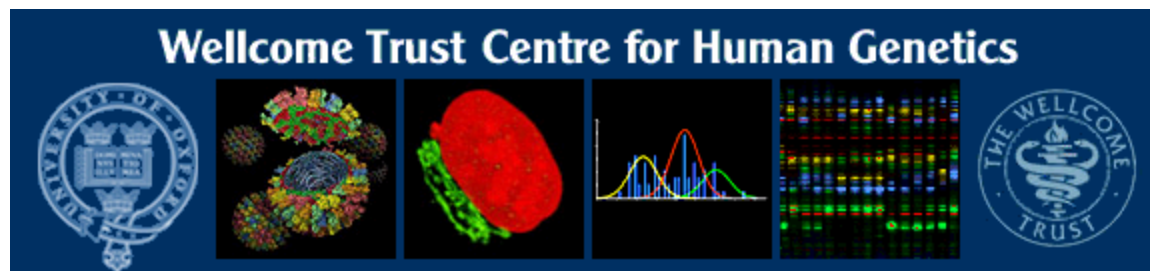




Calculation of IBD probabilities

David Evans and Stacey Cherny

University of Oxford
Wellcome Trust Centre for Human
Genetics





This Session ...

- IBD vs IBS
- Why is IBD important?
- Calculating IBD probabilities
 - Lander-Green Algorithm (MERLIN)
 - Single locus probabilities
 - Hidden Markov Model
 - Other ways of calculating IBD status
 - Elston-Stewart Algorithm
 - MCMC approaches
- MERLIN
- Practical Example
 - IBD determination
 - Information content mapping
 - SNPs vs micro-satellite markers?



Aim of Gene Mapping Experiments

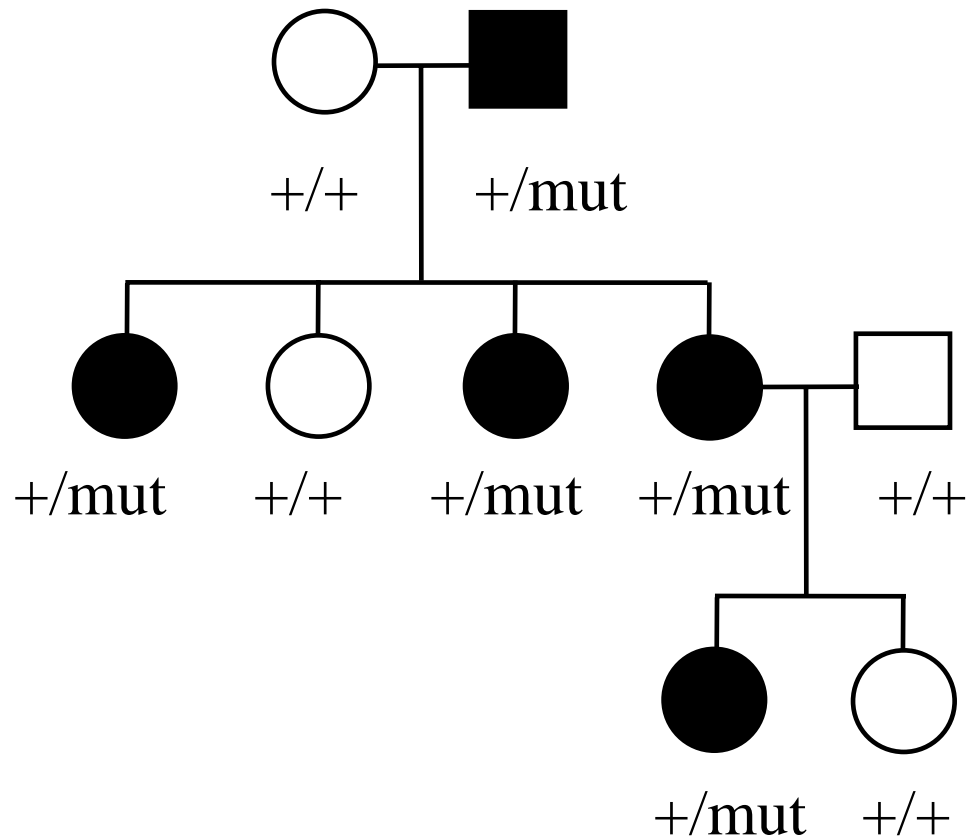
- Identify variants that control interesting traits
 - Susceptibility to human disease
 - Phenotypic variation in the population
- The hypothesis
 - Individuals sharing these variants will be more similar for traits they control
- The difficulty...
 - Testing ~ 10 million variants is impractical...



Identity-by-Descent (IBD)

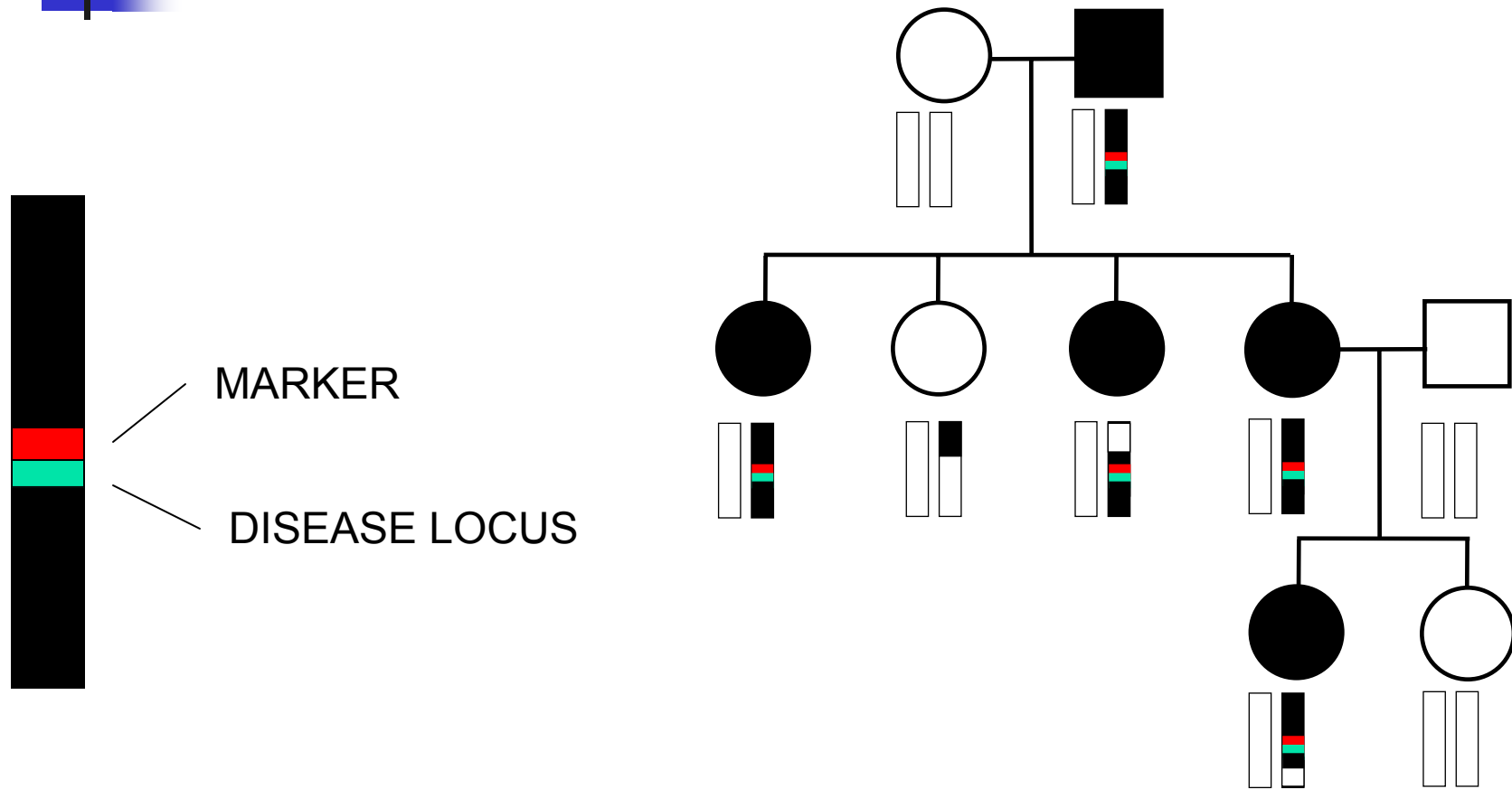
- Two alleles are IBD if they are descended from the same ancestral allele
- If a stretch of chromosome is IBD among a set of individuals, ALL variants within that stretch will also be shared IBD (markers, QTLs, disease genes)
- Allows surveys of large amounts of variation even when a few polymorphisms measured

A Segregating Disease Allele



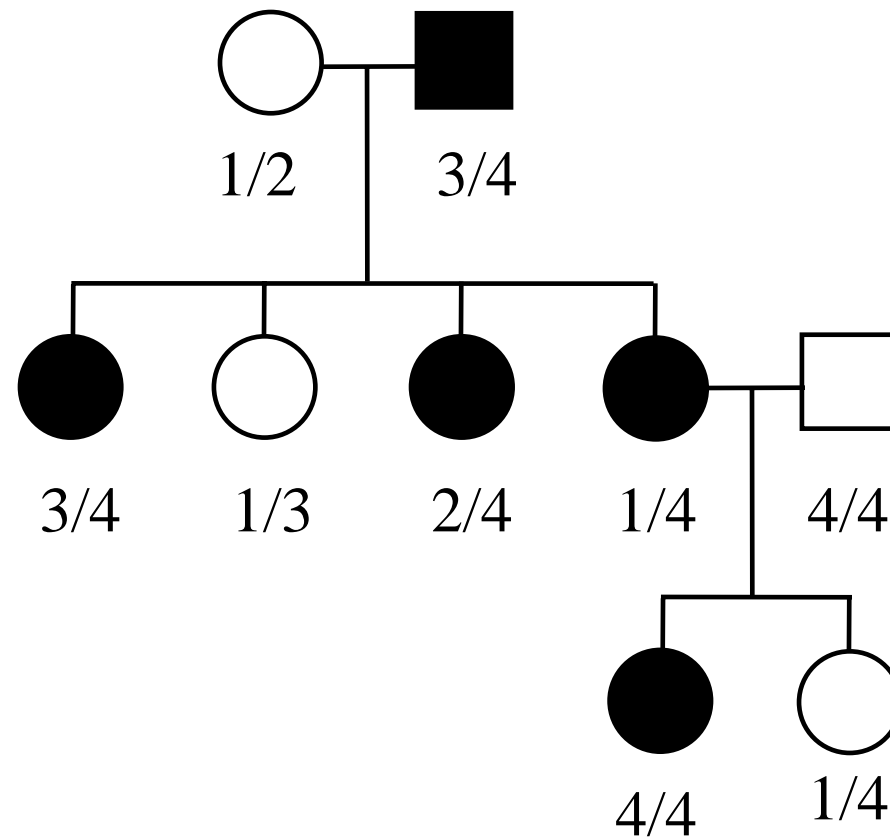
All affected individuals IBD for disease causing mutation

Segregating Chromosomes



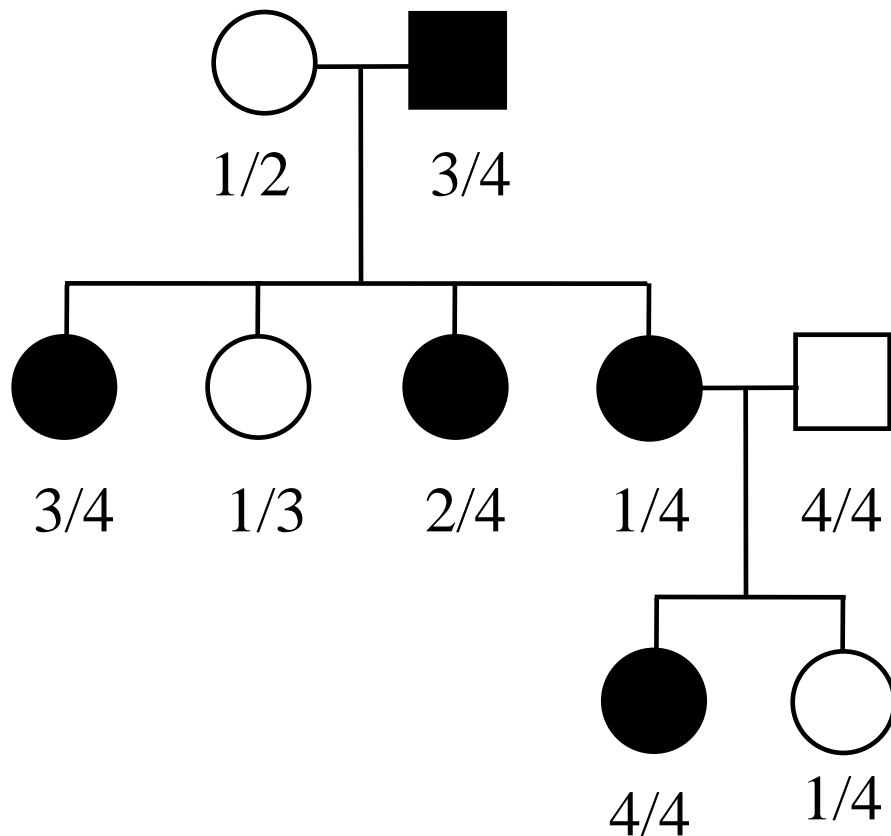
Affected individuals tend to share adjacent areas of chromosome IBD

Marker Shared Among Affecteds



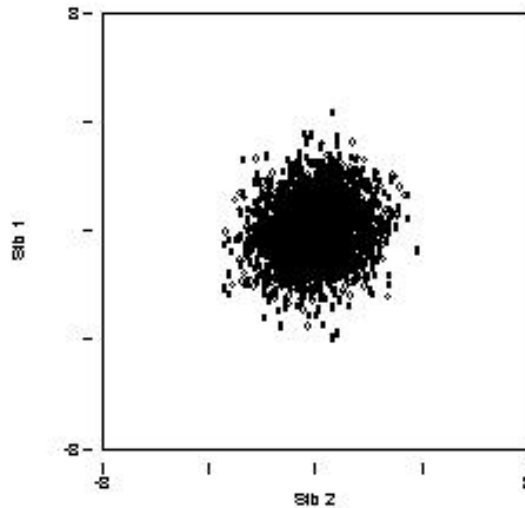
“4” allele segregates with disease

Why is IBD sharing important?



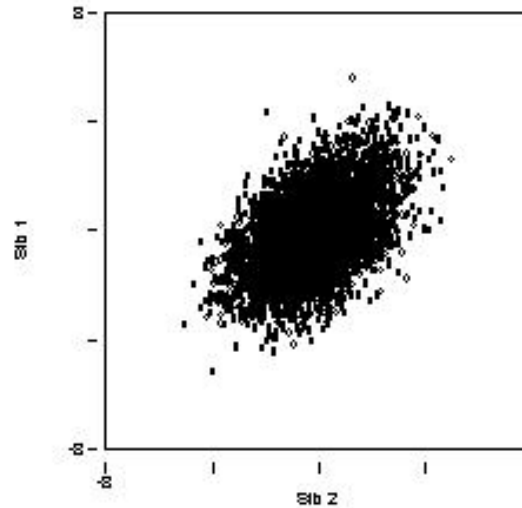
- IBD sharing forms the basis of non-parametric linkage statistics
- Affected relatives tend to share marker alleles close to the disease locus IBD more often than chance

Linkage between QTL and marker

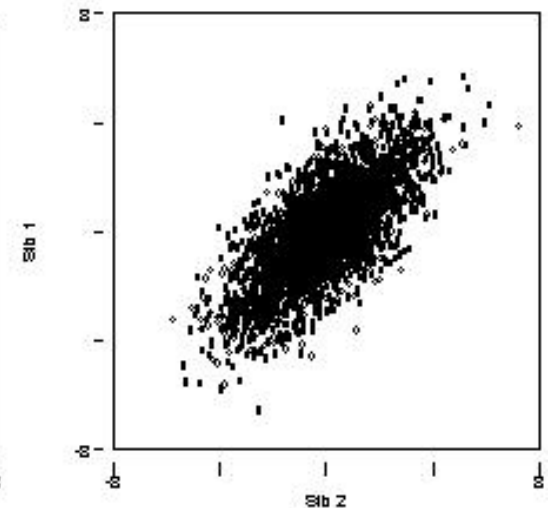


QTL

IBD 0



IBD 1



IBD 2

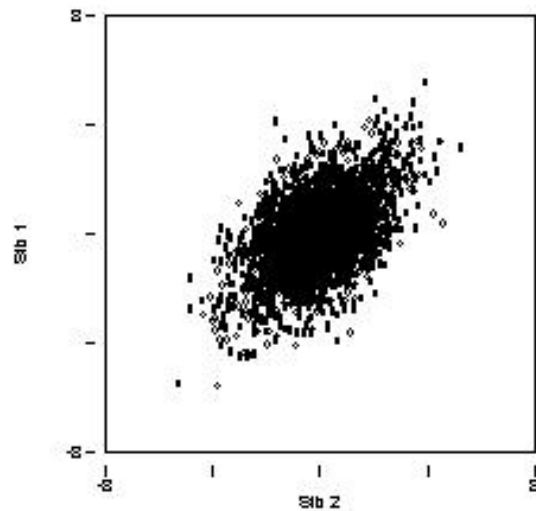
Marker

IBD 0

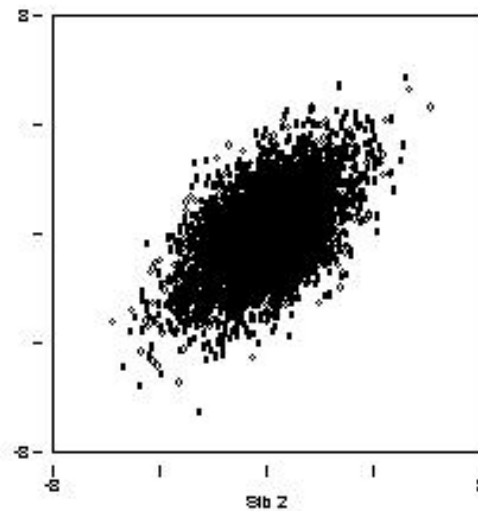
IBD 1

IBD 2

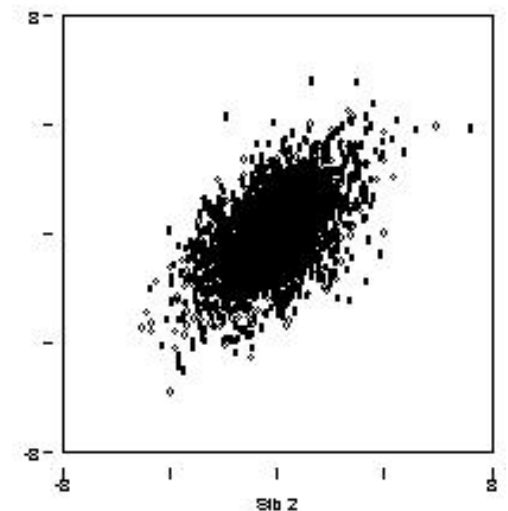
NO Linkage between QTL and marker



IBD 0



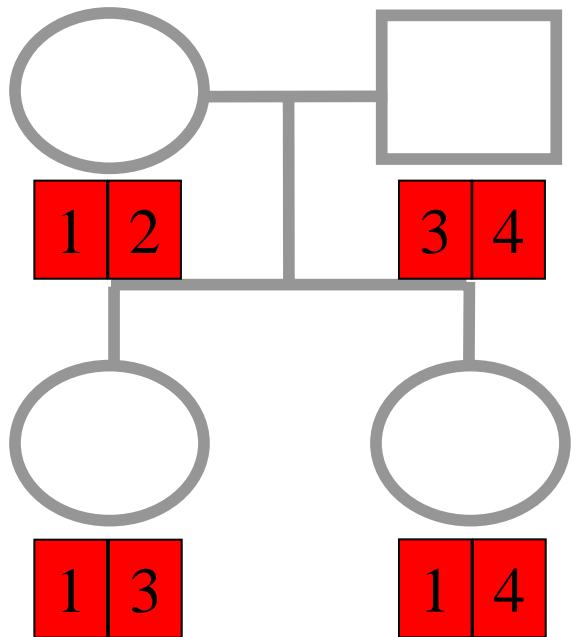
IBD 1



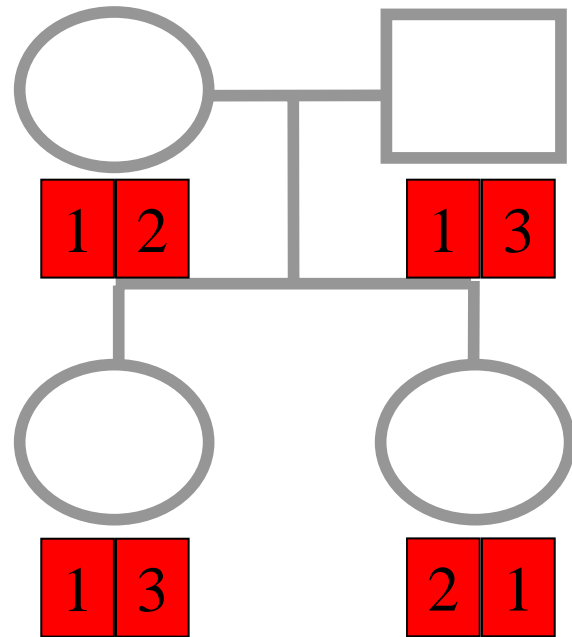
IBD 2

Marker

IBD vs IBS



Identical by Descent
and
Identical by State



Identical by state only



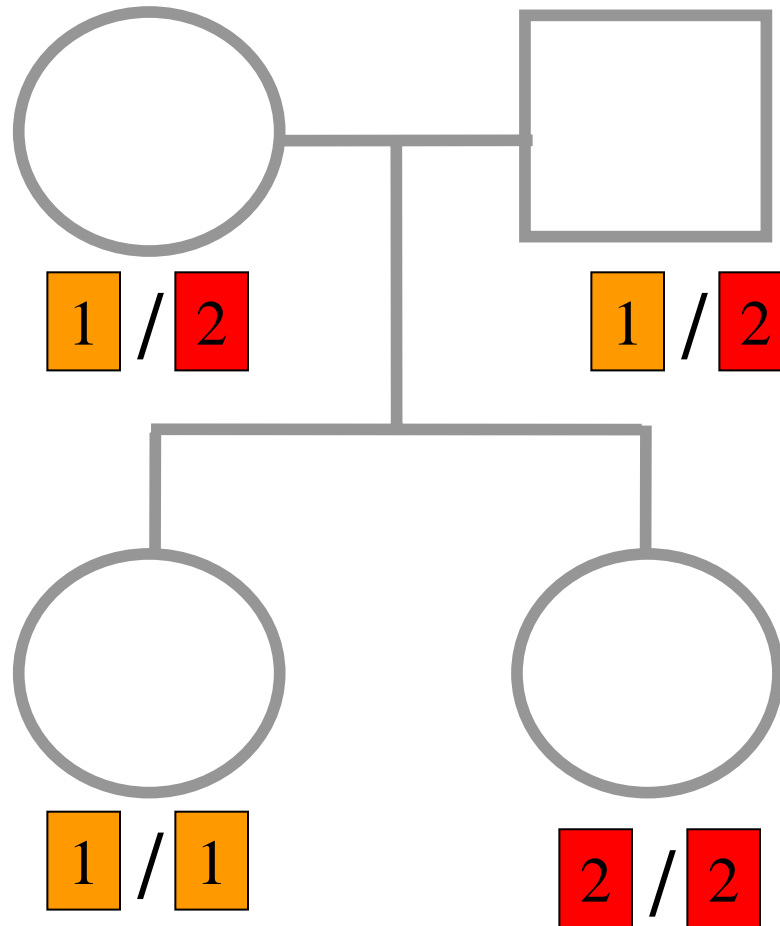
Example: IBD in Siblings

Consider a mating between mother AB x father CD:

		Sib1			
		AC	AD	BC	BD
Sib 2	AC	2	1	1	0
	AD	1	2	0	1
	BC	1	0	2	1
	BD	0	1	1	2

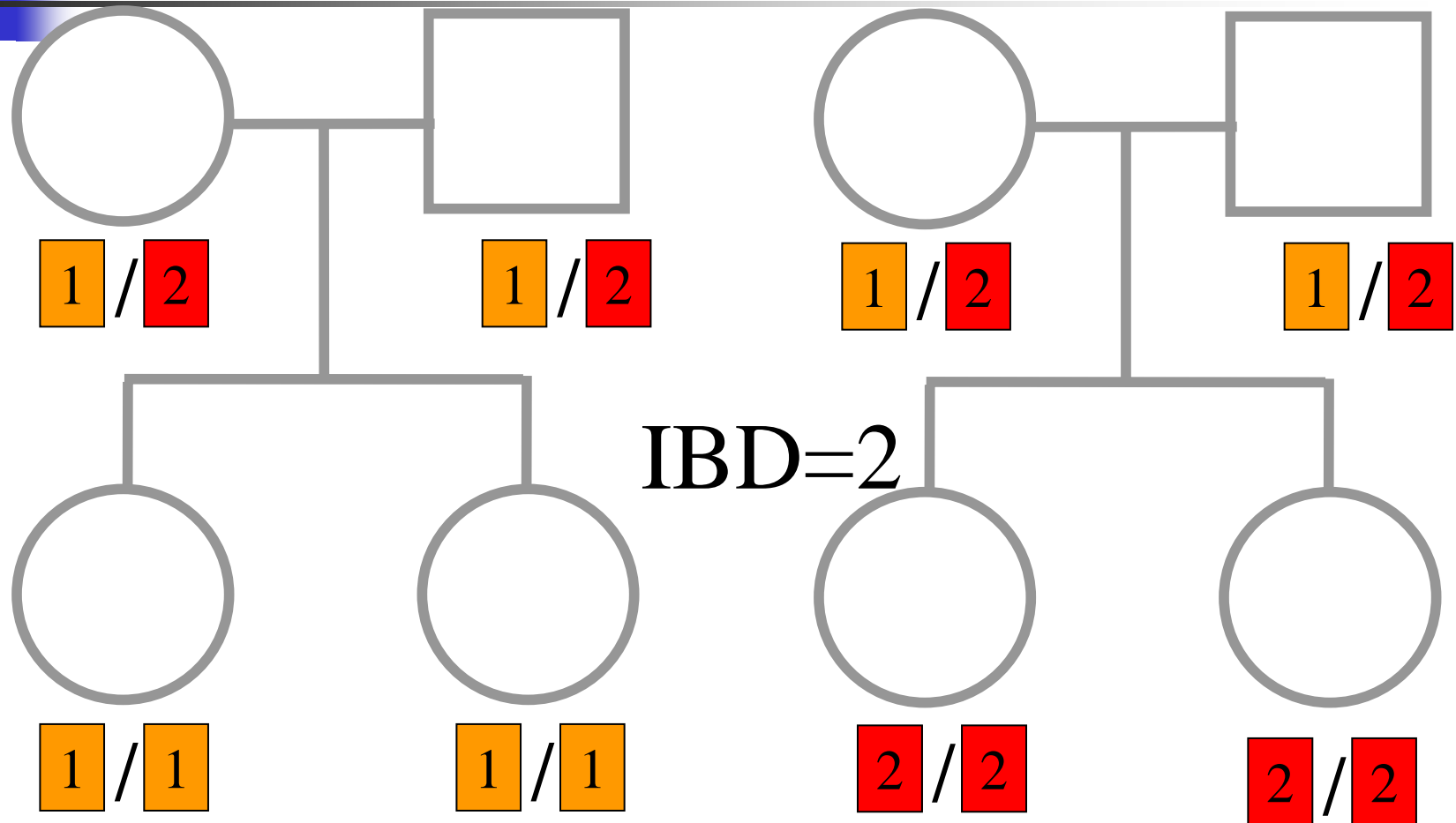
IBD 0 : 1 : 2 = 25% : 50% : 25%

IBD can be trivial...

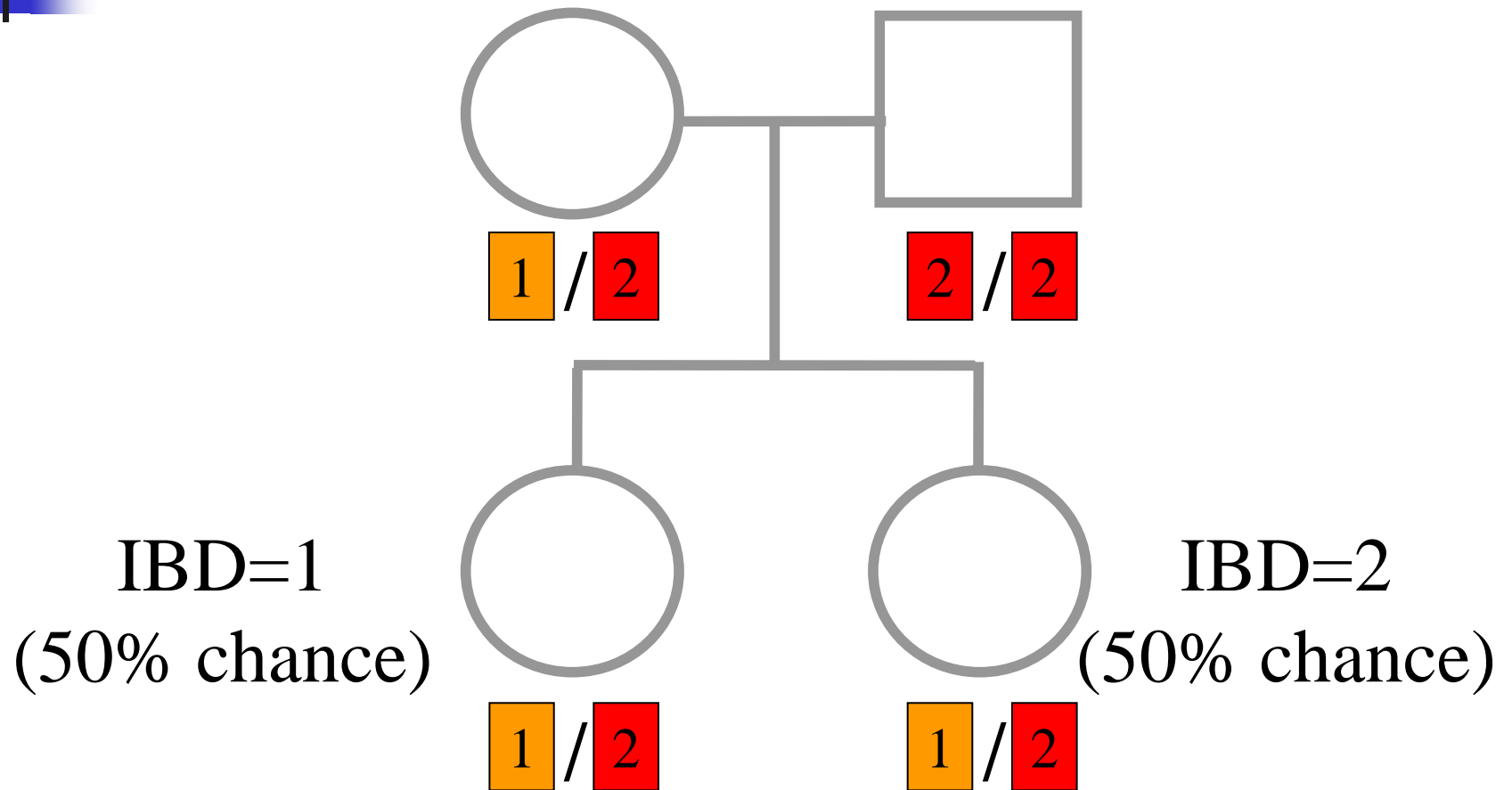


IBD=0

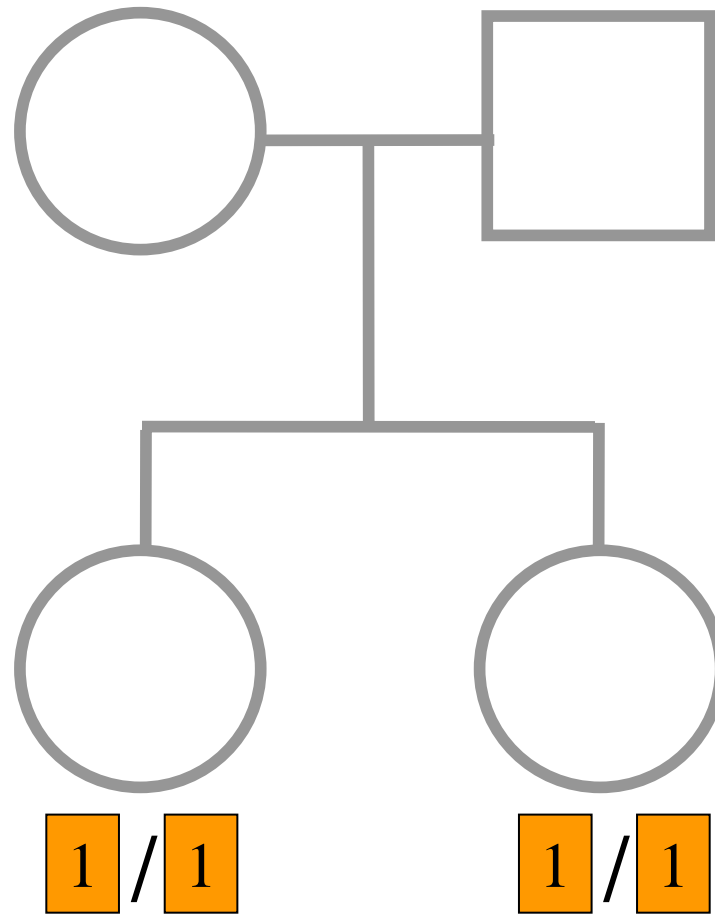
Two Other Simple Cases...



A little more complicated...



And even more complicated...



IBD=?



Bayes Theorem

$$\begin{aligned} P(A_i|B) &= \frac{P(A_i, B)}{P(B)} \\ &= \frac{P(A_i)P(B | A_i)}{P(B)} \\ &= \frac{P(A_i) P(B | A_i)}{\sum_j P(A_j) P(B | A_j)} \end{aligned}$$

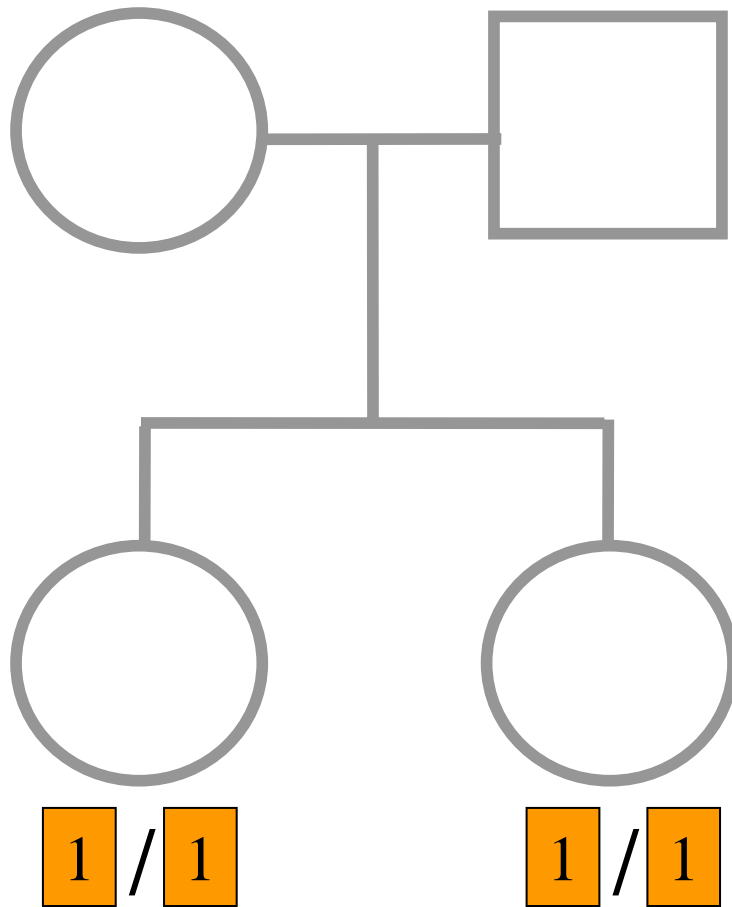
Bayes Theorem for IBD Probabilities

$$\begin{aligned} P(IBD = i | G) &= \frac{P(IBD = i, G)}{P(G)} \\ &= \frac{P(IBD = i)P(G | IBD = i)}{P(G)} \\ &= \frac{P(IBD = i)P(G | IBD = i)}{\sum_j P(IBD = j)P(G | IBD = j)} \end{aligned}$$

P(Marker Genotype|IBD State)

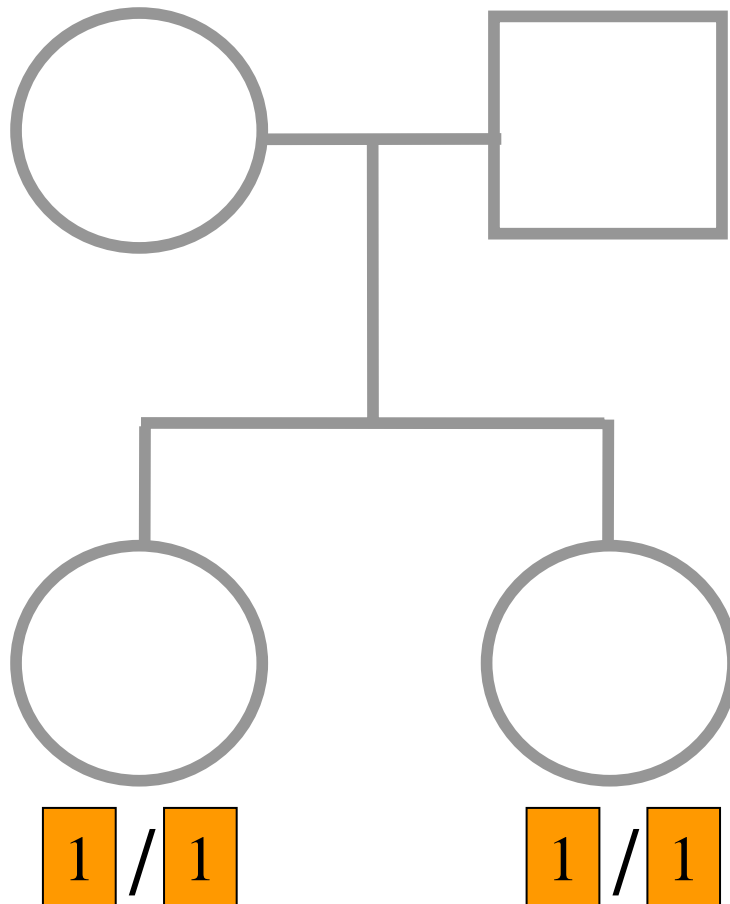
Sib 1	Sib 2	P(observing genotypes / k alleles IBD)		
		$k=0$	$k=1$	$k=2$
A_1A_1	A_1A_1	p_1^4	p_1^3	p_1^2
A_1A_1	A_1A_2	$2p_1^3p_2$	$p_1^2p_2$	0
A_1A_1	A_2A_2	$p_1^2p_2^2$	0	0
A_1A_2	A_1A_1	$2p_1^3p_2$	$p_1^2p_2$	0
A_1A_2	A_1A_2	$4p_1^2p_2^2$	p_1p_2	$2p_1p_2$
A_1A_2	A_2A_2	$2p_1p_2^3$	$p_1p_2^2$	0
A_2A_2	A_1A_1	$p_1^2p_2^2$	0	0
A_2A_2	A_1A_2	$2p_1p_2^3$	$p_1p_2^2$	0
A_2A_2	A_2A_2	p_2^4	p_2^3	p_2^2

Worked Example



$$p_1 = 0.5$$

Worked Example



$$p_1 = 0.5$$

$$P(G | IBD=0) = p_1^4 = \frac{1}{16}$$

$$P(G | IBD=1) = p_1^3 = \frac{1}{8}$$

$$P(G | IBD=2) = p_1^2 = \frac{1}{4}$$

$$P(G) = \frac{1}{4}p_1^4 + \frac{1}{2}p_1^3 + \frac{1}{4}p_1^2 = \frac{9}{64}$$

$$P(IBD=0 | G) = \frac{\frac{1}{4}p_1^4}{P(G)} = \frac{1}{9}$$

$$P(IBD=1 | G) = \frac{\frac{1}{2}p_1^3}{P(G)} = \frac{4}{9}$$

$$P(IBD=2 | G) = \frac{\frac{1}{4}p_1^2}{P(G)} = \frac{4}{9}$$

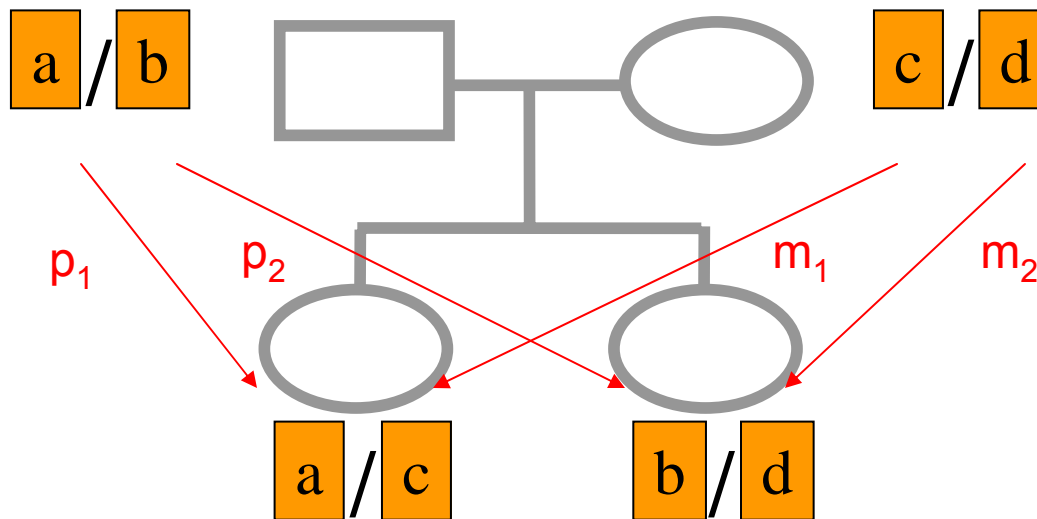
For ANY PEDIGREE the inheritance pattern at every point in the genome can be completely described by a binary inheritance vector:

$$v(x) = (p_1, m_1, p_2, m_2, \dots, p_n, m_n)$$

whose coordinates describe the outcome of the $2n$ paternal and maternal meioses giving rise to the n non-founders in the pedigree

p_i (m_i) is 0 if the grandpaternal allele transmitted

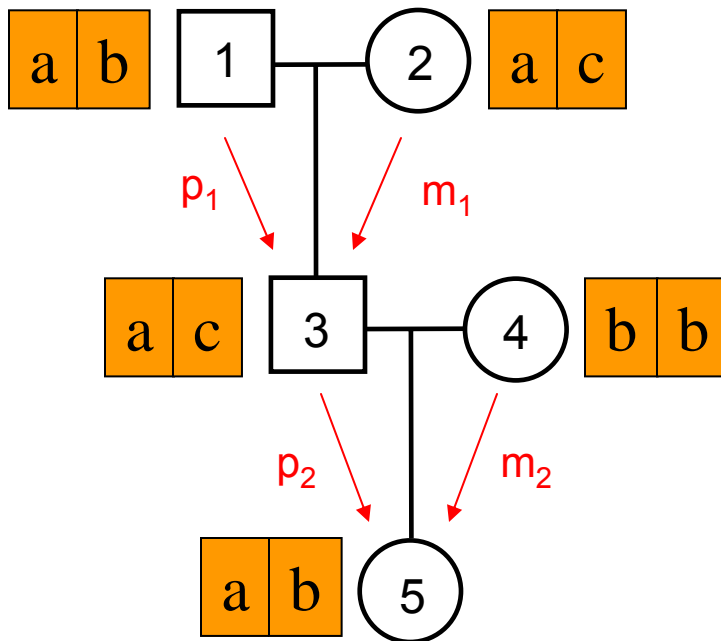
p_i (m_i) is 1 if the grandmaternal allele is transmitted



$$v(x) = [0,0,1,1]$$

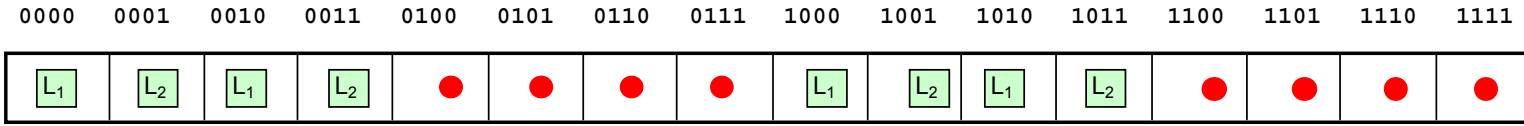
Inheritance Vector

In practice, it is not possible to determine the true inheritance vector at every point in the genome, rather we represent partial information as a probability distribution over the 2^{2n} possible inheritance vectors

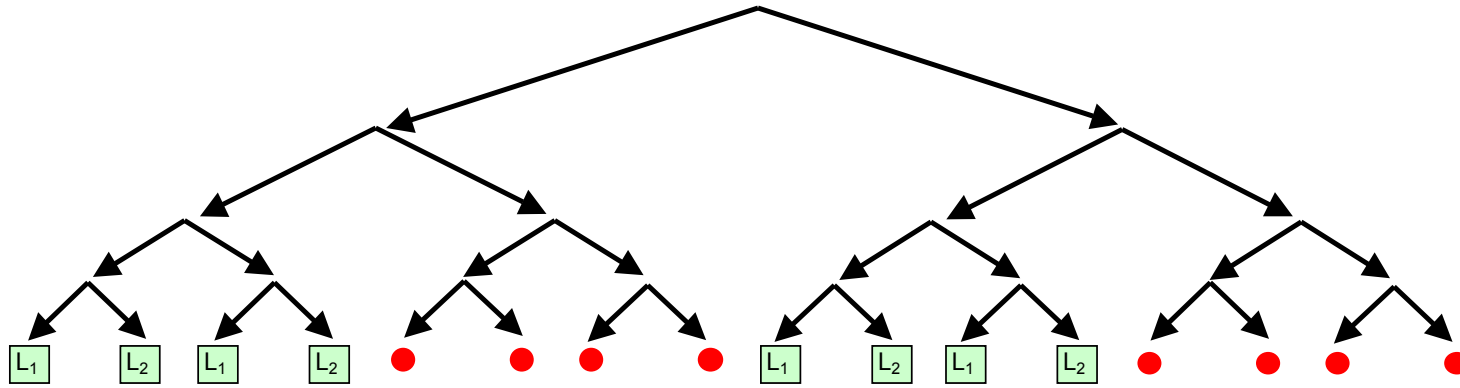


Inheritance vector	Prior	Posterior
0000	1/16	1/8
0001	1/16	1/8
0010	1/16	0
0011	1/16	0
0100	1/16	1/8
0101	1/16	1/8
0110	1/16	0
0111	1/16	0
1000	1/16	1/8
1001	1/16	1/8
1010	1/16	0
1011	1/16	0
1100	1/16	1/8
1101	1/16	1/8
1110	1/16	0
1111	1/16	0

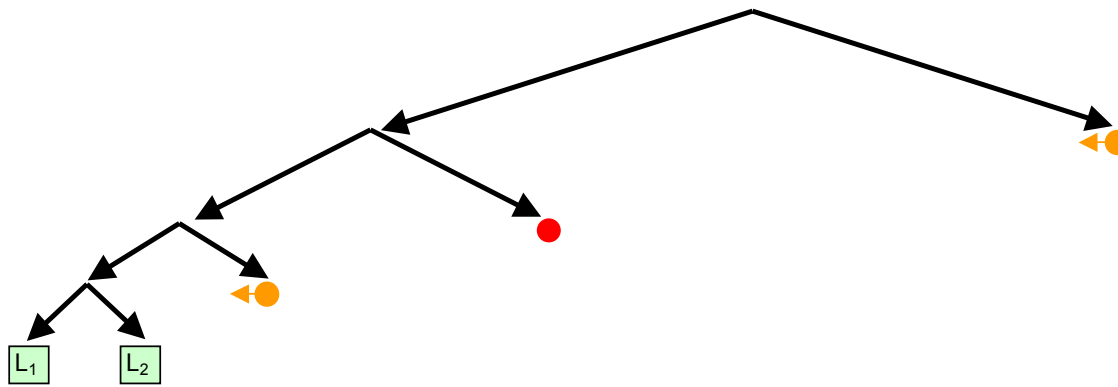
a) bit-indexed array



b) packed tree



c) sparse tree



Legend

- Node with zero likelihood
- ◀● Node identical to sibling
- L₁ L₂ Likelihood for this branch

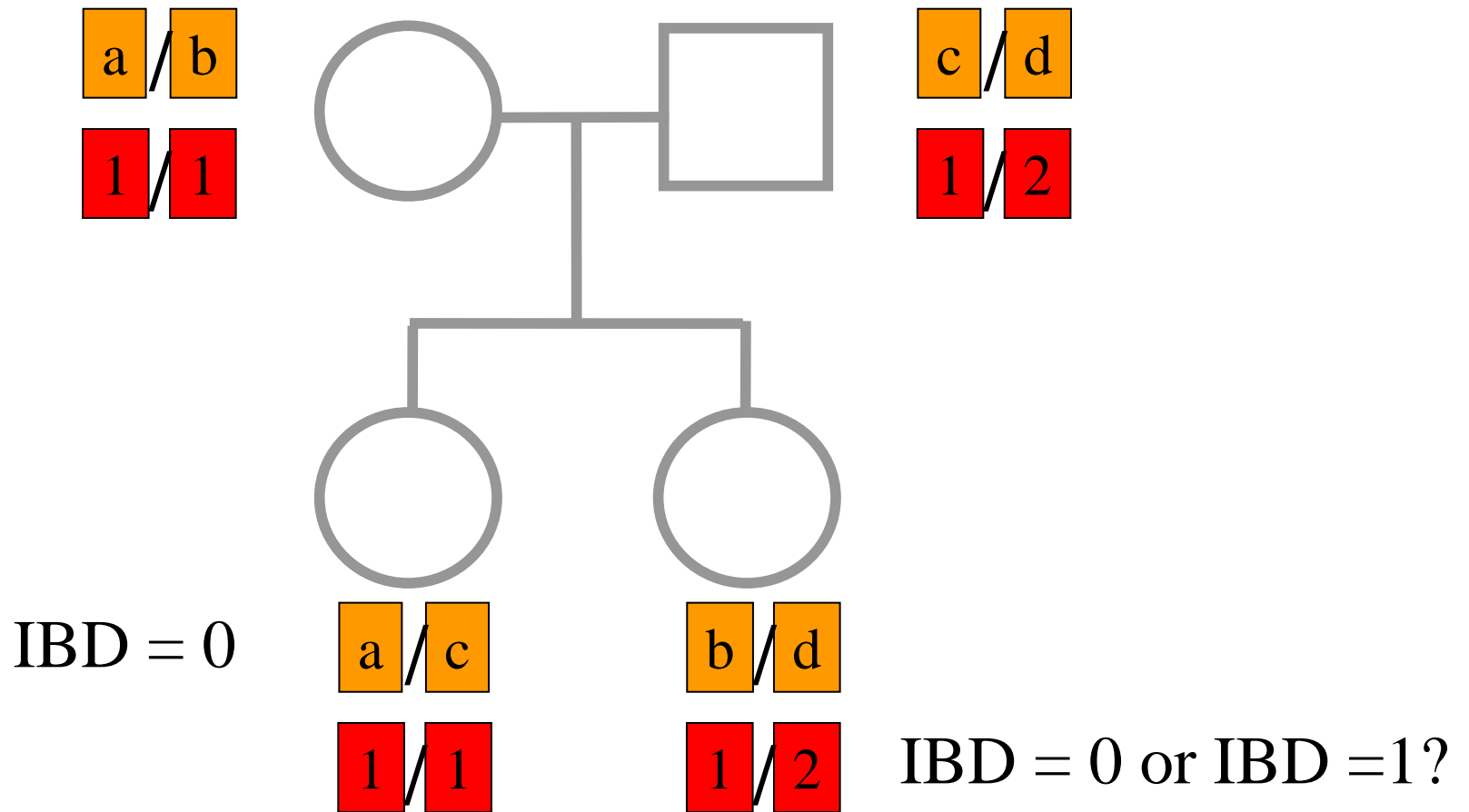
Abecasis et al (2002) *Nat Genet* 30:97-101



Multipoint IBD

- IBD status may not be able to be ascertained with certainty because e.g. the mating is not informative, parental information is not available
- IBD information at uninformative loci can be made more precise by examining nearby linked loci

Multipoint IBD



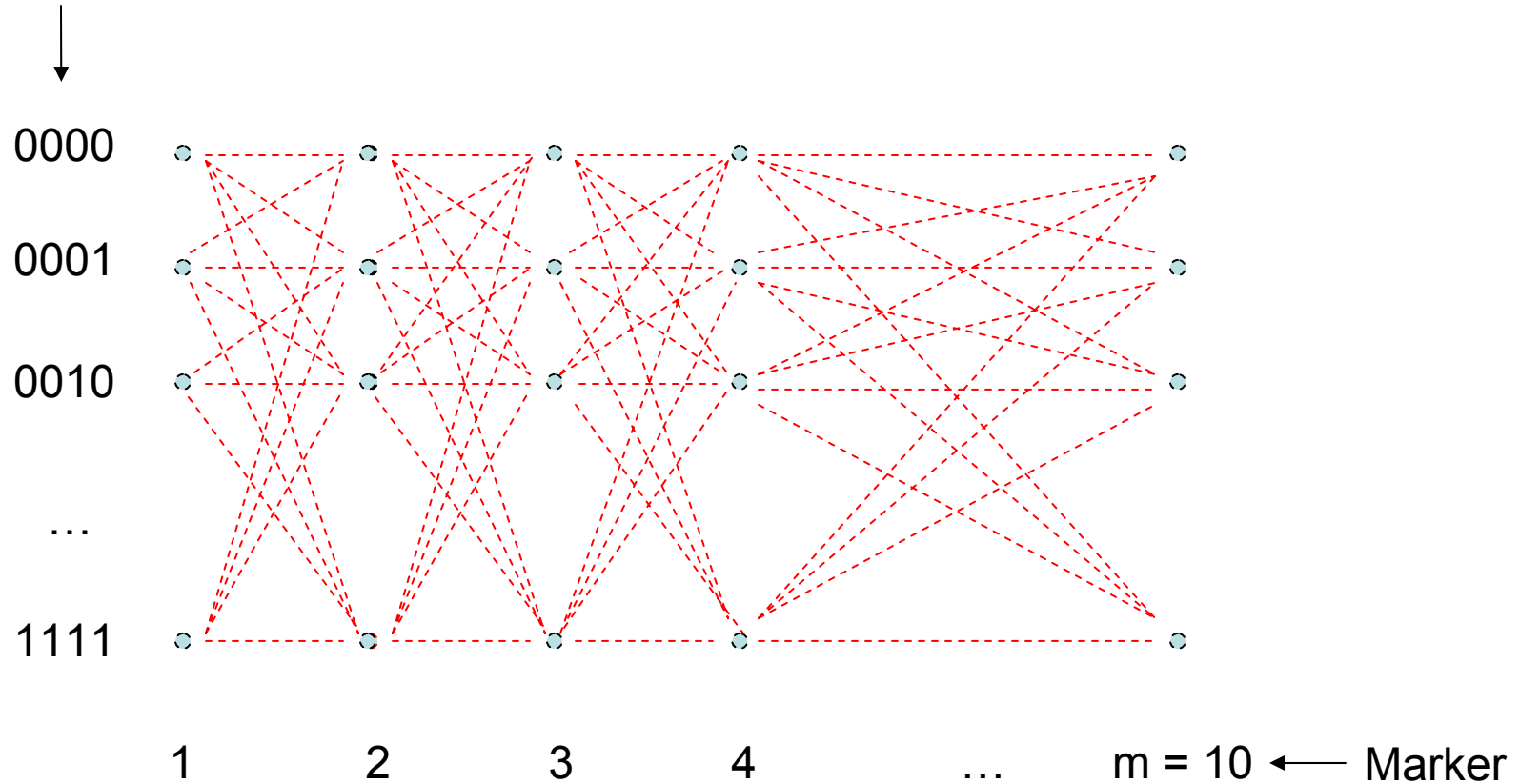


Complexity of the Problem in Larger Pedigrees

- $2n$ meioses in pedigree with n non-founders
 - Each meiosis has 2 possible outcomes
 - Therefore 2^{2n} possibilities for each locus
- For each genetic locus
 - One location for each of m genetic markers
 - Distinct, non-independent meiotic outcomes
- Up to 4^{nm} distinct outcomes!!!

Example: Sib-pair Genotyped at 10 Markers

Inheritance vector



$$(2^{2 \times n})^m = (2^2 \times 2)^{10} = 10^{12} \text{ possible paths !!!}$$



Lander-Green Algorithm

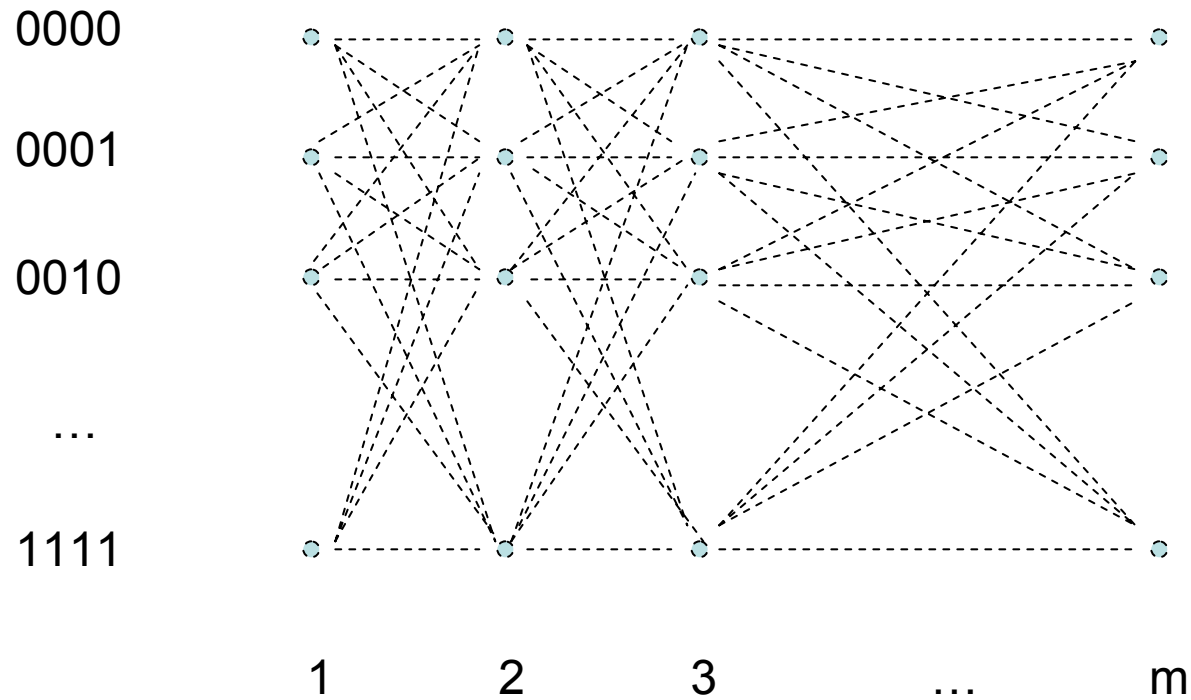
- The inheritance vector at a locus is conditionally independent of the inheritance vectors at all preceding loci given the inheritance vector at the immediately preceding locus (“Hidden Markov chain”)
- The conditional probability of an inheritance vector v_{i+1} at locus $i+1$, given the inheritance vector v_i at locus i is $\theta^j(1-\theta)^{2n-j}$ where θ is the recombination fraction and j is the number of changes in elements of the inheritance vector (“transition probabilities”)

Example:

Locus 1
[0000]

Locus 2
[0001]

Conditional probability = $(1 - \theta)^3\theta$



$$\text{Total Likelihood} = \mathbf{1}' \mathbf{Q}_1 \mathbf{T}_1 \mathbf{Q}_2 \mathbf{T}_2 \dots \mathbf{T}_{m-1} \mathbf{Q}_m \mathbf{1}$$

$$\mathbf{Q}_i = \begin{matrix} & \mathbf{P}[0000] & 0 & 0 & 0 \\ & 0 & \mathbf{P}[0001] & 0 & 0 \\ & 0 & 0 & \dots & 0 \\ & 0 & 0 & 0 & \mathbf{P}[1111] \end{matrix}$$

$2^{2n} \times 2^{2n}$ diagonal matrix of single locus probabilities at locus i

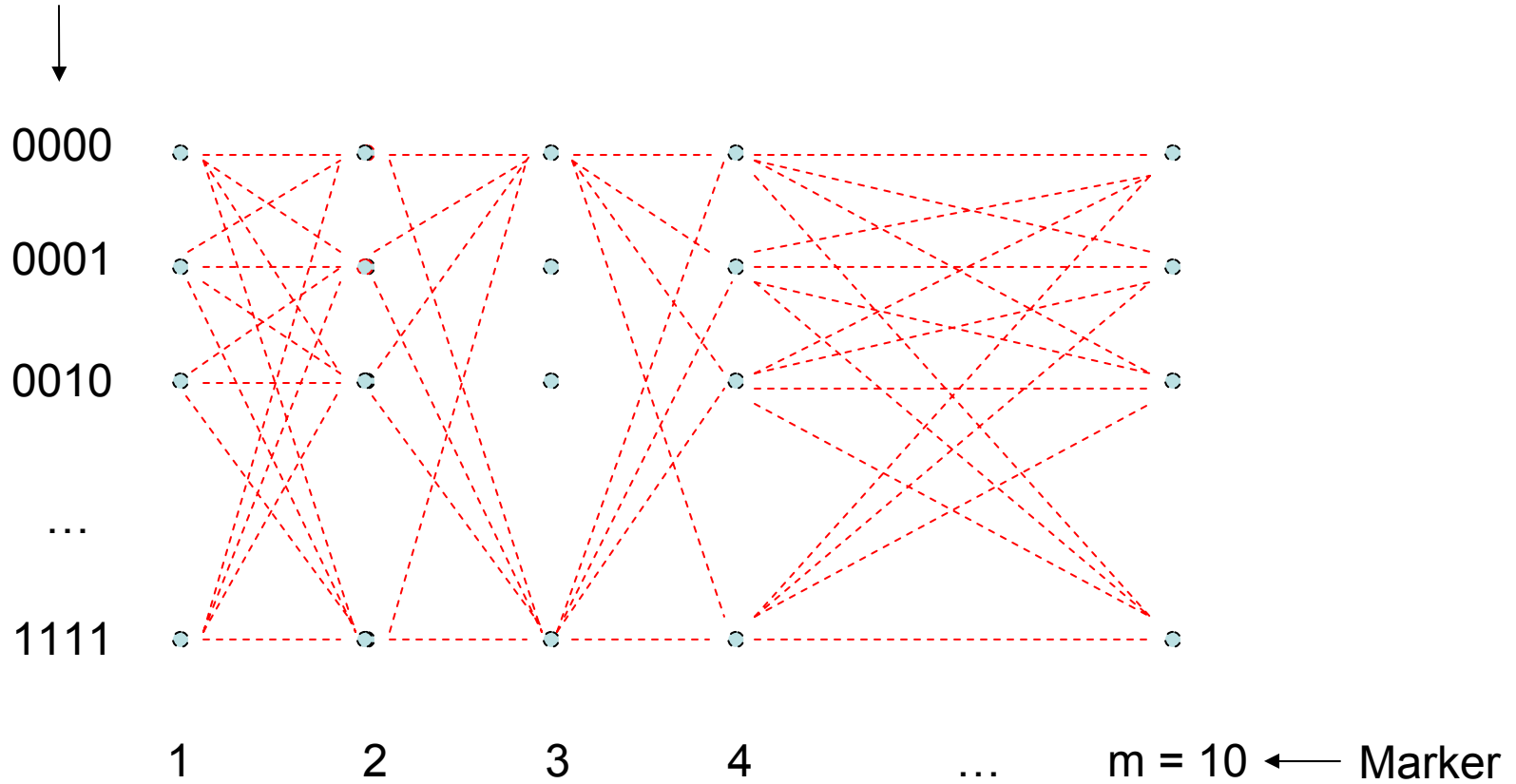
$$\mathbf{T}_i = \begin{matrix} (1-\theta)^4 & (1-\theta)^3\theta & \dots & \theta^4 \\ (1-\theta)^3\theta & (1-\theta)^4 & \dots & (1-\theta)\theta^3 \\ \dots & \dots & \dots & \dots \\ \theta^4 & (1-\theta)\theta^3 & \dots & (1-\theta)^4 \end{matrix}$$

$2^{2n} \times 2^{2n}$ matrix of transitional probabilities between locus i and locus $i+1$

$\sim 10 \times (2^2 \times 2)^2$ operations = 2560 for this case !!!

P(IBD) = 2 at Marker Three

Inheritance vector

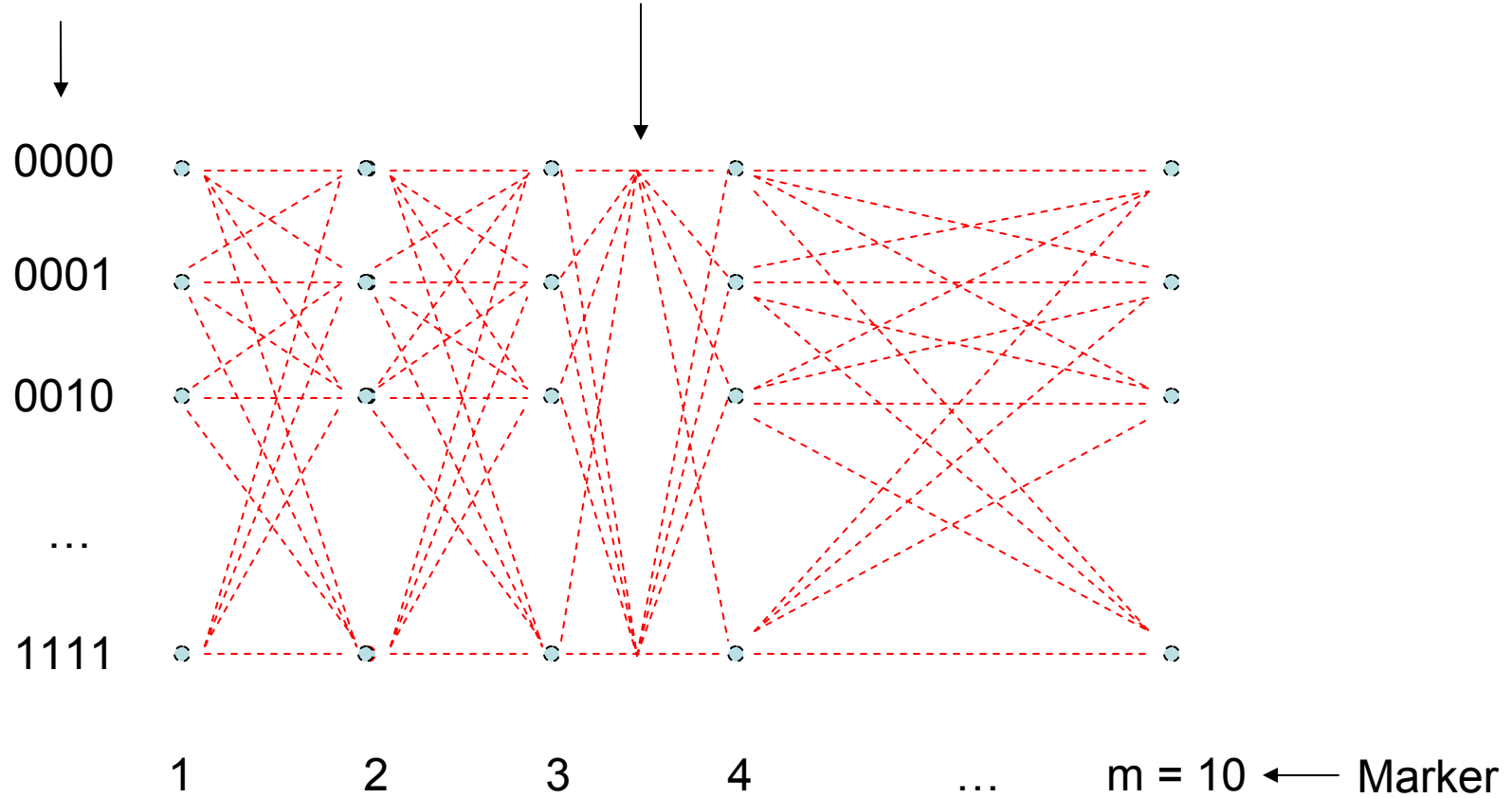


$$L[\text{IBD} = 2 \text{ at marker } 3] / L[\text{ALL}]$$

$$(L[0000] + L[0101] + L[1010] + L[1111]) / L[\text{ALL}]$$

$P(\text{IBD}) = 2$ at arbitrary position on the chromosome

Inheritance vector



$$(L[0000] + L[0101] + L[1010] + L[1111]) / L[\text{ALL}]$$



Further speedups...

- Trees summarize redundant information
 - Portions of inheritance vector that are repeated
 - Portions of inheritance vector that are constant or zero
 - Use sparse-matrix by vector multiplication
- Regularities in transition matrices
 - Use symmetries in divide and conquer algorithm (Idury & Elston, 1997)



Lander-Green Algorithm Summary

- Factorize likelihood by marker
 - Complexity $\propto m \cdot e^n$
- Large number of markers (e.g. dense SNP data)
- Relatively small pedigrees
- MERLIN, GENEHUNTER, ALLEGRO etc



Elston-Stewart Algorithm

- Factorize likelihood by individual
 - Complexity $\propto n \cdot e^m$
- Small number of markers
- Large pedigrees
 - With little inbreeding
- VITESSE etc



Other methods

- Number of MCMC methods proposed
 - \sim Linear on # markers
 - \sim Linear on # people
- Hard to guarantee convergence on very large datasets
 - Many widely separated local minima
- E.g. SIMWALK, LOKI

MERLIN-- Multipoint Engine for Rapid Likelihood Inference

letter

Merlin—rapid analysis of dense genetic maps using sparse gene flow trees

Gonçalo R. Abecasis^{1,2}, Stacey S. Cherny¹, William O. Cookson¹ & Lon R. Cardon¹

Published online: 3 December 2001; DOI: 10.1038/ng786

Efforts to find disease genes using high-density single-nucleotide polymorphism (SNP) maps will produce data sets that exceed the limitations of current computational tools. Here we describe a new, efficient method for the analysis of dense genetic maps in pedigree data that provides extremely fast solutions to common problems such as allele-sharing analysis and haplotyping. We show that sparse binary trees represent patterns of gene flow in general pedigrees in a parsimonious manner, and derive a family of related algorithms for pedigree traversal. With these trees, exact likelihood calculations can be carried out efficiently for single markers or for multiple linked markers. Using an approximate multipoint calculation that ignores the unlikely possibility of a large number of recombinants further improves speed and provides accurate solutions in dense maps with thousands of markers. Our multipoint engine for rapid likelihood inference (Merlin) is a computer program that uses sparse inheritance trees for pedigree analysis. It performs rapid haplotyping, genotype error detection and affected pair linkage analysis and can handle more markers than other pedigree analysis packages.

Linkage and association studies routinely involve analyzing many markers in related individuals to determine phased haplotypes, test for cosegregation of disease and marker loci or identify problems in genotyping. The shift to dense SNP maps^{1,2} poses new problems to pedigree analysis packages³⁻⁷. Packages based on the Elston-Stewart algorithm⁸ can only handle a small number of markers and are not well suited to SNP maps. On the other hand, memory requirements for the Lander-Green algorithm⁹ make analyzing hundreds or thousands of markers a severe challenge in all but the smallest pedigrees. Although Markov-Chain Monte-Carlo (MCMC) sampling methods^{7,10,11} complement some of the deficiencies in these

two approaches, as the number of tightly linked markers increases, it is difficult to guarantee their adequate convergence. Another unresolved issue is undetected genotyping error, which seriously hinders linkage and association studies^{12,13}. As most SNP genotyping errors do not lead to mendelian inconsistencies¹⁴, SNPs require specialized quality-control strategies.

The Lander-Green algorithm⁹ considers each alternative gene flow pattern in a pedigree separately. Allele-sharing statistics for each set of observed phenotypes and likelihoods conditional on observed marker data are calculated and stored in memory⁹, because the pattern of gene flow through a pedigree is fully specified by noting whether the grand-maternal or grand-paternal allele is transmitted in each meiosis, the results of these calculations are typically stored in a bit-indexed array (Fig. 1a), where each index bit indicates the outcome of one meiosis^{6,9,15}. Binary trees provide another natural organization for the results that depend on gene flow patterns. Each level in the tree represents one meiosis, and each branch corresponds to transmission of the grand-maternal or grand-paternal allele (Fig. 1b). Often, many alternative patterns of gene flow have the same outcome, and we reasoned that sparse binary trees might provide an efficient framework for pedigree analysis and extend the scope of the Lander-Green algorithm to very large data sets. These sparse trees are a reduced representation of the full binary tree, where gene flow patterns with identical outcomes are combined into symmetric and premature leaf nodes (Fig. 1c).

We first evaluated the performance of gene flow trees in single marker analyses using simulated replicates of pedigree D (Fig. 2), which includes 40 meioses. Usually the maternal or paternal origin of founder alleles cannot be discerned, and only 2⁴⁰ representative outcomes must be considered⁹. If outcomes were enumerated in

an array, this analysis would exceed the storage capacity of most modern workstations. In comparison, trees describing gene flow pattern likelihoods for SNP markers with equiprobable alleles and 20% missing data have a median size of less than 900 nodes, and are even smaller for more informative markers or smaller amounts of missing data (Table 1). This saves significant amounts of both storage and compiling time, and similar savings result when allele-sharing statistics are calculated for most pedigrees.

Table 1 • Complexity of inheritance tree for pedigree D^a

Missing genotypes	Info ^b	Total nodes			Leaf nodes
		Mean	Median	95% C.I.	
four-allele marker with equiprobable alleles					
—	0.72	154.7	72	64–603	5.2
5%	0.68	245.2	122	64–1,166	9.9
10%	0.64	446.3	171	65–2,409	24.1
20%	0.55	1,747.4	405	69–15,943	107.2
50%	0.28	19,880.6	2,982	154–140,215	2,574.5
two-allele marker with equiprobable alleles					
—	0.42	706.0	151	57–5,447	66.9
5%	0.39	1,299.8	225	57–8,441	159.6
10%	0.36	2,157.7	329	61–15,361	148.9
20%	0.31	8,595.9	872	64–42,592	1,293.9
50%	0.14	55,529.1	4,477	115–183,407	8,173.5

^aPedigree D is represented in Fig. 2. ^bAverage marker informativeness. ^cThe average number of leaf nodes, which correspond to full likelihood evaluations. These statistics summarize 1,000 replicates.



©1998 Jeff Buchbinder



Capabilities

- Linkage Analysis
 - NPL and K&C LOD
 - Variance Components
- Haplotypes
 - Most likely
 - Sampling
 - All
- IBD and info content
- Error Detection
 - Most SNP typing errors are Mendelian consistent
- Recombination
 - No. of recombinants per family per interval can be controlled
- Simulation

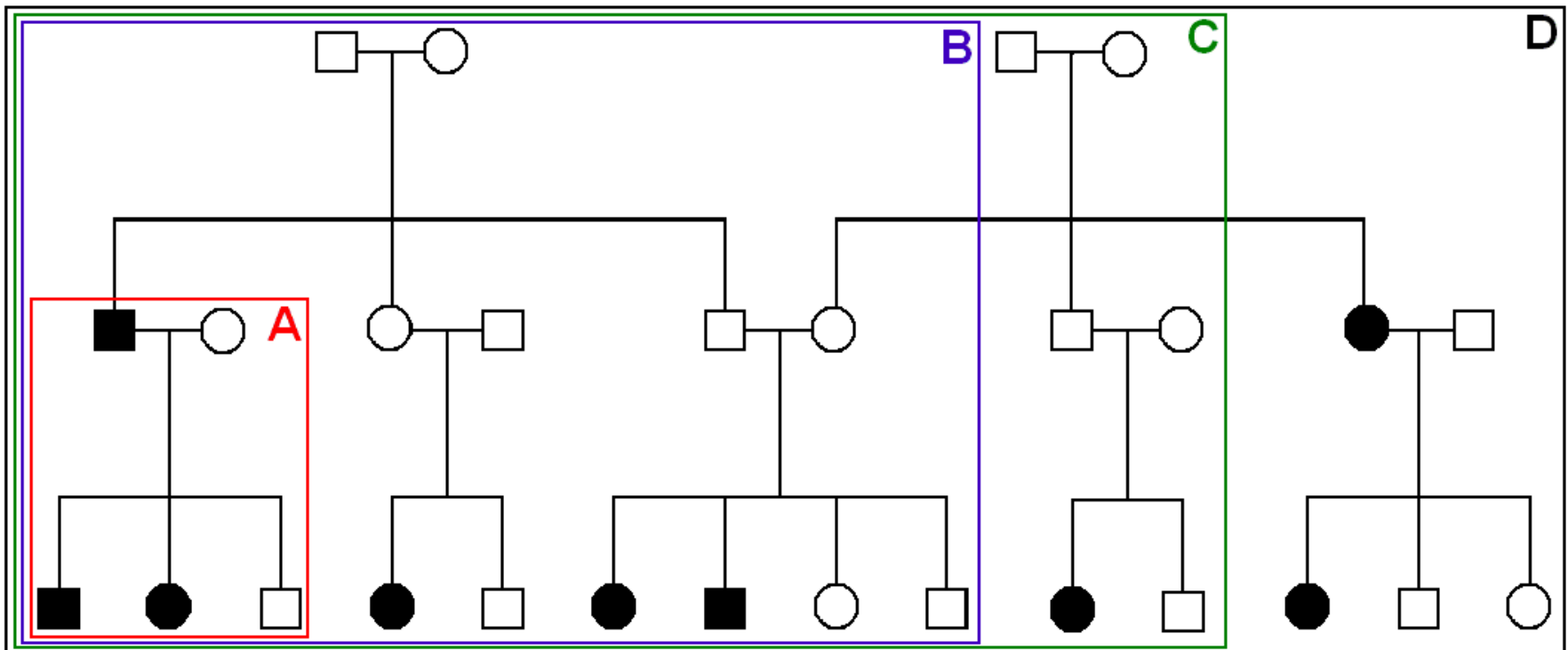


MERLIN Website

www.sph.umich.edu/csg/abecasis/Merlin

- Reference
- FAQ
- Source
- Binaries
- Tutorial
 - Linkage
 - Haplotyping
 - Simulation
 - Error detection
 - IBD calculation

Test Case Pedigrees





Timings – Marker Locations

Top Generation Genotyped				
	A (x1000)	B	C	D
Genehunter	38s	37s	18m16s	*
Allegro	18s	2m17s	3h54m13s	*
Merlin	11s	18s	13m55s	*

Top Generation Not Genotyped				
	A (x1000)	B	C	D
Genehunter	45s	1m54s	*	*
Allegro	18s	1m08s	1h12m38s	*
Merlin	13s	25s	15m50s	*



Intuition: Approximate Sparse T

- Dense maps, closely spaced markers
- Small recombination fractions θ
- Reasonable to set θ^k with zero
 - Produces a very sparse transition matrix
- Consider only elements of \mathbf{v} separated by $< k$ recombination events
 - At consecutive locations



Additional Speedup...

	Time	Memory
Exact	40s	100 MB
No recombination	<1s	4 MB
≤ 1 recombinant	2s	17 MB
≤ 2 recombinants	15s	54 MB
Genehunter 2.1	16min	1024MB

Keavney et al (1998) ACE data, 10 SNPs within gene,
4-18 individuals per family



Input Files

- Pedigree File
 - Relationships
 - Genotype data
 - Phenotype data
- Data File
 - Describes contents of pedigree file
- Map File
 - Records location of genetic markers



Example Pedigree File

<contents of example.ped>

1	1	0	0	1	1	x	3	3	x	x
1	2	0	0	2	1	x	4	4	x	x
1	3	0	0	1	1	x	1	2	x	x
1	4	1	2	2	1	x	4	3	x	x
1	5	3	4	2	2	1.234	1	3	2	2
1	6	3	4	1	2	4.321	2	4	2	2

<end of example.ped>

Encodes family relationships, marker and phenotype information



Example Data File

<contents of example.dat>

T some_trait_of_interest

M some_marker

M another_marker

<end of example.dat>

Provides information necessary to decode
pedigree file



Data File Field Codes

Code	Description
M	Marker Genotype
A	Affection Status.
T	Quantitative Trait.
C	Covariate.
Z	Zygosity.



Example Map File

<contents of example.map>

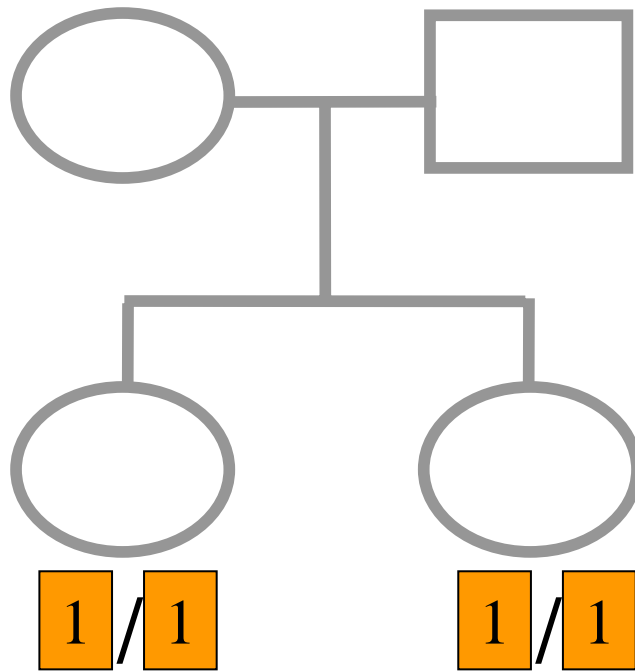
CHROMOSOME	MARKER	POSITION
2	D2S160	160.0
2	D2S308	165.0

...

<end of example.map>

Indicates location of individual markers,
necessary to derive recombination fractions
between them

Worked Example



$$p_1 = 0.5$$

$$P(IBD=0|G) = \frac{1}{9}$$

$$P(IBD=1|G) = \frac{4}{9}$$

$$P(IBD=2|G) = \frac{4}{9}$$

```
merlin -d example.dat -p example.ped -m example.map --ibd
```



Application: Information Content Mapping

- Information content: Provides a measure of how well a marker set approaches the goal of completely determining the inheritance outcome
- Based on concept of entropy
 - $E = -\sum P_i \log_2 P_i$ where P_i is probability of the i th outcome
- $I_E(x) = 1 - E(x)/E_0$
 - Always lies between 0 and 1
 - Does not depend on test for linkage
 - Scales linearly with power

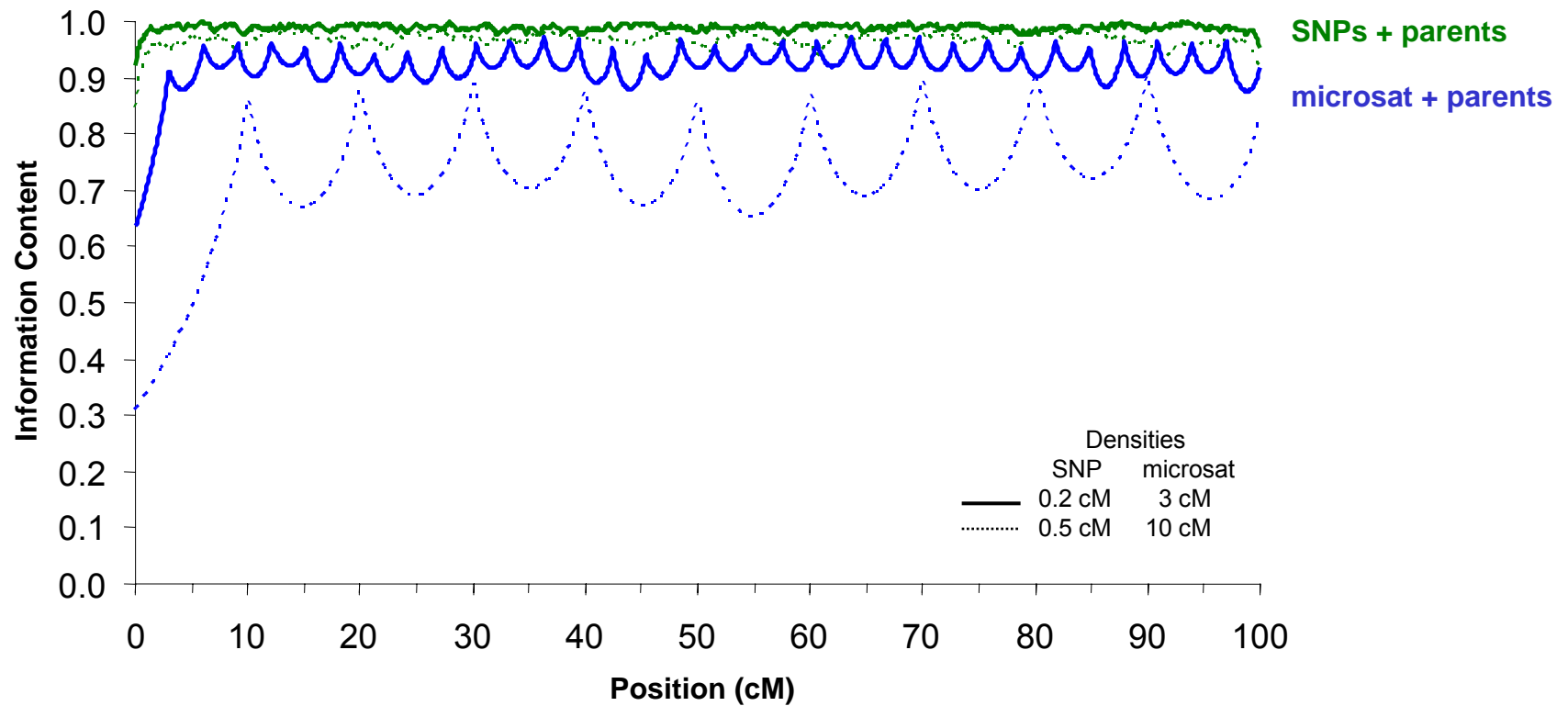


Application: Information Content Mapping

- Simulations (sib-pairs with/out parental genotypes)
 - 1 micro-satellite per 10cM (ABI)
 - 1 microsatellite per 3cM (deCODE)
 - 1 SNP per 0.5cM (Illumina)
 - 1 SNP per 0.2 cM (Affymetrix)
- Which panel performs best in terms of extracting marker information?
- Do the results depend upon the presence of parental genotypes?

```
merlin -d file.dat -p file.ped -m file.map --information --step 1 --markerNames
```

SNPs vs Microsatellites with parents



SNPs vs Microsatellites without parents

