

# Population structure

Fernando Racimo

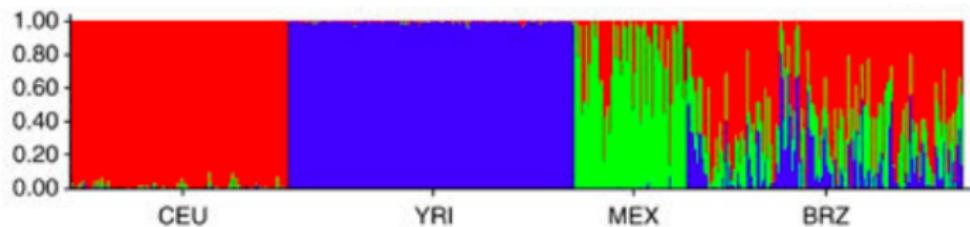
Adelaide, January 2018

# Today

- $F_{ST}$  and population structure
- The mixed-membership clustering model (“Structure”)
- Leveraging haplotype information

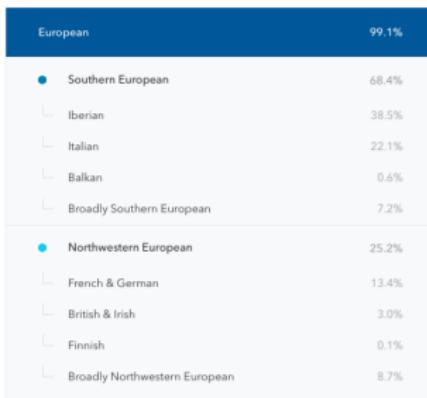
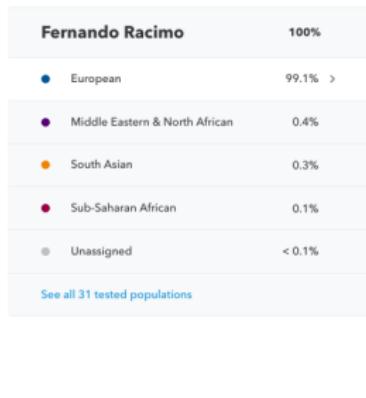
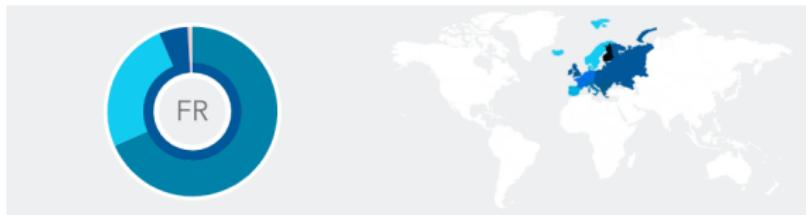
## Questions

- Is there population structure in a population?
- Can we identify subpopulation clusters of shared ancestry?
- Are individuals best modeled as mixtures of ancestral populations?
- How much admixture was passed on from each population?



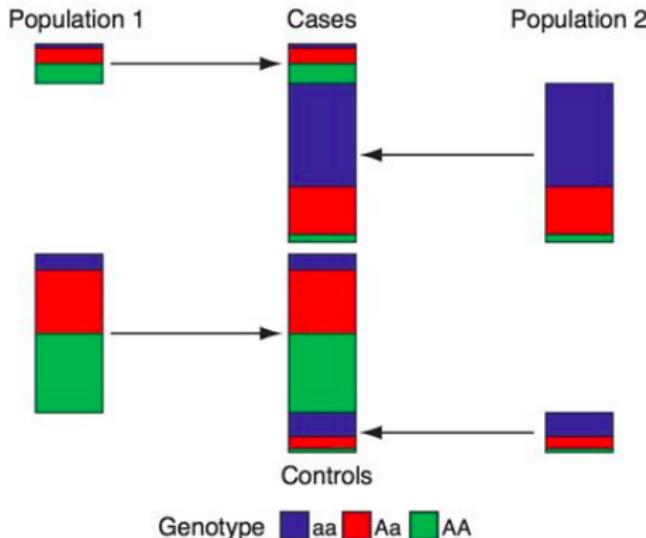
# Objectives

- Learn something about the pasty genetic history of a population under study
- Learn something about ourselves



# Objectives

- Correct for population stratification in downstream analyses (e.g. GWAS)

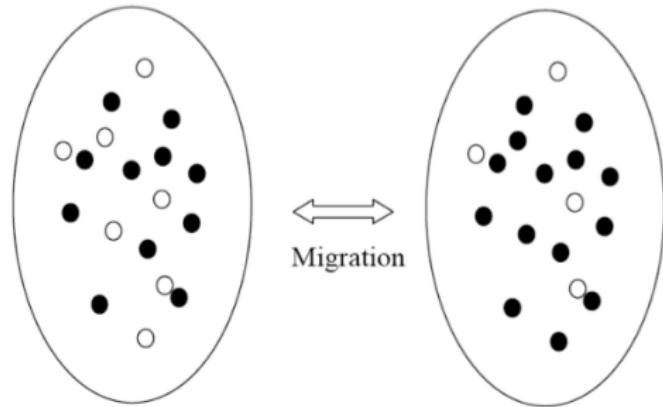


# Today

- $F_{ST}$  and population structure
- The mixed-membership clustering model (“Structure”)
- Leveraging haplotype information

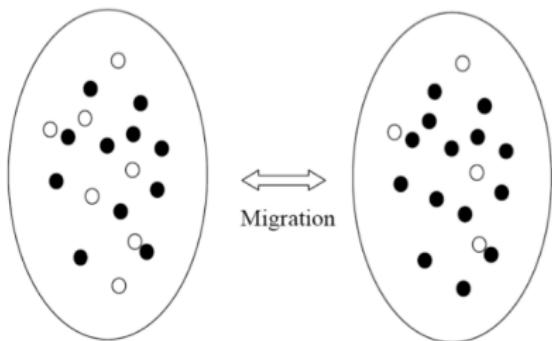
# Population structure

- When a population is not randomly mating due to geographic structure, we say there is **population structure**.
- $F_{st}$  is way to measure how much structure there is in a population.
- Before we define it, let us start by examining a simple model with two subpopulations, each of which contains the same number of individuals and is under HWE.



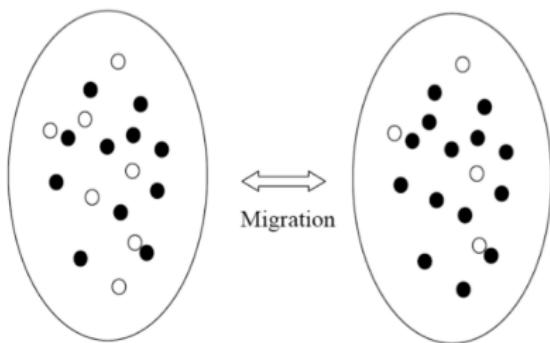
# Population structure

- The frequency of an allele (A) in subpopulation 1 is  $f_{A1}$ . The frequency of the same allele in subpopulation 2 is  $f_{A2}$ .



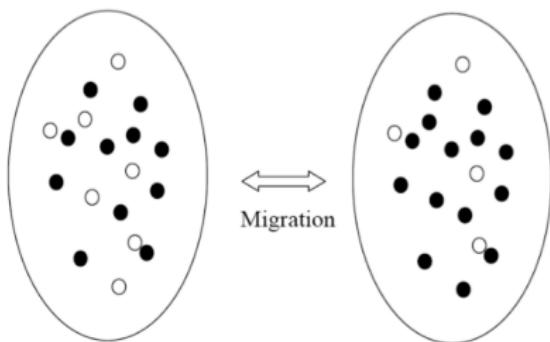
# Population structure

- The frequency of an allele (A) in subpopulation 1 is  $f_{A1}$ . The frequency of the same allele in subpopulation 2 is  $f_{A2}$ .
- The frequency of the allele in the whole population is therefore  $f_A = \frac{Nf_{A1} + Nf_{A2}}{2N} = \frac{f_{A1} + f_{A2}}{2}$ .



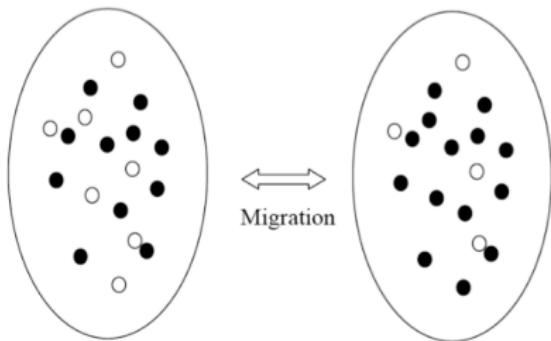
# Population structure

- The frequency of an allele (A) in subpopulation 1 is  $f_{A1}$ . The frequency of the same allele in subpopulation 2 is  $f_{A2}$ .
- The frequency of the allele in the whole population is therefore  $f_A = \frac{Nf_{A1} + Nf_{A2}}{2N} = \frac{f_{A1} + f_{A2}}{2}$ .
- What is the expected heterozygosity under HW?



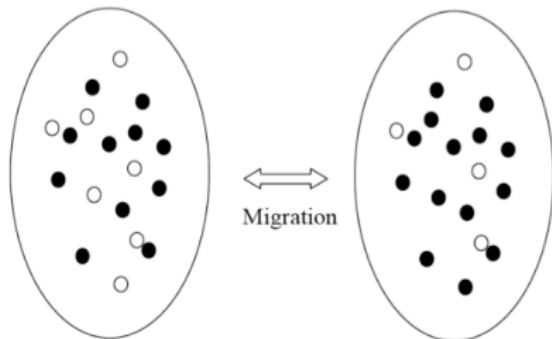
# Population structure

- The frequency of an allele (A) in subpopulation 1 is  $f_{A1}$ . The frequency of the same allele in subpopulation 2 is  $f_{A2}$ .
- The frequency of the allele in the whole population is therefore  $f_A = \frac{Nf_{A1} + Nf_{A2}}{2N} = \frac{f_{A1} + f_{A2}}{2}$ .
- What is the expected heterozygosity under HW?
- Depends...



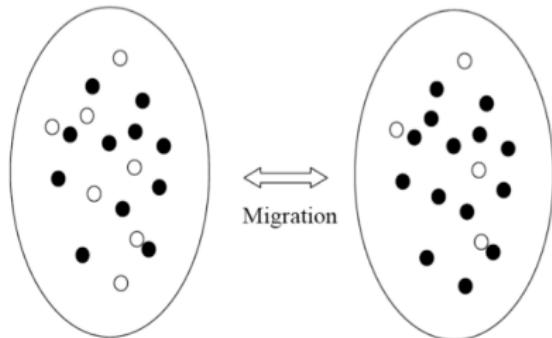
# Population structure

- Recall that for a single population under HW, the expected heterozygosity is  $2p(1-p)$ .



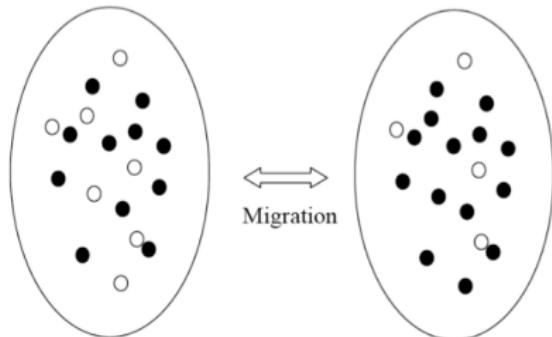
# Population structure

- Recall that for a single population under HW, the expected heterozygosity is  $2p(1-p)$ .
- In this case, the expected heterozygosity is the weighted sum of each subpopulation's heterozygosity...



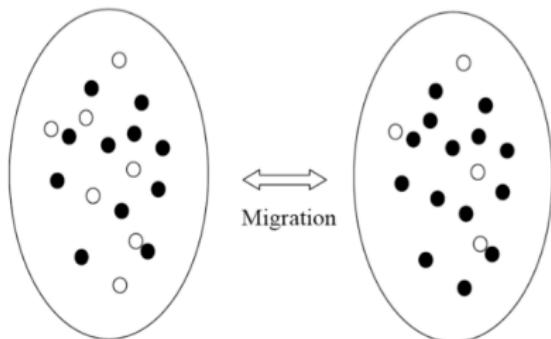
# Population structure

- Recall that for a single population under HW, the expected heterozygosity is  $2p(1-p)$ .
- In this case, the expected heterozygosity is the weighted sum of each subpopulation's heterozygosity...
- $H_S = \frac{H_1+H_2}{2}$



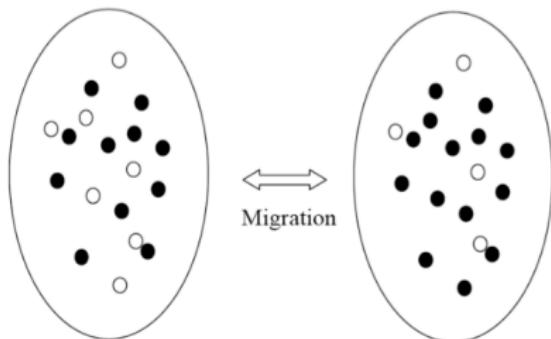
# Population structure

- Recall that for a single population under HW, the expected heterozygosity is  $2p(1-p)$ .
- In this case, the expected heterozygosity is the weighted sum of each subpopulation's heterozygosity...
- $H_S = \frac{H_1 + H_2}{2}$
- $H_1 = 2f_{A1}(1 - f_{A1})$  and  $H_2 = 2f_{A2}(1 - f_{A2})$



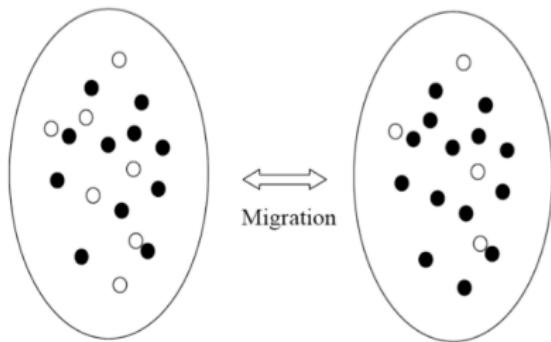
# Population structure

- Recall that for a single population under HW, the expected heterozygosity is  $2p(1-p)$ .
- In this case, the expected heterozygosity is the weighted sum of each subpopulation's heterozygosity...
- $H_S = \frac{H_1 + H_2}{2}$
- $H_1 = 2f_{A1}(1 - f_{A1})$  and  $H_2 = 2f_{A2}(1 - f_{A2})$
- Therefore,  $H_S = (2f_{A1}(1 - f_{A1}) + 2f_{A2}(1 - f_{A2}))/2$



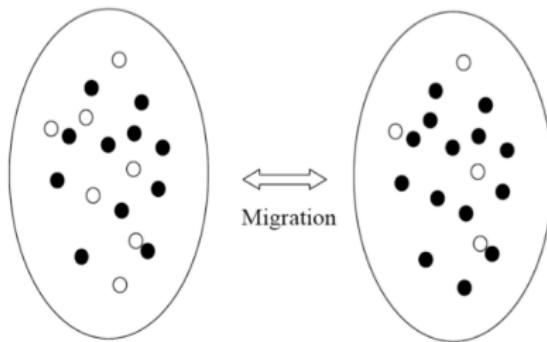
# Population structure

- But if the two subpopulations are fully mixing with each other, then all we would need to calculate the expected heterozygosity is the total frequency of the A allele:  $f_A$ .



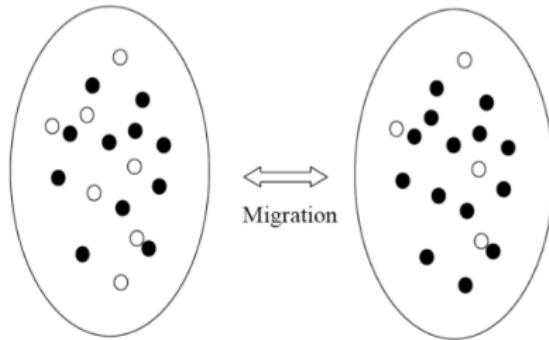
# Population structure

- But if the two subpopulations are fully mixing with each other, then all we would need to calculate the expected heterozygosity is the total frequency of the A allele:  $f_A$ .
- $H_T = 2f_A(1 - f_A) = 2\frac{f_{A1}+f_{A2}}{2}(1 - \frac{f_{A1}+f_{A2}}{2})$



# Wright's Fst

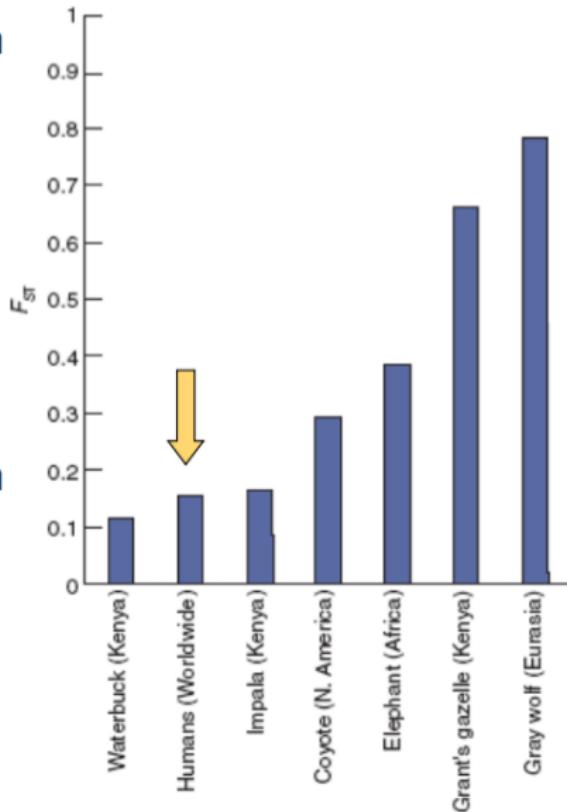
- The fixation index ( $F_{st}$ ) is just the difference between  $H_T$  and  $H_S$ , weighted by  $H_T$ .
- $$F_{st} = \frac{H_T - H_S}{H_T}$$
- Can range from 0 (no differentiation) to 1 (fixation of different alleles in each population).



# Wright's Fst

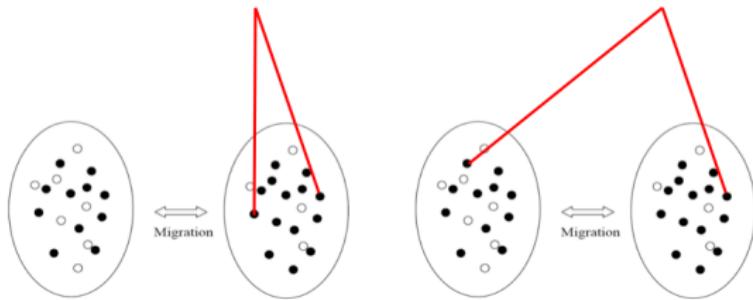
Most variation  
between  
populations

Most variation  
within  
populations



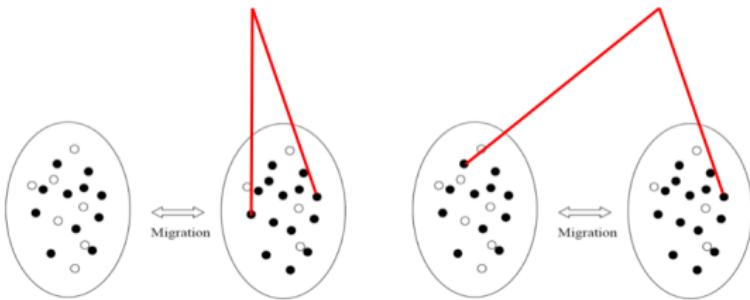
# Interpreting Fst

- $H_T$  is related to the expected time for any two samples to coalesce.
- $H_S$  is related to the expected time for two samples FROM THE SAME SUBPOPULATION to coalesce.



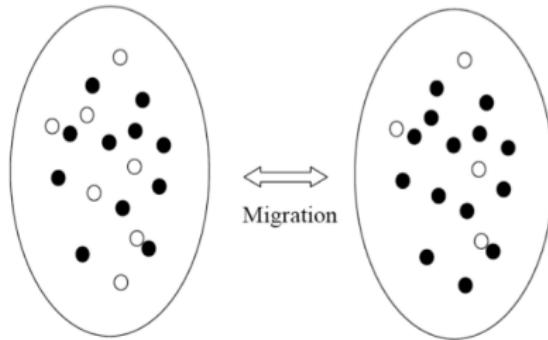
# Interpreting $F_{st}$

- If  $H_S < H_T$ , the expected time for coalescence for two samples **WITHIN THE SAME SUBPOPULATION** is smaller than one would expect if the two populations were randomly mixing, so  $F_{st} \approx 1$ .
- If  $H_S \approx H_T$ , the expected time for coalescence for two samples **WITHIN THE SAME SUBPOPULATION** is approximately what one would expect for any two samples (regardless of which subpopulation they are sampled from), so  $F_{st} \approx 0$ .



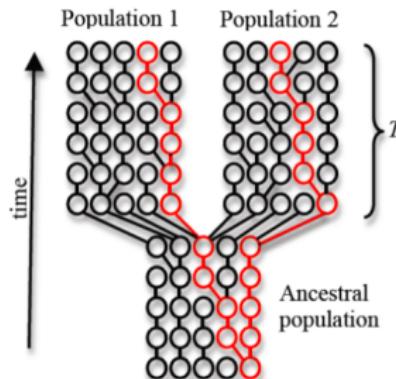
# Interpreting $F_{st}$

- While  $F_{st}$  measures population structure, there are different models of structure that can generate the same  $F_{st}$  values.
- Under the two-population-with-migration model that we've been going over, one can show (Slatkin 1991) that:
- $F_{st} = \frac{1}{1+8Nm}$  where  $m$  is the migration rate between the populations.

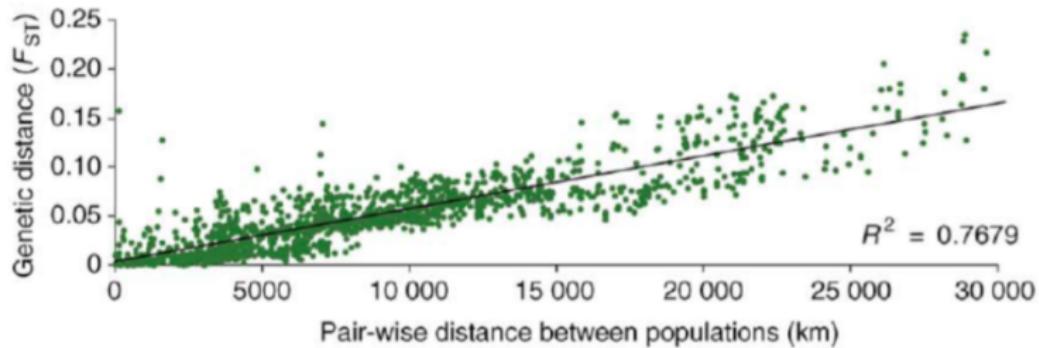


# Interpreting $F_{st}$

- Under a population-split model,  $F_{st}$  can be related to the amount of time since the split:
- $$F_{st} = \frac{T}{T+2}$$
- When  $T = 0$ , there is no population subdivision, and  $F_{st} = 0$ .
- When  $T$  is large,  $F_{st}$  approaches 1.



# Wright's Fst



*Figure 4.7. The relationship between  $FST$  and geographical distance between different pairs of human populations for more than 50 globally distributed populations. The figure is reproduced from Handley et al. 2007.*

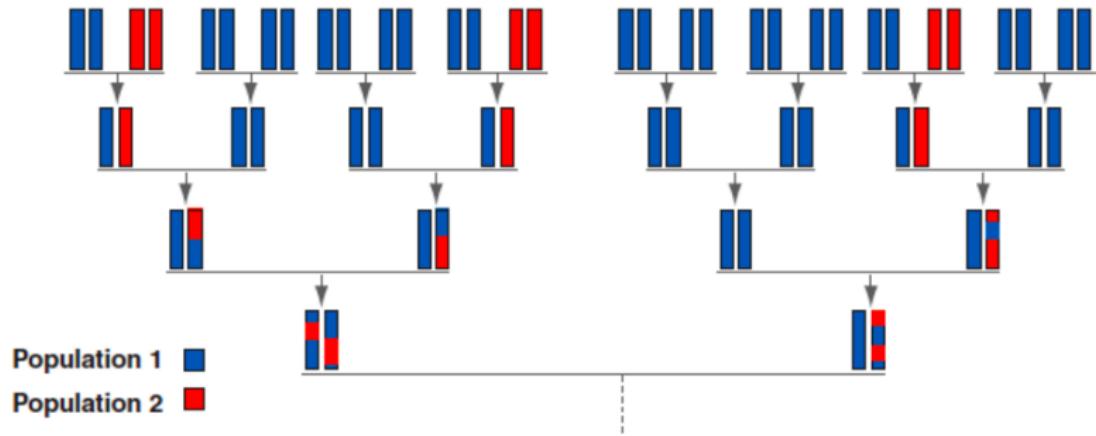
# Today

- $F_{ST}$  and population structure
- **The mixed-membership clustering model (“Structure”)**
- Leveraging haplotype information

## The “Structure” model

- The original model was first proposed by Pritchard et al. (2000)
- **Assumption 1:** each individual can be modeled as a mixture of one or more ancestral “**source populations**”
- **Assumption 2:** each locus is independent
- The proportion of genetic material from each source in each individual is called the “**admixture proportion**”
- **Problem 1:** we don’t know the identity and number of these source populations
- **Problem 2:** we don’t know the admixture proportions
- **Objective:** find best-fitting sources and their proportions

# The “Structure” model



# The “Structure” model

- Known: genotypes (G)
- Unknown:
  - admixture proportions (Q)
  - allele frequencies in source populations (F)
- Need to estimate Q and F, given that we know G.
- Objective: Maximize likelihood function:  $P[G|Q, F]$

**G** (genotypes):

Ind 1	A	T	G	T	T	A	A	T
	T	T	G	C	T	G	T	T
Ind 2	T	T	C	T	T	G	A	G
	T	G	C	T	A	G	T	T
...	T	G	G	T	T	G	A	G
	T	T	C	T	G	T	T	T
Ind 1	A	T	G	T	T	A	A	T
	T	T	G	T	T	G	T	T
Ind 2	T	T	G	T	T	G	A	G
	T	G	G	T	T	G	T	G

**Q** (admixture proportions):

Ind 1	4/16	12/16
Ind 2	6/16	10/16
...	3/16	13/16
	11/16	5/16
	3/16	13/16

**F** (allele frequencies):

Pop 1	2/5	3/4	4/4	3/3	3/3	0/3	2/3	0/2
Pop 2	0/5	4/6	3/6	5/7	5/7	2/7	3/7	6/8

<sup>0</sup>Ida Moltke pers. comm.

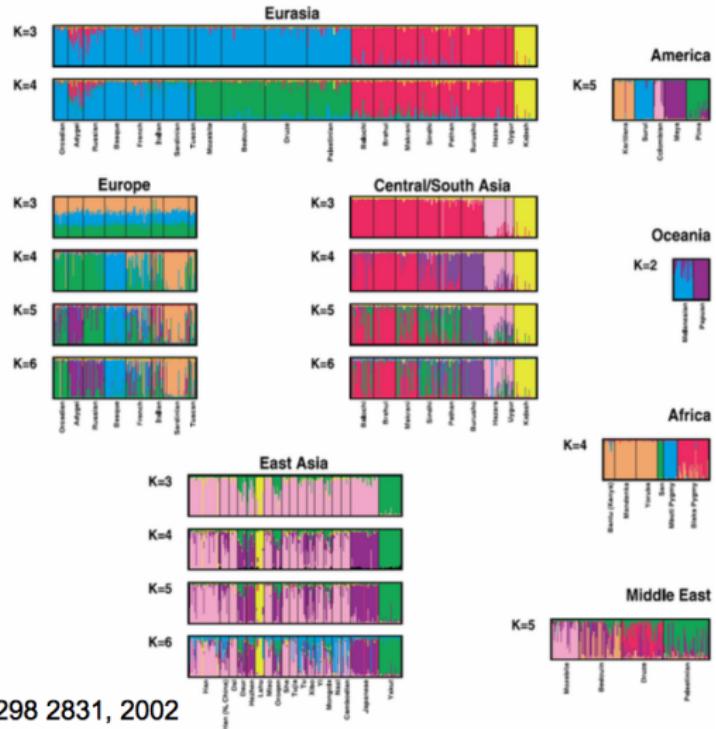
## Likelihood function: one individual, one locus

- A locus  $j$  can have two possible alleles:  $A$  and  $a$
- We'll code a genotype as the number of alleles  $a$ : 0, 1 or 2.
- Let  $G_{i,j}$  be the (diploid) genotype of locus  $j$  in individual  $i$ .
- Let  $F^j = (f^{j,1}, f^{j,2}, \dots, f^{j,K})$  be the allele frequencies of allele  $a$  in the  $K$  "source populations"
- Let  $Q^i = (q^{i,1}, q^{i,2}, \dots, q^{i,K})$  be the  $K$  admixture proportions of individual  $i$ .
- Then, for one of the allele copies in locus  $j$  in individual  $i$ :  
$$P[b|Q, F] = q^{i,1}f^{j,1} + q^{i,2}f^{j,2} + \dots + q^{i,K}f^{j,K} = h^{ij}$$
- Assuming Hardy-Weinberg equilibrium:  
$$P[G_{ij} = 2|Q^i, F^j] = (h^{ij})^2$$
$$P[G_{ij} = 1|Q^i, F^j] = 2(h^{ij})(1 - h^{ij})$$
$$P[G_{ij} = 0|Q^i, F^j] = (1 - h^{ij})^2$$

## Likelihood function: N individuals, M loci

- Assuming loci are independent and individuals are unrelated...
- $P[G|Q, F] = \prod_i^N \prod_j^M P[G_{ij}|Q^i, F_j]$
- $Q$  is a matrix of admixture proportions for each of the  $N$  individuals:  
 $Q = (Q^1, Q^2, \dots, Q^N)$
- $F$  is a matrix of ancestral allele frequencies for each of the  $M$  loci:  
 $F = (F^1, F^2, \dots, F^M)$
- Structure-like methods try to find the **parameters  $Q$  and  $F$  that maximize  $P[G|Q, F]$**

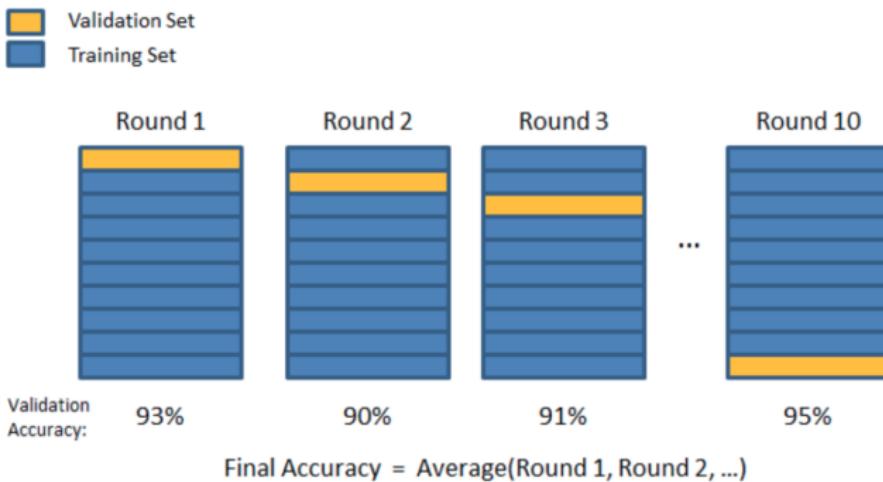
# Structure model applied to human populations



Science 298 2831, 2002

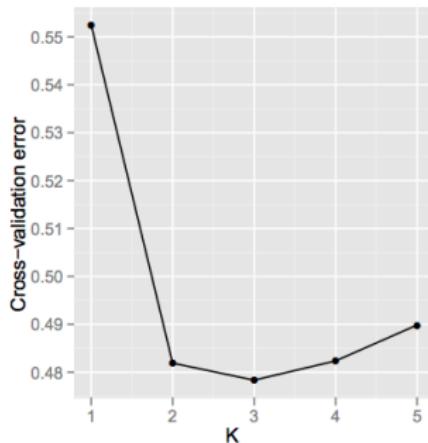
# Choosing K

- We can use cross-validation to find a value of K that does not over-fit
- We leave some genotypes out and predict them based on their estimated ancestries
- **Important:** well-fitting parameter  $\neq$  biologically meaningful parameter



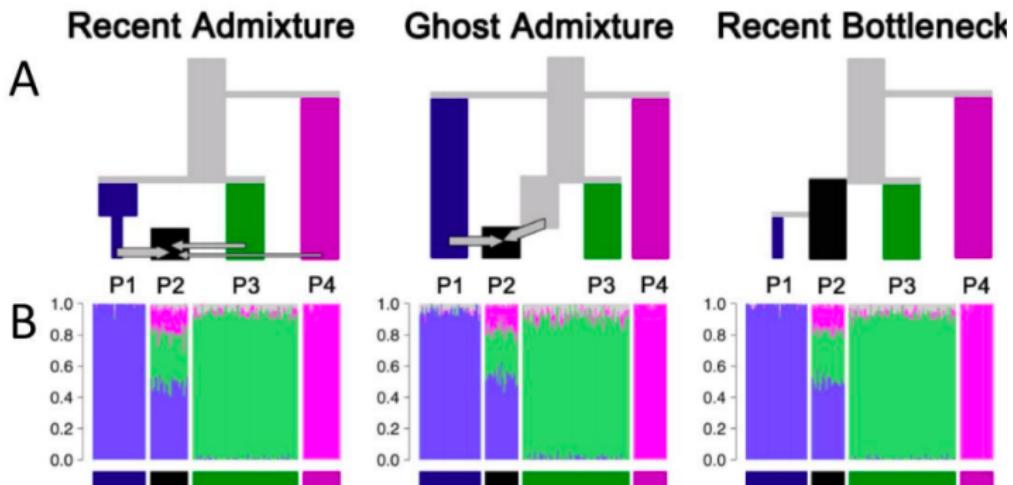
# Choosing K

- We can use cross-validation to find a value of K that does not over-fit
- We leave some genotypes out and predict them based on their estimated ancestries
- **Important:** well-fitting parameter  $\neq$  biologically meaningful parameter



# Over-interpreting Structure results

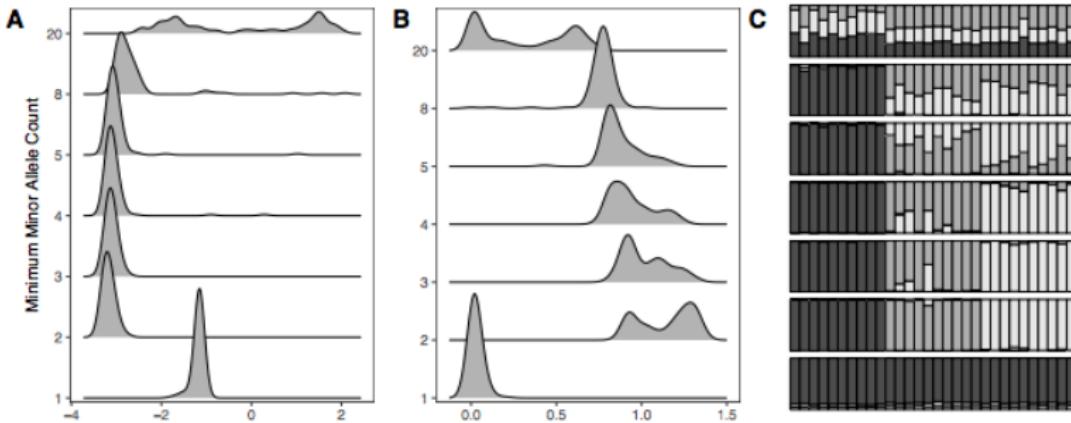
- Structure does not necessarily pick up admixture events!<sup>1</sup>
- “Source populations” need not be real populations that ever existed!
- A population that is highly drifted will be assigned its own cluster at high enough K



<sup>1</sup>Falush et al. 2016

# Minor allele frequency cutoff

- MAF cutoffs can strongly inference under the Structure model<sup>2</sup>

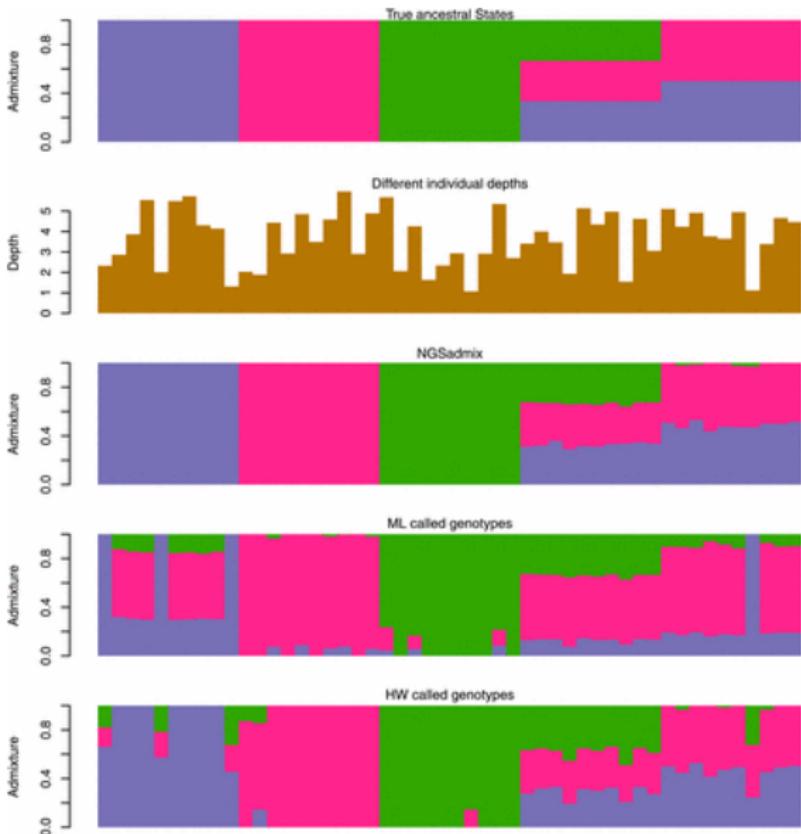


<sup>2</sup>Linck and Battey 2017

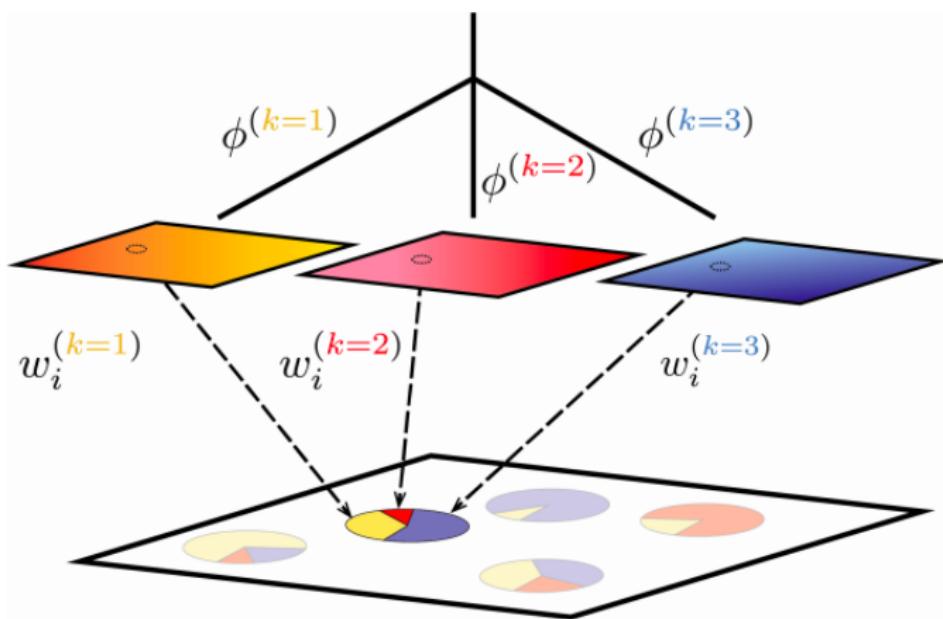
## Variations on a theme...

- Structure (Pritchard et al. 2000): original model; uses Bayesian priors to obtain posterior estimates of  $Q$  and  $F$
- Admixture (Alexander et al. 2011): faster than Structure; uses a maximum likelihood model rather than a Bayesian model; uses cross-validation to choose  $K$
- fastStructure (Raj et al. 2014): faster than Structure; uses variational inference to choose  $K$ ; can detect weak structure
- ngsAdmix (Skotte et al. 2013): can work with genotype likelihoods; better for low coverage data
- Ohana (Cheng et al. 2016): uses Gaussian approximation to model drift in each ancestry component; can detect selection by testing for local deviations from genome-wide model

# ngsAdmix (Skotte et al. 2013)



consStruct: structure + isolation-by-distance (Bradburd et al. 2017)

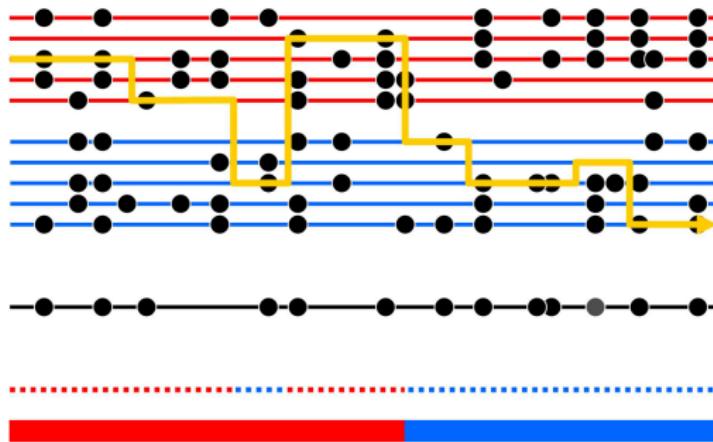


## SNP vs. Haplotype information

- All these methods ignore the spatial distribution of SNPs along the genome
- They require LD pruning: a lot of SNPs sit in the same haplotype and have redundant information
- Advantage: can model each SNP independently (simple model)
- Disadvantage: we ignore haplotype information

# Chromosome painting

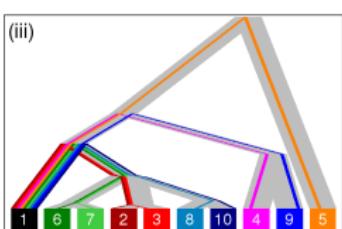
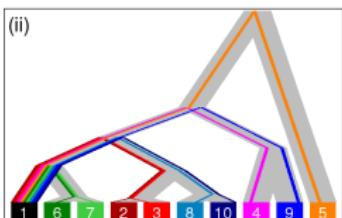
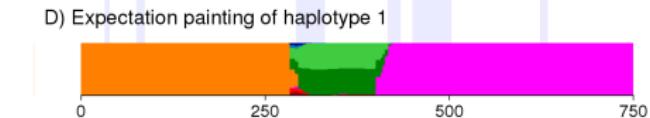
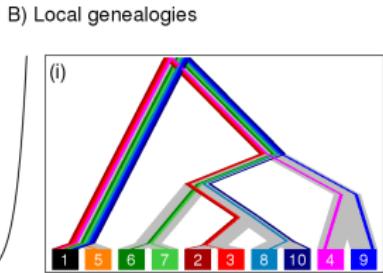
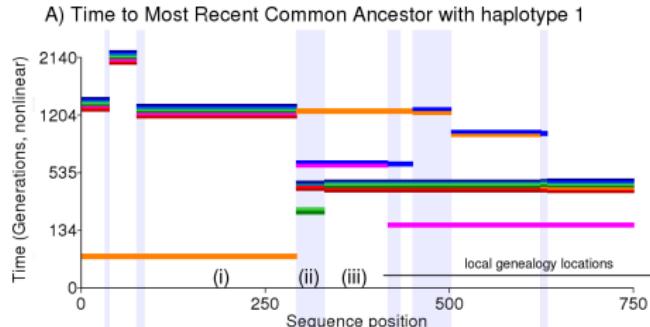
- Painting a chromosome by “copying” segments from other individuals<sup>3</sup>
- Can be done very fast, thanks to Li and Stephens algorithm <sup>4</sup>



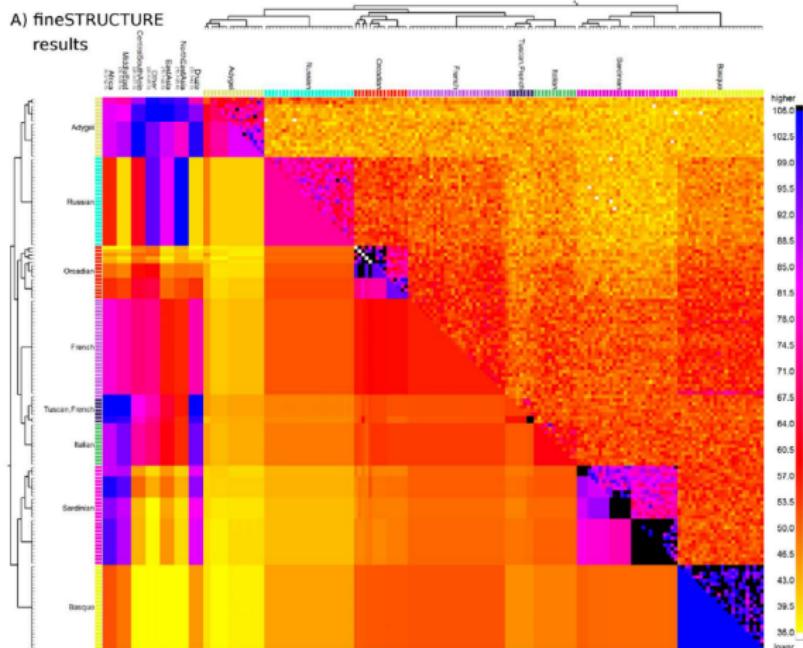
<sup>4</sup>Price et al. 2009

<sup>4</sup>Li and Stephens 2003

# Chromosome painting

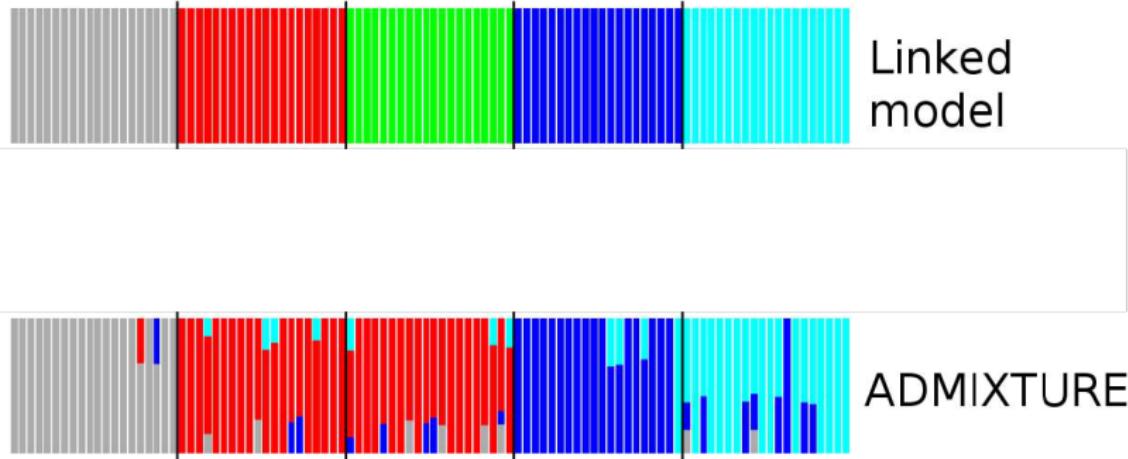


# Coancestry matrix based on painting patterns



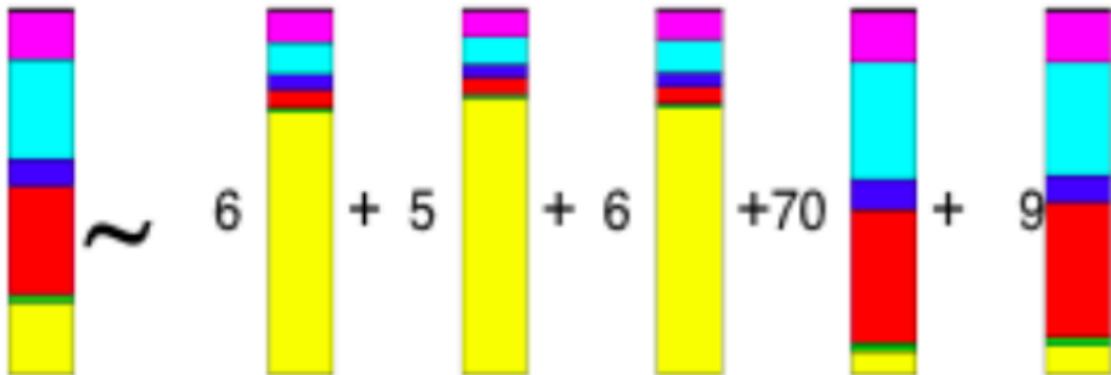
## Chromosome “palettes”

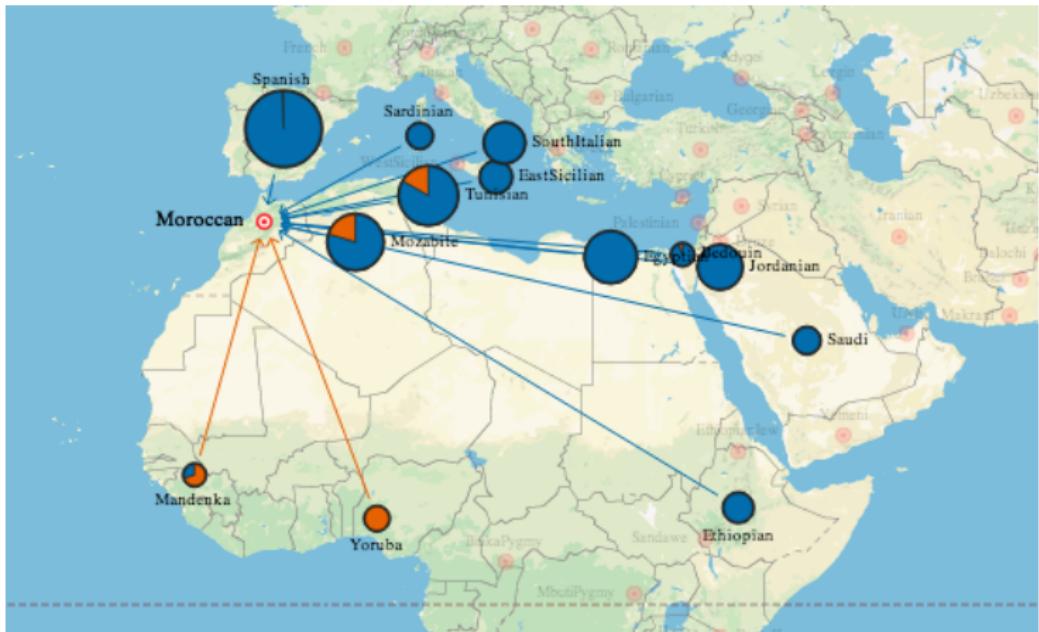
- We can create a chromosome palette by aggregating the ancestry from each of the chromosome segments
- This is comparable (but not the same!) as a Structure / Admixture barplot



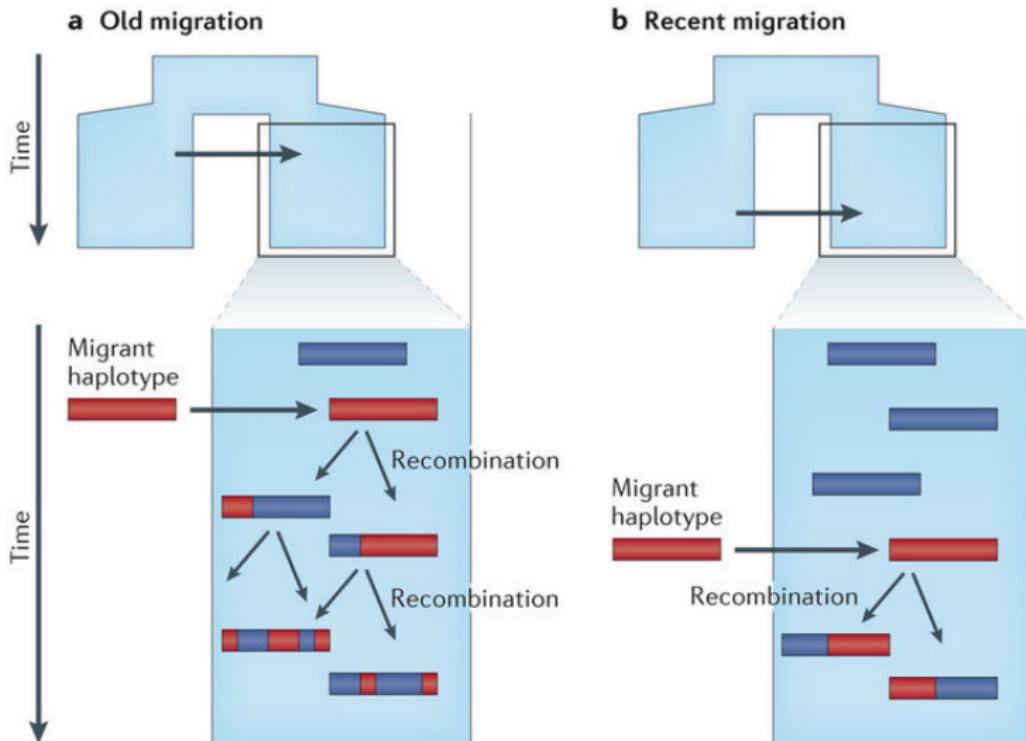
## Chromosome “palettes”

- We can model an individual as a linear mixture of the palettes from other individuals





# Admixture date inference



## Admixture date inference

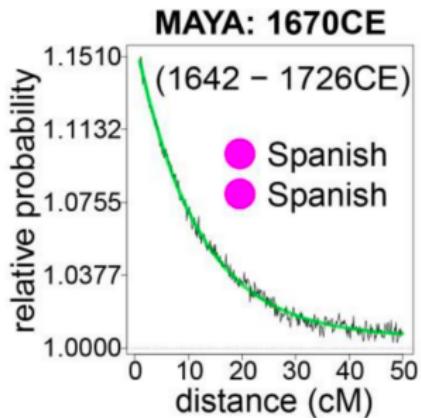
- Assume an admixture event happened at time  $\lambda$  between two populations A and B
- We're interested in modeling the length of tracts from A in an admixed genome
- Assume that the probability of no recombination between two points a distance  $g$  apart since admixture is  $e^{-g\lambda}$  (where  $\lambda$  is the time since admixture scaled by the recombination rate)
- If the fraction of total ancestry from population A in the genome is  $\alpha$ , then the probability that we'll find two loci with ancestry from A a distance  $g$  apart is:
  - $p_{AA}(g) = \alpha(e^{-g\lambda} + (1 - e^{-g\lambda})\alpha) = \alpha^2 + \alpha(1 - \alpha)e^{-g\lambda}$
  - This is just an exponential function of the admixture time  $\lambda!$ <sup>5</sup>



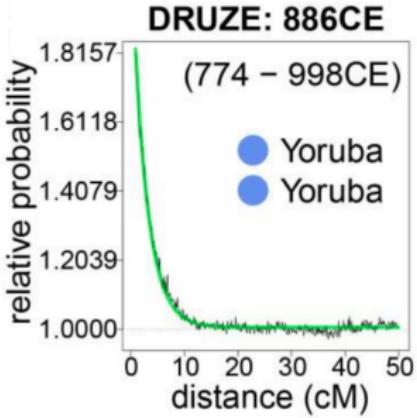
<sup>5</sup>Hellenthal et al. 2014

# Admixture date inference

p\_Spanish,Spanish as a function of distance in a Mayan genome



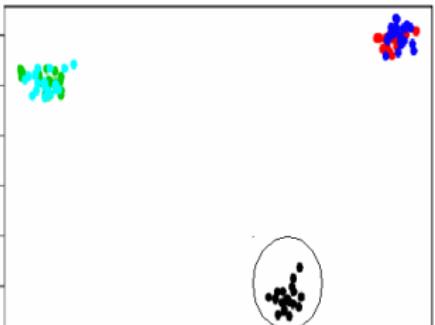
p\_Yoruba,Yoruba as a function of distance in a Druze genome



Haplotype data has more information than SNP data

## Chromo Painter PCA

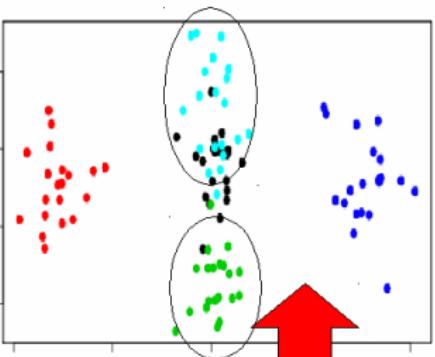
Component 2



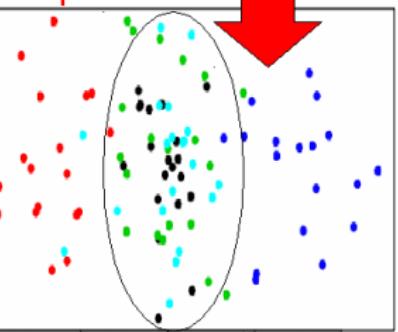
More signal,  
less noise



Component 4



Identify  
Populations



Component 3

## Basic PCA

Component 1

# Haplotype models can capture very recent history

