

Quantifying population structure: beyond PCA

PCA	Dimension-reduction (no population model)	What are the major axes of genetic variation?
STRUCTURE MODEL (STRUCTURE, Admixture, Entropy)	Model-based ancestry inference	What proportion of each individual's genome comes from K ancestral populations?
SNMF	Sparse non-negative matrix factorization	Can we infer ancestry proportions like ADMIXTURE, but faster and with fewer assumptions?

F_{ST} : Quantifying Pop Subdivision/ genetic differentiation

$$F_{ST} = \frac{H_T - H_s}{H_T}$$

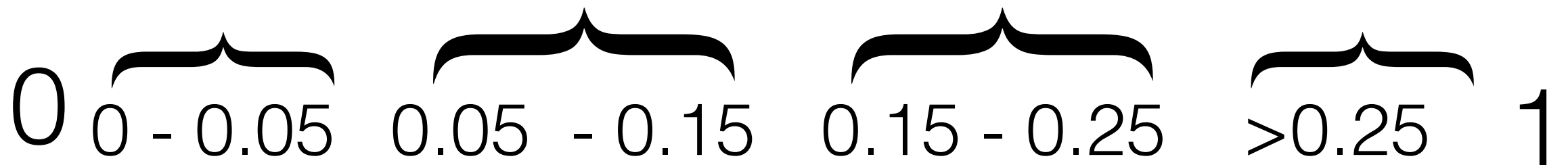
$$F_{ST} = \frac{\text{variation among populations}}{\text{total variation}}$$

little to no
structure

weak
structure

moderate
structure

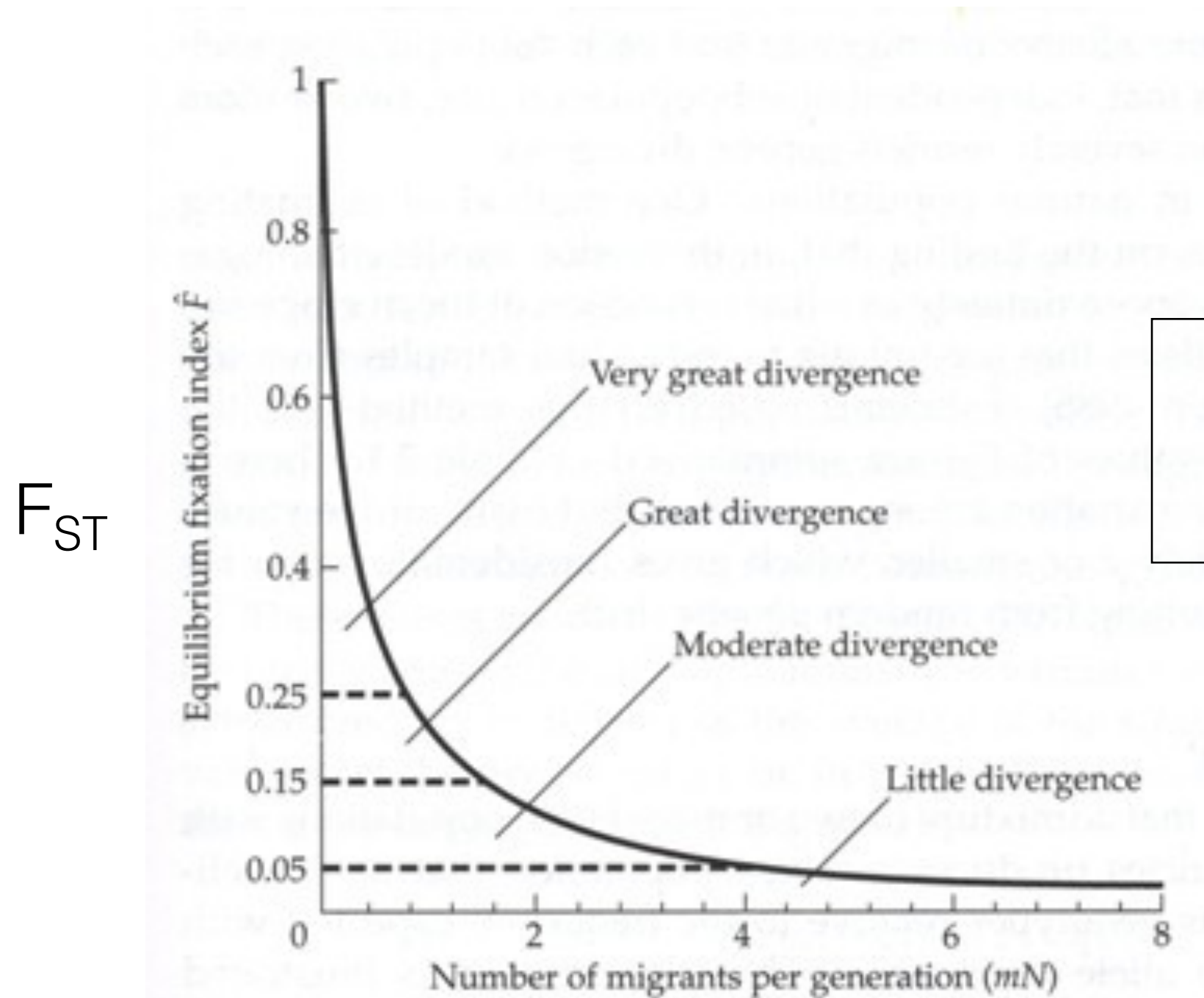
very great
structure



allele frequencies the
same in all populations

different alleles fixed in
different populations

Migration has a huge effect on F_{ST}

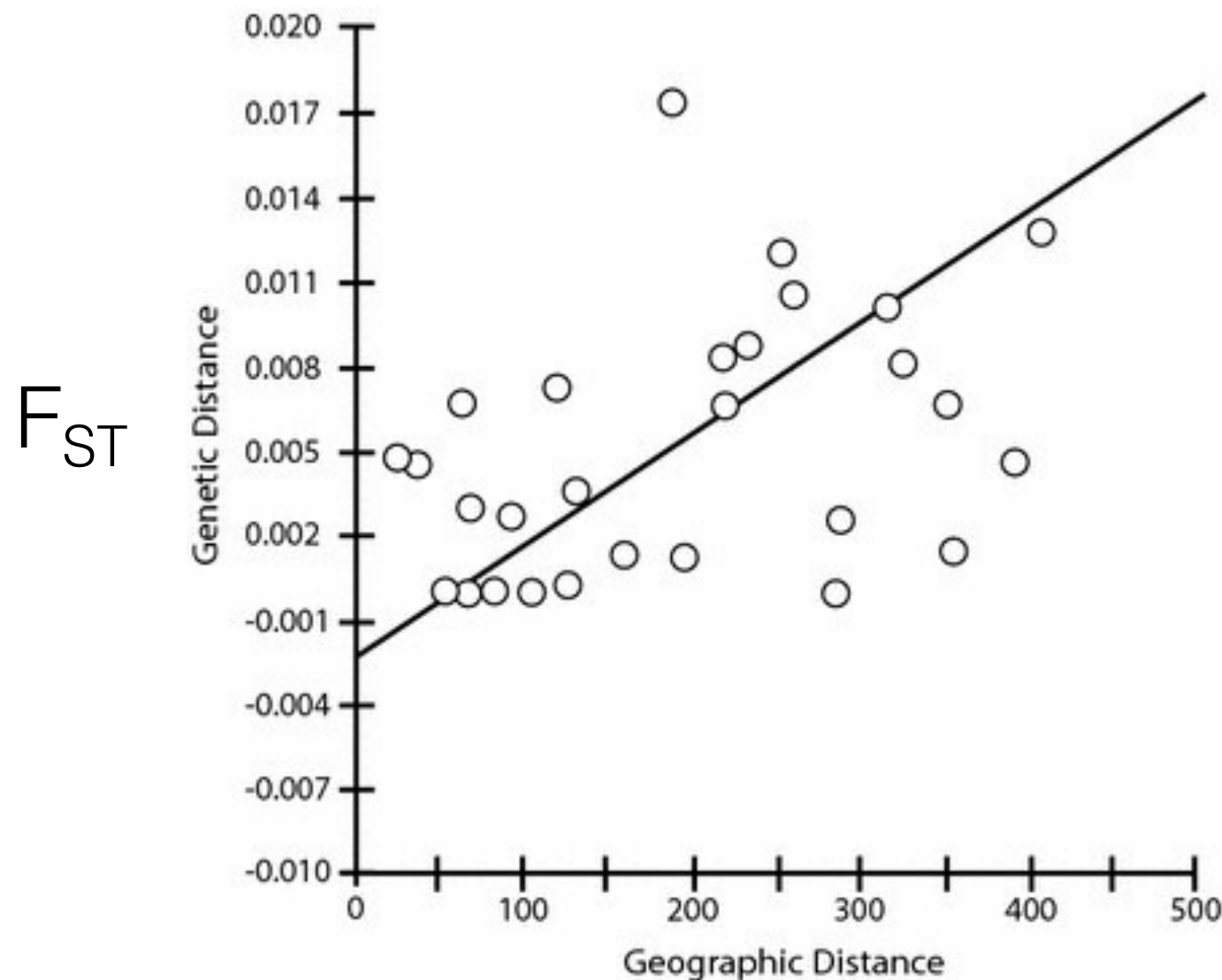


> 1-2 migrants per generation = NO differentiation!

This result extends to many other models

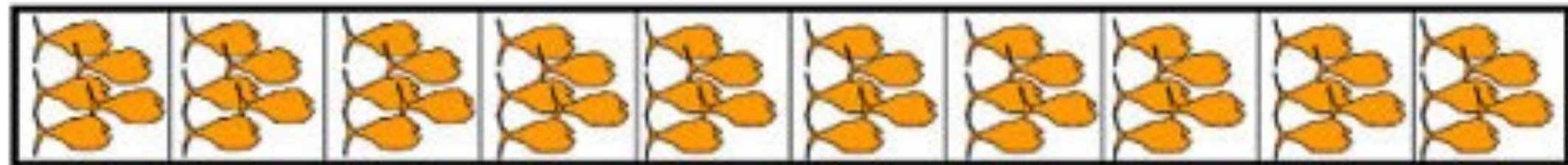
Isolation by Distance

The degree of population subdivision often increases with increasing geographic distance.



Population structure and isolation by distance

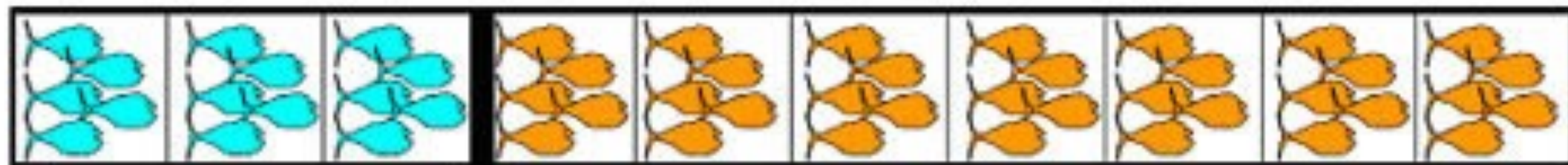
One single population



Unlimited migration

Fishing port

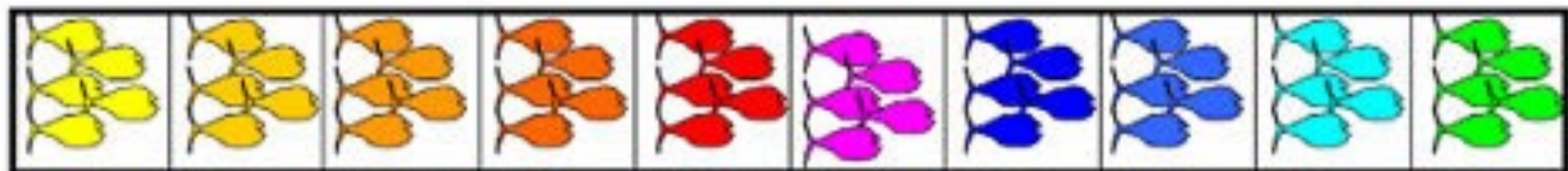
Two distinct populations



Limited migration

Fishing port

Isolation-by-distance (continuous change)



Limited migration between adjacent areas

Fishing port

“One very important aspect of all the F_{ST} -like statistics described above is that they are strongly influenced by within-subpopulation levels of variation (Charlesworth 1998; Jakobsson, Edge, and Rosenberg 2013). Because of this, we refer to them as relative measures of differentiation. In contrast, absolute measures of population differentiation are mostly independent of levels of within-population diversity; absolute measures are also known as genetic distances.”

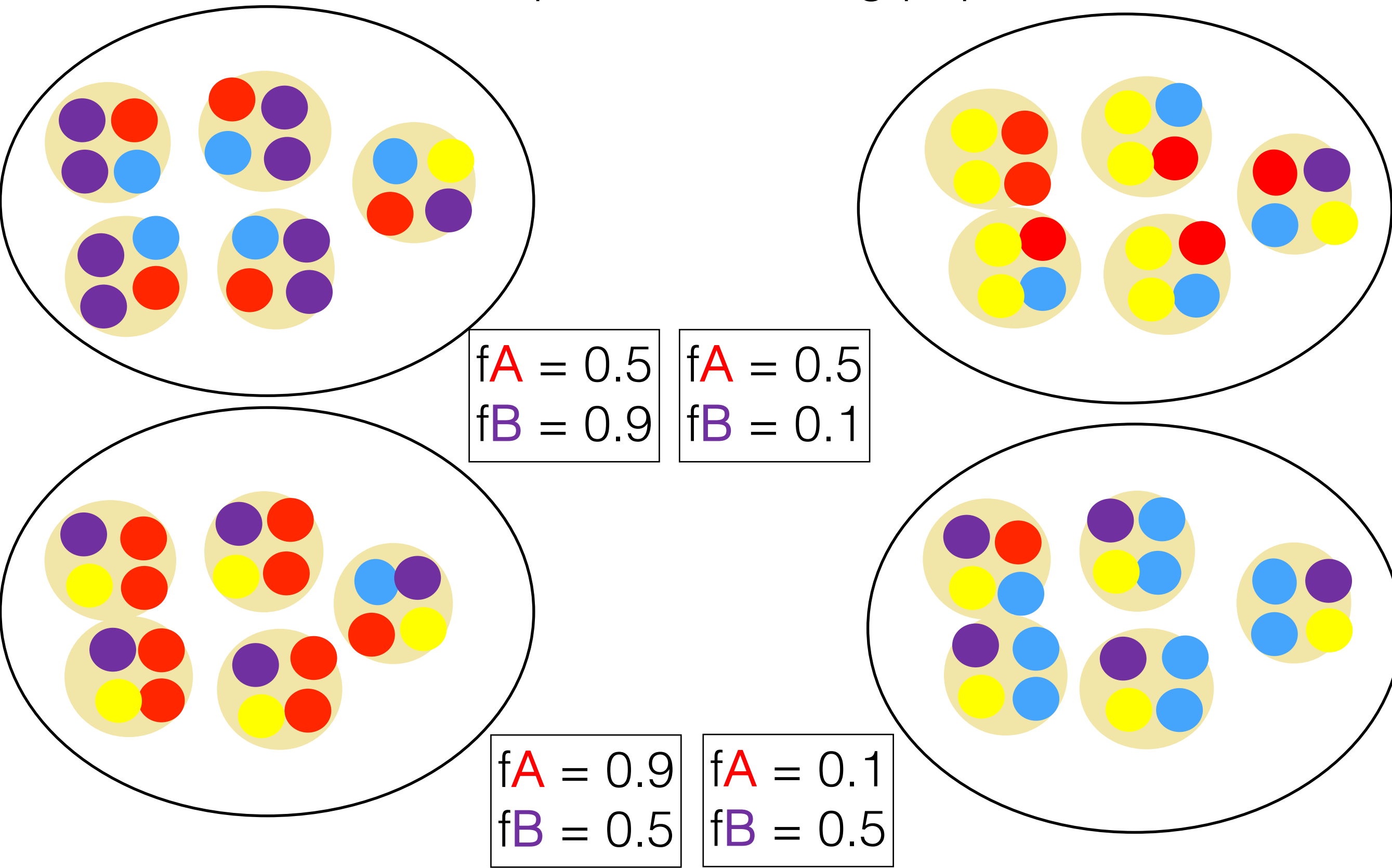
There are several different estimators of F_{st} and these are implemented in R in various packages and scripts.

Think about *how* calculating F_{st} makes sense in your situation — do you want genetic distances between populations? Between individuals? Across loci (as a means of detecting selection)?

When accurately defining populations, we expect to maximize Hardy-Weinburg equilibrium

“...one way to find the best assignment of individuals to populations—or to find the most likely number of populations—is to attempt to minimize the amount of Hardy-Weinberg disequilibrium...The optimal assignment is then the configuration that results in the least amount of Hardy-Weinberg disequilibrium in each population.”

We can calculate the probability that an individual comes from a particular population based on its genotype, using differences in allele frequencies among populations



Inference of Population Structure Using Multilocus Genotype Data

Jonathan K. Pritchard, Matthew Stephens and Peter Donnelly

Department of Statistics, University of Oxford, Oxford OX1 3TG, United Kingdom

Manuscript received September 23, 1999

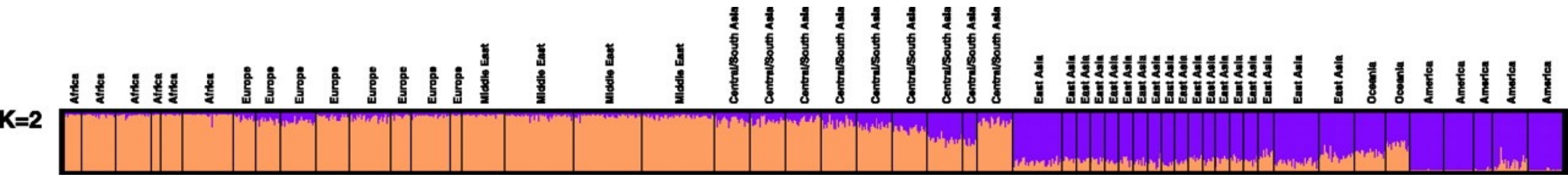
Accepted for publication February 18, 2000

ABSTRACT

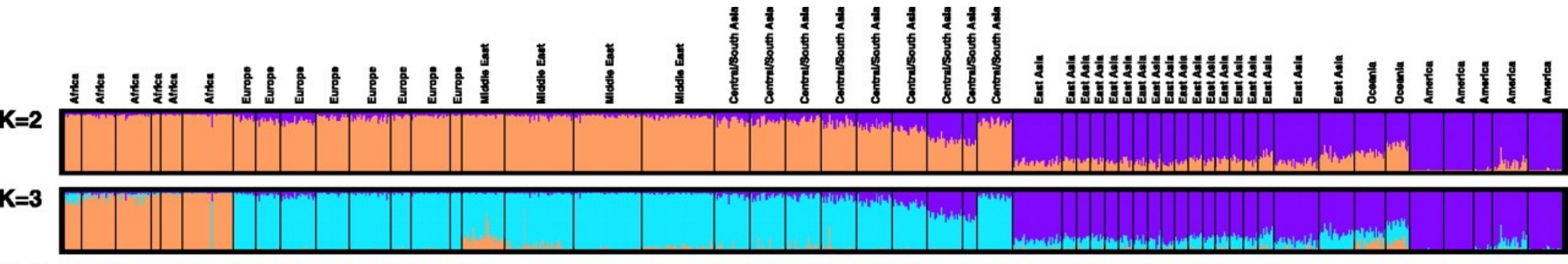
We describe a model-based clustering method for using multilocus genotype data to infer population structure and assign individuals to populations. We assume a model in which there are K populations (where K may be unknown), each of which is characterized by a set of allele frequencies at each locus. Individuals in the sample are assigned (probabilistically) to populations, or jointly to two or more populations if their genotypes indicate that they are admixed. Our model does not assume a particular mutation process, and it can be applied to most of the commonly used genetic markers, provided that they are not closely linked. Applications of our method include demonstrating the presence of population structure, assigning individuals to populations, studying hybrid zones, and identifying migrants and admixed individuals. We show that the method can produce highly accurate assignments using modest numbers of loci—*e.g.*, seven microsatellite loci in an example using genotype data from an endangered bird species. The software used for this article is available from <http://www.stats.ox.ac.uk/~pritch/home.html>.

Prichard 2000 — cited 40,000+ times!

STRUCTURE: proportion of an individual's ancestry that comes from each of k populations

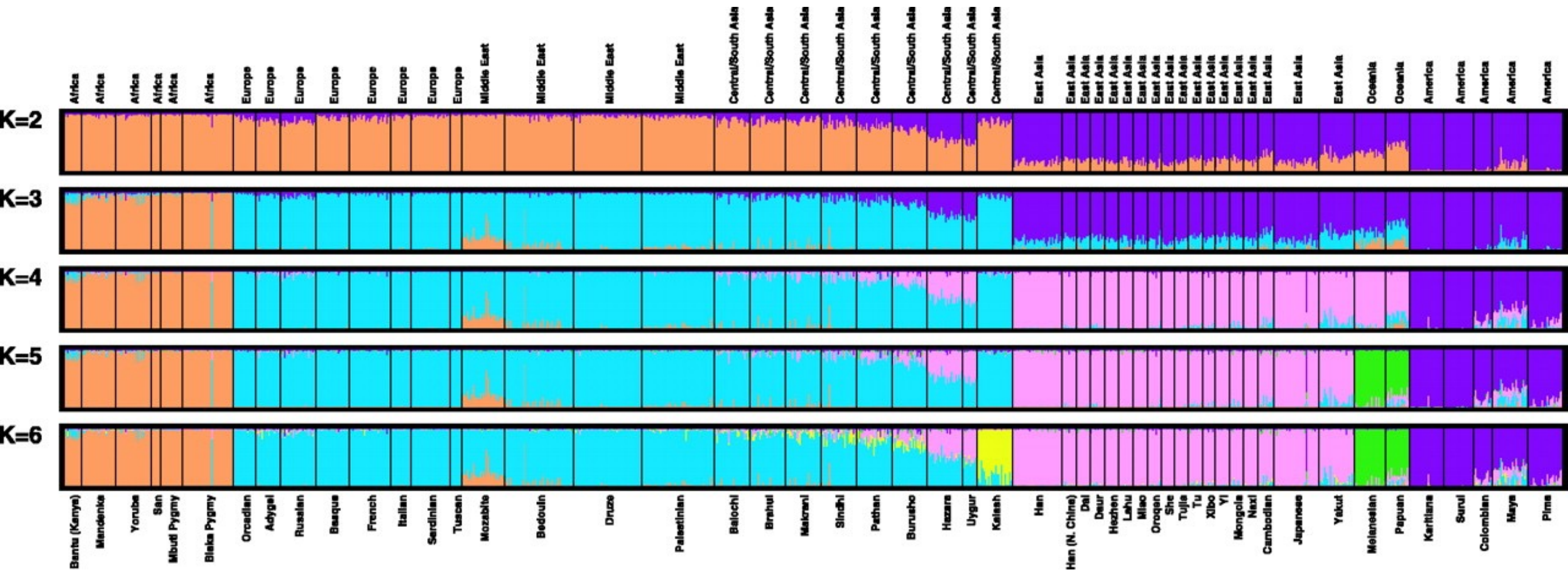


STRUCTURE: proportion of an individual's ancestry that comes from each of k populations



Increasing the value of K can reveal finer-scale population structure

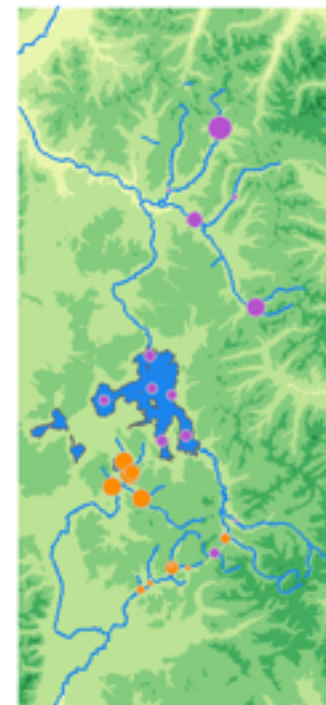
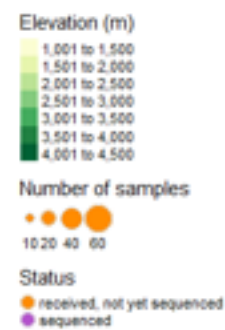
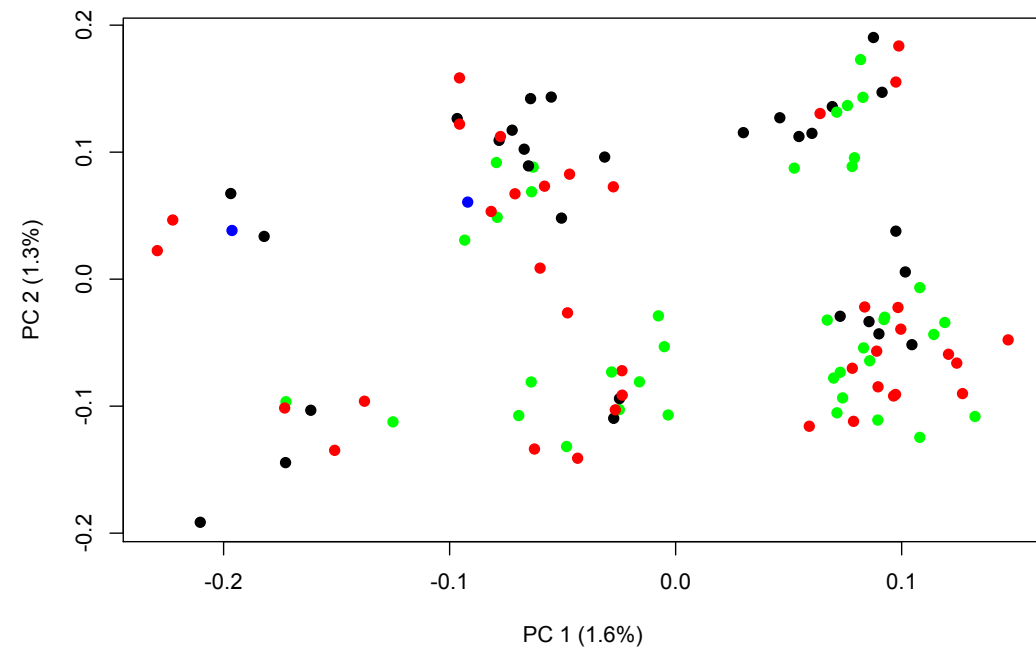
STRUCTURE: proportion of an individual's ancestry that comes from each of k populations



Increasing the value of K can reveal finer-scale population structure

PCA is also useful!

“...Patterson, Price, and Reich (2006) showed that the number of significant eigenvalues in a PCA study is equal to $K - 1$.”



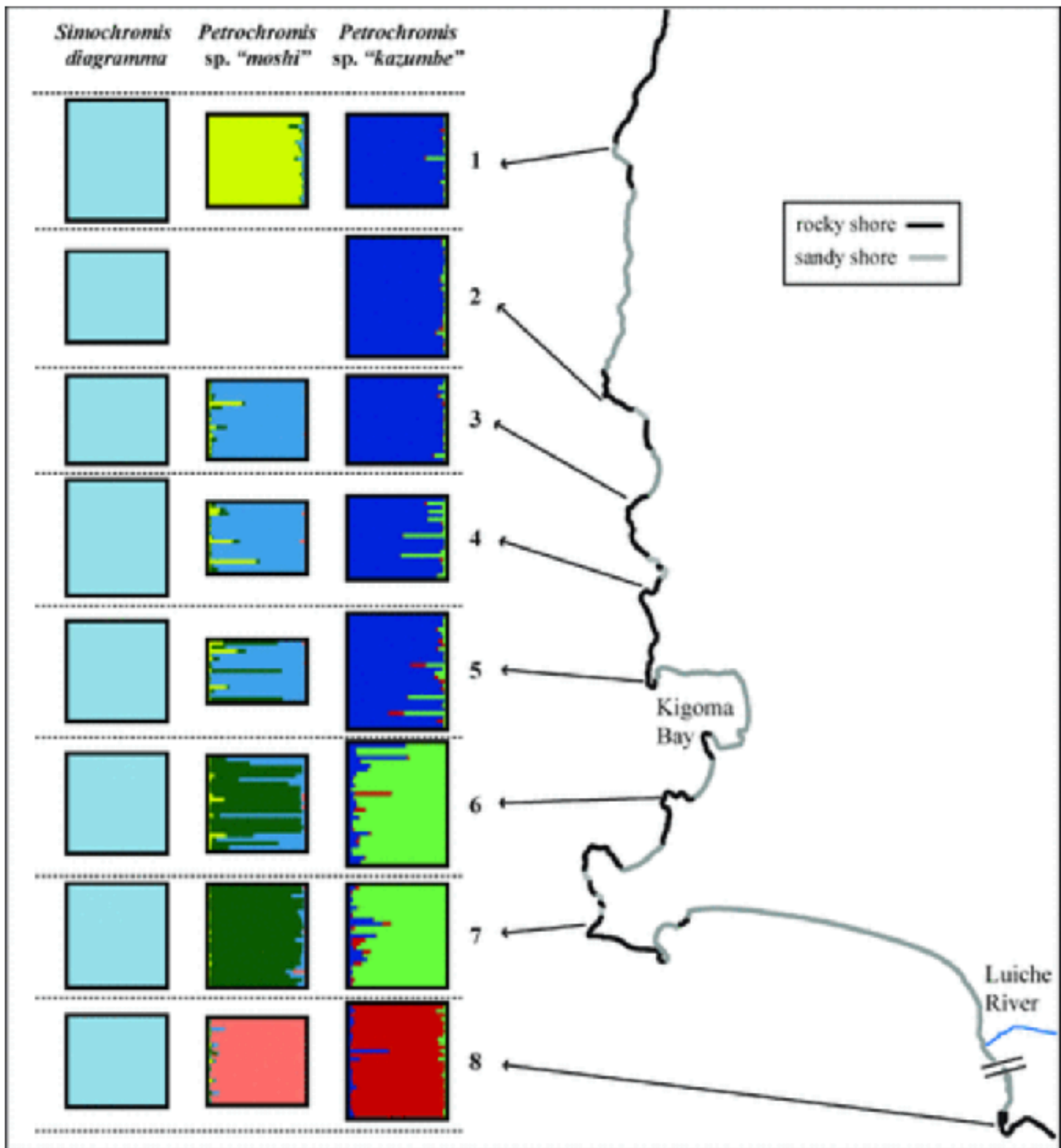
Variations on the STRUCTURE model:

Admixture: Novembre, and K. Lange, 2009 Fast model-based estimation of ancestry in unrelated individuals. Genome research 19(9): 1655–1664

accommodates large genomic datasets

Entropy: Gompert et al 2014, Shastry et al 2021

incorporates genotype likelihoods



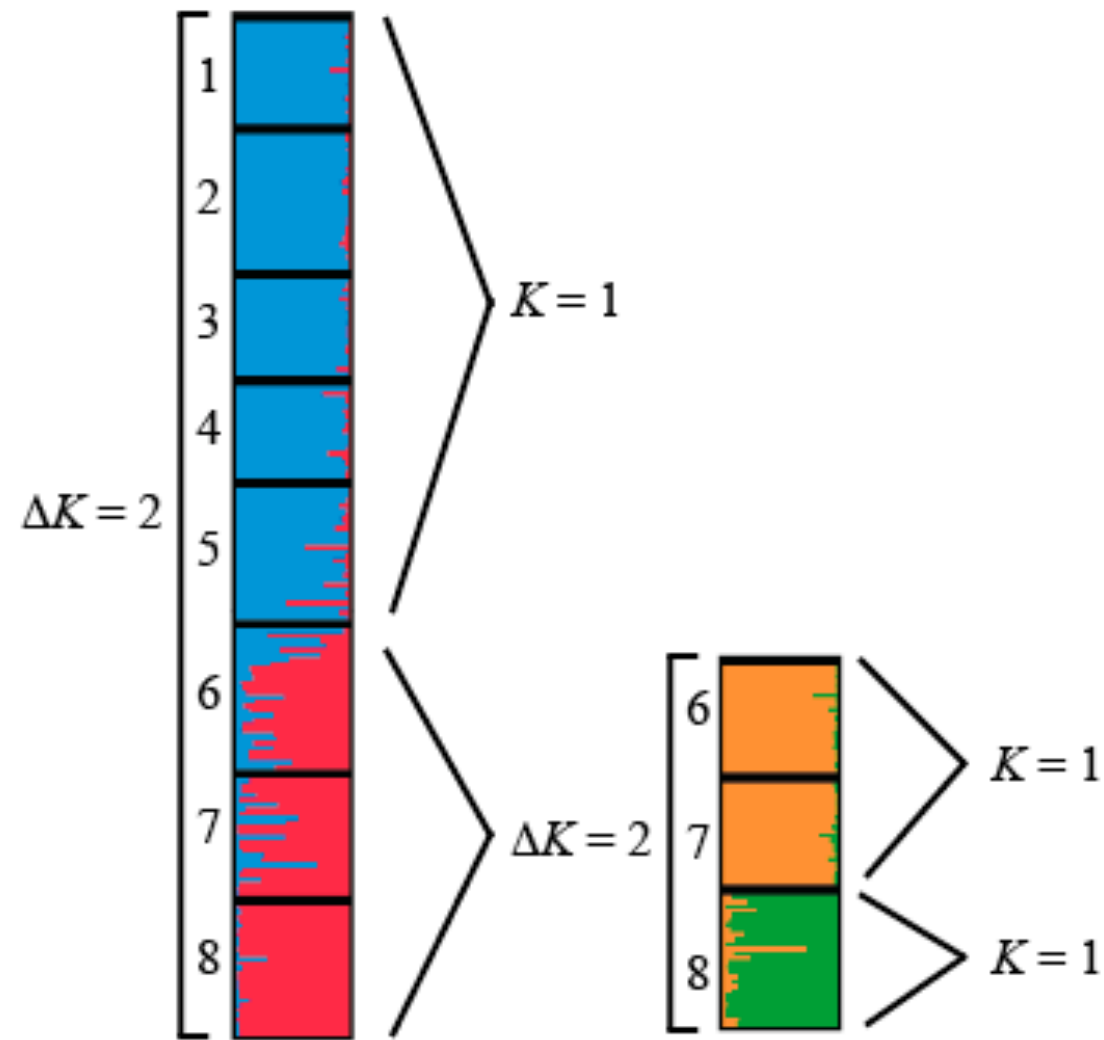
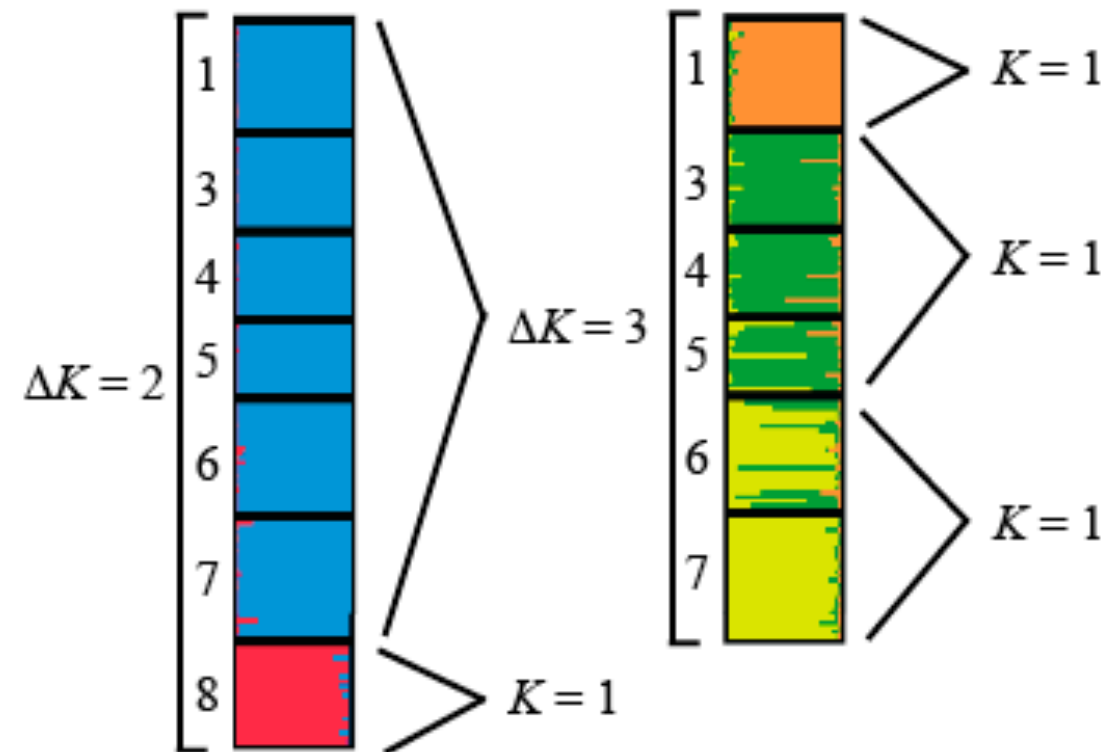
a) *P. kazumbe*b) *P. moshi*

Figure S4. Results of fully hierarchical ΔK STRUCTURE analyses for a) *P. sp.* “kazumbe” and b) *P. sp.* “moshi”. At each round, datasets were divided into K subsets, where K is the number of genetic groups supported by ΔK analyses (Evanno et al 2005) of the dataset from the previous round. Individuals were assigned into subsets if their assignment probabilities were 0.6 or higher for a given group. Assignment probabilities used for subsetting were the consensus of 10 runs at the K indicated by ΔK analyses, generated using CLUMPP (Jakobsson and Rosenberg 2007).

In STRUCTURE, LnP(D) is an approximation of the marginal likelihood

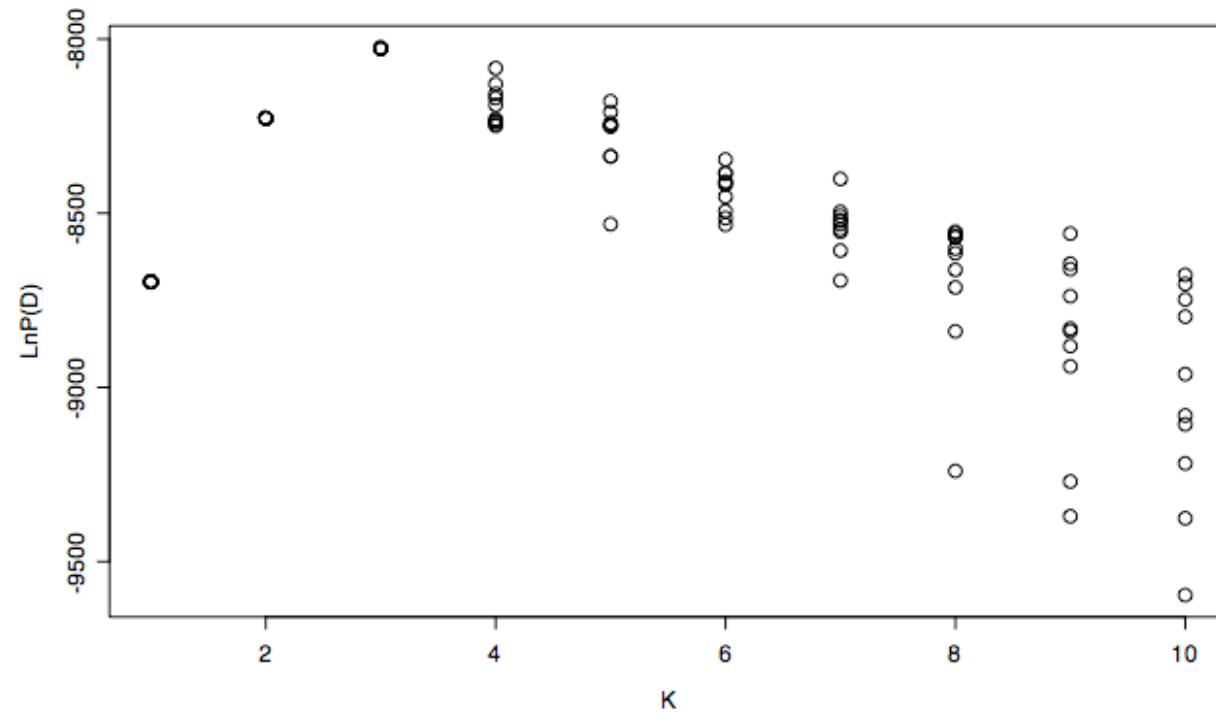


Figure S1. Log probability of the data (LnP(D)) for 10 STRUCTURE runs at each K value, for K 1 through 10 for *P. sp. "kazumbe"*.

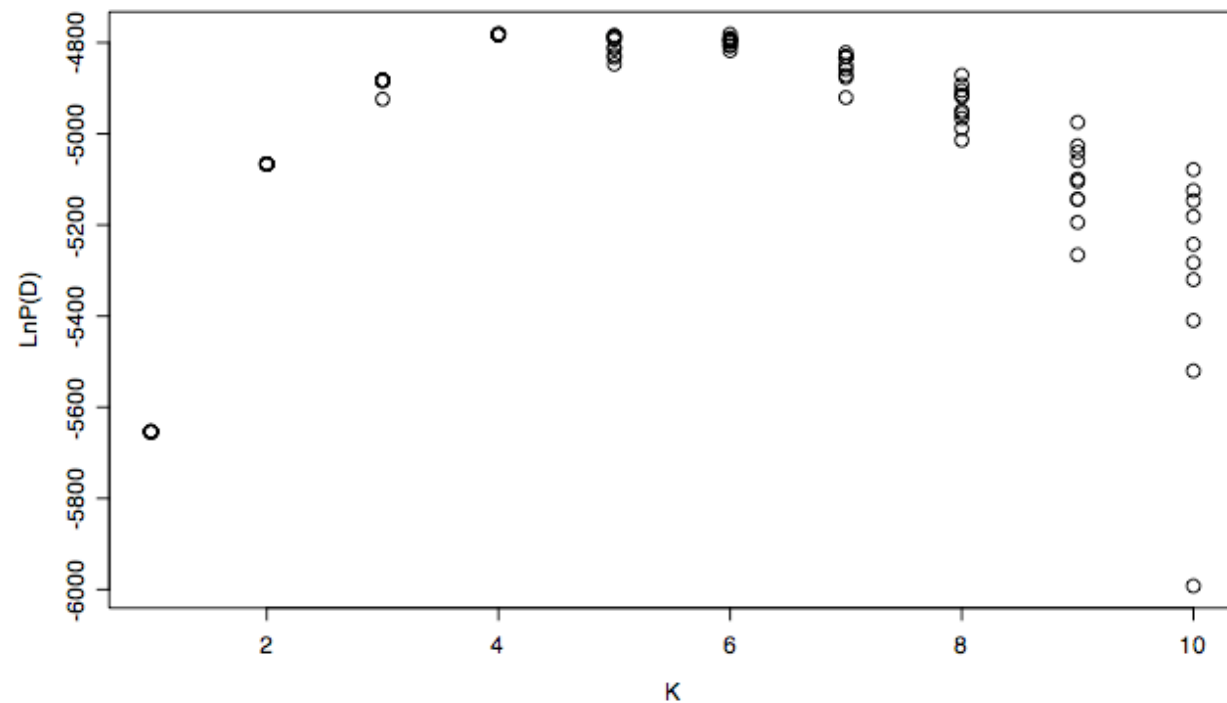


Figure S2. Log probability of the data (LnP(D)) for 10 STRUCTURE runs at each K value, for K 1 through 10 for *P. sp. "moshi"*.

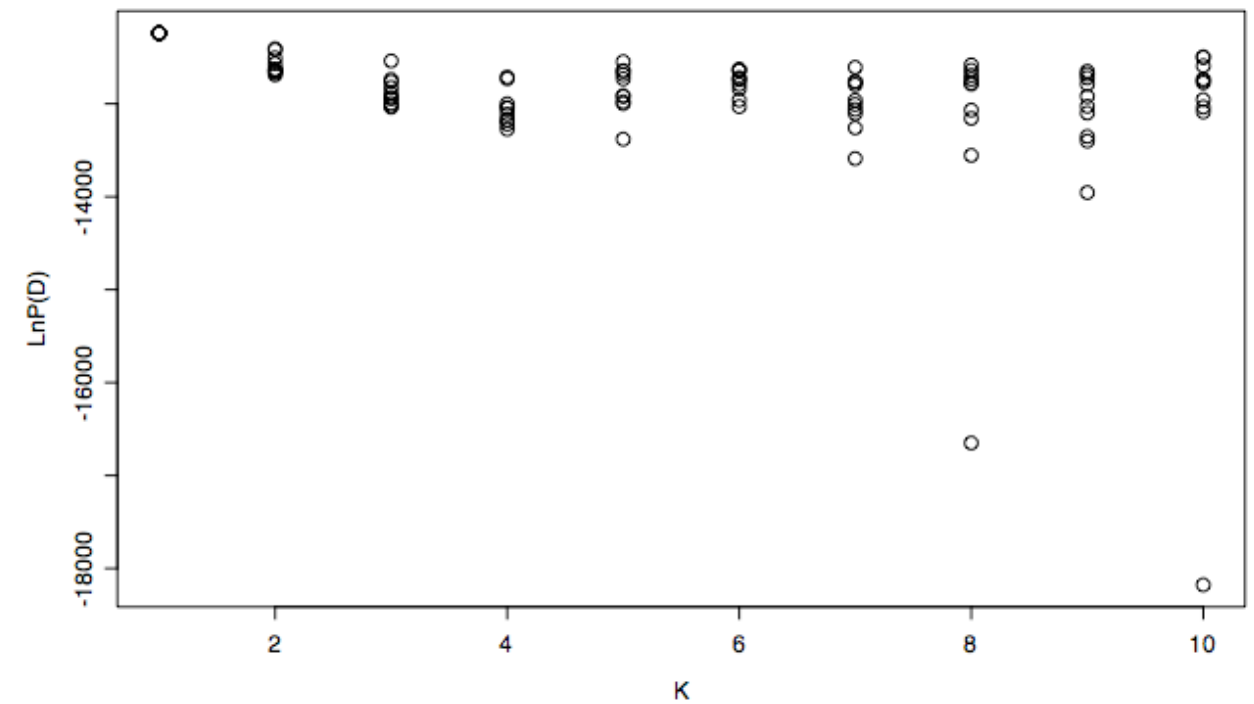
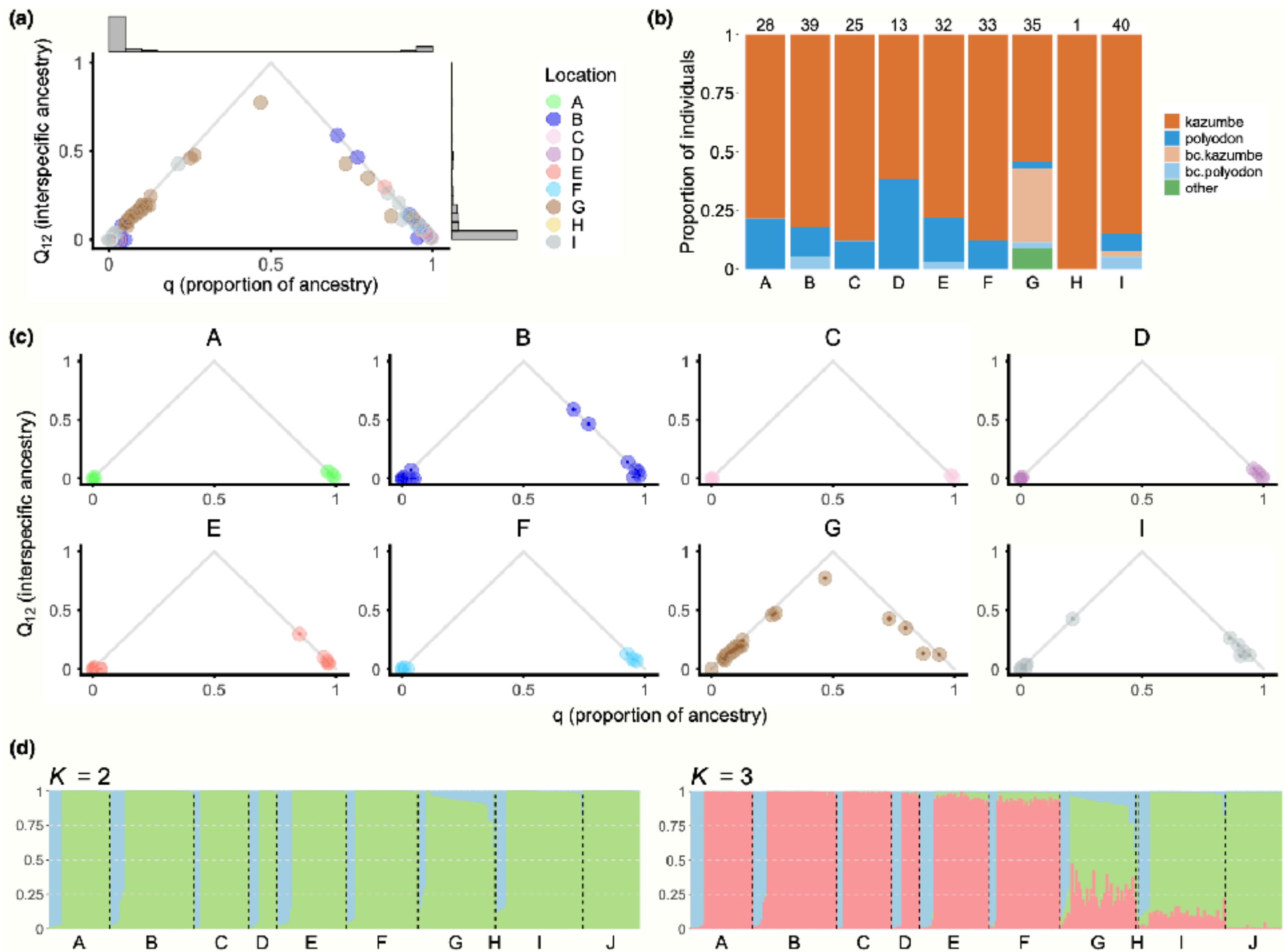
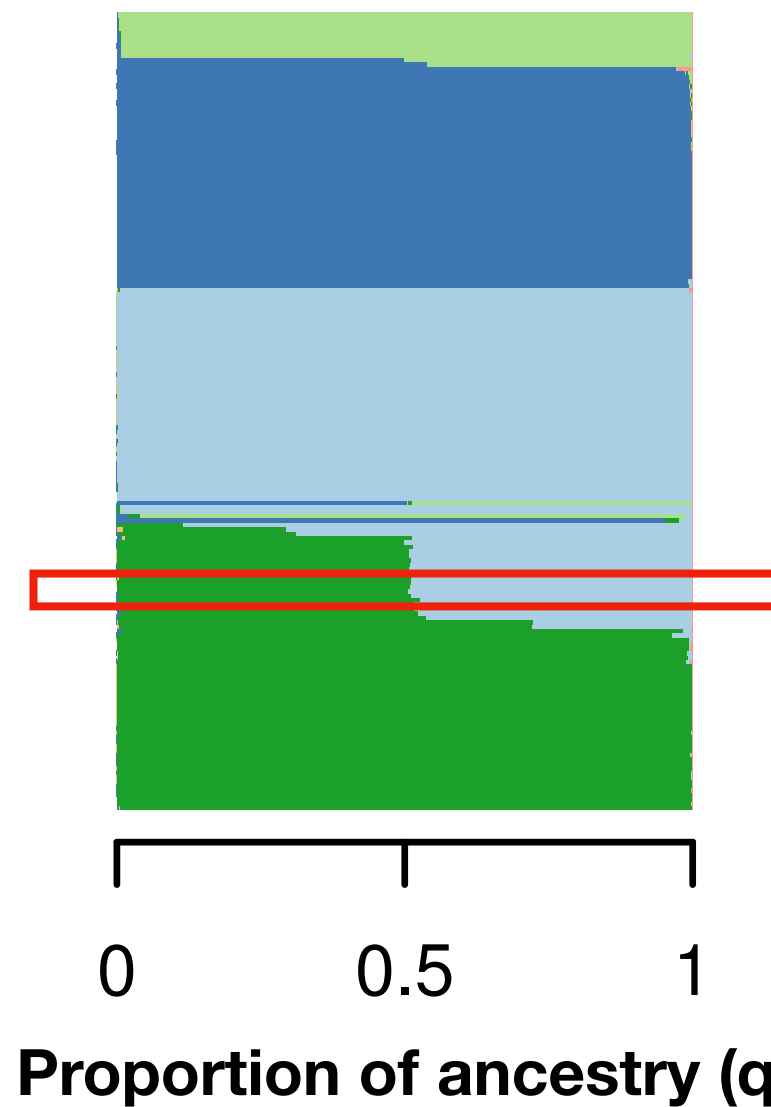
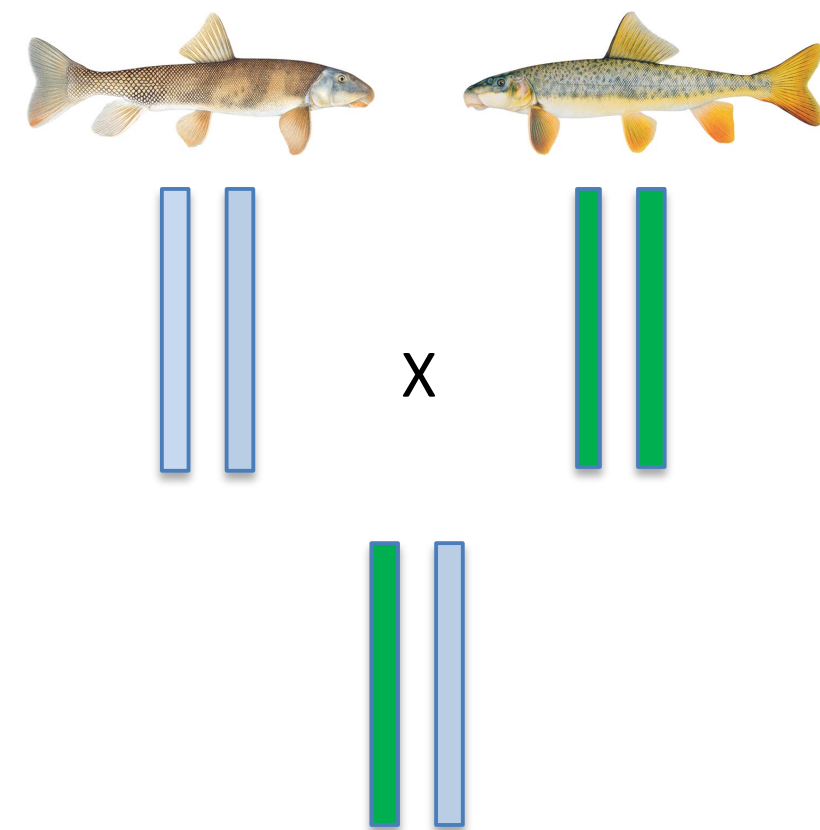


Figure S3. Log probability of the data (LnP(D)) for 10 STRUCTURE runs at each K value, for K 1 through 10 for *S. diagramma*.

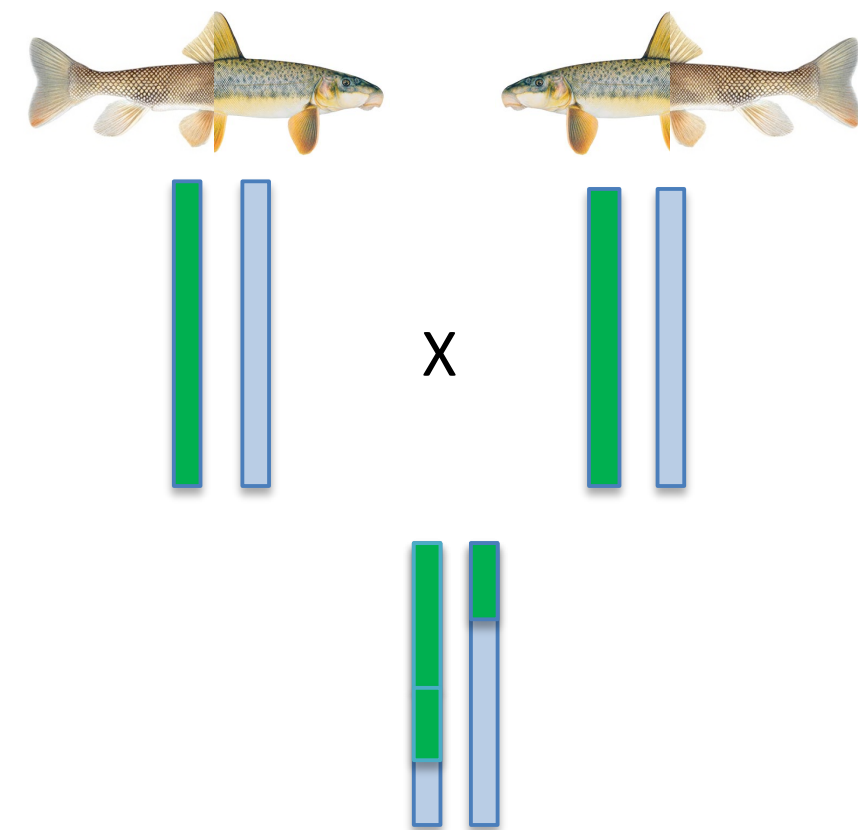


Interspecific ancestry helps differentiate among hybrids

First generation (F_1)

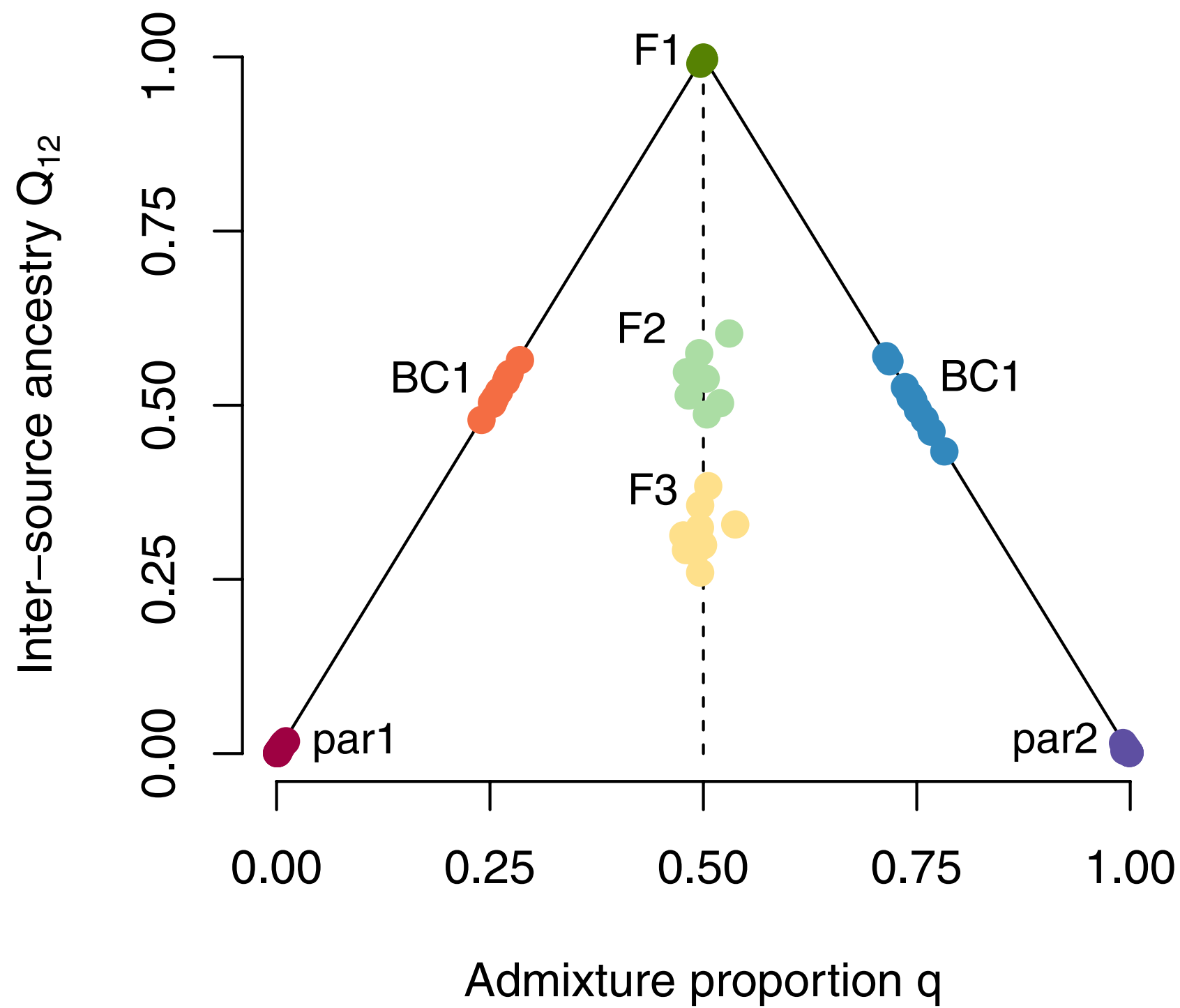


Later-generation ($F_2..F_n$)



Proportion of ancestry=0.5
Interspecific ancestry=1

Proportion of ancestry=0.5
Interspecific ancestry=0.5



Interspecific ancestry reveals history of hybridization

