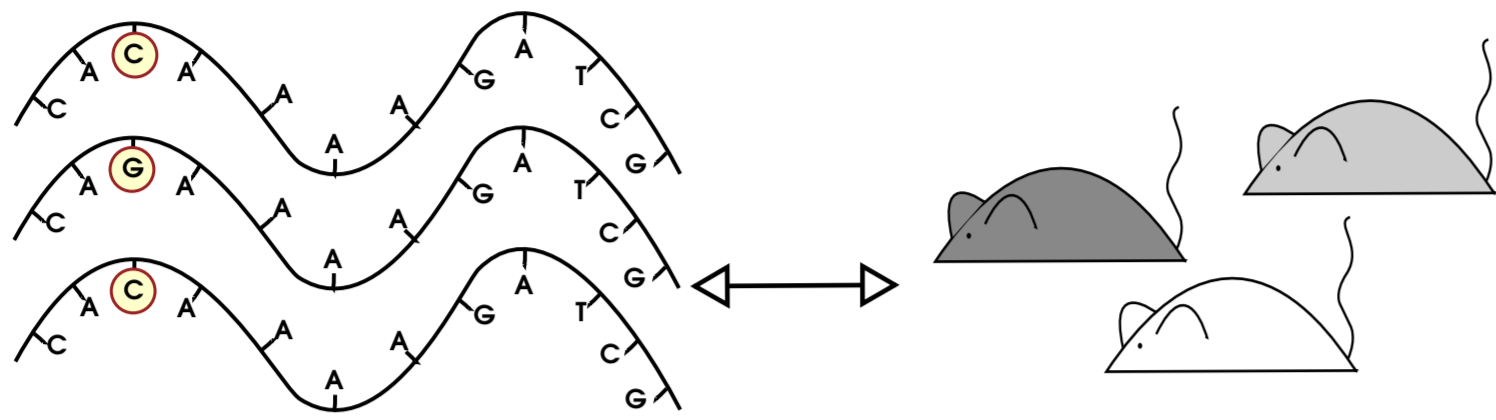




## Genome-Wide Association Studies (GWAS)



- 250 000 – 10 000 000 Single Nucleotide Polymorphisms (SNPs)
- 100 – 10 000 subjects

→ Which SNPs explain the phenotype?

For each SNP individually:

$$\underbrace{y_i}_{\text{Phenotype}} = \alpha_0 + \underbrace{\alpha^\top X_i}_{\text{Covariates}} + \underbrace{\beta g_i}_{\text{Genetics}} + \epsilon_i$$

## Missing heritability

Single-locus GWAS fail to explain most of the heritability of most complex traits

- Rare variants, undetected SNPs with small effect size, ...
- **Joint effects** of multiple SNPs

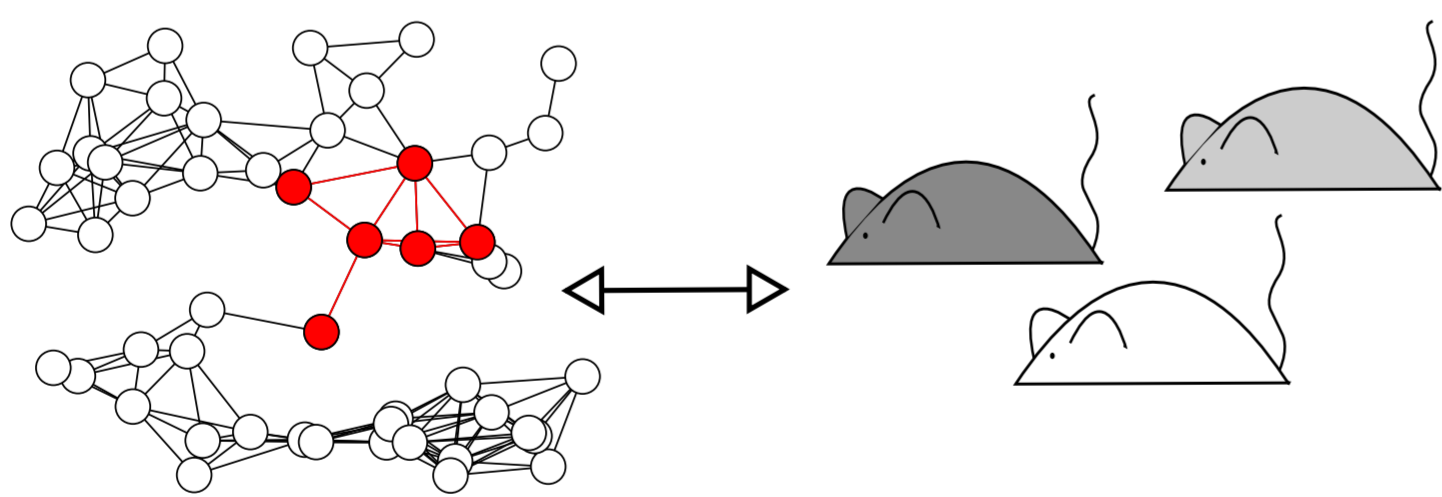
## Multiple-Locus GWAS

Find *multiple* SNPs that *jointly* explain the phenotype

Issue: search space grows exponentially with the number of such SNPs

- **Reduce the search space**
  - Use prior knowledge
  - Space-pruning & sampling techniques
- Limited to **binary / discrete phenotypes / genotypes** and **pairs of SNPs**
- **GPGPU approaches**
- Hardly scale to more than **pairs of SNPs**

## Networks of SNPs



## Selection of graph-structured features

- **Structured sparsity** [Huang et al. 2009]
- **Network-constrained lasso** [Li et al. 2008]
- **Overlapping group-lasso** [Jacob et al. 2009]
- Do not **scale** to hundreds of thousands of SNPs
- **Path-coding penalties** [Mairal and Yu 2012]
- **Directed acyclic graphs** only

## References

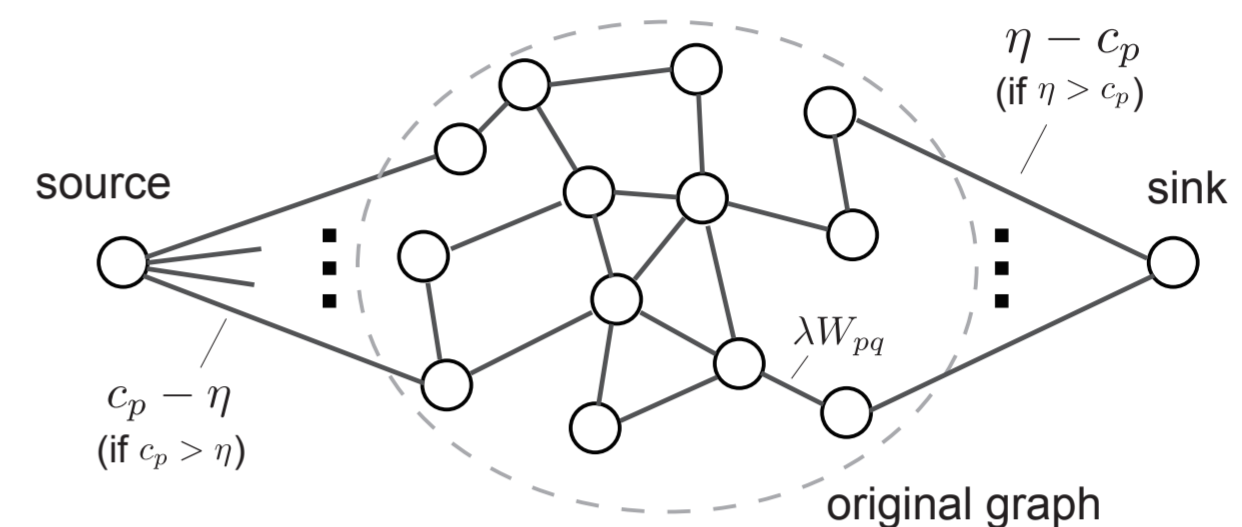
- C.-A. Azencott, D. Grimm, Y. Kawahara, K. Borgwardt. A min-cut solution to mapping phenotypes to networks of genetic markers. *arXiv:1211.2315* 2012
- J. Huang et al. Learning with structured sparsity. *JMLR* 2011
- L. Jacob et al. Group lasso with overlap and graph lasso. *ICML* 2009

## SOS: Subnetworks of Optimal SNPs

Find a *small* number of SNPs, *connected* in a given subnetwork, that *jointly* explain the phenotype

$$\arg \max_{f \in \{0,1\}^n} \underbrace{c^\top f}_{\text{linear association}^*} - \underbrace{\lambda f^\top L f}_{\text{connectivity}} - \underbrace{\eta \|f\|_1}_{\text{sparsity}}$$

Max-flow solution:



\*Sequence Kernel Association Test:

$$\underbrace{y_i}_{\text{Phenotype}} = \alpha_0 + \underbrace{\alpha^\top X_i}_{\text{covariates}} + \underbrace{\beta^\top G_i^S}_{\text{genotypic variation}} + \epsilon_i$$

$$H_0 : \beta = 0 \rightarrow \hat{\mu} = \hat{\alpha}_0 + X_i \hat{\alpha}$$

Variance component test statistic:

$$Q^S = (y - \hat{\mu})^\top K^S (y - \hat{\mu})$$

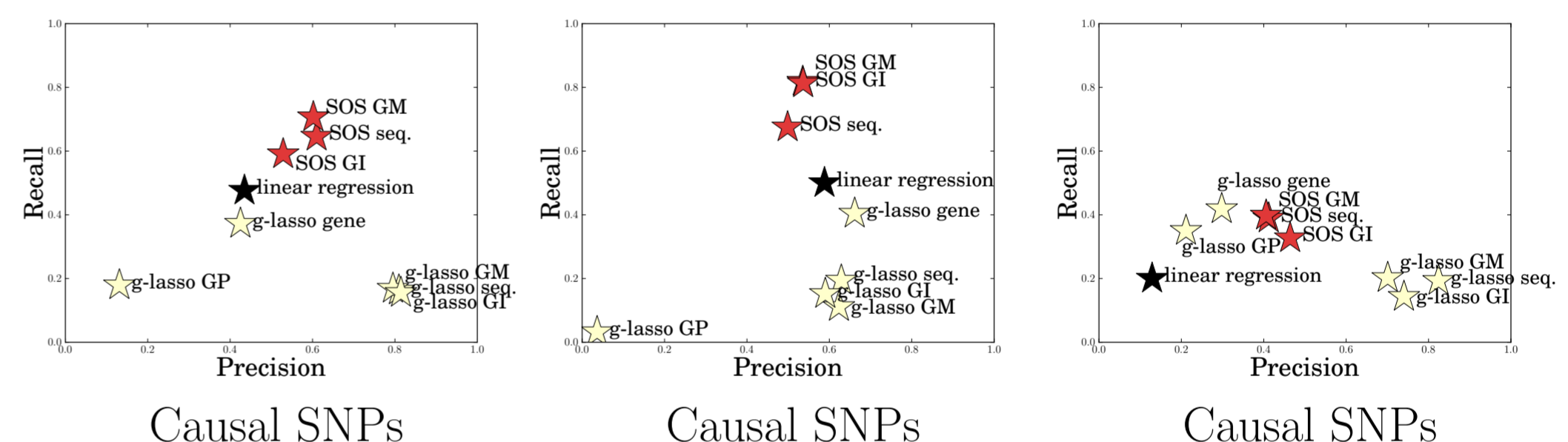
$$K^S = G^S G^{S\top}$$

$$Q(f^S) = c^\top f^S$$

$$c_p = (G^\top (y - \hat{\mu}))_p^2$$

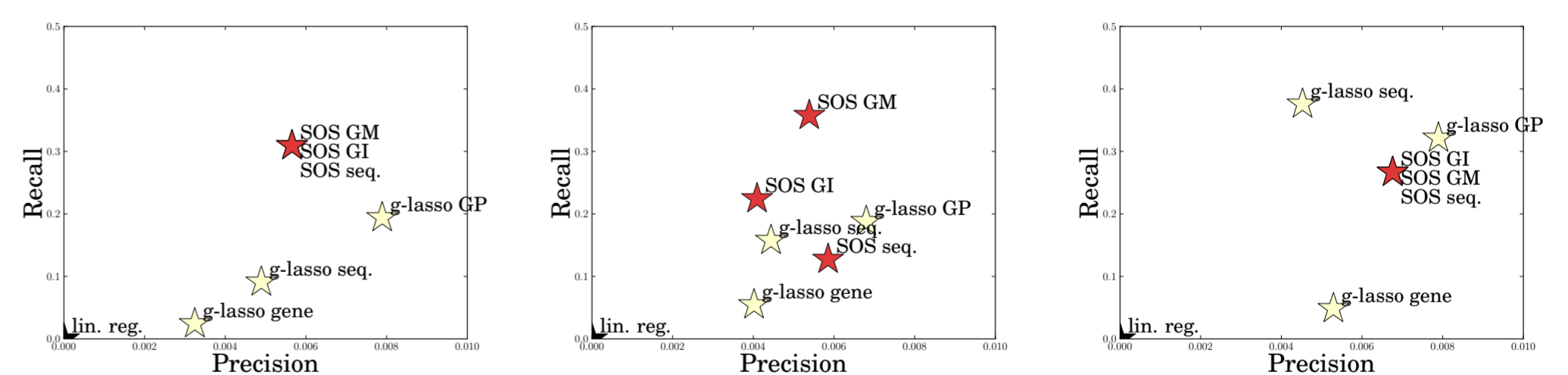
$$f_p^S = \begin{cases} 1 & \text{if } p \in \mathcal{S} \\ 0 & \text{otherwise} \end{cases}$$

## Results on simulated data



“GM”: SNPs in a sequence in a gene are connected; “GI”: groups of SNPs near either of 2 interacting genes are connected; “GP”: groups of SNPs near either of 2 interacting genes. 1,000 SNPs, 500 samples, 10 10-fold cv.

## Arabidopsis thaliana flowering time



4W FTGH LN22  
~ 250,000 SNPs, ~ 800 samples, 10-fold cv. g-lasso GM, GI did not complete after 3 days.  
Prec/recall wrt candidate genes.

## Further research topics

- Generalization to **submodular penalties**
- Extension to **multiple phenotypes**
- Extension to **epistatic effects**
- Definition of the **SNP network**

- C. Li and H. Li. Network-constrained regularization and variable selection for analysis of genomic data. *Bioinformatics* 2008
- J. Mairal and B. Yu. Path-coding penalties for directed acyclic graphs. *arXiv:1205.0079* 2012
- M. C. Wu et al. Rare-variant association testing for sequencing data with the sequence kernel association test. *Am J Hum Genet* 2011