

Genetic Theory

Manuel AR Ferreira

Massachusetts General Hospital

Harvard Medical School

Boston

Boulder, 2007

Outline

1. Aim of this talk
2. Genetic concepts
3. Very basic statistical concepts
4. Biometrical model
5. Introduction to linkage analysis

1. Aim of this talk

Gene mapping

▷ LOCALIZE and then IDENTIFY a locus that regulates a trait



*Linkage
analysis*



*Association
analysis*

Linkage: If a locus regulates a trait, Trait Variance and Covariance between individuals can be expressed as a function of this locus.

Association: If a locus regulates a trait, Trait Mean in the population can be expressed as a function of this locus.

- ▷ Revisit common genetic parameters - such as allele frequencies, genetic effects, dominance, variance components, etc
- ▷ Use these parameters to construct a **biometric genetic model**



Model that expresses the:

(1) Mean

(2) Variance

(3) Covariance between individuals

for a quantitative phenotype as a function of the genetic parameters of a given locus.

- ▷ See how the **biometric model** provides a useful framework for linkage and association methods.

2. Genetic concepts

▷ **A. DNA level**

*DNA structure, organization
recombination*

▷ **B. Population level**

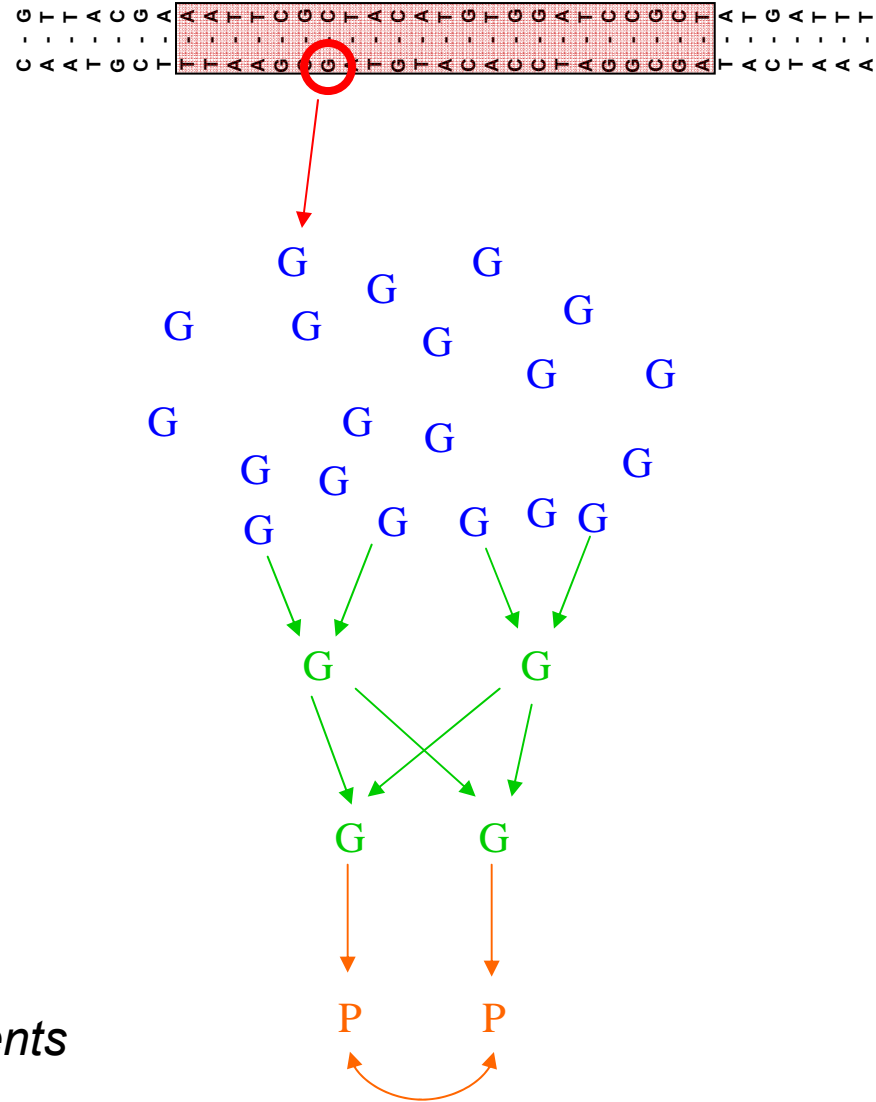
Allele and genotype frequencies

▷ **C. Transmission level**

*Mendelian segregation
Genetic relatedness*

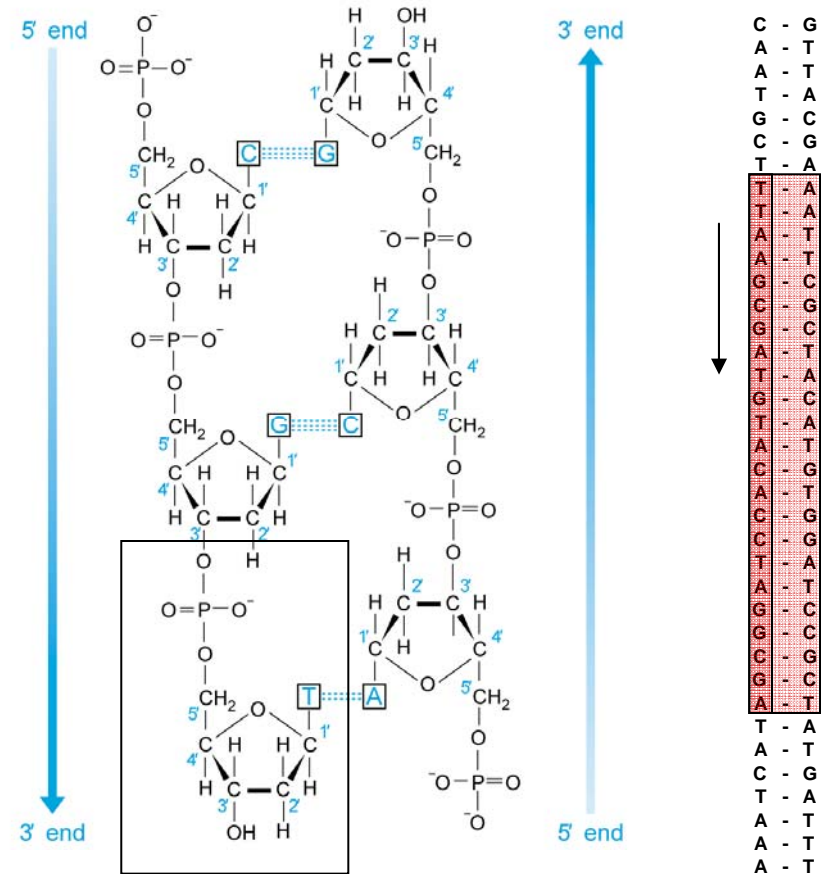
▷ **D. Phenotype level**

*Biometrical model
Additive and dominance components*



A. DNA level

- ▷ 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, $(CA)_n$ repeats, very informative, even, easily automated

▷ **SNPs**

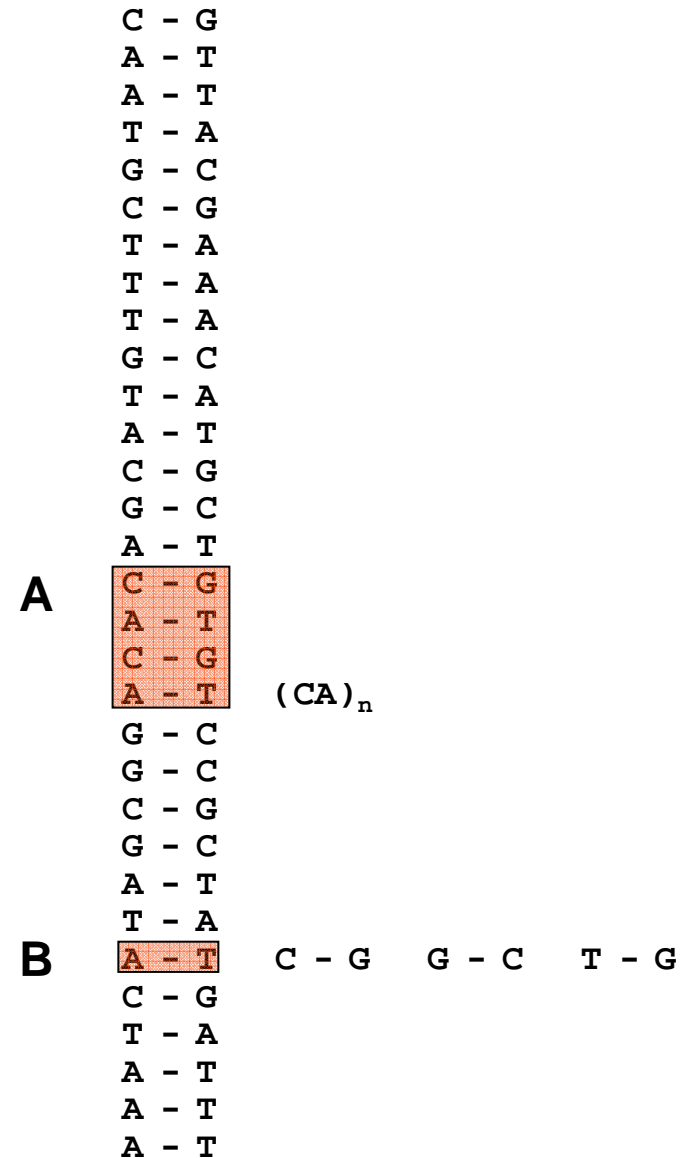
11,961,761 (build 126, 03 Mar '07)

Most with 2 alleles (up to 4), not very informative, even, easily automated

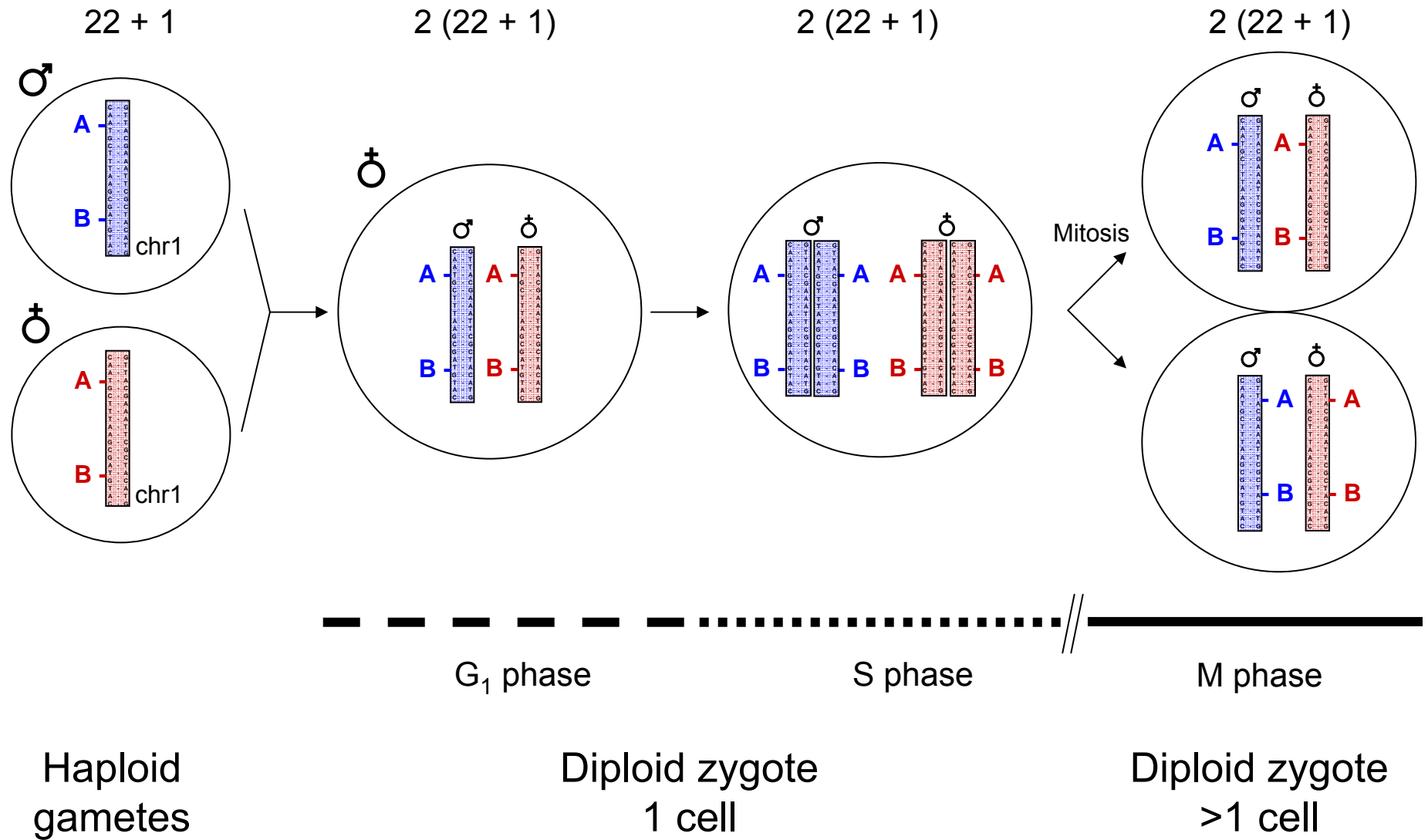
▷ **Copy Number polymorphisms**

~2000-3000 (?)

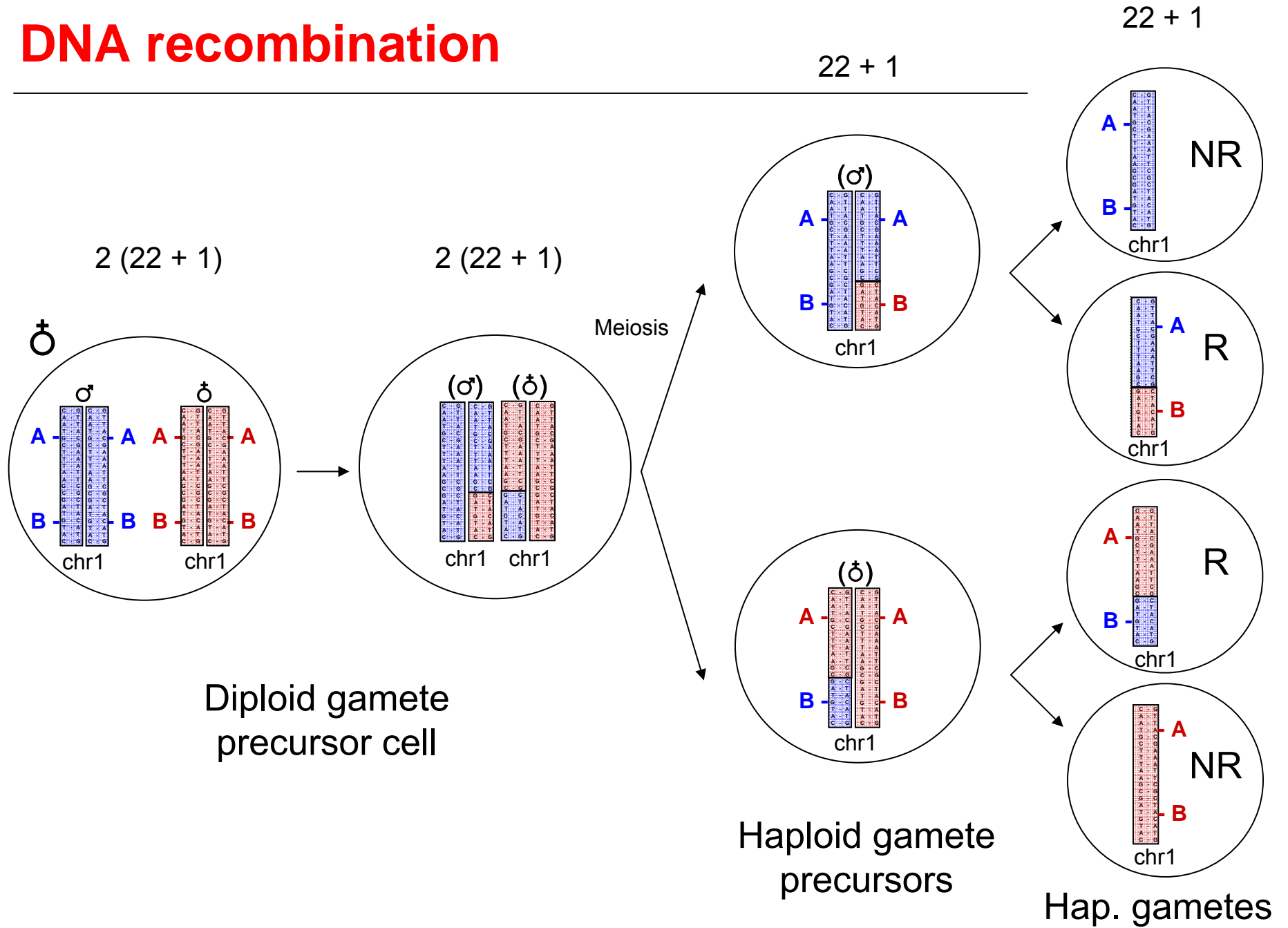
Many alleles, even, not yet automated



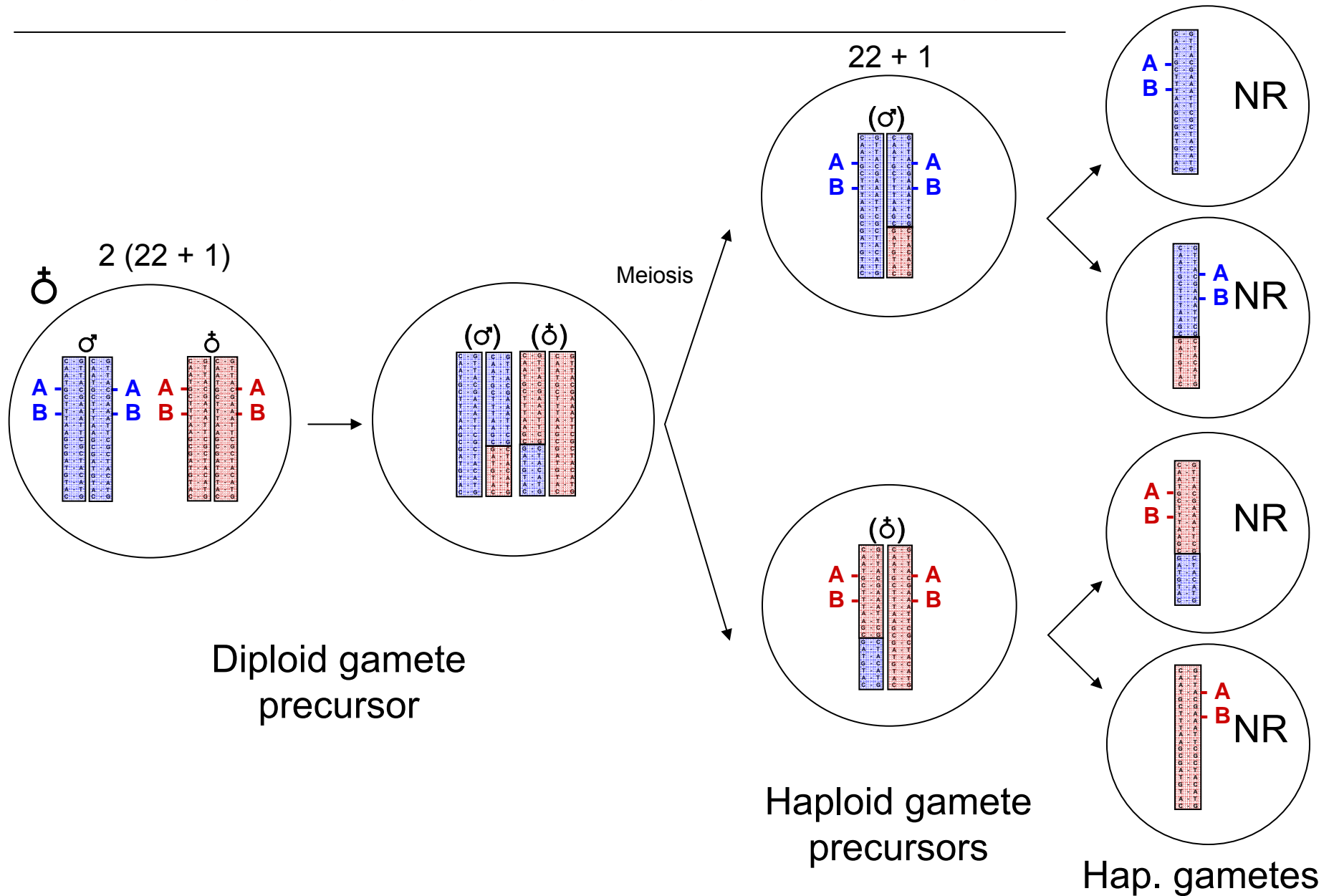
DNA organization



DNA recombination



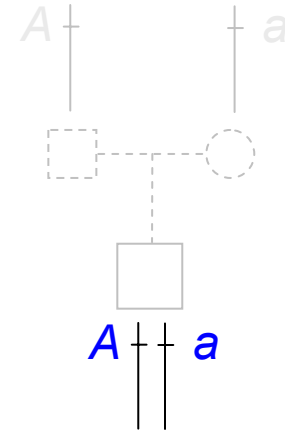
DNA recombination between linked loci



B. Population level

1. Allele frequencies

- ▷ A single locus, with two alleles
 - Biallelic
 - Single nucleotide polymorphism, SNP
- ▷ Alleles **A** and **a**
 - Frequency of **A** is **p**
 - Frequency of **a** is **q = 1 - p**
- ▷ Every individual inherits two alleles
 - A genotype is the combination of the two alleles
 - e.g. **AA**, **aa** (the homozygotes) or **Aa** (the heterozygote)



B. Population level

2. Genotype frequencies (Random mating)

		Allele 1	
		A (p)	a (q)
Allele 2	A (p)	AA (p^2)	Aa (pq)
	a (q)	aA (qp)	aa (q^2)

Hardy-Weinberg Equilibrium frequencies

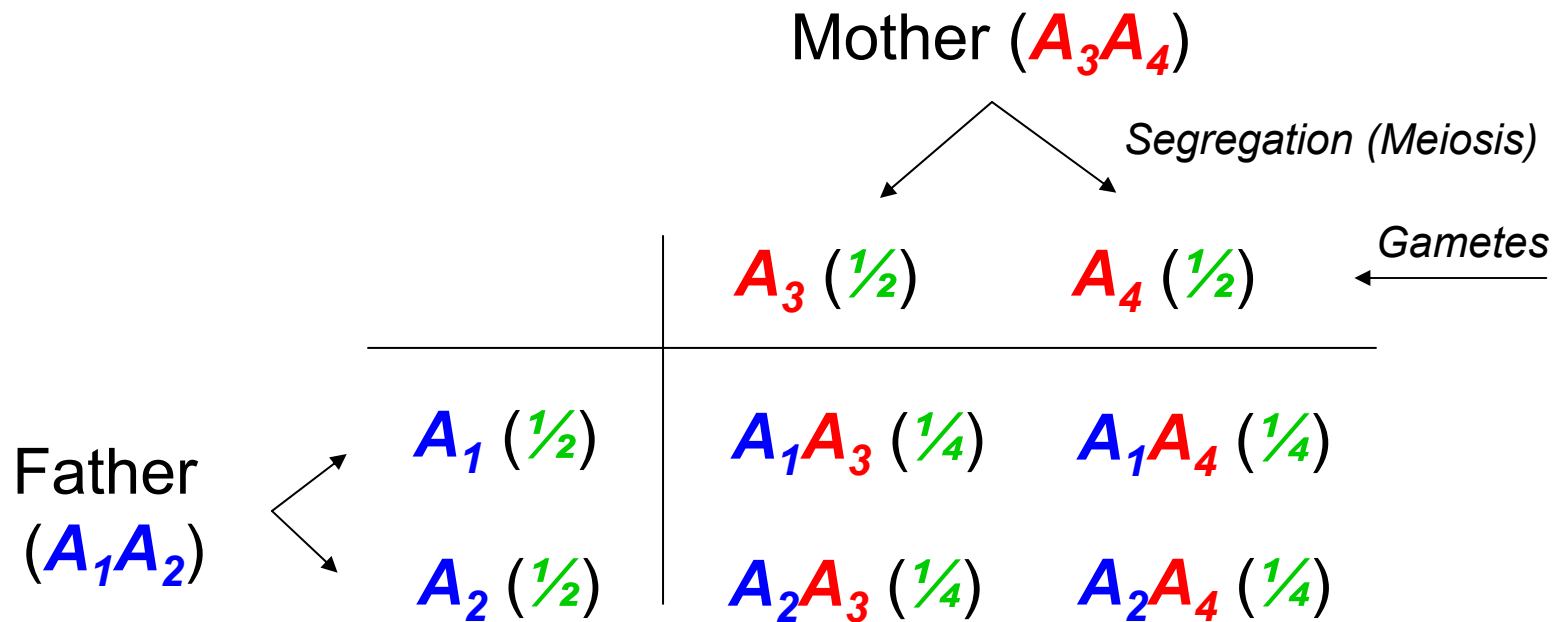
$$P(AA) = p^2$$

$$P(Aa) = 2pq \qquad p^2 + 2pq + q^2 = 1$$

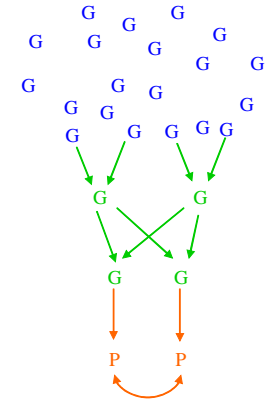
$$P(aa) = q^2$$

C. Transmission level

Mendel's law of segregation



D. Phenotype level



1. Classical Mendelian traits

▷ Dominant trait

- **AA, Aa** **1**
- **aa** **0**

▷ Recessive trait

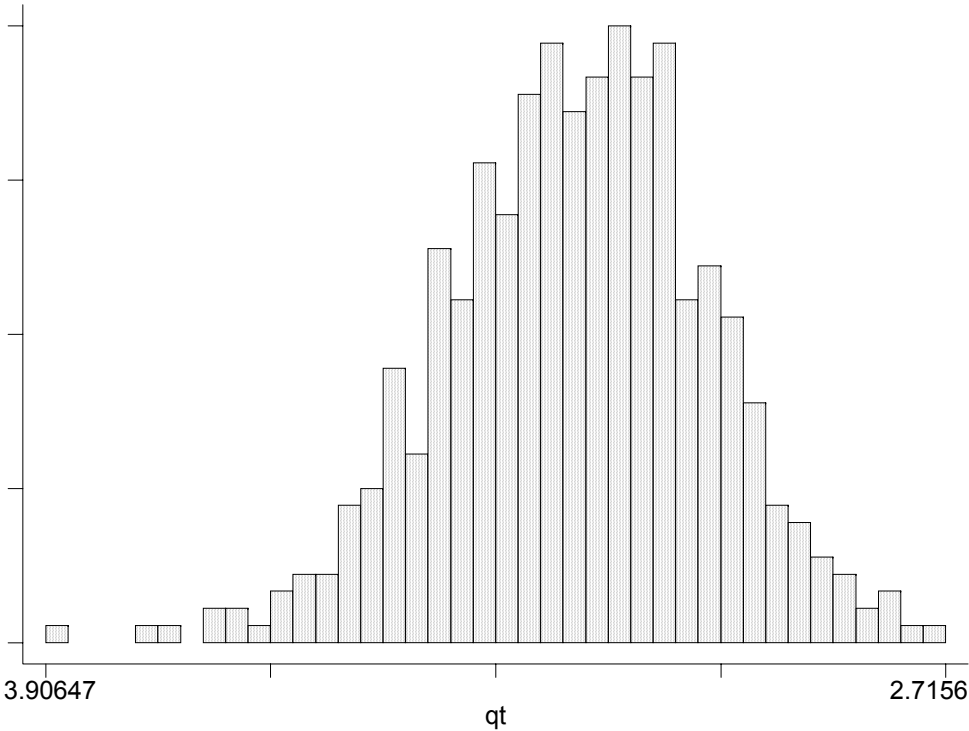
- **AA, Aa** **0**
- **aa** **1**

▷ Codominant trait

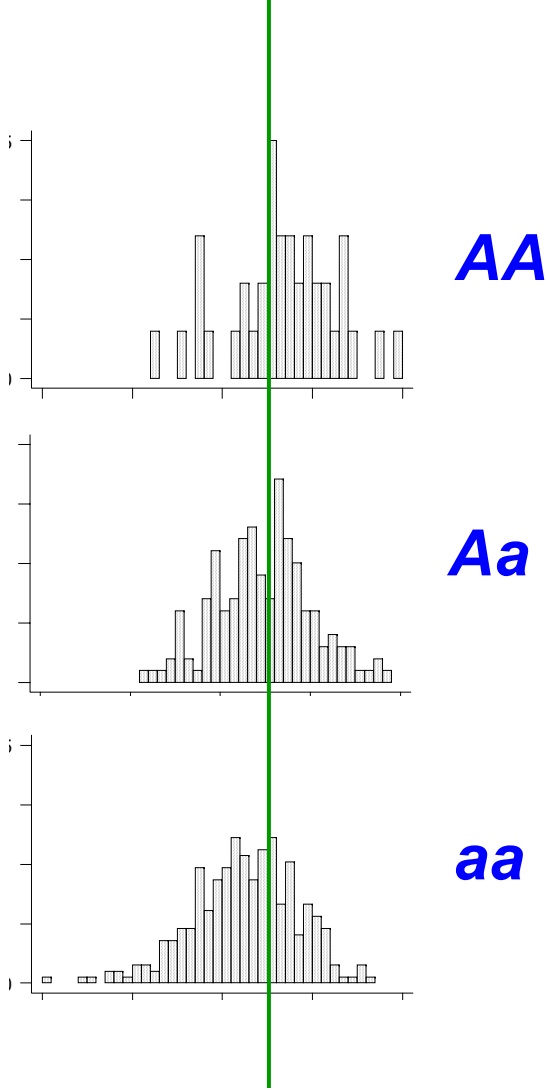
- **AA** **X**
- **Aa** **Y**
- **aa** **Z**

D. Phenotype level

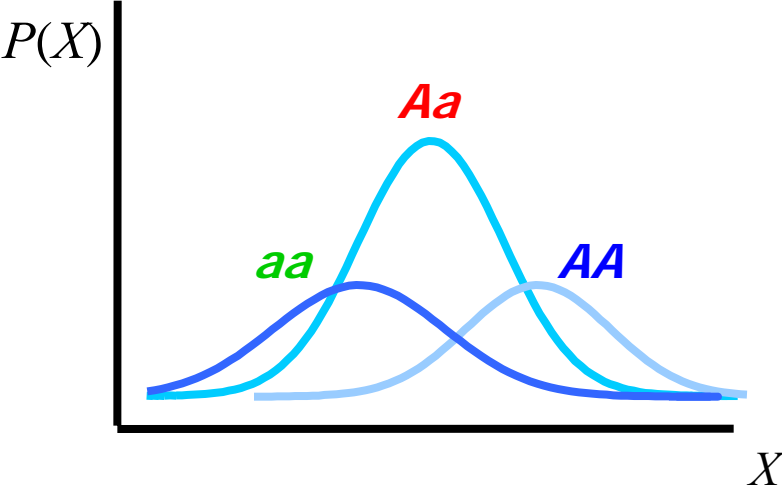
2. Quantitative traits



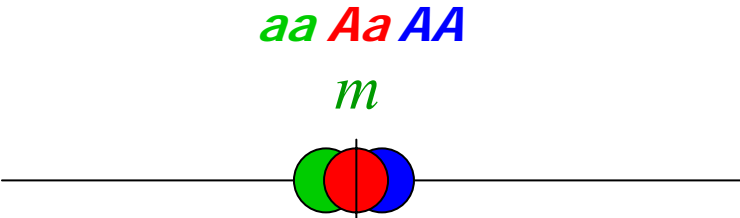
e.g. cholesterol levels



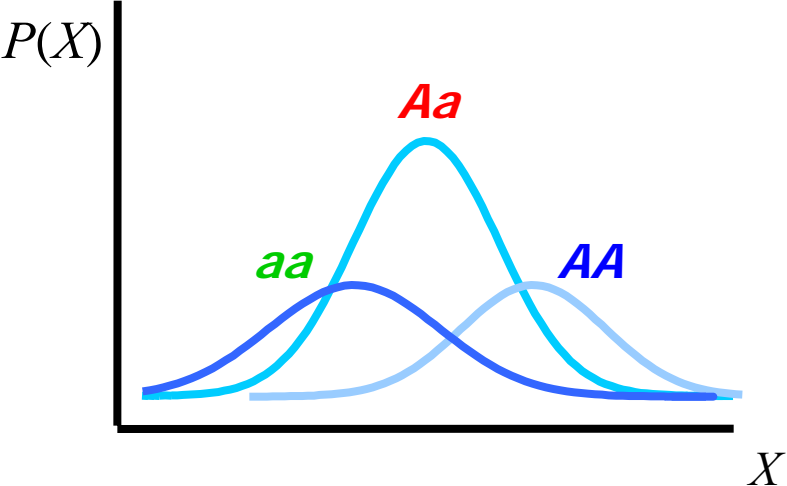
D. Phenotype level



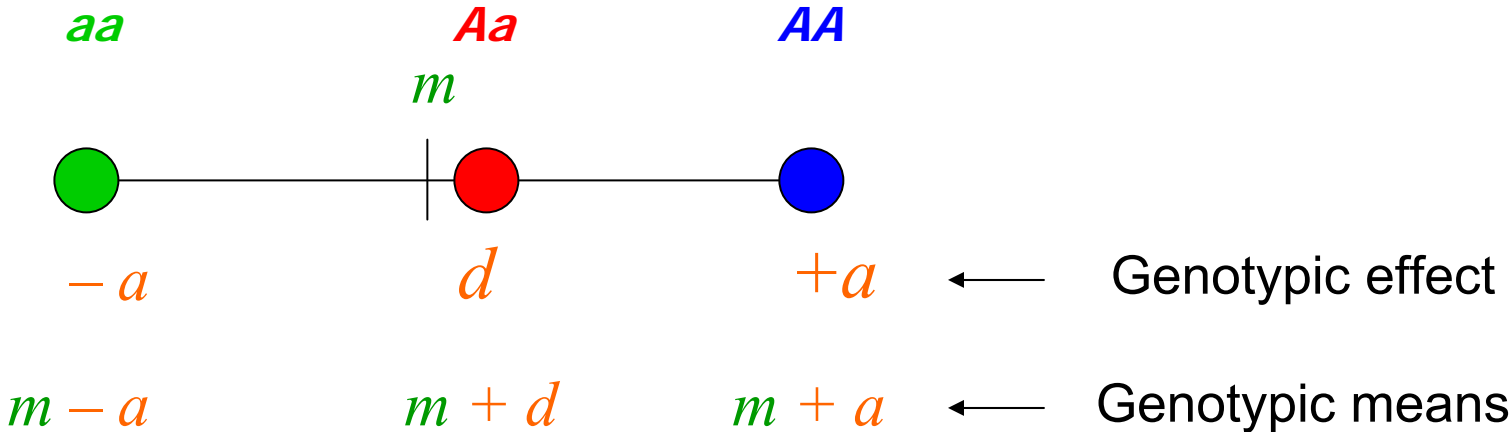
Biometric Model



D. Phenotype level



Biometric Model



3. Very basic statistical concepts

Mean, variance, covariance

1. Mean (X)

$$\mu(X) = \frac{\sum_i x_i}{n}$$
$$= \sum_i x_i f(x_i)$$

X
—
x₁
x₂
x₃
x₄
...
x_n

Mean, variance, covariance

2. Variance (X)

$$\begin{aligned} \text{Var}(X) &= \frac{\sum_i (x_i - \mu)^2}{n-1} \\ &= \sum_i (x_i - \mu)^2 f(x_i) \end{aligned}$$

X	$X-\mu$	$(X-\mu)^2$
x_1	$x_1-\mu$	$(x_1-\mu)^2$
x_2	$x_2-\mu$	$(x_2-\mu)^2$
x_3	$x_3-\mu$	$(x_3-\mu)^2$
x_4	$x_4-\mu$	$(x_4-\mu)^2$
...
x_n	$x_n-\mu$	$(x_n-\mu)^2$

Mean, variance, covariance

3. Covariance (X, Y)

$$\begin{aligned} \text{Cov}(X, Y) &= \frac{\sum_i (x_i - \mu_X)(y_i - \mu_Y)}{n-1} \\ &= \sum_i (x_i - \mu_X)(y_i - \mu_Y) f(x_i, y_i) \end{aligned}$$

X	Y	X- μ_X	Y- μ_Y
x ₁	y ₁	x ₁ - μ_X	y ₁ - μ_Y
x ₂	y ₂	x ₂ - μ_X	y ₂ - μ_Y
x ₃	y ₃	x ₃ - μ_X	y ₃ - μ_Y
x ₄	y ₄	x ₄ - μ_X	y ₄ - μ_Y
...
x _n	y _n	x _n - μ_X	y _n - μ_Y

4. Biometrical model

Biometrical model for single biallelic QTL

- ▷ Biallelic locus
 - Genotypes: **AA, Aa, aa**
 - Genotype frequencies: **$p^2, 2pq, q^2$**
- ▷ Alleles at this locus are transmitted from P-O according to Mendel's law of segregation
- ▷ Genotypes for this locus influence the expression of a quantitative trait X (i.e. locus is a QTL)



Biometrical genetic model that estimates the contribution of this QTL towards the **(1) Mean**, **(2) Variance** and **(3) Covariance between individuals** for this quantitative trait X

Biometrical model for single biallelic QTL

1. Contribution of the QTL to the Mean (X)

e.g. cholesterol levels in the population

$$\mu = \sum_i x_i f(x_i)$$

Genotypes	AA	Aa	aa
Effect, x	a	d	$-a$
Frequencies, $f(x)$	p^2	$2pq$	q^2

$$\text{Mean } (X) = a(p^2) + d(2pq) - a(q^2) = a(p-q) + 2pqd$$

Biometrical model for single biallelic QTL

2. Contribution of the QTL to the Variance (X)

$$Var = \sum_i (x_i - \mu)^2 f(x_i)$$

Genotypes	AA	Aa	aa
Effect, x	a	d	$-a$
Frequencies, $f(x)$	p^2	$2pq$	q^2

$$\begin{aligned} Var(X) &= (a-m)^2 p^2 + (d-m)^2 2pq + (-a-m)^2 q^2 \\ &= V_{QTL} \end{aligned}$$

Broad-sense heritability of X at this locus = V_{QTL} / V_{Total}

Biometrical model for single biallelic QTL

$$\text{Var}(X) = (a-m)^2 p^2 + (d-m)^2 2pq + (-a-m)^2 q^2$$

$$m = a(p-q) + 2pqd \quad = \frac{2pq[a+(q-p)d]^2}{2pq} + \frac{(2pqd)^2}{2pq}$$

$$= V_{A_{QTL}} + V_{D_{QTL}}$$

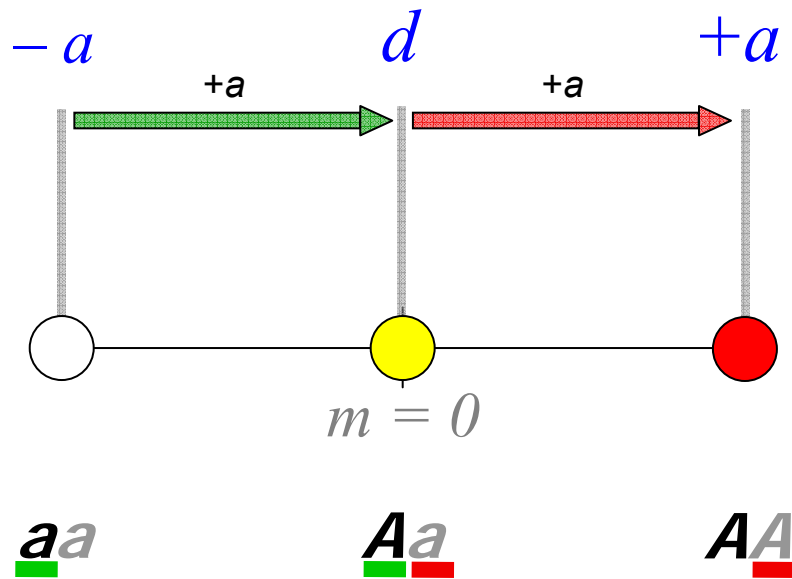
Demonstration: final 3 slides

Additive effects: the main effects of individual alleles

Dominance effects: represent the interaction between alleles

Biometrical model for single biallelic QTL

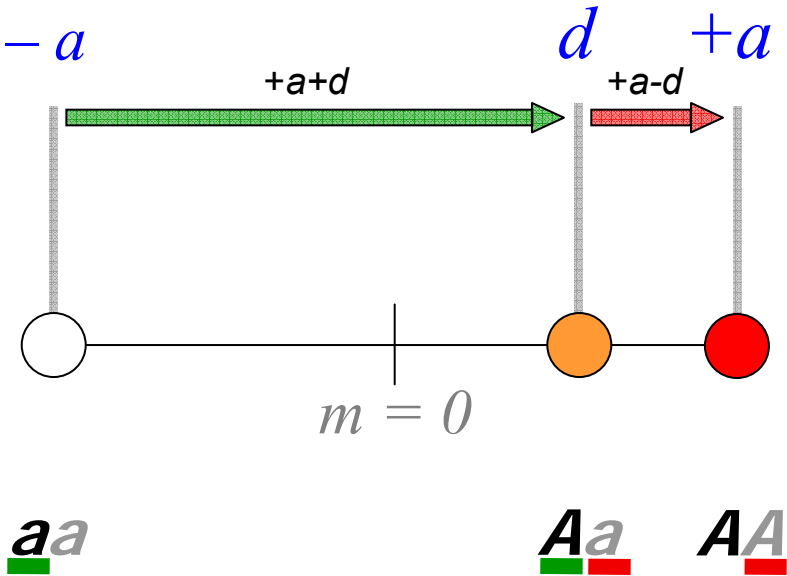
$$d = 0$$



Additive

Biometrical model for single biallelic QTL

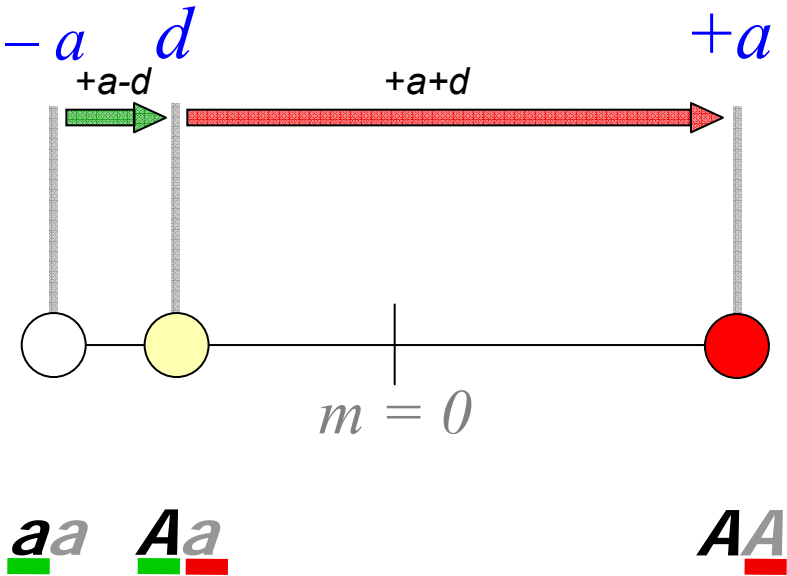
$$d > 0$$



Dominant

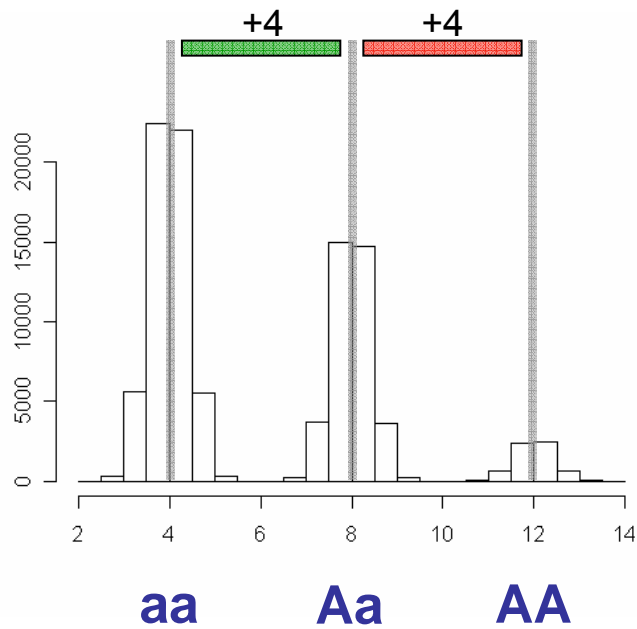
Biometrical model for single biallelic QTL

$d < 0$

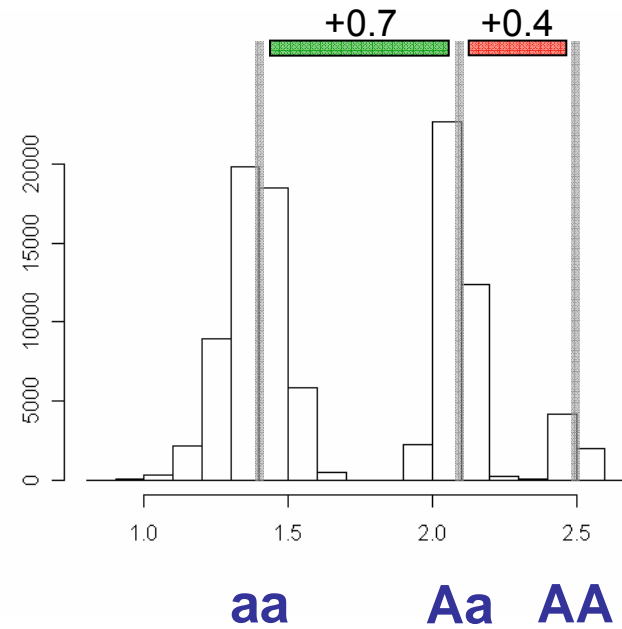


Recessive

Statistical definition of dominance is scale dependent



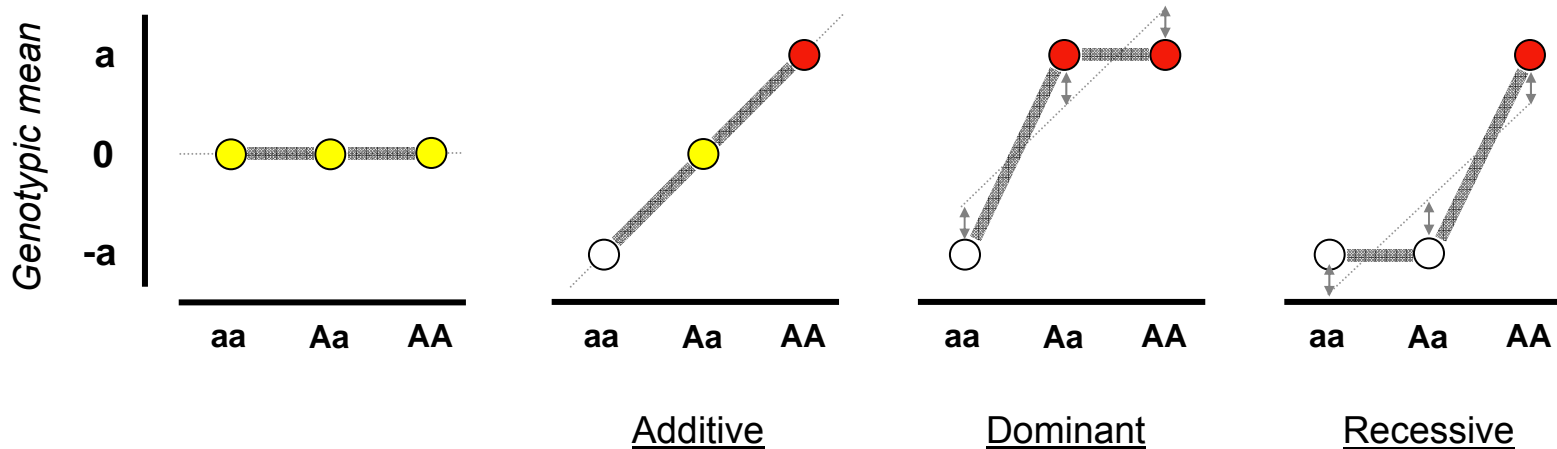
$\log(x)$



No departure from
additivity

Significant departure
from additivity

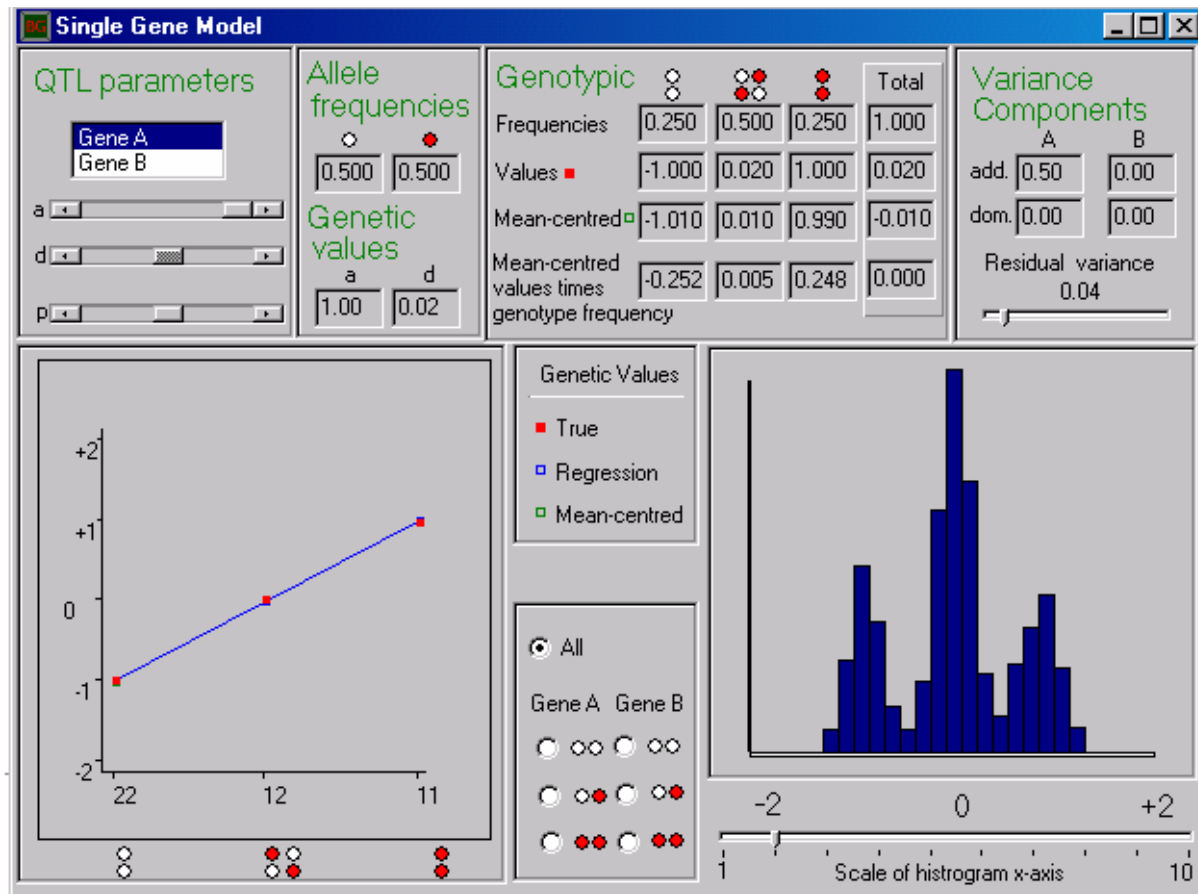
Biometrical model for single biallelic QTL



$$\begin{aligned}\text{Var}(X) &= \text{Regression Variance} + \text{Residual Variance} \\ &= \text{Additive Variance} + \text{Dominance Variance} \\ &= V_{A_{QTL}} + V_{D_{QTL}}\end{aligned}$$

Practical

H:\ferreira\GeneticTheory\sgene.exe



Practical

- ▷ **Aim** Visualize graphically how allele frequencies, genetic effects, dominance, etc, influence trait mean and variance

Ex1

$a=0$, $d=0$, $p=0.4$, Residual Variance = 0.04, Scale = 2.
Vary \underline{a} from 0 to 1.

Ex2

$a=1$, $d=0$, $p=0.4$, Residual Variance = 0.04, Scale = 2.
Vary \underline{d} from -1 to 1.

Ex3

$a=1$, $d=0$, $p=0.4$, Residual Variance = 0.04, Scale = 2.
Vary \underline{p} from 0 to 1.

Look at scatter-plot, histogram and variance components.

Some conclusions

1. Additive genetic variance depends on

allele frequency p

& *additive genetic value* a

as well as

dominance deviation d

2. Additive genetic variance typically greater than dominance variance

Biometrical model for single biallelic QTL

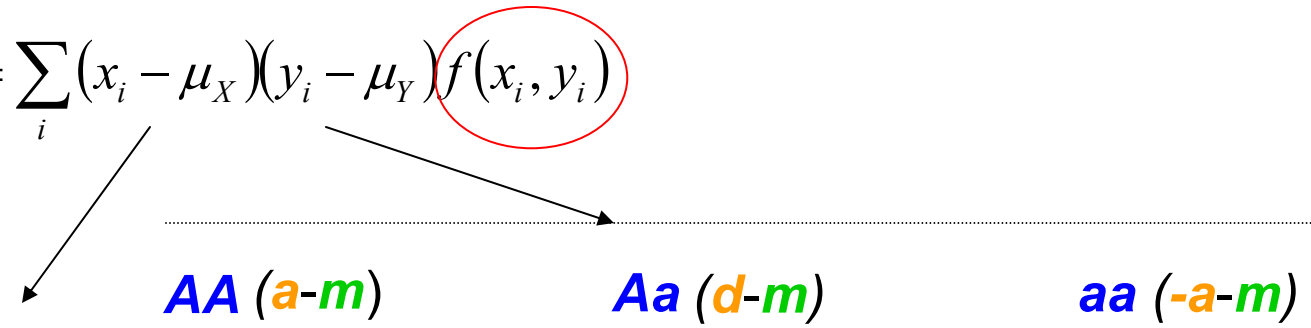
1. Contribution of the QTL to the Mean (X)

2. Contribution of the QTL to the Variance (X)

3. Contribution of the QTL to the Covariance (X, Y)

Biometrical model for single biallelic QTL

3. Contribution of the QTL to the Cov (X, Y)

$$\text{Cov}(X, Y) = \sum_i (x_i - \mu_X)(y_i - \mu_Y) f(x_i, y_i)$$


AA (**a-m**)

Aa (**d-m**)

aa (**-a-m**)

AA (**a-m**)

(**a-m**)²

Aa (**d-m**)

(**a-m**) (**d-m**)

(**d-m**)²

aa (**-a-m**)

(**a-m**) (**-a-m**)

(**d-m**) (**-a-m**)

(**-a-m**)²

Biometrical model for single biallelic QTL

3A. Contribution of the QTL to the Cov (X, Y) – MZ twins

$$\text{Cov}(X, Y) = \sum_i (x_i - \mu_X)(y_i - \mu_Y) f(x_i, y_i)$$

	AA (a-m)	Aa (d-m)	aa (-a-m)
AA (a-m)	$p^2(a-m)^2$		
Aa (d-m)	0 (a-m) (d-m)	$2pq(d-m)^2$	
aa (-a-m)	0 (a-m) (-a-m)	0 (d-m) (-a-m)	$q^2(-a-m)^2$

$$\begin{aligned} \text{Cov}(X, Y) &= (a-m)^2 p^2 + (d-m)^2 2pq + (-a-m)^2 q^2 \\ &= 2pq[a + (q-p)d]^2 + (2pqd)^2 = V_{A_{QTL}} + V_{D_{QTL}} \end{aligned}$$

Biometrical model for single biallelic QTL

3B. Contribution of the QTL to the Cov (X, Y) – Parent-Offspring

	AA $(a-m)$	Aa $(d-m)$	aa $(-a-m)$
AA $(a-m)$	$p^3(a-m)^2$		
Aa $(d-m)$	$p^2q(a-m)(d-m)$	$pq(d-m)^2$	
aa $(-a-m)$	$0(a-m)(-a-m)$	$pq^2(d-m)(-a-m)$	$q^3(-a-m)^2$

- e.g. given an AA father, an AA offspring can come from either $AA \times AA$ or $AA \times Aa$ parental mating types

$AA \times AA$ will occur $p^2 \times p^2 = p^4$
and have AA offspring Prob() $=1$

$AA \times Aa$ will occur $p^2 \times 2pq = 2p^3q$
and have AA offspring Prob() $=0.5$
and have Aa offspring Prob() $=0.5$

$$\begin{aligned} \text{Therefore, P}(AA \text{ father \& } AA \text{ offspring}) &= p^4 + p^3q \\ &= p^3(p+q) \\ &= p^3 \end{aligned}$$

Biometrical model for single biallelic QTL

3B. Contribution of the QTL to the Cov (X, Y) – Parent-Offspring

	AA ($a-m$)	Aa ($d-m$)	aa ($-a-m$)
AA ($a-m$)	$p^3(a-m)^2$		
Aa ($d-m$)	$p^2q(a-m)(d-m)$	$pq(d-m)^2$	
aa ($-a-m$)	$0(a-m)(-a-m)$	$pq^2(d-m)(-a-m)$	$q^3(-a-m)^2$

$$\begin{aligned} \text{Cov}(X, Y) &= (a-m)^2 p^3 + \dots + (-a-m)^2 q^3 \\ &= pq[a + (q-p)d]^2 = \frac{1}{2} V_{A_{QTL}} \end{aligned}$$

Biometrical model for single biallelic QTL

3C. Contribution of the QTL to the Cov (X, Y) – Unrelated individuals

	AA (a-m)	Aa (d-m)	aa (-a-m)
AA (a-m)	$p^4(a-m)^2$		
Aa (d-m)	$2p^3q(a-m)(d-m)$	$4p^2q^2(d-m)^2$	
aa (-a-m)	$p^2q^2(a-m)(-a-m)$	$2pq^3(d-m)(-a-m)$	$q^4(-a-m)^2$

$$\begin{aligned} \text{Cov}(X, Y) &= (a-m)^2 p^4 + \dots + (-a-m)^2 q^4 \\ &= 0 \end{aligned}$$

Biometrical model for single biallelic QTL

3D. Contribution of the QTL to the Cov (X, Y) – DZ twins and full sibs

	$\frac{1}{4}$ genome	$\frac{1}{4}$ genome	$\frac{1}{4}$ genome	$\frac{1}{4}$ genome
# identical alleles inherited from parents	2	1 (father)	1 (mother)	0
	$\frac{1}{4}$ (2 alleles)	+ $\frac{1}{2}$ (1 allele)	+ $\frac{1}{4}$ (0 alleles)	
	<i>MZ twins</i>	<i>P-O</i>	<i>Unrelateds</i>	

$$\begin{aligned}
 \text{Cov}(X, Y) &= \frac{1}{4} \text{Cov}(MZ) + \frac{1}{2} \text{Cov}(P-O) + \frac{1}{4} \text{Cov}(Unrel) \\
 &= \frac{1}{4}(V_{A_{QTL}} + V_{D_{QTL}}) + \frac{1}{2}(\frac{1}{2} V_{A_{QTL}}) + \frac{1}{4}(0) \\
 &= \frac{1}{2} V_{A_{QTL}} + \frac{1}{4} V_{D_{QTL}}
 \end{aligned}$$

Summary so far...

- ▷ Biometrical model predicts contribution of a QTL to the mean, variance and covariances of a trait

$$\text{Mean } (X) = a(p-q) + 2pqd \quad \leftarrow \text{Association analysis}$$

$$\text{Var } (X) = V_{A_{QTL}} + V_{D_{QTL}} \quad \leftarrow \text{Linkage analysis}$$

$$\text{Cov } (MZ) = V_{A_{QTL}} + V_{D_{QTL}}$$

$$\text{Cov } (DZ) = \frac{1}{2}V_{A_{QTL}} + \frac{1}{4}V_{D_{QTL}} \quad \text{On average!}$$

0, 1/2 or 1

0 or 1

For a sib-pair, do the two sibs have 0, 1 or 2 alleles in common?

IBD estimation / Linkage

5. Introduction to Linkage Analysis

For a heritable trait...

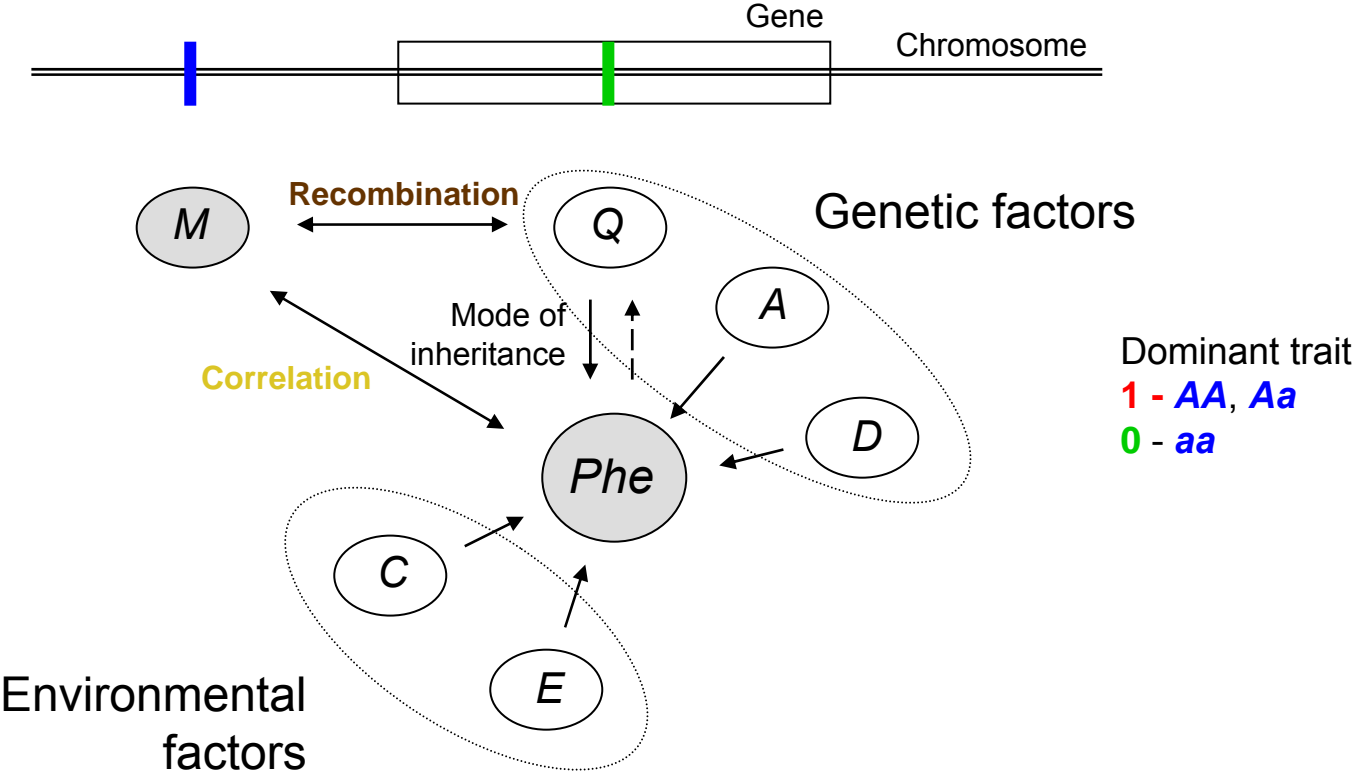
Linkage: localize region of the genome where a QTL that regulates the trait is likely to be harboured

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

Association: identify a QTL that regulates the trait

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

Linkage Analysis: Parametric vs. Nonparametric



Adapted from Weiss & Terwilliger 2000

Approach

▶ Parametric: genotypes marker locus & genotypes trait locus

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

Parameter of interest: θ between marker and trait loci

▶ Nonparametric: genotypes 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 phenotypic similarity and marker genotypic similarity

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

Phenotypic similarity between relatives

▶ Squared trait differences

$$(X_1 - X_2)^2$$

▶ Squared trait sums

$$(X_1 + X_2)^2$$

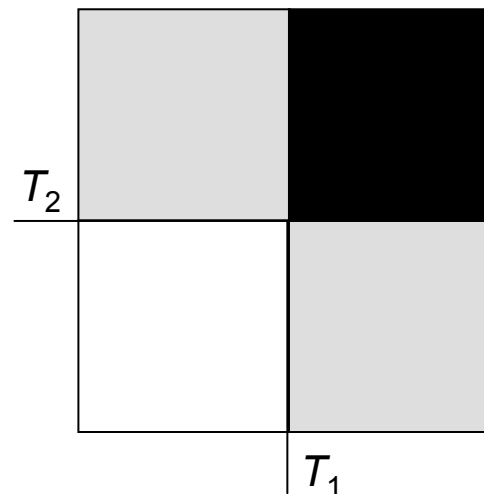
▶ Trait cross-product

$$[(X_1 - \mu) \cdot (X_2 - \mu)]$$

▶ Trait variance-covariance matrix

$$\begin{Bmatrix} \text{Var}(X_1) & \text{Cov}(X_1 X_2) \\ \text{Cov}(X_1 X_2) & \text{Var}(X_2) \end{Bmatrix}$$

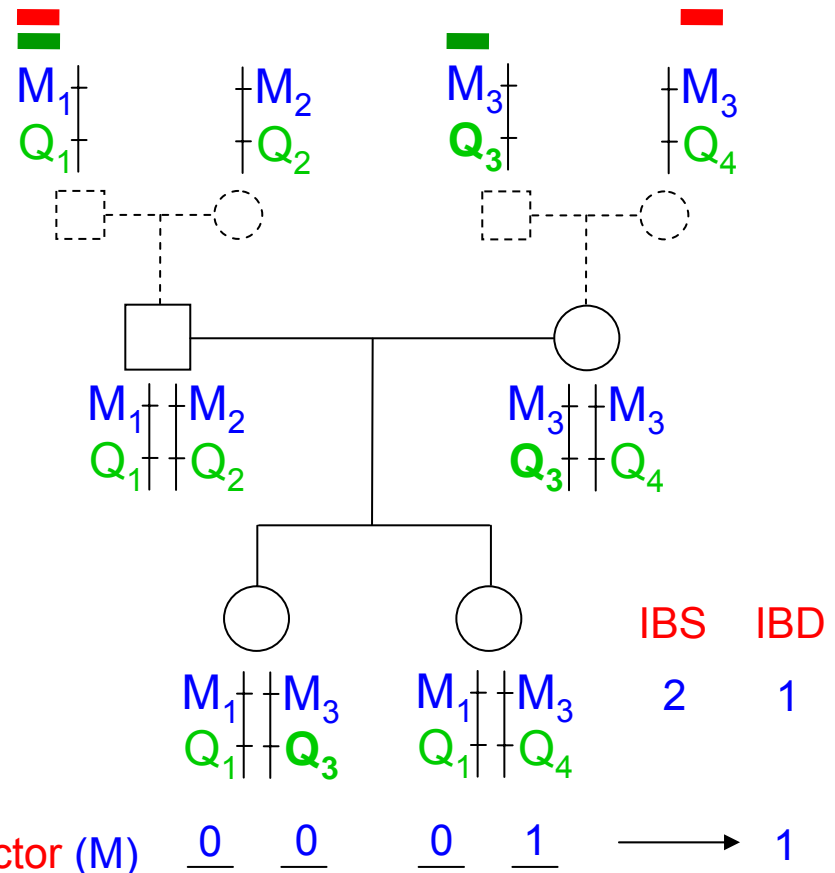
▶ Affection concordance



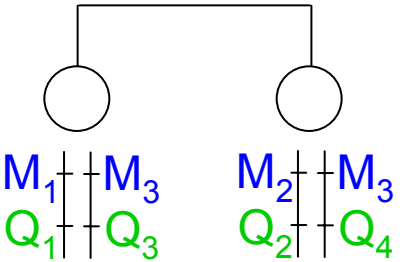
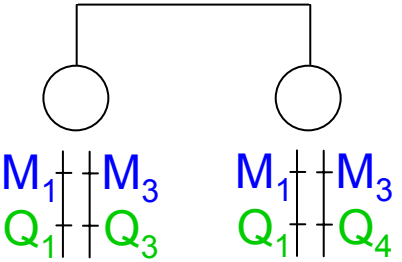
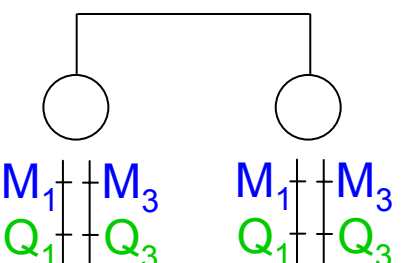
Genotypic similarity between relatives

▶ IBS Alleles shared Identical By State “look the same”, may have the same DNA sequence but they are not necessarily derived from a known common ancestor

▶ IBD Alleles shared Identical By Descent are a copy of the same ancestor allele

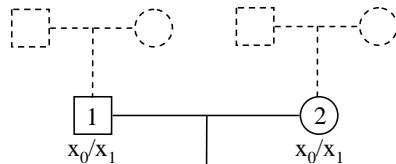


Genotypic similarity between relatives - π

	Inheritance vector (M)	Number of alleles shared IBD	Proportion of alleles shared IBD - π
	$\underline{0}$ $\underline{0}$ $\underline{1}$ $\underline{1}$	0	0
	$\underline{0}$ $\underline{0}$ $\underline{0}$ $\underline{1}$	1	0.5
	$\underline{0}$ $\underline{0}$ $\underline{0}$ $\underline{0}$	2	1

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

A B C D



2^{2n}

		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)

$$\hat{\pi} =$$

$$\text{Var}(X) = V_{A_{QTL}} + V_{D_{QTL}}$$

$$\text{Cov}(MZ) = V_{A_{QTL}} + V_{D_{QTL}}$$

$$\text{Cov}(DZ) = \frac{1}{2}V_{A_{QTL}} + \frac{1}{4}V_{D_{QTL}}$$

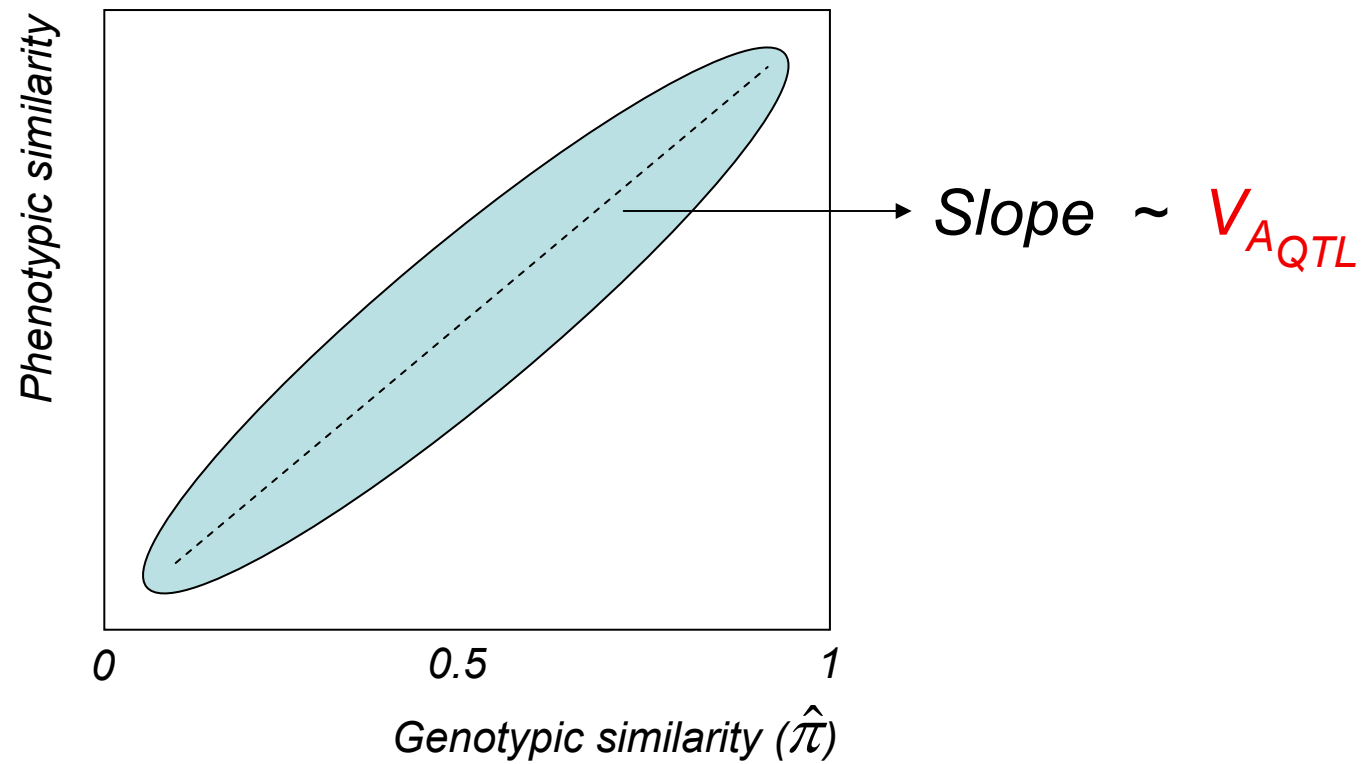
On average!

$$\text{Cov}(DZ) = \hat{\pi} \cdot V_{A_{QTL}} + \pi_2 \cdot V_{D_{QTL}}$$

For a given twin pair

$$\text{Cov}(DZ) = \hat{\pi} \cdot V_{A_{QTL}}$$

$$\text{Cov}(DZ) = V_{A_{QTL}} \cdot \hat{\pi}$$



Statistics that incorporate both phenotypic and genotypic similarities to test V_{QTL}

- ▶ Regression-based methods
Haseman-Elston, MERLIN-regress
- ▶ Variance components methods
Mx, MERLIN, SOLAR, GENEHUNTER

Biometrical model for single biallelic QTL

- ▷ Denote the average allelic effects
 - $\alpha_A = q(a+d(q-p))$
 - $\alpha_a = -p(a+d(q-p))$

- ▷ If only two alleles exist, we can define the *average effect of allele substitution*
 - $\alpha = \alpha_A - \alpha_a$
 - $\alpha = (q-(-p))(a+d(q-p)) = (a+d(q-p))$

- ▷ Therefore:
 - $\alpha_A = q\alpha$
 - $\alpha_a = -p\alpha$