

Topics in Statistical Genetics

Karl W Broman

Department of Biostatistics
Johns Hopkins School of Public Health

Past

- QTLs in experimental crosses
- CEPH data $\left\{ \begin{array}{l} 8 \text{ families } (\sim 130 \text{ people}) \\ >8000 \text{ STRPs} \end{array} \right.$
 - Genetic maps
 - Recombinational variation
 - Autozygosity
 - Crossover interference
- Genotyping/pedigree errors
- Dog genetic maps
- Various disease projects

Traditional genetics: **binary** traits

e.g. disease or no disease

Quantitative traits

e.g. yield of tomato crop

bristles on a fly's abdomen

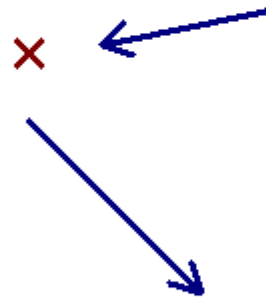
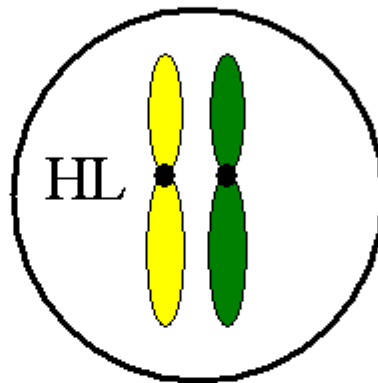
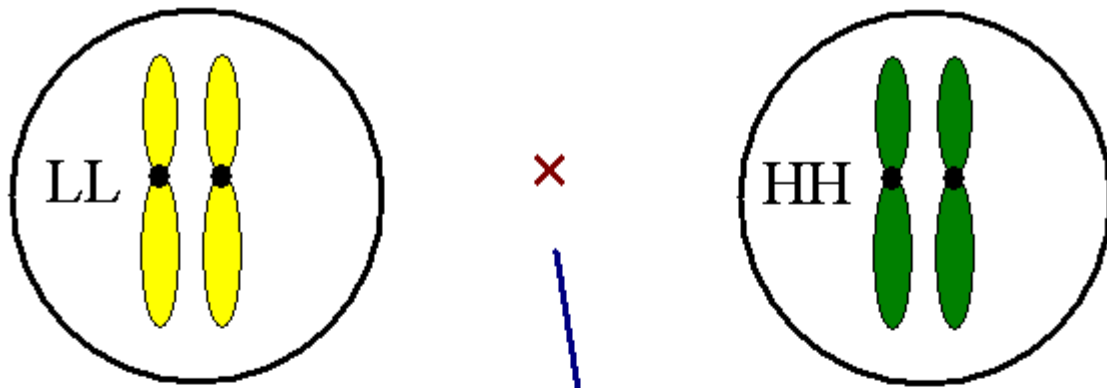
- Many genes
- Environmental variation

Why?

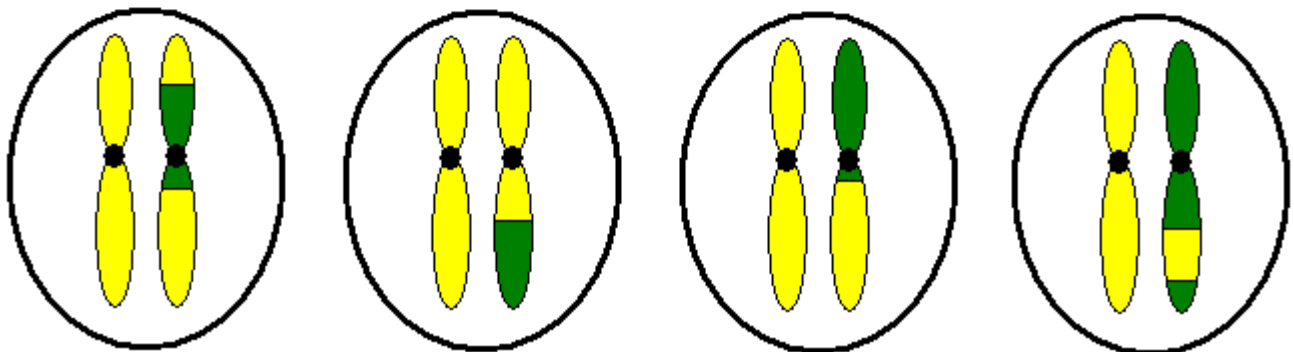
- Biochemical basis of trait
- Selection experiments
- Evolution

Goal: find (some of) the genes

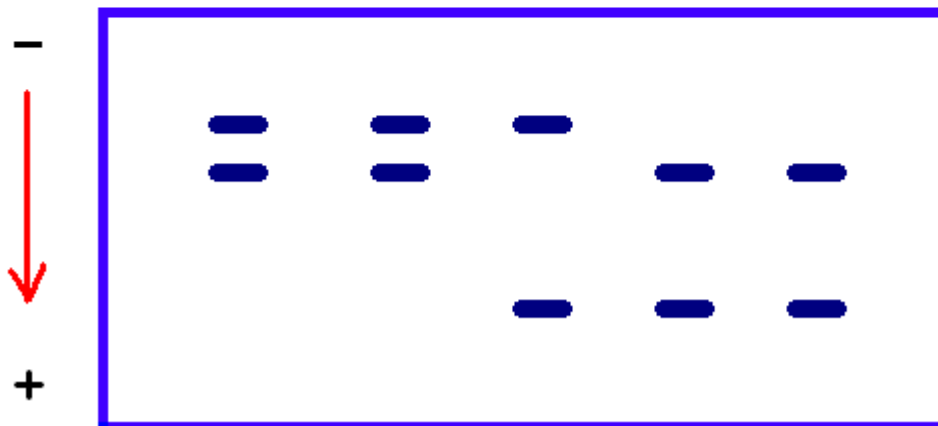
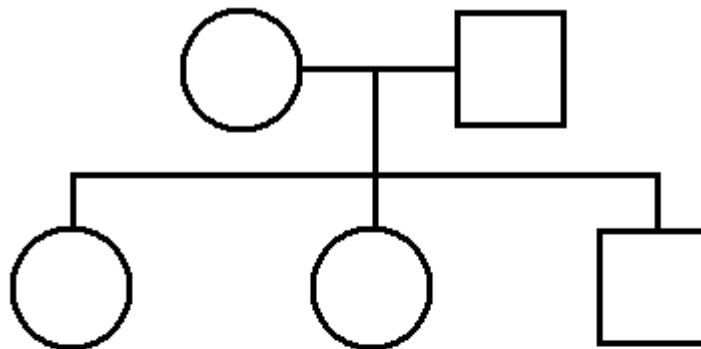
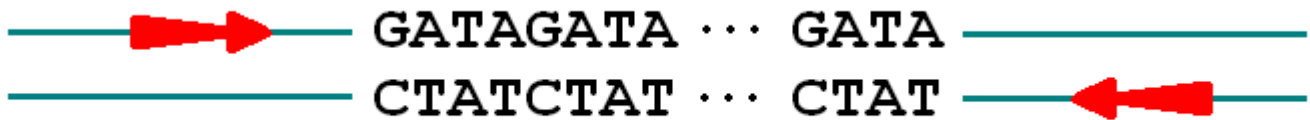
Backcross experiment



Backcross progeny
(LL or HL)



Genetic markers: STRPs or microsatellites



Data

Phenotypes (trait values)

y_i = phenotype for individual i

Marker genotypes

$x_{ij} = 1/0$ if i is HL/LL at marker j

Genetic map

Locations of markers

Models

Recombination: No interference

Phenotype/genotype connection

$$y = \mu + \sum \beta_j z_j + \varepsilon$$

$$\varepsilon \sim \text{Normal}(0, \sigma^2)$$

Problem

100 to 1000 backcross progeny

100 to 400 markers

$$y = \mu + \sum \beta_j x_j + \varepsilon$$

Find the x 's with $\beta_j \neq 0$

Errors:

- Miss important loci
- Include extraneous loci

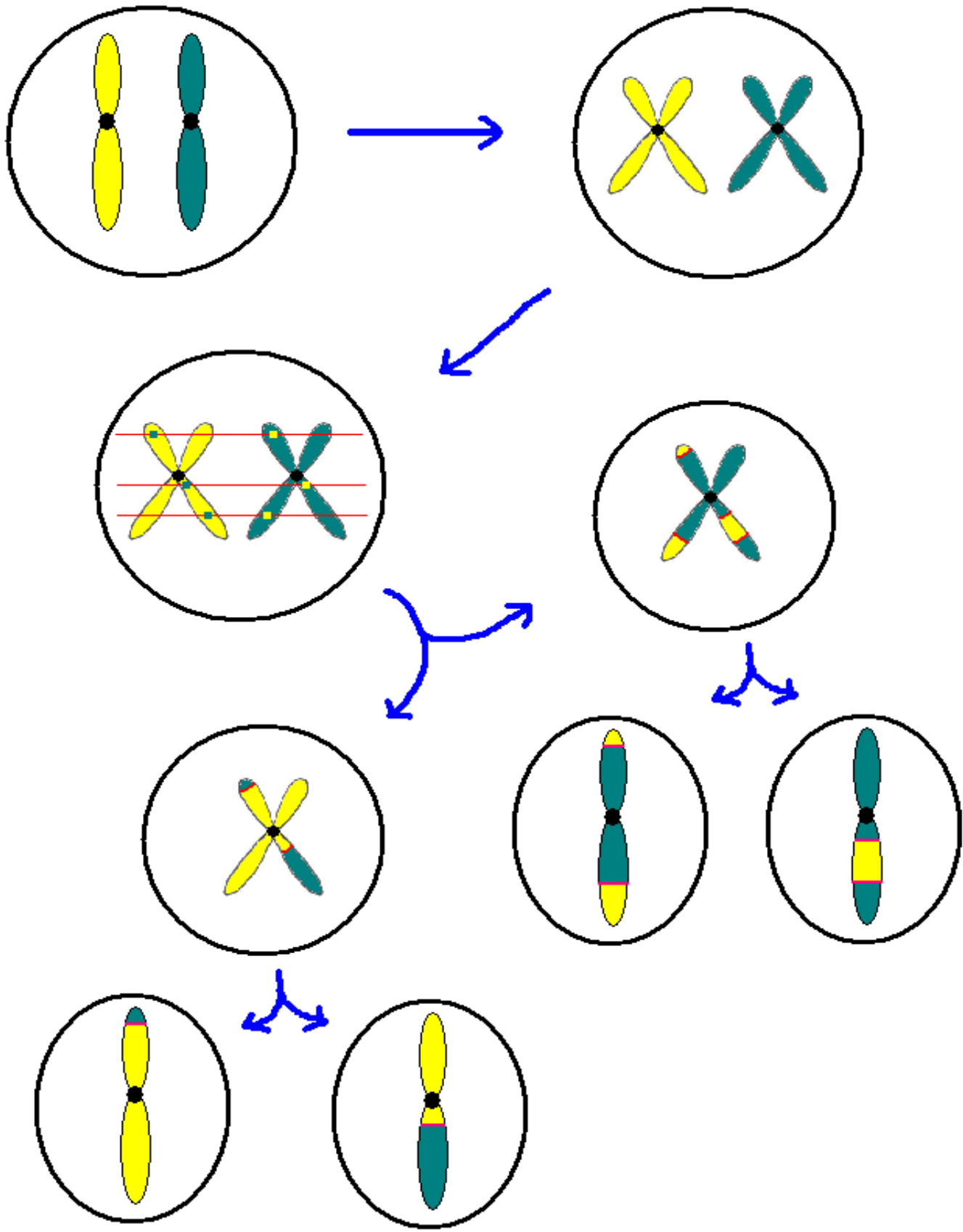
Discussion

- Identifying QTLs is model selection
- Simulation studies are necessary
 - compare procedures
 - understand a procedure's performance
- Different situations will require different procedures

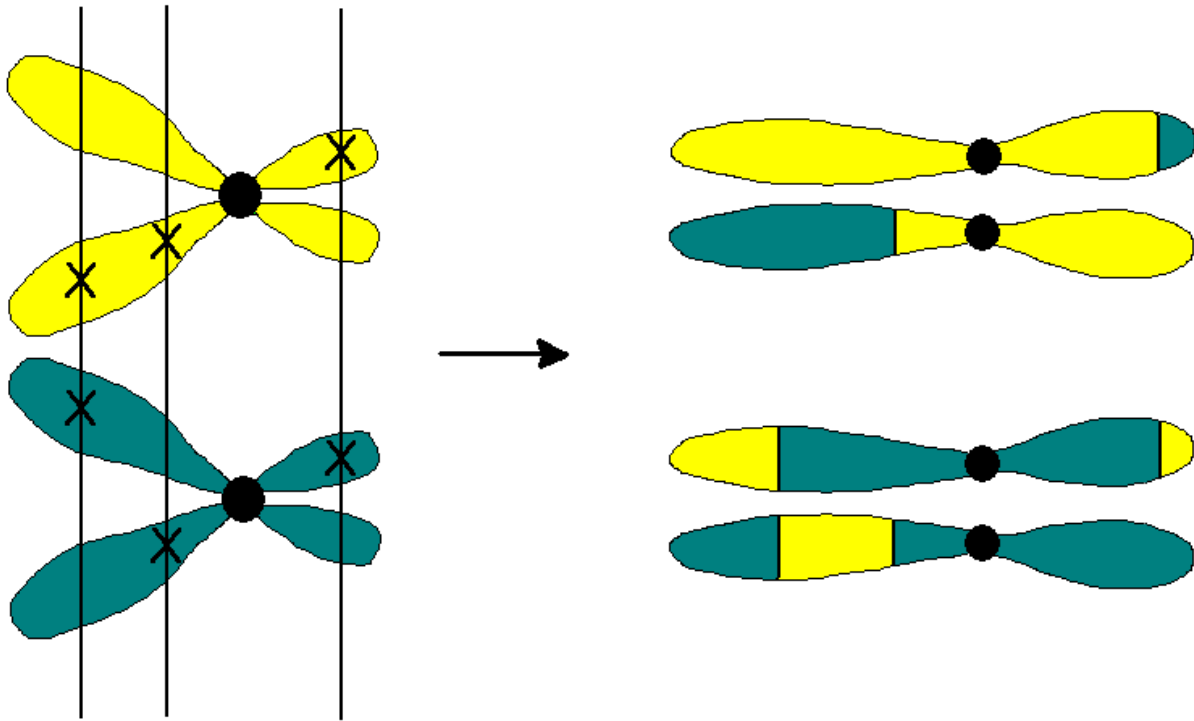
Human data

- www.marshmed.org/genetics
- 8 CEPH families
 - three generations
 - 11 to 15 progeny
 - 92 meioses
- ~8,000 STRP markers
 - 90 ± 7 % typed
- Average spacing
 - female: 0.6 ± 1.2 cM
 - male: 0.4 ± 1.0 cM
 - sex-ave: 0.5 ± 0.9 cM
- Data cleaning
 - Removed 764/964,425 (~0.08%) genotypes resulting in tight double recombinants

Meiosis



Interference



- Strand choice
→ Chromatid interference
- Spacing
→ Chiasma (crossover) interference

Genetic distance

distance (cM) = average # crossovers
in 100 meiotic products

per Morgan $\left\{ \begin{array}{l} 2 \text{ chiasmata on 4-strand bundle} \\ 1 \text{ crossover on meiotic product} \end{array} \right.$

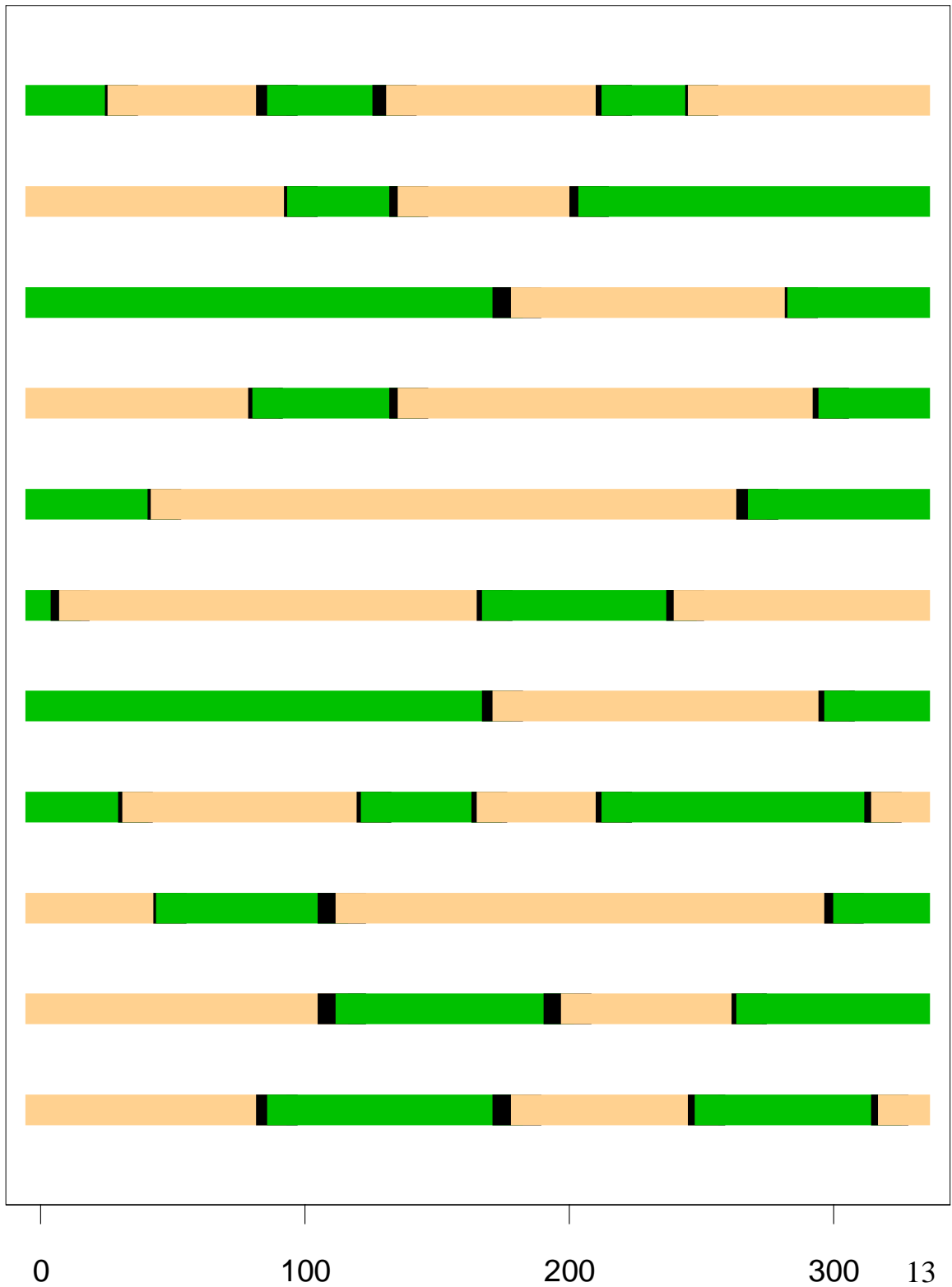
Map function

recombination fraction as a function of genetic distance

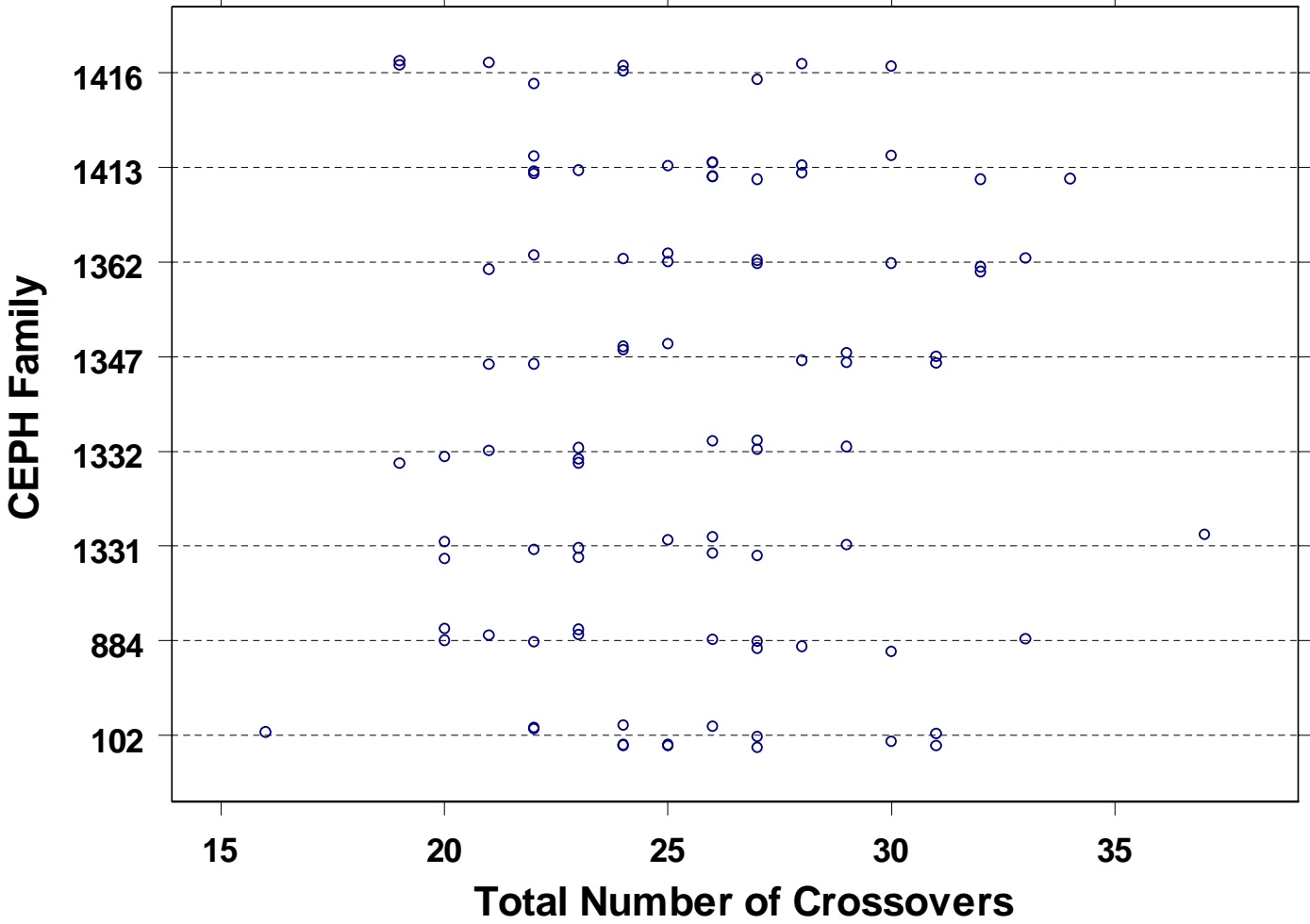
Haldane $r(d) = \frac{1}{2} [1 - \exp(-2d)]$

Kosambi $r(d) = \frac{1}{2} \tanh(2d)$

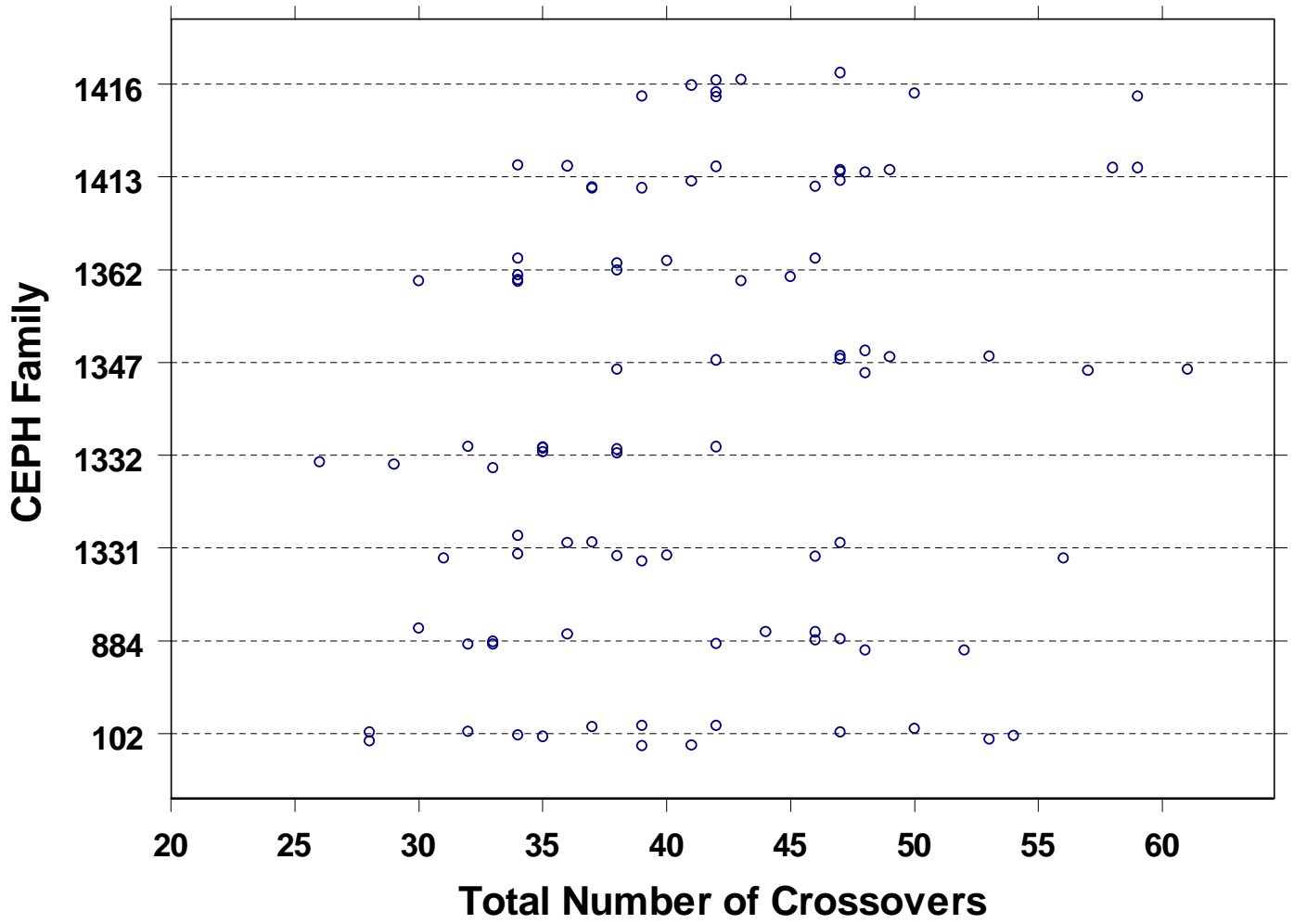
The data



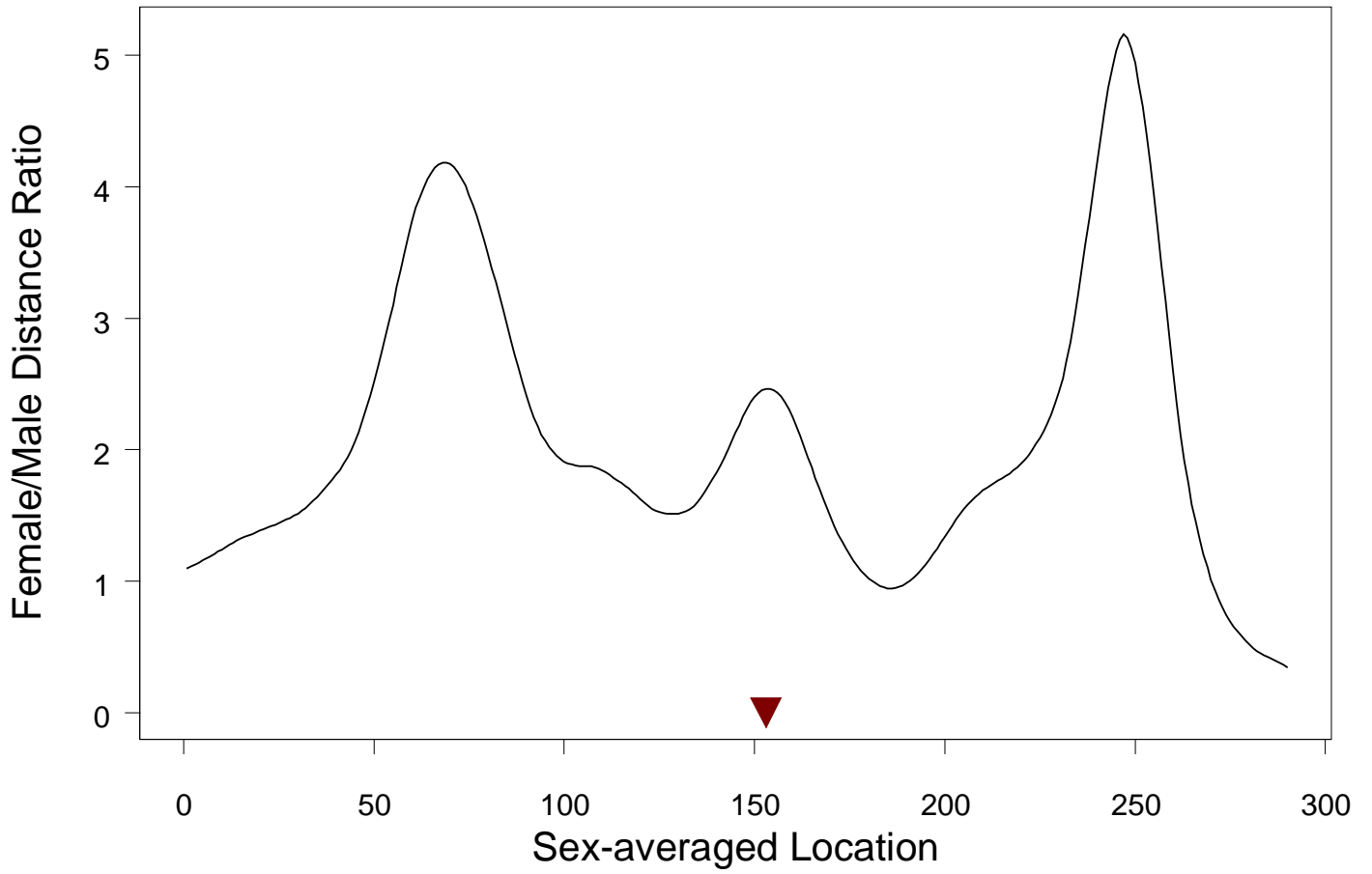
Crossovers in Father



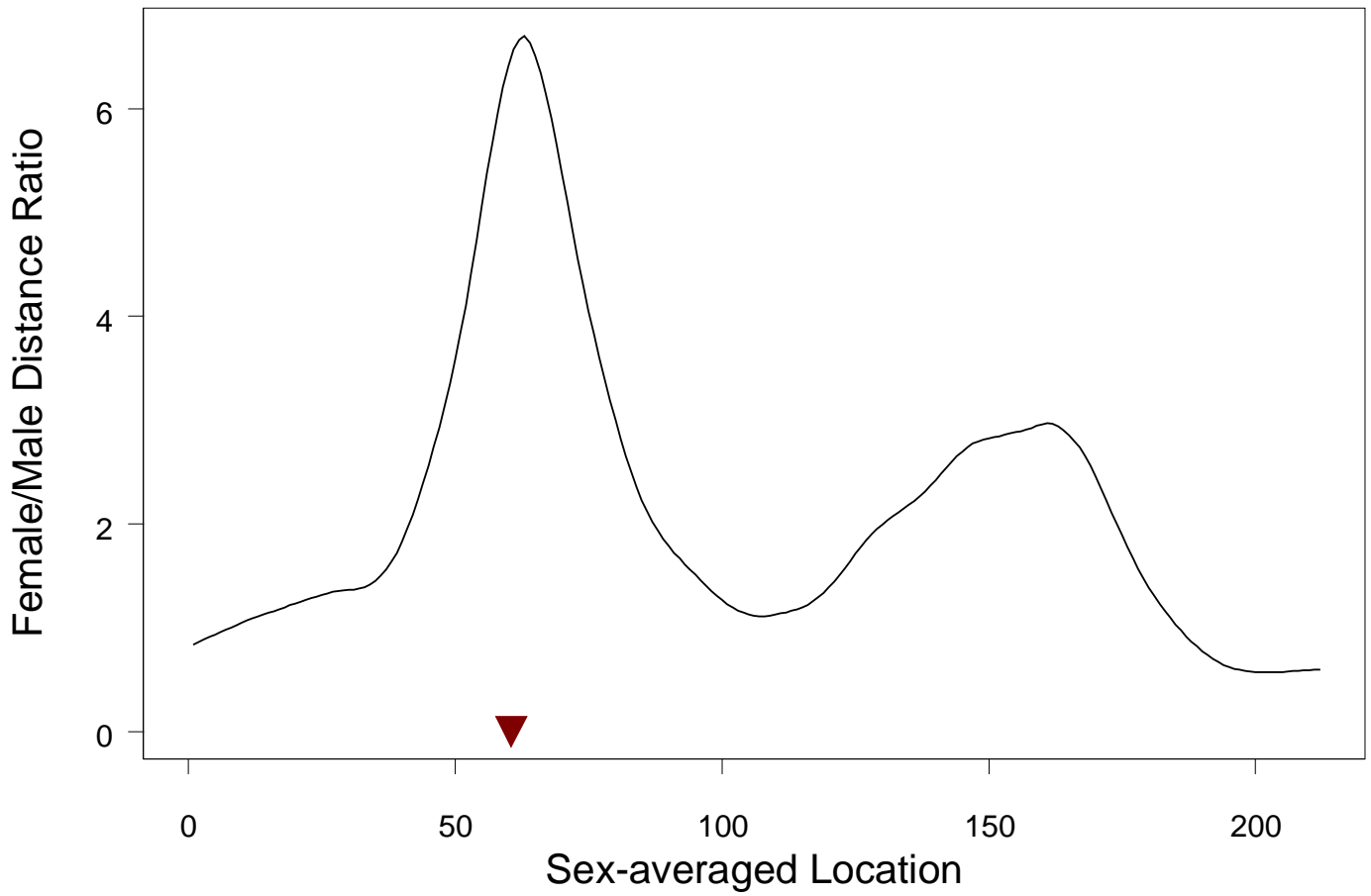
Crossovers in Mother



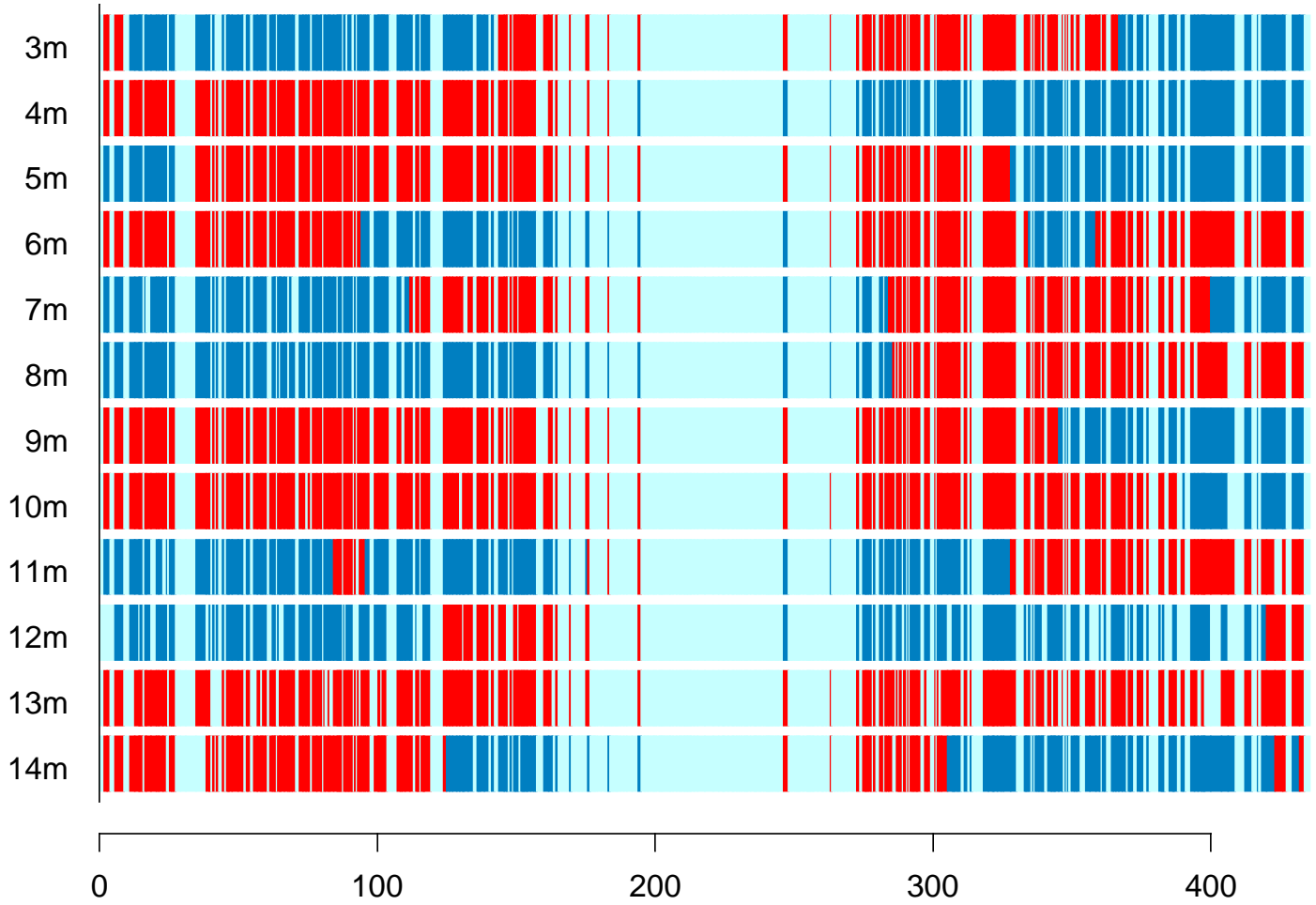
Chromosome 1



Chromosome 4



Chromosome 6 Family 884

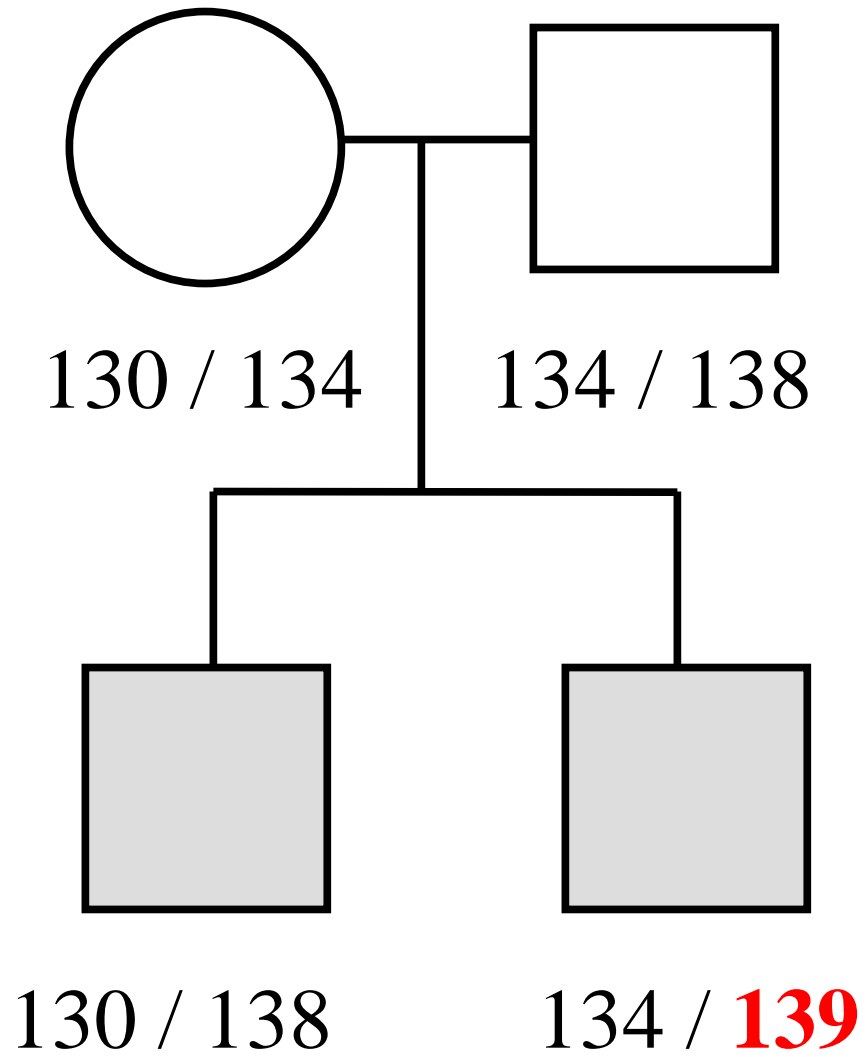


CEPH: homozygous regions

Individual	Number of regions	Length (cM)	Total length (cM)
884-01 (father)	10	2 – 23	116
884-02 (mother)	11	5 – 40	165
884: grandparents	2 7 – 10	7 1 – 29	14 58 – 120
884: 12 progeny	5 – 18	2 – 36	62 – 204
102: 14 progeny	4 – 13	1 – 39	85 – 254
1331-12 (GP)	2	1 – 6	7
1416-14 (GP)	4	6 – 29	76
35/100 of the others	1 – 2	0 – 7	

Present

- Finish the past
- Yellowstone wolves
- Sib pairs where parents are unavailable
- Various disease projects



Future

- Large genetic/epidemiological studies
 - Many people
 - Many phenotypes
 - Many environmental covariates
- Highly computational statistical genetic methods
- Various disease projects