

# Expression QTLs and Mapping of Complex Trait Loci

Paul Schliekelman  
Statistics Department  
University of Georgia

# Definitions: Genes, Loci and Alleles

- A gene codes for a protein. Proteins do everything.
- A gene position on a chromosome is called a locus.
- The alternate forms of the gene at a locus are called alleles.

# Definitions, continued:

- A causative locus for a trait is a locus with DNA sequence variation between individuals that contributes to variation in the trait between them.
- A quantitative trait locus (QTL) is a causative locus for a quantitative trait (e.g. height, IQ).

# Definitions: Genetic Linkage

- Genes on the same chromosome tend to be inherited together. This is called linkage.
- The closer together two genes lie, the more often they are inherited together.

# Genetic Markers

- A genetic marker is a known and detectable variation in a gene or other stretch of DNA.
- The presence of a marker can be detected in an individual organism's genome and its inheritance can be followed.
- A marker that is linked to a locus affecting a trait will appear more often in individuals with that trait.

# Gene Mapping

- Goal of gene mapping is to locate genetic loci that are responsible for variation in traits of interest (e.g. occurrence of disease, intelligence).
- Based on searching for correlations between markers and trait of interest.

# Gene Mapping Successes

- Mendelian traits have inheritance patterns consistent with a single causative locus.
- Hundreds of genes affecting known Mendelian traits (e.g. cystic fibrosis, breast cancer) in humans have been mapped.
- Some cases have led to new treatments and/or screening methods.

# Gene Mapping Failures

- Complex traits do not follow simple one-locus Mendelian expectations.
- Assumed to be caused by mutations at multiple loci.
- Examples include many common diseases: asthma, bipolar disorder, diabetes, prostate cancer, etc.



# Mapping of complex disease genes – little success

Altmuller *et al*, 2001: Reviewed 101 full genome scans of 31 complex diseases.

- 67% did not show significant linkage to any marker.
- of the significant linkages, very few have been reproduced.

# Low power for mapping complex trait loci

- Power to detect genes decreases as the number of loci affecting the trait increases.
- By definition, complex diseases should be harder to map.
- We don't know how many genes underlie complex diseases, so don't know whether to be surprised by lack of success.

The central problem of human disease genetics is to solving this dilemma.

# New Buzzword:

Genetical Genomics: the combination of molecular marker data and genome-wide expression data to elucidate the genetics of complex traits

# Expression QTLs (eQTLs)

- Gene expression levels can be treated as quantitative traits and their quantitative trait loci (QTLs) mapped.
- Transcript: gene whose expression level is being measured.
- An eQTL for that transcript refers to a genetic locus with DNA sequence variation causing variation in the expression level.

# Expression QTLs (eQTLs)

- The eQTL could be at the transcript locus itself, or elsewhere in the genome.
- By combining microarrays and marker data, eQTLs can be simultaneously generated for thousands of gene expression levels

# Recent eQTL surveys

- Numerous recent eQTL surveys have found that expression levels frequently exhibit intermediate to high heritabilities.
- E.g. Schadt *et al*/2003 found 4339 eQTLs over 3701 genes with log-of-odd scores  $>4.3$  and 11,021 genes with eQTL LOD scores  $>3.0$

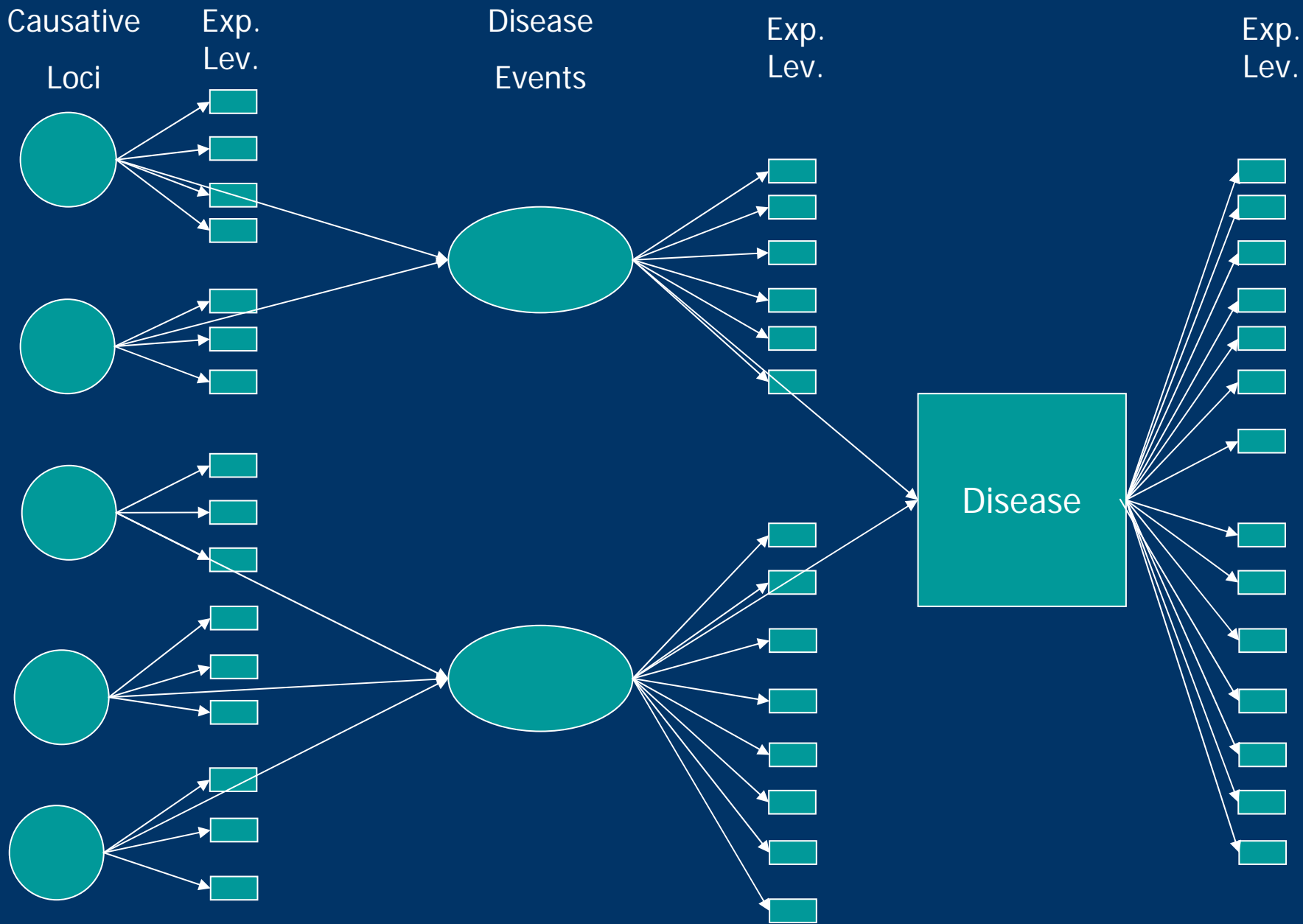
# Recent eQTL surveys

- Thus, it is relatively easy to map QTLs for many expression levels.



# Will eQTLs be useful for dissecting complex physiological traits (e.g. disease)?

- Presumably, disease causative loci are also eQTLs for some gene expression levels.
- Two task must be accomplished:
  1. Show that transcript is related to disease.
  2. Map eQTL for transcript. Then this eQTL should be causative locus for disease.



# Model for Gene Expression

Assume

- $L$  disease causative loci.
- $M$  expression levels ( $X_1 \dots X_M$ ) tested.
- Each expression level  $X_i$  depends on some subset of the  $L$  causative loci (often empty).
- $T$  = indicator variable for disease status.

# Model for Gene Expression

$$P(X | T = 1) = \frac{P(T = 1 | X)P(X)}{P(T = 1)}$$
$$= \frac{\sum_G P(T = 1 | G)P(X | G)P(G)}{K}$$

$K = \textit{Disease prevalence}$

# Expression Levels

- Assume expression levels are normally distributed on some scale (e.g. log scale).
- Mean and variance determined by genotype at one or more of the  $L$  causative loci.
- Expression level distribution is then a mixture of normal distributions with weights determined by genotype distribution conditioned on disease status.

# Multiplicative Model

$$P(T = 1 | G) = P(T_1 = 1 | G_1)P(T_2 = 1 | G_2) \dots P(T_L = 1 | G_L)$$

Can show:

$$P(X | T = 1) = \frac{\sum_{g_1, g_2, \dots, g_c} P(X | g_1, g_2, \dots, g_c) u_1(g_1) u_2(g_2) \dots u_c(g_c) P(g_1) P(g_2) \dots P(g_c)}{K_1 K_2 \dots K_c}$$

$$P(X | T = 0) =$$

$$\frac{\sum_{g_1, g_2, \dots, g_c} P(X | g_1, g_2, \dots, g_c) (1 - u_1(G_1) u_2(G_2) \dots u_c(G_c) K_{c+1} \dots K_L) P(g_1) \dots P(g_c)}{1 - K}$$

Where  $K_i = \sum_{G_i} P(T_i = 1 | G_i) P(G_i)$  and  $u_i(G_i) = P(T_i = 1 | G_i)$

If  $c=1$  (expression controlled by one locus)  
and Hardy-Weinberg Equilibrium (alleles independent):

$$\begin{aligned} P(X | T = 1) &= (1 - p_1)^2 u_1(dd) \phi(X | \mu_{dd}, \sigma_{dd}) \\ &\quad + 2 p_1 (1 - p_1) u_1(Dd) \phi(X | \mu_{Dd}, \sigma_{Dd}) \\ &\quad + p_1^2 u_1(DD) \phi(X | \mu_{DD}, \sigma_{DD}) \end{aligned}$$

Where  $p_1$  is the allele frequency

$D$ =disease allele and  $d$ ="normal" allele

$\phi(X | \mu_{DD}, \sigma_{DD})$  is normal pdf

# Additive Model

$$P(T = 1 | G) = P(T_1 = 1 | G_1) + P(T_2 = 1 | G_2) + \dots + P(T_L = 1 | G_L)$$

Can show:

$$P(X | T = 1) = \frac{\sum_{G_1} P(X | G_1)(u_1(G_1) + K_2 + \dots + K_L)P(G_1)}{K}$$

$$P(X | T = 0) = \frac{\sum_{G_1} P(X | G_1)(1 - u_1(G_1) - K_2 - K_3 - \dots - K_L)P(G_1)}{1 - K}$$



# Expression Level Mean

- Expression level mean assumed to have either multiplicative or additive dependence on genotype:

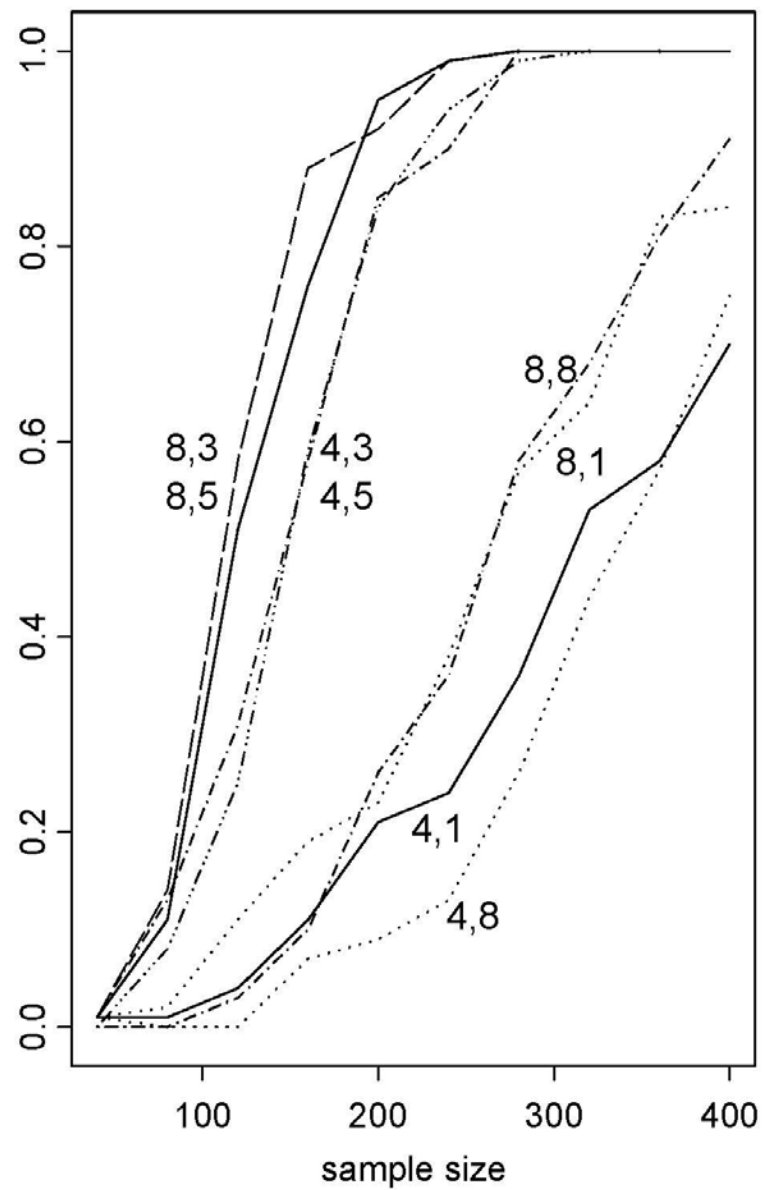
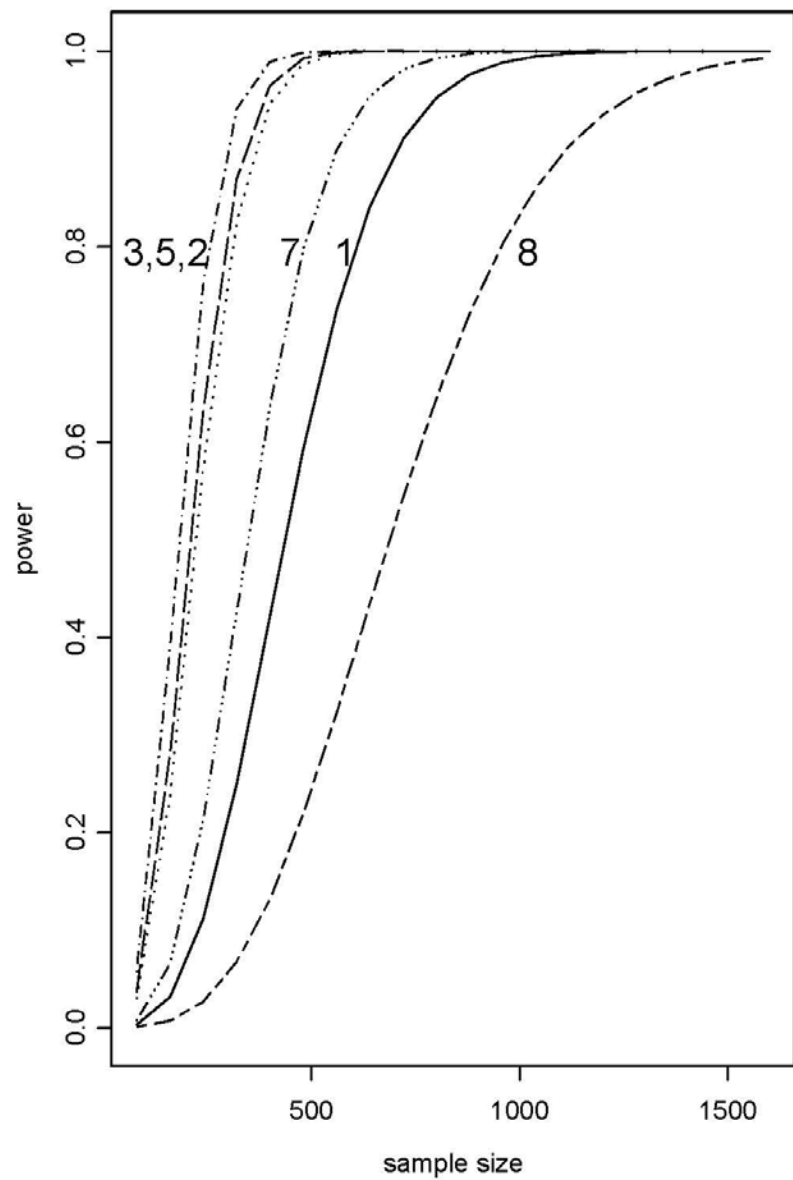
$$\mu(g_1, \dots, g_c) = \mu_1(g_1) \mu_2(g_2) \dots \mu(g_c)$$

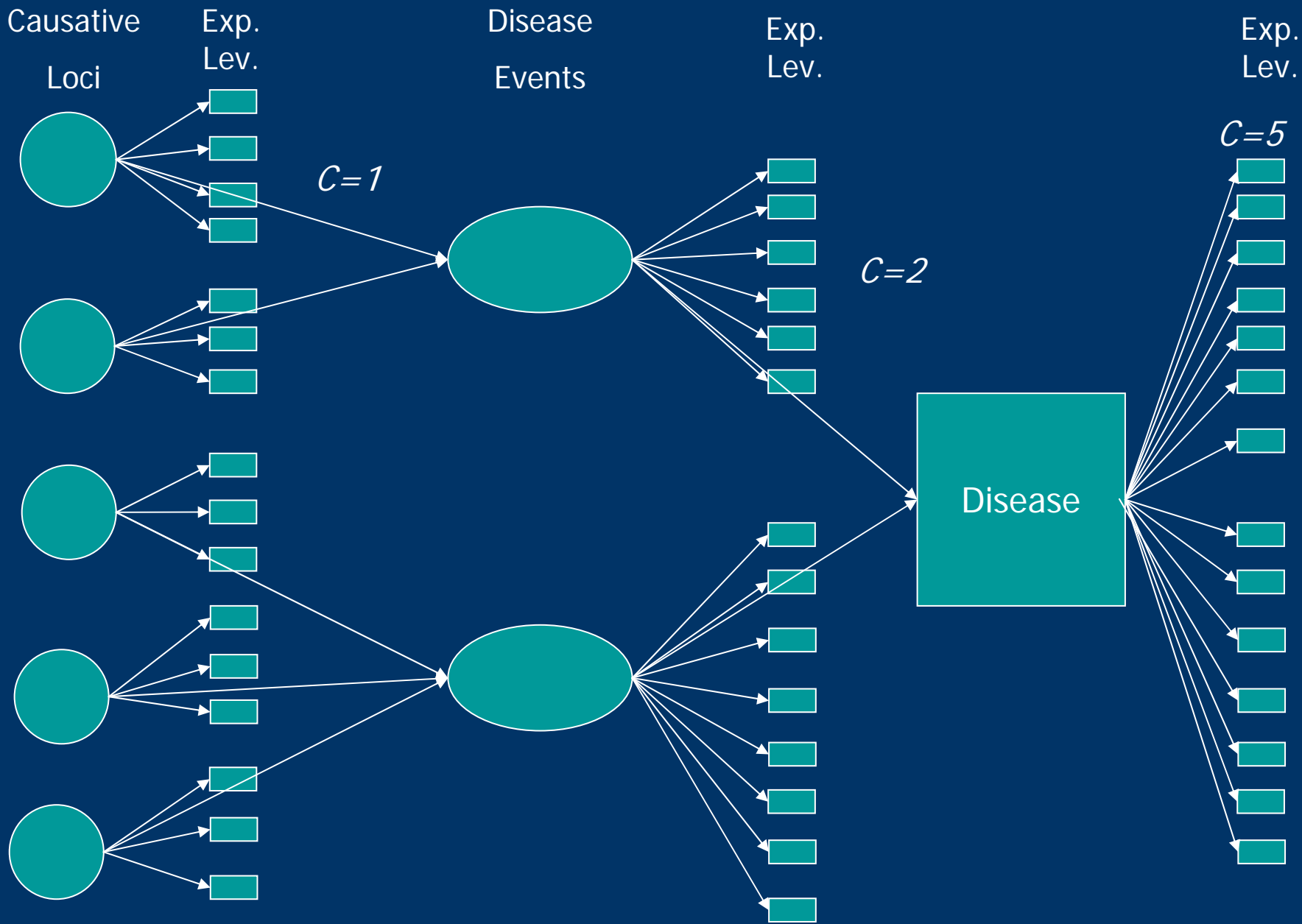
*or*

$$\mu(g_1, \dots, g_c) = \mu_1(g_1) + \mu_2(g_2) + \dots + \mu(g_c)$$

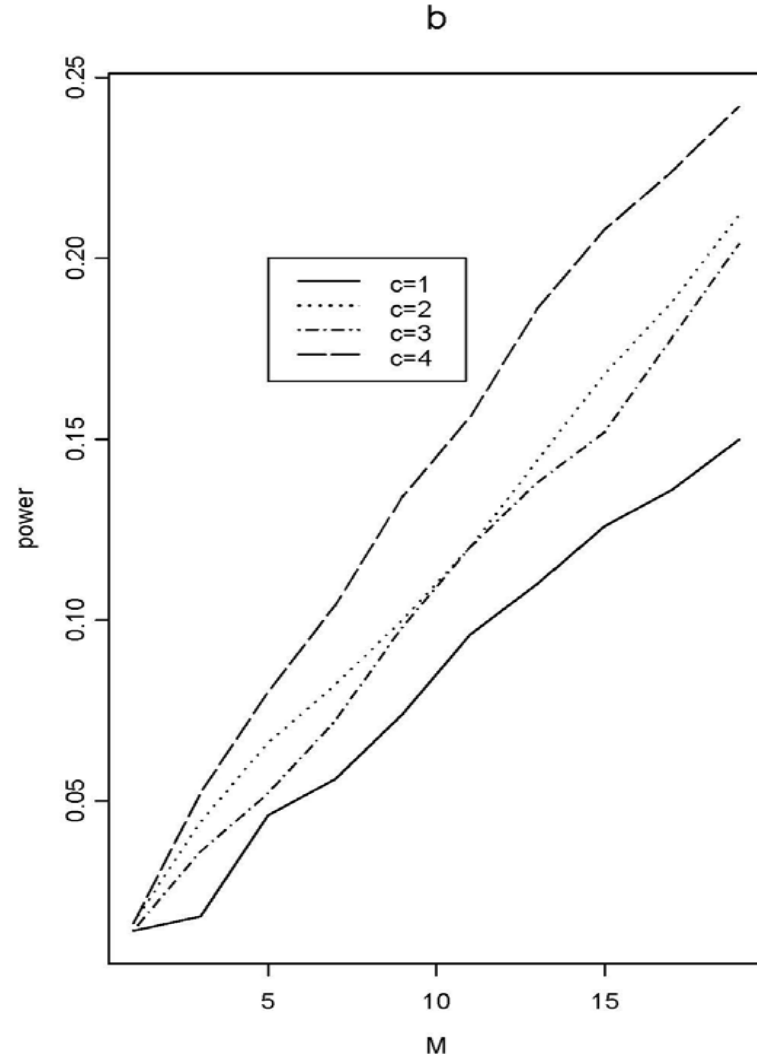
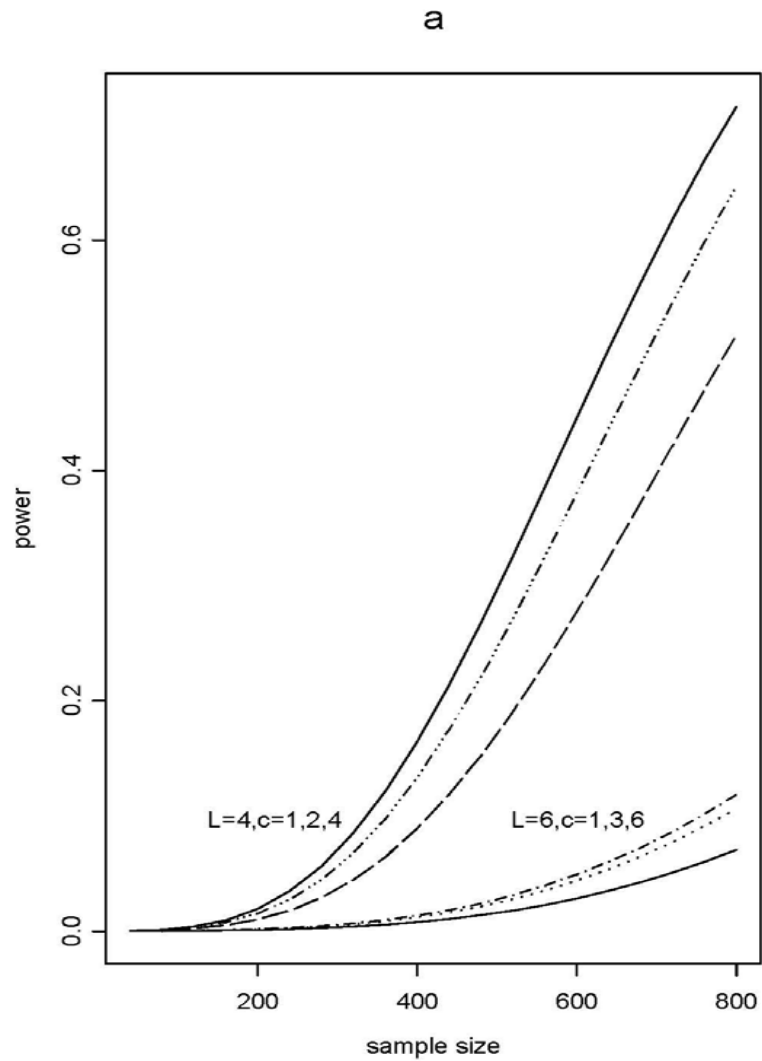
# Power Calculations

- Two groups: disease affected and unaffected.
- t-test for group difference conducted for each gene.
- Calculate power to detect differential expression.
- Interested in relationship between power and genetic model.





# Additive Model



# Detecting expression level differences: Conclusions

- Power to detect expression level differences is poor for a multiplicative model with  $c=1$ , but reasonable for  $c=2-5$ .
- Power to detect expression level differences is very poor for an additive model.

# Power for Mapping eQTLs

Power for mapping QTLs in natural populations deteriorates quickly as the number of QTLs increases. Power poor even for  $c=2$ .

# Conclusions – Natural Populations

- Power to detect expression level differences extremely poor if disease probability is additive.
- For multiplicative model there is no value of  $c$  (number of controlling loci) where power is good for detecting both linkage and expression level differences.



# Many Unknowns

- Distribution of relationships between expression level and causative loci.
- Dominance relationships.
- Form of gene interactions.
- Allele frequencies.

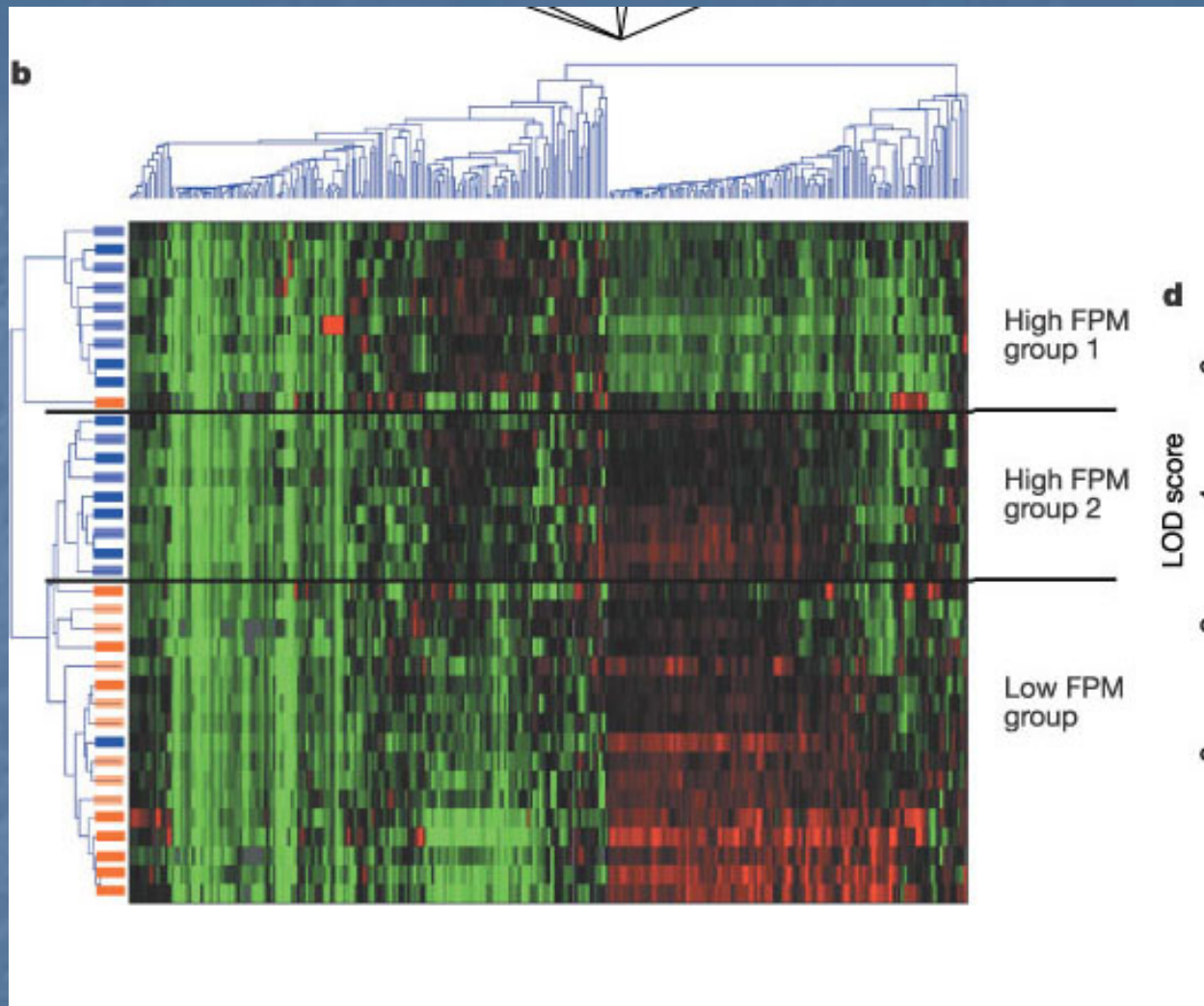
# Much Potential For Improved Power

- There is much information in the joint behavior of expression levels.
- Better experimental designs are possible (e.g. utilize family structure).

# How can we use joint expression information?

Schadt et al 2003:

- Crossed two mouse strains: one susceptible to diabetes and weight gain on high fat diet and the other not.
- 111 offspring from cross were put on high fat diet for 4 months.
- Liver tissue from each was profiled using microarrays.
- Genotypes for several hundred markers were obtained for each mouse.



From Schadt *et al.*, 2003 Genetics of gene expression surveyed in maize, mouse and man. Nature 422: 297-302.

# Results:

- Microarray data splits high FPM data into two groups.
- Gene mapping comparing separate high-FPM groups to low-FPM produced significant linkage for a obesity locus.
- Gene mapping comparing combined high-FPM groups to low-FPM groups do not find significant linkage.

# Genetic Heterogeneity

	Affected Individuals													
Locus	1	2	3	4	5	6	7	8	9	10	11	12	13	14
1		x				x								x
2				x				x	x			x		
3					x		x							
4			x								x			
5				x	x			x						
6										x				x

Major variation in which disease mutations are carried by affected individuals



Goal: A method for identifying the genetic heterogeneity using genomewide expression data

Gene mapping power could be dramatically improved if this heterogeneity could be accounted for



# Proposed EM Algorithm Based Method:

Data:

- $N_A$  affected and  $N_U$  unaffected individuals.
- Microarray for each individual
- Disease status for each individual

		Transcript					
Sample Individ.	Disease Status	1	2	3	4	5	6
1	$Y_1$	$X_{11}$	$X_{12}$	$X_{13}$	$X_{14}$	$X_{15}$	$X_{16}$
2	$Y_2$	$X_{21}$	$X_{22}$	$X_{23}$	$X_{24}$	$X_{25}$	$X_{26}$
3	$Y_3$	$X_{31}$	$X_{32}$	$X_{33}$	$X_{34}$	$X_{35}$	$X_{36}$
4	$Y_4$	$X_{41}$	$X_{42}$	$X_{43}$	$X_{44}$	$X_{45}$	$X_{46}$
5	$Y_{-5}$	$X_{51}$	$X_{52}$	$X_{53}$	$X_{54}$	$X_{55}$	$X_{56}$

# Unobserved Variables

$Z_i$ =controlling locus for transcript  $i$ ,  $i=1$  to  $M$

$V_{jk}$ =genotype on controlling locus  $j$  in individual  $k$ .

Controlling loci 1 to  $L$  are the disease loci.

Controlling loci  $L+1$  to  $L+c$  are "null loci" that do not affect disease.

More specific goal:

Cluster sample individuals by their values of  $V_{jk}$  (that is, their genotype on the disease loci. ).

# Likelihood

$$L(\vec{\theta} | \vec{X}, \vec{Y}, \vec{Z}, \vec{V}) =$$

$$P(\vec{X} | \vec{Z}, \vec{V}, \vec{\theta}) P(\vec{Y} | \vec{V}, \vec{\theta}) P(\vec{Z}, \vec{V} | \vec{\theta})$$

$$= \prod_{i=1}^N \prod_{j=1}^M P(X_{ij} | Z_j, V_{iZ_j}, \vec{\theta}) \prod_{k=1}^N P(Y_k | \vec{V}_k, \vec{\theta}) P(\vec{Z}, \vec{V} | \vec{\theta})$$

$P\left(X_{ij} \mid Z_j, V_{iZ_j}, \vec{\theta}\right) = \text{Expression level distn. for gene } j$   
*in indiv. } i*

$P\left(Y_k \mid \vec{V}_k, \vec{\theta}\right) = \text{Disease probability for indiv. } k$

$P\left(\vec{Z}, \vec{V} \mid \vec{\theta}\right) = \text{Controlling locus and genotype probs.}$

# Parameter Estimates:

Expression level mean for transcript  $a$  with controlling locus genotype  $b$ :

$$\hat{\mu}_{ab} = \frac{\sum_{t=1}^L \sum_{i=1}^N x_{ia} P\left(Z_a = t, V_{it} = b \mid \vec{X}, \vec{Y}, \vec{\theta}_0\right)}{\sum_{t=1}^L \sum_{i=1}^N P\left(Z_a = t, V_{it} = b \mid \vec{X}, \vec{Y}, \vec{\theta}_0\right)}$$

Expression level variance for transcript  $a$  with controlling locus genotype  $b$ :

$$\hat{\sigma}_{ab}^2 = \frac{\sum_{t=1}^L \sum_{i=1}^N (x_{ia} - \mu_{ab})^2 P(Z_a = t, V_{it} = b | \bar{X}, \bar{Y}, \bar{\theta}_0)}{\sum_{t=1}^L \sum_{i=1}^N P(Z_a = t, V_{it} = b | \bar{X}, \bar{Y}, \bar{\theta}_0)}$$



Probability that the expression of gene  $r$  has controlling locus  $q$ )

$$\hat{\alpha}_{rq} = P\left(Z_r = q \mid \vec{X}, \vec{Y}, \vec{\theta}_0\right)$$

Probability that controlling locus  $s$   
has genotype  $t$

$$\hat{\lambda}_{st} = \frac{\sum_{j=1}^N P(V_{js} = t \mid \vec{X}, \vec{Y}, \vec{\theta}_0)}{N}$$

# Disease penetrance parameters

The genotype specific risk parameters are found by numerical maximization.

# Assigning Genotypes:

$$P(v_{ij} = k | Y_i, \vec{W}_{ij}) = \frac{P(Y_i | v_{ij} = k) P(\vec{W}_{ij} | v_{ij} = k)}{P(Y_i, \vec{W}_{ij})}$$

$\vec{W}_{ij}$  = vector containing the values of expression levels in individual  $i$  that are controlled by locus  $j$ .

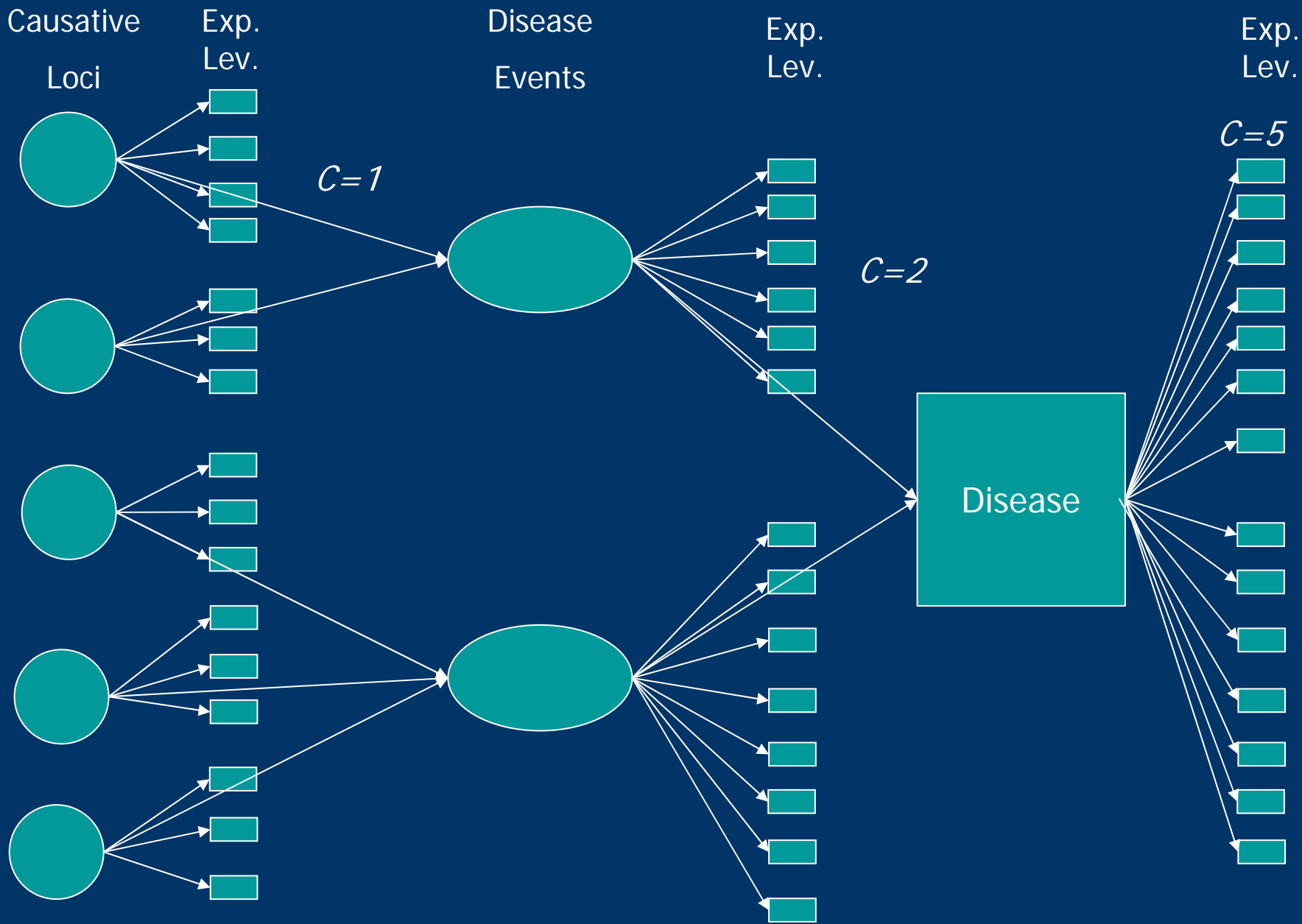
- Transcripts assigned to same controlling locus if their expression is correlated in a way that is consistent with genetic model.
- Sample individuals assigned to same genotype at controlling locus if controlled expression levels are similar.

# How well will it work?

- Potentially much information about genotype available in expression levels.

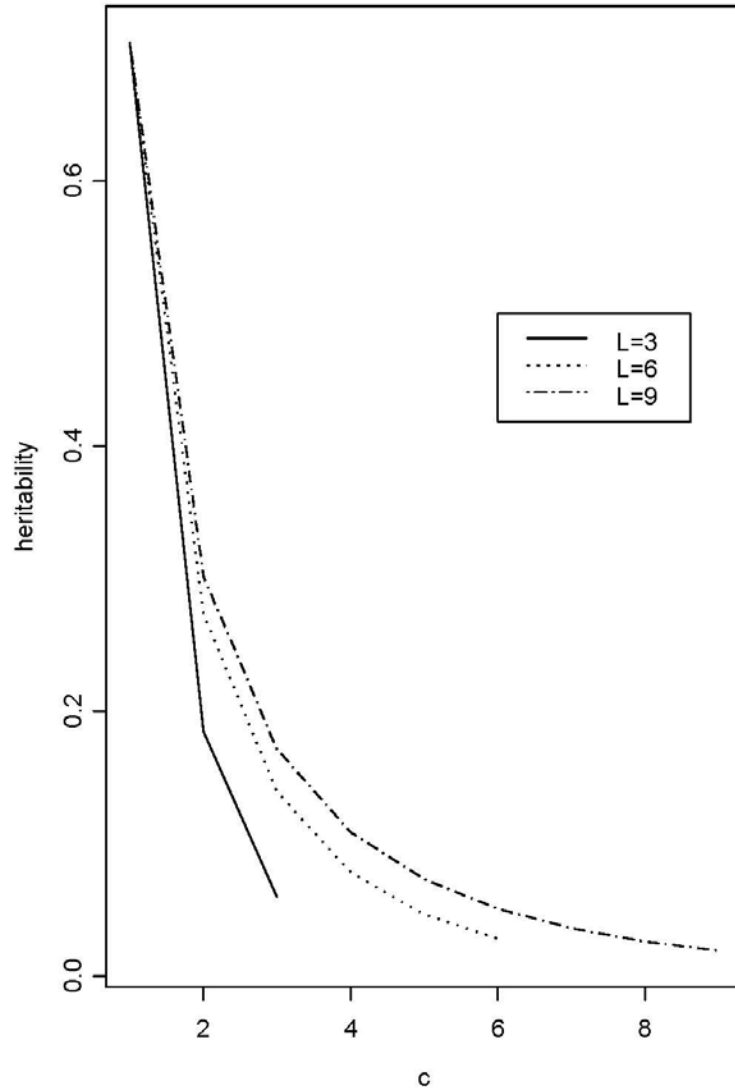
On the other hand:

- Highly assumption laden.
- Many parameters to estimate.
- High-level transcripts may overwhelm analysis.



# Mapping Power - multiplicative

a



b

