

# A model selection approach for the identification of quantitative trait loci in experimental crosses

---

Karl W Broman

Department of Biostatistics, Johns Hopkins University

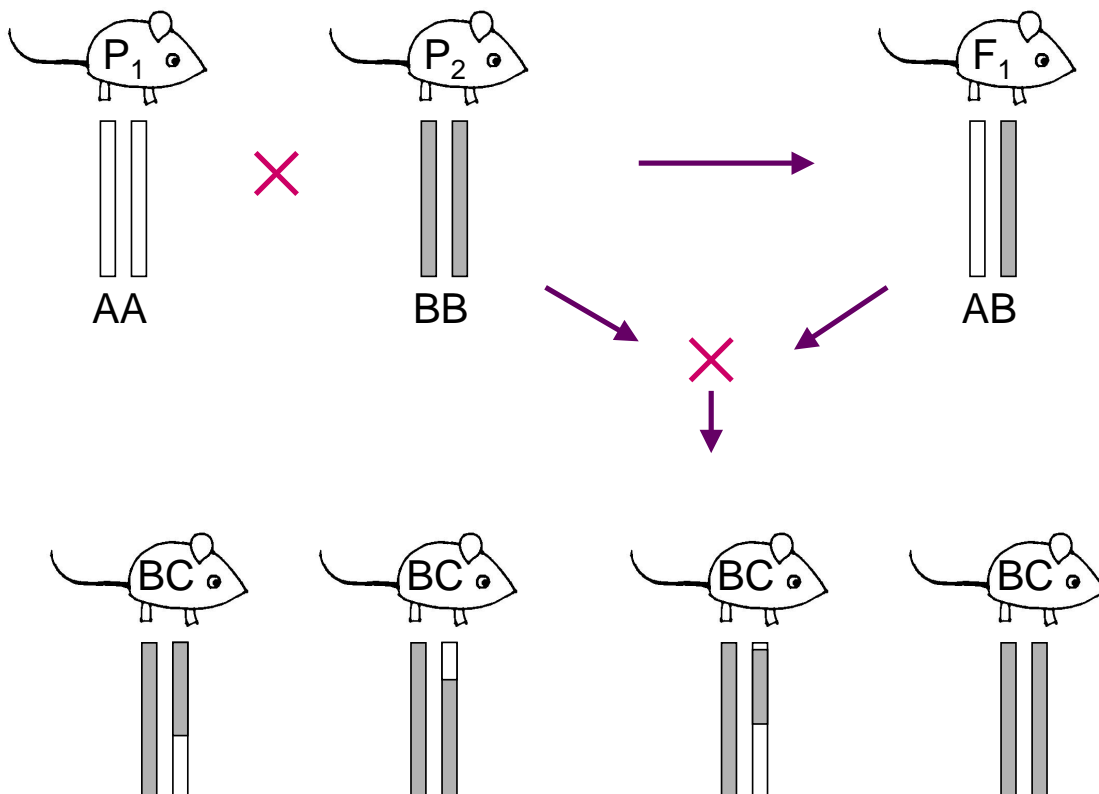
Terry Speed

Department of Statistics, University of California, Berkeley

Walter and Eliza Hall Institute (Melbourne, Australia)

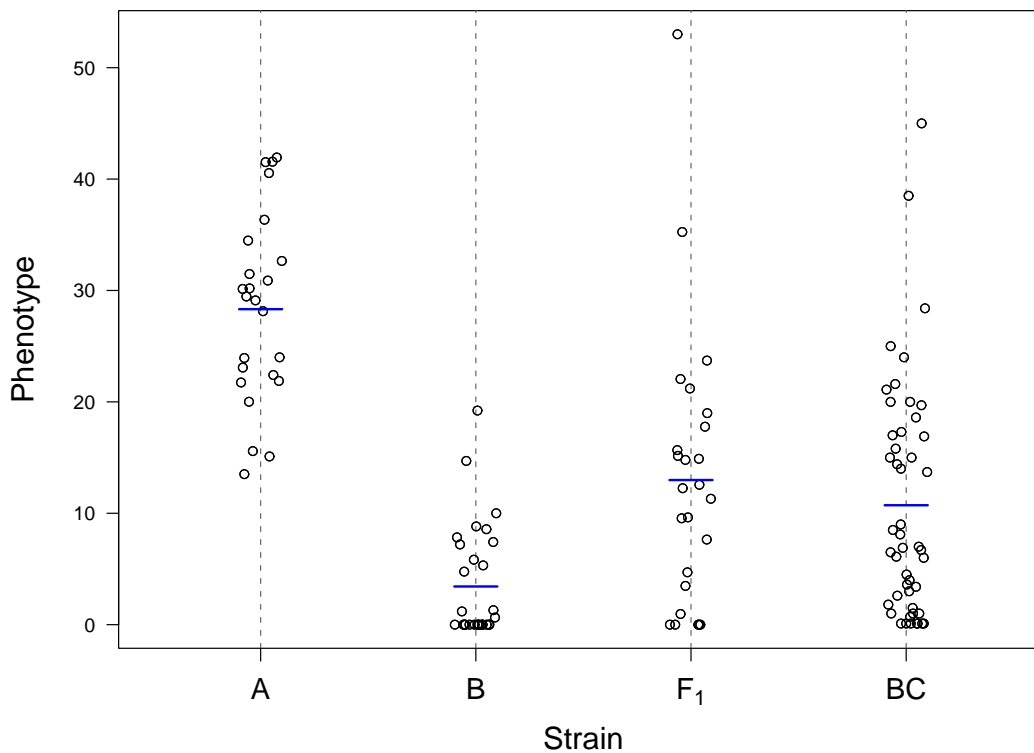
## Backcross experiment

---



# Trait distributions

---



## Data and Goals

---

**Phenotypes:**

$y_i$  = trait value for mouse  $i$

**Genotypes:**

$x_{ij}$  = 1/0 if mouse  $i$  is BB/AB at marker  $j$   
(for a backcross)

**Genetic map:**

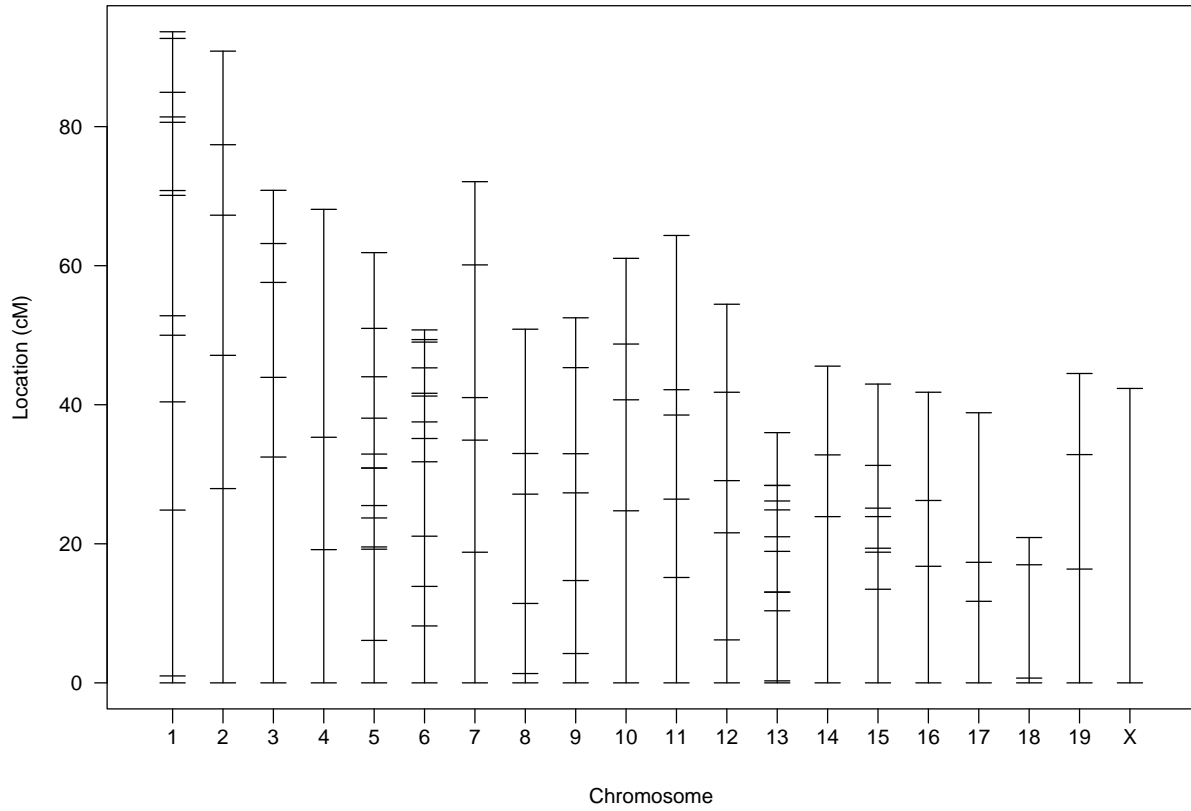
Locations of markers

---

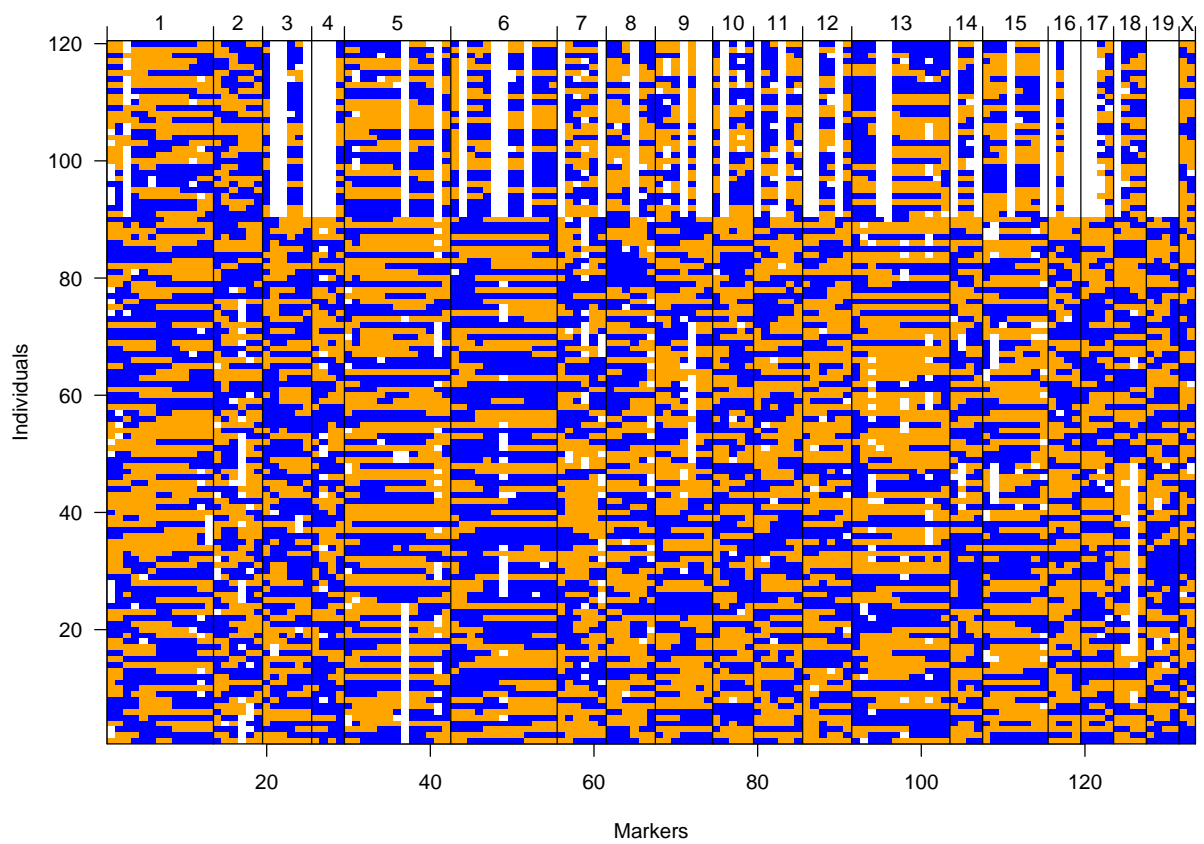
**Goals:**

- Identify the (or at least one) genomic regions (QTLs) that contribute to variation in the trait.
- Form confidence intervals for QTL locations.
- Estimate QTL effects.

## Genetic map



## Genotype data



# Models: Recombination

---

We assume no crossover interference.

⇒ Points of exchange (crossovers) are according to a Poisson process.

⇒ The  $\{x_{ij}\}$  (marker genotypes) form a Markov chain

## Models: Genotype $\longleftrightarrow$ Phenotype

---

Let  $y$  = phenotype  
 $g$  = whole genome genotype

Imagine a small number of QTLs with genotypes  $g_1, \dots, g_p$ .  
( $2^p$  distinct genotypes)

$$E(y|g) = \mu_{g_1, \dots, g_p} \quad \text{var}(y|g) = \sigma_{g_1, \dots, g_p}^2$$

**Homoscedasticity** (constant variance):  $\sigma_g^2 \equiv \sigma^2$

**Normally distributed residual variation:**  $y|g \sim N(\mu_g, \sigma^2)$ .

**Additivity:**  $\mu_{g_1, \dots, g_p} = \mu + \sum_{j=1}^p \Delta_j g_j$  ( $g_j = 1$  or  $0$ )

**Epistasis:** Any deviations from additivity.

# Abstractions / simplifications

---

- Complete marker data
  - QTLs are at the marker loci
  - QTLs act additively
- This work is not **useful in practice** but serves to **illustrate** the key issues.

## The problem

---

n backcross mice; M markers

$x_{ij}$  = genotype (1/0) of mouse  $i$  at marker  $j$

$y_i$  = phenotype (trait value) of mouse  $i$

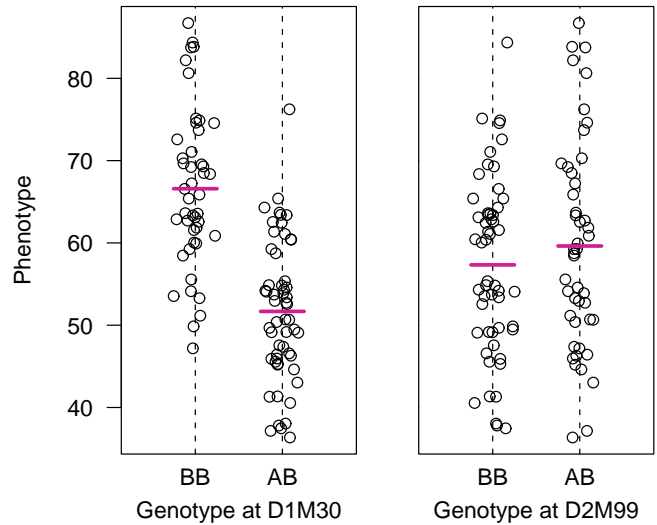
$$y_i = \mu + \sum_{j=1}^M \Delta_j x_{ij} + \epsilon_i \quad \text{Which } \Delta_j \neq 0?$$

- Errors:**
- Miss important loci
  - Include extraneous loci

# The simplest method: ANOVA

---

- Split mice into groups according to genotype at a marker.
- Do a t-test / ANOVA.
- Repeat for each marker.
- Adjust for multiple testing



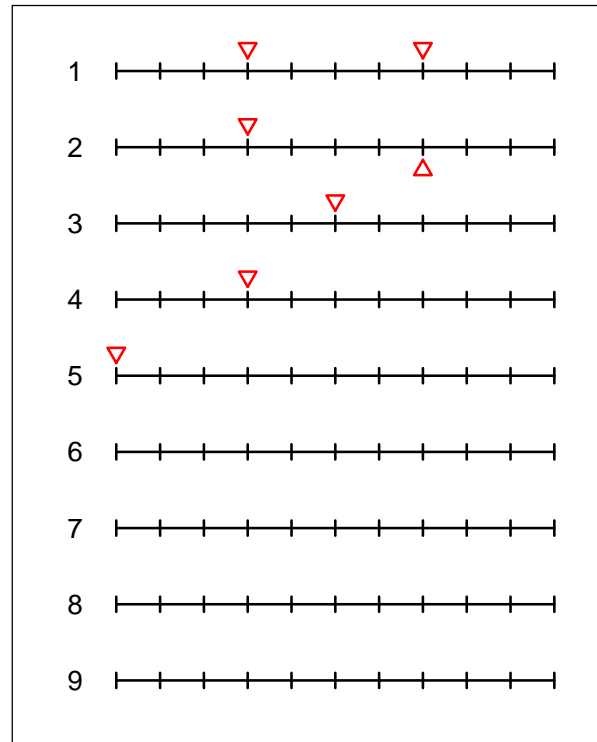
## Model selection

---

- **Select a class of models**
  - Additive models
  - Up to pairwise interactions
  - Regression trees
- **Compare models**
  - Estimated prediction error
  - $\text{BIC}_\delta(\gamma) = \log \text{RSS}_\gamma + \delta|\gamma| \log n/n$
  - Permutation tests
- **Search model space**
  - Forward selection
  - Backward elimination
  - Stepwise selection
  - MCMC
- **Assess the performance of a procedure**

# Simulations

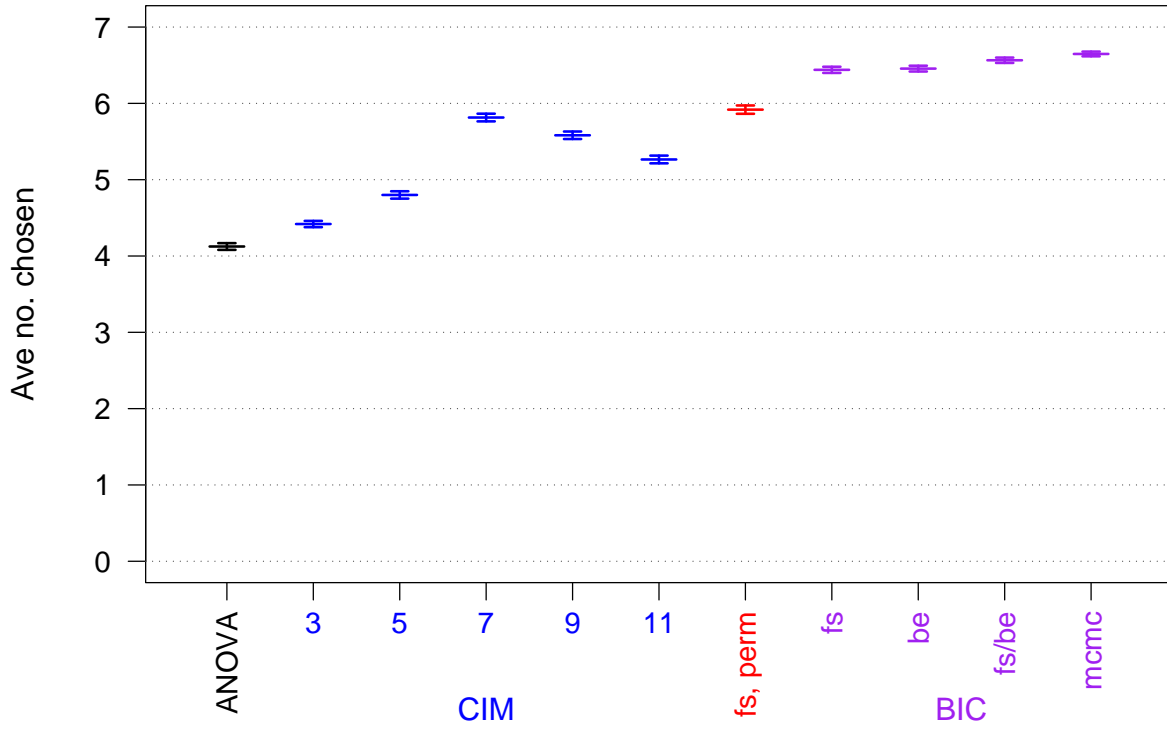
- Backcross with  $n=250$
- No crossover interference
- 9 chr, each 100 cM
- Markers at 10 cM spacing; complete genotype data
- 7 QTLs
  - One pair in **coupling**
  - One pair in **repulsion**
  - Three unlinked QTLs
- **Heritability** = 50%
- 2000 simulation replicates



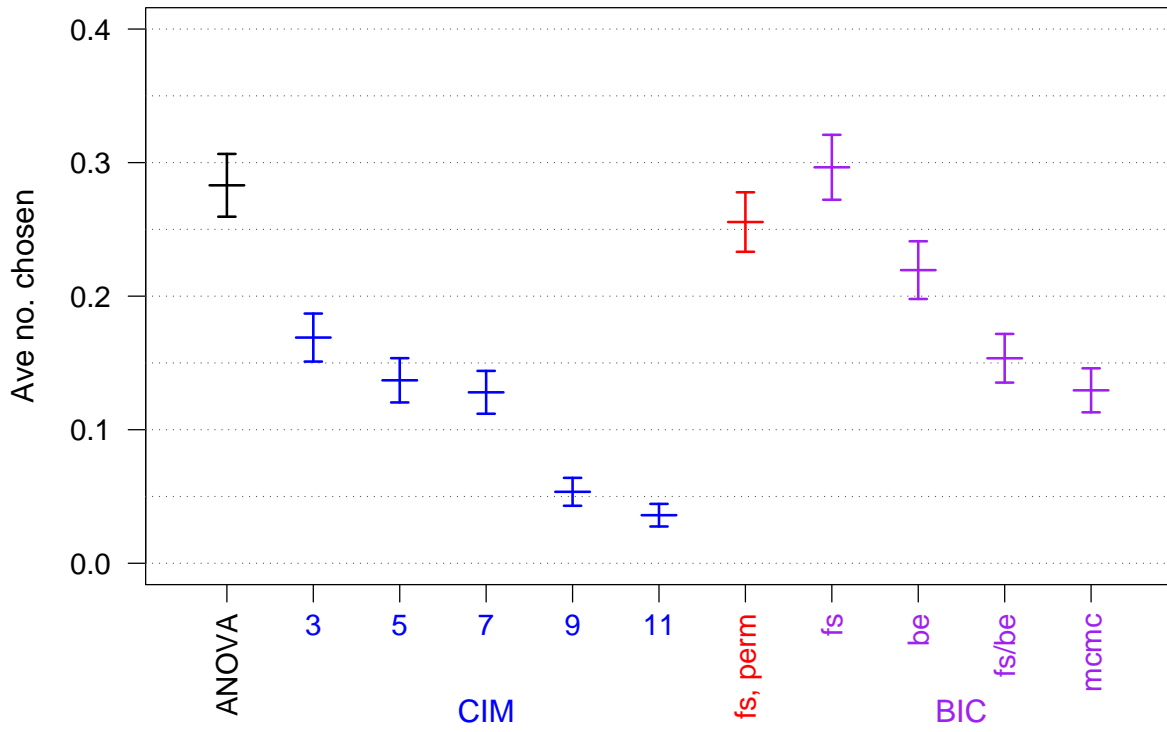
# Methods

- ANOVA at marker loci
  - Composite interval mapping (CIM)
  - Forward selection with permutation tests
  - Forward selection with  $BIC_{\delta}$
  - Backward elimination with  $BIC_{\delta}$
  - **FS followed by BE with  $BIC_{\delta}$**
  - MCMC with  $BIC_{\delta}$
- A **selected marker** is deemed **correct** if it is within 10 cM of a QTL (i.e., correct or adjacent)

Correct

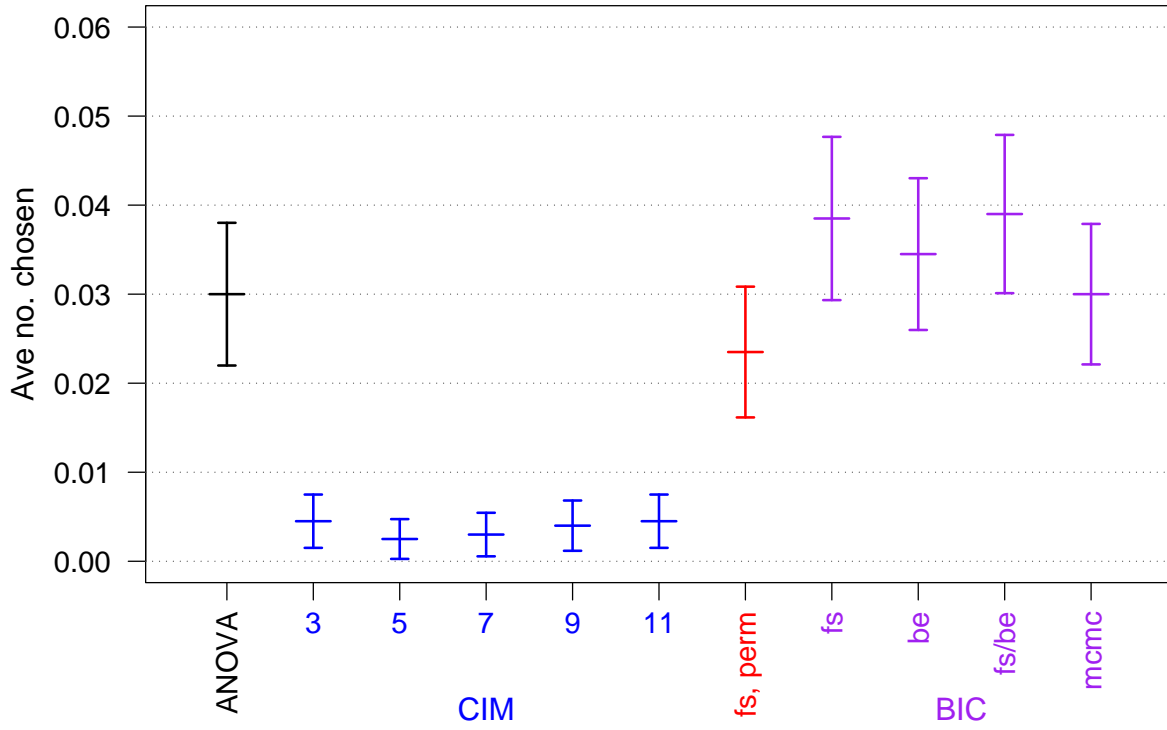


Extraneous linked

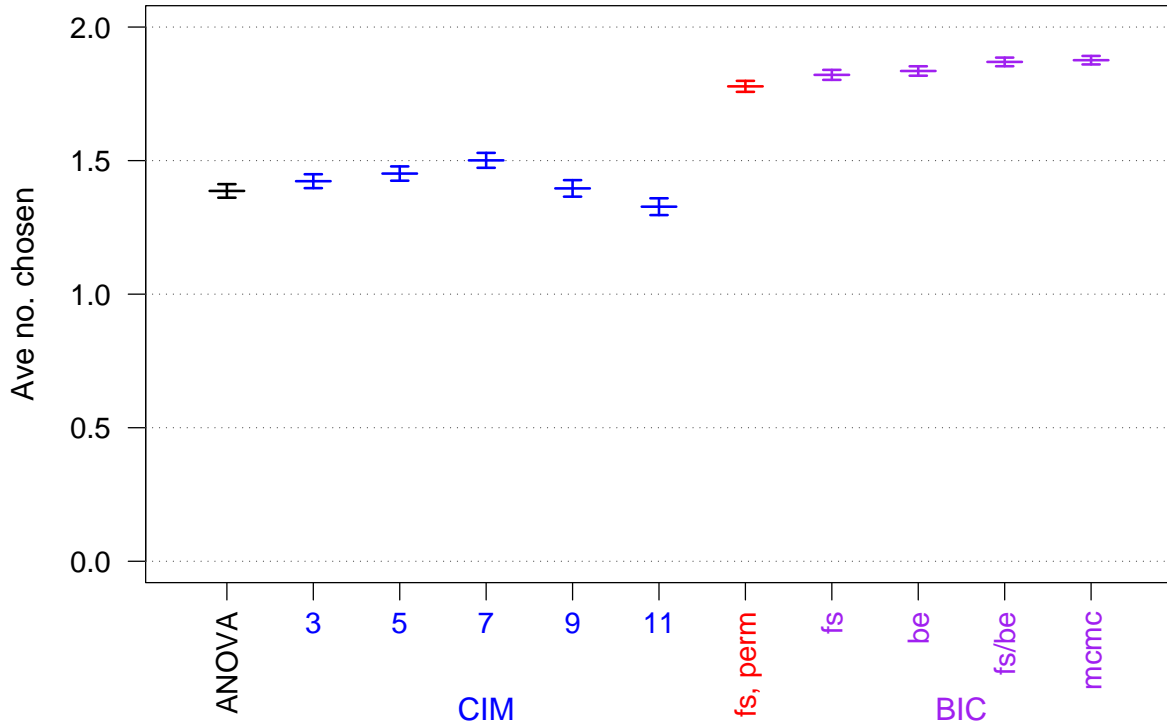




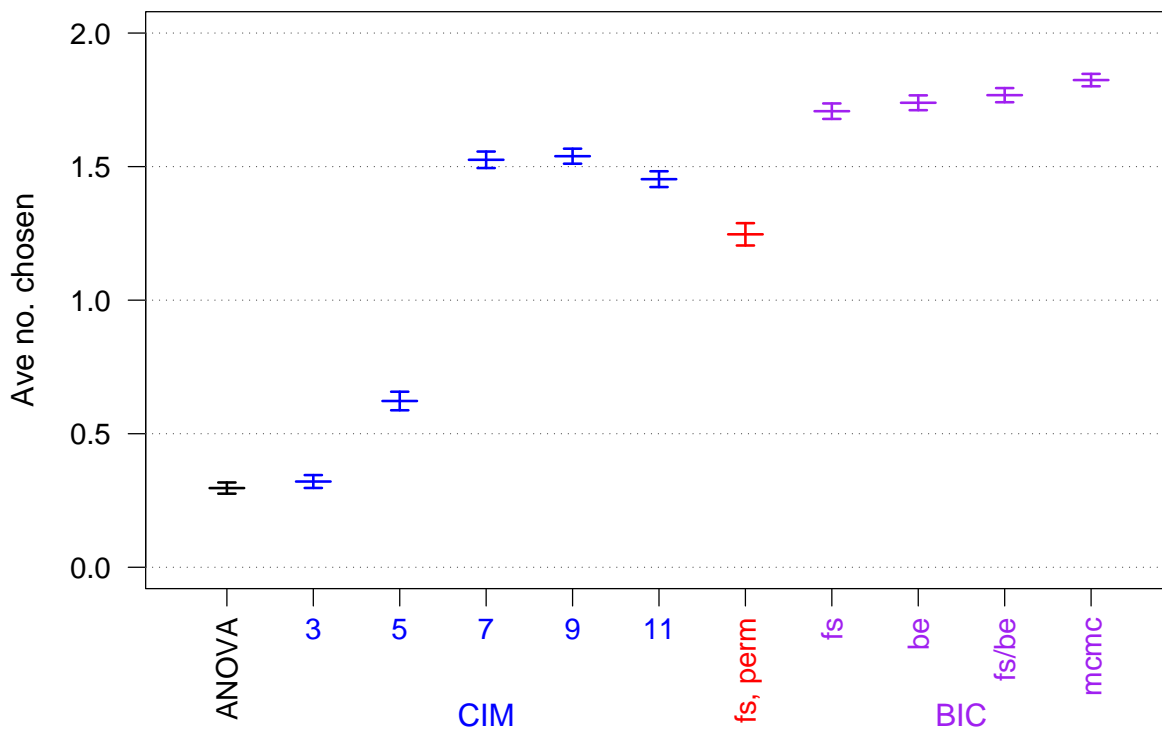
### Extraneous unlinked



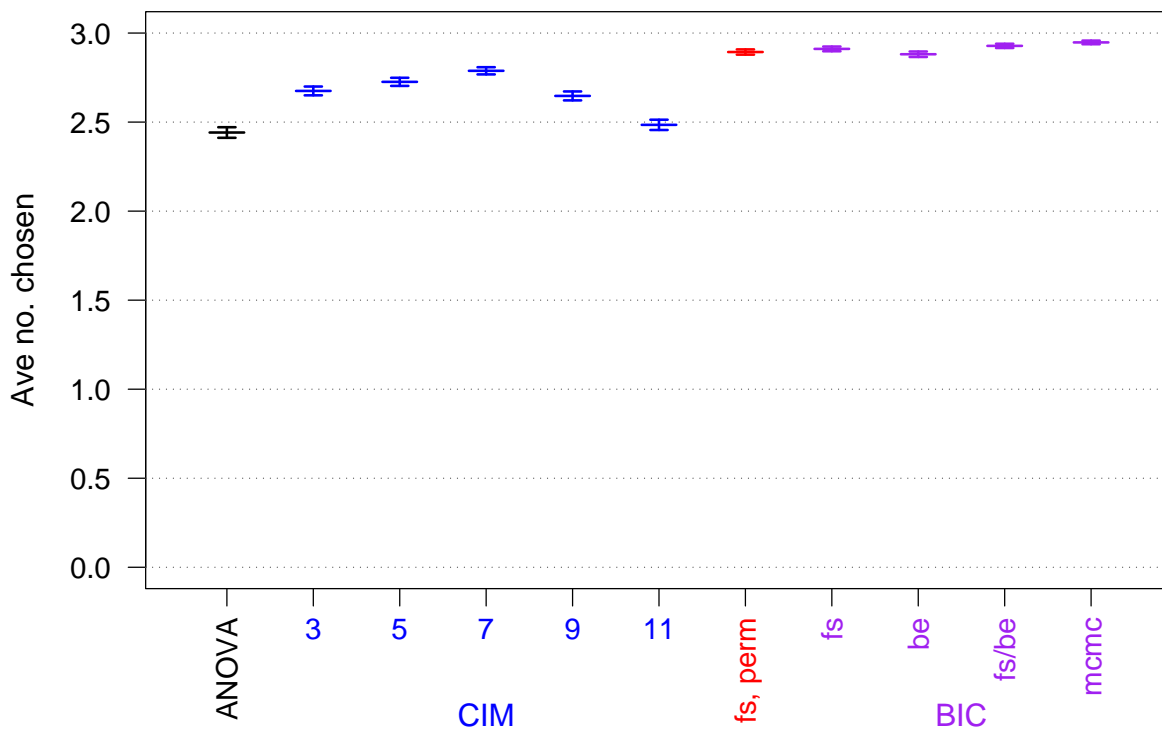
### QTLs linked in coupling



QTLs linked in repulsion



Other QTLs



# Summary

---

- QTL mapping is a **model selection** problem.
- Key issue: **the comparison of models**.
- Large-scale simulations are important.
- More refined procedures do not necessarily give improved results.
- **$BIC_\delta$**  with forward selection followed by backward elimination works quite well.