

Bayesian variable selection logistic regression: multivariate metaanalysis in GWAS

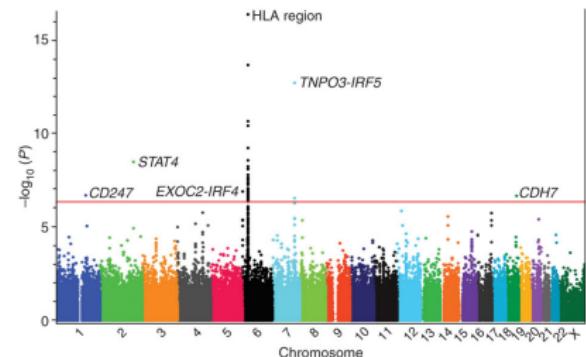
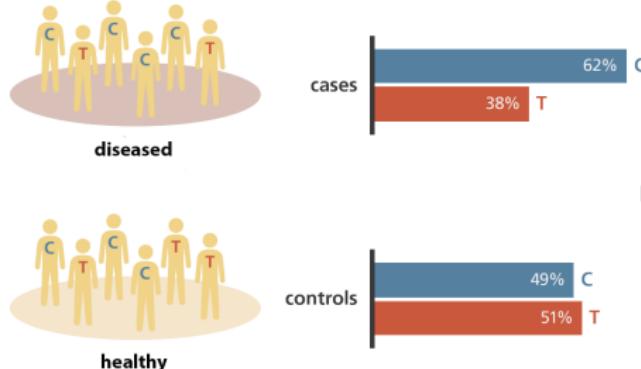
Saikat Banerjee

Max Planck Institute for Biophysical Chemistry

FEBRUARY 2, 2017



Genome-wide association studies (GWAS)



- ▶ Discovered thousands of variants associated with complex diseases

Association tests in GWAS

# of samples	Genotype (x)						Phenotype
	..	GA	TT	AT	GC	AC	..
..	GG	TT	AT	GG	AA	..	y_2
..	AA	TT	AA	GC	AA	..	y_3
..	GA	TC	AT	GC	CC	..	y_4
..	GG	TT	AT	CC	AC	..	y_5
..	AA	TC	TT	CC	AC	..	y_6

Association tests in GWAS

# of samples	Genotype (x)					Phenotype	
	GA	TT	AT	GC	AC	..	y_1
..	GG	TT	AT	GG	AA	..	y_2
..	AA	TT	AA	GC	AA	..	y_3
..	GA	TC	AT	GC	CC	..	y_4
..	GG	TT	AT	CC	AC	..	y_5
..	AA	TC	TT	CC	AC	..	y_6



Association tests in GWAS

# of samples	Genotype (x)					Phenotype	
	..	GA	TT	AT	GC	AC	..
..	GG	TT	AT	GG	AA	..	y_2
..	AA	TT	AA	GC	AA	..	y_3
..	GA	TC	AT	GC	CC	..	y_4
..	GG	TT	AT	CC	AC	..	y_5
..	AA	TC	TT	CC	AC	..	y_6



Association tests in GWAS

	Genotype (x)						Phenotype	
# of samples	..	GA	TT	AT	GC	AC	..	y_1
	..	GG	TT	AT	GG	AA	..	y_2
	..	AA	TT	AA	GC	AA	..	y_3
	..	GA	TC	AT	GC	CC	..	y_4
	..	GG	TT	AT	CC	AC	..	y_5
	..	AA	TC	TT	CC	AC	..	y_6



Association tests in GWAS

# of samples	Genotype (x)						Phenotype	
	..	GA	TT	AT	GC	AC	..	y_1
	..	GG	TT	AT	GG	AA	..	y_2
	..	AA	TT	AA	GC	AA	..	y_3
	..	GA	TC	AT	GC	CC	..	y_4
	..	GG	TT	AT	CC	AC	..	y_5
	..	AA	TC	TT	CC	AC	..	y_6



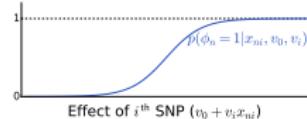
$$y_n = v_0 + v_i x_{ni} + \epsilon, \quad \text{with} \quad \epsilon \sim \mathcal{N}(0, \sigma^2) \qquad \text{Quantitative phenotype}$$

Association tests in GWAS

	Genotype (x)						Phenotype	
# of samples	..	GA	TT	AT	GC	AC	..	y_1
	..	GG	TT	AT	GG	AA	..	y_2
	..	AA	TT	AA	GC	AA	..	y_3
	..	GA	TC	AT	GC	CC	..	y_4
	..	GG	TT	AT	CC	AC	..	y_5
	..	AA	TC	TT	CC	AC	..	y_6

$$y_n = v_0 + v_i x_{ni} + \epsilon, \quad \text{with} \quad \epsilon \sim \mathcal{N}(0, \sigma^2) \quad \text{Quantitative phenotype}$$

$$p(y_n = 1 | x_{ni}, v_0, v_i) = \text{lf}(v_0 + v_i x_{ni}) \quad \text{Binary phenotype}$$



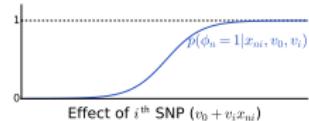
Association tests in GWAS

# of samples	Genotype (x)						Phenotype	
	..	GA	TT	AT	GC	AC	..	y_1
	..	GG	TT	AT	GG	AA	..	y_2
	..	AA	TT	AA	GC	AA	..	y_3
	..	GA	TC	AT	GC	CC	..	y_4
	..	GG	TT	AT	CC	AC	..	y_5
	..	AA	TC	TT	CC	AC	..	y_6

$$y_n = v_0 + v_i x_{ni} + \epsilon, \quad \text{with} \quad \epsilon \sim \mathcal{N}(0, \sigma^2) \quad \text{Quantitative phenotype}$$

$$p(y_n = 1 | x_{ni}, v_0, v_i) = \text{lf}(v_0 + v_i x_{ni}) \quad \text{Binary phenotype}$$

- ▶ Is the coefficient v_i significantly different from 0? \Rightarrow P-values



Strengths

- Straightforward
- Computationally fast
- Conservative
- Easy to interpret

Strengths

- Straightforward
- Computationally fast
- Conservative
- Easy to interpret

Challenges

- Linkage disequilibrium
- Genetic networks
- Low effect sizes

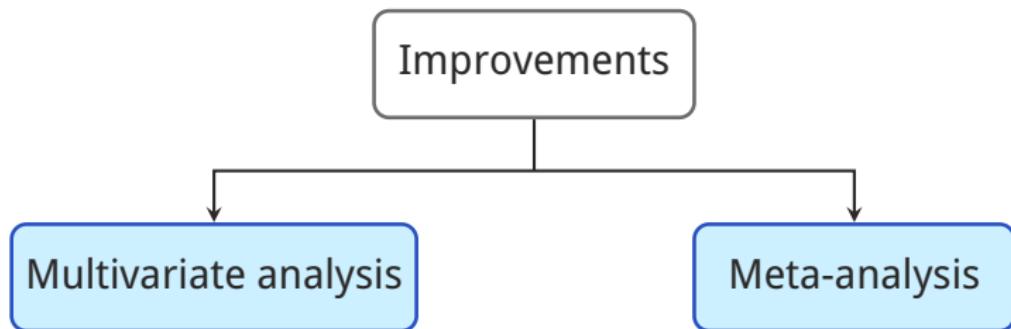
Univariate methods

Strengths

- Straightforward
- Computationally fast
- Conservative
- Easy to interpret

Challenges

- Linkage disequilibrium
- Genetic networks
- Low effect sizes

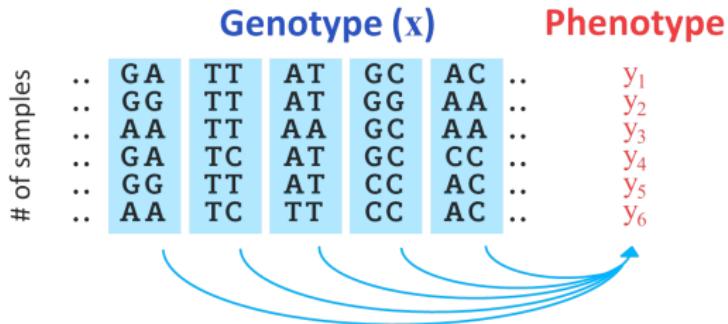


Multivariate methods

# of samples	Genotype (x)					Phenotype	
	..	GA	TT	AT	GC	AC	..
	..	GG	TT	AT	GG	AA	y_1
	..	AA	TT	AA	GC	AA	y_2
	..	GA	TC	AT	GC	CC	y_3
	..	GG	TT	AT	CC	AC	y_4
	..	AA	TC	TT	CC	AC	y_5
							y_6



Multivariate methods



Multivariate methods

# of samples	Genotype (x)					Phenotype		
	..	GA	TT	AT	GC	AC	..	y_1
	..	GG	TT	AT	GG	AA	..	y_2
	..	AA	TT	AA	GC	AA	..	y_3
	..	GA	TC	AT	GC	CC	..	y_4
	..	GG	TT	AT	CC	AC	..	y_5
	..	AA	TC	TT	CC	AC	..	y_6



The diagram illustrates the relationship between genotype and phenotype. It shows a grid of genotype data (x) with rows labeled by sample index (..). Each row contains five genotype entries (e.g., GA, TT, AT, GC, AC). To the right of the grid, six phenotype values (y1 through y6) are listed. A series of blue curved arrows originates from the right side of each row and points towards the corresponding phenotype value, indicating the mapping from genotype to phenotype.

$$y_n = v_0 + \sum_i v_i x_{ni} + \epsilon, \quad \text{with} \quad \epsilon \sim \mathcal{N}(0, \sigma^2) \quad \text{Quantitative phenotype}$$

Multivariate methods

	Genotype (x)					Phenotype		
# of samples	..	GA	TT	AT	GC	AC	..	y_1
..	GG	TT	AT	GG	AA	AA	..	y_2
..	AA	TT	AA	GC	AA	AA	..	y_3
..	GA	TC	AT	GC	CC	CC	..	y_4
..	GG	TT	AT	CC	AC	AC	..	y_5
..	AA	TC	TT	CC	AC	AC	..	y_6



$$y_n = v_0 + \sum_i v_i x_{ni} + \epsilon, \quad \text{with} \quad \epsilon \sim \mathcal{N}(0, \sigma^2) \quad \text{Quantitative phenotype}$$

$$p(y_n = 1 | x_{ni}, v_0, v_i) = \text{lf}\left(v_0 + \sum_i v_i x_{ni}\right) \quad \text{Binary phenotype}$$

Multivariate methods

	Genotype (x)					Phenotype		
# of samples	..	GA	TT	AT	GC	AC	..	y_1
..	GG	TT	AT	GG	AA	AA	..	y_2
..	AA	TT	AA	GC	AA	AA	..	y_3
..	GA	TC	AT	GC	CC	CC	..	y_4
..	GG	TT	AT	CC	AC	AC	..	y_5
..	AA	TC	TT	CC	AC	AC	..	y_6



$$y_n = v_0 + \sum_i v_i x_{ni} + \epsilon, \quad \text{with} \quad \epsilon \sim \mathcal{N}(0, \sigma^2) \quad \text{Quantitative phenotype}$$

$$p(y_n = 1 | x_{ni}, v_0, v_i) = \text{lf}\left(v_0 + \sum_i v_i x_{ni}\right) \quad \text{Binary phenotype}$$

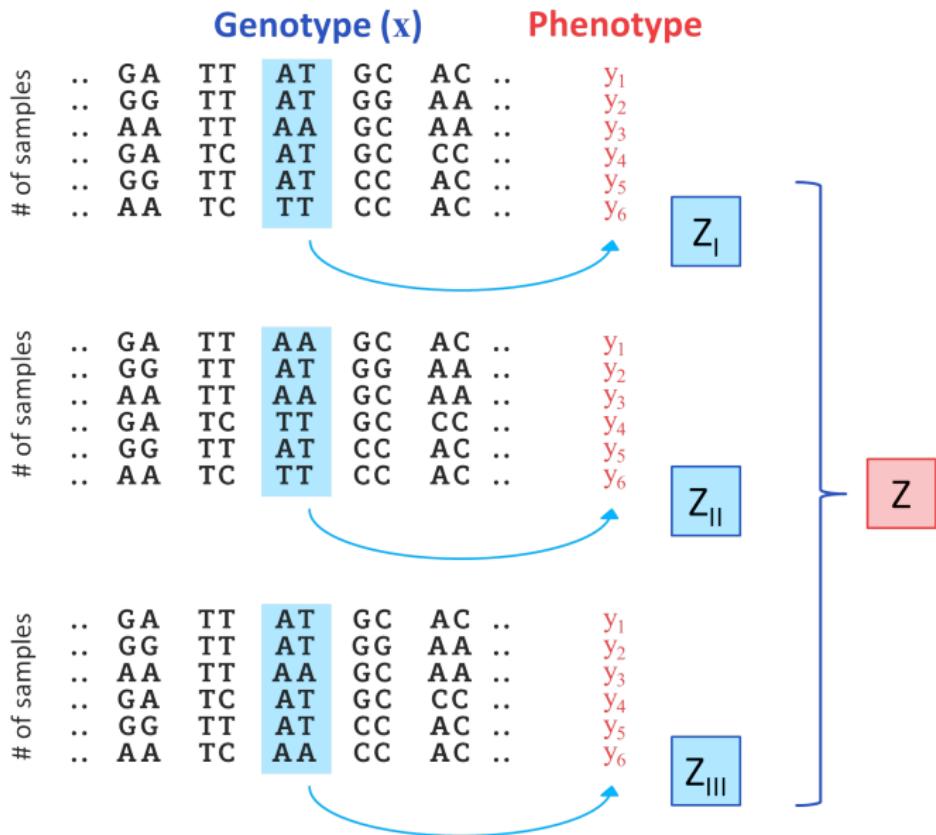
- ▶ Multivariate methods perform better than univariate methods

Meta-analysis

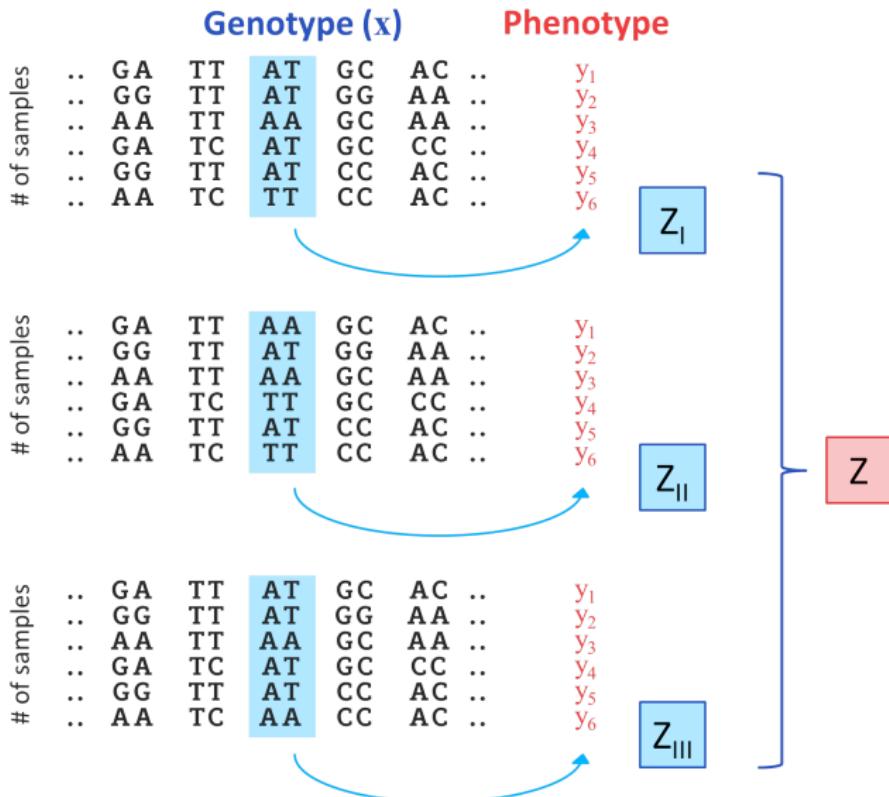
# of samples	Genotype (x)						Phenotype	
	..	GA	TT	AT	GC	AC	..	y_1
..	GG	TT	AT	GG	AA	..	y_2	
..	AA	TT	AA	GC	AA	..	y_3	
..	GA	TC	AT	GC	CC	..	y_4	
..	GG	TT	AT	CC	AC	..	y_5	
..	AA	TC	TT	CC	AC	..	y_6	



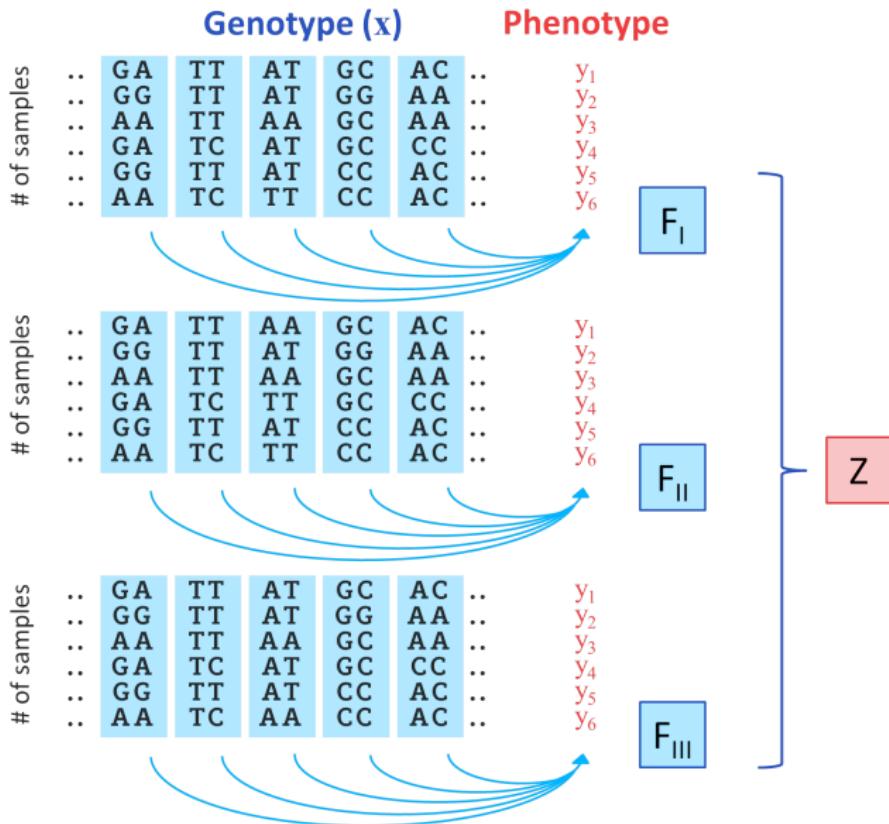
Meta-analysis



Goal of our method



Goal of our method



Bayesian variable selection regression (BVSR)

BIMBAM

Servin and Stephens, *PLoS Genetics* 2007
Guan and Stephens, *Ann. Appl. Stats.* 2011

$$y_n = v_0 + \sum_i v_i x_{ni} + \epsilon, \quad \text{with} \quad \epsilon \sim \mathcal{N}(0, \tau^{-1}) \quad \text{Quantitative phenotype}$$

Bayesian variable selection regression (BVSR)

BIMBAM

Servin and Stephens, *PLoS Genetics* 2007
Guan and Stephens, *Ann. Appl. Stats.* 2011

$$y_n = v_0 + \sum_i v_i x_{ni} + \epsilon, \quad \text{with} \quad \epsilon \sim \mathcal{N}(0, \tau^{-1}) \quad \text{Quantitative phenotype}$$

- ▶ Likelihood for N patients:

$$p(\mathbf{y} | \mathbf{x}, \mathbf{v}, \tau) = \mathcal{N}(\mathbf{y} | \mathbf{x}^\top \mathbf{v}, \tau^{-1} \mathbb{I})$$

Bayesian variable selection regression (BVSR)

BIMBAM

Servin and Stephens, *PLoS Genetics* 2007
Guan and Stephens, *Ann. Appl. Stats.* 2011

$$y_n = v_0 + \sum_i v_i x_{ni} + \epsilon, \quad \text{with} \quad \epsilon \sim \mathcal{N}(0, \tau^{-1}) \quad \text{Quantitative phenotype}$$

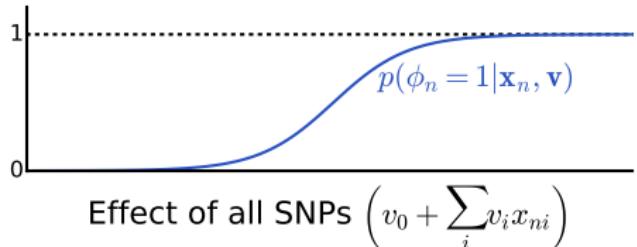
- ▶ Likelihood for N patients:

$$p(\mathbf{y} | \mathbf{x}, \mathbf{v}, \tau) = \mathcal{N}(\mathbf{y} | \mathbf{x}^T \mathbf{v}, \tau^{-1} \mathbb{I})$$

- ▶ Number of SNPs \gg samples → Overfitting
- ▶ Effective priors on \mathbf{v} for sparsity

Bayesian variable selection logistic regression (BVSLR)

$$p(\phi_n = 1 | \mathbf{x}_n, \mathbf{v}) = \text{lf}\left(v_0 + \sum_i v_i x_{ni}\right)$$

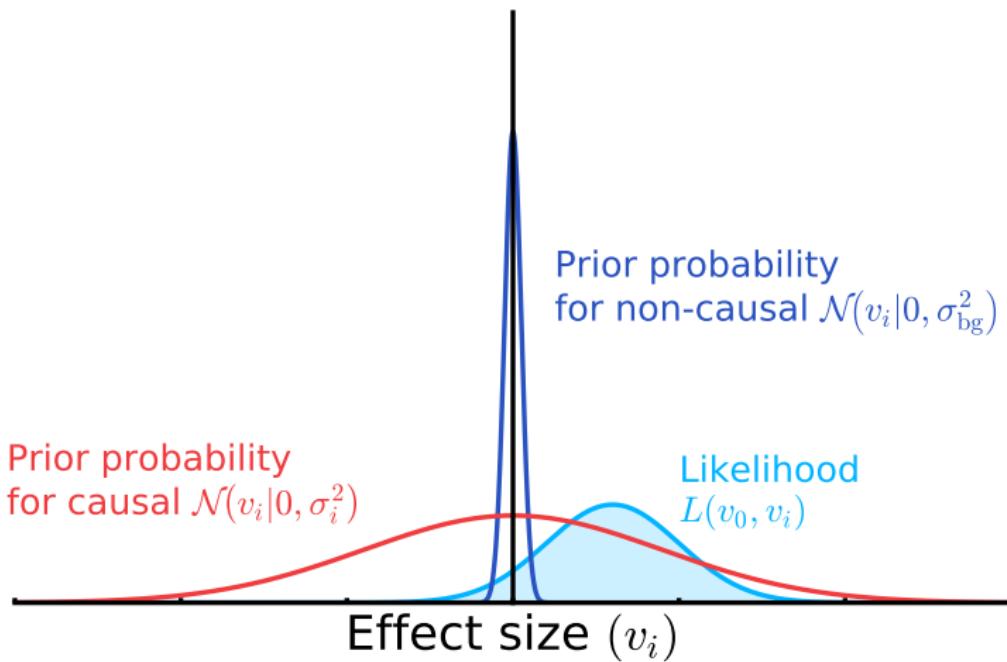


- Likelihood for N patients:

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

Sparsity in BVSR / BVSLR

Looks at both **null hypothesis** and **alternate hypothesis**



Simulation details

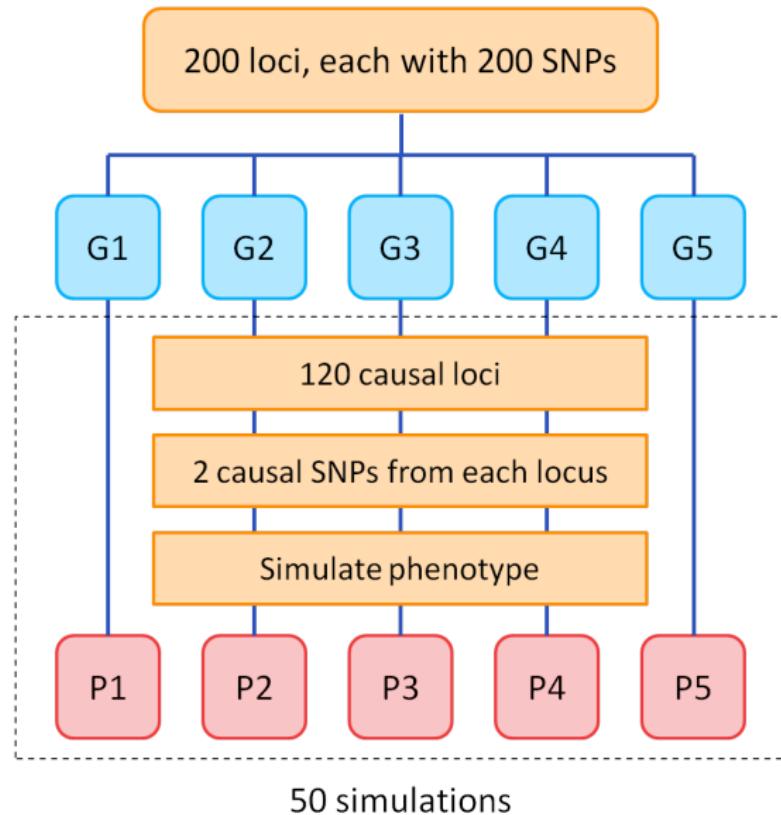
- ▶ Genotype: German Myocardial Infarction Family Study (GERMIFS)
- ▶ Five cohorts : G1, G2, G3, G4, G5
- ▶ Phenotype simulation:

$$\text{Disease liability} \quad Y_n = \sum_i v_i x_{ni} + \varepsilon_n$$

$$\text{Var}\left(\sum_i v_i x_{ni}\right) = h_g^2 = 0.4 \text{ and } \text{Var}(\varepsilon) = 0.6$$

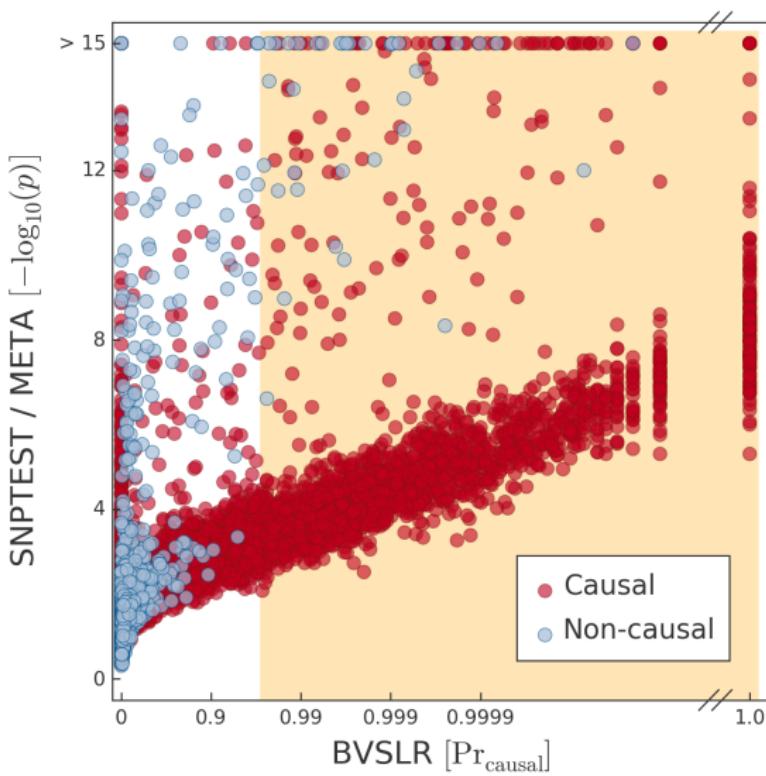
Cases are sampled from Y_n exceeding the threshold of normal distribution truncating the proportion of k (disease prevalence)

Simulation details



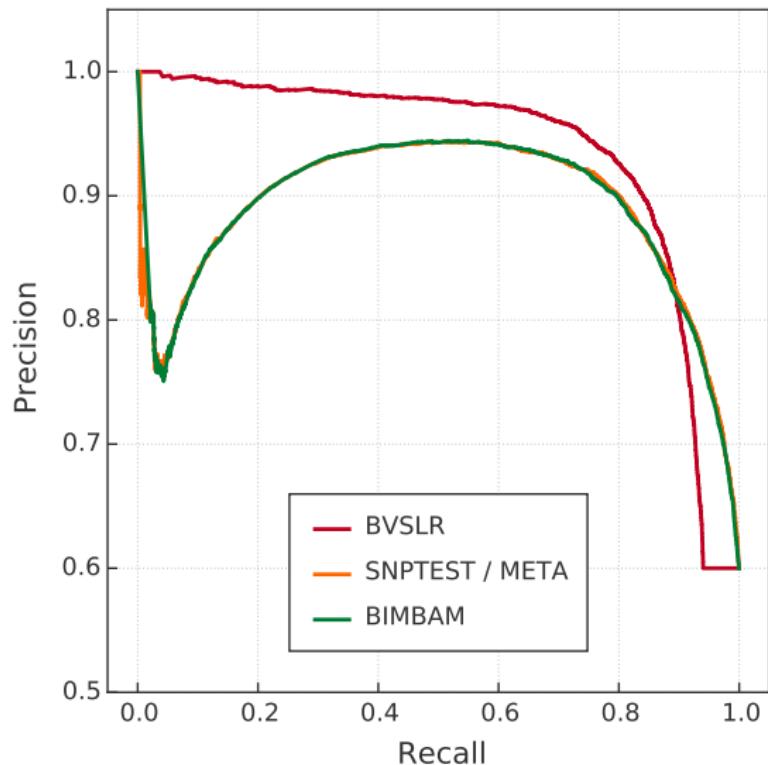
- ▶ BVSLR
- ▶ BIMBAM
- ▶ SNPTEST / META
- ▶ PAINTOR

Prediction of causal loci



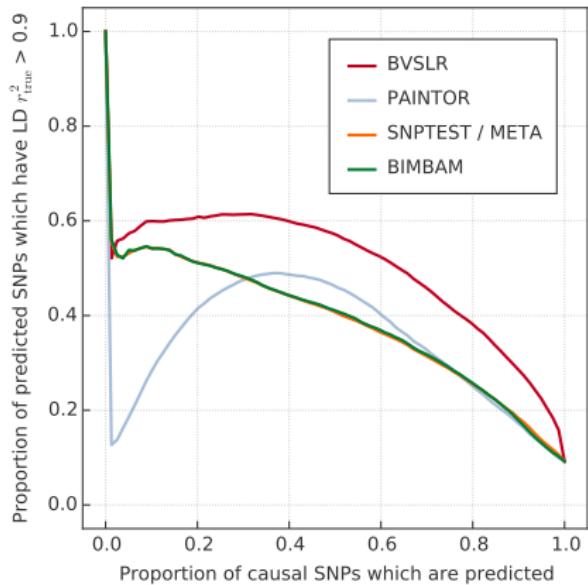
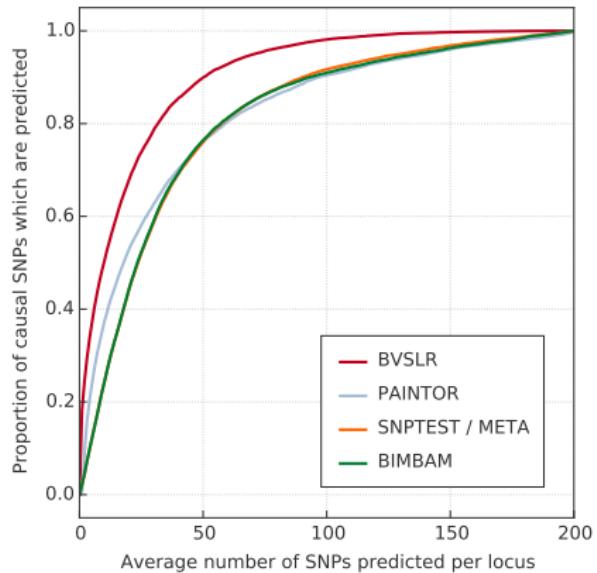
- ▶ 6000 causal loci (120 from each of the 50 simulations)
- ▶ 4000 non-causal loci (80 from each of the 50 simulations)

Prediction of causal loci



$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$
$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

Finemapping causal variants



Challenges of BVSLR

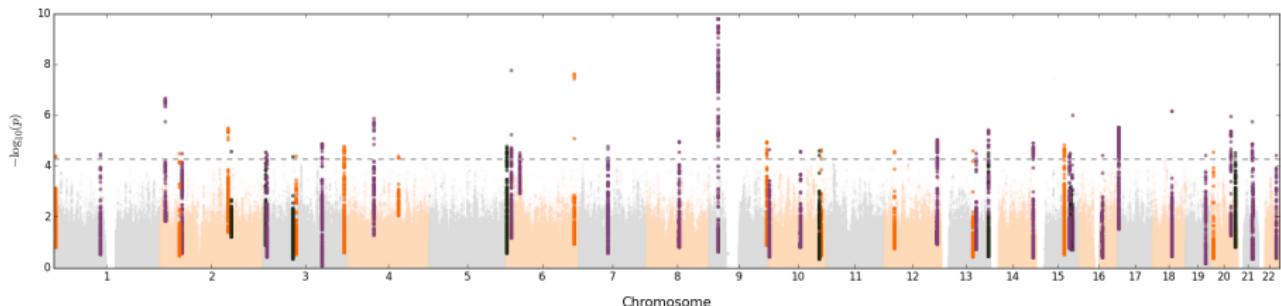
- ▶ Computational time (summing over z-states)
- ▶ New summary statistics
- ▶ Prior assumptions

Association with coronary artery diseases (CAD)

- ▶ 5 GERMIFS cohorts
- ▶ 6228 cases, 6854 controls
- ▶ Imputed with 1000G Phase 1

Association with coronary artery diseases (CAD)

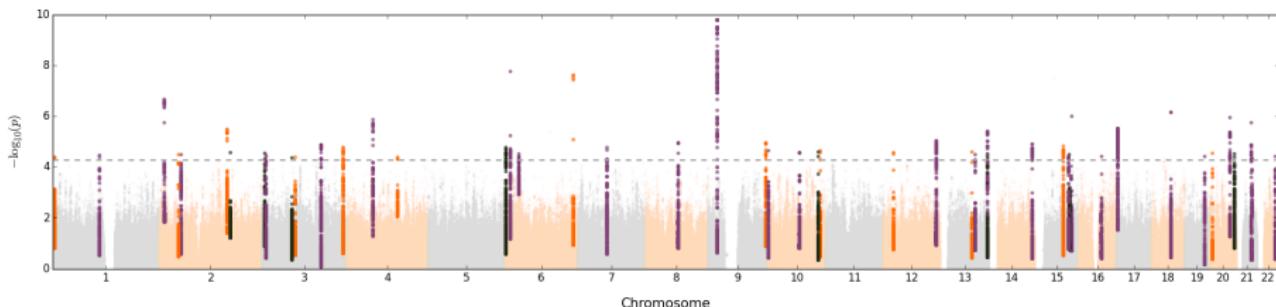
- ▶ 5 GERMIFS cohorts
- ▶ 6228 cases, 6854 controls
- ▶ Imputed with 1000G Phase 1



GWAS using SNPTEST / META

Association with coronary artery diseases (CAD)

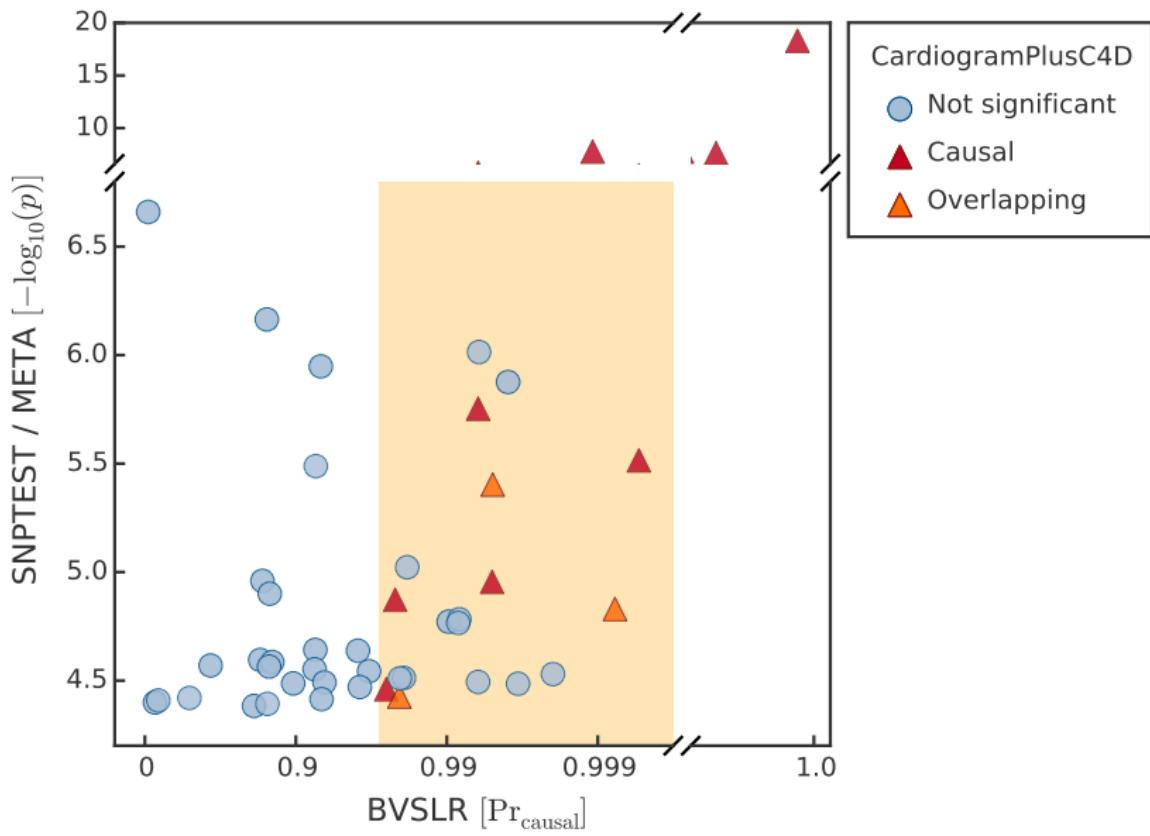
- ▶ 5 GERMIFS cohorts
- ▶ 6228 cases, 6854 controls
- ▶ Imputed with 1000G Phase 1



GWAS using SNPTEST / META

- ▶ Applied BVSLR on these 45 loci, selecting 400 SNPs at each locus.

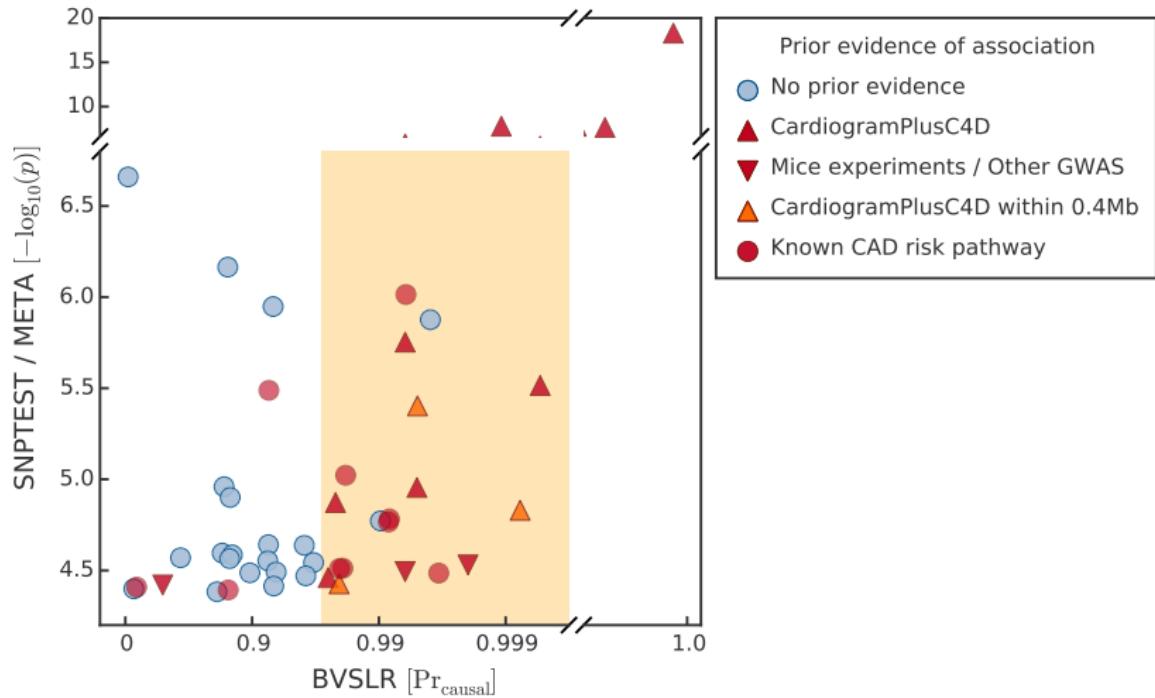
BVSLR predictions



Top BVSLR predicted loci

Region	<i>Pr</i>	Gene	Comments
6p21.3	0.998	C6orf10-BTNL2	GWAS for CAD in Han Chinese, 2012
15q25	0.997	IL-16	GWAS for CAD in Han Chinese, 2012
4q13.1	0.996	desert	-
15q25	0.994	AKAP13	C. hypertrophy (mice) / GWAS (BP in Koreans, 2011)
2p16	0.994	NRXN1	GWAS for CAD in OHGS1 + WTCCC2
3q28	0.992	IL1RAP	Involved in risk pathway
6p25	0.992	SERPINB	patented as biomarker for CVD
7q11.22	0.990	AUTS2	-
12q24	0.982	ZNF664	GWAS hit for HDL-C, TG
13q21.1	0.981	ARHGEF1	Controls vascular tone and BP

Literature-based classification



- ▶ Novel Bayesian method for GWAS
- ▶ Multivariate analysis in meta studies
- ▶ Precision
- ▶ Predicts new associations in CAD

Many thanks to ...



Johannes Söding



Heribert Schunkert



Jeanette Erdmann

Many thanks to ...



Johannes Söding



Heribert Schunkert



Jeanette Erdmann



AG Söding

Many thanks to ...



Johannes Söding



Heribert Schunkert



Jeanette Erdmann



AG Söding

Thank you!