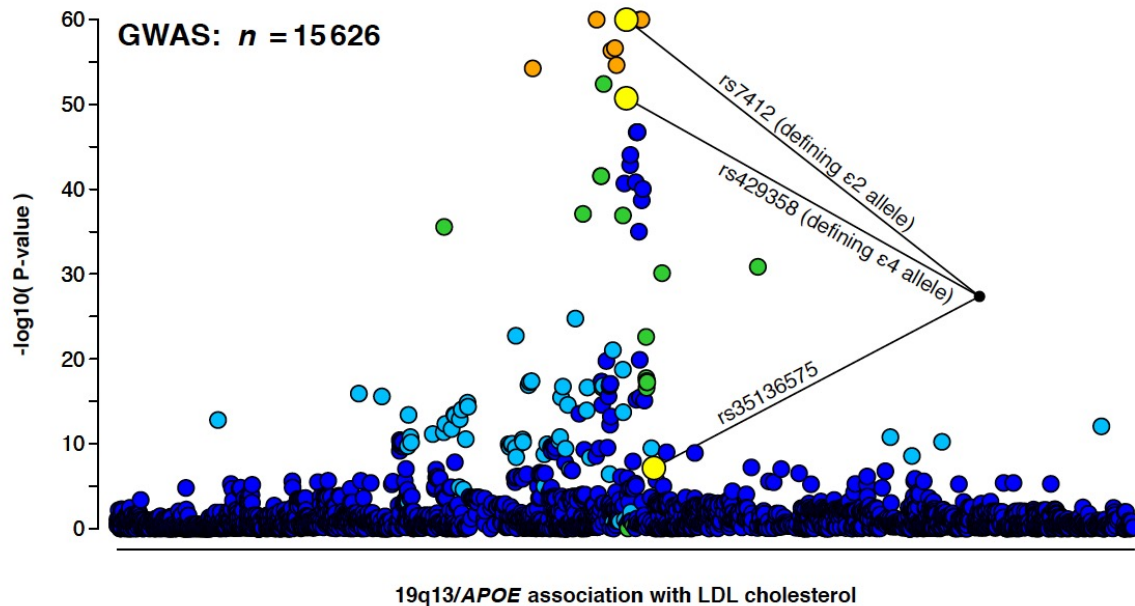


Example: Variable selection in genetics



➤ Variants near each other are often highly correlated

■ $|r| > 0.90$ very common

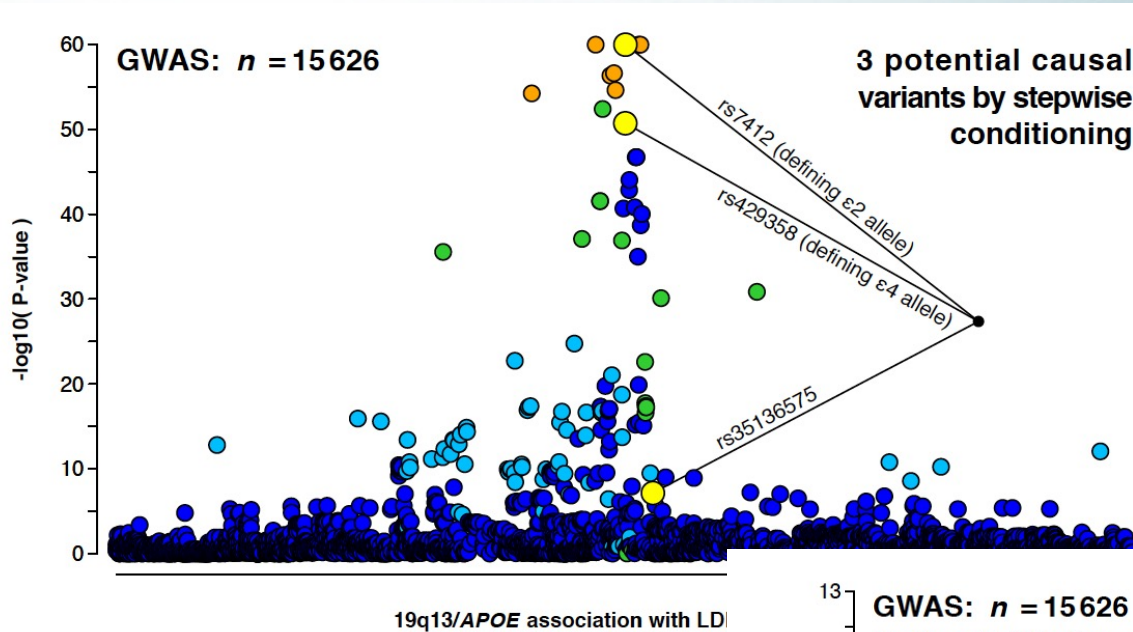
➤ Which ones are causal and which are just passengers?

$n \sim 10^5$
(samples = individuals)

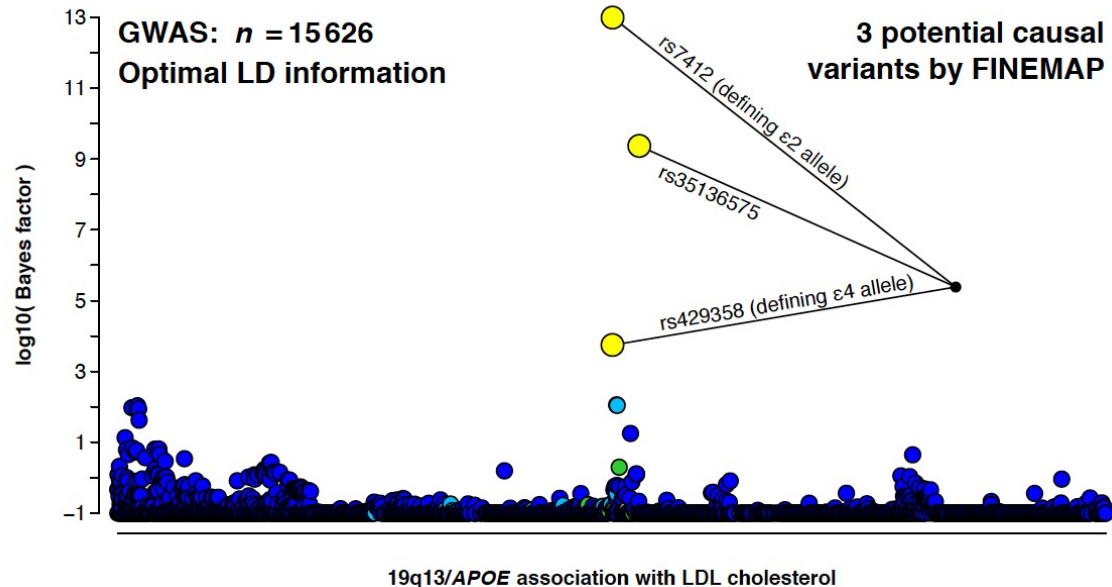
IND1	M	A/A	A/C	A/C	C/C	A/A	A/C	A/C	C/C	C/C
IND2	M	A/A	A/C	A/A	C/C	A/C	A/C	A/C	C/C	A/C
IND3	M	A/C	C/C	A/C	C/C	A/A	C/C	C/C	C/C	C/C
IND4	F	A/A	A/C	A/A	C/C	A/C	A/C	A/C	C/C	C/C
IND5	F	A/A	C/C	A/C	C/C	A/A	A/A	C/C	A/A	A/C
IND6	M	A/C	A/C	C/C	C/C	A/A	A/C	A/C	C/C	C/C
IND7	F	A/A	A/A	A/C	C/C	A/A	C/C	A/C	C/C	A/C
IND8	M	A/C	A/C	C/C	C/C	C/C	A/C	C/C	C/C	C/C
IND9	F	A/A	A/A	A/C	C/C	A/A	A/C	C/C	A/C	C/C

$p = 10^3 \dots 10^4$ (predictors = variants)

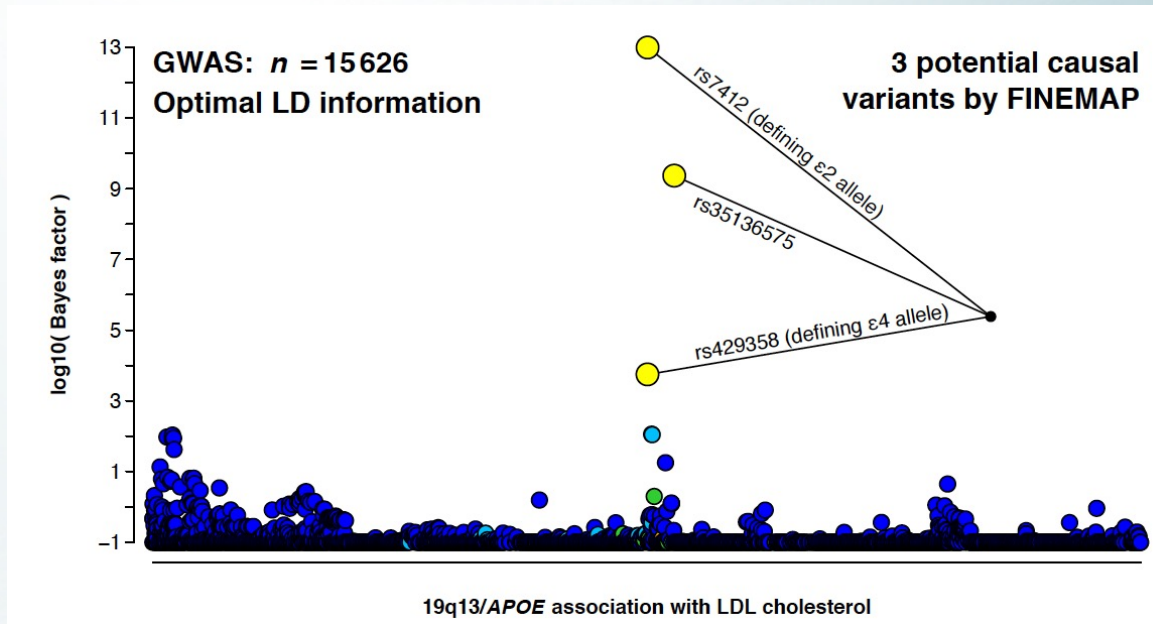
Variable selection = “fine-mapping”



- Variants near each other are often highly correlated
- Which ones are causal and which are just passengers?



Final results are probabilities



- › What is the probability for each configuration of variants being causal?
- › What is the probability for each variant being causal?

rank	config	config_prob	config_log10bf
1	rs15,rs47	0.59	44.6
2	rs15,rs42,rs47	0.02	44.9
3	rs15,rs34,rs47	0.01	44.7

index	snp	snp_prob	snp_log10bf
15	rs15	1.00	11.3
47	rs47	1.00	10.6
42	rs42	0.03	-0.22

Goal of probabilistic variable selection

To provide, for each predictor,

(1) a measure of being important (“causal”)

- We want variable selection **with estimates of uncertainty** that elastic net does not give (by default)

(2) by accounting for correlation between the predictors

- Elastic net gives only an optimum but does not output other possible solutions that could include some of the highly correlated predictors instead of the chosen ones. We want a longer list of most probable configurations.

Three pieces of efficient variable selection



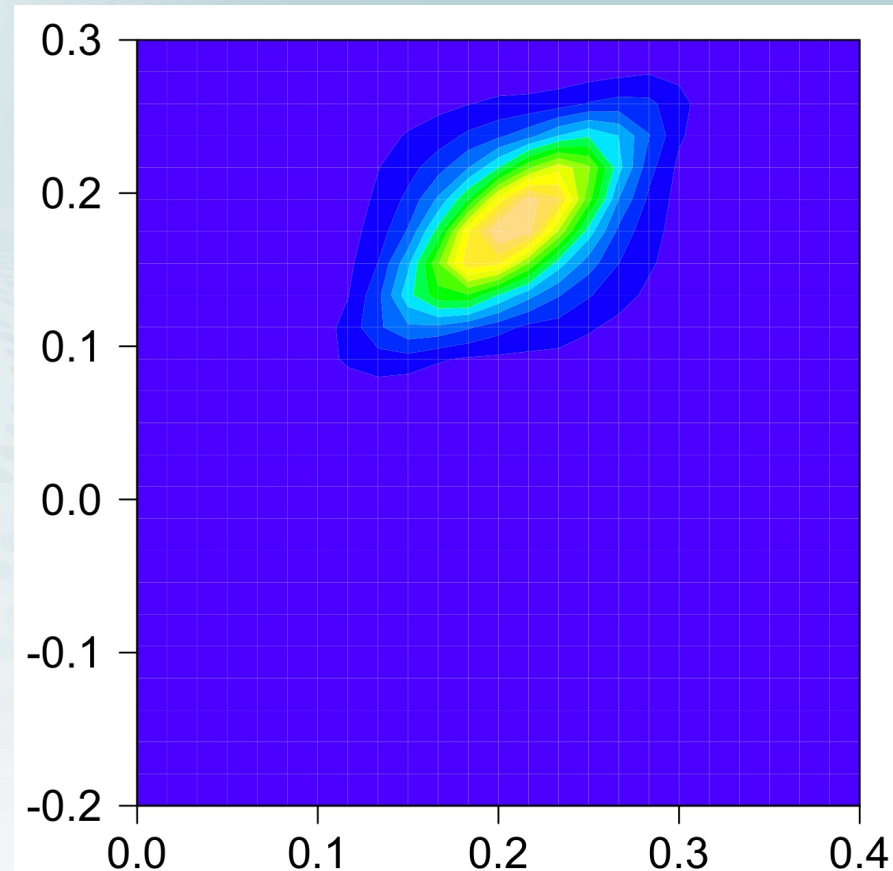
Univariate “betas”

$$\mathbf{Y}_{n \times 1} = \mathbf{X}_l \beta_l + \boldsymbol{\varepsilon}_{n \times 1}$$

$n \times 1 \quad 1 \times 1$

$$\hat{\beta}_l = (\mathbf{X}_l^T \mathbf{X}_l)^{-1} \mathbf{X}_l^T \mathbf{y}$$

These are simple regression coeffs between one predictor and outcome. Computed efficiently as in Exercise 1.5.



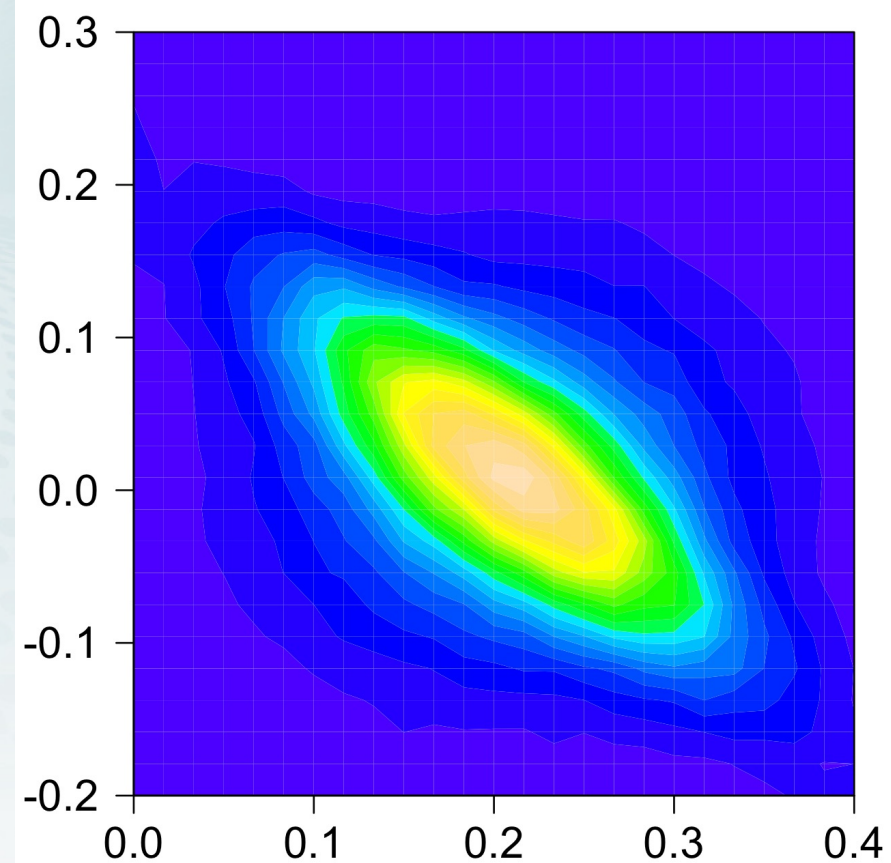
- › $\text{Cor}(X_1, X_2) = 0.85$
- › X_1 has effect 0.2, X_2 is null (=0)

Multiple regression “lambdas”

$$\mathbf{Y}_{n \times 1} = \mathbf{X}_{n \times p} \boldsymbol{\lambda}_{p \times 1} + \boldsymbol{\varepsilon}_{n \times 1}$$

$$\hat{\boldsymbol{\lambda}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

These are the direct effects and account for other predictors but are complicated to compute for large p , especially when $p > n$.



- › $\text{Cor}(X_1, X_2) = 0.85$
- › X_1 has effect 0.2, X_2 is null (=0)

Betas and lambdas

Assuming standardized predictors

$$\beta = \frac{1}{n} (\mathbf{X}^T \mathbf{X}) \lambda = \mathbf{R} \lambda$$

where \mathbf{R} is pairwise correlation matrix of predictors.

$$\beta = \begin{bmatrix} 1 & 0.85 \\ 0.85 & 1 \end{bmatrix} \begin{bmatrix} 0.2 \\ 0 \end{bmatrix} = \begin{bmatrix} 0.2 \\ 0.17 \end{bmatrix}$$

Summary data

- › Computation for multiple regression model is possible using **summary data**: univariate z-scores and correlation matrix of predictors (R matrix)

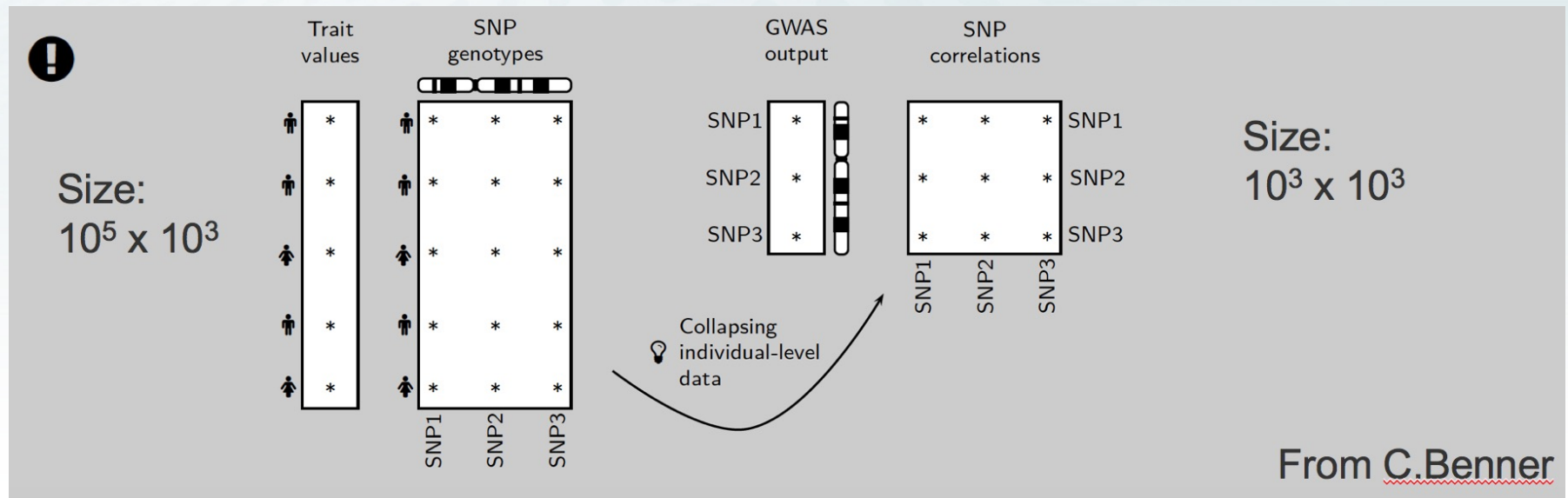
$$\underset{n \times 1}{\mathbf{Y}} = \underset{n \times p}{\mathbf{X}} \underset{p \times 1}{\boldsymbol{\lambda}} + \underset{n \times 1}{\boldsymbol{\varepsilon}}$$

$$\hat{\boldsymbol{\lambda}} = \underbrace{(\mathbf{X}^T \mathbf{X})^{-1}}_{n\mathbf{R}} \underbrace{\mathbf{X}^T \mathbf{Y}}_{\sqrt{n}\sigma_{\varepsilon}\mathbf{Z}}$$

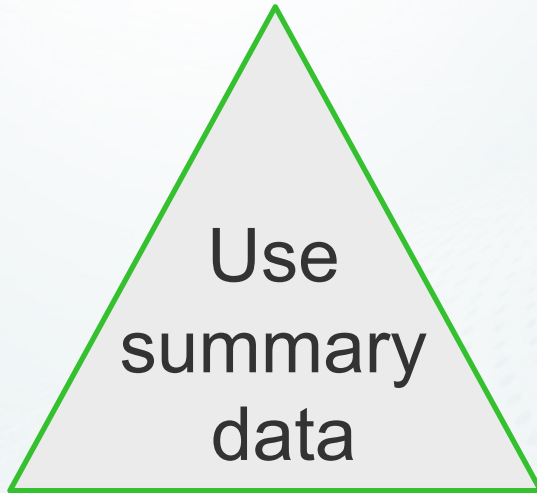
\mathbf{R} , predictors' correlations
 \mathbf{Z} , univariate z-scores

Use summary data to make multiple regression possible from univariate results

- › Working with less data but with full information
 - For $p=1,000$ and $n=100,000$, data reduction is **100 fold**



Three pieces of efficient variable selection



Assumption: true configuration is sparse

- › Joint MLE (or ridge regression) of all predictors is not our final answer to variable selection since it does not lead to sparse solutions
- › Bayesian answer:
 - Define a prior probabilities for configurations
 - Define a prior distribution for regression coefficients of a configuration
 - Integrate (prior x likelihood) leading to *marginal likelihood* for the configuration

0	1	0	1	0	0	0	1	0	0	Causal configuration γ
0	2.1	0	0.1	0	0	0	3.1	0	0	Causal effects λ
1.3	2.0	0.7	0.2	1.5	0.3	0.2	3.2	2.9	0.1	MLE $\hat{\lambda}$

“causal effects” = “direct effects” = “multiple regression coefficients”

Bayesian model for variable selection

- › Define a configuration γ as a binary vector over predictors

$$\gamma = \begin{array}{|c|c|c|c|c|c|c|c|} \hline 0 & 0 & 1 & 0 & 0 & 0 & 1 & 0 \\ \hline \end{array}$$

- › This configuration represents model where predictors 3 and 7 are allowed to have non-zero effects while the other predictors have effect size 0

Bayesian model for variable selection

- › Define a configuration γ as a binary vector for predictors

$$\gamma = \begin{bmatrix} 0 & 0 & 1 & 0 & 0 & 0 & 1 & 0 \end{bmatrix}$$

- › In total there are 2^p configurations on p predictors, but we will assume that only **sparse** configurations are plausible, say those with < 10 non-zero predictors
 - This is similar idea to LASSO that sets many coefficients to 0
- › Ultimate goal is to compute probability for each configuration, given the observed data
 - This is much more challenging than the LASSO optimization

Bayesian model for variable selection

- › Define a configuration γ as a binary vector for predictors
- › Each non-zero predictor picks its effect from $N(0, s^2)$
 - This is the slab part of the spike and slab prior
 - This is similar prior as in ridge regression but now the model is sparse which is different from ridge regression

$$p(\lambda|\gamma) = \mathcal{N}(\mathbf{0}, s^2 \mathbf{\Delta}_\gamma)$$

Causal configuration γ

1	0	0	0	0	0	0	0	0	0
---	---	---	---	---	---	---	---	---	---

$$\mathbf{\Delta}_\gamma = \text{diag}(\gamma) = \begin{bmatrix} 1 & & & & & & & & & \\ & 0 & & & & & & & & \\ & & \ddots & & & & & & & \\ & & & \ddots & & & & & & \\ & & & & 0 & & & & & \end{bmatrix}$$

Bayesian model for variable selection

- › Define causal configuration γ as a binary vector for predictors
- › Each non-zero predictor picks its effect from $N(0, s^2)$
- › For each configuration, compute the Bayes factor (BF), i.e., how well the configuration explains the data relative to the null model

$$\text{BF}_{\gamma} = \frac{P(\text{DATA}|\gamma)}{P(\text{DATA}|\text{NULL})}$$

- › How to compute the numerator?

Marginal likelihood for a configuration

$$\begin{aligned}\mathcal{L}(\gamma) &= \int p(\mathbf{y}|\boldsymbol{\lambda}, \mathbf{X}) p(\boldsymbol{\lambda}|\gamma) d\boldsymbol{\lambda} && \text{(Likelihood x prior)} \\ &= \int \mathcal{N}(\hat{\boldsymbol{\lambda}}|\boldsymbol{\lambda}, \sigma^2(\mathbf{X}^T \mathbf{X})^{-1}) \mathcal{N}(\boldsymbol{\lambda}|0, s_\lambda^2 \sigma^2 \boldsymbol{\Delta}_\gamma) d\boldsymbol{\lambda} \\ &= \mathcal{N}(\hat{\boldsymbol{\lambda}}|0, \sigma^2(n\mathbf{R})^{-1} + s_\lambda^2 \sigma^2 \boldsymbol{\Delta}_\gamma) \\ &= \mathcal{N}(\hat{\mathbf{z}}|0, \mathbf{R} + \mathbf{R}\boldsymbol{\Delta}_\gamma^* \mathbf{R}) && \text{Depends on X and Y only through} \\ &&& \text{Summary data Z and R!}\end{aligned}$$

$$\hat{\mathbf{z}} = \hat{\boldsymbol{\beta}}/\text{SE}_\beta = \frac{\sqrt{n}}{\sigma} \hat{\boldsymbol{\beta}} \text{ and } \boldsymbol{\Delta}_\gamma^* = s_\lambda^2 \boldsymbol{\Delta}_\gamma$$

Marginal likelihood for a configuration

$$\mathcal{L}(\gamma) = \mathcal{N}(\hat{\mathbf{z}}|0, \mathbf{R} + \mathbf{R}\mathbf{\Delta}_{\gamma}^*\mathbf{R})$$

- ✦ Depends on data only through univariate summary statistics and correlation matrix \mathbf{R} , i.e.,
summary statistics
Thus, we do not need access to original \mathbf{X} and \mathbf{Y} !
- Dimension is p , the number of predictors that can be 10,000s, which makes evaluation of many configurations impossible since each config requires decomposition of a $p \times p$ matrix and this is $O(p^3)$

Using only causal predictors

Consider configuration γ
Divide predictors into causal
(C) and non-causal (N)

$$\mathbf{z} = \begin{bmatrix} \mathbf{z}_C \\ \mathbf{z}_N \end{bmatrix} \quad \mathbf{R} = \begin{bmatrix} \mathbf{R}_{CC} & \mathbf{R}_{CN} \\ \mathbf{R}_{NC} & \mathbf{R}_{NN} \end{bmatrix}$$

Using only causal variants

Consider configuration γ
Divide predictors into causal
(C) and non-causal (N)

$$\mathbf{z} = \begin{bmatrix} \mathbf{z}_C \\ \mathbf{z}_N \end{bmatrix} \quad \mathbf{R} = \begin{bmatrix} \mathbf{R}_{CC} & \mathbf{R}_{CN} \\ \mathbf{R}_{NC} & \mathbf{R}_{NN} \end{bmatrix}$$

Cond. distr of $\mathbf{z}_N | \mathbf{z}_C$ is the same for configuration γ as it is for null model !

$$\begin{aligned} \text{BF}(\gamma : \text{NULL}) &= \frac{\mathcal{N}(\mathbf{z}|0, \mathbf{R} + \mathbf{R} \mathbf{\Delta}_\gamma^* \mathbf{R})}{\mathcal{N}(\mathbf{z}|0, \mathbf{R})} \\ &= \frac{\mathcal{N}(\mathbf{z}_C|0, \mathbf{R}_{CC} + \mathbf{R}_{CC} \mathbf{\Delta}_{CC}^* \mathbf{R}_{CC}) \mathcal{N}(\mathbf{z}_N|\mathbf{z}_C)}{\mathcal{N}(\mathbf{z}_C|0, \mathbf{R}_{CC}) \mathcal{N}(\mathbf{z}_N|\mathbf{z}_C)} \\ &= \frac{\mathcal{N}(\mathbf{z}_C|0, \mathbf{R}_{CC} + \mathbf{R}_{CC} \mathbf{\Delta}_{CC}^* \mathbf{R}_{CC})}{\mathcal{N}(\mathbf{z}_C|0, \mathbf{R}_{CC})} \end{aligned}$$

Benner et al. 2016

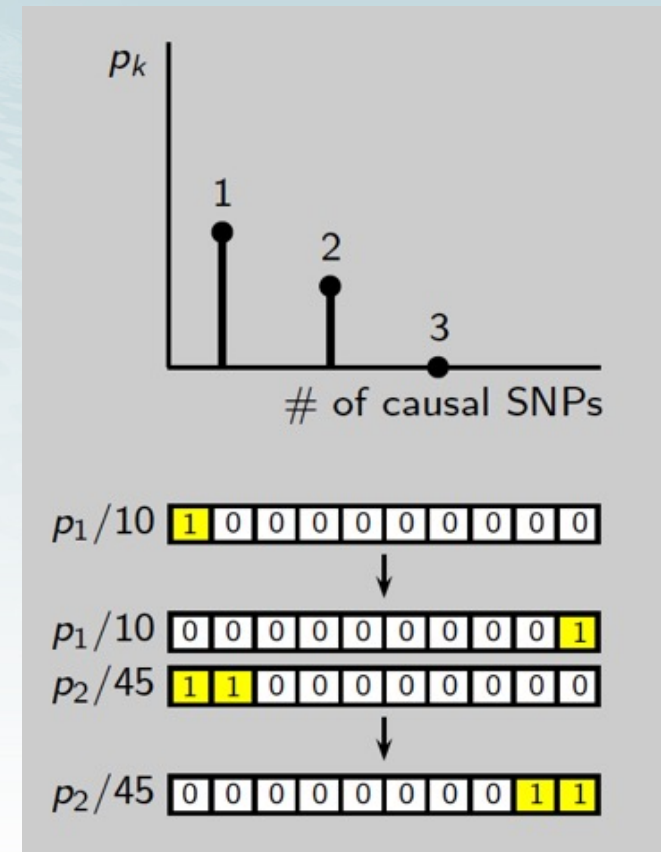
Bayesian model for variable selection

- › Define causal configuration γ as a binary vector for predictors
- › Each non-zero predictor picks its effect from $N(0, s^2)$
- › For each configuration compute the Bayes factor (BF), i.e., how well the configuration explains the data relative to the null model
- › By combining BFs with prior probabilities of configurations we get the posterior probabilities

$$p_{\gamma} = P(\gamma|\text{DATA}) \propto \text{prior}_{\gamma} \times \text{BF}_{\gamma}$$

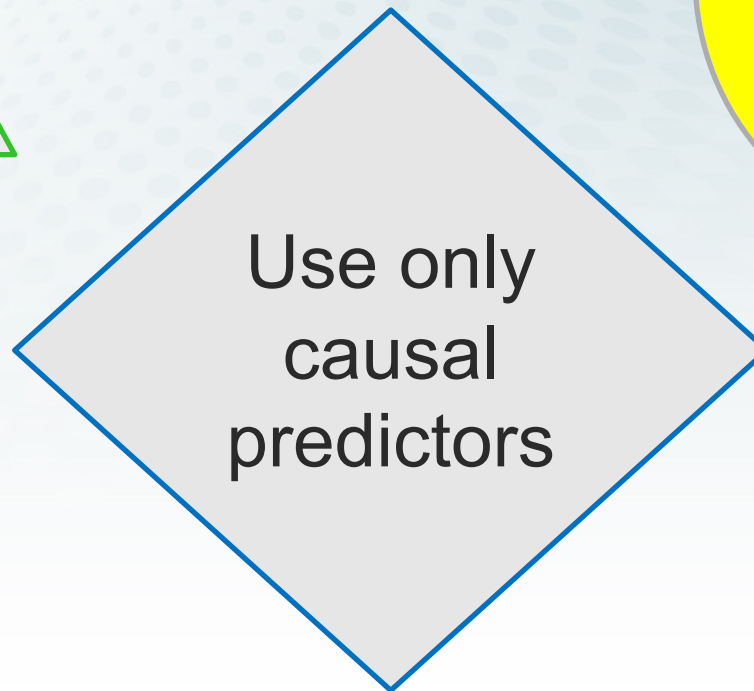
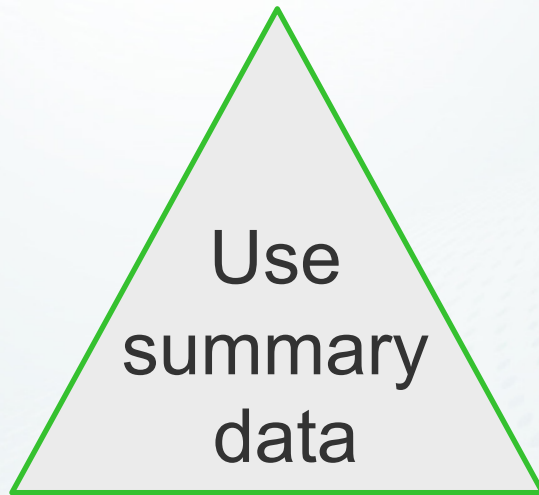
Prior on configurations to enforce sparsity

- › Specify probability p_k that there are k non-zero predictors
- › Divide that probability equally between all configurations having k non-zero predictors
- › This prior could be learned from data or remain as an ad-hoc choice



C. Benner

Three pieces of efficient fine-mapping



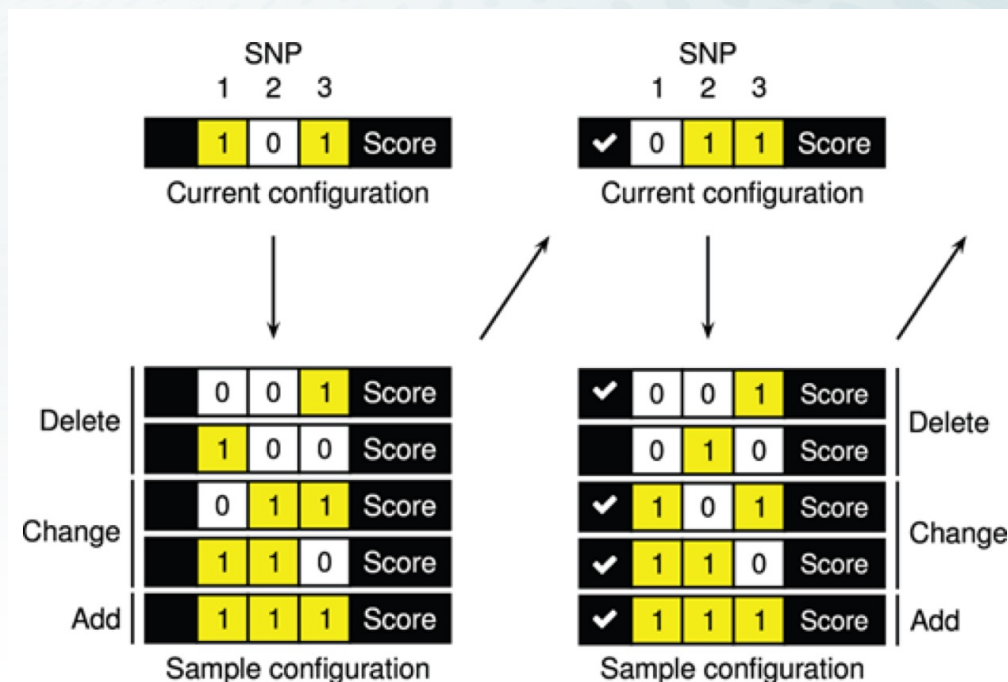
Space of causal configurations is huge, 2^p

› Best subset algorithms evaluate every configuration

- Can allow at most 3 causal predictors when 1000s of predictors are available
- Experimenting with genetic data: On average only about 100 configs out of 70,000,000 already covered 95% of posterior in setting: $p=750$, 5 causal predictors (Benner et al. 2016)
 - Can be different in some other application fields!

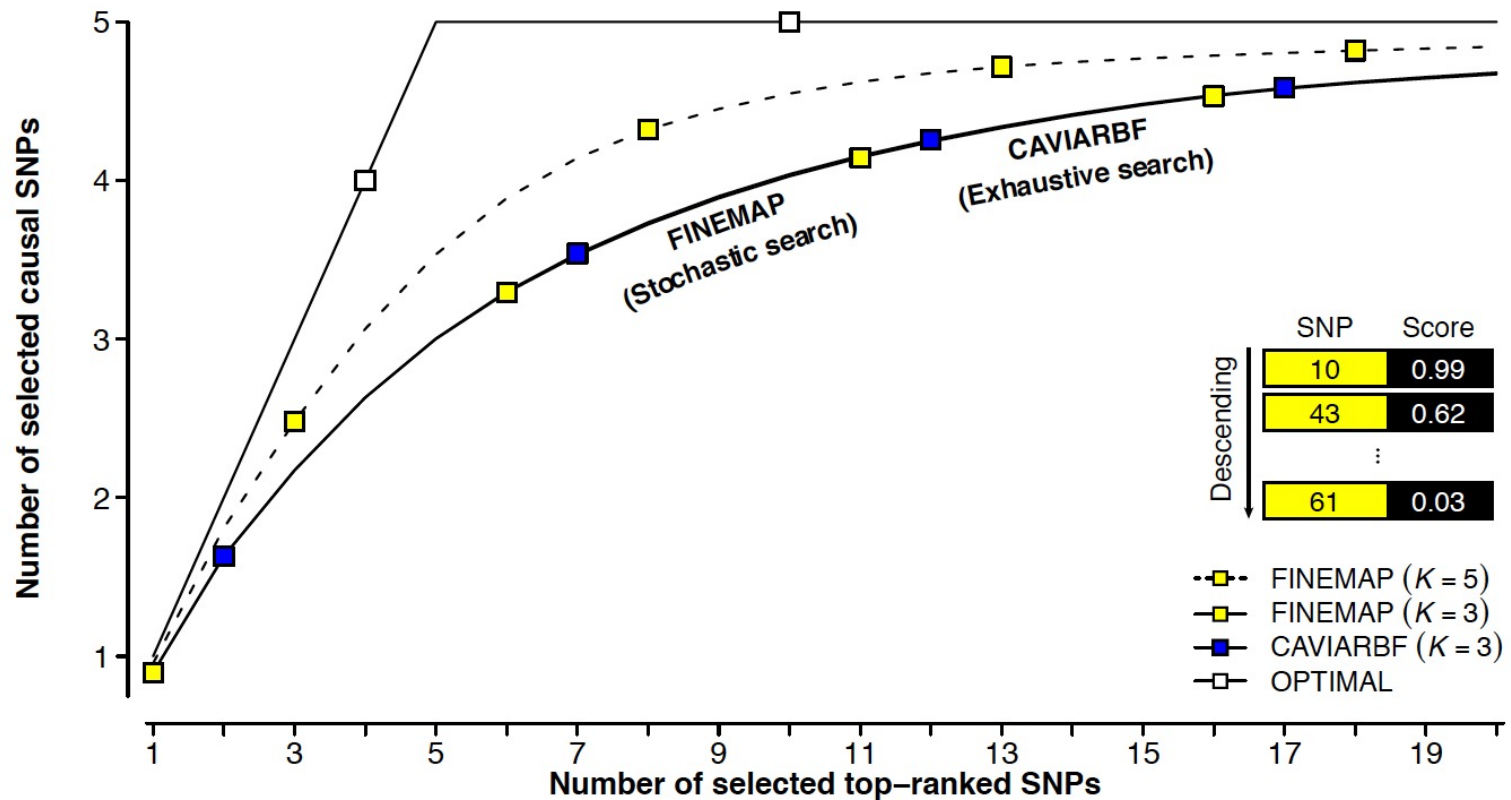
Shotgun stochastic search algorithm

- › Collect configurations from a high probability region using Shotgun stochastic search (Hans et al. 2007)
 - Memorize BFs of all those configurations seen during the search
 - Stop once not much new probability mass is found
 - Renormalize posteriors with respect to the configurations visited

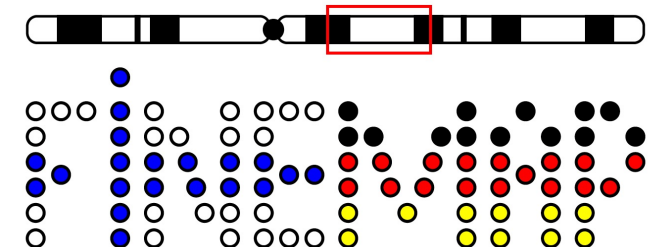
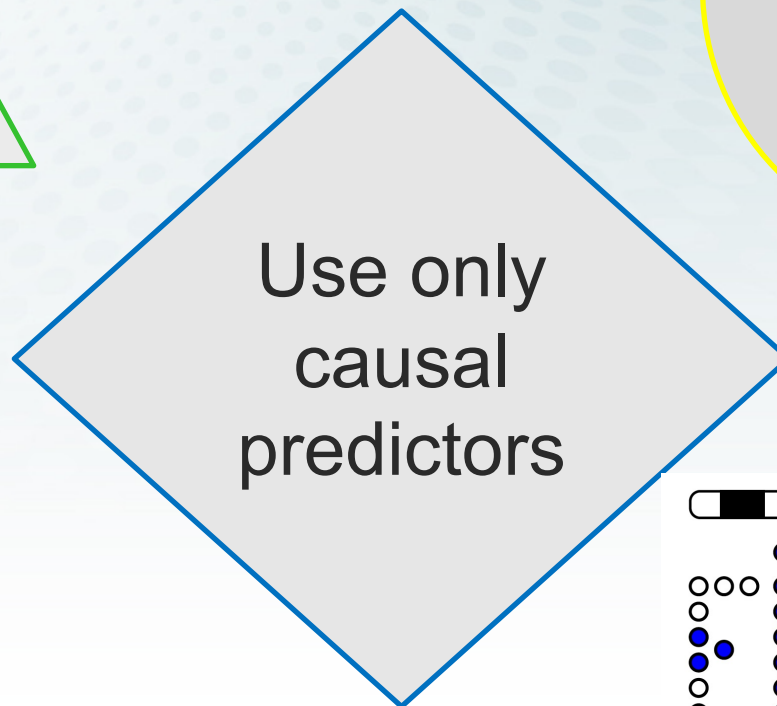
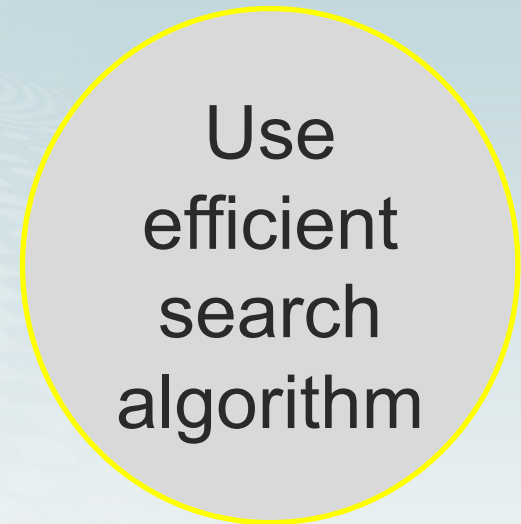
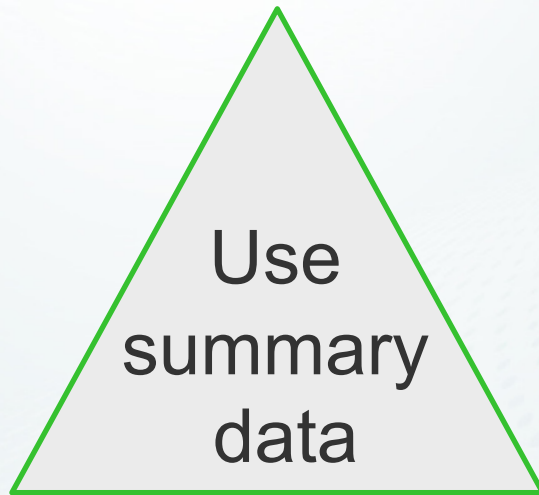


Example: FINEMAP software

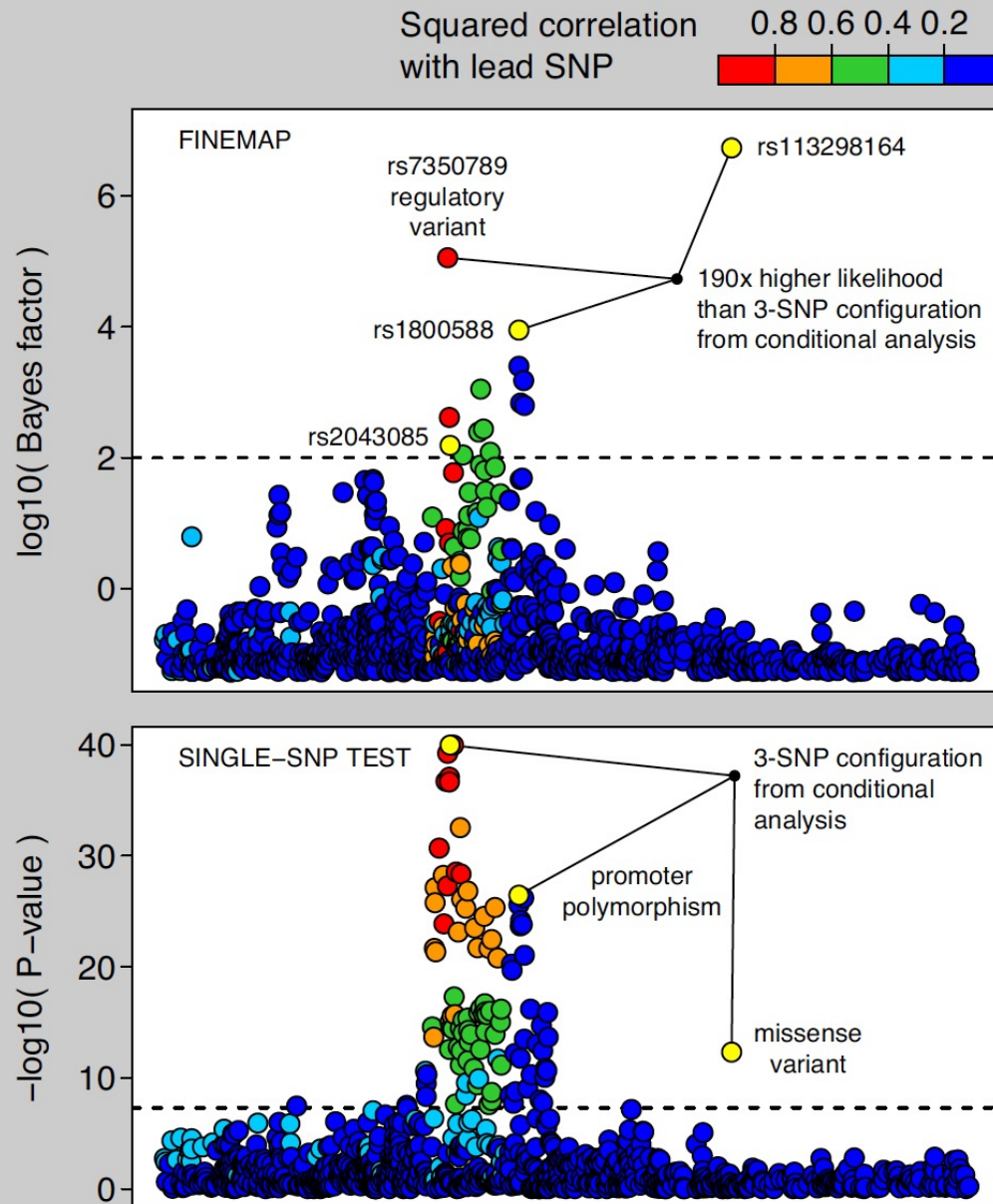
- › Simulations with $p=1500$ of which 5 are truly non-zero
 - FINEMAP runs in a few seconds
 - Enumeration is impossible in practice



Three pieces of efficient variable selection



15q21/*LIPC* association with HDL cholesterol



6 Mb region
8612 variants



**NATIONAL INSTITUTE FOR
HEALTH AND WELFARE**
FINLAND

FINRISK STUDY
20000 individuals

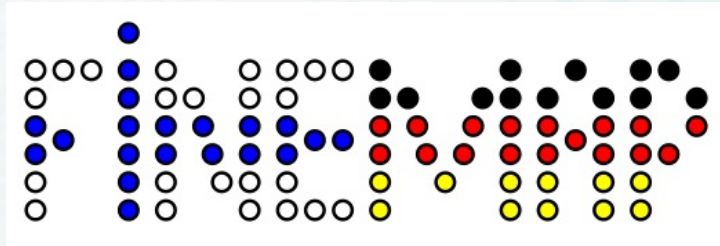
Benner et al. 2016

Surakka et al.
Nat. Genet. 2015

Acknowledgements



Christian Benner (www.finemap.me)



ACADEMY OF FINLAND



NATIONAL INSTITUTE
FOR HEALTH AND WELFARE