



**Questions we are addressing**

- How reproductively isolated are a group of populations?
- How many populations are in sample?

**Allele Frequencies**

- Allele frequencies (gene frequencies) = proportion of all alleles in an all individuals in the group in question which are a particular type
- Allele frequencies:
  - $p + q = 1$
- Expected genotype frequencies:
  - $p^2 + 2pq + q^2$

**Hardy-Weinberg Equilibrium**

- Null Model = population is in HW Equilibrium
  - Useful
  - Often predicts genotype frequencies well

**Hardy-Weinberg Theorem**

if only random mating occurs, then allele frequencies remain unchanged over time.

After one generation of random-mating, genotype frequencies are given by

AA	Aa	aa
$p^2$	$2pq$	$q^2$

$p = \text{freq (A)}$   
 $q = \text{freq (a)}$

**Expected Genotype Frequencies**

- The possible range for an allele frequency or genotype frequency therefore lies between (0 – 1)
- with 0 meaning complete absence of that allele or genotype from the population (no individual in the population carries that allele or genotype)
- 1 means complete **fixation** of the allele or genotype (fixation means that every individual in the population is homozygous for the allele -- i.e., has the same genotype at that locus).

**ASSUMPTIONS**

- 1) diploid organism
- 2) sexual reproduction
- 3) generations are non-overlapping
- 4) mating occurs at random
- 5) large population size
- 6) migration = 0
- 7) mutation = 0
- 8) no selection on genes

For a population to be in Hardy Weinberg Equilibrium, the observed genotype frequencies must match those predicted by the equation  $p^2 + 2pq + q^2$

Finding "p": equals frequency of AA + 1/2(frequency of Aa)  
 Finding "q": equals frequency of aa + 1/2(frequency of Aa)

Then  $p^2$  = predicted genotype frequency of AA  
 $q^2$  = predicted genotype frequency of aa  
 $2pq$  = predicted genotype frequency of Aa

	$f_{AA}$	$f_{Aa}$	$f_{aa}$
Population 1	0.3	0.0	0.7
Population 2	0.2	0.2	0.6
Population 3	0.1	0.4	0.5

In each group, the allele frequency is actually the same

Population 1     $A = 0.3$   
                      $a = 0.7$

Population 2     $A = 0.2 + .5(0.2) = 0.3$   
                      $a = 0.6 + .5(0.2) = 0.7$

Population 3     $A = 0.1 + .5(0.4) = 0.3$   
                      $a = 0.5 + .5(0.4) = 0.7$

Only **population 3** is in Hardy Weinberg equilibrium (need to do chi-square test)  
 $p^2 + 2pq + q^2 = 1$   
 $(0.3)^2 + 2(0.3)(0.7) + (0.7)^2 = 1$   
 $0.09 + 0.42 + 0.49 = 1$

### IMPORTANCE OF HW THEOREM

If the only force acting on the population is random mating, allele frequencies remain unchanged and genotypic frequencies are constant.

Mendelian genetics implies that genetic variability can persist indefinitely, unless other evolutionary forces act to remove it

### Population Substructure

- Many species naturally subdivide themselves into herds, flocks, colonies, schools etc.
- Patchy environments can also cause subdivision
- Subdivision decreases heterozygosity and generates genetic differentiation via:
  - Natural selection
  - Genetic drift

### Wright's Fixation Index

- Equals the reduction in heterozygosity expected with random mating at one level of population hierarchy relative to another more inclusive level.

$$F_{ST} = (H_T - H_S) / H_T$$

Step 1: Calculate mean heterozygosities at each population level

- Heterozygosity = mean percentage of heterozygous individuals per locus
- Assuming H-W, heterozygosity ( $H$ ) =  $2pq$  where  $p$  and  $q$  represent mean allele frequencies
- $H_S$  = sum of all subpopulation heterozygosities divided by the total number of subpopulations

	Pop1	Pop2	Pop3
Sample size	20	20	20
AA	10	5	0
Aa	4	10	8
aa	6	5	12

	Pop1	Pop2	Pop3
Freq			
p	$(20 + 1/2*8)/40 = 0.60$	$(10 + 1/2*20)/40 = .50$	$(0 + 1/2*16)/40 = 0.20$
q	$(12 + 1/2*8)/40 = 0.40$	$(10 + 1/2*20)/40 = .50$	$(24 + 1/2*16)/40 = 0.80$

- Calculate  $H_S$  ( $2pq$ )
  - Pop1:  $2*0.60*0.40 = 0.48$
  - Pop2:  $2*0.50*0.50 = 0.50$
  - Pop3:  $2*0.20*0.80 = 0.32$

Pop	$H_S$	p	q	$H_T$	$F_{ST}$
1	0.48	0.60	0.40		
2	0.50	0.50	0.50		
3	0.32	0.20	0.80		
Mean	0.43	0.43	0.57	0.49	0.12

### Interpreting $F_{ST}$

- Can range from 0 (no genetic differentiation) to 1 (fixation of alternative alleles).
- Wright's Guidelines:
  - 0 - 0.05, little differentiation
  - 0.05 - 0.15, moderate
  - 0.15 - 0.25, great
  - > 0.25, very great

## $F_{ST}$ for various organisms

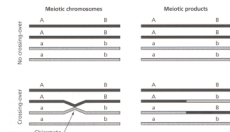
Organism	Number of Populations	Number Loci	$H_B$	$H_S$	$F_{ST}$
Human (major races)	37		0.131	0.121	0.085
Yanomama Indian Villages	37	15	0.039	0.036	0.077
Hawai moale	4	40	0.007	0.006	0.113
Jumping rodent	9	16	0.007	0.012	0.676
Fruit fly	5	27	0.201	0.179	0.109
Horseshoe crab	4	25	0.068	0.061	0.114
Lycobid plant	4	19	0.071	0.051	0.282

## Multiple loci

- Should be in HW equilibrium at all loci
- If recombination occurs freely, all combinations of alleles should be found equally

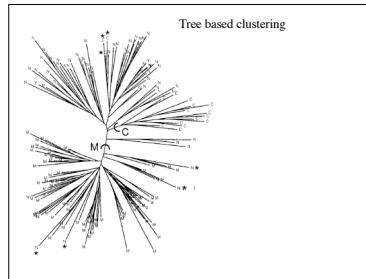
## Linkage Disequilibrium

- $D = P_{AB} - P_A P_B$
- $D = 0$  when population is in equilibrium



## Inferring numbers of populations

- Phylogenetic/distance methods
  - Generate tree from population data, and look for obvious structure.



## Inferring numbers of populations: clustering methods

- Minimize departures from HW
- Minimize linkage disequilibrium
- Assign individuals to K populations so achieve linkage equilibrium and HW equilibrium

## Program Structure is a model based clustering method to infer population structure

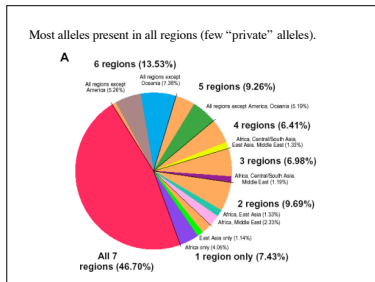
- Main parameter is K, the number of populations
- Run models several times with different K values and compare likelihoods.

## Inferring numbers of populations: clustering methods

- Minimize departures from HW
- Minimize linkage disequilibrium
- Assign individuals to K populations so achieve linkage equilibrium and HW equilibrium

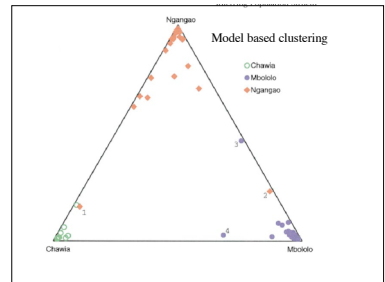
## Why care about population structure?





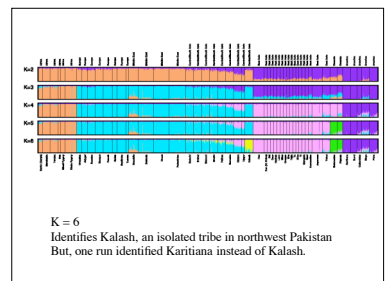
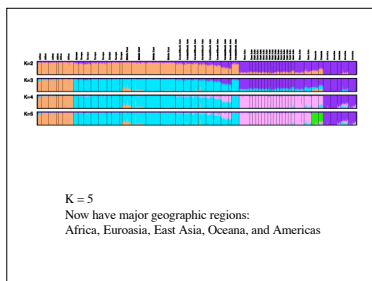
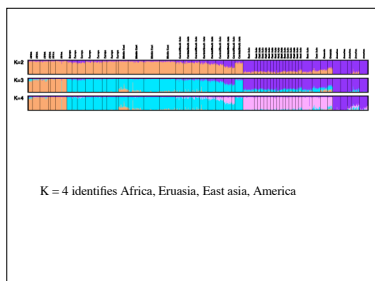
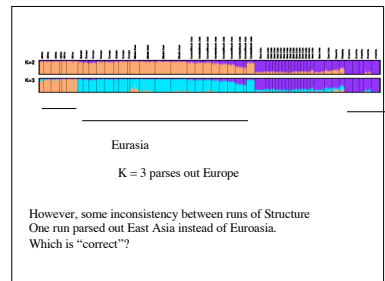
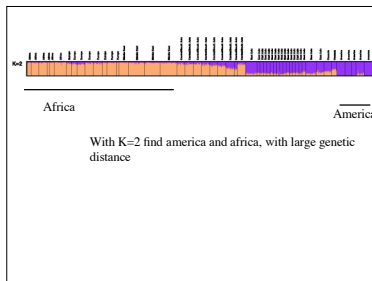
With so little variation between population, how to identify populations

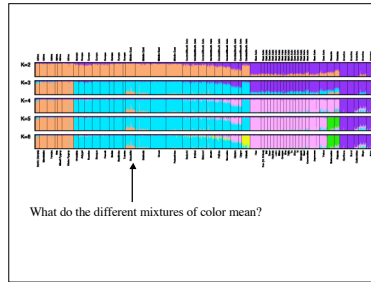
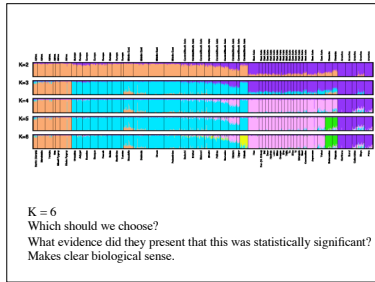
- Look at allele frequencies.
- Cluster to identify groups with distinctive allele frequencies
  - Minimize departures from H-W equilibrium
  - Minimize LD



First cluster entire populations

- Run at a variety of K populations





Supplementary Table 1. Genetic distances between regional groups. The consistency or  $F_{st}$  distance (32, p. 104) was estimated using GDA (21).

	Africa	Europe	Middle East	Central/South Asia	East Asia	Oceania
Europe	0.010					
Middle East	0.032	0.025				
Central/South Asia	0.027	0.026	0.028			
East Asia	0.054	0.028	0.026	0.026		
Oceania	0.068	0.041	0.029	0.029	0.047	
America	0.101	0.070	0.081	0.068	0.060	0.102

Example output from Structure

Estimated Ln Prob of Data = -4018.0  
Mean value of ln likelihood = -3874.2  
Variance of ln likelihood = 287.4  
Mean value of alpha = 0.0432

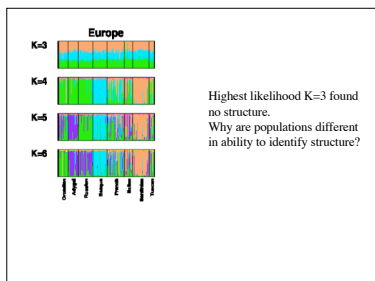
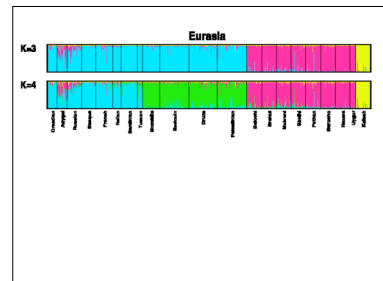
Mean value of  $F_{st_1}$  = 0.0703  
Mean value of  $F_{st_2}$  = 0.0643  
Mean value of  $F_{st_3}$  = 0.1004

Inferred ancestry of individuals:  
Label (Miss) : Inferred clusters (and 90% probability intervals)

1	1	(0) : 0.414 0.566 0.018	(0.000,1.000)	(0.000,1.000)	(0.000,0.116)
2	2	(0) : 0.633 0.356 0.011	(0.000,1.000)	(0.000,1.000)	(0.000,0.061)
3	3	(0) : 0.378 0.610 0.012	(0.000,1.000)	(0.000,1.000)	(0.000,0.070)
4	4	(0) : 0.707 0.285 0.008	(0.000,1.000)	(0.000,1.000)	(0.000,0.039)
5	5	(0) : 0.314 0.650 0.036	(0.000,1.000)	(0.000,1.000)	(0.000,0.264)
6	6	(0) : 0.649 0.342 0.009	(0.000,1.000)	(0.000,1.000)	(0.000,0.053)
7	7	(0) : 0.659 0.330 0.010	(0.000,1.000)	(0.000,1.000)	(0.000,0.059)
8	8	(0) : 0.491 0.490 0.019	(0.000,1.000)	(0.000,1.000)	(0.000,0.131)
9	9	(0) : 0.587 0.375 0.038	(0.000,1.000)	(0.000,1.000)	(0.000,0.230)

Now, can we break down each population further

- Analyze individual populations for structures
- Why not just keep increasing K, like up to 52?



Genetic drift occurs rapidly in small populations, particularly those isolated quickly accrue “distinctive” allele frequencies. This is the signal structure is good at identifying.

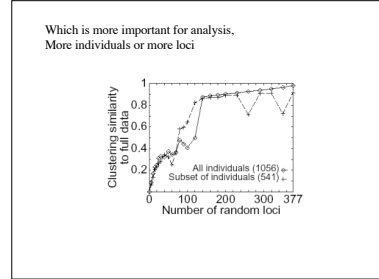
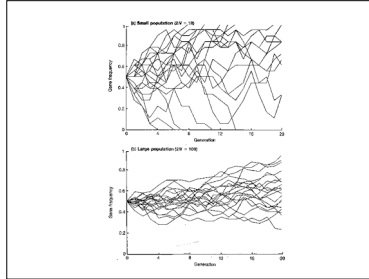
Principles of Neutral Theory

- If a population contains a neutral allele with frequency  $P_0$ , then the probability of fixation of the allele is  $P_0$ 
  - Since initial frequency is  $1/2N$ , the probability of fixation is  $1/2N$
- Can also calculate rates of fixation and loss
  - Both depend on effective population size

Table 1. Average time to fixation or loss of neutral alleles.

$N^*$	Average generations until loss <sup>a</sup>	Average generations until fixation <sup>b</sup>
50	7.4	160
100	8.5	320
250	9.9	800
500	11.1	1,600
1,000	12.2	3,200
5,000	14.7	16,000
10,000	15.8	32,000
50,000	18.4	160,000

<sup>a</sup> Tabulated values assume that  $N_e = 0.8N$ .  
<sup>b</sup> Averages are for selectively neutral alleles destined to be lost or fixed, respectively.



### Major conclusions

- Most variation is within-population differences
- Identified six main genetic clusters, five of which correspond to major geographic regions (?)
- General agreement between genetic and self-reported ancestry

### Future studies

- Models of human migration
- Admixture in human origins
  - Did humans breed with other hominids during evolution?
- Application to other organisms