# Introduction to QTL mapping

Manuel Ferreira

*Boulder Introductory Course 2006*

# Outline

1. Aim

2. The Human Genome

3. Principles of Linkage Analysis

4. Parametric Linkage Analysis

5. Nonparametric Linkage Analysis

# 1. Aim

# QTL mapping

▷ <u>LOCALIZE</u> and then <u>IDENTIFY</u> a <span style="color:red">locus that regulates a trait (QTL)</span>

*Nucleotide or sequence of nucleotides with variation in the population, with different variants associated with different trait levels.*

For a heritable trait...

**Linkage:**   <u>localize</u> region of the genome where a QTL that regulates the trait is likely to be harboured

<u>Family-specific phenomenon:</u>
Affected individuals in a family share the same ancestral predisposing DNA segment at a given QTL

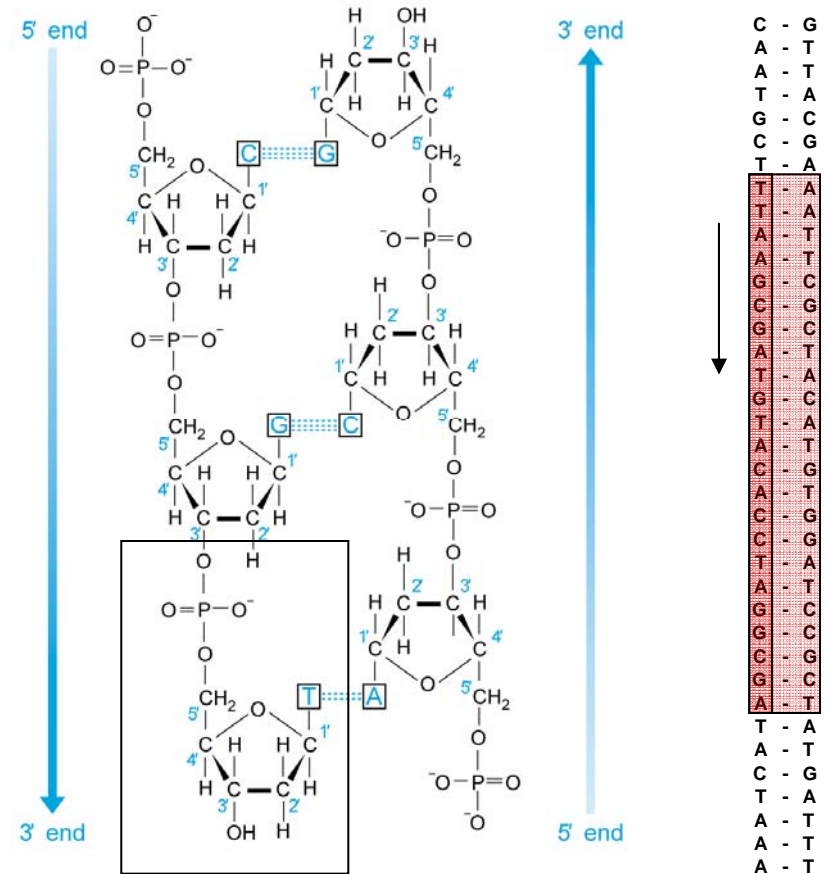**Association:** <u>identify</u> a QTL that regulates the trait

<u>Population-specific phenomenon:</u>
Affected individuals in a population share the same ancestral predisposing DNA segment at a given QTL

# 2. Human Genome

# DNA structure

▷ A DNA molecule is a linear backbone of alternating sugar residues and phosphate groups

▷ Attached to carbon atom 1' of each sugar is a nitrogenous base: A, C, G or T

▷ Two DNA molecules are held together in anti-parallel fashion by hydrogen bonds between bases [Watson-Crick rules] *Antiparallel double helix*

▷ A gene is a segment of DNA which is transcribed to give a protein or RNA product

▷ Only one strand is read during gene transcription

▷ Nucleotide: 1 phosphate group + 1 sugar + 1 base
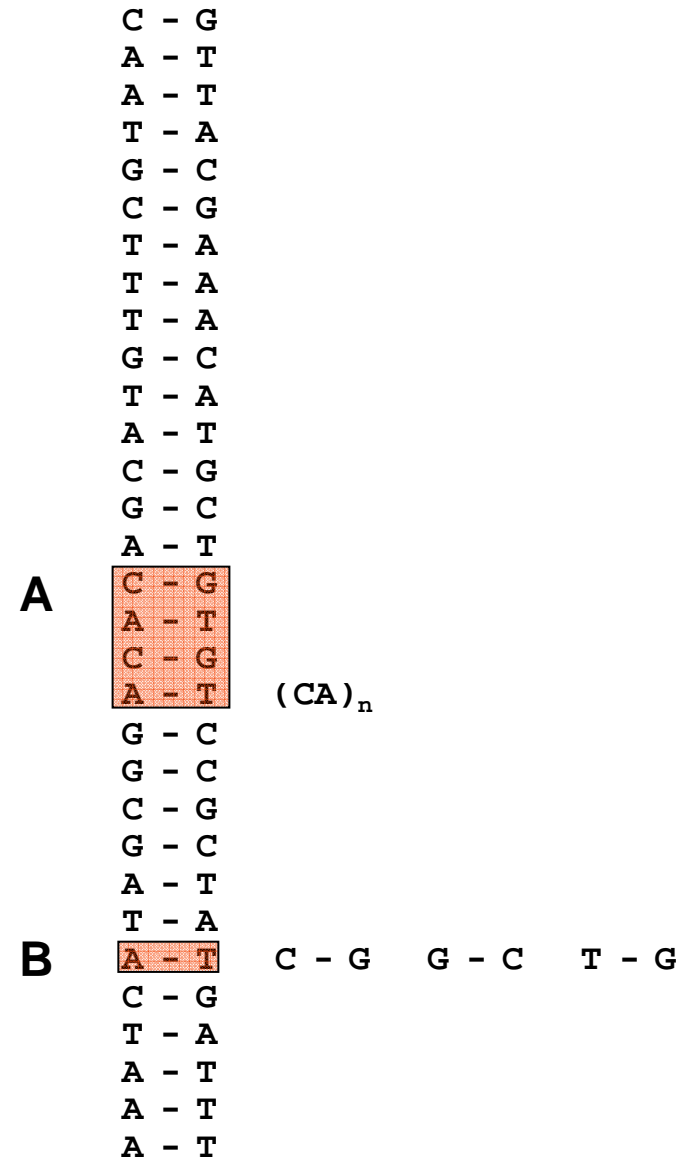
# DNA polymorphisms

▷ Microsatellites
>100,000
Many alleles, $(CA)_n$, very
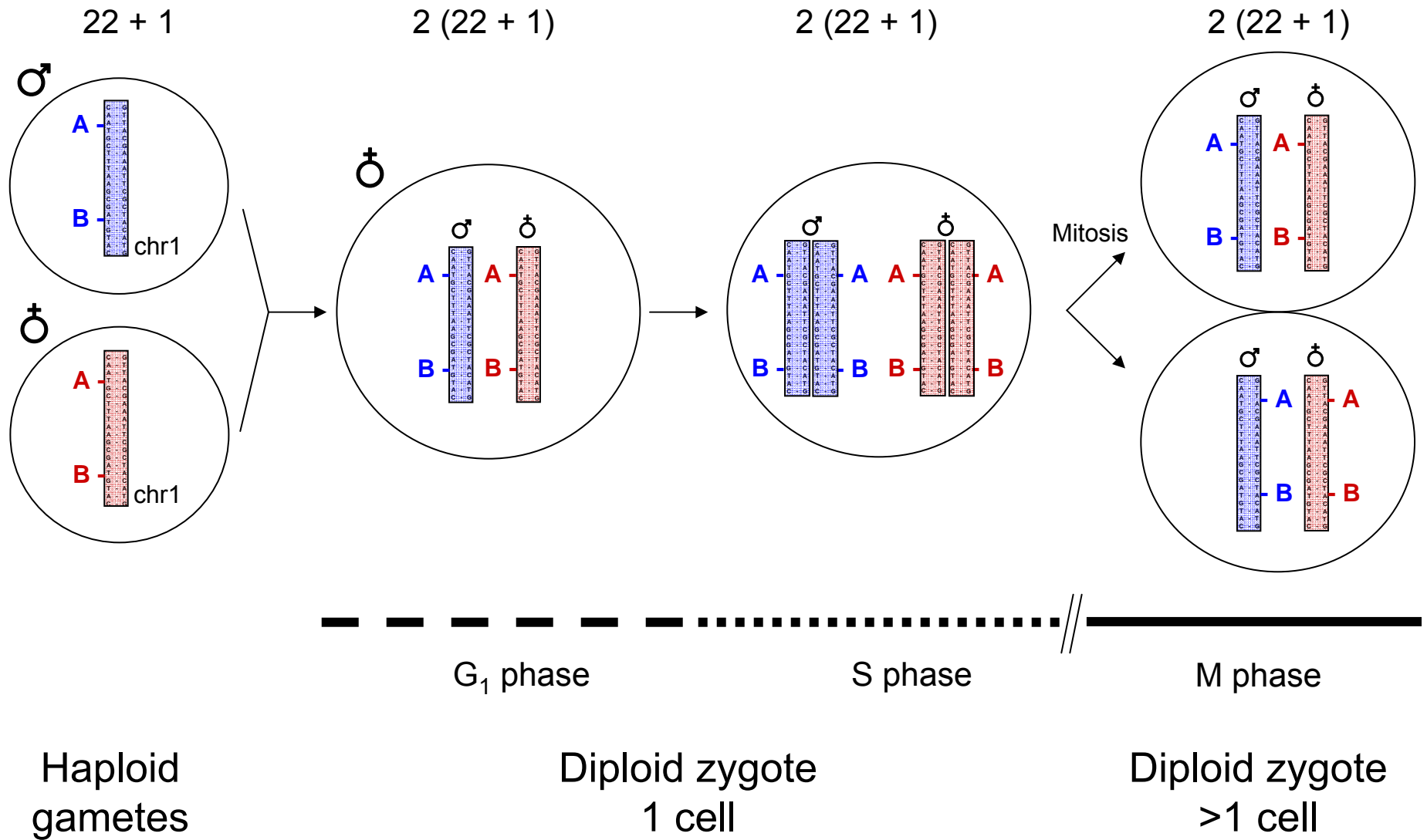informative, even, easily automated

▷ SNPs
10,054,521 (25 Jan '05)
10,430,753 (11 Mar '06)
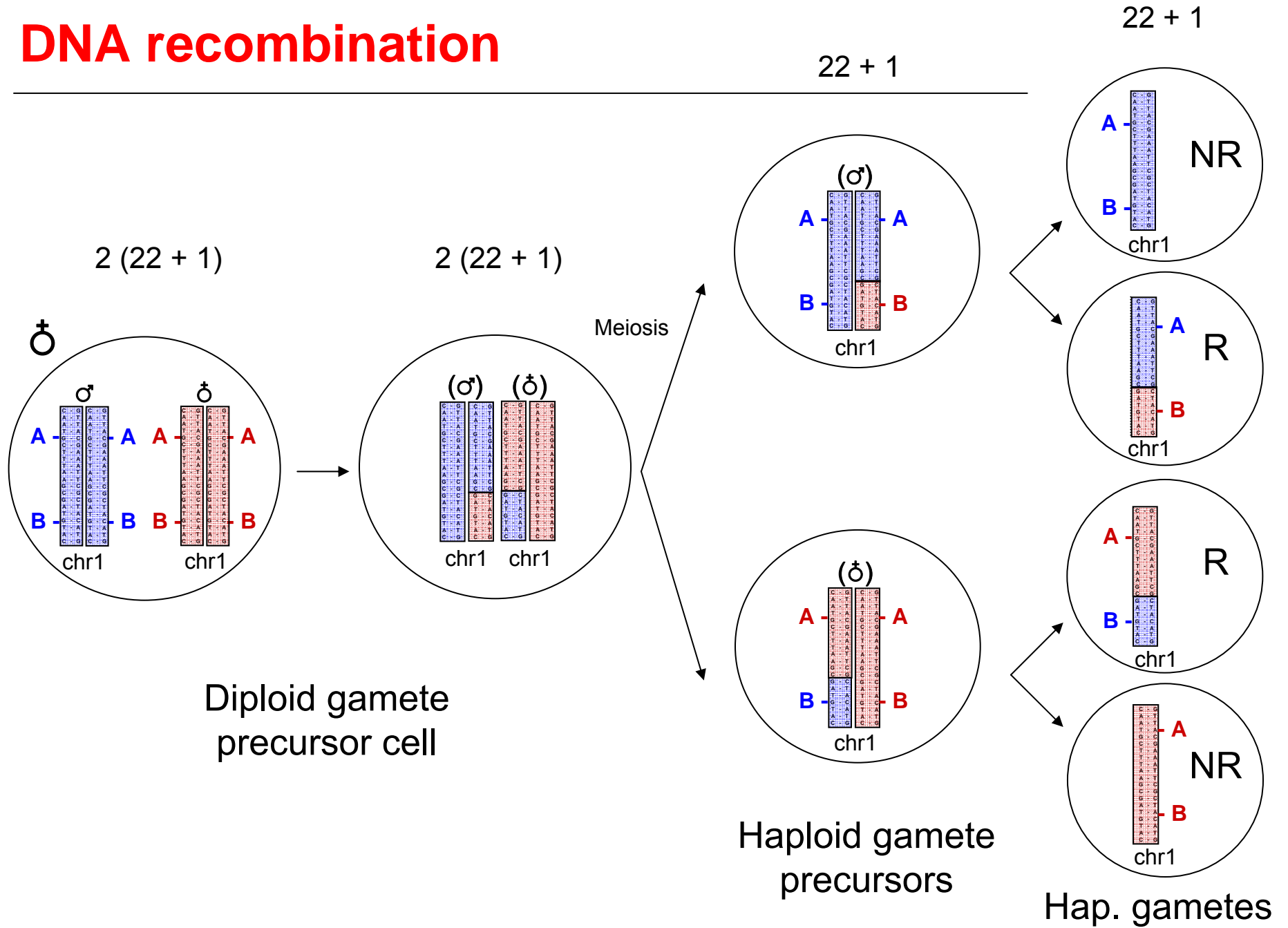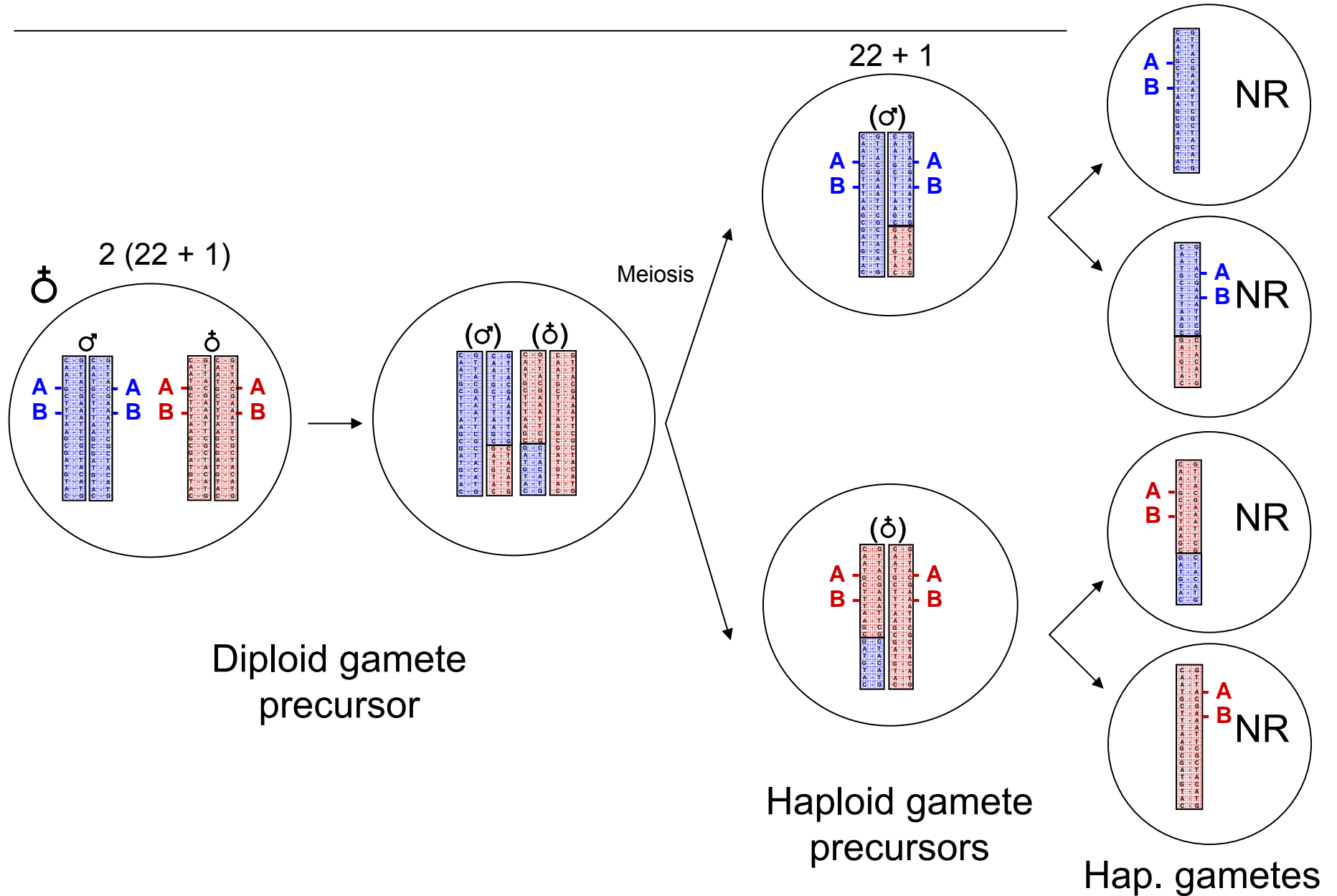Most with 2 alleles (up to 4), not very
informative, even, easily automated

```
    C - G
    A - T
    A - T
    T - A
    G - C
    C - G
    T - A
    T - A
    T - A
    G - C
    T - A
    A - T
    C - G
    G - C
    A - T
A   C - G
    A - T
    C - G
    A - T      (CA)n
    G - C
    G - C
    C - G
    G - C
    A - T
    T - A
B   A - T      C - G   G - C   T - G
    C - G
    T - A
    A - T
    A - T
    A - T
```

# DNA organization

# DNA recombination

# DNA recombination between linked loci

# Human Genome - summary

▷ <u>DNA is a linear sequence of nucleotides partitioned into 23 chromosomes</u>

Two copies of each chromosome (2x22 autosomes + XY), from

paternal and maternal origins. During meiosis in gamete precursors,

recombination can occur between maternal and paternal homologs

▷ <u>Recombination fraction between loci A and B ($\theta$)</u>

Proportion of gametes produced that are recombinant for A and B

If A and B are very far apart: 50%R:50%NR  -  $\theta$ = 0.5

If A and B are very close together: <50%R   -  0 ≤ $\theta$ < 0.5

▷ <u>Recombination fraction ($\theta$) can be converted to genetic distance (cM)</u>

Haldane:   $cM = 100 \cdot \left[ -0.5 \cdot \ln(1 - 2 \cdot \theta) \right]$                eg. θ=0.17, cM=20.8
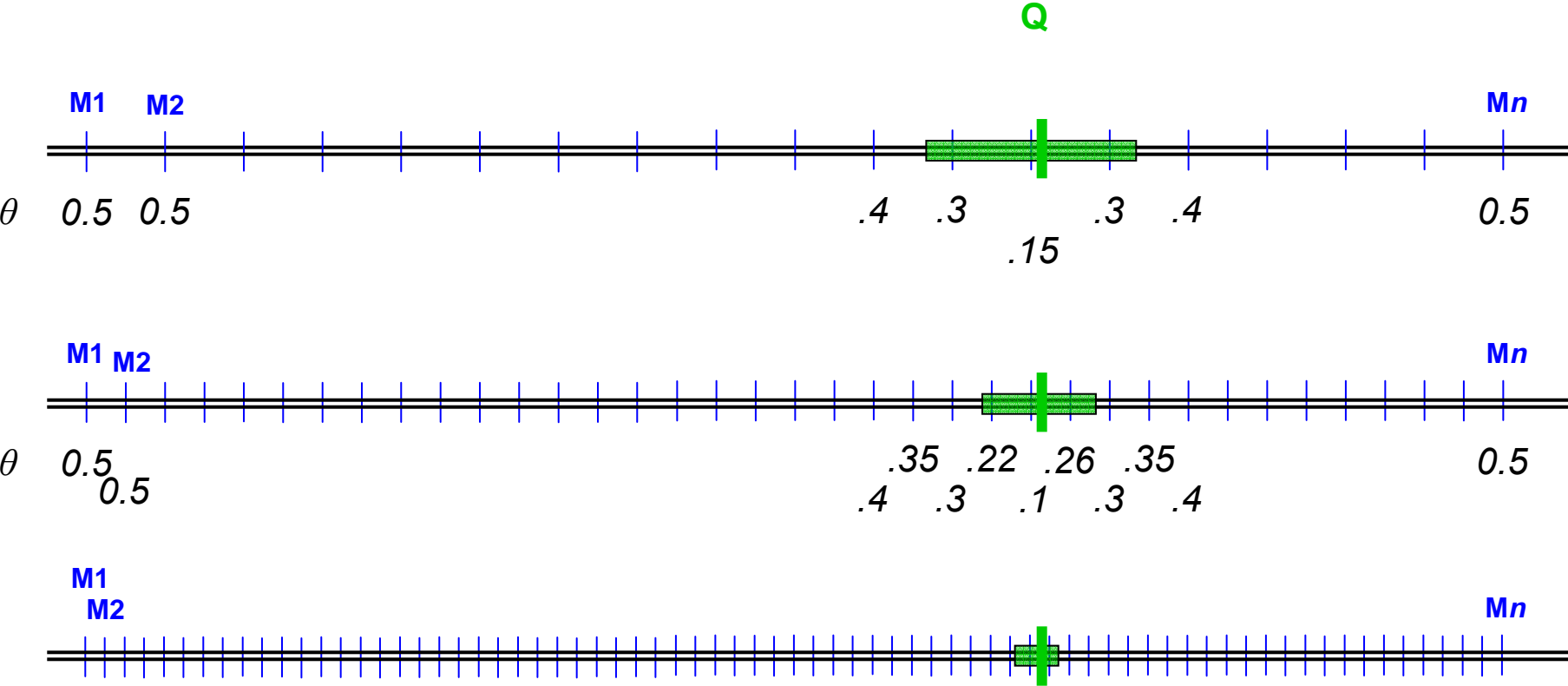
Kosambi:   $cM = 100 \cdot \left[ 0.25 \cdot \ln\left( (1 + 2 \cdot \theta)/(1 - 2 \cdot \theta) \right) \right]$  eg. θ=0.17, cM=17.7

# 3. Principles of Linkage Analysis

# Linkage Analysis requires genetic markers

# Linkage Analysis: Parametric vs. Nonparametric



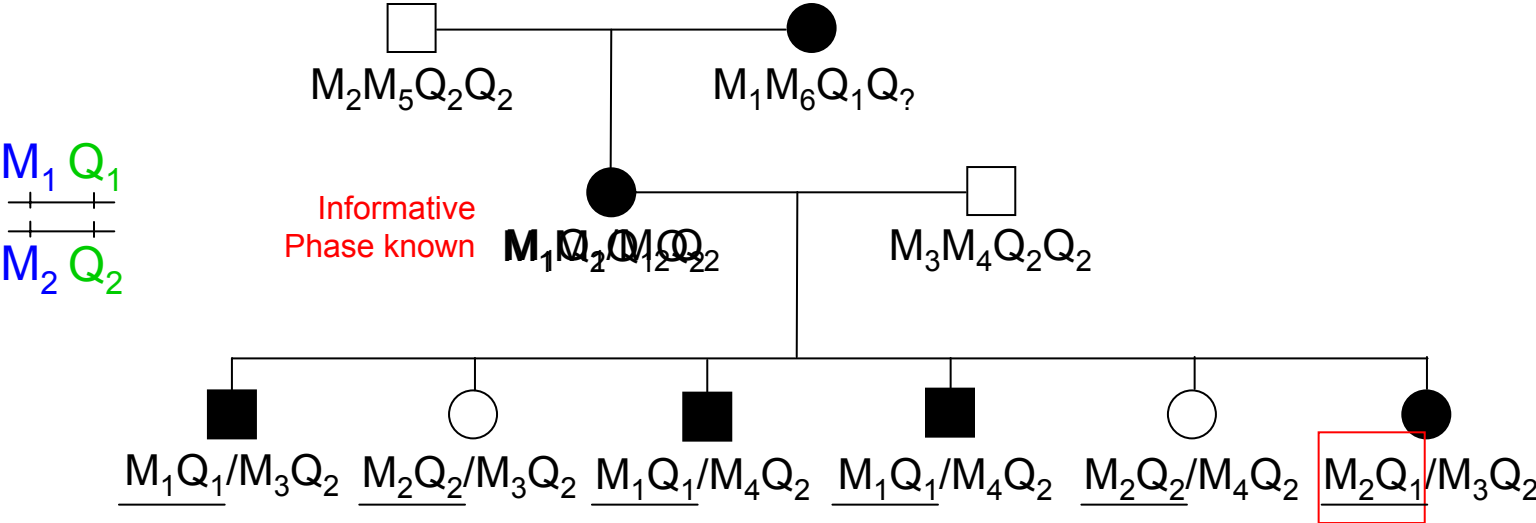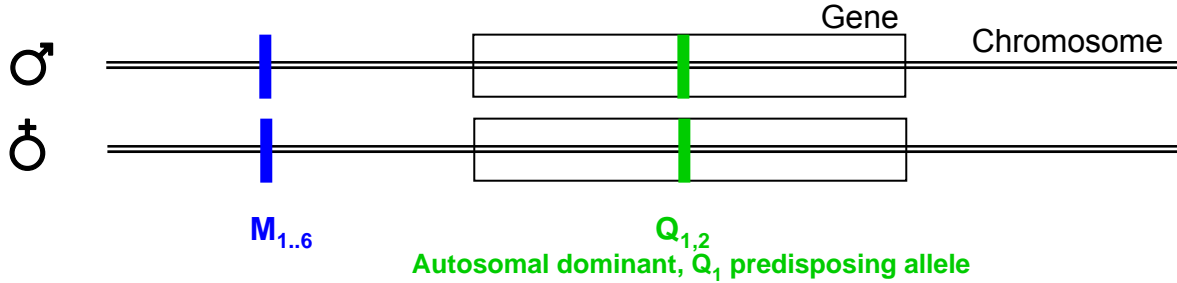*Adapted from Weiss & Terwilliger 2000*

# 4. Parametric Linkage Analysis

# Linkage with informative phase known meiosis



Gene

Chromosome

♂

♂

$M_{1..6}$

$Q_{1,2}$

**Autosomal dominant, $Q_1$ predisposing allele**

$M_2M_5Q_2Q_2$          $M_1M_6Q_1Q_?$

$M_1\ Q_1$
____
$M_2\ Q_2$

Informative
Phase known  $M_1M_2Q_{12}Q_{22}$          $M_3M_4Q_2Q_2$

$M_1Q_1/M_3Q_2$  $M_2Q_2/M_3Q_2$  $M_1Q_1/M_4Q_2$  $M_1Q_1/M_4Q_2$  $M_2Q_2/M_4Q_2$  $M_2Q_1/M_3Q_2$

NR: $M_1Q_1$
NR: $M_2Q_2$
R: $M_1Q_2$
R: $M_2Q_1$

$\theta_{MQ} = 1/6 = 0.17$   (~20.8 cM)

# Linkage with informative phase unknown meiosis



$M_1 Q_1$     $M_1 Q_2$

$M_2 Q_2$     $M_2 Q_1$

**Informative Phase <u>unknown</u>**

| $M_1 Q_1/M_2 Q_2$ | | $P$ | $N$ |
|---|---|---|---|
| NR: | $M_1 Q_1$ | $\frac{1}{2}(1-\theta)$ | 3 |
| NR: | $M_2 Q_2$ | $\frac{1}{2}(1-\theta)$ | 2 |
| R: | $M_1 Q_2$ | $\frac{1}{2}\theta$ | 0 |
| R: | $M_2 Q_1$ | $\frac{1}{2}\theta$ | 1 |

| $M_1 Q_2/M_2 Q_1$ | | $P$ | $N$ |
|---|---|---|---|
| R: | $M_1 Q_1$ | $\frac{1}{2}\theta$ | 3 |
| R: | $M_2 Q_2$ | $\frac{1}{2}\theta$ | 2 |
| NR: | $M_1 Q_2$ | $\frac{1}{2}(1-\theta)$ | 0 |
| NR: | $M_2 Q_1$ | $\frac{1}{2}(1-\theta)$ | 1 |

$$L(X\mid\theta)= \frac{1}{2}\cdot\left[\theta^1\cdot(1-\theta)^5\right] + \frac{1}{2}\cdot\left[\theta^5\cdot(1-\theta)^1\right]$$

$$L(X\mid\theta=0.5)= \frac{1}{2}\cdot\left[0.5^1\cdot(1-0.5)^5\right] + \frac{1}{2}\cdot\left[0.5^5\cdot(1-0.5)^1\right] = (0.5)^6$$
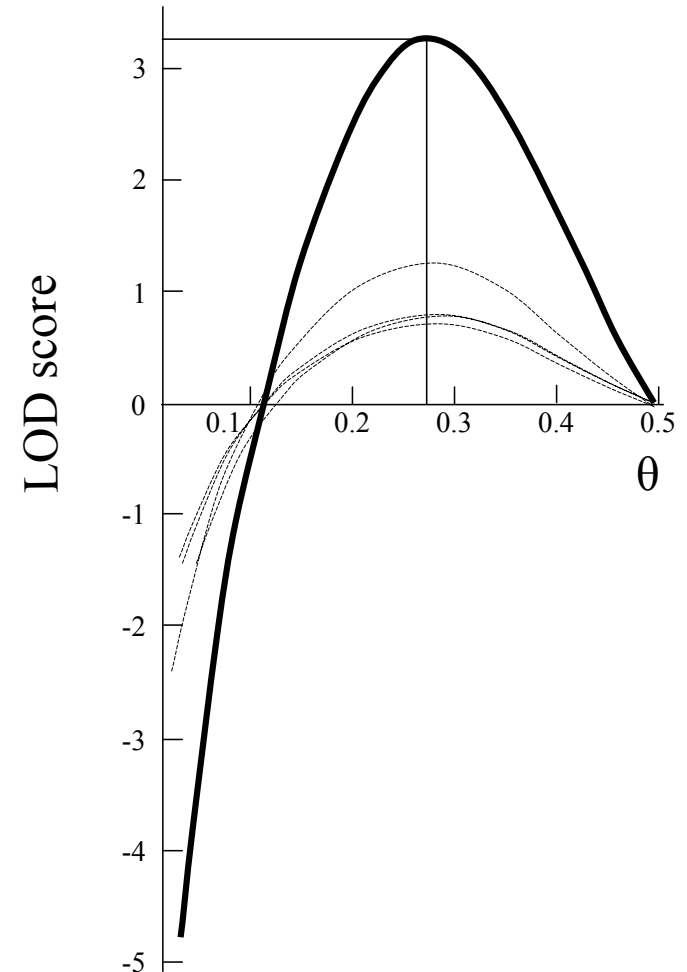
# Parametric LOD score calculation

$$OD = \frac{L(X \mid \theta)}{L(X \mid \theta = 0.5)} \qquad LOD = \log_{10} \frac{L(X \mid \theta)}{L(X \mid \theta = 0.5)}$$

$$LOD = \log_{10} \frac{\frac{1}{2} \cdot \left[\theta^1 \cdot (1-\theta)^5\right] + \frac{1}{2} \cdot \left[\theta^5 \cdot (1-\theta)^1\right]}{(0.5)^6}$$

$$OD = \prod_{i=1}^{n} \frac{L(X_i \mid \theta)}{L(X_i \mid \theta = 0.5)}$$

$$LOD = \log_{10} \left( \prod_{i=1}^{n} \frac{L(X_i \mid \theta)}{L(X_i \mid \theta = 0.5)} \right)$$

$$LOD = \sum_{i=1}^{n} \log_{10} \left( \frac{L(X_i \mid \theta)}{L(X_i \mid \theta = 0.5)} \right) = \sum_{i=1}^{n} LOD_i$$
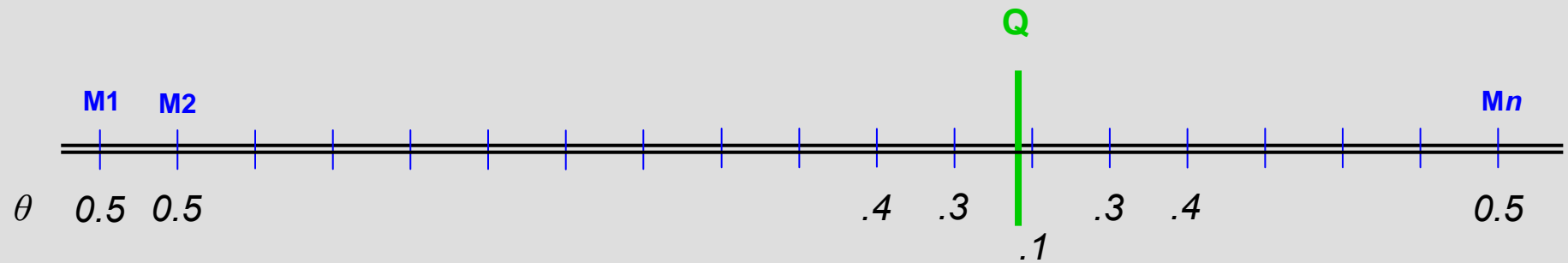


▷ Overall LOD score for a given $\theta$ is the sum of all family LOD scores at $\theta$

eg. LOD=3 for $\theta$=0.28

# Parametric Linkage Analysis - summary

Q

M1  M2                                                                    M*n*

$\theta$   0.5  0.5                                   .4   .3       .3   .4              0.5

.1

▷ For each marker, estimate the $\theta$ that yields highest LOD score across all families

▷ This $\theta$ (and the LOD) will depend upon the mode of inheritance assumed

   MOI determines the genotype at the trait locus Q and thus determines the number of meiosis which are recombinant or nonrecombinant. Limited to Mendelian diseases.

▷ Markers with a significant parametric LOD score (>3) are said to be linked to the trait locus with recombination fraction $\theta$

# Outline

1. Aim

2. The Human Genome

3. Principles of Linkage Analysis

4. Parametric Linkage Analysis

5. Nonparametric Linkage Analysis

# 5. Nonparametric Linkage Analysis

# Approach

▷ [Parametric: genotype marker locus & genotype trait locus](#)

(latter inferred from phenotype according to a specific disease model)

Parameter of interest: $\theta$ between marker and trait loci

▷ [Nonparametric: genotype marker locus & phenotype](#)

If a trait locus truly regulates the expression of a phenotype, then two relatives with similar phenotypes should have similar genotypes at a marker in the vicinity of the trait locus, and vice-versa.

Interest: correlation between <u>phenotypic similarity</u> and marker <u>genotypic similarity</u>

No need to specify mode of inheritance, allele frequencies, etc...

# Phenotypic similarity between relatives

▷ Squared trait differences $\qquad (X_1 - X_2)^2$

▷ Squared trait sums $\qquad (X_1 + X_2)^2$

▷ Trait cross-product $\qquad [(X_1 - \mu) \cdot (X_2 - \mu)]$

▷ Trait variance-covariance matrix $\qquad \begin{Bmatrix} Var(X_1) & Cov(X_1 X_2) \\ Cov(X_1 X_2) & Var(X_2) \end{Bmatrix}$
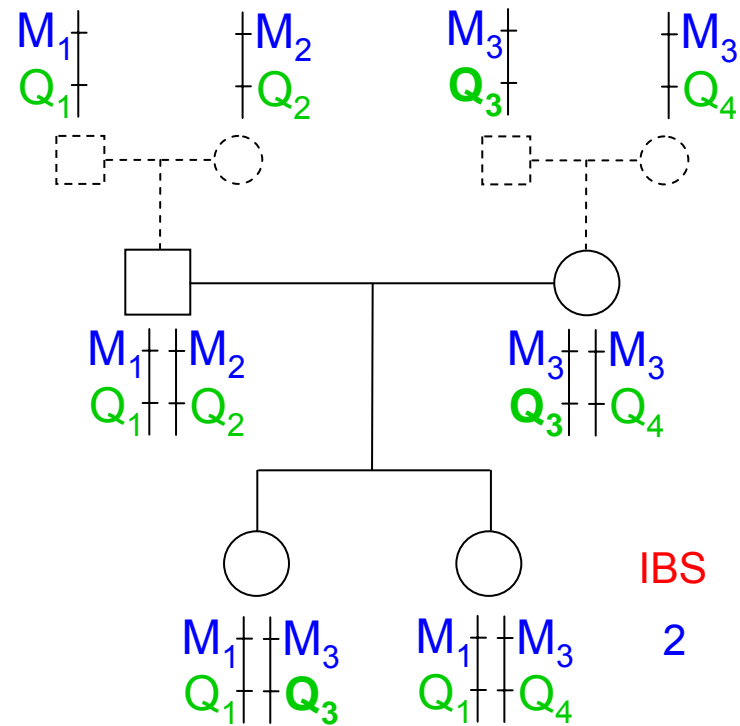
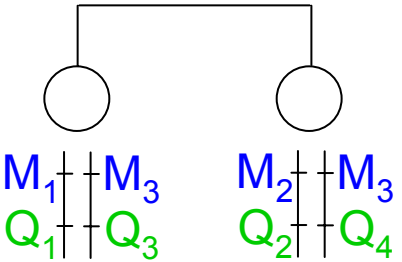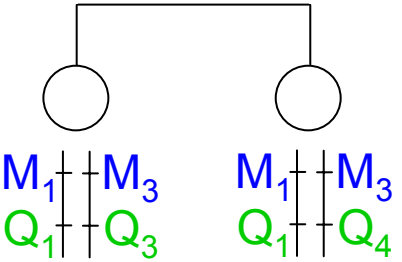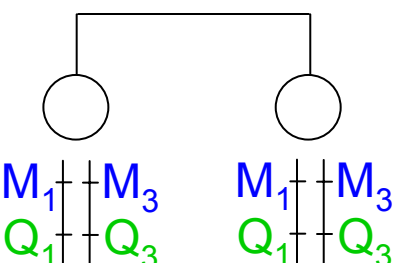▷ Affection concordance

# Genotypic similarity between relatives

▷ **IBS**   Alleles shared <u>Identical By State</u> "look the same", may have the same DNA sequence but they are not necessarily derived from a known common ancestor

▷ **IBD**   Alleles shared <u>Identical By Descent</u> are a copy of the same ancestor allele
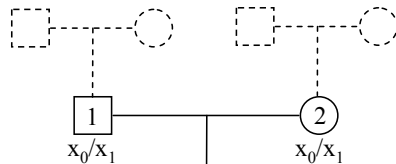
# Genotypic similarity between relatives - $\pi$



| | Inheritance vector (M) | Number of alleles shared IBD | Proportion of alleles shared IBD - $\pi$ |
|---|---|---|---|
| $M_1$ $M_3$ $\quad$ $M_2$ $M_3$ $\quad$ $Q_1$ $Q_3$ $\quad$ $Q_2$ $Q_4$ | 0 $\quad$ 0 $\quad$ 1 $\quad$ 1 | 0 | 0 |
| $M_1$ $M_3$ $\quad$ $M_1$ $M_3$ $\quad$ $Q_1$ $Q_3$ $\quad$ $Q_1$ $Q_4$ | 0 $\quad$ 0 $\quad$ 0 $\quad$ 1 | 1 | 0.5 |
| $M_1$ $M_3$ $\quad$ $M_1$ $M_3$ $\quad$ $Q_1$ $Q_3$ $\quad$ $Q_1$ $Q_3$ | 0 $\quad$ 0 $\quad$ 0 $\quad$ 0 | 2 | 1 |

# Genotypic similarity between relatives - $\hat{\pi}$

A      B      C      D

$2^{2n}$

| 3 | 4 | Inheritance vector | IBD |
|---|---|---|---|
| $x_0/x_0$ | $x_0/x_0$ | 0000 | 2 |
| $x_0/x_0$ | $x_0/x_1$ | 0001 | 1 |
| $x_0/x_0$ | $x_1/x_0$ | 0010 | 1 |
| $x_0/x_0$ | $x_1/x_1$ | 0011 | 0 |
| $x_0/x_1$ | $x_0/x_0$ | 0100 | 1 |
| $x_0/x_1$ | $x_0/x_1$ | 0101 | 2 |
| $x_0/x_1$ | $x_1/x_0$ | 0110 | 0 |
| $x_0/x_1$ | $x_1/x_1$ | 0111 | 1 |
| $x_1/x_0$ | $x_0/x_0$ | 1000 | 1 |
| $x_1/x_0$ | $x_0/x_1$ | 1001 | 0 |
| $x_1/x_0$ | $x_1/x_0$ | 1010 | 2 |
| $x_1/x_0$ | $x_1/x_1$ | 1011 | 1 |
| $x_1/x_1$ | $x_0/x_0$ | 1100 | 0 |
| $x_1/x_1$ | $x_0/x_1$ | 1101 | 1 |
| $x_1/x_1$ | $x_1/x_0$ | 1110 | 1 |
| $x_1/x_1$ | $x_1/x_1$ | 1111 | 2 |

**P (IBD=0)**
**P (IBD=1)**
**P (IBD=2)**

$$\hat{\pi} =$$

# Practical

▷ **Aim** (1) Estimate IBD with MERLIN; (2) IBD estimation can be influenced by genotyped individuals and allele frequencies; (3) compute $\hat{\pi}$

`H:\manuel  -    Copy folder "Linkage" to C:\`

    1. Open with Notepad: pr1.ped  pr1.dat  pr1.map  pr1.freq
    2. Start>Run>C:/Linkage/pfe32.exe
    3. Run Command Prompt
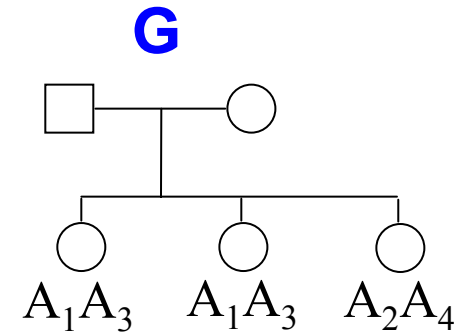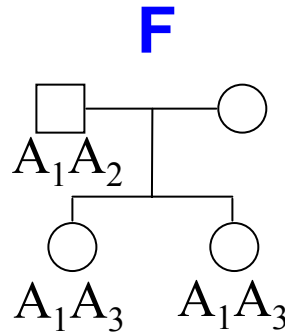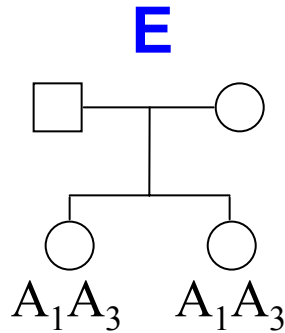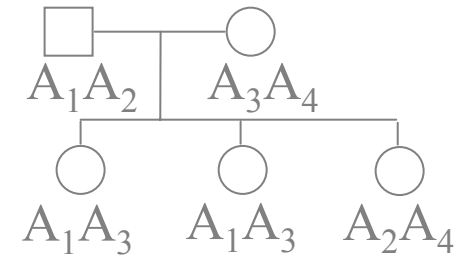    4. Keep a File Explorer window open

## Exercice1

(1) Estimate IBD for pedigrees **A**, **B** and **C** in the previous slide

(2) Change allele frequencies (pr1.freq) from 0.25 0.25 0.25 0.25 to
       (i) 0.45 0.25 0.25 0.05 and
       (ii) 0.05 0.25 0.25 0.45

# Practical



$A_1A_2$ | $A_3A_4$

$A_1A_3$  $A_1A_3$  $A_2A_4$

## Exercice 2

(1) Modify pr1.ped and estimate IBD probabilities and $\hat{\pi}$ between twin 1 and twin 2 for pedigrees **E**, **F** and **G**:
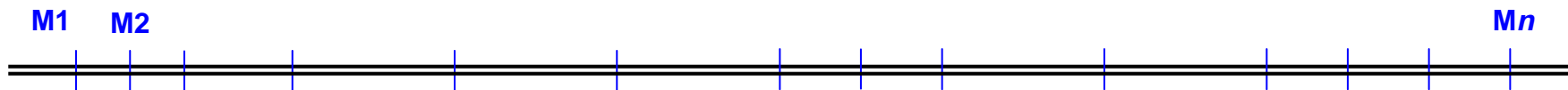
**E**



$A_1A_3$  $A_1A_3$

**F**

$A_1A_2$



$A_1A_3$  $A_1A_3$

**G**



$A_1A_3$  $A_1A_3$  $A_2A_4$

---
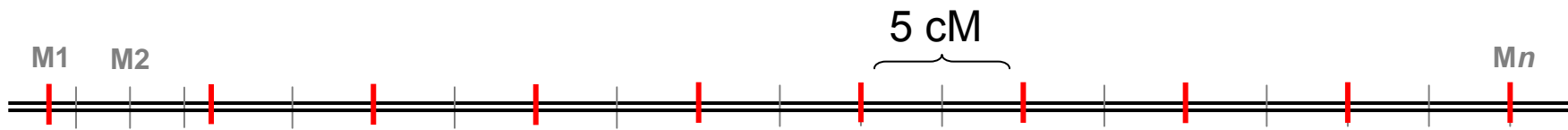
$P(\text{IBD}=0)$

$P(\text{IBD}=1)$

$P(\text{IBD}=2)$

$\hat{\pi}$

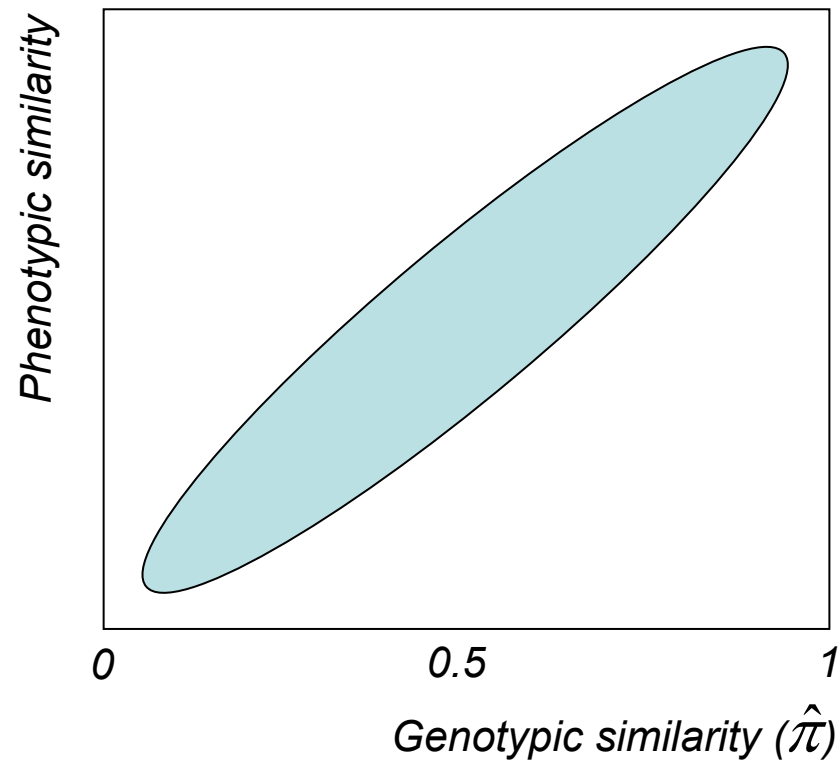Allele frequencies on pr1.freq: 0.25 0.25 0.25 0.25
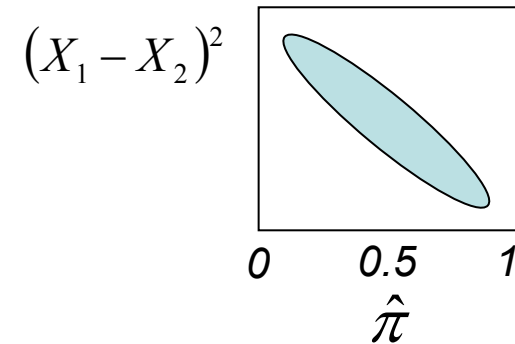
**IBD at a marker**
Singlepoint IBD

5 cM

**IBD at a 'grid'**
Multipoint IBD

# Statistics that incorporate both phenotypic and genotypic similarities

# Haseman-Elston regression – Quantitative traits

$$E\left[(X_1 - X_2)^2 \mid \hat{\pi}\right]$$

$$= E\left[(X_1^2 + X_2^2 - 2 \cdot X_1 \cdot X_2) \mid \hat{\pi}\right]$$

$$= Var(X_1) + Var(X_2) - 2Cov(X_1 X_2 \mid \hat{\pi})$$

$$(X_1 - X_2)^2$$



$$0 \qquad 0.5 \qquad 1$$

$$\hat{\pi}$$

$$Var(X_1) = Var(X_2) = V_Q + V_A + V_C + V_E$$

$$Cov(X_1, X_2 \mid \hat{\pi}) = \hat{\pi} \cdot V_Q + \frac{1}{2} \cdot V_A + V_C$$

| | $X_1$ | $X_2$ | $(X_1-X_2)^2$ | $\hat{\pi}$ |
|---|---|---|---|---|
| 1 | 2.2 | 2.1 | 0.01 | 0.9 |
| 2 | 1.9 | 2.3 | 0.16 | 0.6 |
| 3 | 2.3 | 2.6 | 0.09 | 0.7 |
| 4 | 3.4 | 1.6 | 3.24 | 0.1 |
| 5 | 2.5 | 2.3 | 0.04 | 0.8 |
| ... | | | | |
| 1000 | 2.4 | 2.4 | 0 | 0.9 |

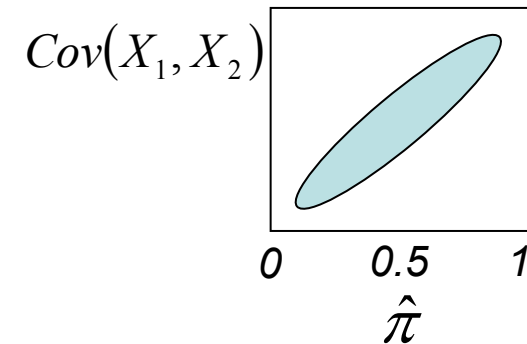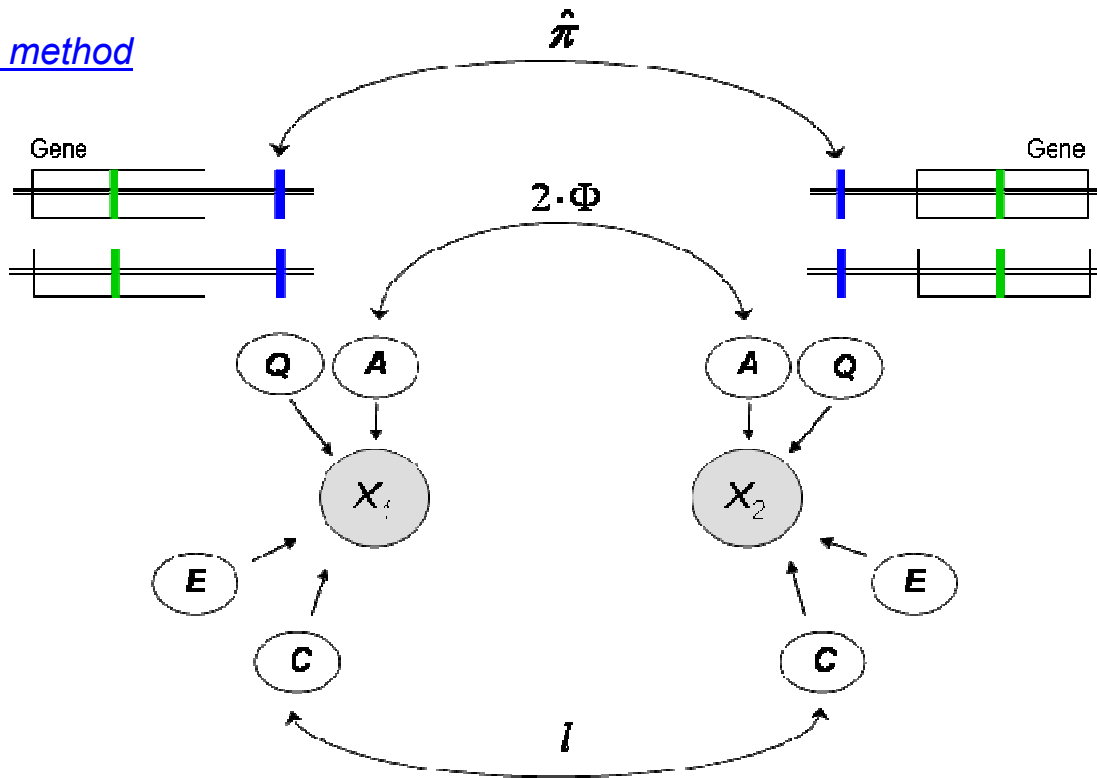$$E\left[(X_1 - X_2)^2 \mid \hat{\pi}\right] = -2 \cdot V_Q \cdot \hat{\pi} + 2 \cdot V_Q + V_A + 2 \cdot V_E$$

*Phenotypic dissimilarity*  =  **b ×** *Genotypic similarity*  **+**  **c**

# VC ML – Quantitative & Categorical traits

$\hat{\pi}$ _method_



$$H_1: \quad Var(X_1) = Var(X_2) = V_Q + V_A + V_C + V_E$$
$$Cov(X_1, X_2 \mid \hat{\pi}) = \hat{\pi} \cdot V_Q + 2 \cdot \Phi \cdot V_A + l \cdot V_C$$

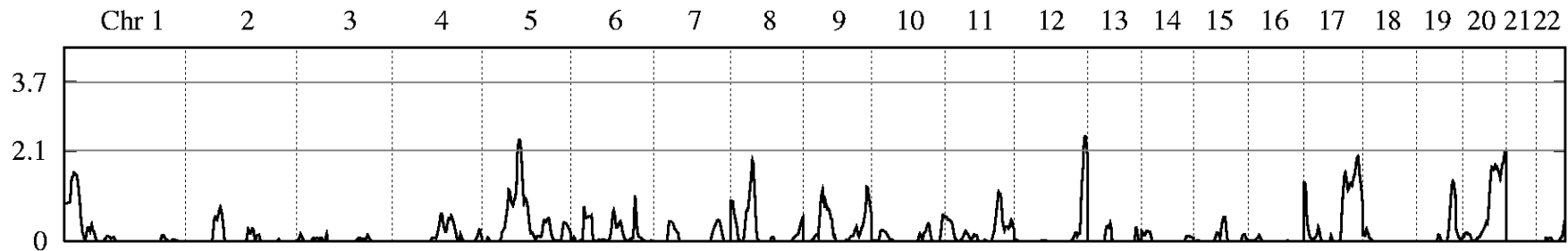$$H_0: \quad Var(X_1) = Var(X_2) = \qquad V_A + V_C + V_E$$
$$Cov(X_1, X_2 \mid \hat{\pi}) = \qquad 2 \cdot \Phi \cdot V_A + l \cdot V_C$$

$$LOD = \log_{10} \frac{L(H_1)}{L(H_0)}$$

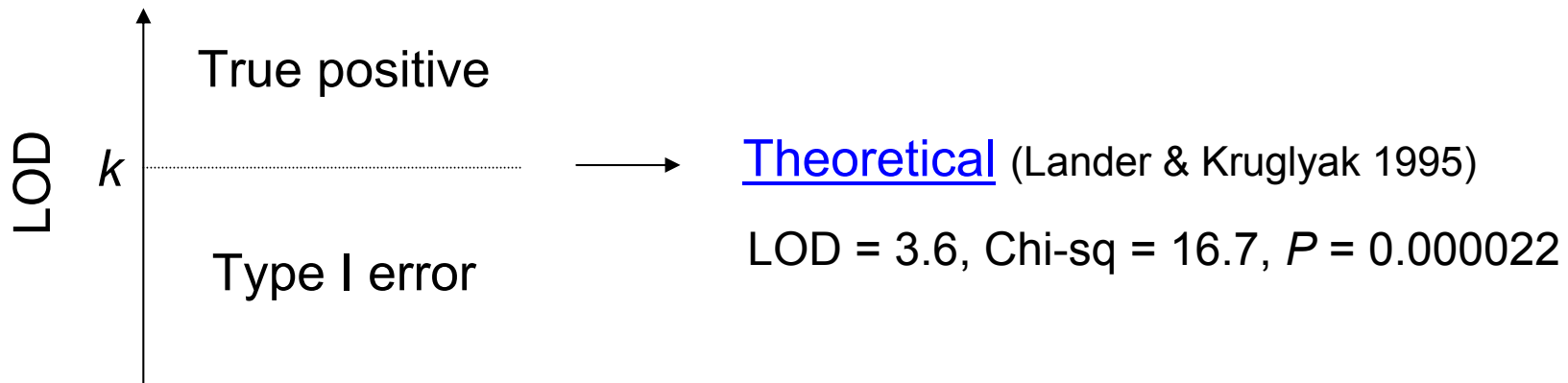_e.g. LOD=3_

# Genome-wide linkage analysis (e.g. VC)



Individual LOD scores can be expressed as $P$ values (Pointwise)

$$LOD \xrightarrow{(x4.6)} \text{Chi-sq (n-df)} \longrightarrow P \text{ value}$$

| 2.1 | 9.67 | 0.0009 |

LOD

True positive

$k$ $\longrightarrow$ Theoretical (Lander & Kruglyak 1995)

LOD = 3.6, Chi-sq = 16.7, $P$ = 0.000022

Type I error

# Nonparametric Linkage Analysis - summary

▷ No need to specify mode of inheritance

▷ Models phenotypic and genotypic similarity of relatives

▷ Expression of phenotypic similarity, calculation of IBD

▷ HE and VC are the most popular statistics used for linkage of quantitative traits

▷ Other statistics available, specially for affection traits