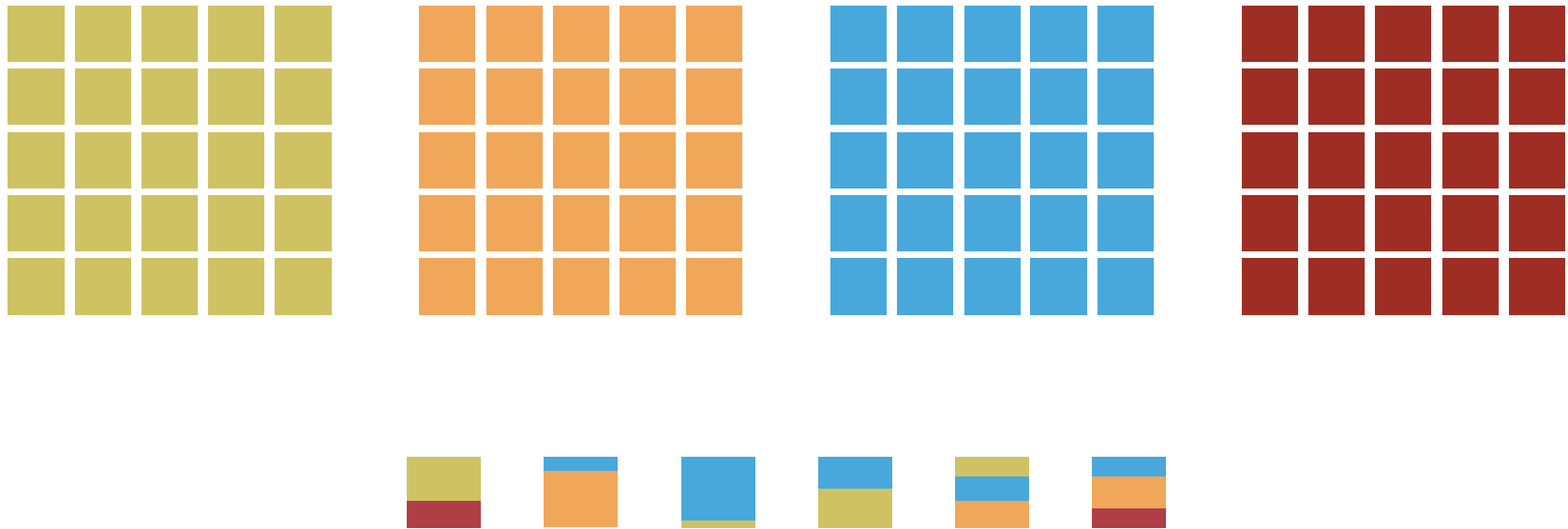


Reference Populations (known or unknown)



Identify **ancestry proportions** for individuals with **admixed** ancestry

Approaches: Structure (MCMC, Bayesian)
Or ADMIXTURE (quadratic programming)

Genetic Structure of Human Populations

Noah A. Rosenberg,^{1*} Jonathan K. Pritchard,² James L. Weber,³
Howard M. Cann,⁴ Kenneth K. Kidd,⁵ Lev A. Zhivotovsky,⁶
Marcus W. Feldman⁷

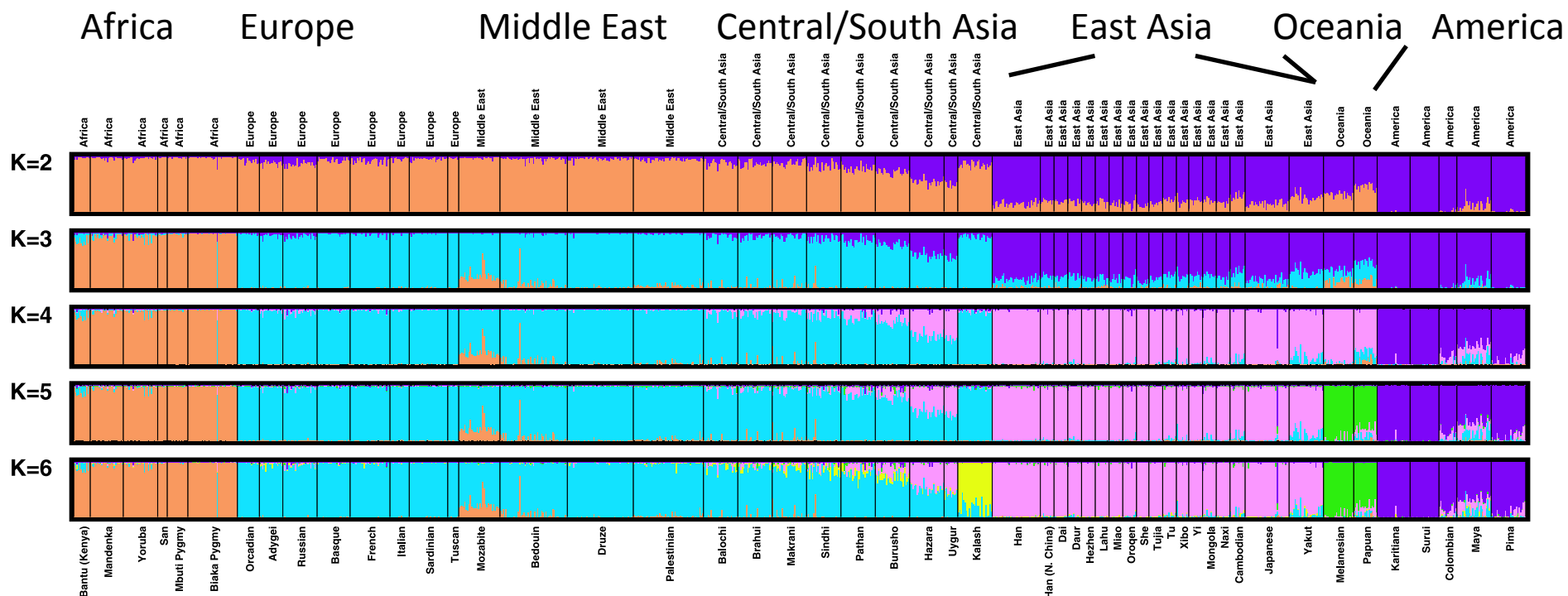
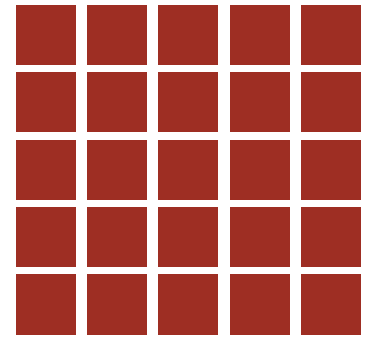
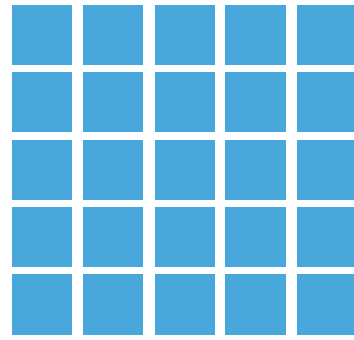
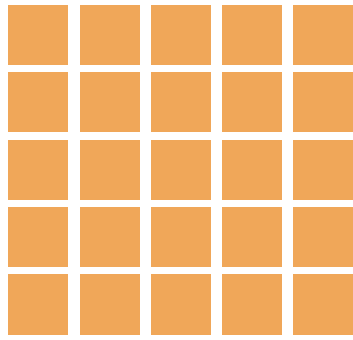
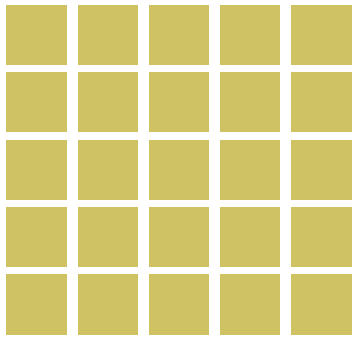


Fig. 1. Estimated population structure. Each individual is represented by a thin vertical line, which is partitioned into K colored segments that represent the individual's estimated membership fractions in K clusters. Black lines separate individuals of different populations. Populations are labeled below the figure, with their regional affiliations above it. Ten *structure* runs at each

K produced nearly identical individual membership coefficients, having pairwise similarity coefficients above 0.97, with the exceptions of comparisons involving four runs at $K = 3$ that separated East Asia instead of Eurasia, and one run at $K = 6$ that separated Karitiana instead of Kalash. The figure shown for a given K is based on the highest probability run at that K .

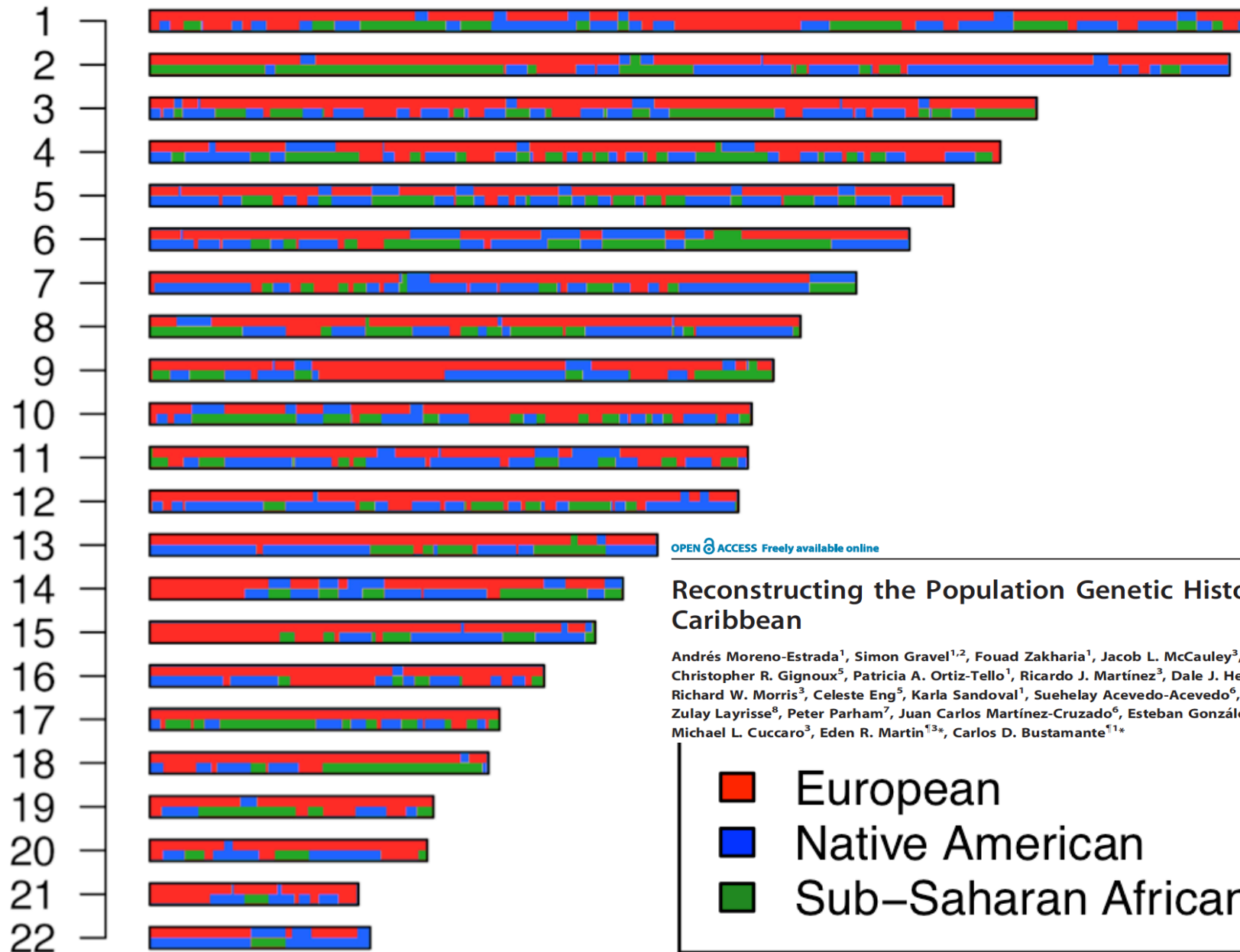
Reference Populations



Identify origins of **chromosomal segments** in individuals of **admixed** ancestry

Approaches: Based on Hidden Markov Models

Local ancestry



OPEN ACCESS Freely available online

PLOS GENETICS

Reconstructing the Population Genetic History of the Caribbean

Andrés Moreno-Estrada¹, Simon Gravel^{1,2}, Fouad Zakharia¹, Jacob L. McCauley³, Jake K. Byrnes^{1,4}, Christopher R. Gignoux⁵, Patricia A. Ortiz-Tello¹, Ricardo J. Martínez³, Dale J. Hedges³, Richard W. Morris³, Celeste Eng⁵, Karla Sandoval¹, Suehelay Acevedo-Acevedo⁶, Paul J. Norman⁷, Zulay Layrisse⁸, Peter Parham⁷, Juan Carlos Martínez-Cruzado⁶, Esteban González Burchard⁵, Michael L. Cuccaro³, Eden R. Martin^{1,3*}, Carlos D. Bustamante^{1,1*}

■ European
■ Native American
■ Sub-Saharan African

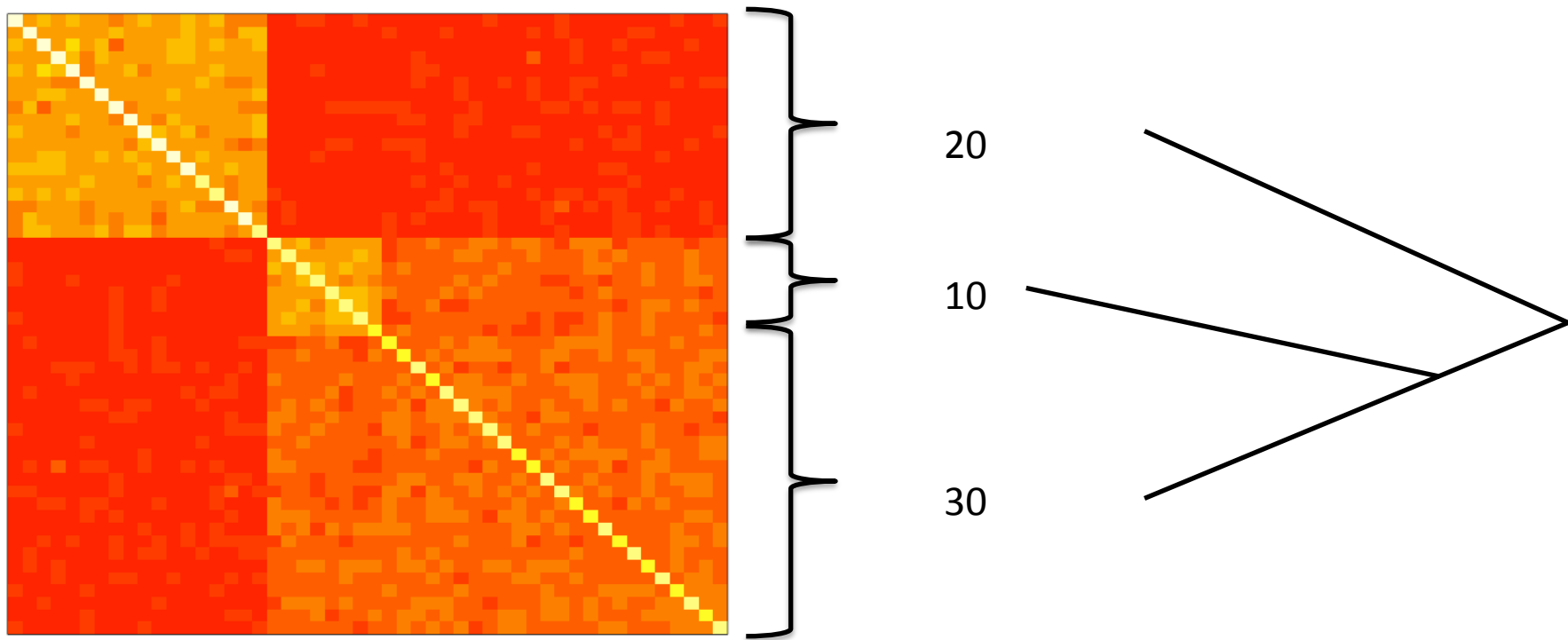
Principal components analysis

Example

Section 1.7.2 of notes

(Simulated data, N=50 individuals, L=1000 SNPs)

Relatedness matrix R



i^{th} and j^{th} entry = average over loci (l) of $(X_{li} - \bar{X}_l)(X_{lj} - \bar{X}_l)$

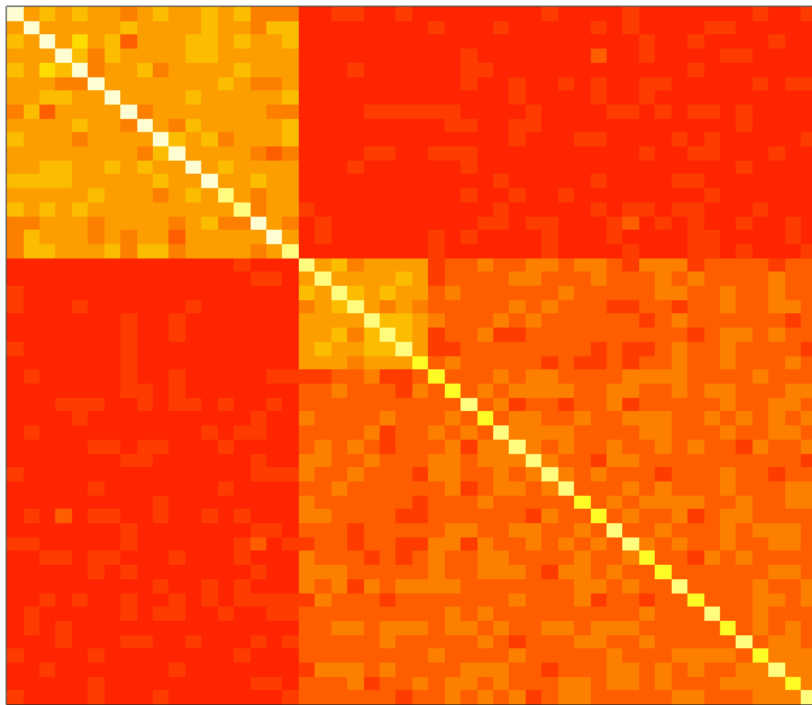
Where X_l is mean freq. of the l^{th} locus.

Modified from slide by Gavin Band

Principal components analysis Example

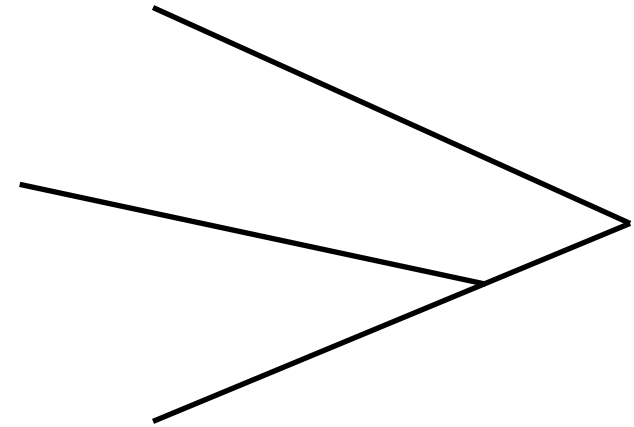
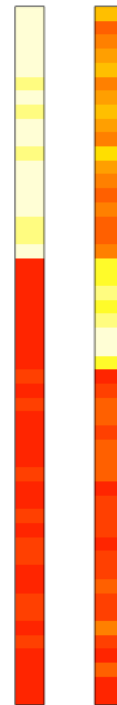
(Simulated data, 50 individuals, 1000 SNPs)

Relatedness matrix R



Principal components
(Eigen-vectors)

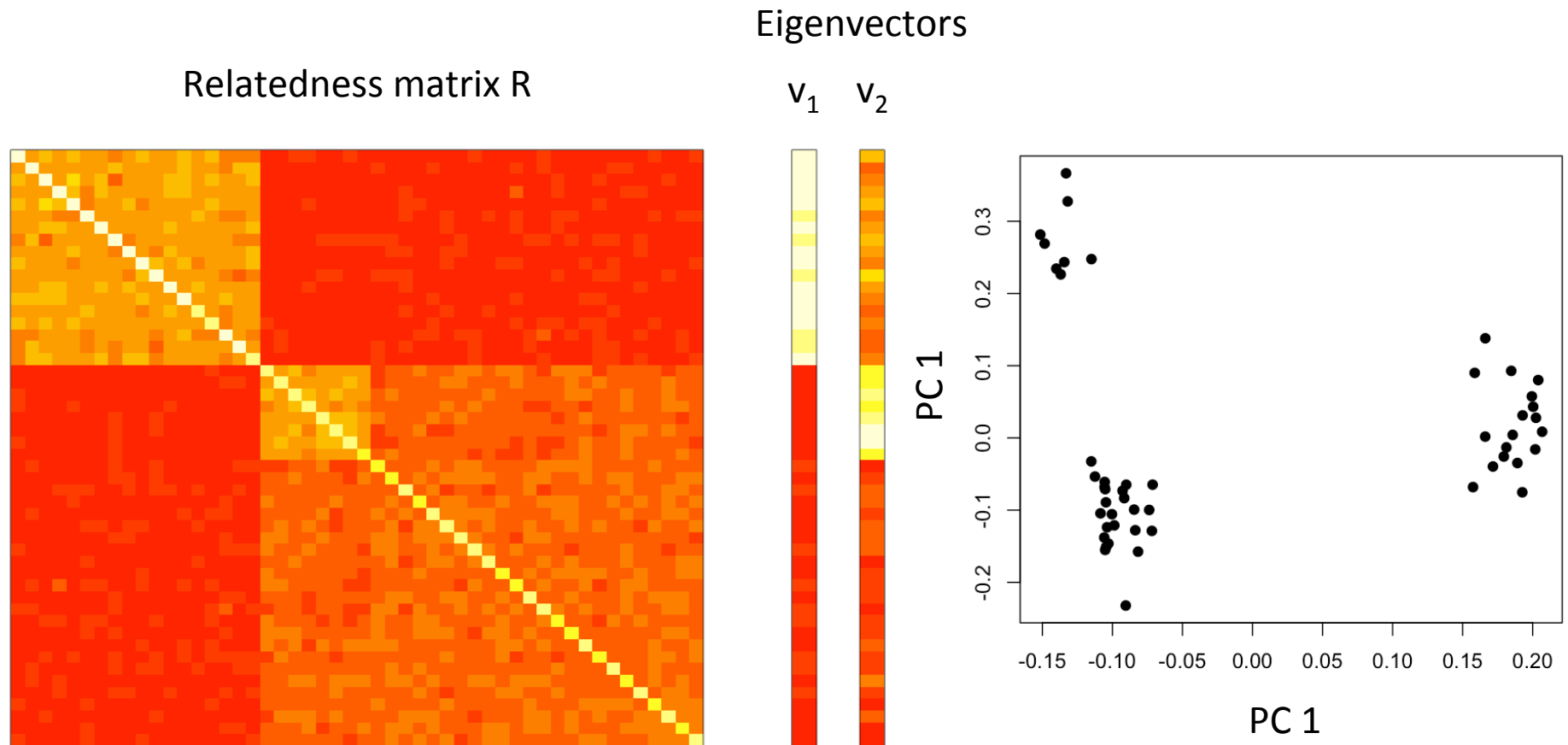
v_1 v_2



Modified from slide by Gavin Band

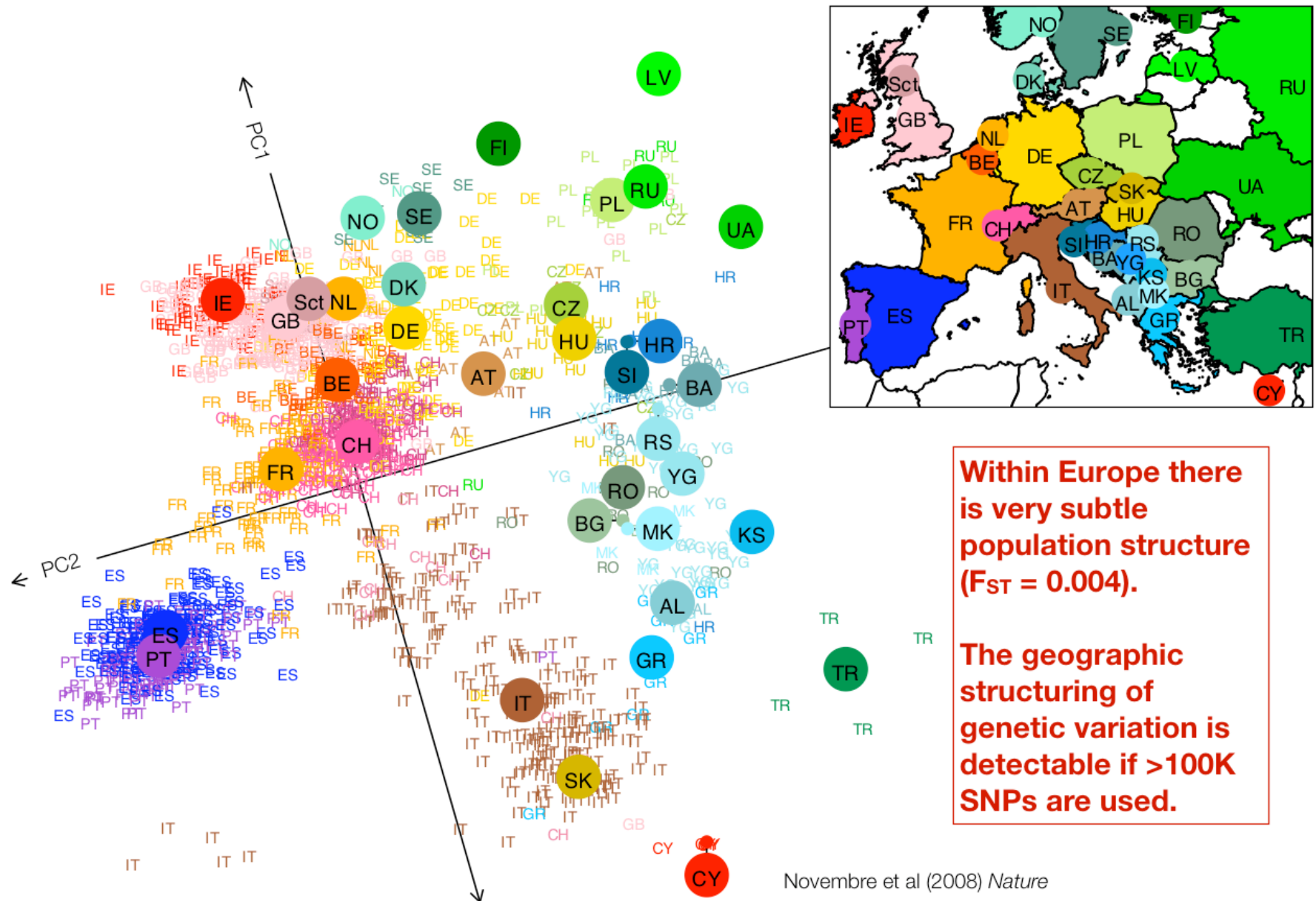
Principal components analysis Example

(Simulated data, 50 individuals, 1000 SNPs)



Modified from slide by Gavin Band

Principal Component Analysis of Europeans

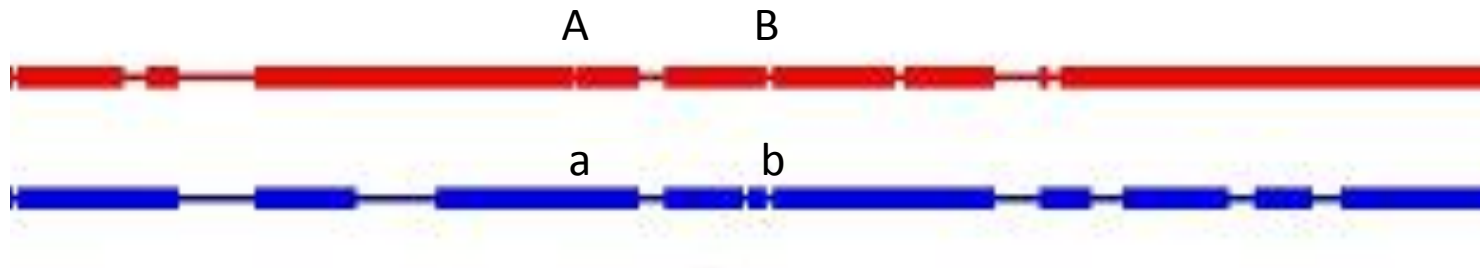


Within Europe there is very subtle population structure ($F_{ST} = 0.004$).

The geographic structuring of genetic variation is detectable if >100K SNPs are used.

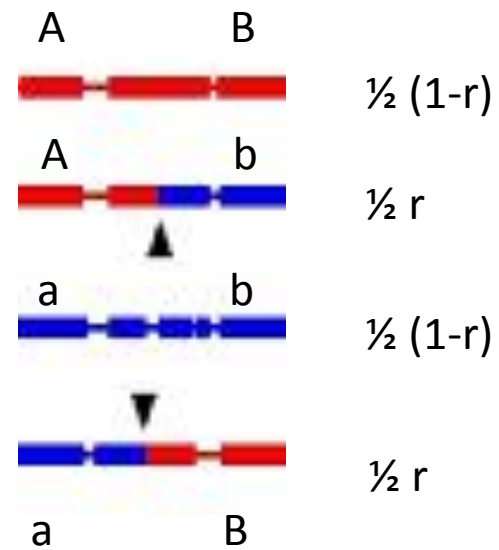
Novembre et al (2008) *Nature*

Recombination and Linkage Disequilibrium (LD)



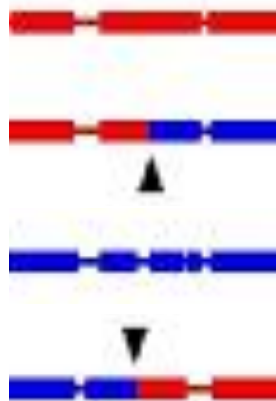
r = recombination fraction
 probability of an odd
 Number of crossovers occur
 Between our markers

$$0 < r < \frac{1}{2}$$



Linkage disequilibrium: The non-random association of alleles at different sites in the genome in a population.

If independent the expected



p_{AB} = frequency of AB

p_{ab} = frequency of ab

p_{Ab} = frequency of Ab

p_{aB} = frequency of aB

frequency of gametes (haplotypes)

$$p_A \times p_B$$

$$p_a \times p_b$$

$$p_A \times p_b$$

$$p_a \times p_B$$

Define “D”

$$D_{AB} = p_{AB} - p_A p_B$$

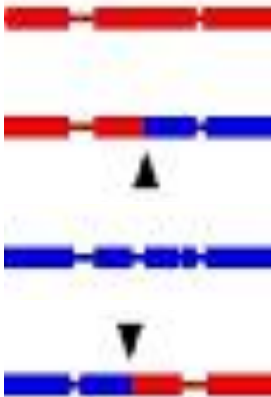
$$D_{ab} = p_{ab} - p_a p_b$$

$$D_{Ab} = p_{Ab} - p_A p_b$$

$$D_{aB} = p_{aB} - p_a p_B$$

The covariance of A and B.

Linkage disequilibrium: The non-random association of alleles at different sites in the genome.



Define “D”

$$D_{AB} = p_{AB} - p_A p_B$$

$$D_{ab} = p_{ab} - p_a p_b$$

$$D_{Ab} = p_{Ab} - p_A p_b$$

$$D_{aB} = p_{aB} - p_a p_B$$

$$D_{AB} = - D_{Ab}$$

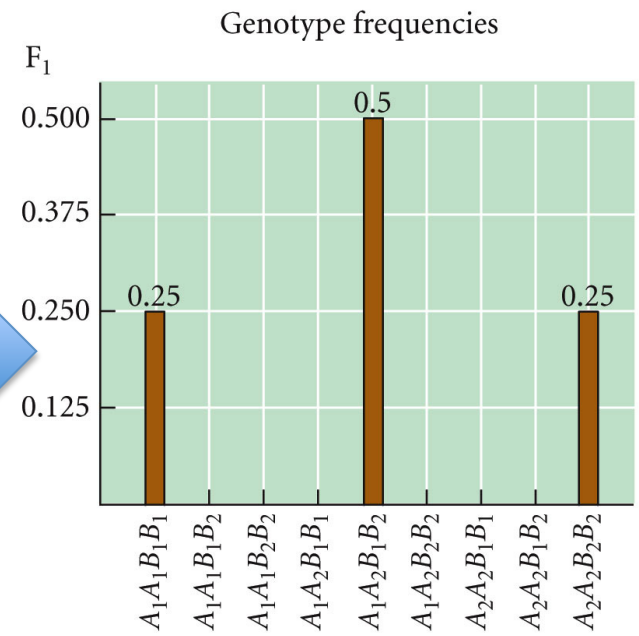
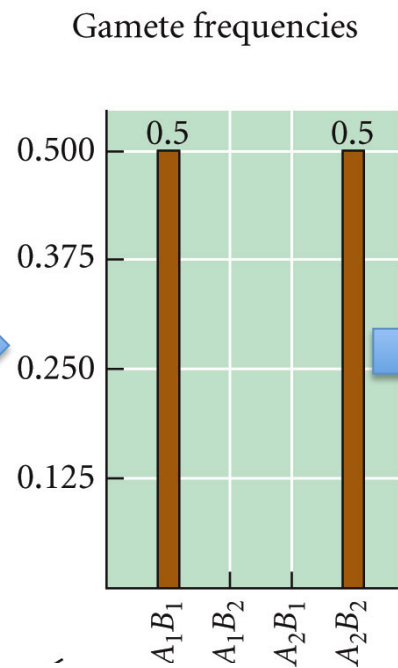
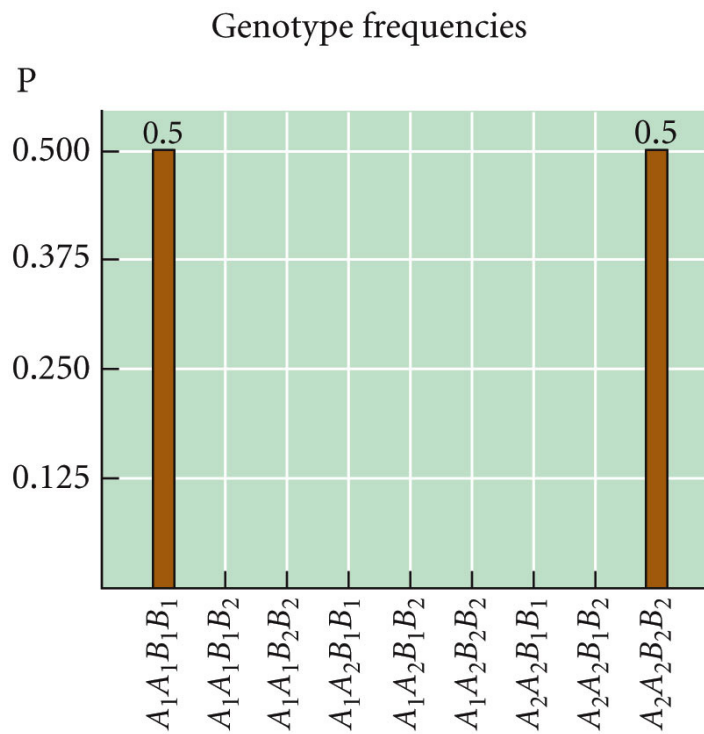
$$D_{AB} = D_{ab} \text{ and } D_{Ab} = D_{aB}$$

(so, knowing D_{AB} is enough - call this “D”)

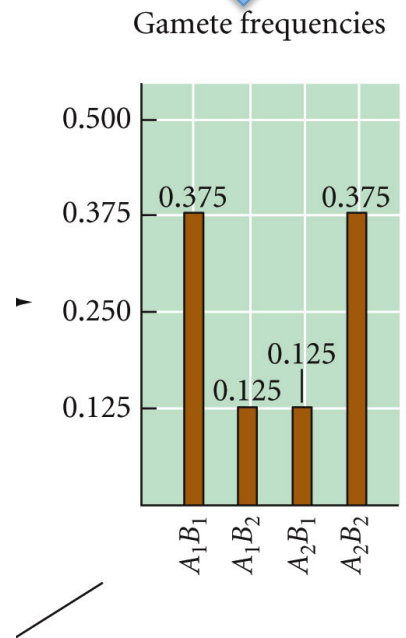
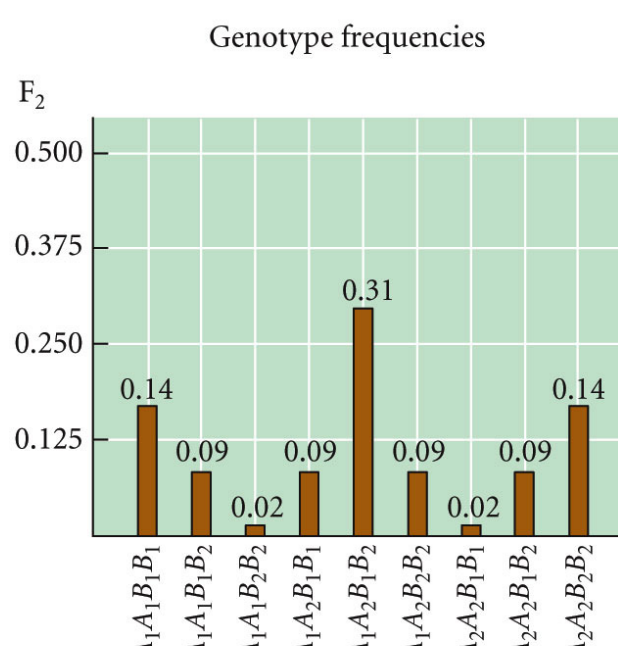
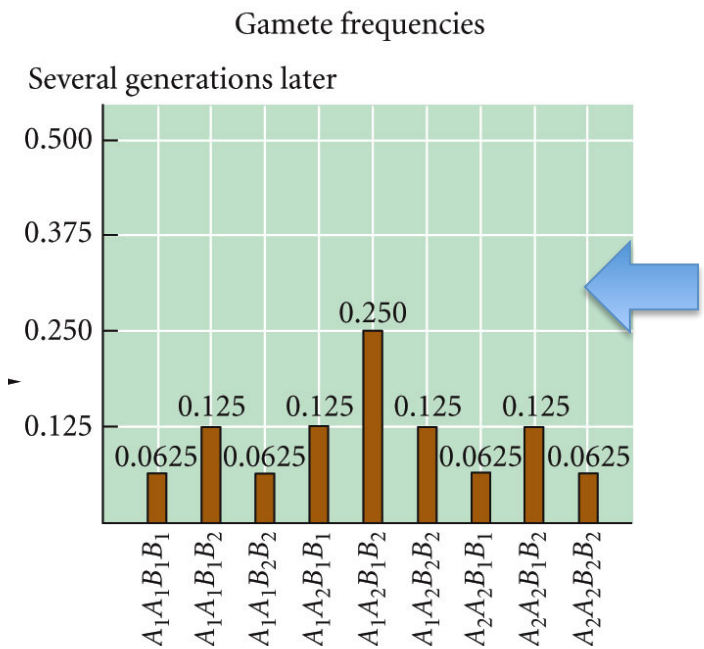
If $O = E$, then $D = 0$

If $D > 0$ (or $D < 0$) then there is “linkage disequilibrium (LD)”

Note: you can also write $p_{AB} = p_A p_B + D$



EVOLUTION 2e, Figure 9.17



Decay of LD in a very large boring
randomly mating population

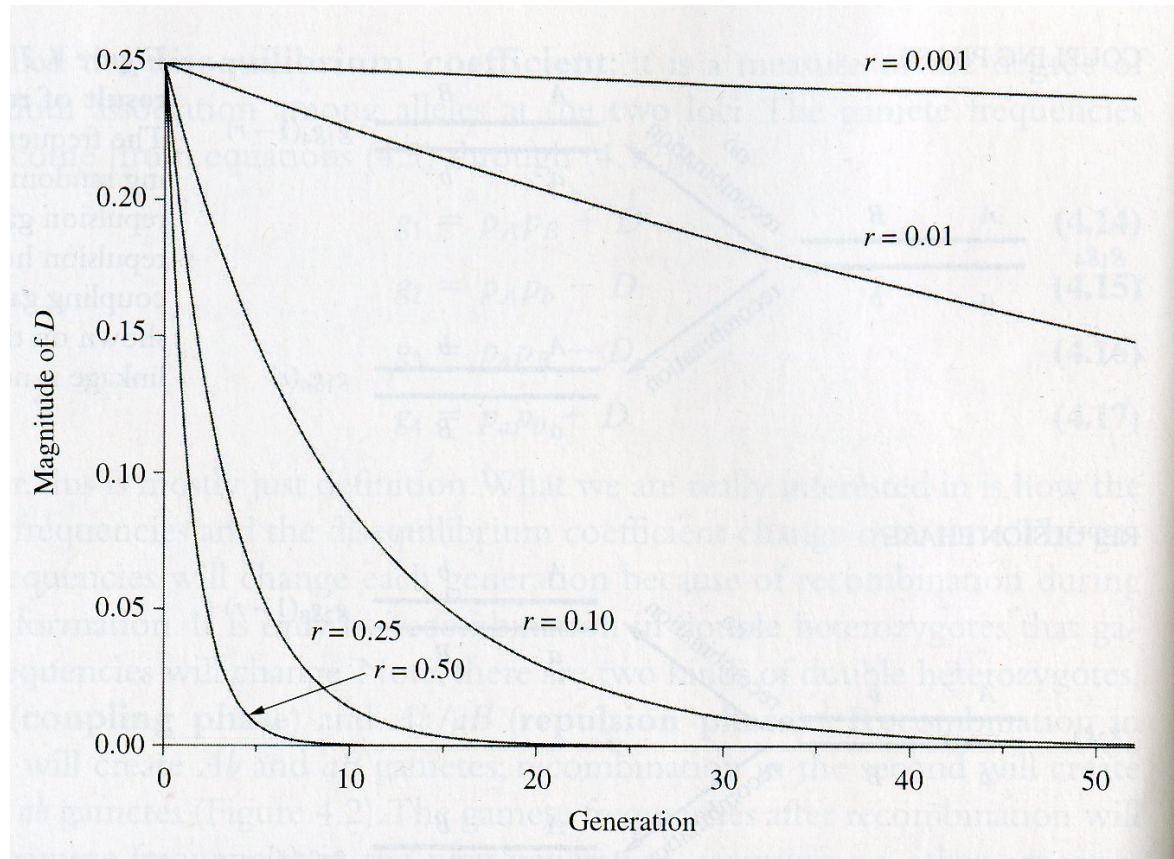
$$D_t = (1 - r)^t D_o$$

With inbreeding coefficient f replace r with $r(1-f)$

linkage disequilibrium

How does LD change over time due to recombination?

$$D_t = (1 - r)^t D_o$$



Note: more distant markers recombine more!

So eventually recombination leads to $D=0$.

Even with free recombination ($r=0.5$), it isn't instantaneous