_{\text{bg}}^2\right)$$

$$\pi_i = \frac{1}{1 + \exp(-\boldsymbol{\xi}_i^\intercal \boldsymbol{\beta}_\pi)}$$

$z_i \in \{0,1\} \Rightarrow$ Indicator variable of causality

$\boldsymbol{\xi}_i \Rightarrow$ vector of local genomic features



Prior probability for non-causal $\mathcal{N}(v_i \mid 0, \sigma_{\text{bg}}^2)$
Prior probability for causal $\mathcal{N}(v_i \mid 0, \sigma_i^2)$
Likelihood $L(v_0, v_i)$
Effect size $(v_i)$

- ♦ $z_i = 1$    SNP $i$ is causal
- ♦ $z_i = 0$    SNP $i$ is non-causal

## 3 OPTIMIZATION

*Evidence approximation*: maximizing the marginal likelihood

$$mL(\boldsymbol{\theta}) = p(\boldsymbol{\phi} \mid \mathbf{x}, \boldsymbol{\theta}) = \sum_{\mathbf{z}} p(\mathbf{z} \mid \boldsymbol{\theta}) \int p(\boldsymbol{\phi} \mid \mathbf{x}, \mathbf{v}) \mathcal{N}(\mathbf{v} \mid \boldsymbol{\mu}_{\mathbf{z}}, \text{diag}(\boldsymbol{\sigma}_{\mathbf{z}}^2)) d\mathbf{v} \to \max$$

**Quasi-Laplace approximation**:

$$p(\boldsymbol{\phi} \mid \mathbf{x}, \mathbf{v}) \mathcal{N}(\mathbf{v} \mid \boldsymbol{\mu}_{\mathbf{z}}, \text{diag}(\boldsymbol{\sigma}_{\mathbf{z}}^2)) = \underbrace{p(\boldsymbol{\phi} \mid \mathbf{x}, \mathbf{v}) \mathcal{N}(\mathbf{v} \mid \tilde{\boldsymbol{\mu}}, \text{diag}(\tilde{\boldsymbol{\sigma}}^2))}_{\propto \mathcal{N}(\mathbf{v} \mid \tilde{\mathbf{v}}, \tilde{\boldsymbol{\Lambda}}^{-1})} \times \frac{\mathcal{N}(\mathbf{v} \mid \boldsymbol{\mu}_{\mathbf{z}}, \text{diag}(\boldsymbol{\sigma}_{\mathbf{z}}^2))}{\mathcal{N}(\mathbf{v} \mid \tilde{\boldsymbol{\mu}}, \text{diag}(\tilde{\boldsymbol{\sigma}}^2))}$$

The optimization can be done over multiple studies,

$$mL(\boldsymbol{\theta}) = p(\boldsymbol{\phi} \mid \mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_S, \boldsymbol{\theta}) = \int p(\boldsymbol{\phi} \mid \mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_S, \mathbf{v}) p(\mathbf{v} \mid \boldsymbol{\theta}) d\mathbf{v} \to \max$$

assuming that the quasi-Laplace approximation holds for each individual study
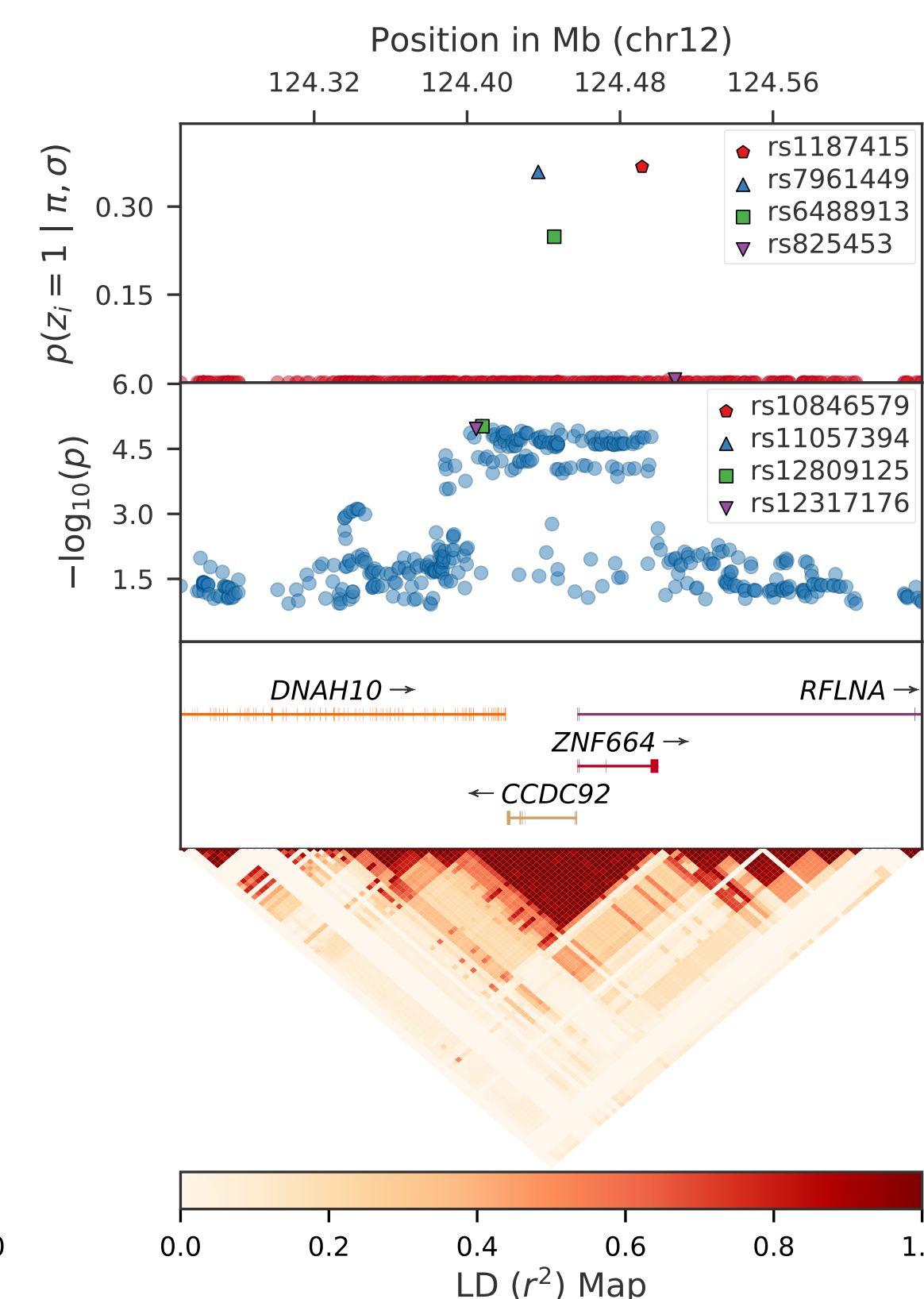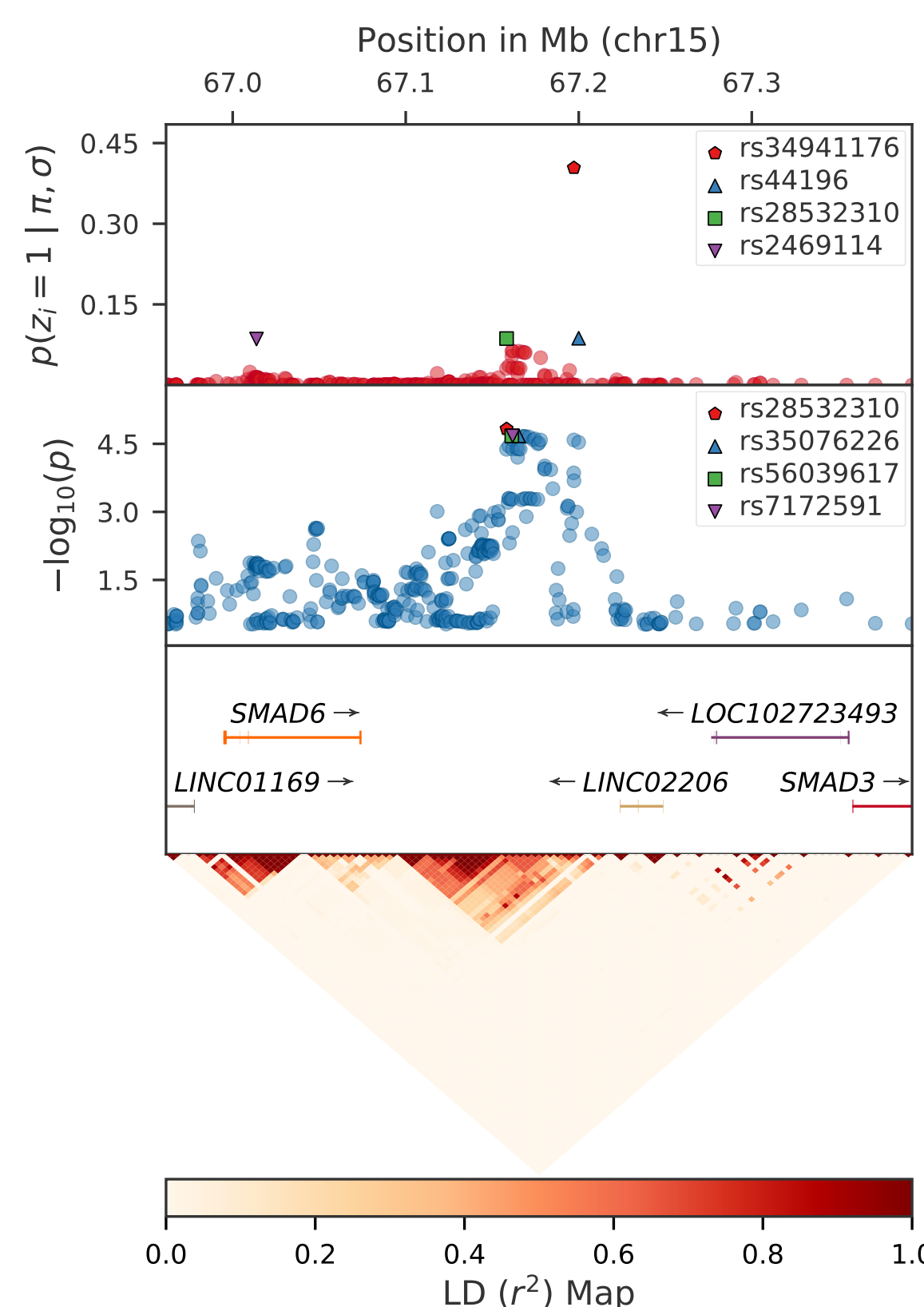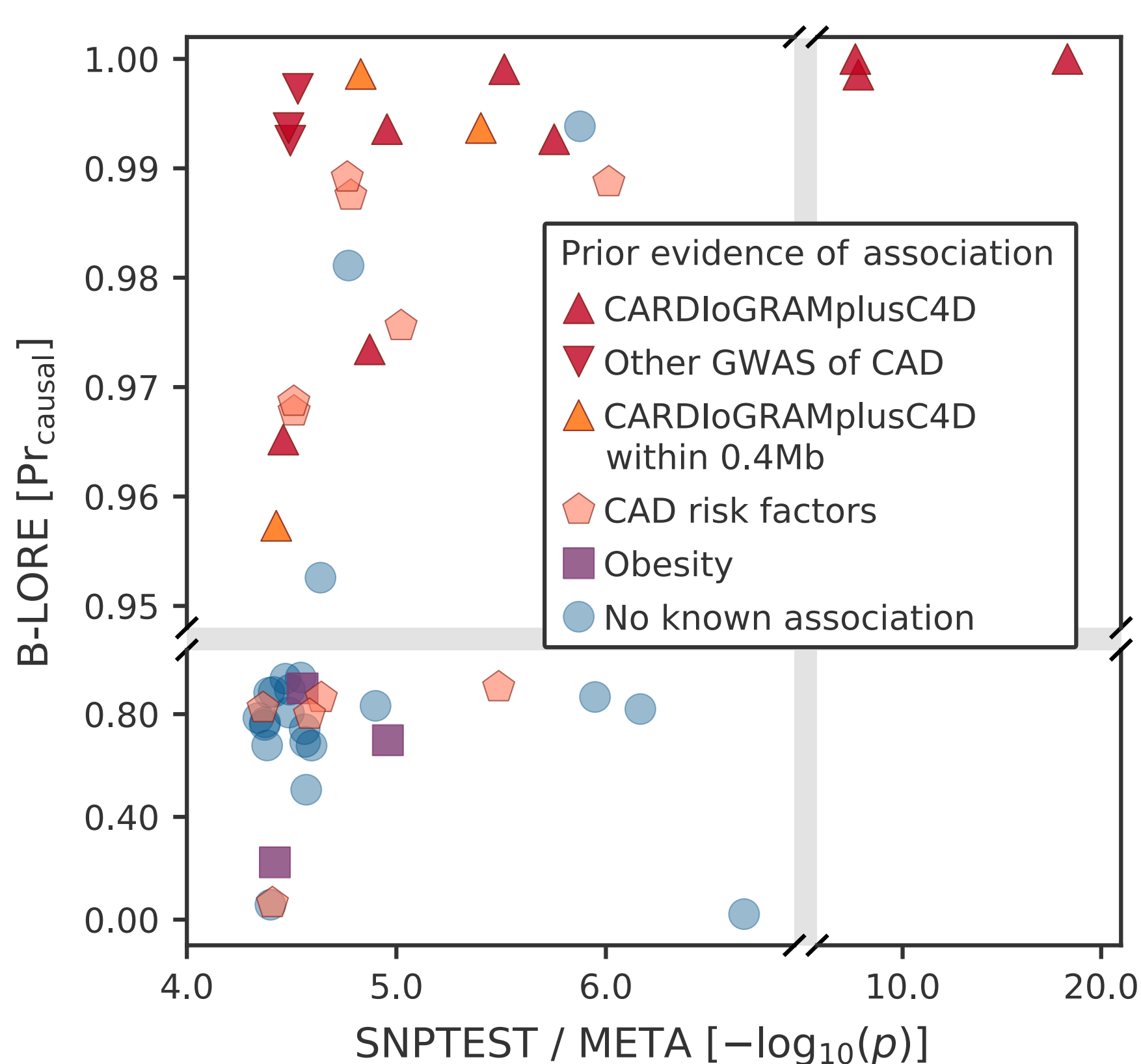
$$\prod_{s=1}^{S} \left[ p(\boldsymbol{\phi} \mid \mathbf{x}_s, \mathbf{v}) \mathcal{N}(\mathbf{v} \mid \tilde{\boldsymbol{\mu}}_{\mathbf{z},s}, \text{diag}(\tilde{\boldsymbol{\sigma}}_{\mathbf{z},s}^2)) \right] \propto \prod_{s=1}^{S} \mathcal{N}\left(\mathbf{v} \mid \tilde{\mathbf{v}}_s, \tilde{\boldsymbol{\Lambda}}_s^{-1}\right) = \mathcal{N}\left(\mathbf{v} \mid \tilde{\mathbf{v}}, \tilde{\boldsymbol{\Lambda}}^{-1}\right)$$

where $\tilde{\boldsymbol{\Lambda}} = \sum_{s=1}^{S} \tilde{\boldsymbol{\Lambda}}_s$ and $\tilde{\mathbf{v}} = \tilde{\boldsymbol{\Lambda}}^{-1} \sum_{s=1}^{S} \tilde{\boldsymbol{\Lambda}}_s \tilde{\mathbf{v}}_s$.

**B-LORE schema**

1. Two optimizations at each cohort to estimate $(\tilde{\boldsymbol{\mu}}_s, \tilde{\boldsymbol{\sigma}}_s)$ and $(\tilde{\mathbf{v}}_s, \tilde{\boldsymbol{\Lambda}}_s)$.

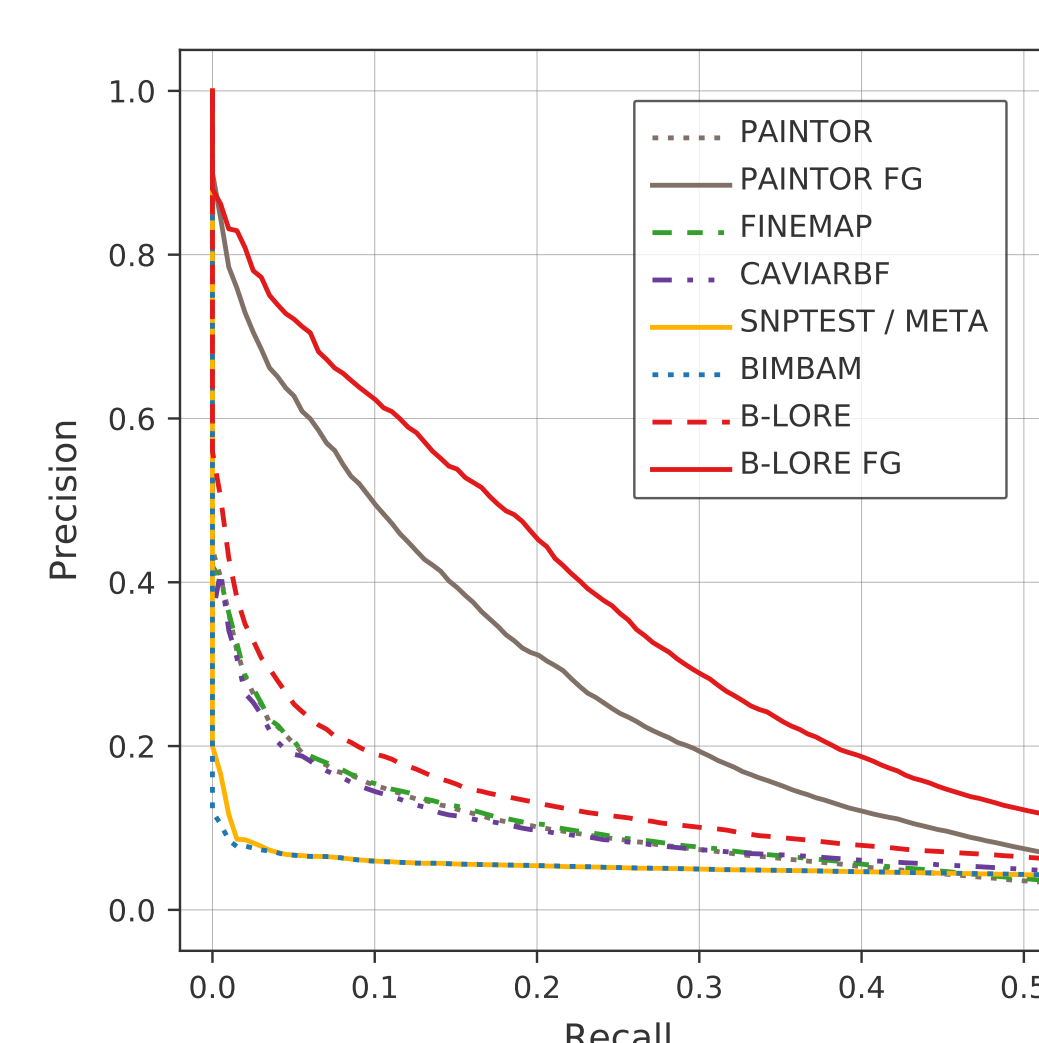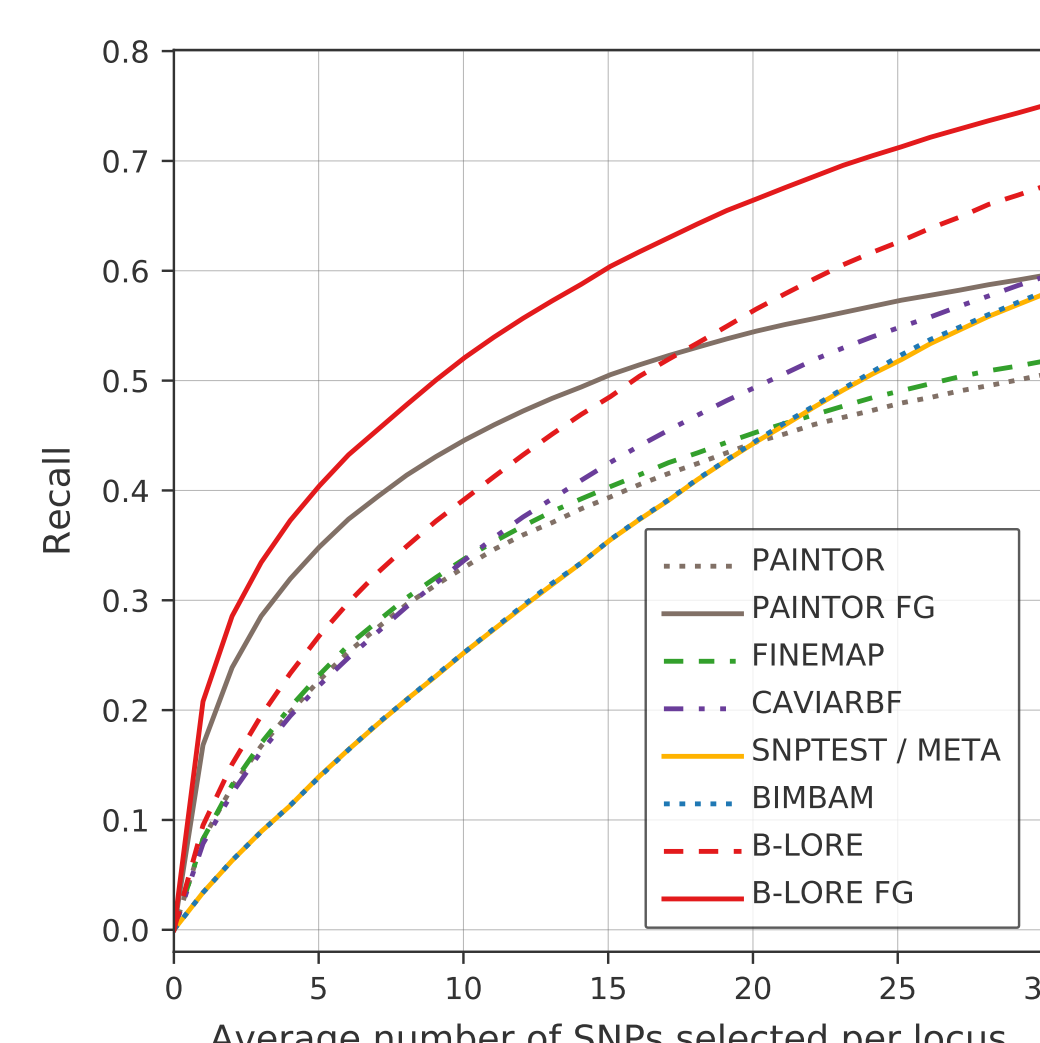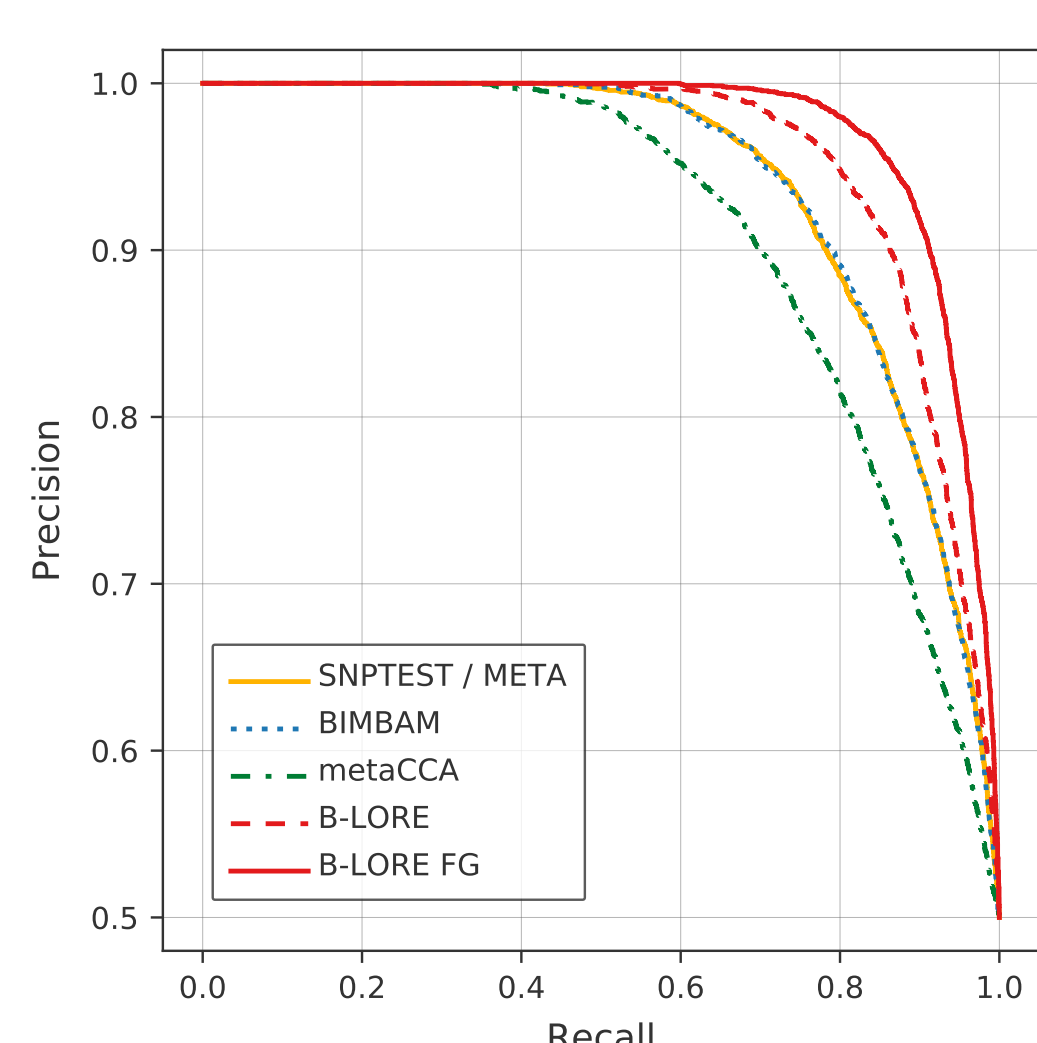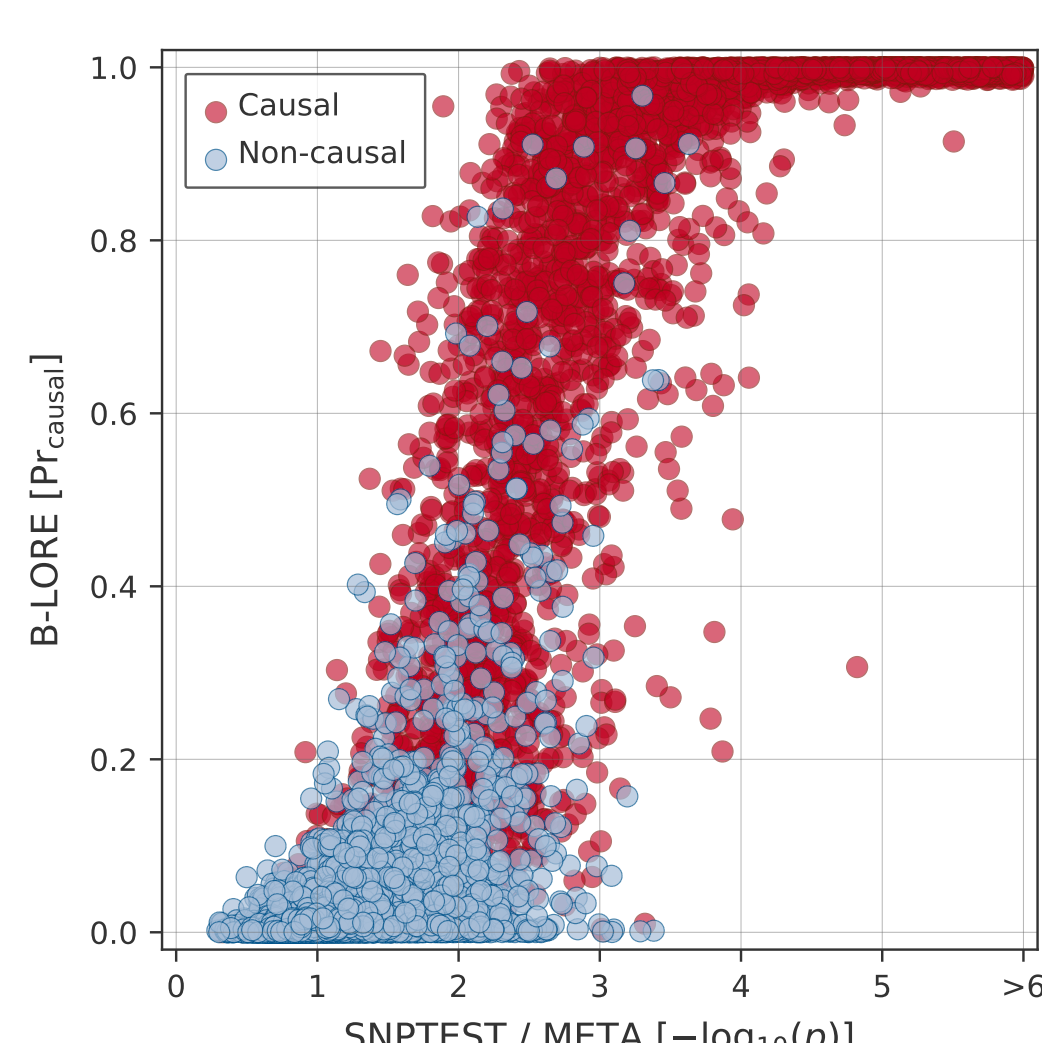2. Optimize summary statistics to estimate the hyperparameters.

## 4 INFERENCE

*Prediction of causality of each locus.*

The probability for a locus to be causally associated with the disease is

$$\text{Pr}_{\text{causal}} = p\left(\text{locus is causal} \mid \boldsymbol{\phi}, \mathbf{X}, \hat{\boldsymbol{\theta}}\right)$$
$$= 1 - p\left(\mathbf{z} = 0 \mid \boldsymbol{\phi}, \mathbf{X}, \hat{\boldsymbol{\theta}}\right)$$

*Statistical finemapping of causal variants.*

The posterior probability for SNP $i$ to be causal is

$$p\left(z_i = 1 \mid \boldsymbol{\phi}, \mathbf{X}, \hat{\boldsymbol{\theta}}\right)$$

## 5 APPLICATION ON CORONARY ARTERY DISEASE



Prior evidence of association
- ▲ CARDIoGRAMplusC4D
- ▼ Other GWAS of CAD
- ▲ CARDIoGRAMplusC4D within 0.4Mb
- ⬠ CAD risk factors
- ■ Obesity
- ● No known association

B-LORE [$\text{Pr}_{\text{causal}}$] vs SNPTEST / META [$-\log_{10}(p)$]



Position in Mb (chr15) / Position in Mb (chr12)
$p(z_i = 1 \mid \pi, \sigma)$
$-\log_{10}(p)$
LD ($r^2$) Map

Meta-analysis of 5 cohorts (Germal Myocardial Infarction Family Study, GerMIFS I-V) with a total of 6234 cases and 6848 controls from white European ancestry. We pre-selected the top 50 loci with SNPTEST / META, and applied B-LORE, using 112 functional genomics features for each SNP from DNase-seq data of the ENCODE project.

## 6 META-ANALYSES WITH REAL GENOTYPE AND SIMULATED PHENOTYPE

We selected 200 loci each with 200 SNPs from GerMIFS. We used 112 functional genomics features. We randomly selected 3 features as significant and simulated binary phenotype for $\sim 13000$ patients. Each simulation had 100 causal loci and $\sim 450$ causal SNPs with a total heritability of 0.25.



## 7 REFERENCES

1. Banerjee *et al.* bioRxiv 2017, doi:10.1101/198911
2. Servin *et al.* PLOS Genet 2007, doi:10.1371/journal.pgen.0030114
3. Guan *et al.* Ann Appl Stat 2011, doi:10.1214/11-AOAS455
4. Schunkert *et al.*, Nat Genet 2011, doi:10.1038/ng.784
5. CARDIoGRAMplusC4D Nat Genet 2015, doi:10.1038/ng.3396

## 8 ACKNOWLEDGEMENT