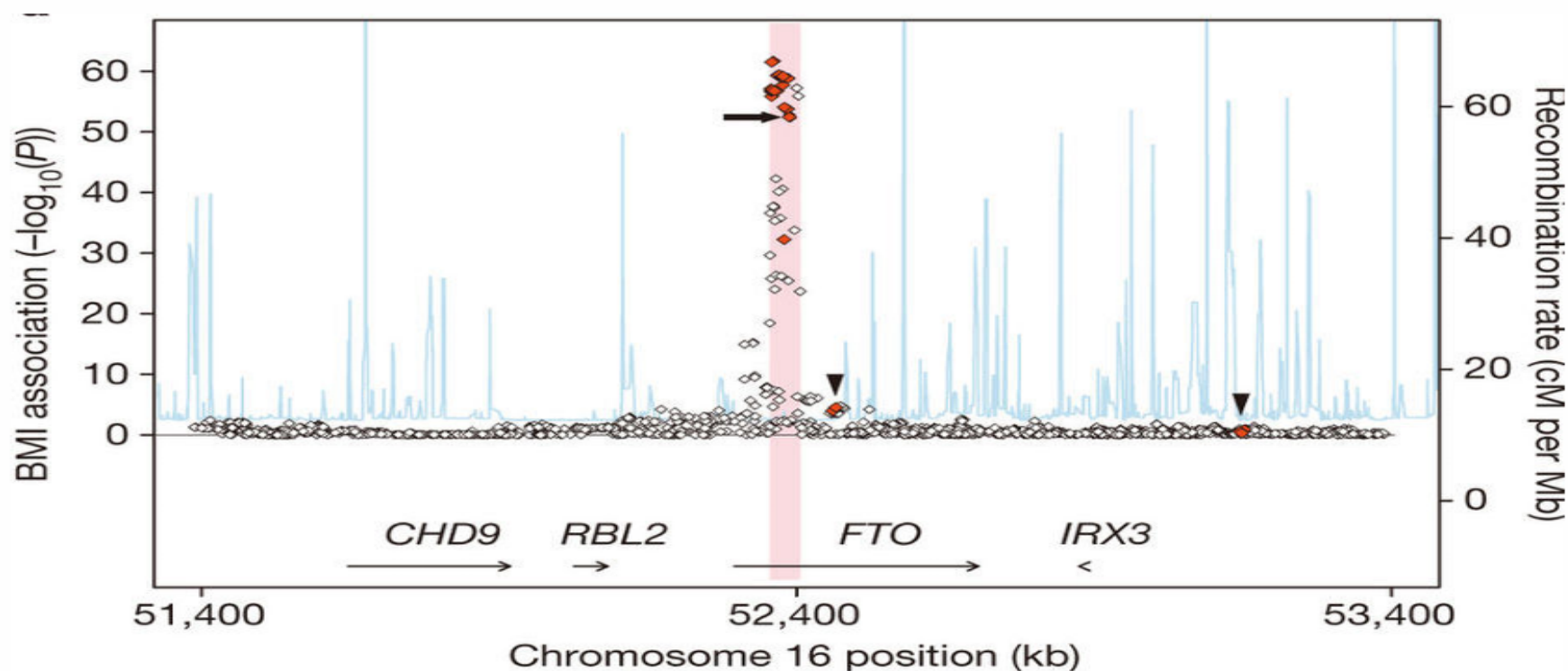


Finemapping and its application in Coloc

- A Simple New Approach to Variable Selection in Regression, with Application to Genetic Fine Mapping
- A Bayesian framework for multiple trait colocalization from summary association statistics
- A more accurate method for colocalisation analysis allowing for multiple causal variants

Single causal variant can drive association signals in multiple SNPs

- Causal variant: the variant with a true effect on the phenotype.



Lead SNPs are often not causal variants

- Because of sampling errors: a SNP in close LD with the causal SNP may have similar association statistics.
- Simulations with 1000 cases and 1000 controls: at effect size 1.1 and AF 5%, causal variant has 2.4% chance of being the lead SNP.
- The vast majority of variants discovered by GWAS have small effect sizes $\log -\text{OR} < 1.1$.

Lead SNPs are often not causal variants

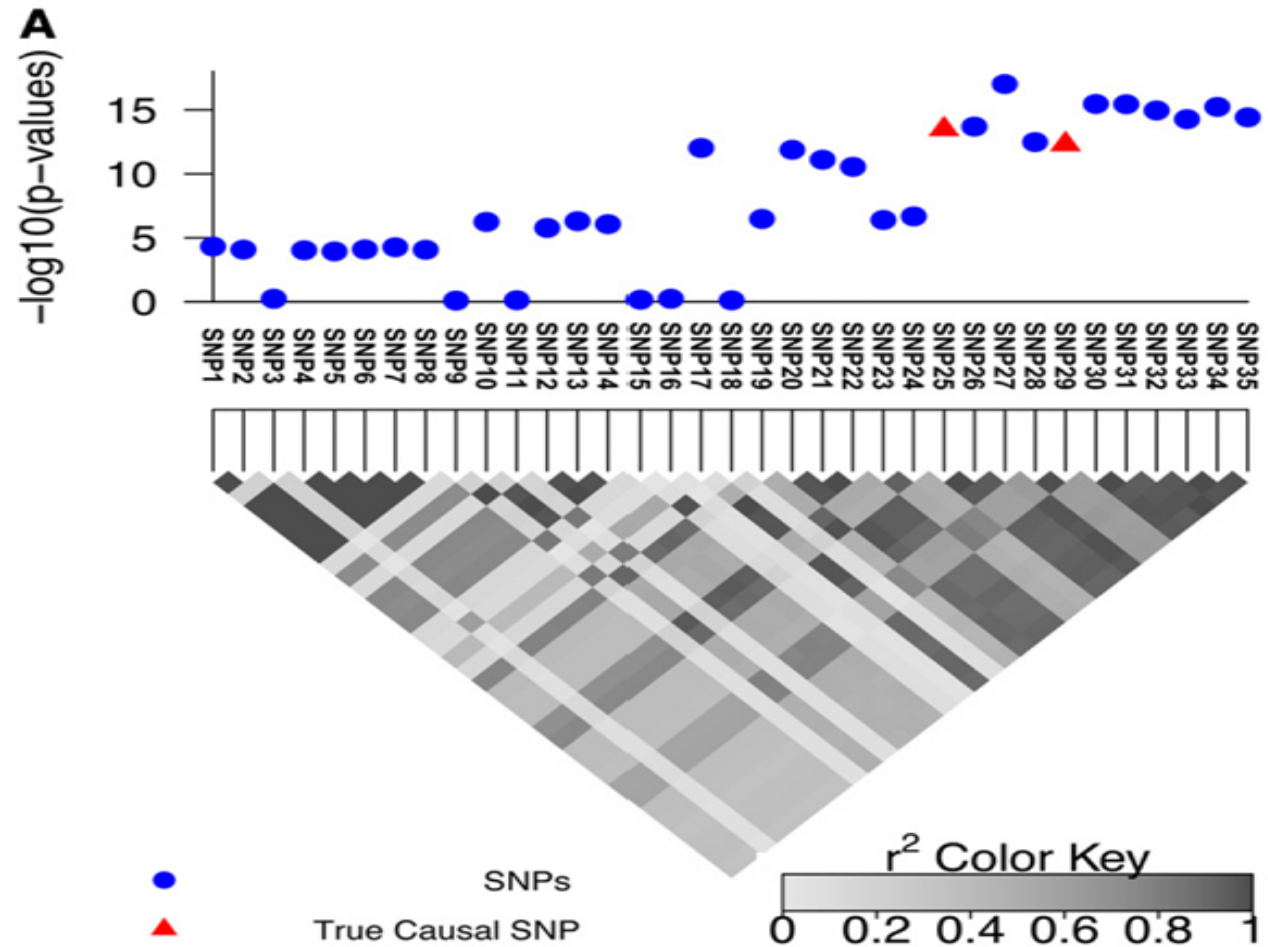


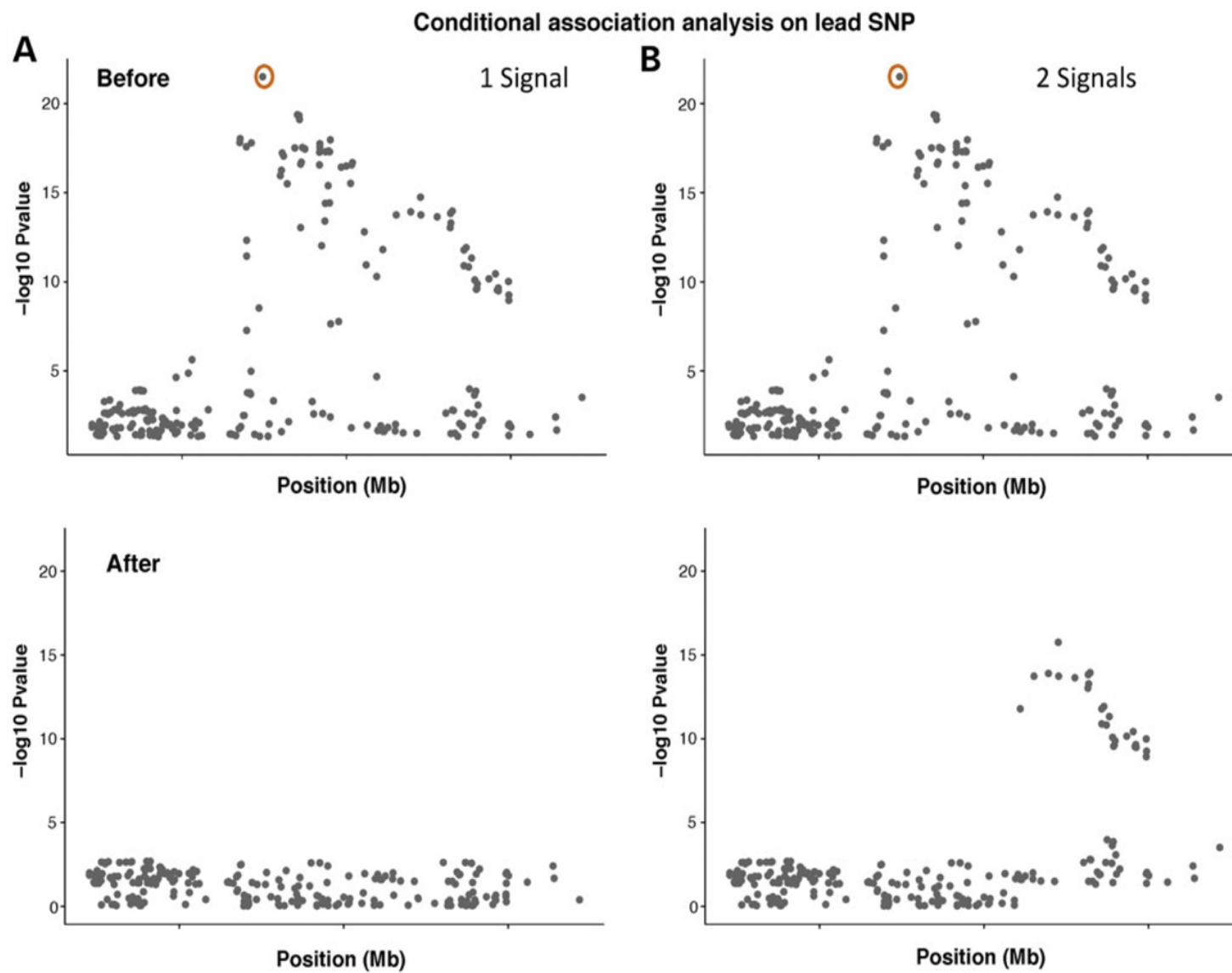
Figure: Simulated data with 2 causal variants (red) [Hormozdiari et al, Genetics, 2014]

Conditional Regression

- Regression analysis on a SNP j by conditioning on the lead SNP:

$$Y = G_{\text{lead}} \cdot \beta_{\text{lead}} + G_j \beta_j + \epsilon$$

- A SNP passing the threshold will be chosen. ($p\text{val} < 5 \cdot 10^{-8}$)
- Repeat this analysis: at each step, condition on all SNPs chosen at previous steps.



Conditional Regression

- Lead SNPs may not be causal: wrong decision at the beginning.
- It is unclear how to account for multiple testing and choose the threshold: at each step, many hypothesis are tested.
- Low power of detecting the secondary SNP.

Finemapping

- Finding a small number of variants that explain all the associations in a region.
- PIP: the posterior probability that a SNP is causal, summing over all possible configurations.

$$P(\gamma_j = 1 \mid D) = \sum_{\gamma} \underbrace{I(\gamma_j = 1)}_{\substack{1 \text{ if } \gamma_j=1; 0 \text{ otherwise}}} \cdot P(\gamma \mid D)$$

- Credible set: the minimum set of variant that contains the causal variant with probability α (typically, 95%).

Define the "confidence level" of a variant set S . If there is a single causal variant, then it is simply the sum of PIP of all the variants in S . Ex. $S = \{A, B, C\}$, its confidence level is:

$$\rho = \text{PIP}_A + \text{PIP}_B + \text{PIP}_C$$

SuSiE

$$\mathbf{y} = \mathbf{X}\mathbf{b} + \mathbf{e} \quad \mathbf{e} \sim N_n(0, \sigma^2 I_n)$$

- variables j with non-zero effects ($b_j \neq 0$) as 'effect variables'.
- Assume now that exactly two variables are effect variables-variables 1 and 4, and that these two effect variables are each completely correlated with another non-effect variable, say $x_1 = x_2$ and $x_3 = x_4$. Further suppose that no other pairs of variables are correlated. Here, because the effect variables are completely correlated with other variables, it is impossible to select the correct variables confidently, even when n is very large. However, given sufficient data it should be possible to conclude that there are (at least) two effect variables, and that

$$(b_1 \neq 0 \text{ or } b_2 \neq 0) \text{ and } (b_3 \neq 0 \text{ or } b_4 \neq 0)$$

SuSiE

'single-effect regression' (SER) model

$$\begin{aligned}\mathbf{y} &= \mathbf{X}\mathbf{b} + \mathbf{e} \\ \mathbf{e} &\sim N_n(0, \sigma^2 I_n) \\ \mathbf{b} &= b\boldsymbol{\gamma} \\ \boldsymbol{\gamma} &\sim \text{Mult}(1, \boldsymbol{\pi}) \\ b &\sim N_1(0, \sigma_0^2)\end{aligned}$$

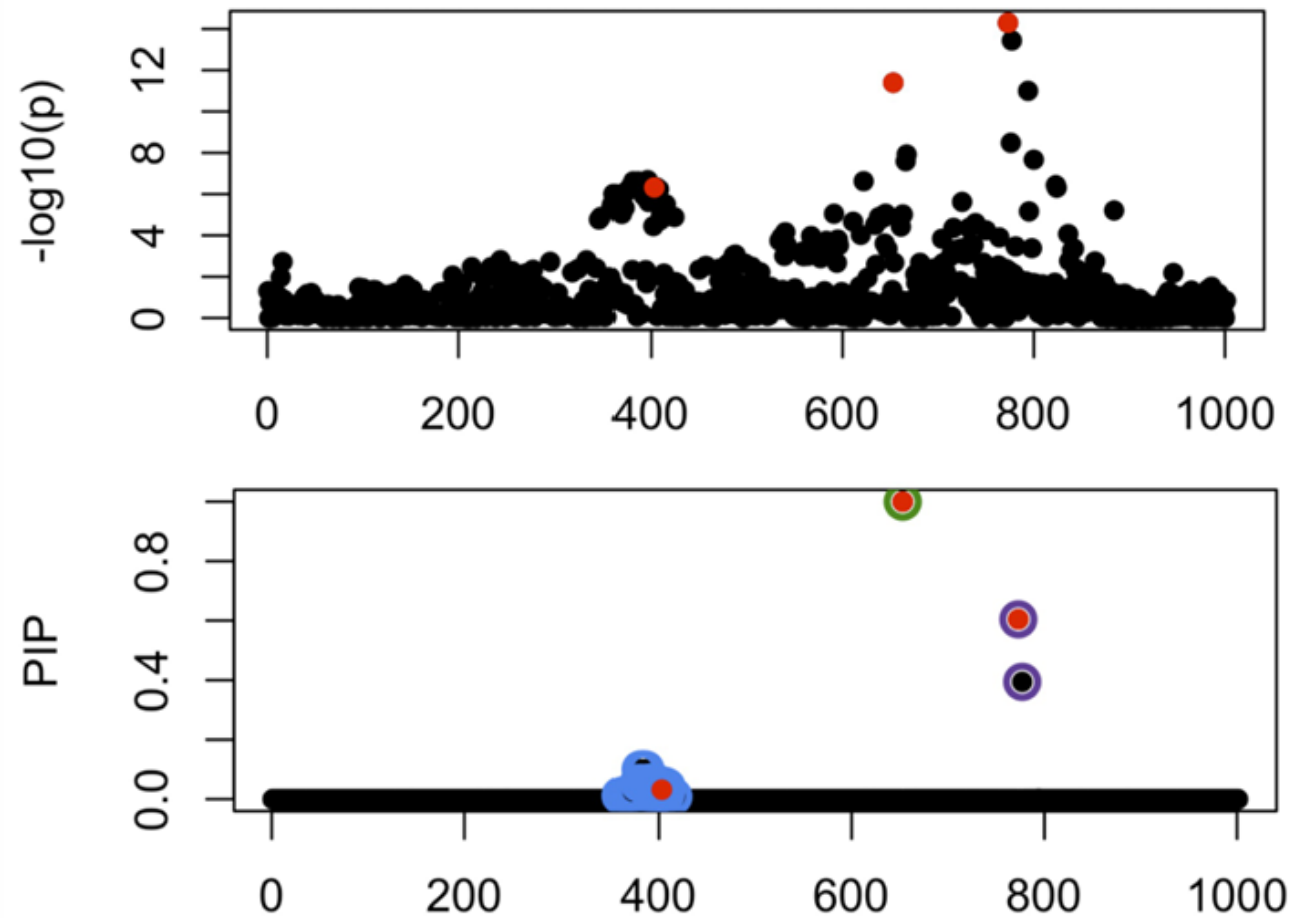
Mult($m, \boldsymbol{\pi}$) denotes the multinomial distribution on class counts that is obtained when m samples are drawn with class probabilities given by $\boldsymbol{\pi}$.

the prior variance of the non-zero effect, σ_0^2 , and the prior inclusion probabilities $\boldsymbol{\pi} = (\pi_1, \dots, \pi_p)$, in which π_j gives the prior probability that variable j is the effect variable are given

SuSiE

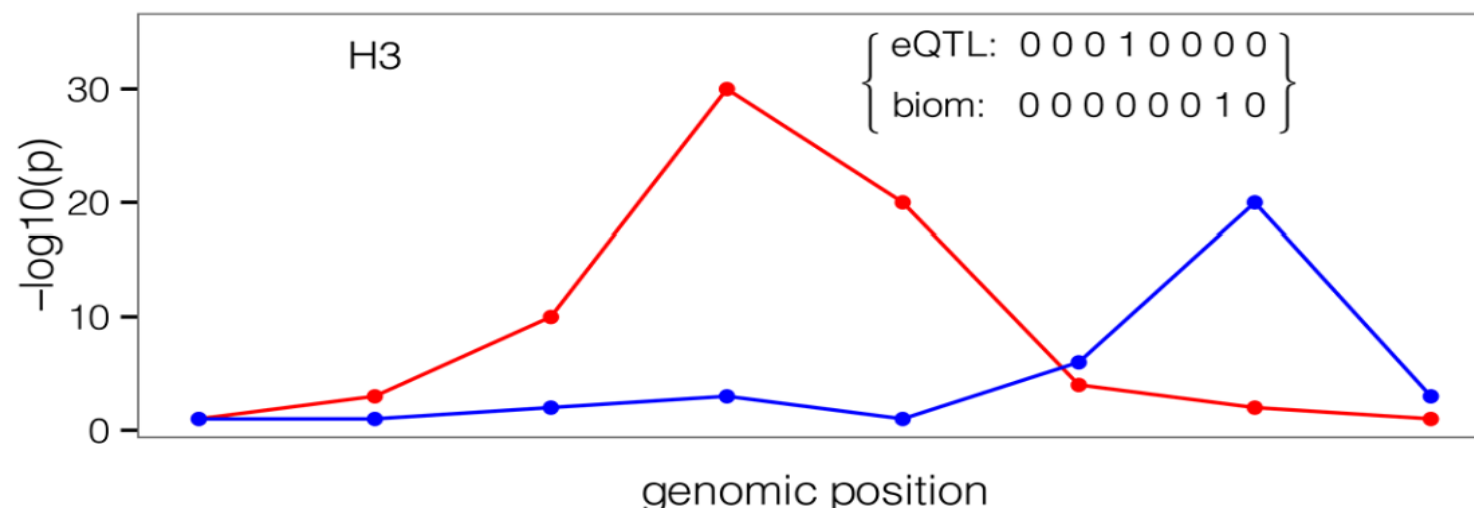
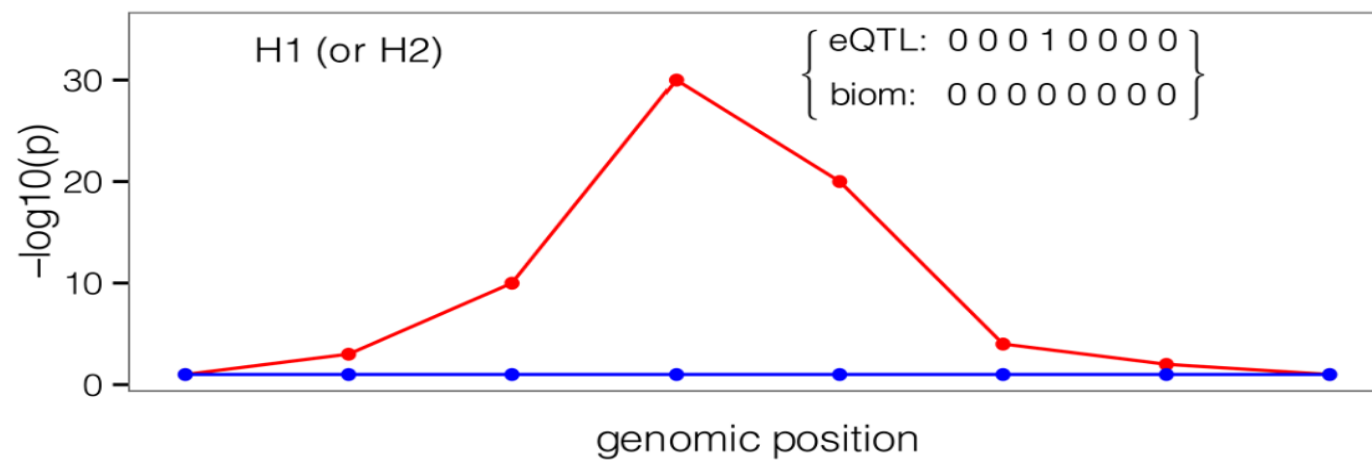
The key idea is simple: introduce multiple single-effect vectors $\mathbf{b}_1, \dots, \mathbf{b}_L$ and construct the overall effect vector \mathbf{b} as the sum of these single effects. We call this the 'sum of single effects' regression model SuSiE:

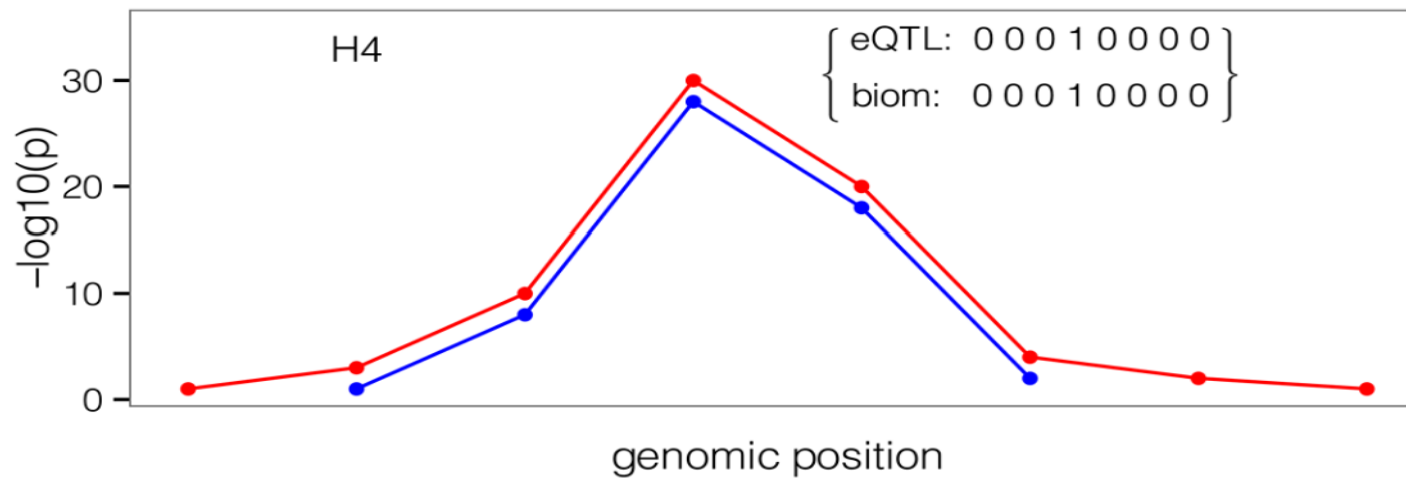
$$\begin{aligned}\mathbf{y} &= \mathbf{X}\mathbf{b} + \mathbf{e} \\ \mathbf{e} &\sim N_n(0, \sigma^2 I_n) \\ \mathbf{b} &= \sum_{l=1}^L \mathbf{b}_l \\ \mathbf{b}_l &= \gamma_l b_l \\ \gamma_l &\sim \text{Mult}(1, \pi) \\ b_l &\sim N_1(0, \sigma_{0l}^2)\end{aligned}$$



Fine-mapping (SuSiE) will uncover several credible sets , each capturing one causal variant.

COLOC





H_0 : no association with either trait in the region

H_1 : association with trait 1 only

H_2 : association with trait 2 only

H_3 : both traits are associated, but have different single causal variants

H_4 : both traits are associated and share the same single causal variant

Coloc

Assumption

- Firstly, that the causal variant is included in the set of Q variants, either directly typed or well imputed
- Secondly, that at most one association is present for each trait in the genomic region of interest.

Input

- Summary statistics
- Imputation

Coloc

$$\begin{aligned} PP4 &= P(H_4 \mid D) \\ &= \frac{P(H_4 \mid D)}{P(H_0 \mid D) + P(H_1 \mid D) + P(H_2 \mid D) + P(H_3 \mid D) + P(H_4 \mid D)} \\ &= \frac{\frac{P(H_4|D)}{P(H_0|D)}}{1 + \frac{P(H_1|D)}{P(H_0|D)} + \frac{P(H_2|D)}{P(H_0|D)} + \frac{P(H_3|D)}{P(H_0|D)} + \frac{P(H_4|D)}{P(H_0|D)}} \end{aligned}$$

- The ratios in the numerator and denominator are:

$$\frac{P(H_h \mid D)}{P(H_0 \mid D)} = \sum_{S \in S_h} \frac{P(D \mid S)}{P(D \mid S_0)} \times \frac{P(S)}{P(S_0)}$$

Coloc

ABF

- Asymptotically, that is as n increase, the MLE $\hat{\theta}$ has the normal distribution $N(\theta, V)$.
- Combining this "likelihood" with a normal prior, $N(0, W)$, on the log relative risk, θ , gives the asymptotic Bayes factor

$$\text{ABF} = \sqrt{\frac{V + W}{V}} \exp \left(-\frac{z^2}{2} \frac{W}{(V + W)} \right)$$

where $z^2 = \hat{\theta}^2 / V$ is the usual Wald statistic. High/low values of the asymptotic Bayes factor occur when z^2 is small/large and correspond to evidence for/against the null hypothesis.

Coloc

- If $S \in S_0$, then $\frac{P(S)}{P(S_0)} = \frac{p_0^Q}{p_0^Q} = 1$
- If $S \in S_1$, then $\frac{P(S)}{P(S_0)} = \frac{p_0^{Q-1}}{p_0^Q} \times p_1 = \frac{p_1}{p_0} \approx p_1$
- If $S \in S_2$, then $\frac{P(S)}{P(S_0)} = \frac{p_0^{Q-1}}{p_0^Q} \times p_2 = \frac{p_2}{p_0} \approx p_2$
- If $S \in S_3$, then $\frac{P(S)}{P(S_0)} = \frac{p_0^{Q-2}}{p_0^Q} \times p_1 \times p_2 = \frac{p_1}{p_0} \times \frac{p_2}{p_0} \approx p_1 \times p_2$
- If $S \in S_4$, then $\frac{P(S)}{P(S_0)} = \frac{p_0^{Q-1}}{p_0^Q} \times p_{12} = \frac{p_{12}}{p_0} \approx p_{12}$

Coloc

- $\frac{P(H_0|D)}{P(H_0|D)} = 1$
 - $\frac{P(H_1|D)}{P(H_0|D)} = p_1 \times \sum_{j=1}^Q ABF_j^1$
 - $\frac{P(H_2|D)}{P(H_0|D)} = p_2 \times \sum_{j=1}^Q ABF_j^2$
 - $\frac{P(H_3|D)}{P(H_0|D)} = p_1 \times p_2 \times \sum_{j,k,j \neq k} ABF_j^1 ABF_k^2$
 - $\frac{P(H_4|D)}{P(H_0|D)} = p_{12} \times \sum_{j=1}^Q ABF_j^1 \times ABF_j^2$
-
- Superscript 1,2 : trait1,2
 - Subscript ij : snp ij

Coloc

Coloc vs MR

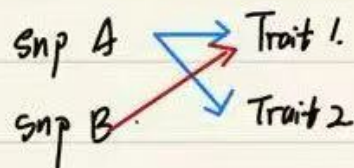
- MR Asymmetric in the traits: one trait is the exposure, the other is the outcome
- Coloc Symmetric in the traits: the traits are treated equivalently in the analysis
- Coloc is often used as sensitivity analysis after MR

Advantage & Disadvantage

- Robust
- Cumbersome if there are more than 2 traits (MOLOC)
- Difficult to scan the whole genome

Coloc & SuSiE

Case 1.



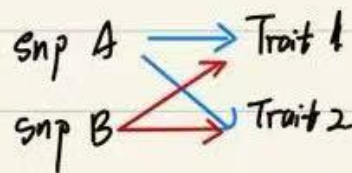
step 1. "SuSiE" → $\begin{cases} \text{Trait 1 } CS_{11} & CS_{21} \\ \text{Trait 2 } CS_{12} & \end{cases} \Rightarrow \text{选择每个CS中PIP最大的SNP}$

step 2. "merge" $\begin{matrix} \nearrow \text{Trait 1} \\ \text{SNP A} \\ \searrow \text{Trait 2} \end{matrix}$

step 3 "label" $\begin{matrix} \nearrow A: r^2(V_i, A) > 0.5 \text{ 且 } r^2(V_i, A) > r^2(V_i, B) \\ \text{SNP}_i \searrow B: r^2(V_i, B) > 0.5 \text{ 且 } r^2(V_i, B) > r^2(V_i, A) \end{matrix}$

step 4 "coloc" PP4?

Case 2



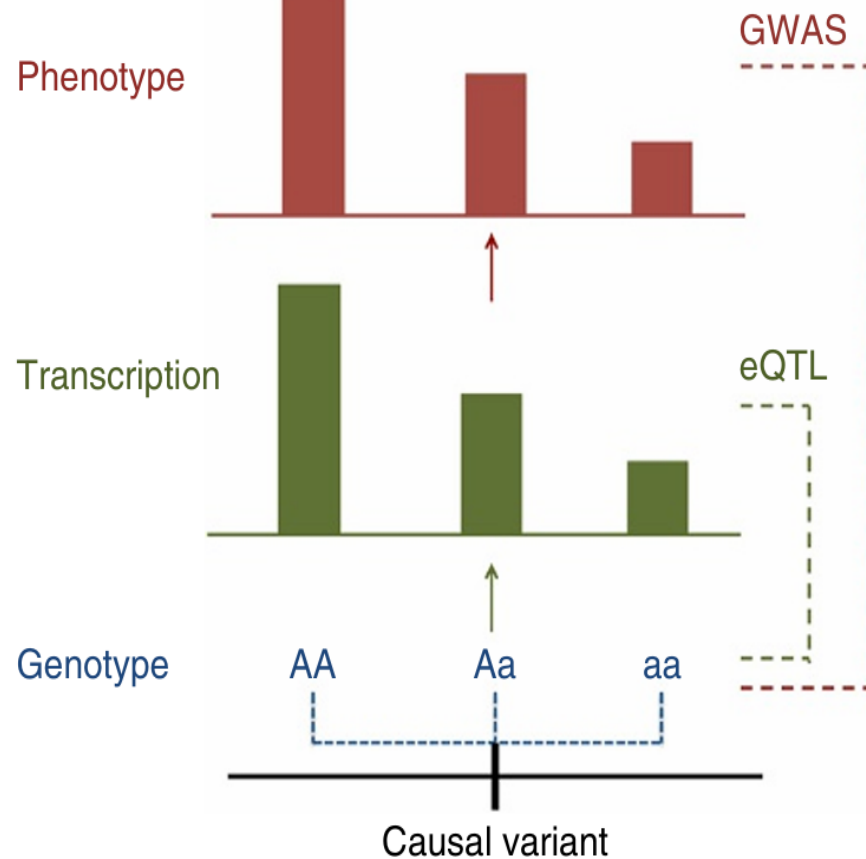
- pairwise coloc

- $H_4: AA$

$H_4: BB$

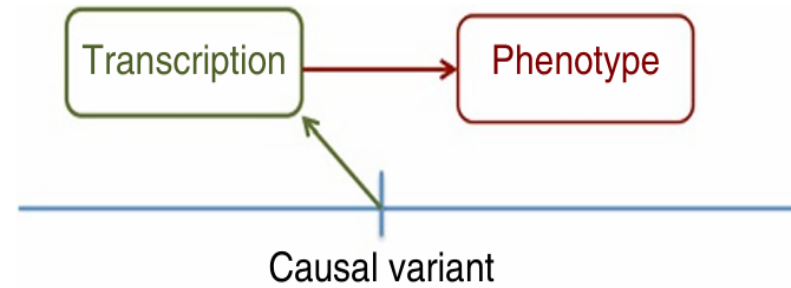
SMR & HEIDI

a

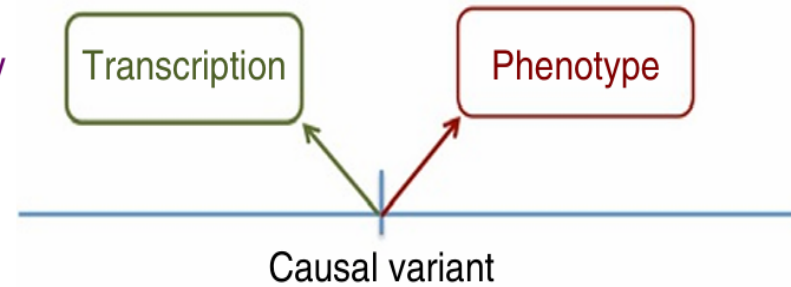


b

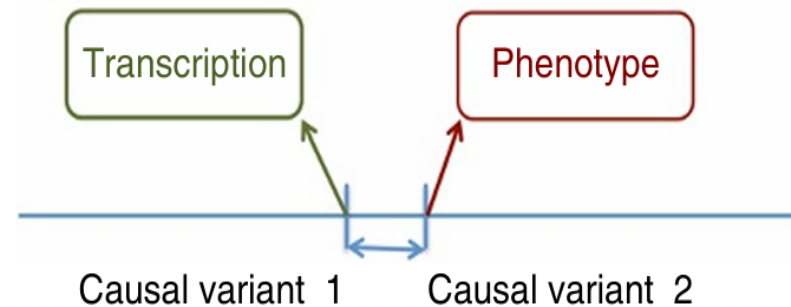
Causality



Pleiotropy



Linkage



HEIDI

- Analogue to Cochran Q test in MR_IVW
- Distinguish functional association from linkage.

If a trait and gene expression are affected by the same causal variant (pleiotropy), then b_{xy} calculated using any SNP in LD with the causal variant is identical. This is because, under Hardy-Weinberg equilibrium, for any SNP i

$$b_{xy(i)} = \frac{b_{zy(i)}}{\beta_{zx(i)}} = \frac{b_{zy(0)} r_{0i} \sqrt{h_0/h_i}}{\beta_{zx(0)} r_{0i} \sqrt{h_0/h_i}} = \frac{b_{zy(0)}}{\beta_{zx(0)}} = b_{xy(0)}$$

HEIDI

- Under the null hypothesis that there is no heterogeneity, that is, where $\mathbf{d} = 0$, we have a vector of standard normal variables $\mathbf{z}_d = \{z_{d(1)}, z_{d(2)}, \dots, z_{d(m)}\}$ with $z_{d(i)} = \hat{d}_i / \sqrt{\text{var}(\hat{d}_i)}$ and $\mathbf{z}_d \sim \text{MVN}(0, \mathbf{R})$, where \mathbf{R} is the correlation matrix with the ij th element being $r(z_{d(i)}, z_{d(j)}) = \text{cov}(\hat{d}_i, \hat{d}_j) / \sqrt{\text{var}(\hat{d}_i) \text{var}(\hat{d}_j)}$
- HEIDI test statistic: $T_{\text{HEIDI}} = \mathbf{z}_d \mathbf{z}_d^T$, that is, $T_{\text{HEIDI}} = \sum_i^m z_{d(i)}^2$

Multi-HEIDI

$$\hat{d}_{x_{ij}(s)} = \hat{b}_{x_i(s)} - \hat{b}_{x_i(\text{top } i)}$$

- where the subscript "i" represents the i th molecular phenotype, s represents an instrument from a set of x SNPs present in all the xQTL datasets and passing an xQTL p value threshold, and the subscript "top;" represents the lead xSNP for molecular phenotype i . The x SNPs are required to be instrument SNPs (i.e., $P_{XQTL} < 1.6 \times 10^{-3}$) for all tested molecular phenotypes.
- The null hypothesis of the multi-exposure HEIDI test is that all $d_{x_{iy}(s)} = 0$ and the alternative hypothesis of the multi-exposure HEIDI test is that any $d_{x_{iy}(s)} \neq 0$. If $\hat{d}_{x_{iy}(s)}$ is computed for two xSNPs

Remark

- HEIDI+COJO

(部分) 解决了这个问题, 但是整体而言, 这个方法相较coloc对LD是比较敏感的, 且在有多个exposure时表现较差

- COLOC+SuSiE

利用了finemapping的结果, 但是对于哪些snp是某个casual snp的“LD代理”的做法比较草率, PIP最大的也未必是causal variant

$$\mathbf{A}r^2(v_i, A) > 0.5 \wedge r^2(v_i, A) > r^2(v_i, B)$$

$$\mathbf{B}r^2(v_i, B) > 0.5 \wedge r^2(v_i, B) > r^2(v_i, A)$$

Remark

- HEIDI
- Low power in detecting linkage when two variants are in high LD
- Single causal variant assumption
- Multi-HEIDI
- Reject "pleiotropy" when two causal variants are in LD
- Heterogeneity in different QTLs datasets