

# Seattle Summer Institute 2009

## Advanced QTL

Brian S. Yandell, UW-Madison  
[www.stat.wisc.edu/~yandell/statgen](http://www.stat.wisc.edu/~yandell/statgen)

- overview: multiple QTL approaches
- Bayesian QTL mapping & model selection
- data examples in detail
- software demos: R/qtl and R/qtlbim

*Real knowledge is to know the extent of one's ignorance.*  
Confucius (on a bench in Seattle)

# Overview of Multiple QTL

1. what is the goal of multiple QTL study?
2. gene action and epistasis
3. Bayesian vs. classical QTL
4. QTL model selection
5. QTL software options

# 1. what is the goal of QTL study?

- uncover underlying biochemistry
  - identify how networks function, break down
  - find useful candidates for (medical) intervention
  - epistasis may play key role
  - statistical goal: maximize number of correctly identified QTL
- basic science/evolution
  - how is the genome organized?
  - identify units of natural selection
  - additive effects may be most important (Wright/Fisher debate)
  - statistical goal: maximize number of correctly identified QTL
- select “elite” individuals
  - predict phenotype (breeding value) using suite of characteristics (phenotypes) translated into a few QTL
  - statistical goal: minimize prediction error

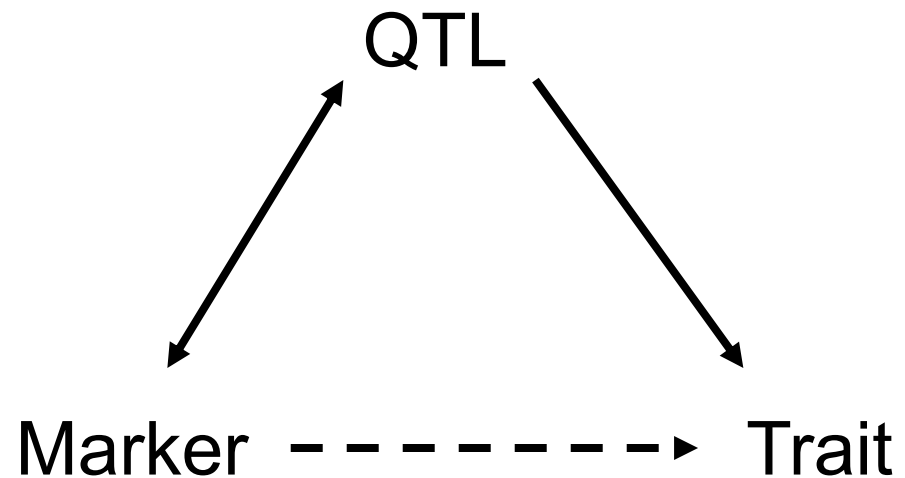
cross two inbred lines

→ linkage disequilibrium

→ associations

→ linked segregating QTL

(after Gary Churchill)



# problems of single QTL approach

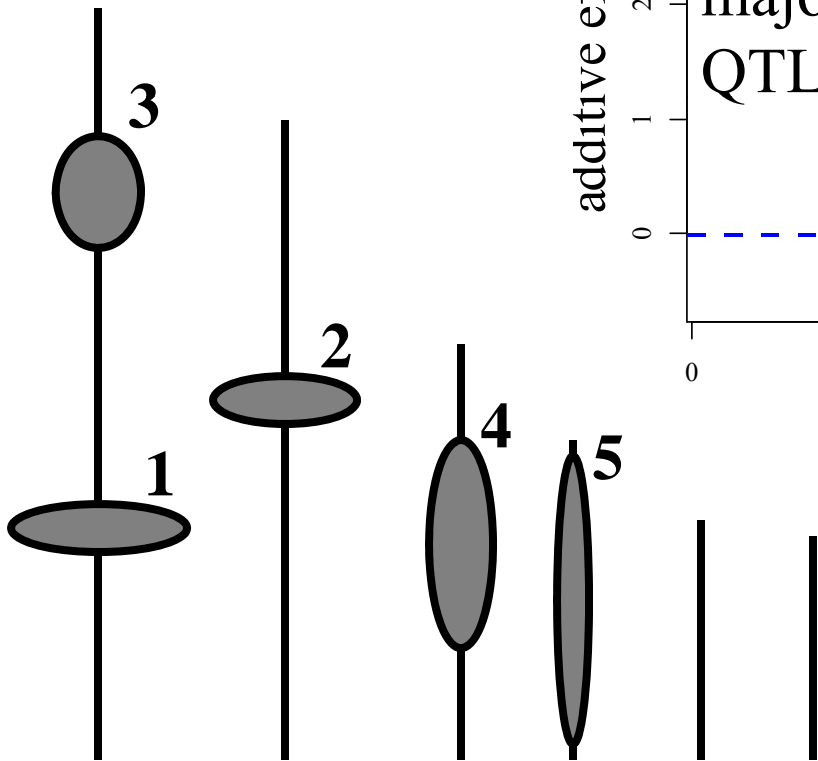
- wrong model: biased view
  - fool yourself: bad guess at locations, effects
  - detect ghost QTL between linked loci
  - miss epistasis completely
- low power
- bad science
  - use best tools for the job
  - maximize scarce research resources
  - leverage already big investment in experiment

# advantages of multiple QTL approach

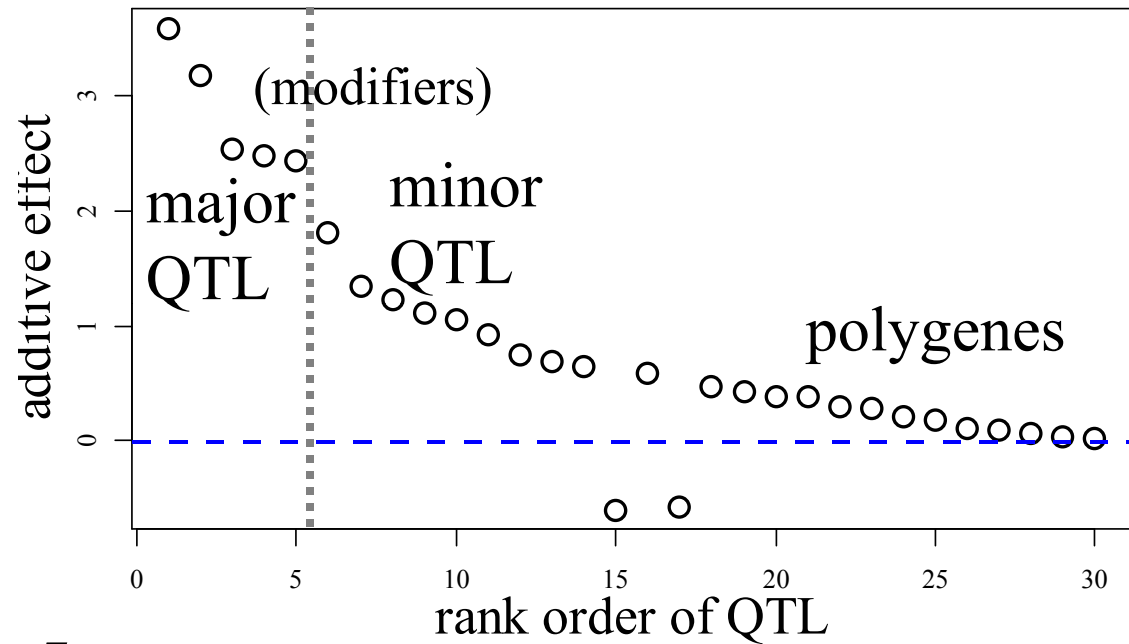
- improve statistical power, precision
  - increase number of QTL detected
  - better estimates of loci: less bias, smaller intervals
- improve inference of complex genetic architecture
  - patterns and individual elements of epistasis
  - appropriate estimates of means, variances, covariances
    - asymptotically unbiased, efficient
  - assess relative contributions of different QTL
- improve estimates of genotypic values
  - less bias (more accurate) and smaller variance (more precise)
  - mean squared error =  $MSE = (\text{bias})^2 + \text{variance}$

# Pareto diagram of QTL effects

major QTL on linkage map



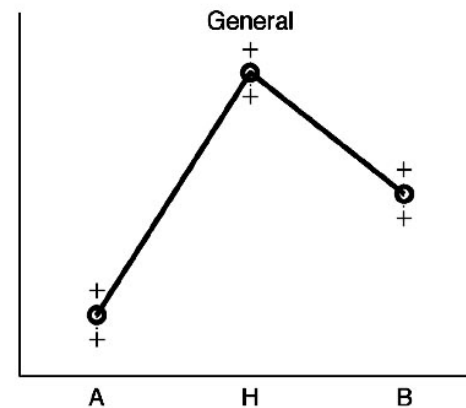
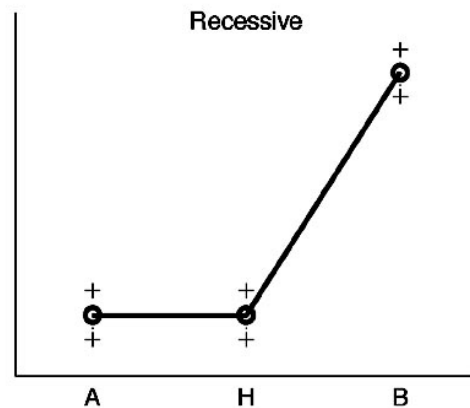
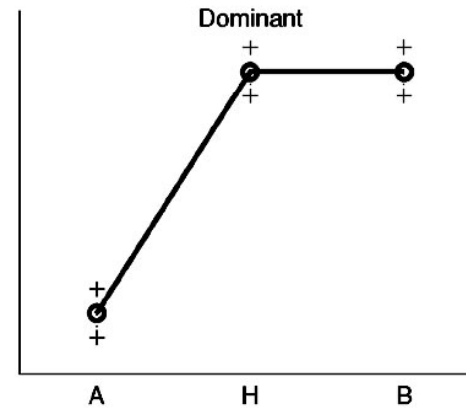
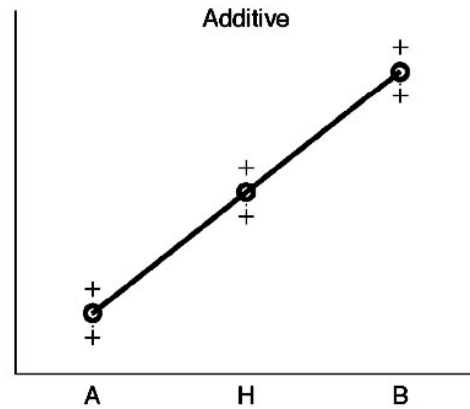
QTL 2: Overview



Seattle SISG: Yandell © 2009

## 2. Gene Action and Epistasis

additive, dominant, recessive, general effects  
of a single QTL (Gary Churchill)

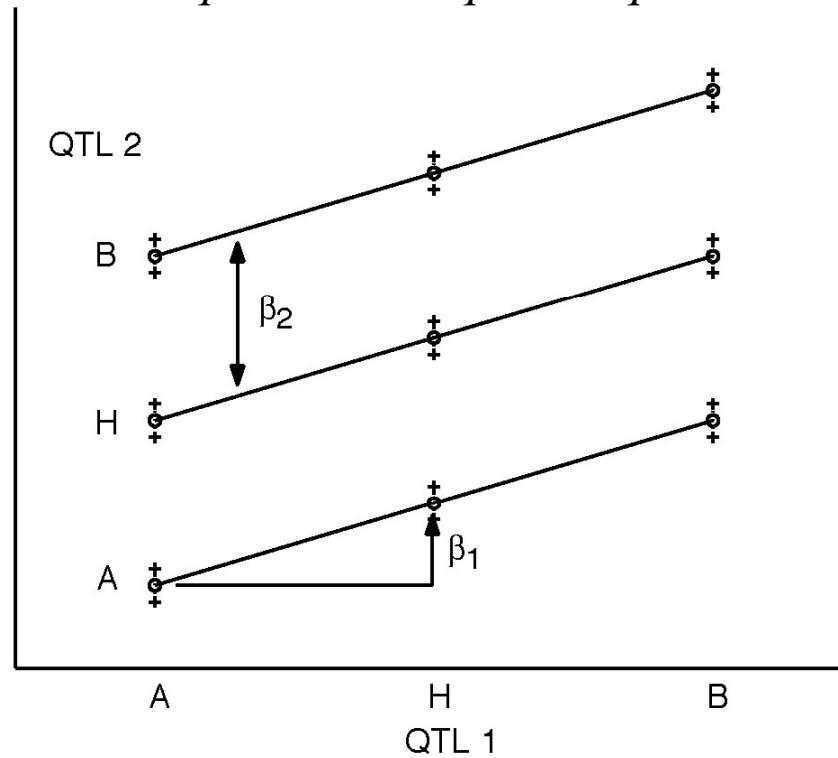


QTL 2: Overview



# additive effects of two QTL (Gary Churchill)

$$\mu_q = \mu + \beta_{q1} + \beta_{q2}$$



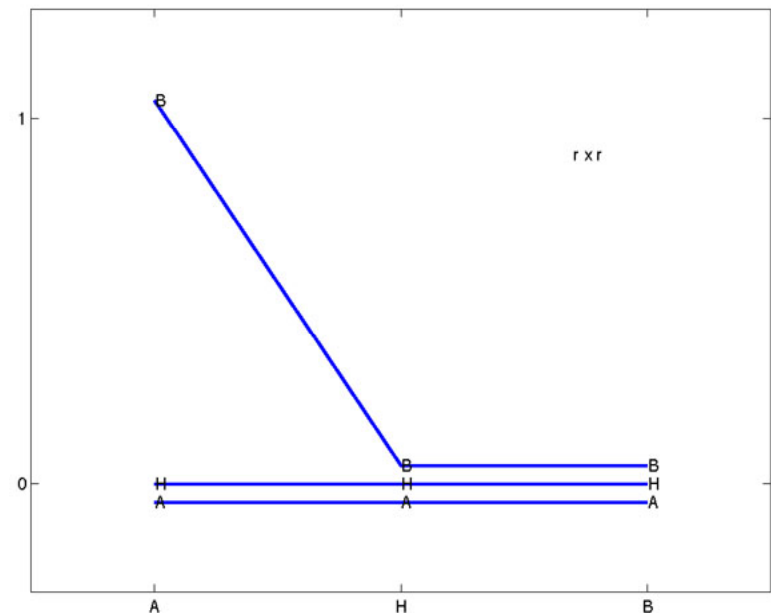
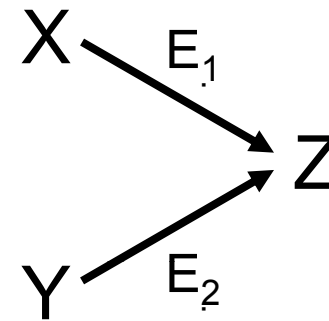
# Epistasis (Gary Churchill)

The allelic state at one locus can mask or uncover the effects of allelic variation at another.

- W. Bateson, 1907.

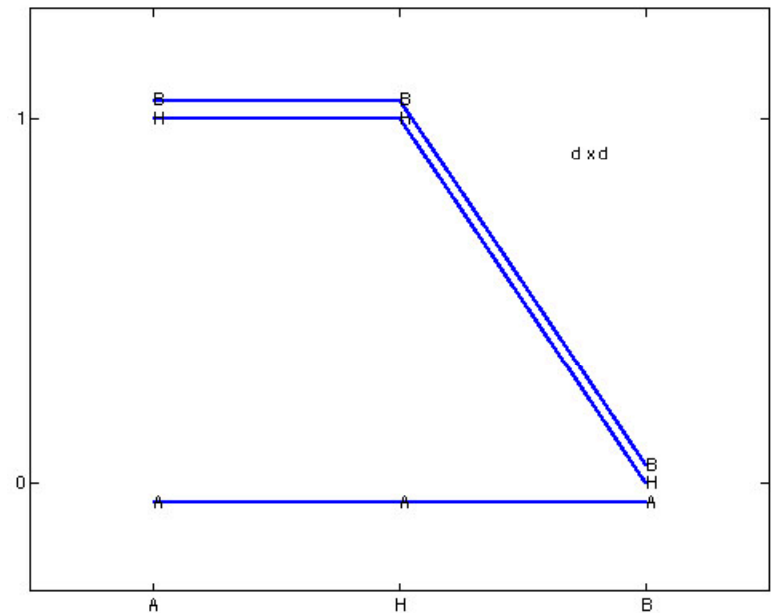
# epistasis in parallel pathways (GAC)

- Z keeps trait value low
- neither  $E_1$  nor  $E_2$  is rate limiting
- loss of function alleles are segregating from parent A at  $E_1$  and from parent B at  $E_2$



# epistasis in a serial pathway (GAC)

- Z keeps trait value high
- **either**  $E_1$  **or**  $E_2$  is rate limiting
- loss of function alleles are segregating from parent B at  $E_1$  **or** from parent A at  $E_2$



# epistatic interactions

- model space issues
  - 2-QTL interactions only?
    - or general interactions among multiple QTL?
  - partition of effects
    - Fisher-Cockerham or tree-structured or ?
- model search issues
  - epistasis between significant QTL
    - check all possible pairs when QTL included?
    - allow higher order epistasis?
  - epistasis with non-significant QTL
    - whole genome paired with each significant QTL?
    - pairs of non-significant QTL?
- see papers of Nengjun Yi (2000-7) in *Genetics*

# limits of epistatic inference

- power to detect effects
  - epistatic model sizes grow quickly
    - $|A| = 3^{n.qtl}$  for general interactions
  - power tradeoff
    - depends sample size *vs.* model size
    - want  $n / |A|$  to be fairly large (say  $> 5$ )
    - 3 QTL,  $n = 100$  F2:  $n / |A| \approx 4$
- rare genotypes may not be observed
  - $aa/BB$  &  $AA/bb$  rare for linked loci
  - empty cells mess up balance
    - adjusted tests (type III) are wrong
  - confounds main effects & interactions

2 linked QTL  
empty cell  
with  $n = 100$

	<i>bb</i>	<i>bB</i>	<i>BB</i>
<i>aa</i>	6	15	0
<i>aA</i>	15	25	15
<i>AA</i>	3	15	6

# limits of multiple QTL?

- limits of statistical inference
  - power depends on sample size, heritability, environmental variation
  - “best” model balances fit to data and complexity (model size)
  - genetic linkage = correlated estimates of gene effects
- limits of biological utility
  - sampling: only see some patterns with many QTL
  - marker assisted selection (Bernardo 2001 *Crop Sci*)
    - 10 QTL ok, 50 QTL are too many
    - phenotype better predictor than genotype when too many QTL
    - increasing sample size may not give multiple QTL any advantage
  - hard to select many QTL simultaneously
    - $3^m$  possible genotypes to choose from

# QTL below detection level?

- problem of selection bias
  - QTL of modest effect only detected sometimes
  - effects overestimated when detected
  - repeat studies may fail to detect these QTL
- think of probability of detecting QTL
  - avoids sharp in/out dichotomy
  - avoid pitfalls of one “best” model
  - examine “better” models with more probable QTL
- rethink formal approach for QTL
  - directly allow uncertainty in genetic architecture
  - QTL model selection over genetic architecture



### 3. Bayesian vs. classical QTL study

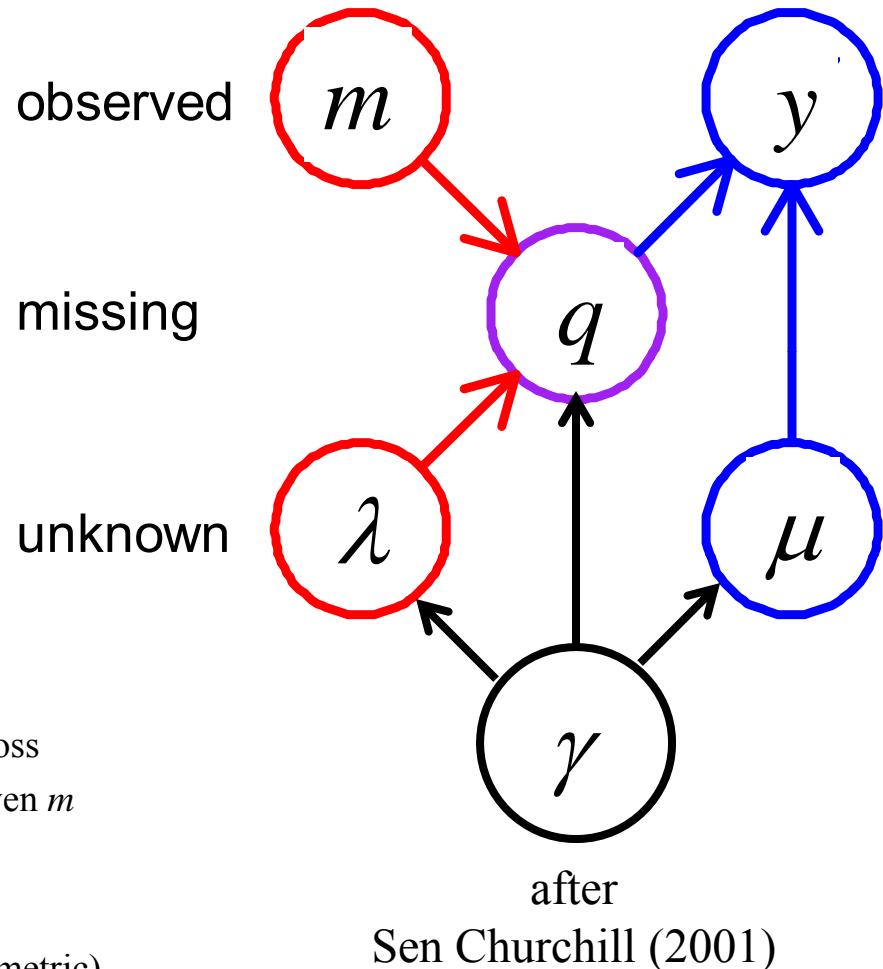
- classical study
  - *maximize* over unknown effects
  - *test* for detection of QTL at loci
  - model selection in stepwise fashion
- Bayesian study
  - *average* over unknown effects
  - *estimate* chance of detecting QTL
  - sample all possible models
- both approaches
  - average over missing QTL genotypes
  - scan over possible loci

# Bayesian idea

- Reverend Thomas Bayes (1702-1761)
  - part-time mathematician
  - buried in Bunhill Cemetary, Moongate, London
  - famous paper in 1763 *Phil Trans Roy Soc London*
  - was Bayes the first with this idea? (Laplace?)
- basic idea (from Bayes' original example)
  - two billiard balls tossed at random (uniform) on table
  - where is first ball if the second is to its left?
    - prior: anywhere on the table
    - posterior: more likely toward right end of table

# QTL model selection: key players

- observed measurements
  - $y$  = phenotypic trait
  - $m$  = markers & linkage map
  - $i$  = individual index  $(1, \dots, n)$
- missing data
  - missing marker data
  - $q$  = QT genotypes
    - alleles QQ, Qq, or qq at locus
- unknown quantities
  - $\lambda$  = QT locus (or loci)
  - $\mu$  = phenotype model parameters
  - $\gamma$  = QTL model/genetic architecture
- $\text{pr}(q|m, \lambda, \gamma)$  genotype model
  - grounded by linkage map, experimental cross
  - recombination yields multinomial for  $q$  given  $m$
- $\text{pr}(y|q, \mu, \gamma)$  phenotype model
  - distribution shape (assumed normal here)
  - unknown parameters  $\mu$  (could be non-parametric)



# Bayes posterior vs. maximum likelihood

- LOD: classical Log ODds
  - maximize likelihood over effects  $\mu$
  - R/qt1 scanone/scantwo: method = "em"
- *LPD*: Bayesian *Log Posterior Density*
  - average posterior over effects  $\mu$
  - R/qt1 scanone/scantwo: method = "imp"

$$\text{LOD}(\lambda) = \log_{10} \{ \max_{\mu} \text{pr}(y | m, \mu, \lambda) \} + c$$

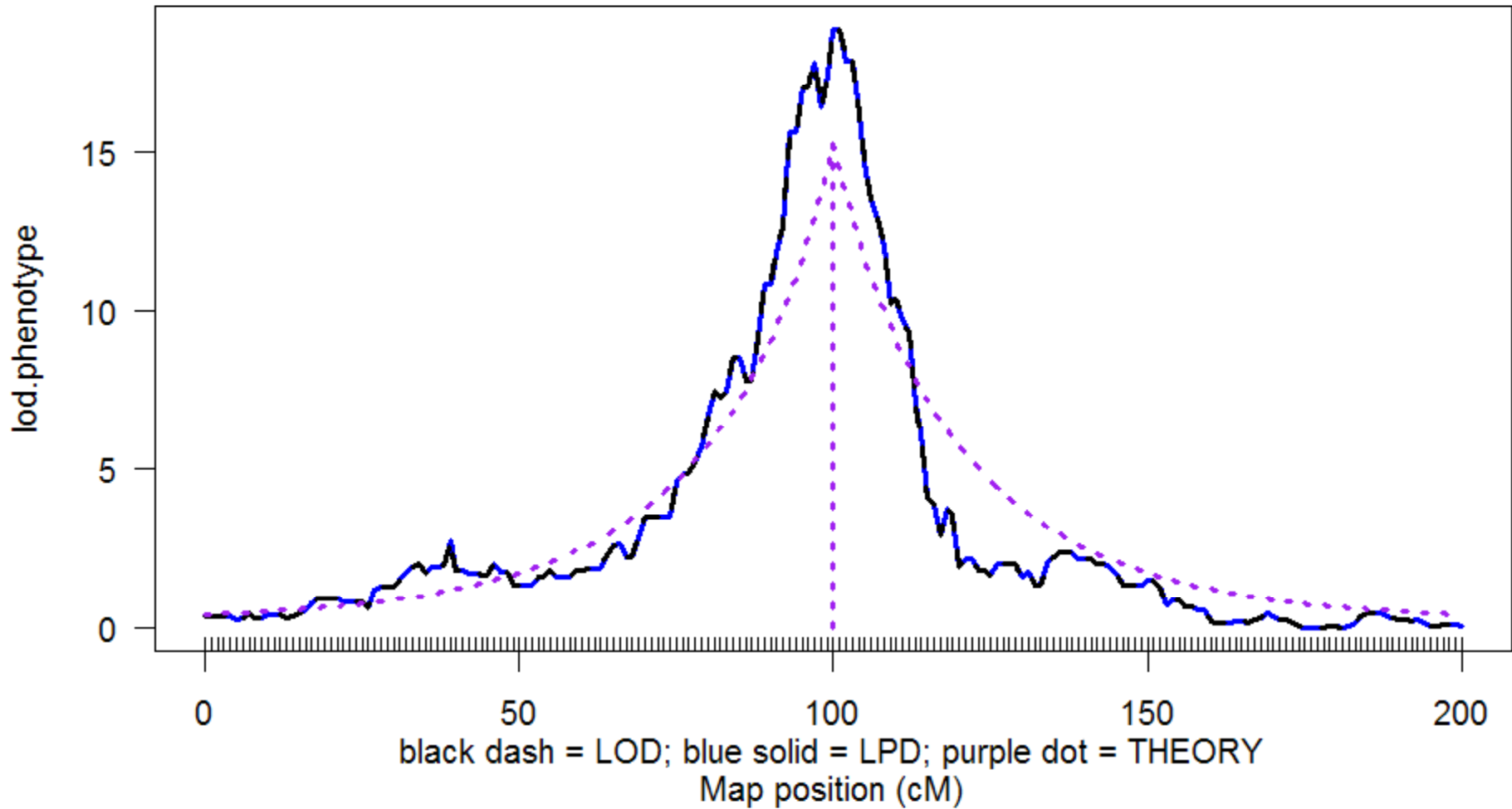
$$\text{LPD}(\lambda) = \log_{10} \{ \text{pr}(\lambda | m) \int \text{pr}(y | m, \mu, \lambda) \text{pr}(\mu) d\mu \} + C$$

likelihood mixes over missing QTL genotypes :

$$\text{pr}(y | m, \mu, \lambda) = \sum_q \text{pr}(y | q, \mu) \text{pr}(q | m, \lambda)$$

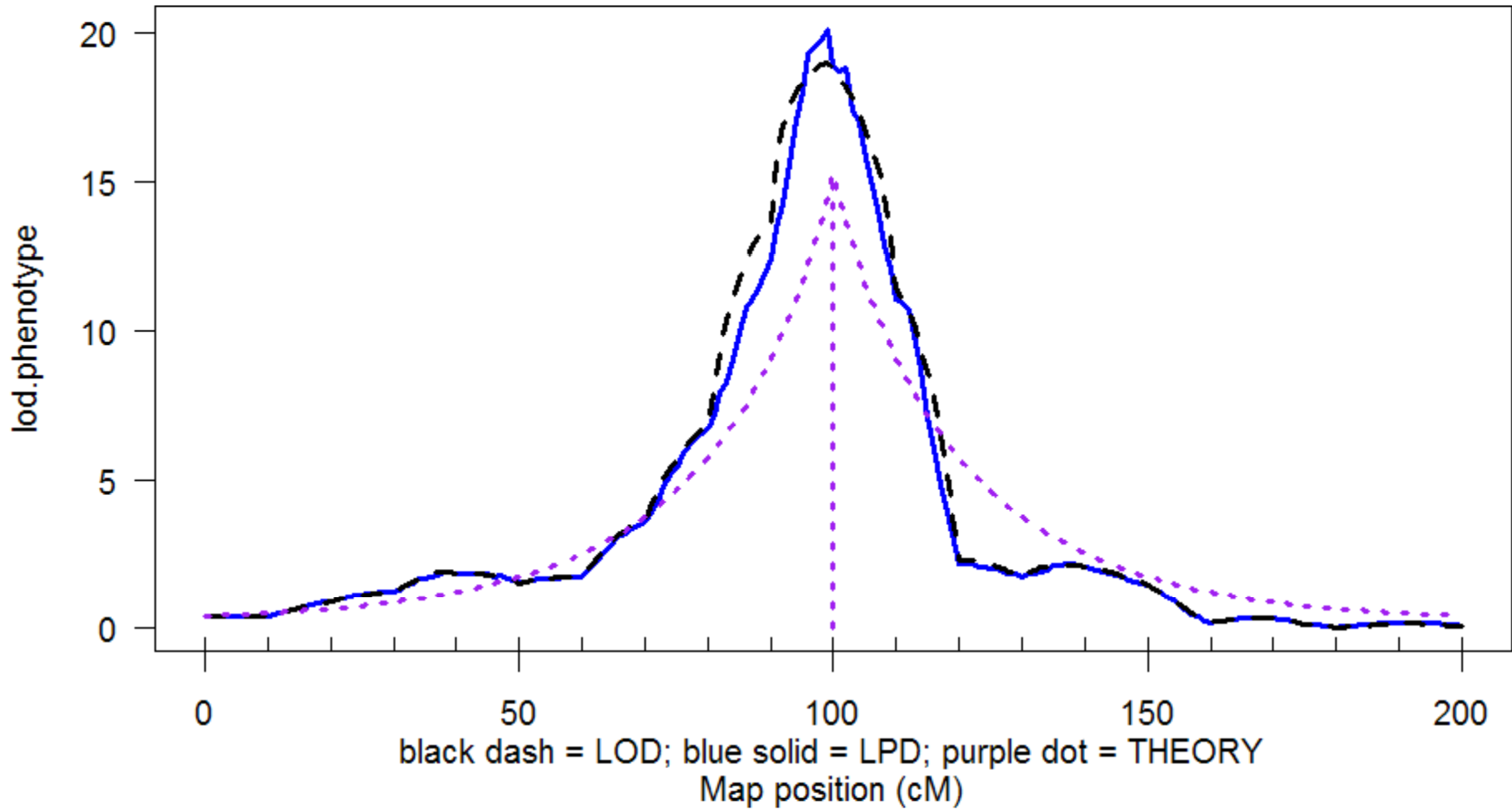
# LOD & LPD: 1 QTL

n.ind = 100, 1 cM marker spacing



# LOD & LPD: 1 QTL

n.ind = 100, 10 cM marker spacing



# marginal LOD or LPD

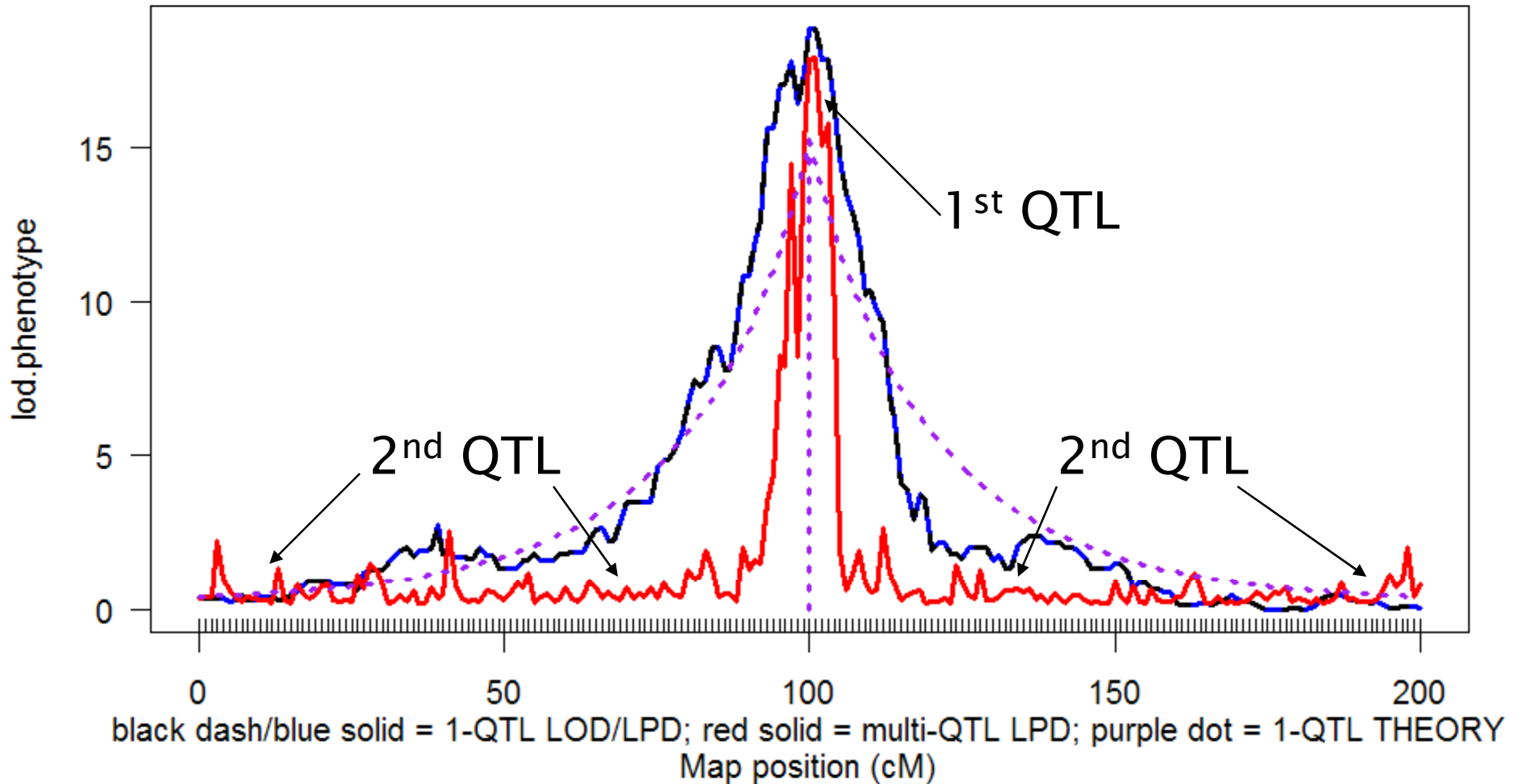
- compare two genetic architectures ( $\gamma_2, \gamma_1$ ) at each locus
  - with ( $\gamma_2$ ) or without ( $\gamma_1$ ) another QTL at locus  $\lambda$ 
    - preserve model hierarchy (e.g. drop any epistasis with QTL at  $\lambda$ )
  - with ( $\gamma_2$ ) or without ( $\gamma_1$ ) epistasis with QTL at locus  $\lambda$
  - $\gamma_2$  contains  $\gamma_1$  as a sub-architecture
- allow for multiple QTL besides locus being scanned
  - architectures  $\gamma_1$  and  $\gamma_2$  may have QTL at several other loci
  - use marginal LOD, LPD or other diagnostic
  - posterior, Bayes factor, heritability

$$\text{LOD}(\lambda | \gamma_2) - \text{LOD}(\lambda | \gamma_1)$$

$$\text{LPD}(\lambda | \gamma_2) - \text{LPD}(\lambda | \gamma_1)$$

# LPD: 1 QTL vs. multi-QTL

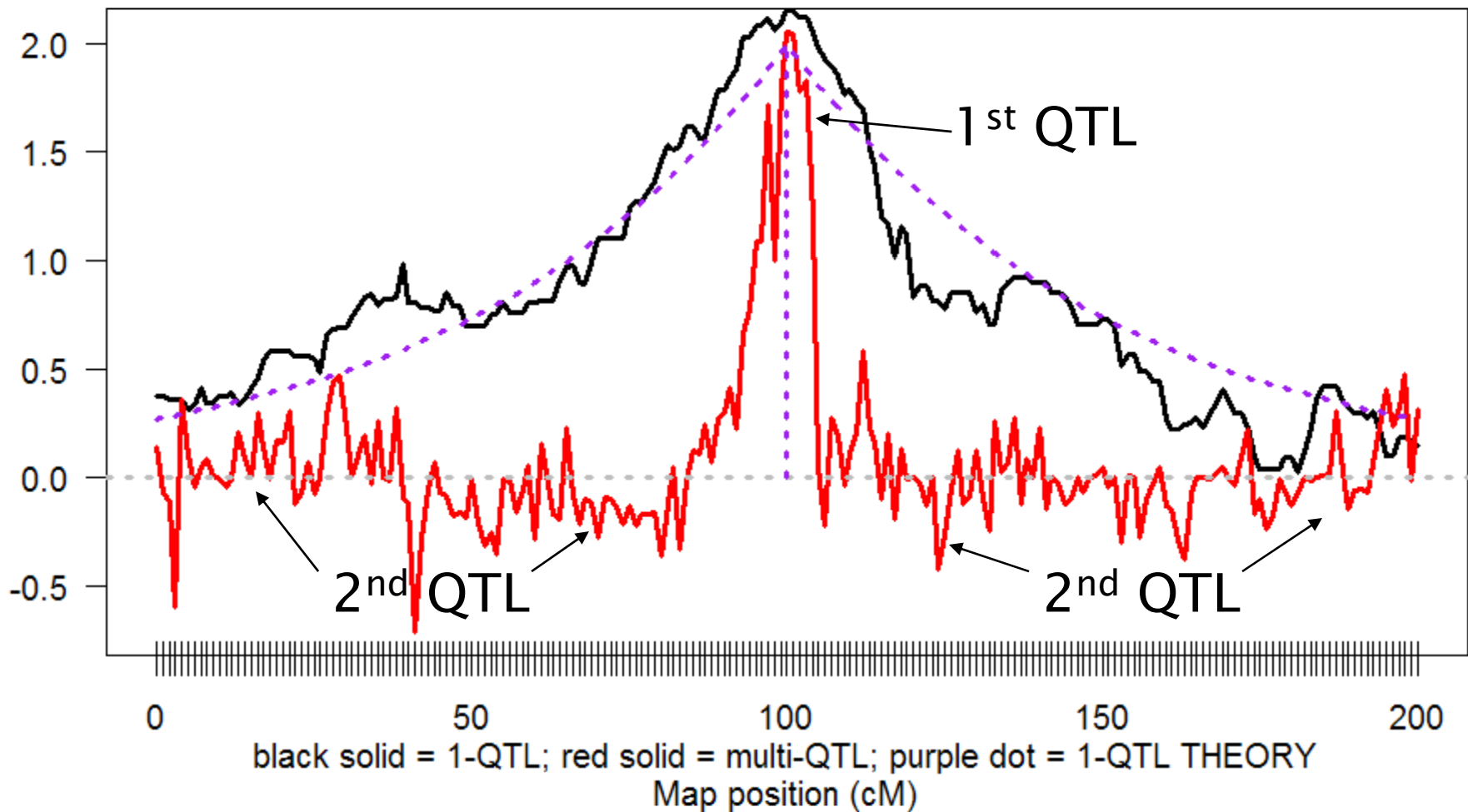
marginal contribution to LPD from QTL at  $\lambda$





# substitution effect: 1 QTL vs. multi-QTL

single QTL effect vs. marginal effect from QTL at  $\lambda$



# why use a Bayesian approach?

- first, do *both* classical and Bayesian
  - always nice to have a separate validation
  - each approach has its strengths and weaknesses
- classical approach works quite well
  - selects large effect QTL easily
  - directly builds on regression ideas for model selection
- Bayesian approach is comprehensive
  - samples most probable genetic architectures
  - formalizes model selection within one framework
  - readily (!) extends to more complicated problems

## 4. QTL model selection

- select class of models
  - see earlier slides above
- decide how to compare models
  - (Bayesian interval mapping talk later)
- search model space
  - (Bayesian interval mapping talk later)
- assess performance of procedure
  - see Kao (2000), Broman and Speed (2002)
  - Manichaukul, Moon, Yandell, Broman (in prep)
  - be wary of HK regression assessments

# pragmatics of multiple QTL

- evaluate some objective for model given data
  - classical likelihood
  - Bayesian posterior
- search over possible genetic architectures (models)
  - number and positions of loci
  - gene action: additive, dominance, epistasis
- estimate “features” of model
  - means, variances & covariances, confidence regions
  - marginal or conditional distributions
- art of model selection
  - how select “best” or “better” model(s)?
  - how to search over useful subset of possible models?

# comparing models

- balance model fit against model complexity
  - want to fit data well (maximum likelihood)
  - without getting too complicated a model

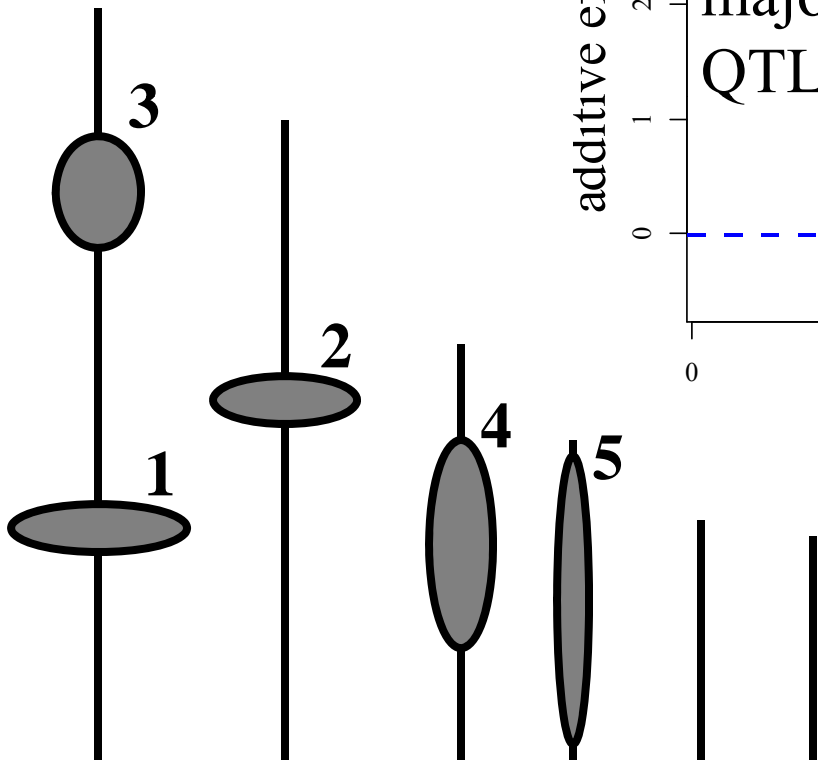
	<b>smaller model</b>	<b>bigger model</b>
<b>fit model</b>	miss key features	fits better
<b>estimate phenotype</b>	may be biased	no bias
<b>predict new data</b>	may be biased	no bias
<b>interpret model</b>	easier	more complicated
<b>estimate effects</b>	low variance	high variance

# Bayesian model averaging

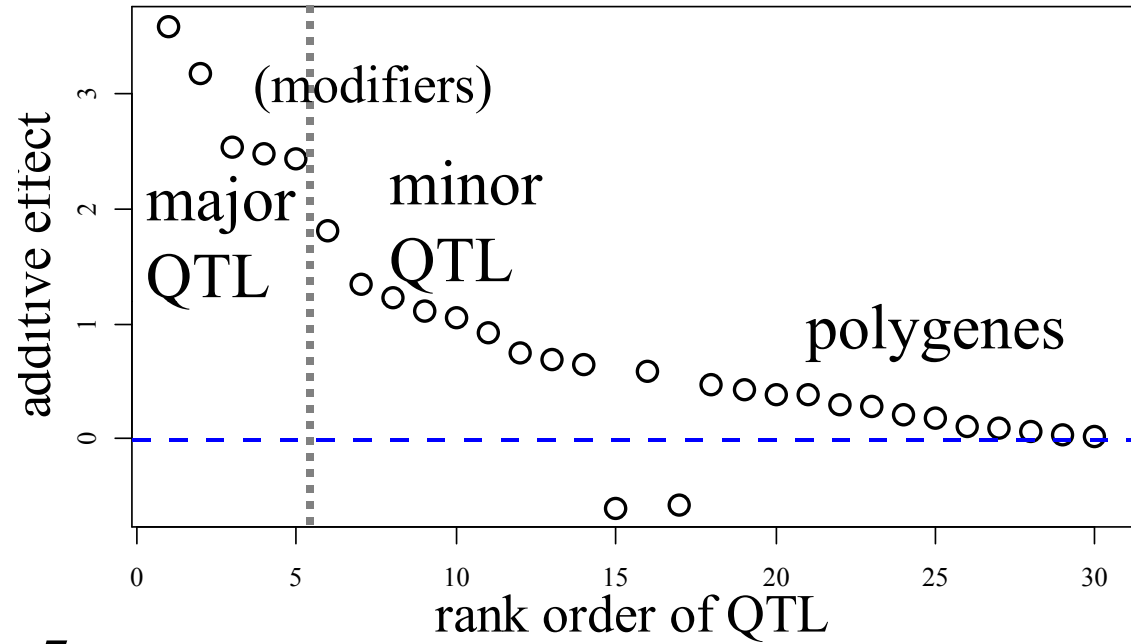
- average summaries over multiple architectures
- avoid selection of “best” model
- focus on “better” models
- examples in data talk later

# Pareto diagram of QTL effects

major QTL on linkage map



QTL 2: Overview



Seattle SISG: Yandell © 2009

# 5. QTL software options

- methods
  - approximate QTL by markers
  - exact multiple QTL interval mapping
- software platforms
  - MapMaker/QTL (obsolete)
  - QTLCart ([statgen.ncsu.edu/qtlcart](http://statgen.ncsu.edu/qtlcart))
  - R/qtl ([www.rqtl.org](http://www.rqtl.org))
  - R/qtlbim ([www.qtlbim.org](http://www.qtlbim.org))
  - Yandell, Bradbury (2007) book chapter



# approximate QTL methods

- marker regression
  - locus & effect confounded
  - lose power with missing data
- Haley-Knott (least squares) regression
  - correct mean, wrong variance
  - biased by pattern of missing data (Kao 2000)
- extended HK regression
  - correct mean and variance
  - minimizes bias issue (R/qtl “ehk” method)
- composite interval mapping (QTLCart)
  - use markers to approximate other QTL
  - properties depend on marker spacing, missing data

# exact QTL methods

- interval mapping (Lander, Botstein 1989)
  - scan whole genome for single QTL
  - bias for linked QTL, low power
- multiple interval mapping (Kao, Zeng, Teasdale 1999)
  - sequential scan of all QTL
  - stepwise model selection
- multiple imputation (Sen, Churchill 2001)
  - fill in (impute) missing genotypes along genome
  - average over multiple imputations
- Bayesian interval mapping (Yi et al. 2005)
  - sample most likely models
  - marginal scans conditional on other QTL

# QTL software platforms

- QTLCart ([statgen.ncsu.edu/qtlcart](http://statgen.ncsu.edu/qtlcart))
  - includes features of original MapMaker/QTL
    - not designed for building a linkage map
  - easy to use Windows version WinQTLCart
  - based on Lander-Botstein maximum likelihood LOD
    - extended to marker cofactors (CIM) and multiple QTL (MIM)
    - epistasis, some covariates (GxE)
    - stepwise model selection using information criteria
  - some multiple trait options
  - OK graphics
- R/qtl ([www.rqtl.org](http://www.rqtl.org))
  - includes functionality of classical interval mapping
  - many useful tools to check genotype data, build linkage maps
  - excellent graphics
  - several methods for 1-QTL and 2-QTL mapping
    - epistasis, covariates (GxE)
  - tools available for multiple QTL model selection

# Bayesian QTL software options

- Bayesian Haley-Knott approximation: no epistasis
  - Berry C (1998)
    - R/bqtl ([www.r-project.org](http://www.r-project.org) contributed package)
- multiple imputation: epistasis, mostly 1-2 QTL but some multi-QTL
  - Sen and Churchill (2000)
    - matlab/pseudomarker ([www.jax.org/staff/churchill/labsite/software](http://www.jax.org/staff/churchill/labsite/software))
  - Broman et al. (2003)
    - R/qlt ([www.rqtl.org](http://www.rqtl.org))
- Bayesian interval mapping via MCMC: no epistasis
  - Satagopan et al. (1996); Satagopan, Yandell (1996) Gaffney (2001)
    - R/bim ([www.r-project.org](http://www.r-project.org) contributed package)
    - WinQTLCart/bmapqtl ([statgen.ncsu.edu/qltcart](http://statgen.ncsu.edu/qltcart))
  - Stephens & Fisch (1998): no code release
  - Sillanpää Arjas (1998)
    - multimapper ([www.rni.helsinki.fi/~mjs](http://www.rni.helsinki.fi/~mjs))
- Bayesian interval mapping via MCMC: epistasis
  - Yandell et al. (2007)
    - R/qltbim ([www.qltbim.org](http://www.qltbim.org))
- Bayesian shrinkage: no epistasis
  - Wang et al. Xu (2005): no code release

# R/qtlbim: [www.qtlbim.org](http://www.qtlbim.org)

- Properties
  - cross-compatible with R/qtl
  - new MCMC algorithms
    - Gibbs with loci indicators; no reversible jump
  - epistasis, fixed & random covariates, GxE
  - extensive graphics
- Software history
  - initially designed (Satagopan Yandell 1996)
  - major revision and extension (Gaffney 2001)
  - R/bim to CRAN (Wu, Gaffney, Jin, Yandell 2003)
  - R/qtlbim to CRAN (Yi, Yandell et al. 2006)
- Publications
  - Yi et al. (2005); Yandell et al. (2007); ...

# many thanks

## U AL Birmingham

Nengjun Yi

Tapan Mehta

Samprit Banerjee

Daniel Shriner

Ram Venkataraman

David Allison

## Jackson Labs

Gary Churchill

Hao Wu

Hyuna Yang

Randy von Smith

## Alan Attie

Jonathan Stoehr

Hong Lan

Susie Clee

Jessica Byers

Mark Gray-Keller

## Tom Osborn

David Butruille

Marcio Ferrera

Josh Udahl

Pablo Quijada

## UW-Madison Stats

### Yandell lab

Jaya Satagopan

Fei Zou

Patrick Gaffney

Chunfang Jin

Elias Chaibub

W Whipple Neely

Jee Young Moon

Elias Chaibub

### Michael Newton

Karl Broman

Christina Kendziorski

Daniel Gianola

Liang Li

Daniel Sorensen

USDA Hatch, NIH/NIDDK (Attie), NIH/R01s (Yi, Broman)