

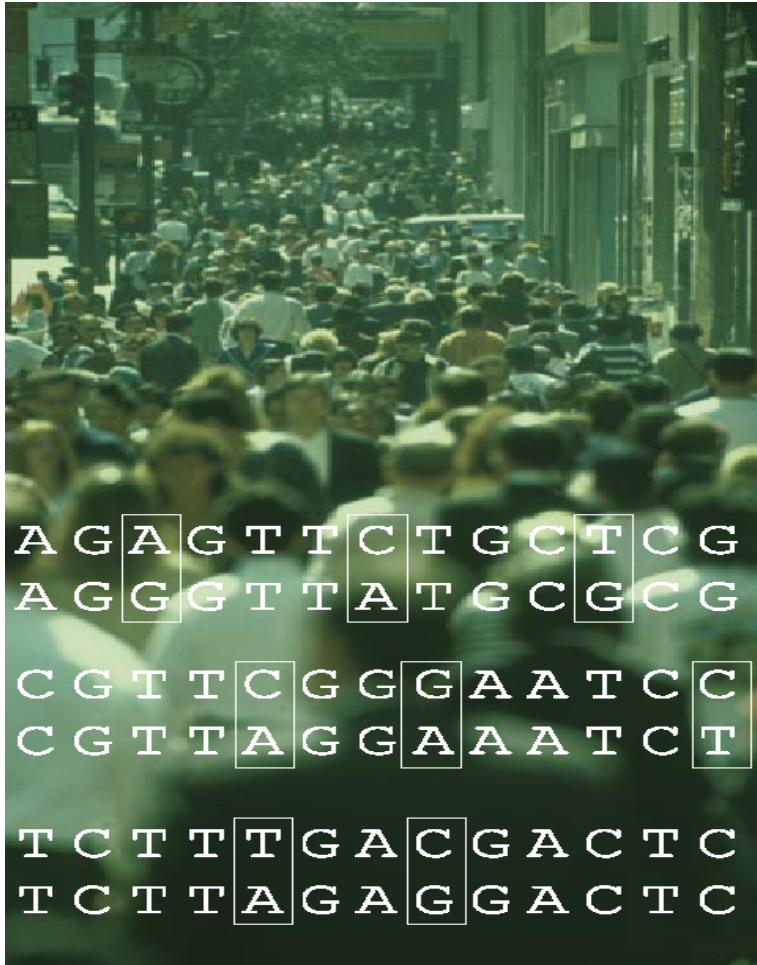
Linkage Analysis

02-710 Computational Genomics

Seyoung Kim

Genome Polymorphisms

Genetic Variation



Phenotypic Variation

The ABO Blood System

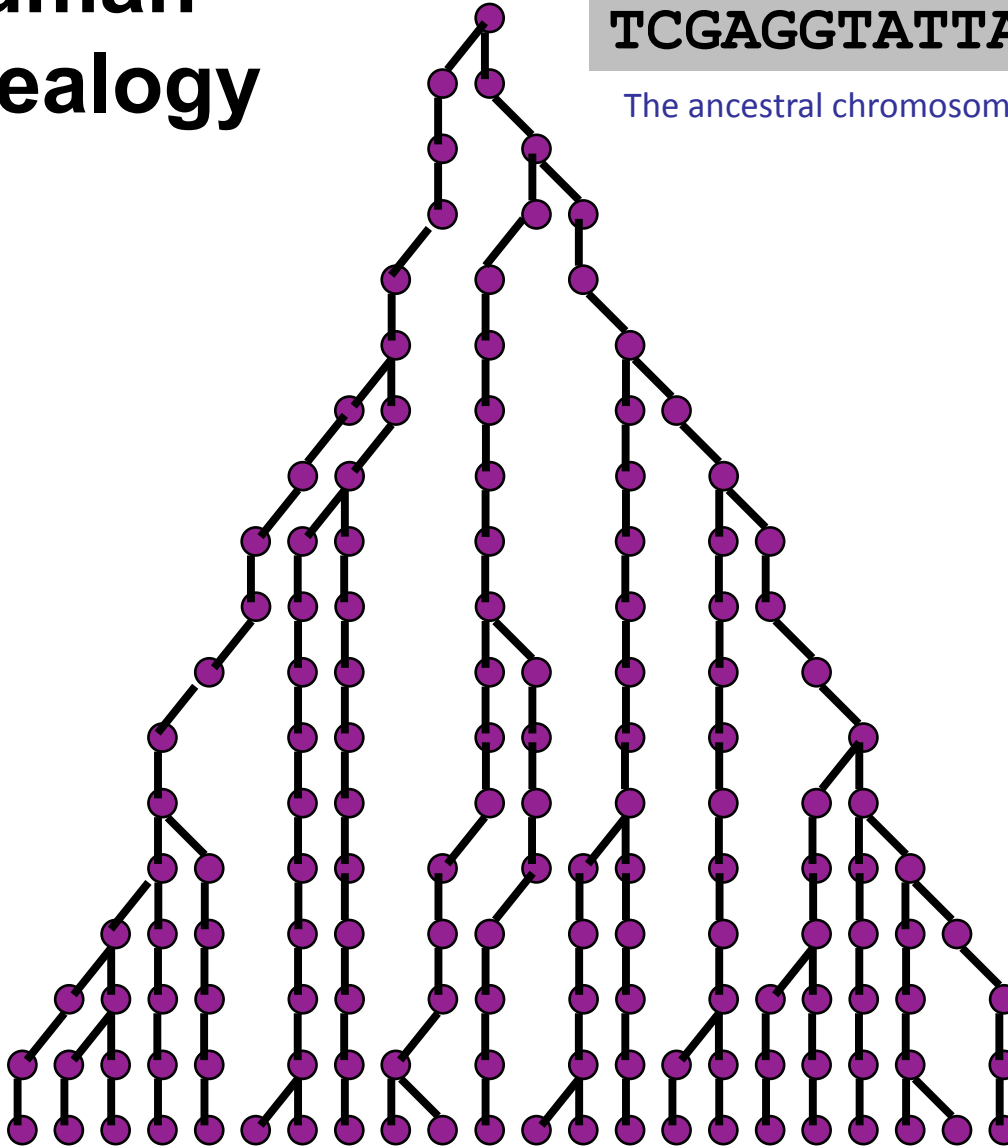
Blood Type (genotype)	Type A (AA, AO)	Type B (BB, BO)	Type AB (AB)	Type O (OO)
Red Blood Cell Surface Proteins (phenotype)	A agglutinogens only	B agglutinogens only	A and B agglutinogens	No agglutinogens
Plasma Antibodies (phenotype)	b agglutinin only	a agglutinin only	NONE	a and b agglutini



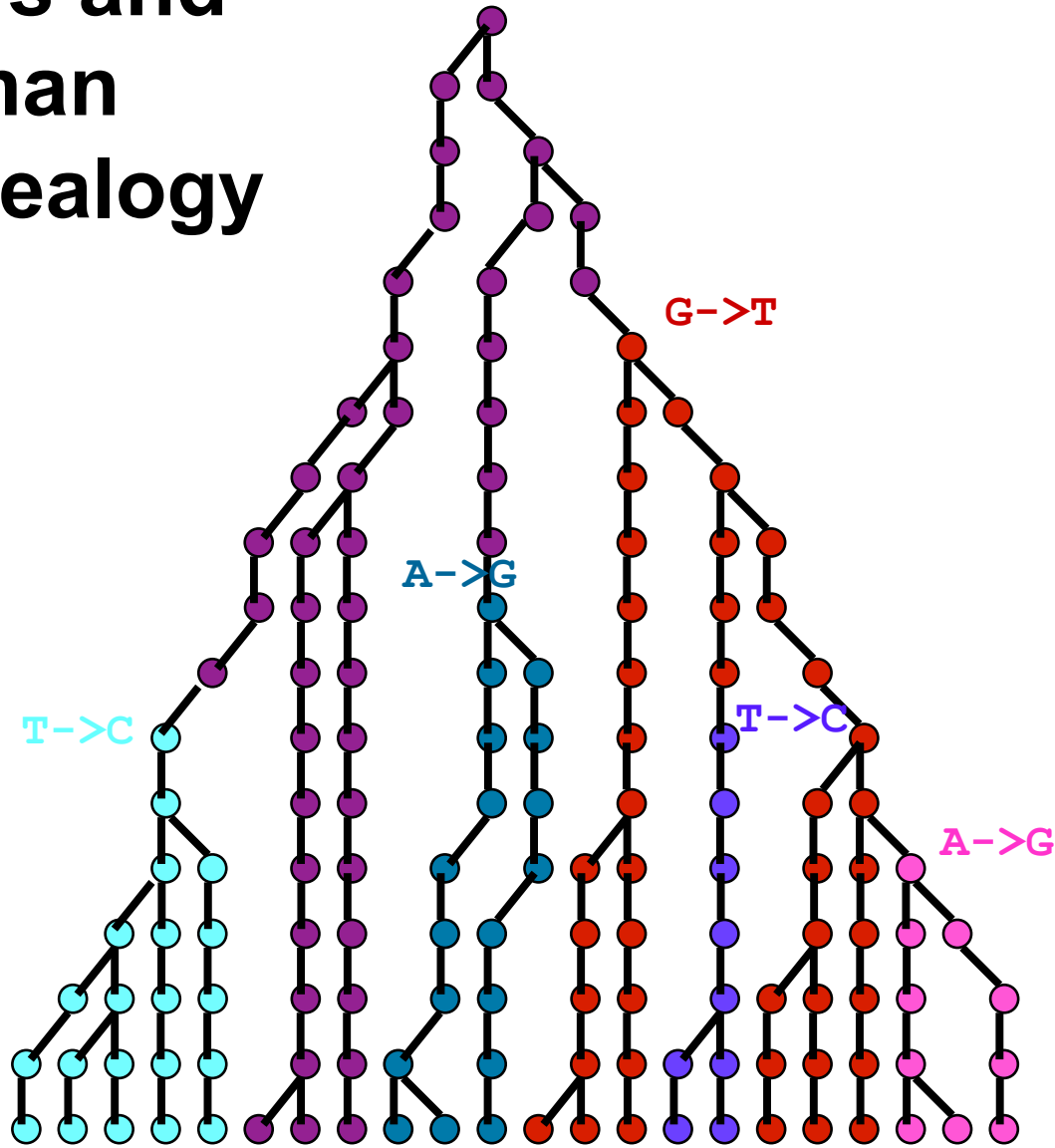
A Human Genealogy

TCGAGGTATTAAC

The ancestral chromosome



SNPs and Human Genealogy



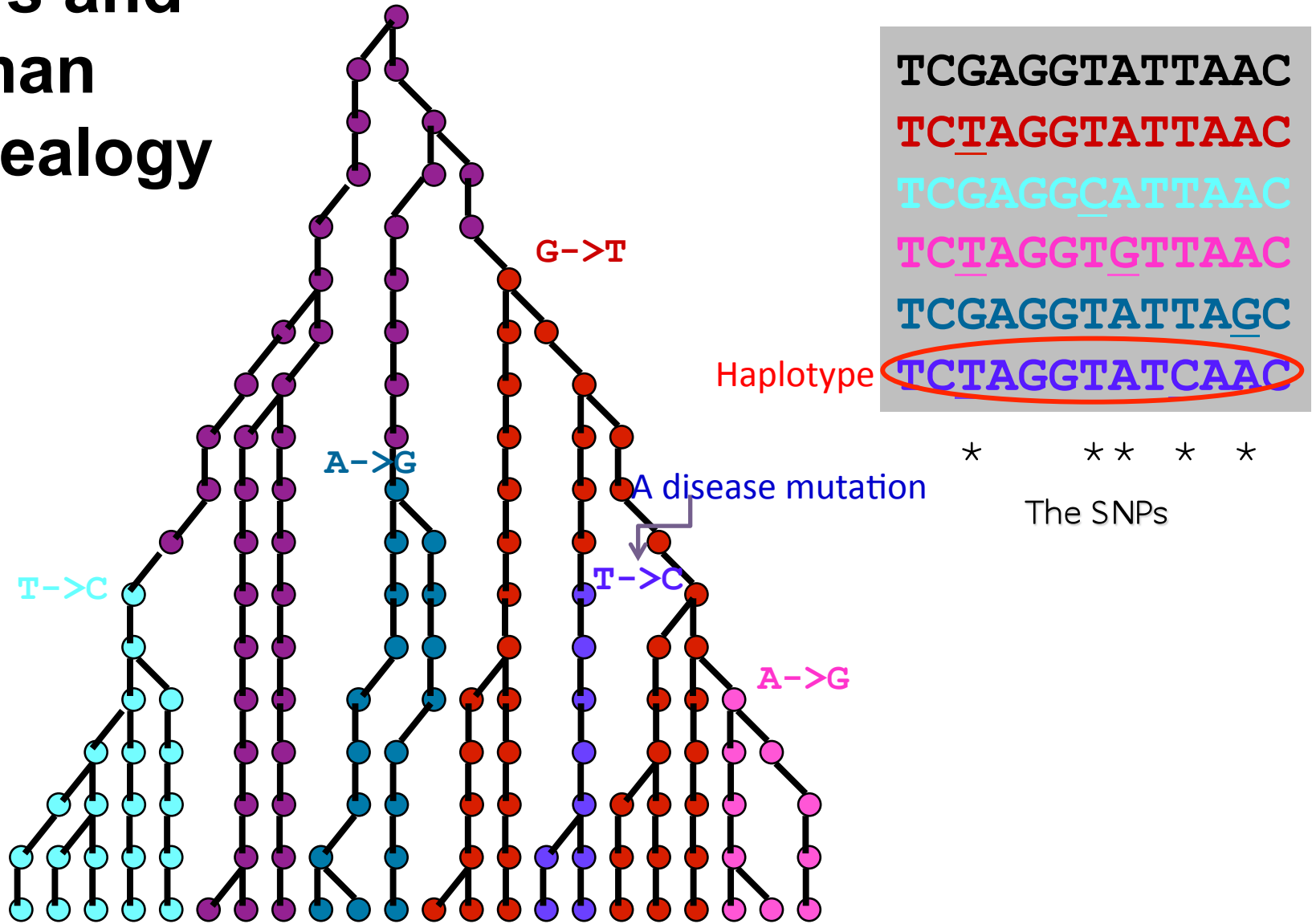
```

TCGAGGTATTAAC
TCTAGGTATTAAC
TCGAGGCATTAAC
TCTAGGTGTTAAC
TCGAGGTATTAGC
TCTAGGTATCAAC
  
```

* * * * *

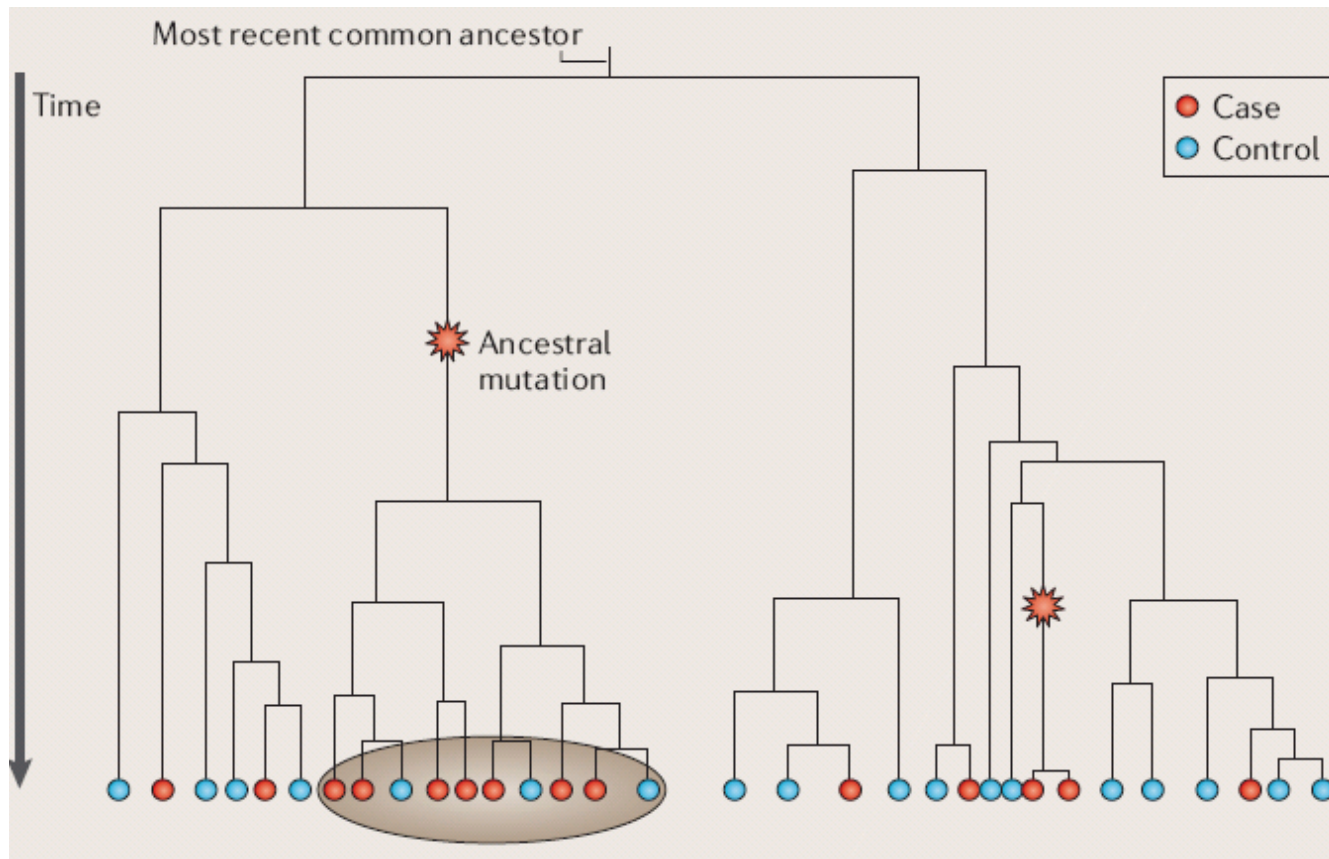
The SNPs

SNPs and Human Genealogy



Identifying Disease Loci

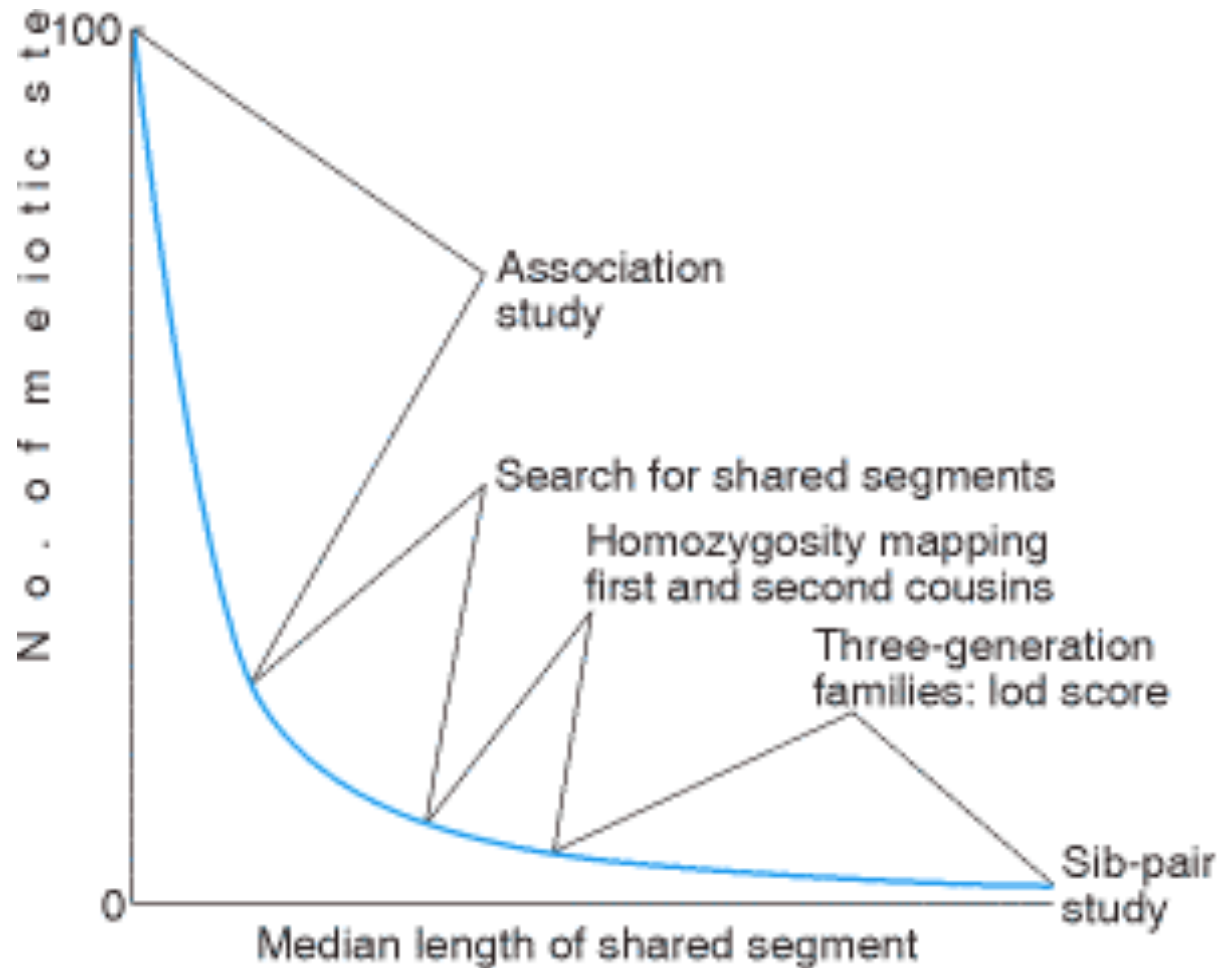
- All individuals are related if we go back far enough in the ancestry



Overview

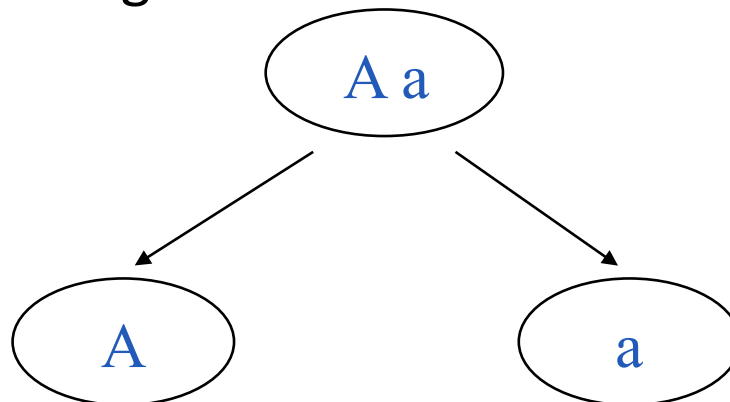
- How can we identify the genetic loci responsible for determining phenotypes?
 - Linkage analysis
 - Data are collected for family members
 - Difficult to collect data on a large number of families
 - Effective for rare diseases
 - Low resolution on the genomes due to only few recombinations
 - » a large region of linkage
 - Genome-wide association studies
 - Data are collected for unrelated individuals
 - Easier to find a large number of affected individuals
 - Effective for common diseases, compared to family-based method
 - Relatively high resolution for pinpointing the locus linked to the phenotype

Linkage Analysis vs. Association Analysis



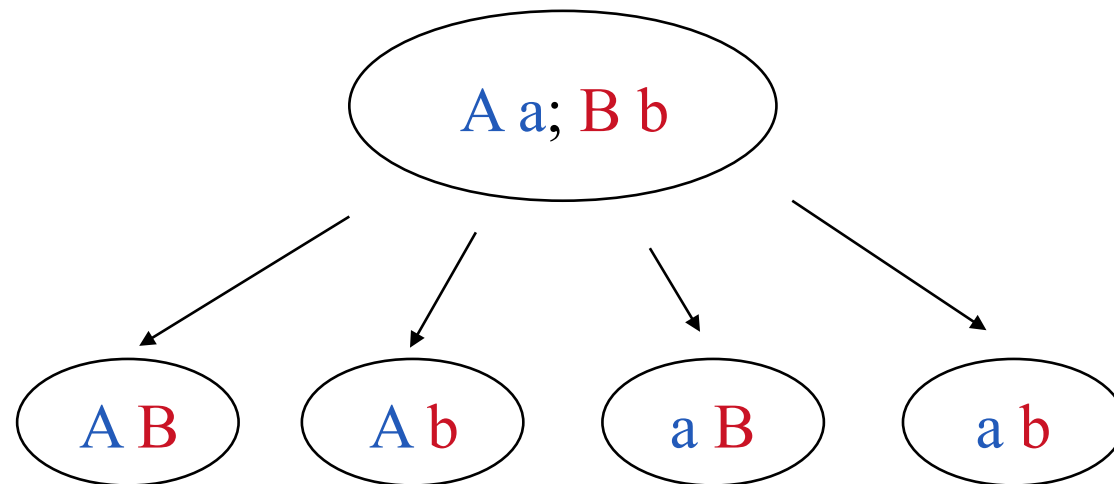
Mendel's two laws

- Modern genetics began with Mendel's experiments on garden peas. He studied seven contrasting pairs of characters, including:
 - The form of ripe seeds: round, wrinkled
 - The color of the seed albumen: **yellow**, **green**
 - The length of the stem: long, short
- **Mendel's first law:** Characters are controlled by pairs of genes which separate during the formation of the reproductive cells (meiosis)



Mendel's two laws

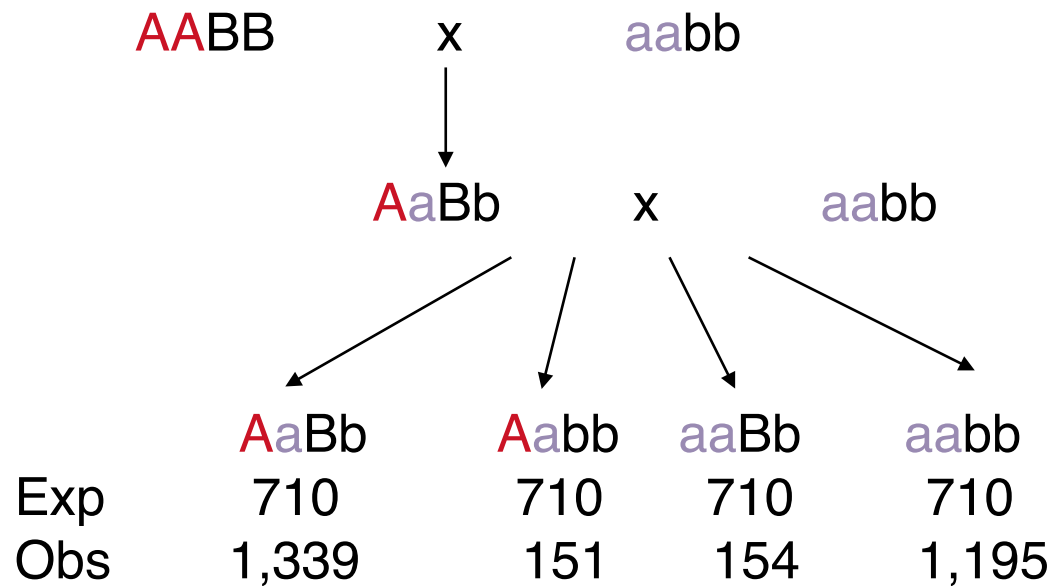
- **Mendel's second law:** When two or more pairs of genes segregate simultaneously, they do so independently.



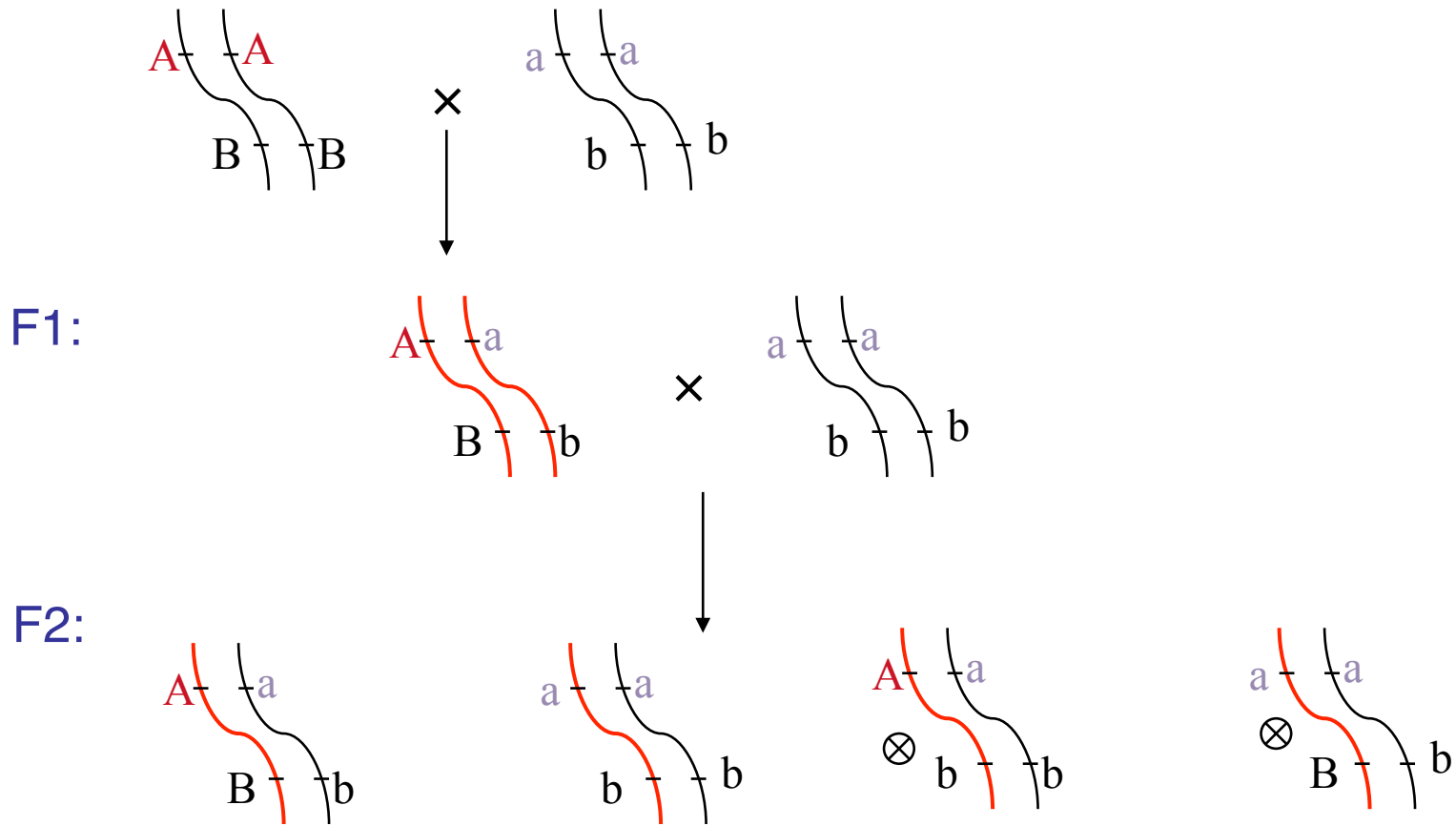
“Exceptions” to Mendel’s Second Law

Morgan’s fruitfly data (1909): 2,839 flies

Eye color **A**: red a: purple
 Wing length **B**: normal b: vestigial



Morgan's explanation



⊗ *Crossover has taken place*

Recombination

- *Parental types*: AaBb, aabb
- *Recombinants*: Aabb, aaBb
 - The proportion of recombinants between the two genes (or characters) is called the **recombination fraction** between these two genes.
- **Recombination fraction** It is usually denoted by r or θ . For Morgan's traits:

$$r = (151 + 154)/2839 = 0.107$$

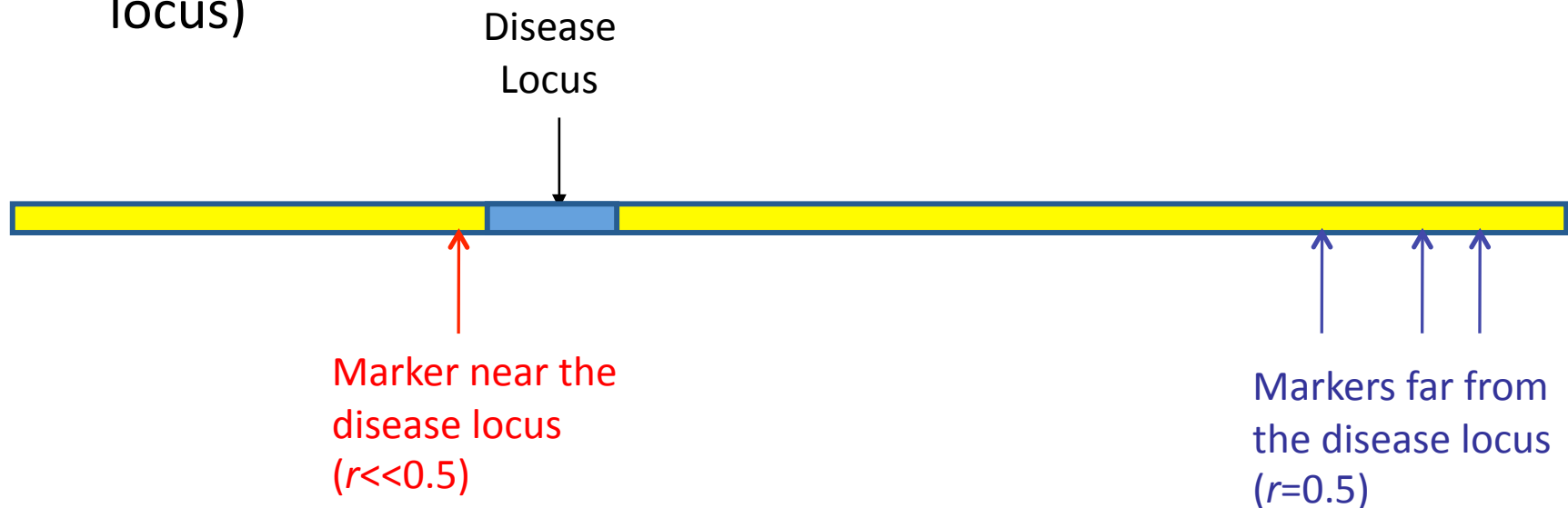
If $r < 1/2$: two genes are said to be **linked**.

If $r = 1/2$: independent segregation (Mendel's second law).

Now we move on to (small) pedigrees.

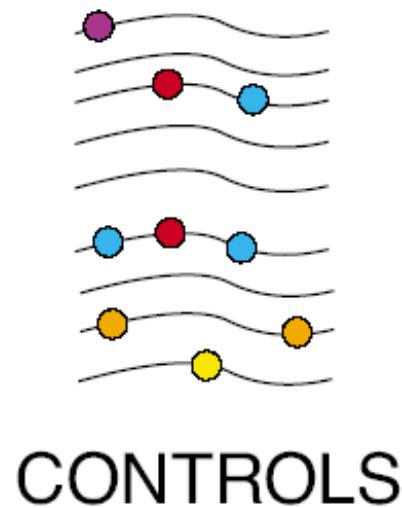
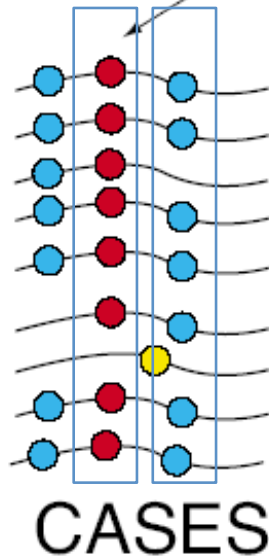
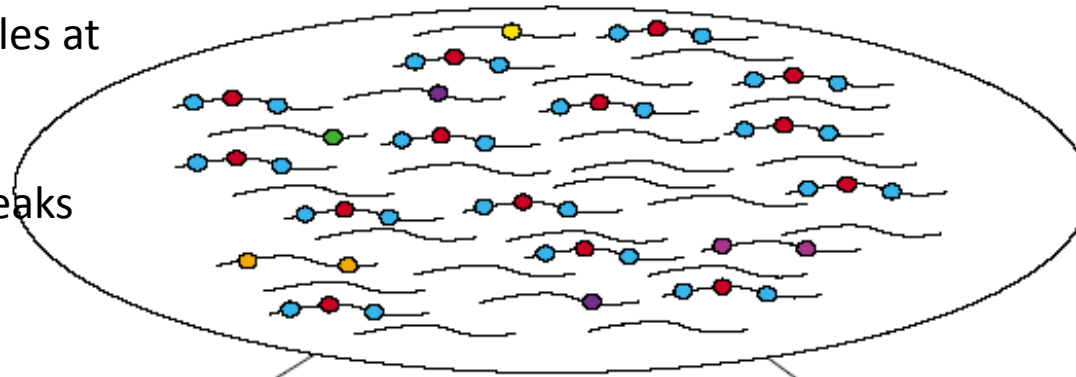
Linkage Analysis

- Goal: Identify the unknown disease locus
- Idea: Given pedigree data and a map of genetic markers, let's look for the markers that are linked to the unknown disease locus (i.e. linkage between the disease locus and the marker locus)



Linkage Disequilibrium in Gene Mapping

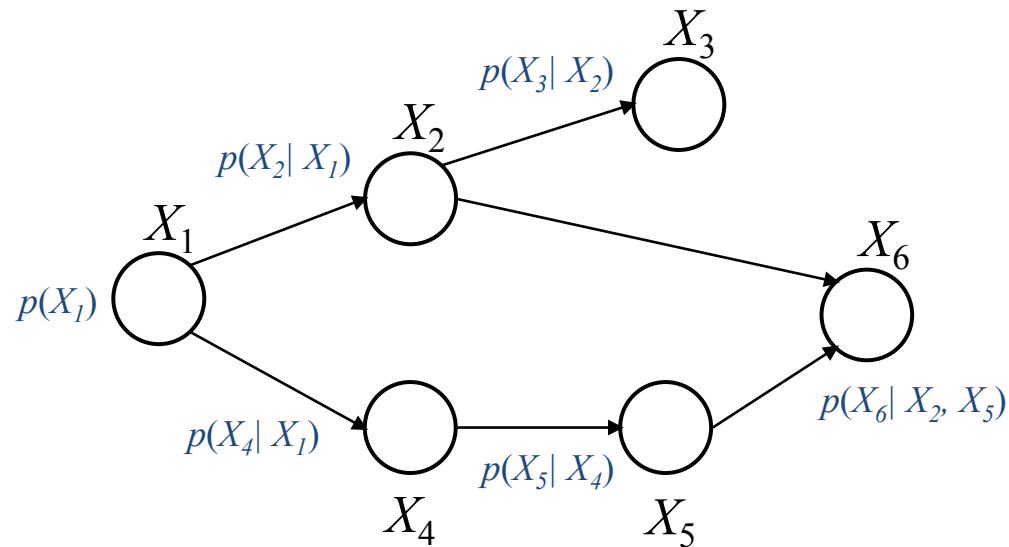
- LD is the non-random association of alleles at different loci
- Genetic recombination breaks down LD



Linkage Analysis

- Parametric Linkage Analysis
 - Need to specify the disease model
 - Compute LOD-score based on the model for each marker
 - Markers with the high LOD-scores are considered as linked to disease locus
 - Highly effective for Mendelian disease caused by a single locus
 - Usually based on a large pedigree

Probabilistic Graphical Models

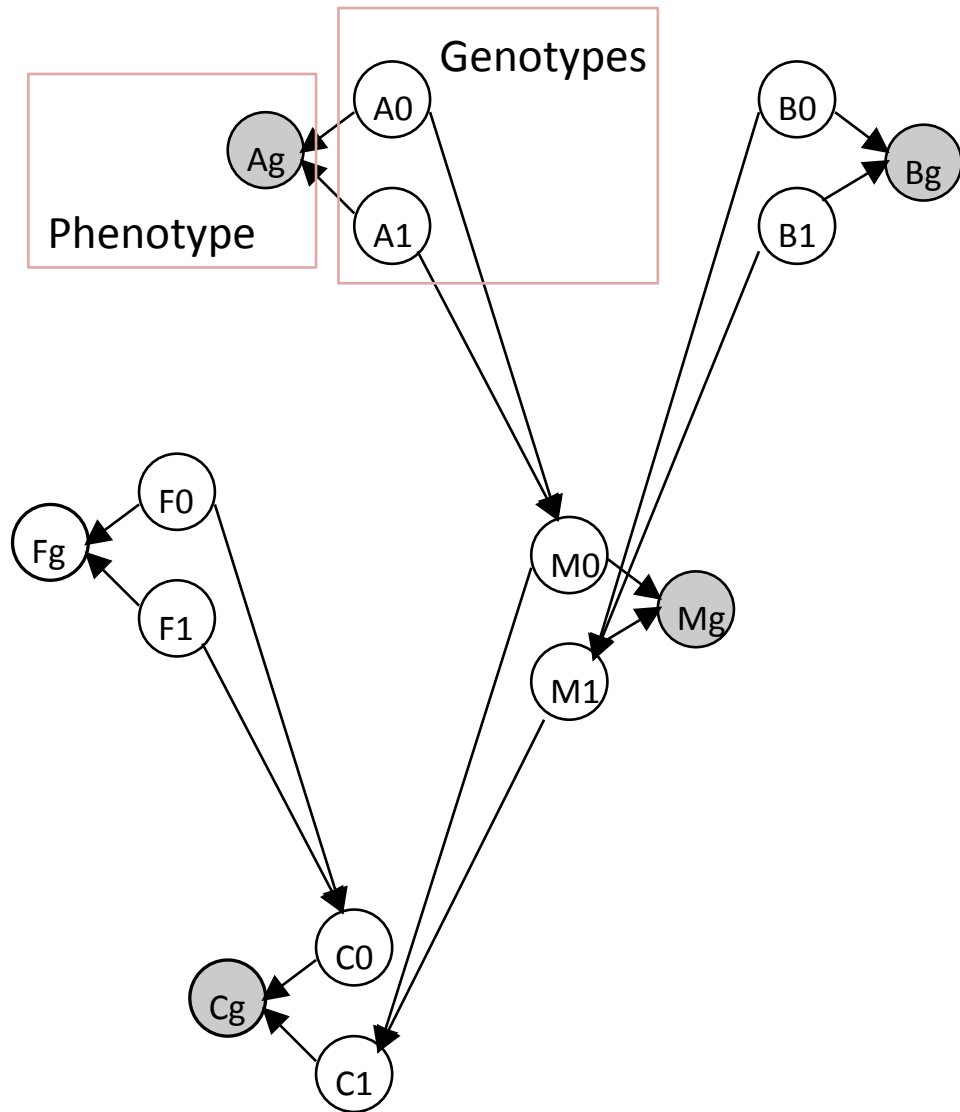
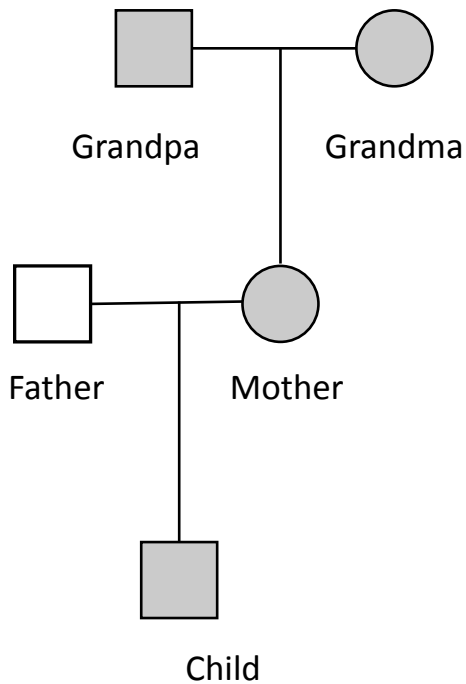


- The joint distribution on (X_1, X_2, \dots, X_N) factors according to the “parent-of” relations defined by the edges E :

$$p(X_1, X_2, X_3, X_4, X_5, X_6) = p(X_1) p(X_2|X_1) p(X_3|X_2) p(X_4|X_1) p(X_5|X_4) p(X_6|X_2, X_5)$$

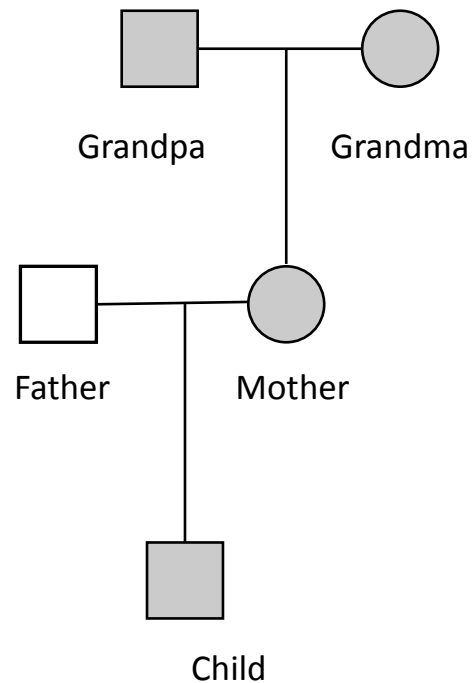
Pedigree as Graphical Models: the Allele Network

Shaded means affected,
blank means unaffected.



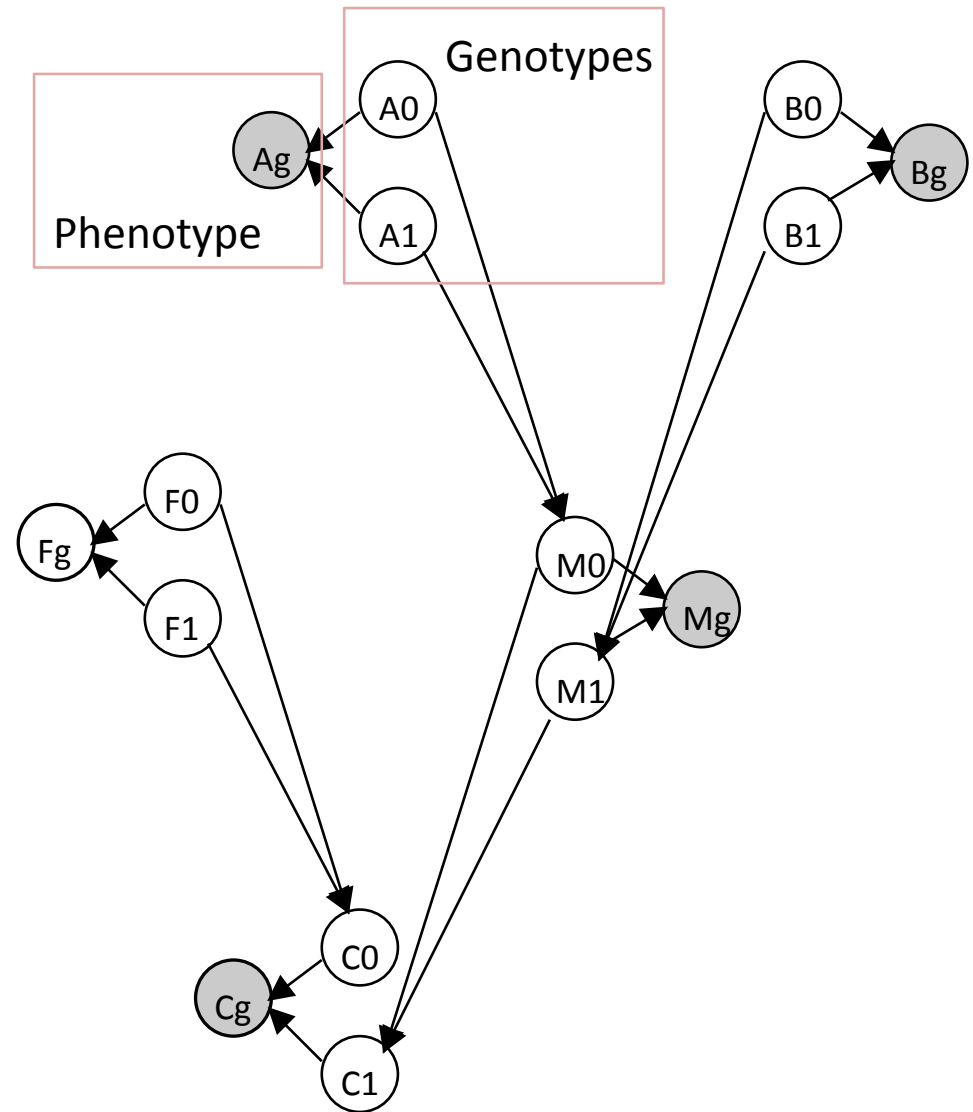
Founders and Non-founders

- **Founders:** individuals whose parents are not in the pedigree.
- **Non-founders:** individuals whose parents are not in the pedigree.



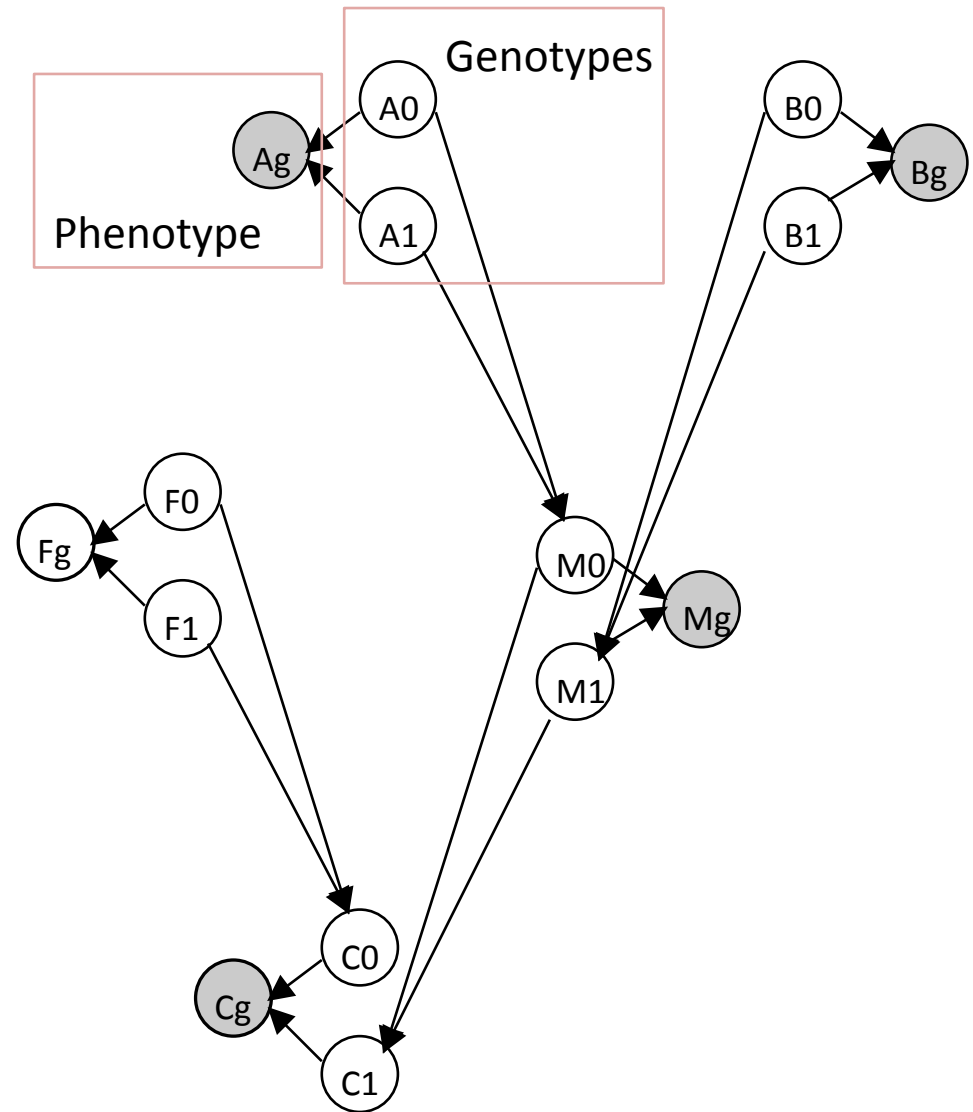
Probability Models over Pedigree

- **Founder genotype probabilities**
- **Transmission probabilities:**
 $P(\text{child's genotype} \mid \text{father's genotype, mother's genotype})$
- **Penetrance model:**
 $P(\text{phenotype} \mid \text{genotype})$ for each individual



Probability Models over Pedigree

- Genotype probabilities are independent
 - across different founders
 - Across siblings of the same parents
- Phenotype probability of each individual is independent of all other individuals genotypes, conditional on their own genotype



One Locus: Founder Genotype Probabilities

- **Assign founder probabilities** to their genotypes, assuming **Hardy-Weinberg equilibrium**
 - Example: If the frequency of D is .01, *HWE* says

$$\begin{array}{c} \boxed{1} \\ Dd \end{array}$$

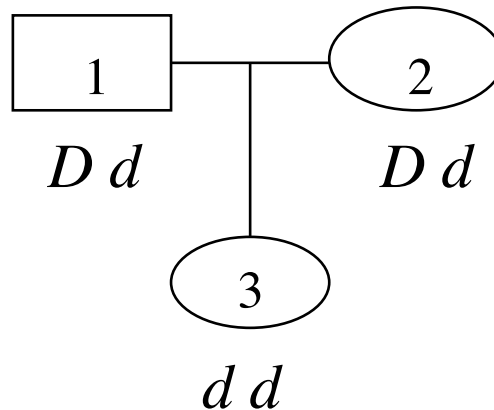
$$P(Dd) = 2 \times .01 \times .99$$

$$\begin{array}{c} \textcircled{2} \\ dd \end{array}$$

$$P(dd) = (.99)^2$$

One Locus: Transmission Probabilities

- Children get their genes from their parents' genes, independently, according to **Mendel's laws**;

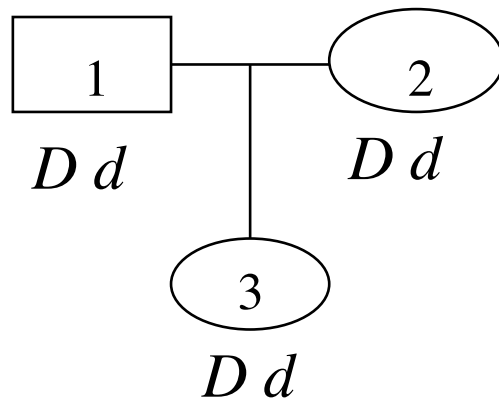


$$\begin{aligned} P(G_{\text{ch3}} = dd \mid G_{\text{pop1}} = Dd, G_{\text{mom2}} = Dd) \\ = 1/2 \times 1/2 \end{aligned}$$

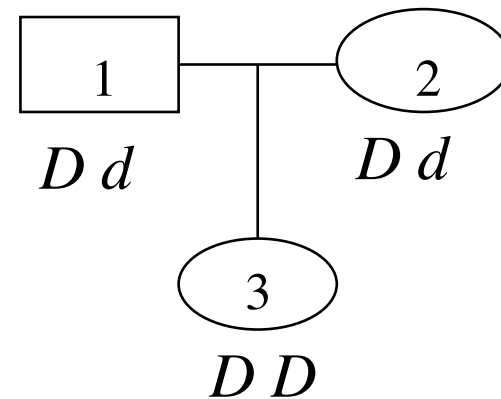
- The inheritances are independent for different children.

One Locus: Transmission Probabilities

- Children get their genes from their parents' genes, independently, according to **Mendel's laws**;



$$P(G_{\text{ch3}} = Dd \mid G_{\text{pop1}} = Dd, G_{\text{mom2}} = Dd) \\ = (1/2 \times 1/2) \times 2$$



$$P(G_{\text{ch3}} = DD \mid G_{\text{pop1}} = Dd, G_{\text{mom2}} = Dd) \\ = 1/2 \times 1/2$$

The factor 2 comes from summing over the two mutually exclusive and equiprobable ways Child3 can get a D and a d .

One Locus: Penetrance Probabilities

- Complete penetrance:

$$P(\text{Ph} = \text{affected} \mid G=DD) = 1$$

DD



- Incomplete penetrance:

$$P(\text{Ph} = \text{affected} \mid G=DD) = .8$$

DD



- Independent Penetrance Model:

- Pedigree analyses usually suppose that, given the genotype at all loci, and in some cases age and sex, the chance of having a particular phenotype depends only on genotype at one locus, and is independent of all other factors: genotypes at other loci, environment, genotypes and phenotypes of relatives, etc.

One Locus: Penetrance Probabilities

- Age and sex-dependent penetrance:



DD (45 years old)

$$P(\text{Ph} = \text{affected} \mid G = DD, \text{sex} = \text{male}, \text{age} = 45 \text{ y.o.}) = .6$$

One Locus: Putting it All Together

- The overall pedigree likelihood is given as

$$L = \prod_{f \in \text{Founders}} P(G_f) \prod_{ch \in \text{Nonfounders}} P(G_{ch} | G_{pop}, G_{mom}) \prod_{i \in \{\text{Founders}, \text{Nonfounders}\}} P(Ph_i | G_i)$$

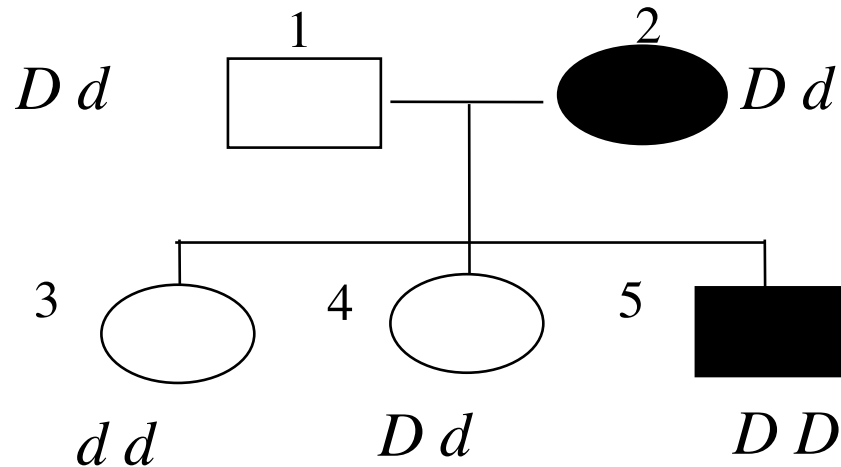
- If founder or non-founder genotypes are **unavailable/missing**, we sum over all possible genotypes for those individuals with missing genotypes to obtain the likelihood

$$L = \sum_{G_f, G_{ch}} \prod_{f \in \text{Founders}} P(G_f) \prod_{ch \in \text{Nonfounders}} P(G_{ch} | G_{pop}, G_{mom}) \prod_{i \in \{\text{Founders}, \text{Nonfounders}\}} P(Ph_i | G_i)$$

One Locus: LOD Score

- Null hypothesis: the disease locus is **unlinked** to the given marker locus being tested
- Alternative hypothesis: the disease locus is **linked** to the given marker locus being tested
- LOD Score = $\text{Log}_{10}(\text{Likelihood under the alternative hypothesis}) - \text{Log}_{10}(\text{Likelihood under the null hypothesis})$
 - Likelihood under the null hypothesis can be obtained by summing the pedigree likelihood over all possible genotypes of the all pedigree individuals: Computationally expensive but efficient algorithms exist

One Locus Example



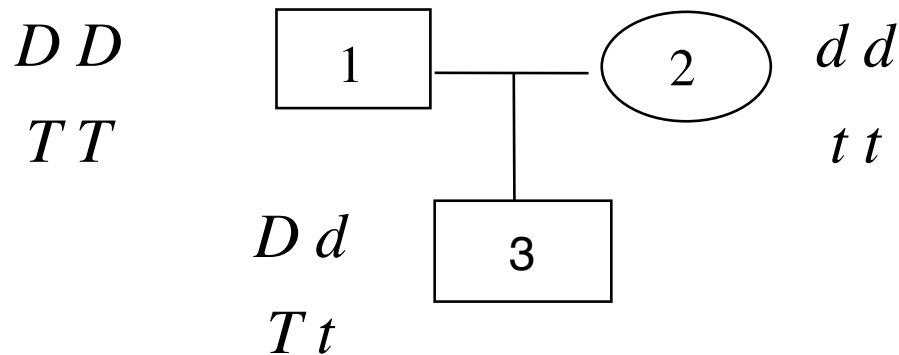
- Assume
 - Penetrances: $P(\text{affected} \mid dd) = .1$, $P(\text{affected} \mid Dd) = .3$, $P(\text{affected} \mid DD) = .8$.
 - Allele D has frequency .01.
- The probability of this pedigree is given as

$$(2 \times .01 \times .99 \times .7) \times (2 \times .01 \times .99 \times .3) \times (1/2 \times 1/2 \times .9) \times (2 \times 1/2 \times 1/2 \times .7) \times (1/2 \times 1/2 \times .8)$$

One Locus Analysis

- Two algorithms:
 - The general strategy of beginning with founders, then non-founders, and multiplying and summing as appropriate, has been codified in what is known as the **Elston-Stewart algorithm** for calculating probabilities over pedigrees. It is one of the two widely used approaches.
 - The other is called the **Lander-Green algorithm** and takes a quite different approach. Lander-Green algorithm uses hidden Markov models to model multiple loci

Two Loci: Linkage and Recombination



- Son 3 produces sperms with $D-T$, $D-t$, $d-T$ or $d-t$ in proportions:

	T	t	
D	$(1-\theta)/2$	$\theta/2$	$1/2$
d	$\theta/2$	$(1-\theta)/2$	$1/2$
	$1/2$	$1/2$	

no recomb.

Two Loci: Linkage and Recombination

- Son produces sperm with DT , Dt , dT or dt in proportions:

	T	t	
D	$(1-\theta)/2$	$\theta/2$	1/2
d	$\theta/2$	$(1-\theta)/2$	1/2
	1/2	1/2	

$\theta = 1/2$: independent assortment (*cf* Mendel) unlinked loci

$\theta < 1/2$: linked loci

$\theta \approx 0$: tightly linked loci

Note: $\theta > 1/2$ is never observed

Two Loci: Linkage and Recombination

- Son produces sperm with DT , Dt , dT or dt in proportions:

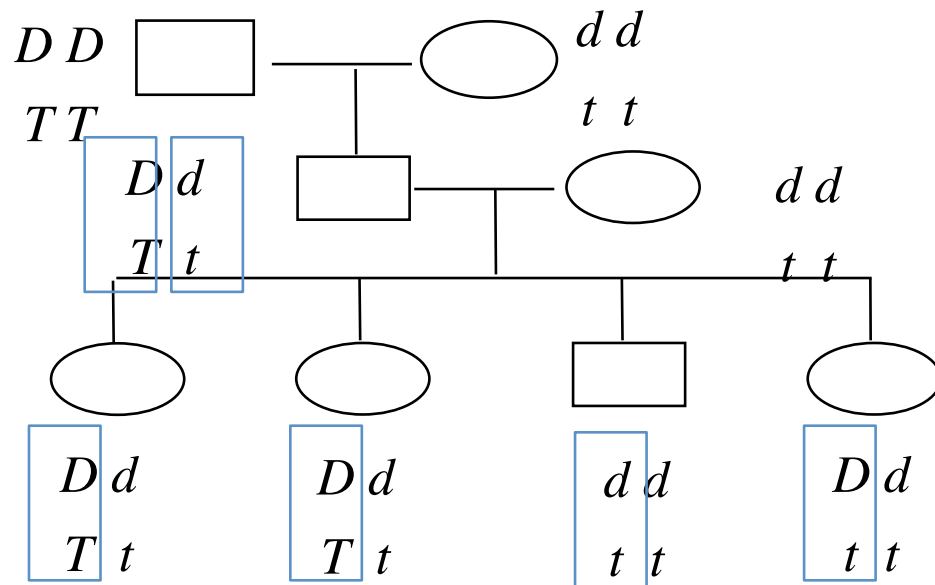
	T	t	
D	$(1-\theta)/2$	$\theta/2$	1/2
d	$\theta/2$	$(1-\theta)/2$	1/2
	1/2	1/2	

- If the loci are linked,
 - $D-T$ and $d-t$ are **parental** haplotypes
 - $D-t$ and $d-T$ are **recombinant** haplotypes

Phase

- Phase is **known** for an individual if you **can** tell whether the gamete was **parental or recombinant**
- Phase is **unknown** if you **cannot** tell whether the gamete was **parental or recombinant**

Two Loci: Phase Known Pedigree



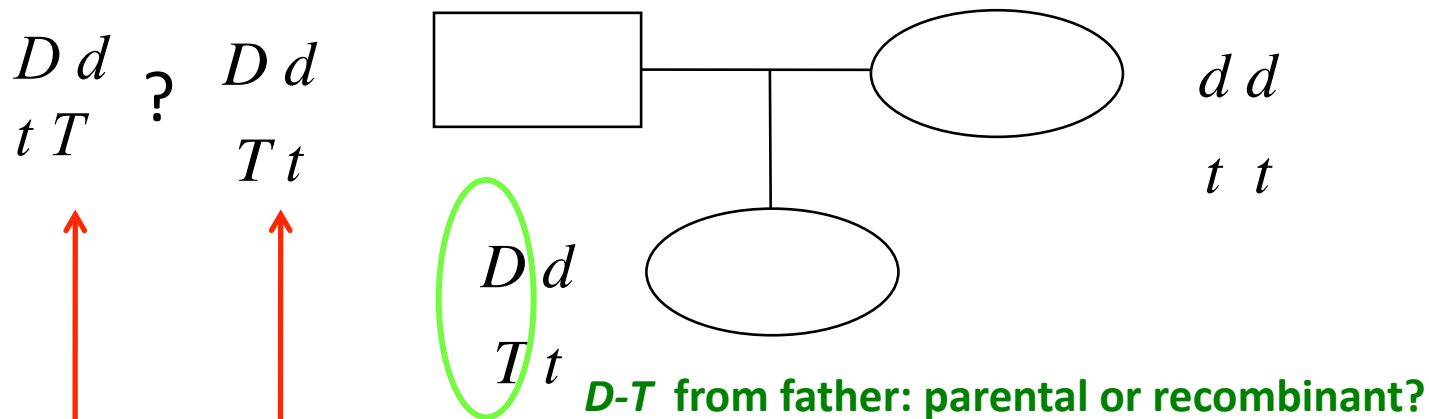
Recombination only discernible in the father. Here $\hat{\theta} = 1/4$ (why?)

This is called the *phase-known double backcross* pedigree.

What if the grandparents' genotypes are not known?

Two Loci: Phase Unknown Pedigree

- Suppose the grandparents' genotypes are unavailable in the double backcross pedigree



If father got $D-T$ from one parent and $d-t$ from the other, the daughter's **paternally derived haplotype** is **parental**.

If father got $D-t$ from one parent and $d-T$ from the other, daughter's **paternally derived haplotype** would be **recombinant**.

Two Loci: Dealing with Phase

- Phase is usually regarded as unknown genetic information
- Sometimes, but not always, phase can be inferred with certainty from genotype data on parents, multiple children, relatives.
- In practice, probabilities must be calculated under all phases compatible with the observed data, and added together: **computationally intensive**, especially with multilocus analyses.

Two Loci: Founder Probabilities

- Assume **linkage equilibrium**, i.e. independence of genotypes across the two loci.

- Allele frequencies at locus one: $D = .01$, and $d = .99$

Allele frequencies at locus two: $T = .25$ and $t = .75$


- Haplotype frequencies

- $DT = .01 \times .25$
- $Dt = .01 \times .75$
- $dT = .99 \times .25$
- $dt = .99 \times .75$

- Together with Hardy-Weinberg, this implies that

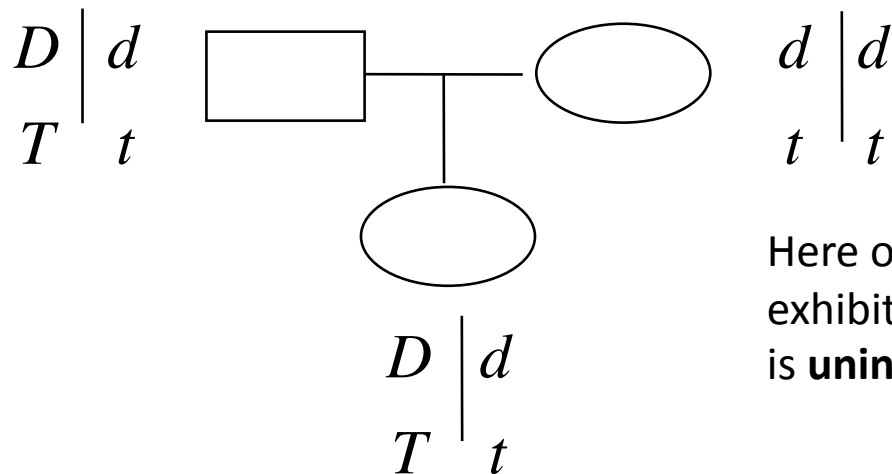
$$\begin{aligned} P(G = DdTt) &= (2 \times .01 \times .99) \times (2 \times .25 \times .75) \\ &= 2 \times (.01 \times .25) \times (.99 \times .75) + 2 \times (.01 \times .75) \times (.99 \times .25). \end{aligned}$$

adds haplotype pair probabilities.

Dd
 Tt 

Two Loci: Transmission Probabilities

- For a given haplotype inheritance:



Here only the father can exhibit recombination: mother is **uninformative**.

$$\begin{aligned}
 &P(G_{\text{ch}} = DT/dt \mid G_{\text{pop}} = DT/dt, G_{\text{mom}} = dt/dt) \\
 &= P(G_{\text{ch}} = DT \mid G_{\text{pop}} = DT/dt) \times P(G_{\text{ch}} = dt \mid G_{\text{mom}} = dt/dt) \\
 &= (1-\theta)/2 \times 1.
 \end{aligned}$$

- Sum the probabilities over all possible phases/haplotypes.

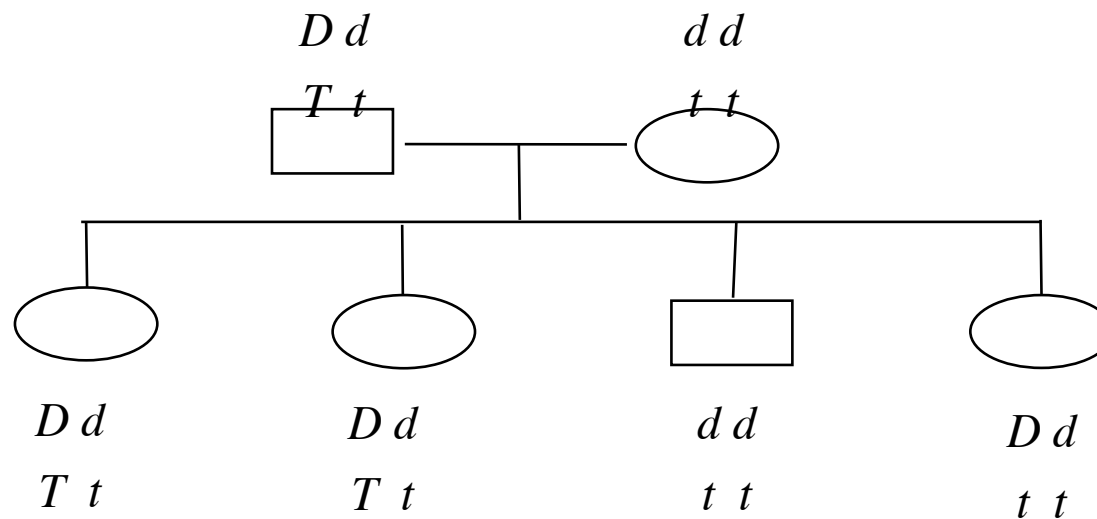
Two Loci: Penetrance

- In all standard linkage programs, different parts of phenotype are conditionally independent given all genotypes, and two-loci penetrances split into products of one-locus penetrances.
- Assuming the penetrances for DD, Dd and dd given earlier, and that T,t are two alleles at a **co-dominant** marker locus.

$$\begin{aligned} & P(Ph_1 = \text{affected}, Ph_2 = Tt \mid G_1 = DD, G_2 = Tt) \\ &= \Pr(Ph_1 = \text{affected} \mid G_1 = DD, G_2 = Tt) \times \Pr(Ph_2 = Tt \mid G_1 = DD, G_2 = Tt) \\ &= 0.8 \times 1 \end{aligned}$$

Two Loci: Phase Unknown Double Backcross

- We assume below pop is as likely to be DT / dt as Dt / dT .



$$\begin{aligned}
 & P(\text{all data} \mid \theta) \\
 &= P(\text{parents' data} \mid \theta) \times P(\text{kids' data} \mid \text{parents' data}, \theta) \\
 &= P(\text{parents' data}) \times \left\{ \left[\frac{((1-\theta)/2)^3 \times \theta/2}{2} + \frac{(\theta/2)^3 \times (1-\theta)/2}{2} \right] \right\}
 \end{aligned}$$

This is then maximised in θ , in this case numerically. Here $\hat{\theta} = 0.25$

Log (base 10) Odds or LOD Scores

- Suppose $P(\text{data} \mid \theta)$ is the likelihood function of a recombination fraction θ generated by some 'data', and $P(\text{data} \mid 1/2)$ is the same likelihood when $\theta = 1/2$.
- This can equally well be done with $\text{Log}_{10}L$, i.e.
$$\text{LOD}(\theta^*) = \text{Log}_{10}P(\text{data} \mid \theta^*) - \text{Log}_{10}P(\text{data} \mid 1/2)$$
measures the relative strength of the data for $\theta = \theta^*$ (optimal θ) rather than $\theta = 1/2$.

Facts about/interpretation of LOD scores

1. Positive LOD scores suggests stronger support for θ^* than for $1/2$, negative LOD scores the reverse.
2. Higher LOD scores means stronger support, lower means the reverse.
3. LODs are additive across independent pedigrees, and under certain circumstances can be calculated sequentially.
4. For a single two-point linkage analysis, the threshold $\text{LOD} \approx 3$ has become the de facto standard for "establishing linkage", i.e. rejecting the null hypothesis of no linkage.
5. When more than one locus or model is examined, the remark in 4 must be modified, sometimes dramatically.

Assumptions underpinning most 2-point human linkage analyses

- ***Founder Frequencies***: Hardy-Weinberg, random mating at each locus. Linkage equilibrium across loci, **known** allele frequencies; founders independent.
- ***Transmission***: Mendelian segregation, no mutation.
- ***Penetrance***: single locus, no room for dependence on relatives' phenotypes or environment. **Known** (including phenocopy rate).
- ***Implicit***: phenotype and genotype data **correct**, marker order and location correct
- ***Comment***: Some analyses are *robust*, others can be *very sensitive* to violations of some of these assumptions.