

B-LORE

Bayesian multiple logistic regression for case-control GWAS

Saikat Banerjee¹, Lingyao Zeng², Heribert Schunkert² and Johannes Söding¹

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

<https://github.com/soedinglab/b-lore>



1 MOTIVATION

Genetic variants in genome-wide association studies (GWAS) are tested for disease association mostly using simple regression, one variant at a time. In post-GWAS analyses, such as finemapping, **MULTIPLE REGRESSION** use multiple SNPs with Bayesian variable selection, in which a sparsity-enforcing prior on effect sizes is used to avoid overtraining. The effect sizes are integrated out for posterior inference.

MULTIPLE LOGISTIC REGRESSION has not yielded clear improvements over the linear model for binary traits in case-control GWAS.

MCMC SAMPLING has proved to be costly and technically challenging to perform the integration.

LINEAR APPROXIMATION of the logistic function is often used for case-control data.

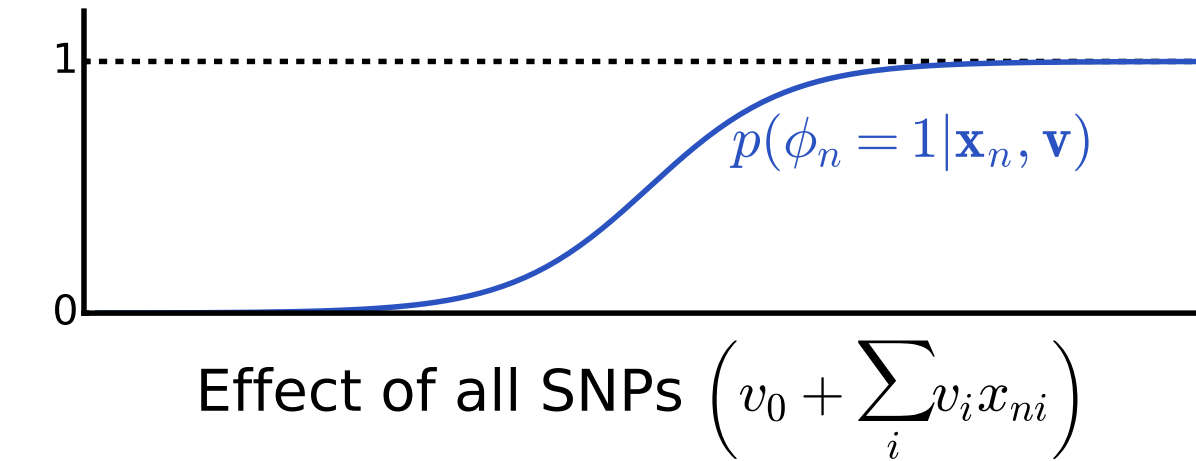
How to perform multiple logistic regression more accurately and faster?

In B-LORE, we introduce the **quasi-Laplace approximation** to analytically integrate over variant effect sizes. B-LORE improves finemapping with increasing number of controls.

2 MODEL AND PRIORS

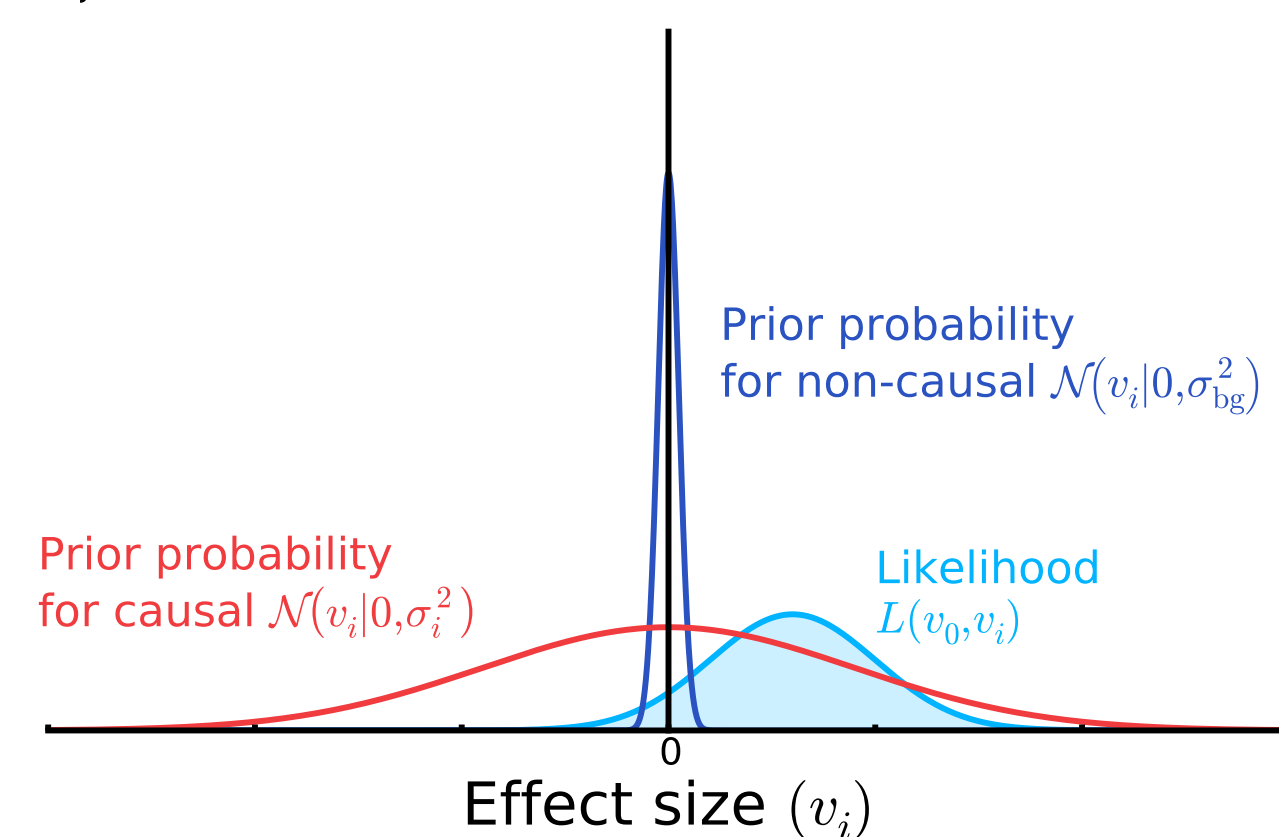
Probability of n^{th} individual with genotype \mathbf{x}_n to be diseased:

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



Prior on effect sizes given hyperparameters π and σ , $p(v_i | \pi, \sigma)$

$$\begin{aligned} &= \pi \mathcal{N}(v_i | 0, \sigma^2) + (1 - \pi) \delta_0 \\ &= \sum_{z_i=0,1} \pi^{z_i} (1 - \pi)^{(1-z_i)} \mathcal{N}(v_i | 0, \text{diag}(\sigma_{\mathbf{z},i}^2)) \\ &= \sum_{z_i=0,1} p(\mathbf{z} | \pi) \mathcal{N}(v_i | 0, \text{diag}(\sigma_{\mathbf{z},i}^2)) \end{aligned}$$



where, $\sigma_{\mathbf{z},i}^2 = z_i \sigma^2$

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

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

3 OPTIMIZATION

Evidence approximation: maximizing the marginal likelihood

$$m\mathcal{L}(\pi, \sigma) := p(\phi | \mathbf{x}, \pi, \sigma) = \sum_{\mathbf{z}} p(\mathbf{z} | \pi) \int p(\phi | \mathbf{x}, \mathbf{v}) \mathcal{N}(\mathbf{v} | \mathbf{0}, \text{diag}(\sigma_{\mathbf{z}}^2)) d\mathbf{v} \rightarrow \max$$

Quasi-Laplace approximation:

$$p(\phi | \mathbf{x}, \mathbf{v}) \mathcal{N}(\mathbf{v} | \mathbf{0}, \text{diag}(\sigma_{\mathbf{z}}^2)) = \underbrace{p(\phi | \mathbf{x}, \mathbf{v}) \mathcal{N}(\mathbf{v} | \mathbf{0}, \tilde{\sigma}^2 \mathbf{I})}_{\propto \mathcal{N}(\mathbf{v} | \tilde{\mathbf{v}}, \tilde{\Lambda}^{-1})} \frac{\mathcal{N}(\mathbf{v} | \mathbf{0}, \text{diag}(\sigma_{\mathbf{z}}^2))}{\mathcal{N}(\mathbf{v} | \mathbf{0}, \tilde{\sigma}^2 \mathbf{I})}$$

Benefits:

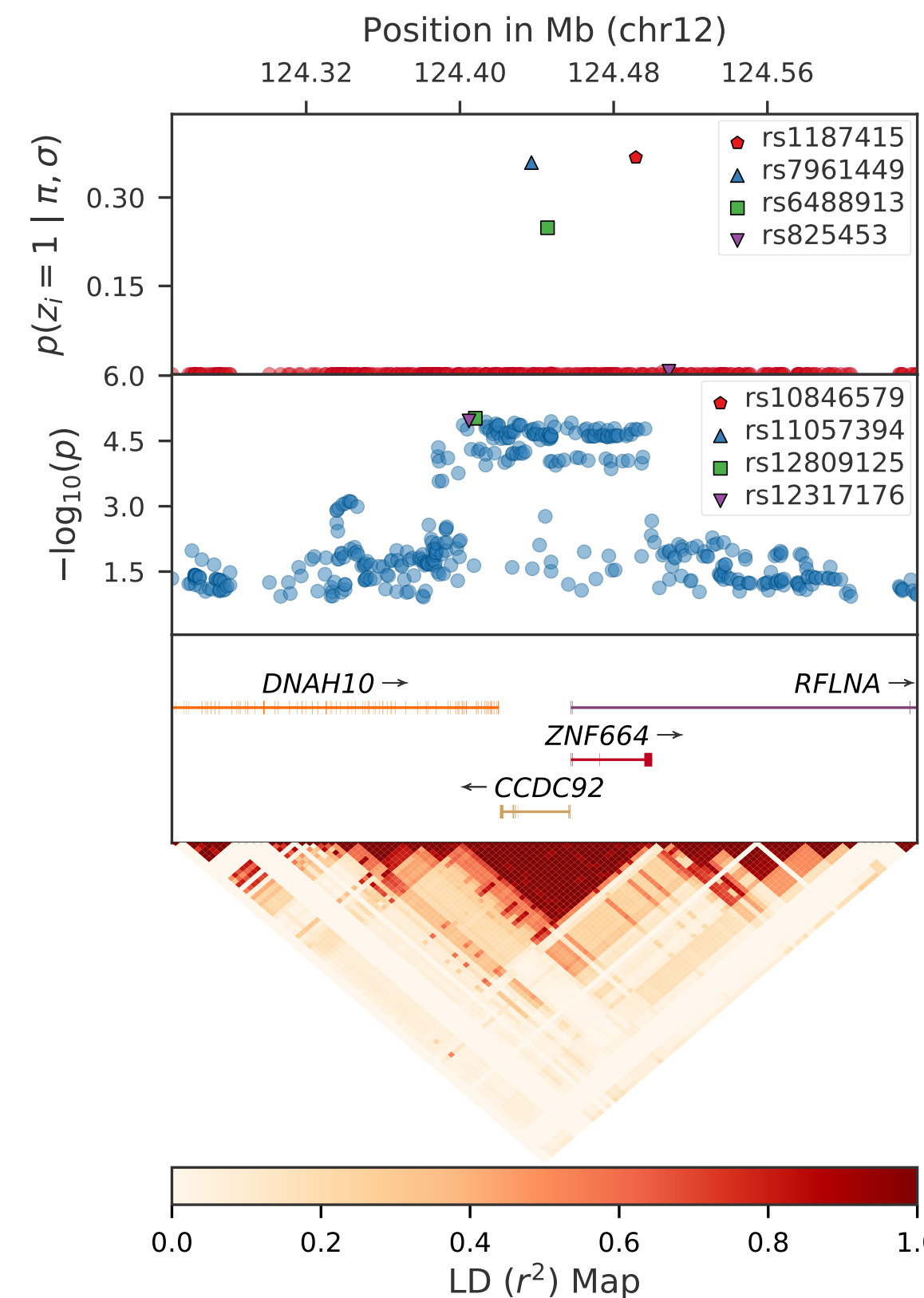
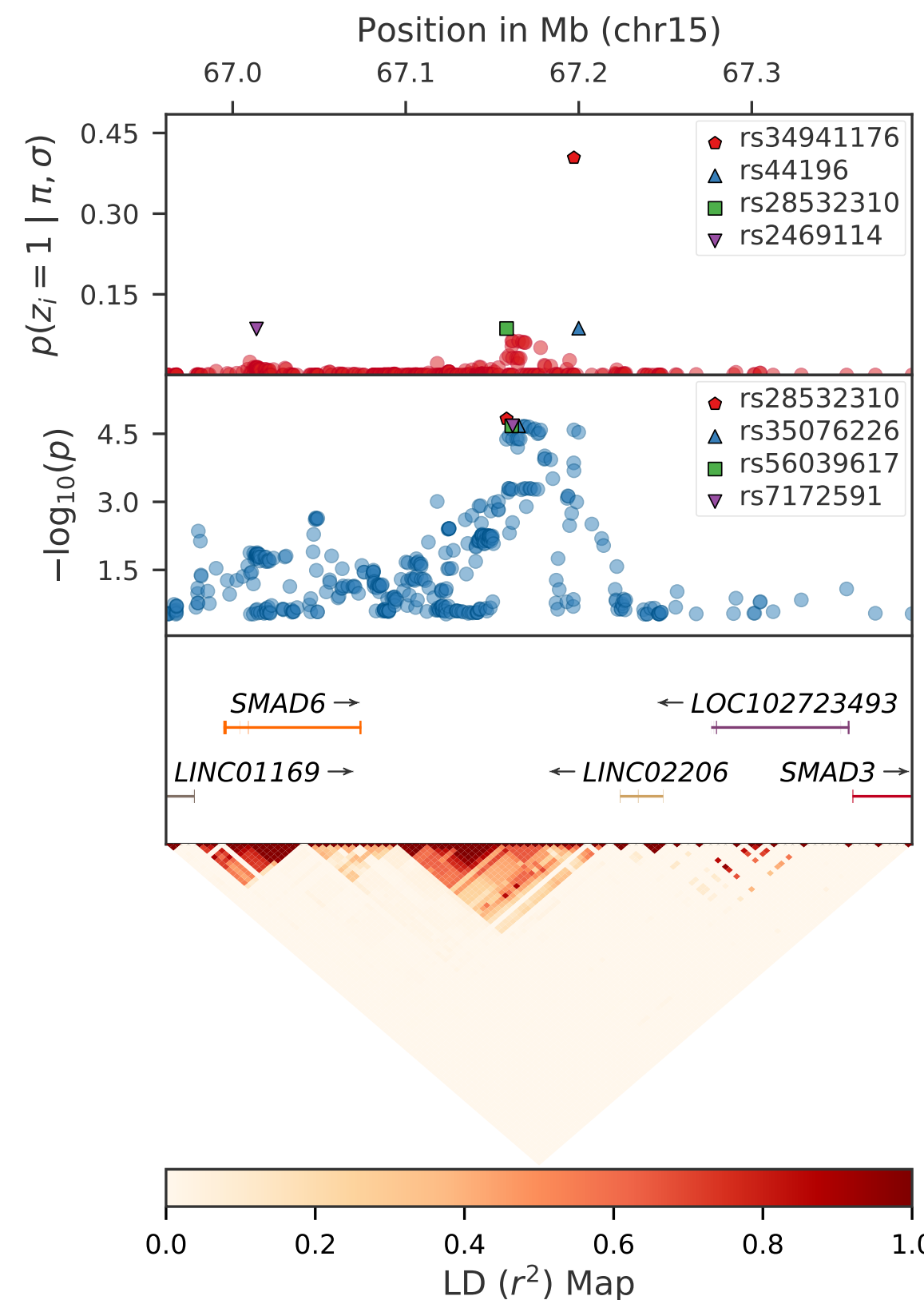
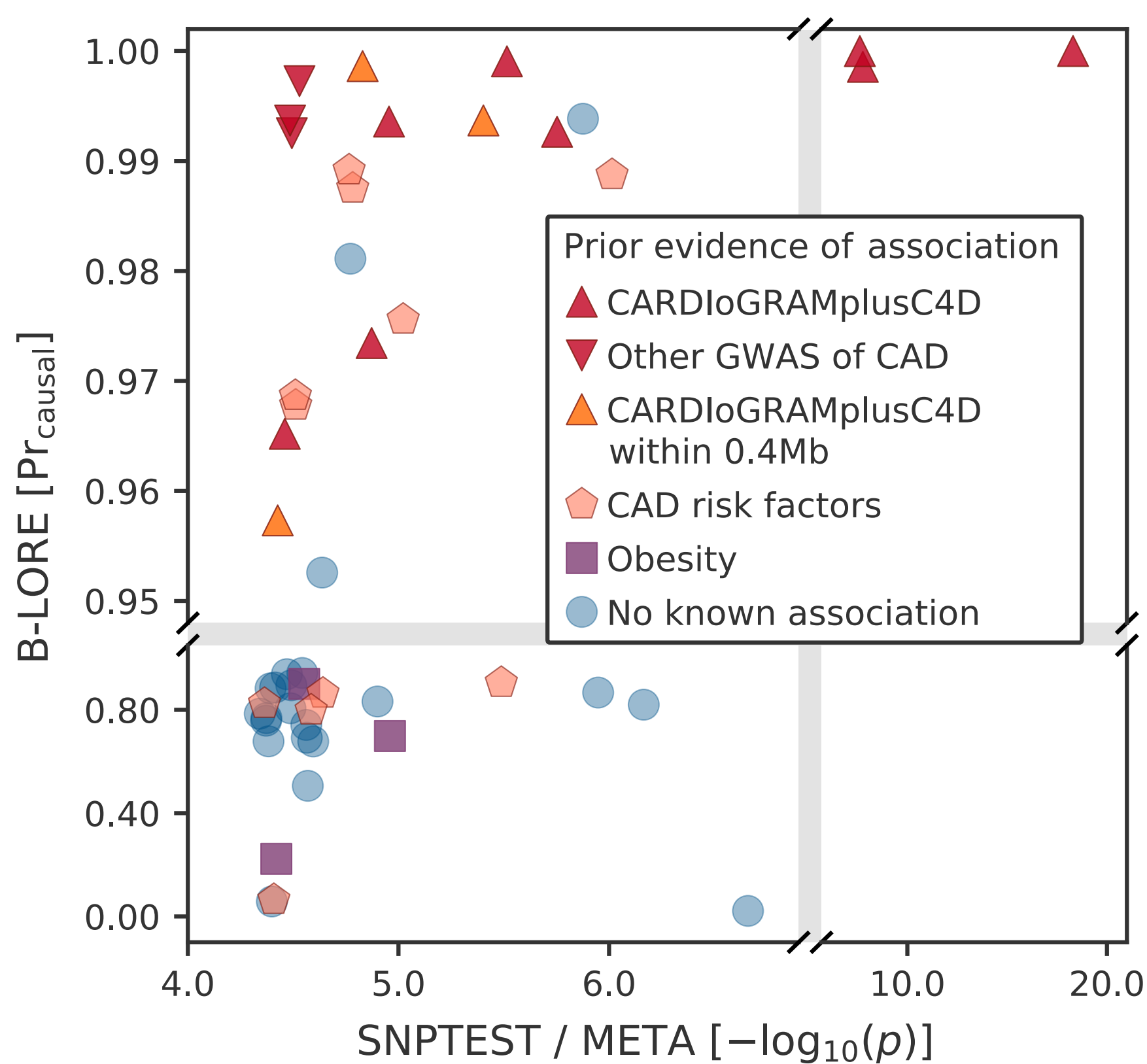
- The regularizer pulls the maximum of the regularized likelihood near to the mode of the integral, making it more accurate than Laplace approximation.
- Can be extended to multiple studies.
- Fast gradient-descent optimization.

B-LORE schema

1. Two-step optimization at each cohort to estimate $\tilde{\sigma}$ and $(\tilde{\mathbf{v}}, \tilde{\Lambda})$.

2. Estimation of hyperparameters (π, σ) .

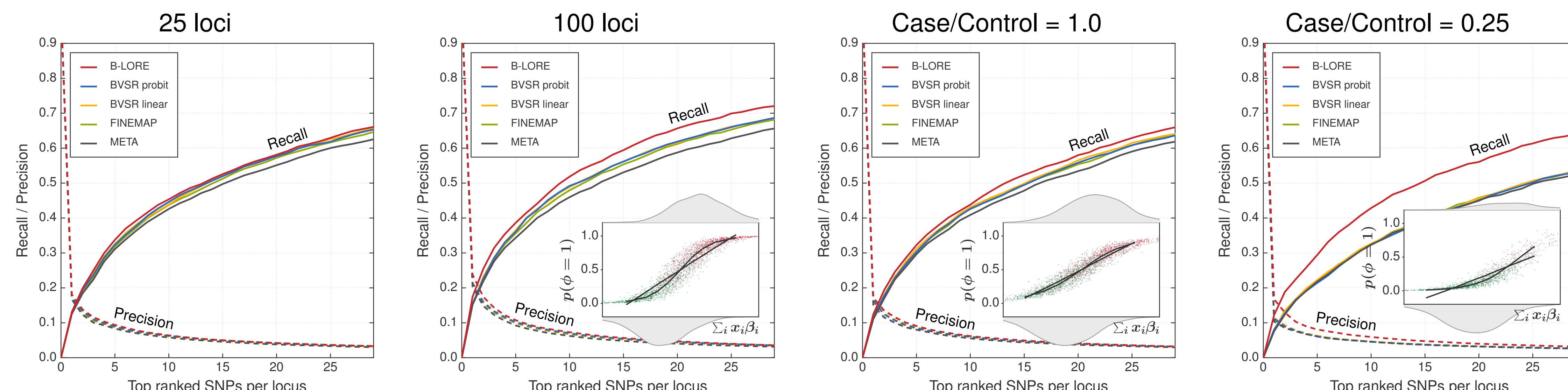
5 APPLICATION ON CORONARY ARTERY DISEASE



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 SNPTTEST / META, and applied B-LORE.

6 META-ANALYSIS WITH REAL GENOTYPE AND SIMULATED PHENOTYPE

We simulated 13082 phenotypes for 5 cohorts using 100 loci of ~200 SNPs, using one or more causal SNPs in each locus.



4 INFERENCE

Prediction of causality of each locus.

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

$$\begin{aligned} \text{Pr}_{\text{causal}} &= p(\text{locus is causal} | \phi, \mathbf{X}, \hat{\pi}, \hat{\sigma}) \\ &= 1 - p(\mathbf{z} = 0 | \phi, \mathbf{X}, \hat{\pi}, \hat{\sigma}) \end{aligned}$$

Statistical finemapping of causal variants.

The posterior probability for SNP i to be causal is

$$p(z_i = 1 | \phi, \mathbf{X}, \hat{\pi}, \hat{\sigma})$$

7 REFERENCES

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

8 ACKNOWLEDGEMENT

We thank Prof. Dr. Jeanette Erdmann and all members of the Söding lab for helpful suggestions and discussions. This work was supported by the German Federal Ministry of Education and Research (BMBF) within the framework of the e:Med research and funding concept (grant 01ZX1313A-2014).

