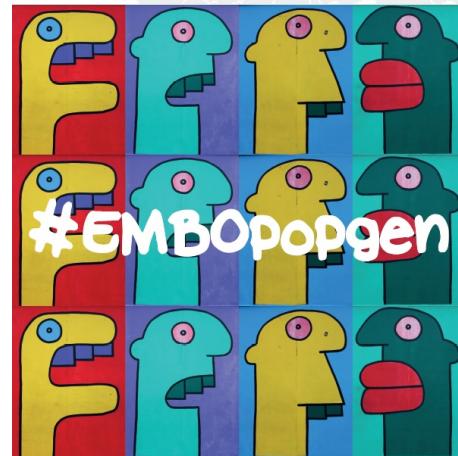


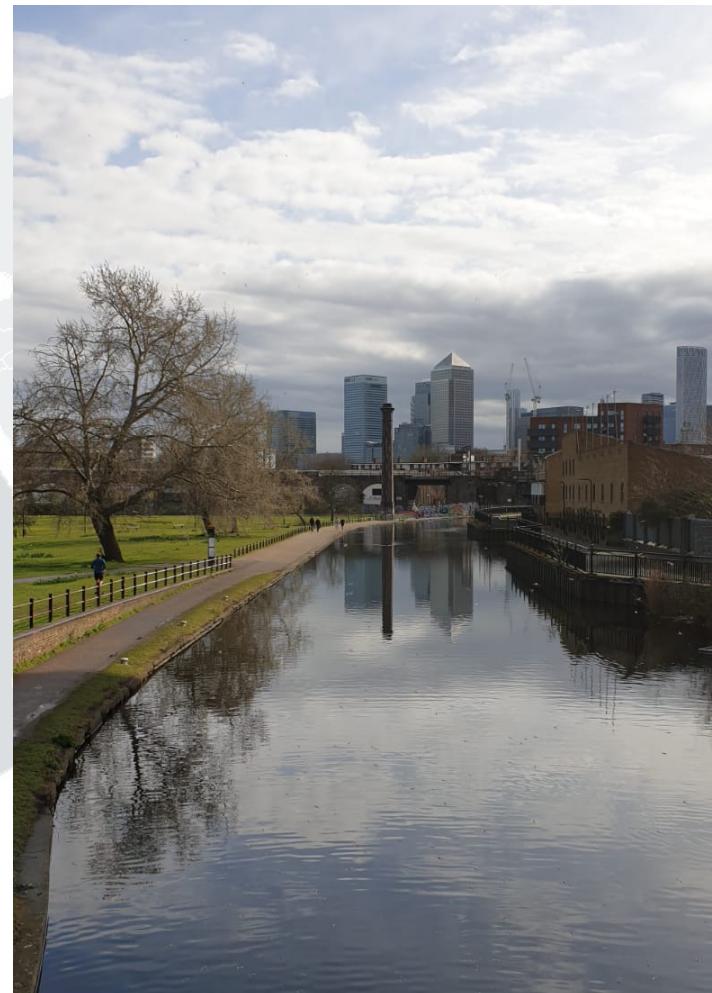
Inferring population structure and admixture histories



Lucy van Dorp

UCL Genetics Institute

lucy.dorp.12@ucl.ac.uk



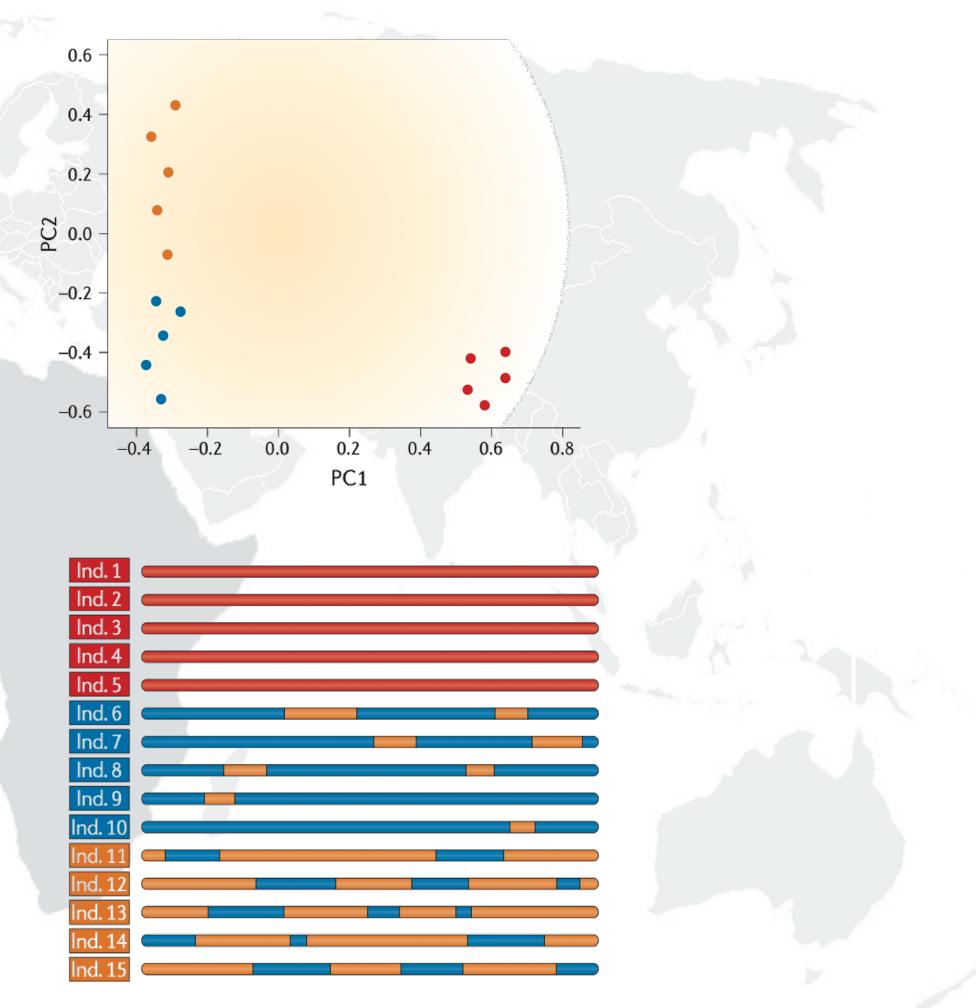
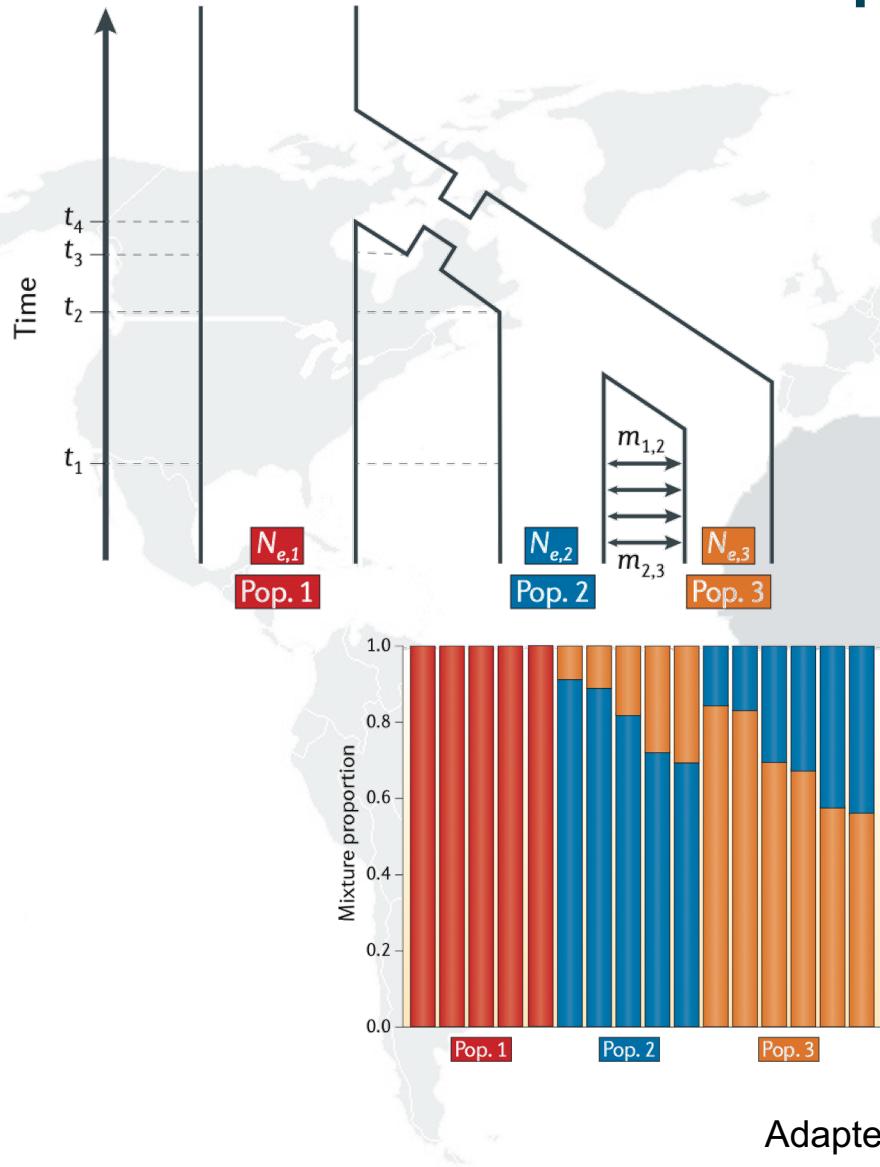
Genetic diversity in populations

- Genome-wide patterns of variation across individuals provide a powerful source of data
- Genetic diversity is underpinned by the key ancestral processes
 - Genetic Drift
 - Mutation
 - Recombination
 - Selection
- There are many different ways of measuring population differentiation and also many different types of data, from **allele frequencies** at single polymorphic sites (or multiple unlinked sites) to full **haplotypic sequences** in coding or noncoding regions.

Determination of population structure

Name	Data type	Inference	Notes	Ref
STRUCTURE	Unlinked multi-allelic genotypes	Population structure, admixture	User friendly GUI; can be computationally demanding	Pritchard, Stephens & Donnelly. Genetics. 2000.
ADMIXTURE	Unlinked bi-allelic SNPs	Population structure, admixture	Estimates the number of populations via cross-validation error	Alexander, Novembre, & Lange. Genome Res. 2009.
fineSTRUCTURE	Phased haplotypes	Population structure, admixture, chromosome painting	Can be used to identify the number and identity of populations	Lawson, Hellenthal, Myers & Falush. PLoS Genet. 2012.
GLOBETROTTER	Phased haplotypes	Population structure, admixture, chromosome painting	Estimate unsampled ancestral populations and admixture times	Hellenthal et al. Science. 2014.

Determination of population structure

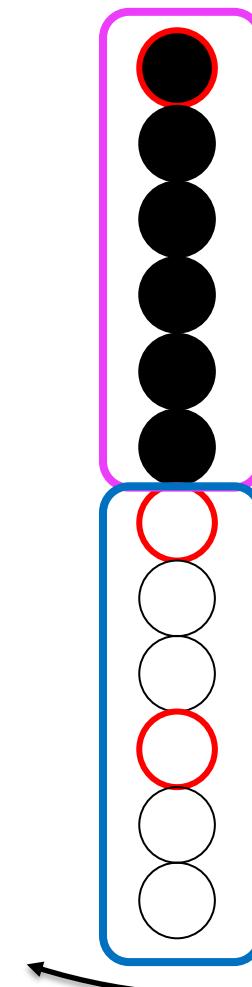
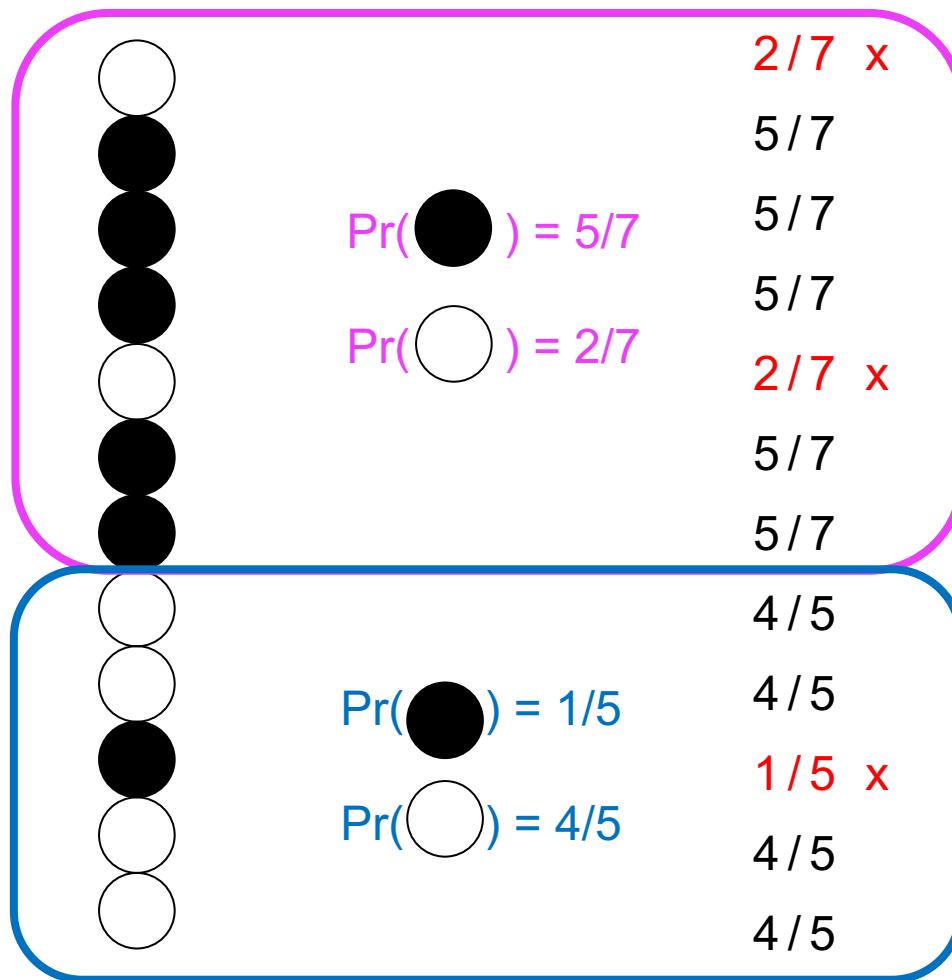


Adapted from Schraiber & Akey. *Nature Review Genetics*. 2015.

Allele frequency-based clustering

S_{i1}, \dots, S_{iL}

$\Pr(S_{il} = j) = pk_{jl}$, where pk_{jl} is the frequency of allele j at locus l in cluster k



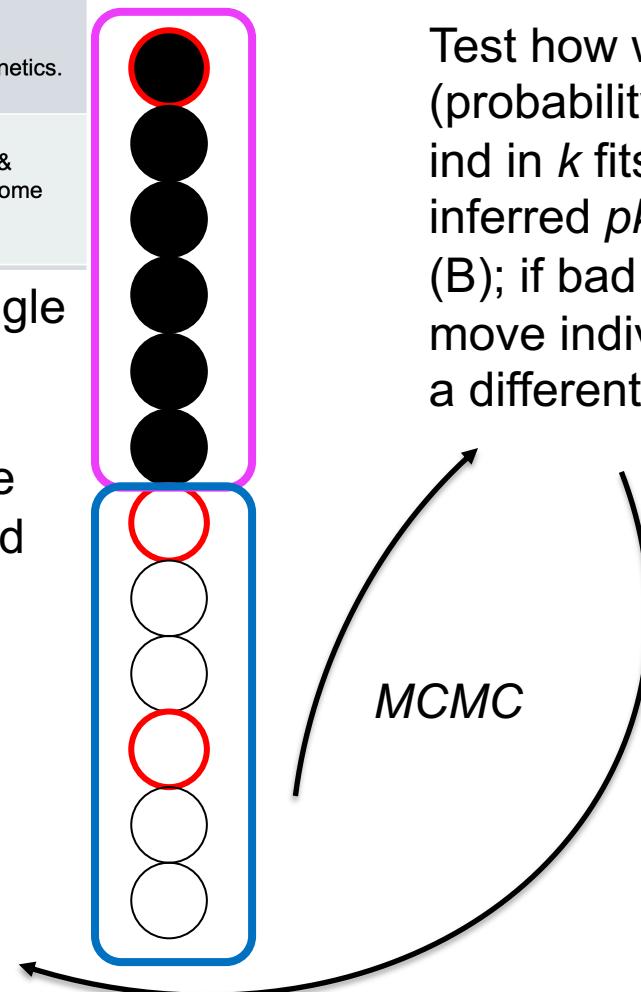
Test how well (probability) each ind in k fits the inferred pk_{jl} from (B); if bad fit → move individual to a different cluster

LD prune

Allele frequency-based clustering

Name	Data type	Inference	Notes	Ref
STRUCTURE	Unlinked multi-allelic genotypes	Population structure, admixture	User friendly GUI; can be computationally demanding	Pritchard, Stephens & Donnelly. Genetics. 2000.
ADMIXTURE	Unlinked bi-allelic SNPs	Population structure, admixture	Estimates the number of populations via cross-validation error	Alexander, Novembre, & Lange. Genome Res. 2009.

- “no admixture model” – assign each ind i to single cluster k
- “admixture model” – assign each ind to multiple clusters (i.e. infer % of ind i ’s genome assigned to clusters $1, \dots, K$)
- “linkage model” – can identify regions of ind i assigned to each cluster (Falush et al 2003, *Genetics* **164**:1567)
- Spatial priors on cluster membership.

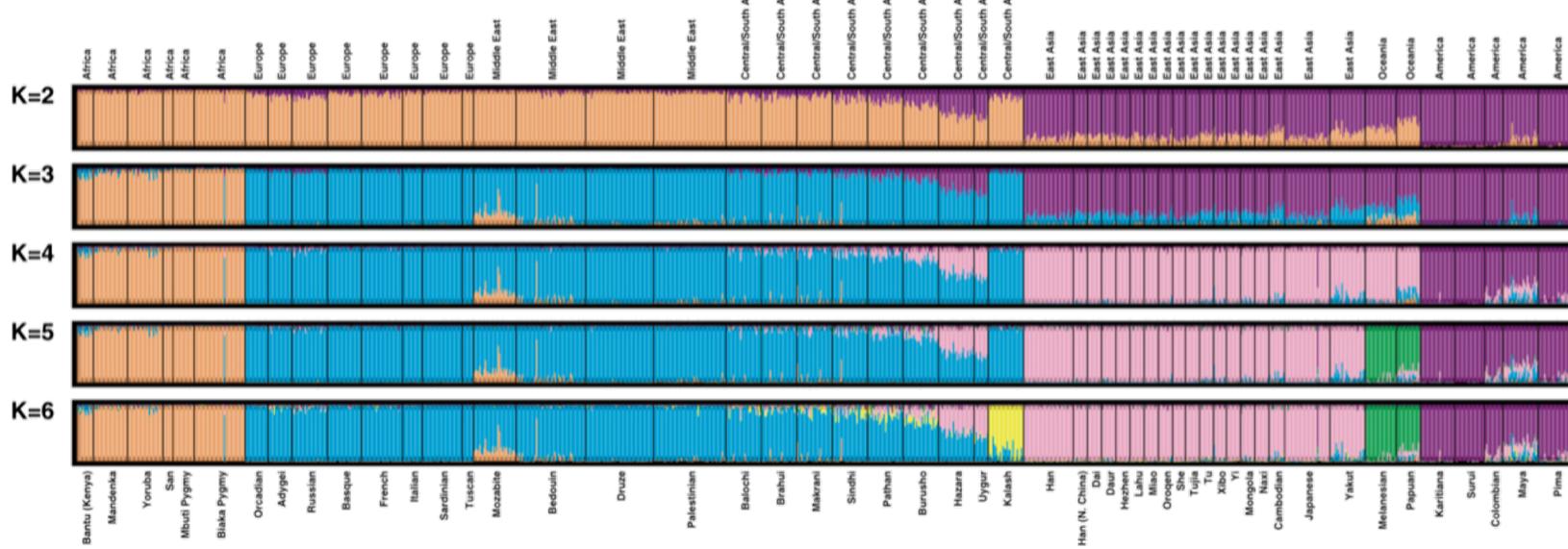


Test how well (probability) each ind in k fits the inferred $p_{kj|l}$ from (B); if bad fit → move individual to a different cluster

Allele frequency-based clustering

Human Genome Diversity Panel

Rosenberg et al. *Science*. (2002), 298, 5602, 2381-2385.

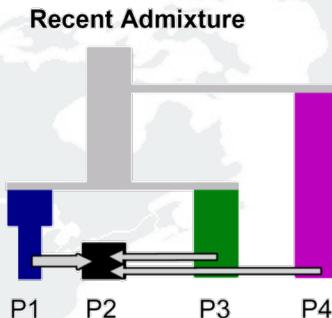


- Accuracy of population identification depends on the number of individuals included, the number of loci used and the allele frequency differences among those populations.
- Population identification can be strongly affected by the amount of admixture within individuals and the proportion of admixed individuals in the dataset – all populations are admixed....

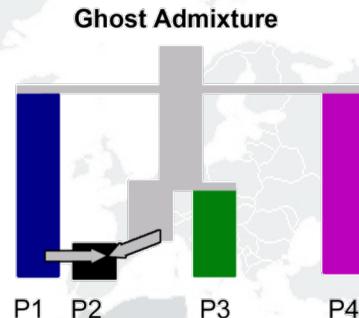
Allele frequency-based clustering

Different combinations of drift, admixture and ancient relatedness can give similar signals

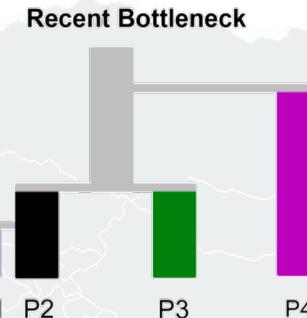
a



Recent admixture of P1,
P3 & P4
Some parallels to
African Americans



P2 50:50 admixture
of P1 and
unsampled
population most
closely related to
P3

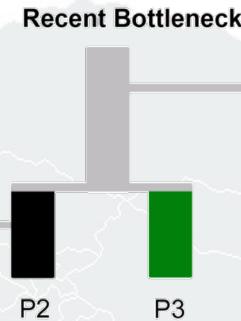
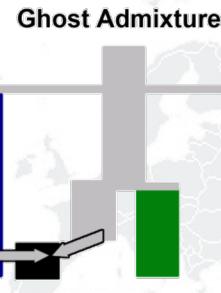
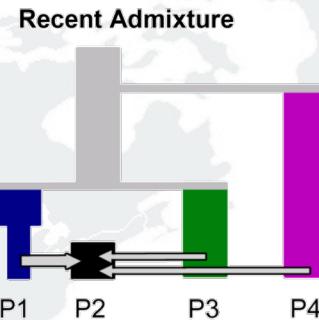


P1 – a sister
population to P2
that underwent a
strong recent
bottleneck

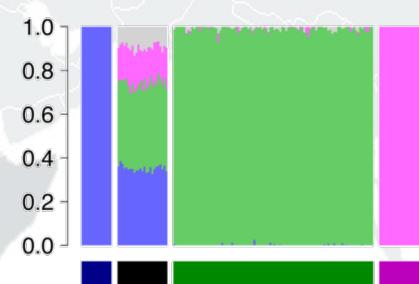
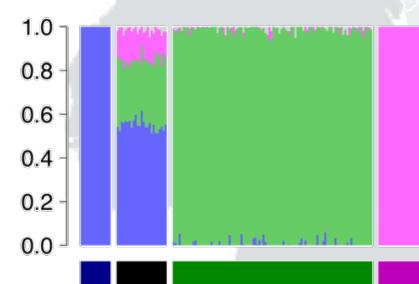
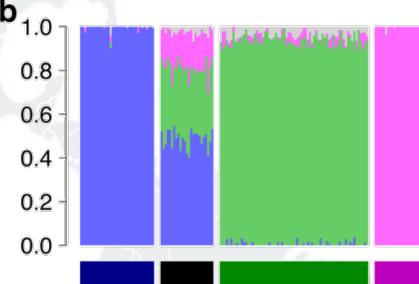
Allele frequency-based clustering

Different combinations of drift, admixture and ancient relatedness can give similar signals

a

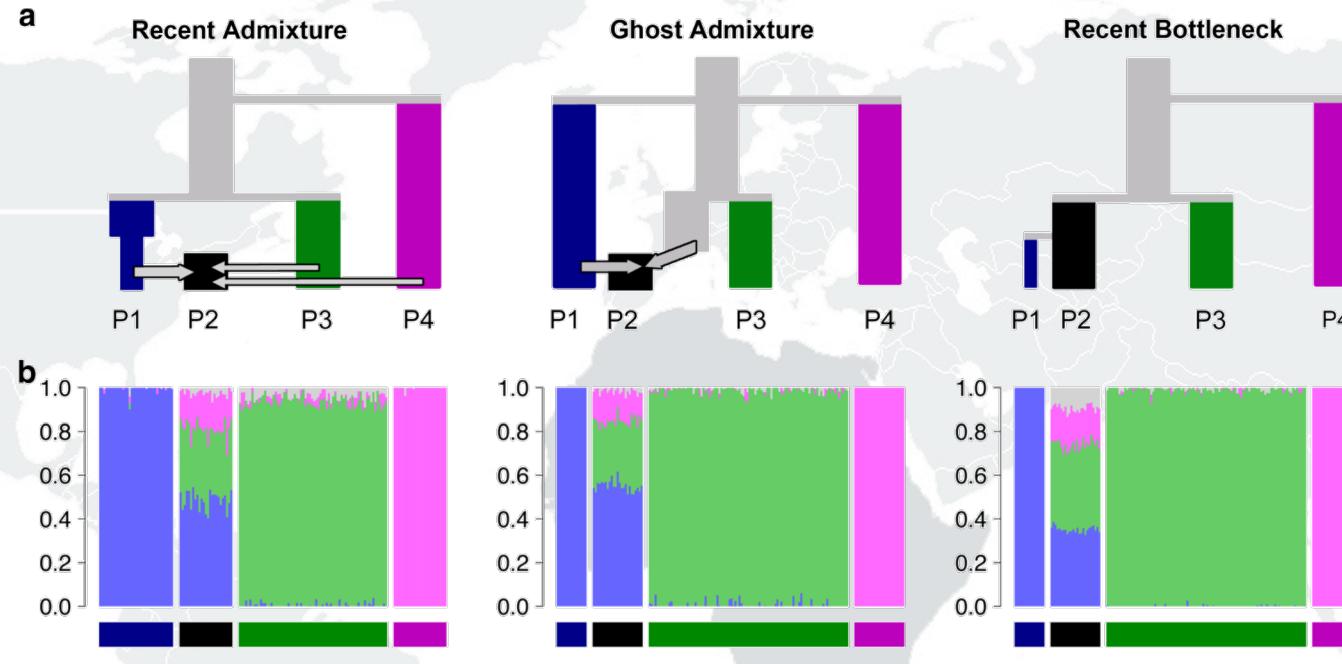


b



Allele frequency-based clustering

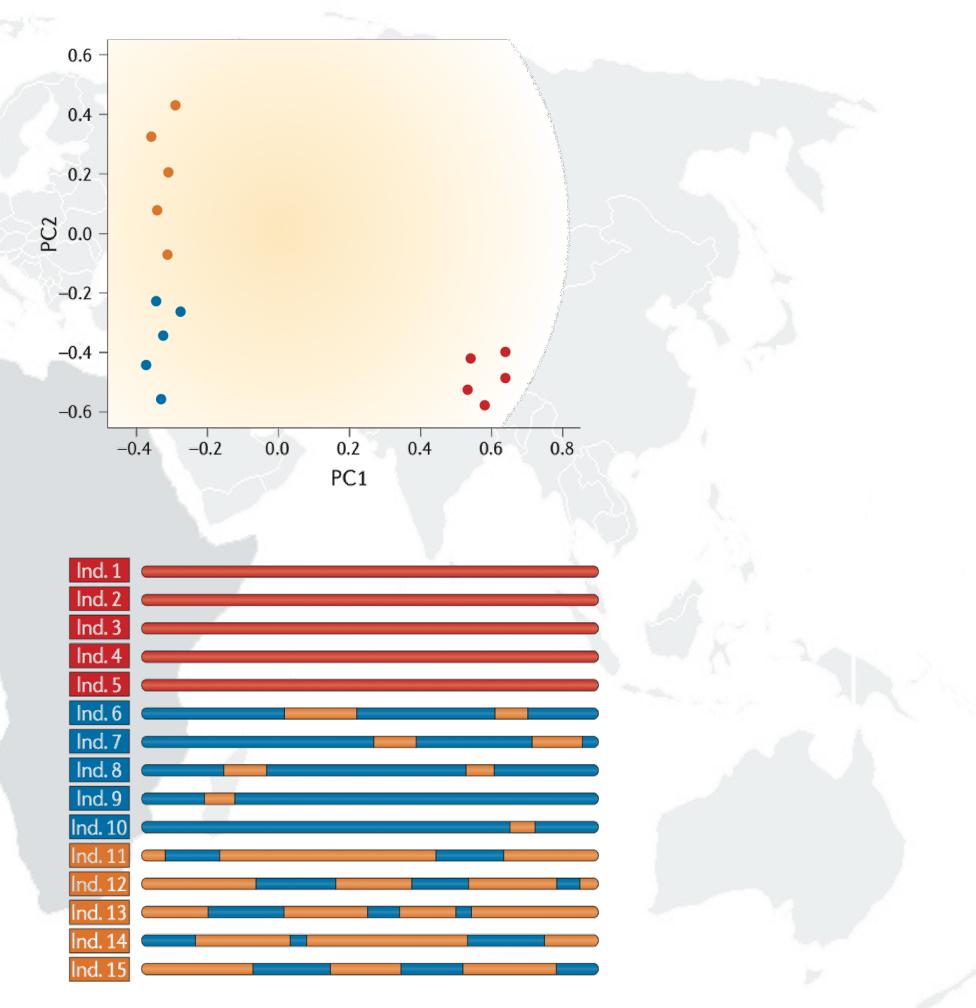
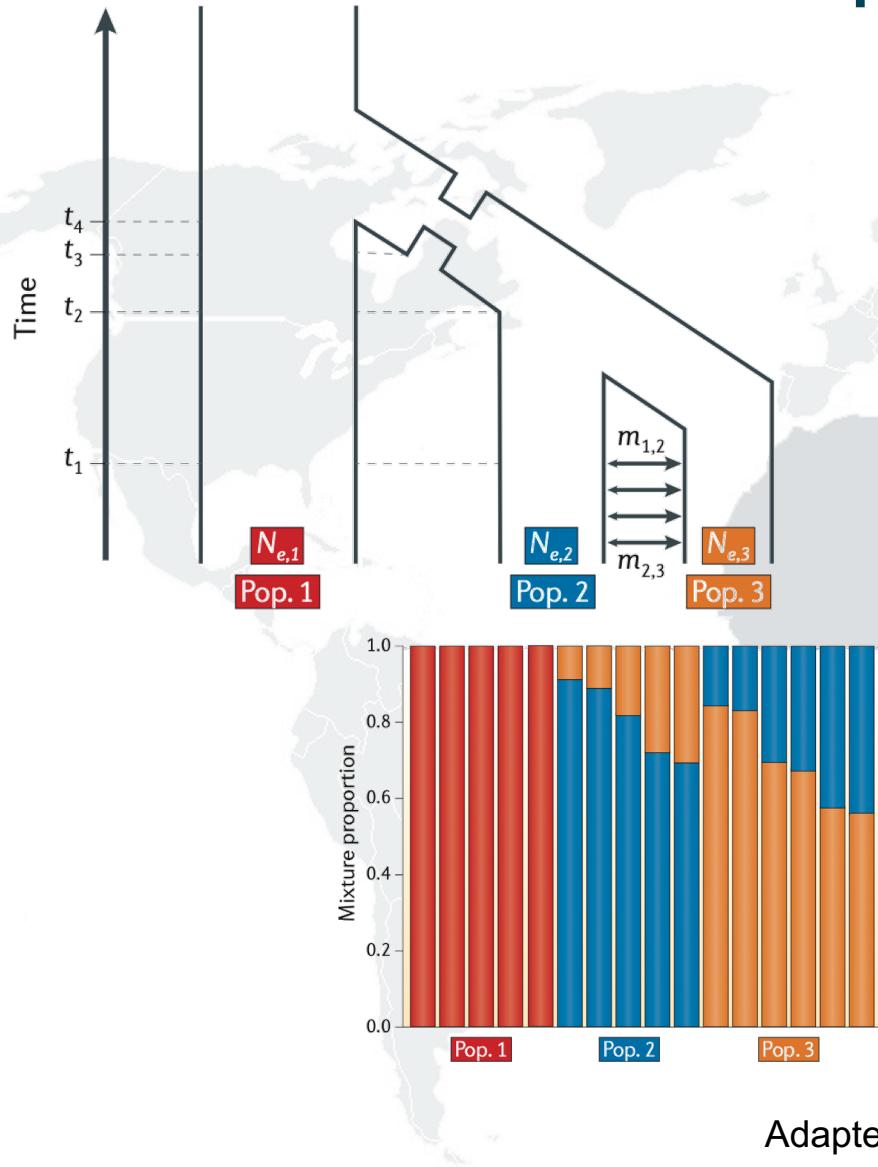
Different combinations of drift, admixture and ancient relatedness can give similar signals



Lawson, van Dorp & Falush. *Nat Comms.* (2018).

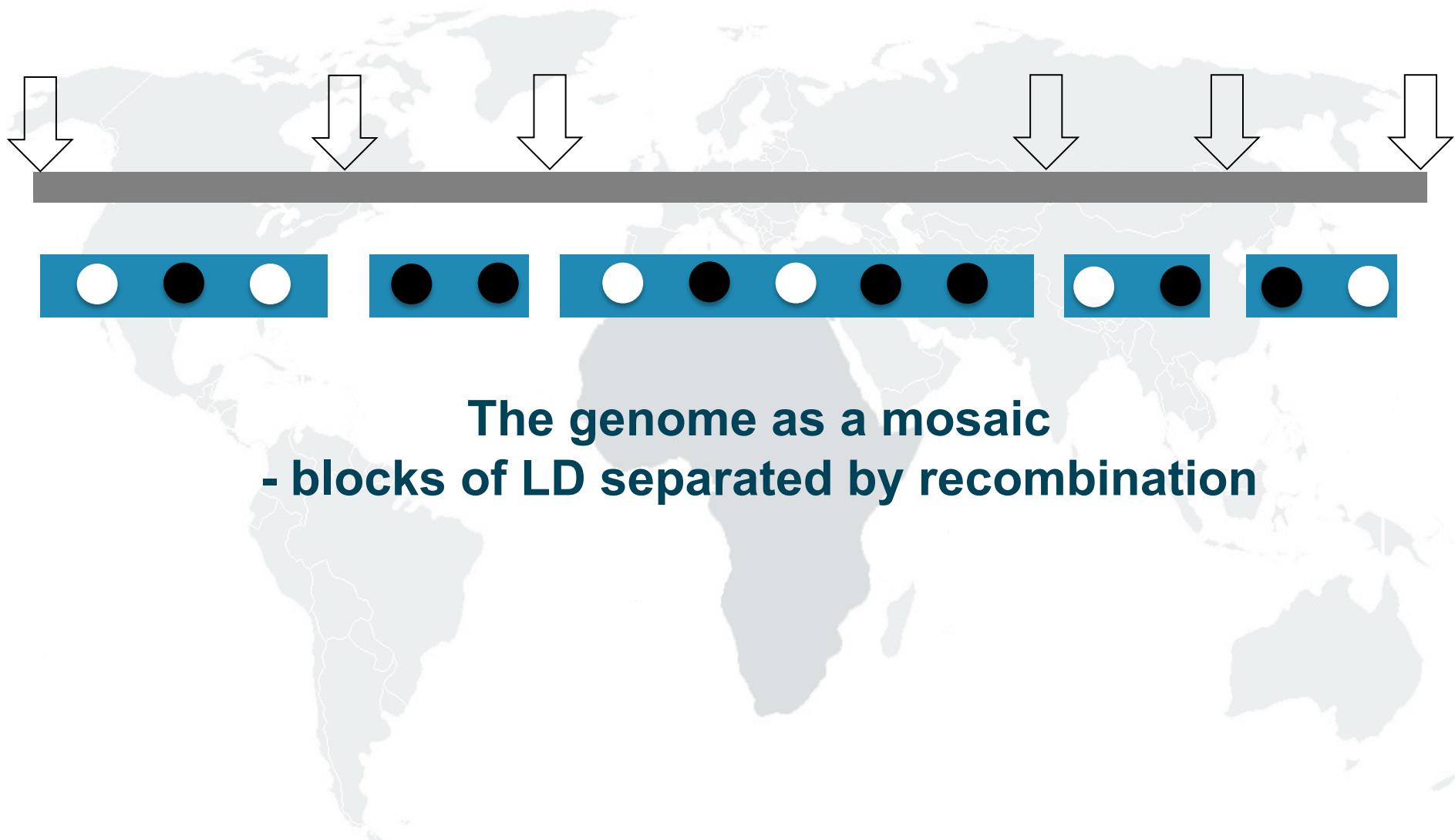
Ancestry assignments in STRUCTURE/ADMIXTURE often termed admixture proportions
 they do not identify where admixture has occurred → they do not alone indicate a history of
 admixture.

Determination of population structure



Adapted from Schraiber & Akey. *Nature Review Genetics*. 2015.

Incorporating haplotype information

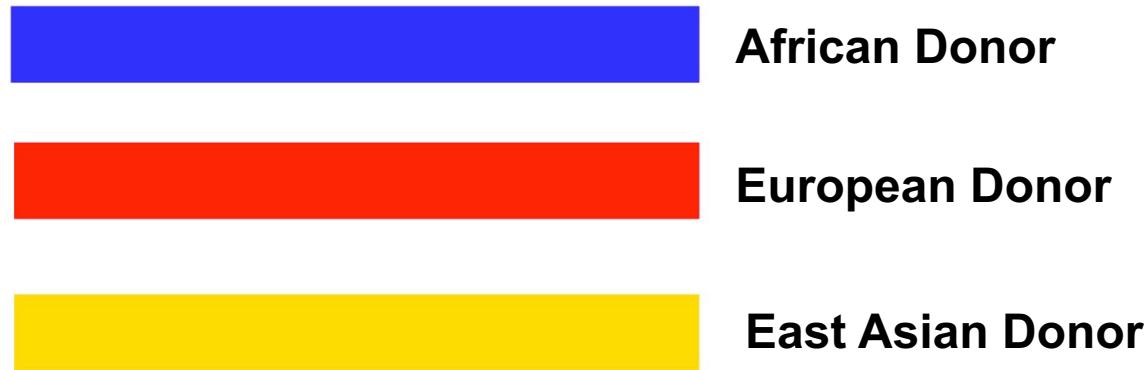


Incorporating haplotype information

Name	Data type	Inference	Notes	Ref
HAPMIX	Phased haplotypes, reference panel	Chromosome painting	Requires populations to be pre-specified	Price et al. PLoS Genetics. 2009
LAMP	Phased haplotypes	Chromosome painting	Identifies local ancestry in windows, rather than using an HMM	Sankararaman et al. Am. J. Hum. Gen. 2008.
PCAdmix	Phased haplotypes	Chromosome painting, population structure	Uses PCA in small chunks followed by HMM to estimate local ancestry	Brisbin et al. Hum. Biol. 2012.
fineSTRUCTURE	Phased haplotypes	Population structure, admixture, chromosome painting	Can be used to identify the number and identity of populations	Lawson, Hellenthal, Myers & Falush. PLoS Genet. 2012.
GLOBETROTTER	Phased haplotypes	Population structure, admixture, chromosome painting	Estimate unsampled ancestral populations and admixture times	Hellenthal et al. Science. 2014.

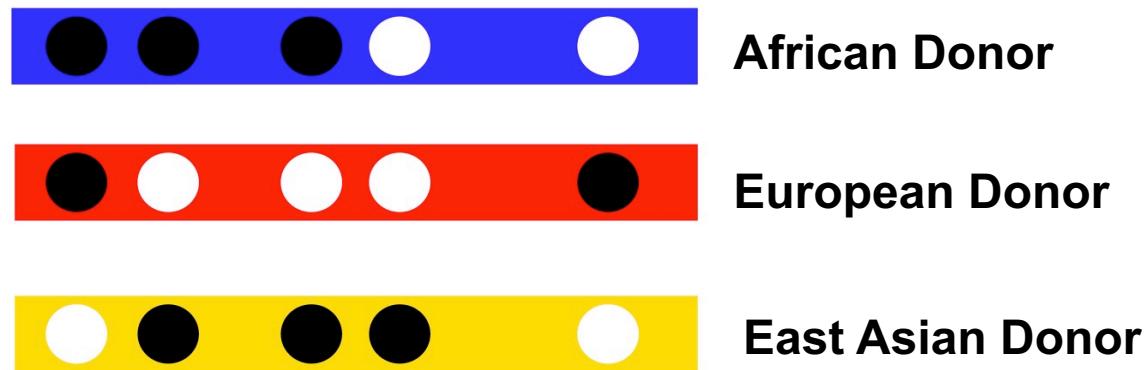
Incorporating haplotype information

Chromosome painting



Incorporating haplotype information

Chromosome painting

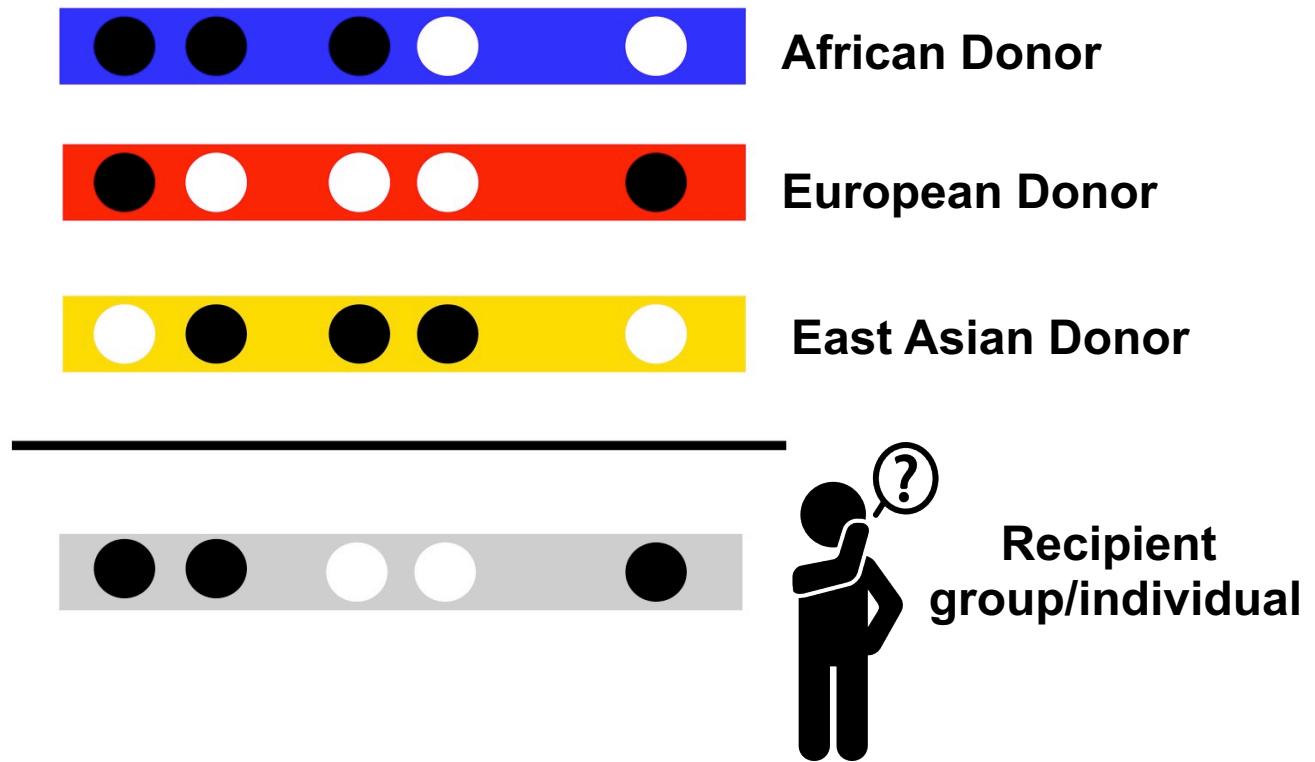


Li & Stephens, *Genetics*, 2003

Lawson et al, *PLoS Genetics*, 2012

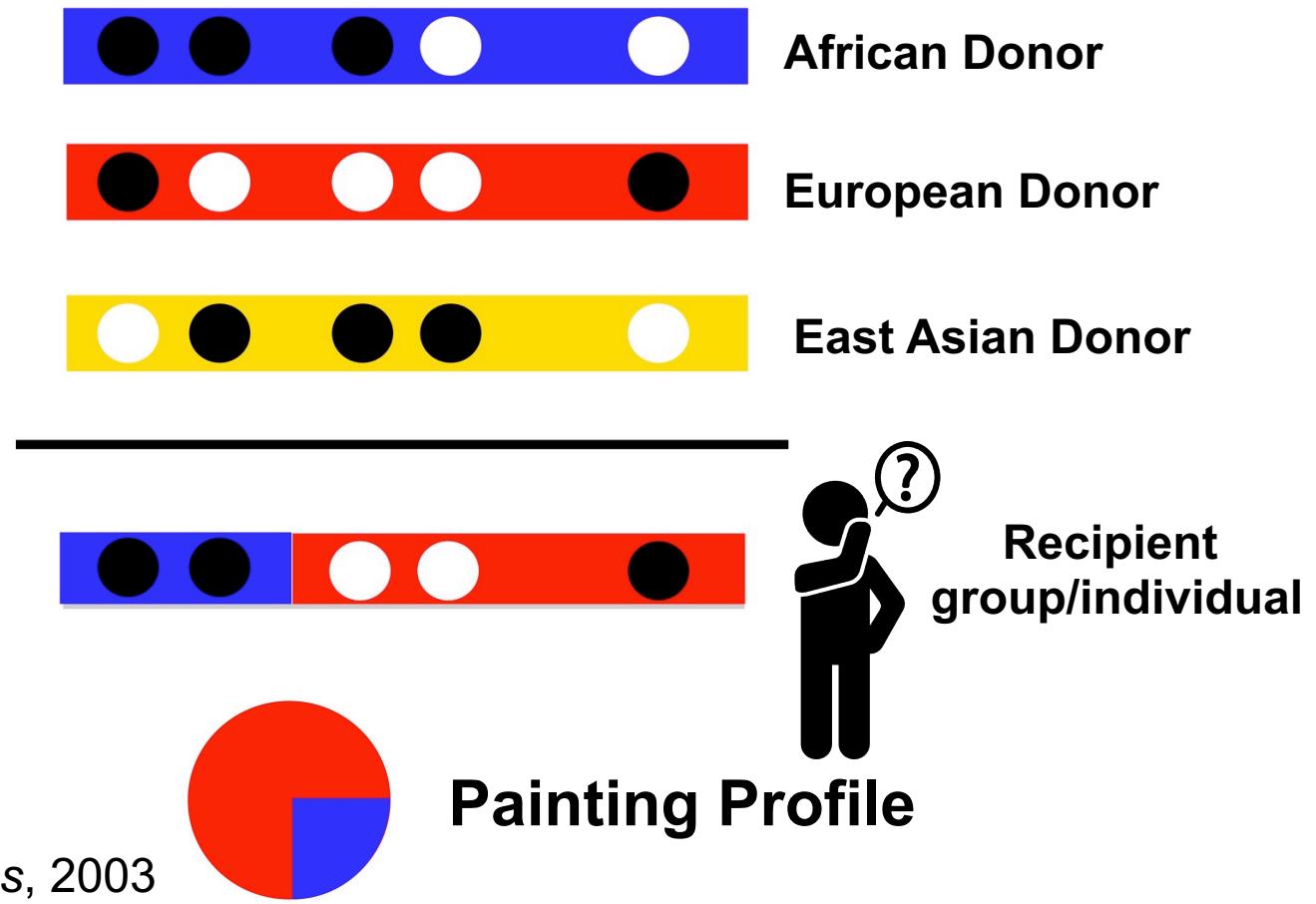
Incorporating haplotype information

Chromosome painting



Incorporating haplotype information

Chromosome painting

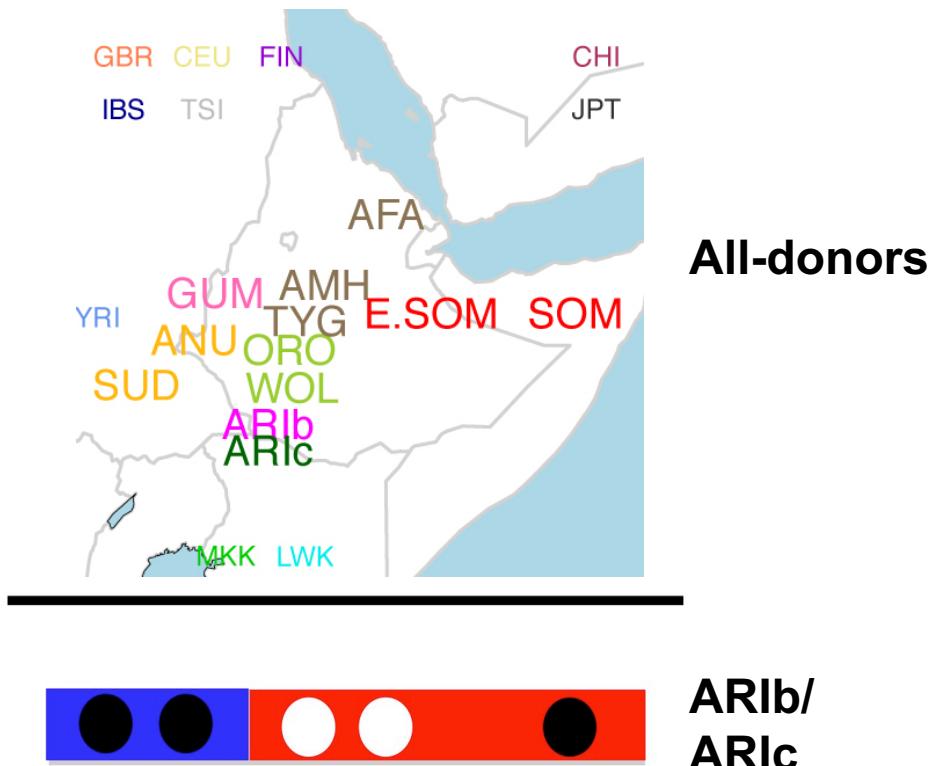


Li & Stephens, *Genetics*, 2003

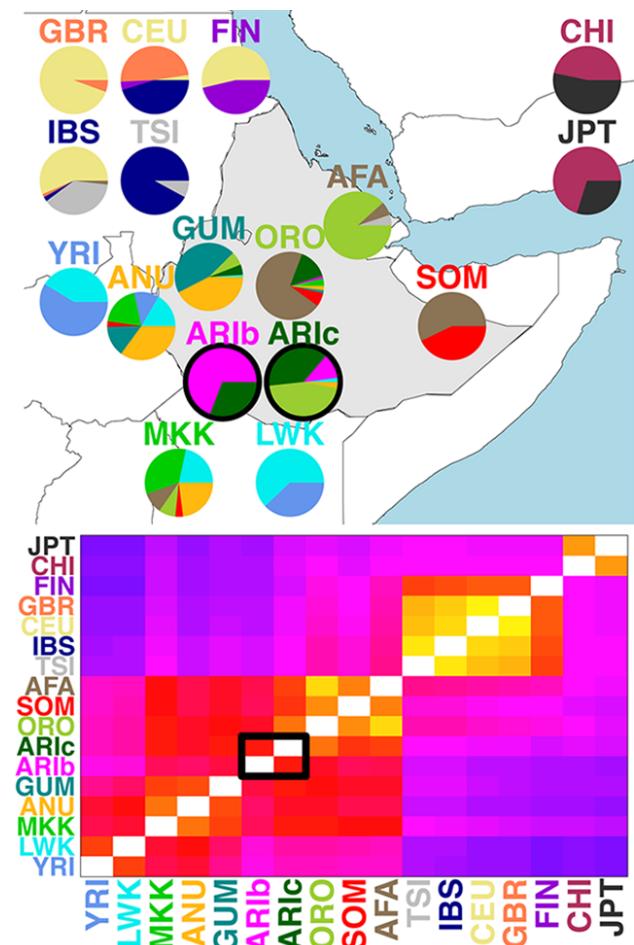
Lawson et al, *PLoS Genetics*, 2012

Incorporating haplotype information

Chromosome painting

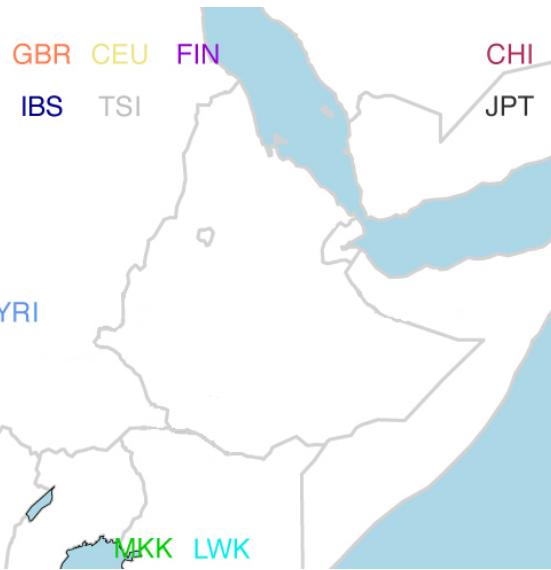


Average haplotype length shared = 2.76cM



Incorporating haplotype information

Chromosome painting



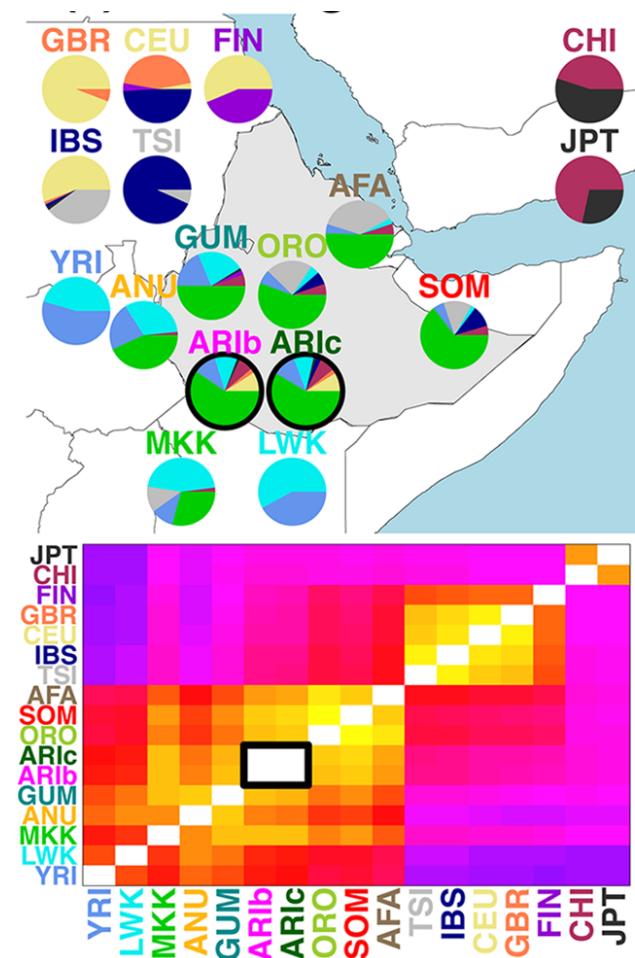
No Ethiopia



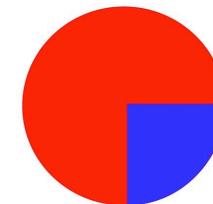
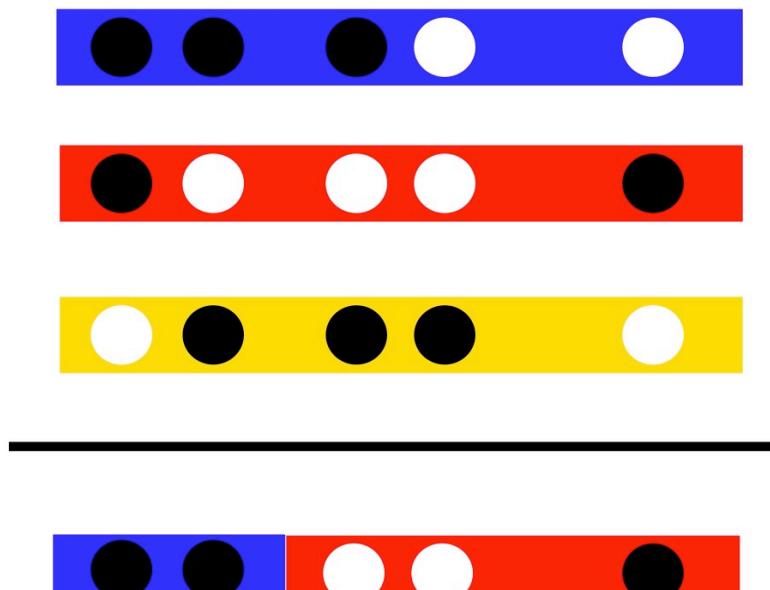
ARIb/
ARIC

Average haplotype length shared = 0.65cM

van Dorp et al, *PLoS Genetics*, 2015



Haplotype clustering

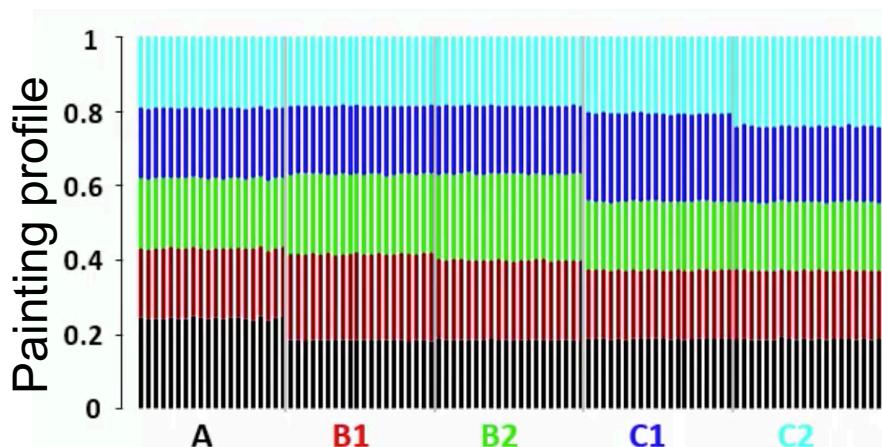
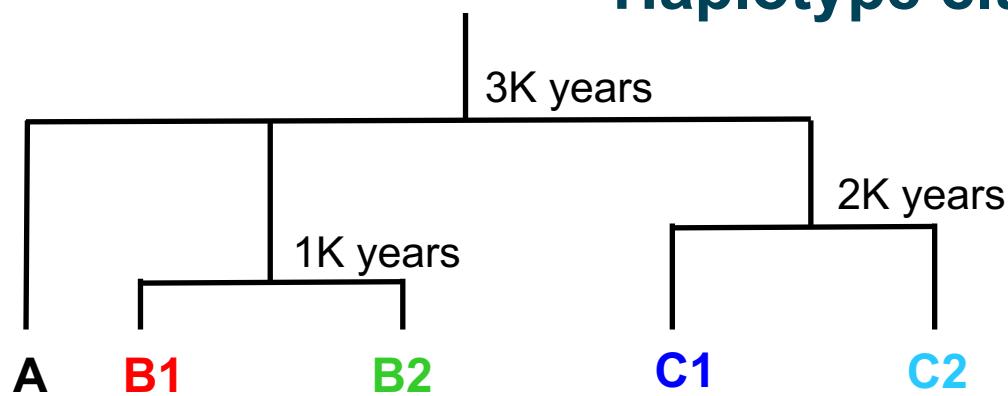


Clustering
fineSTRUCTURE



Painting Profile

Haplotype clustering

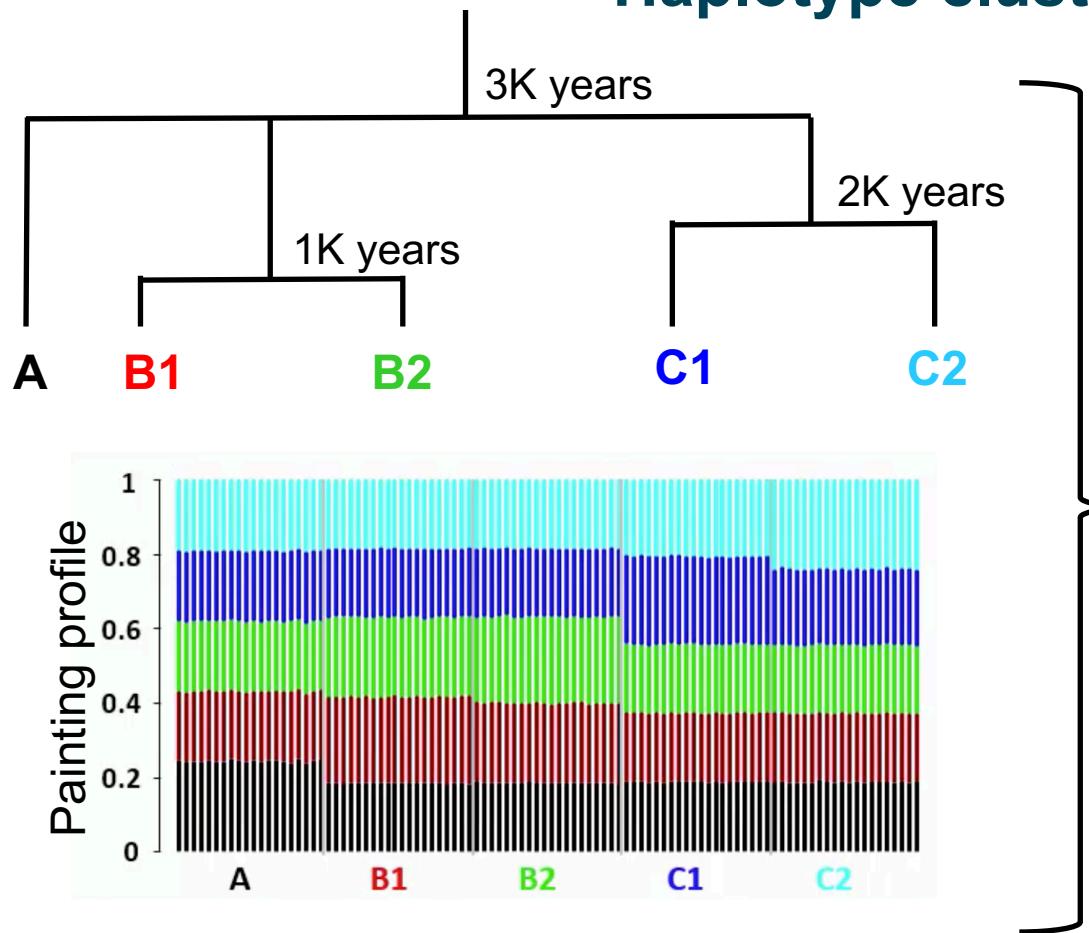


Donor individual = j

Recipient individual = i

Number of chunks of DNA by which individual i is painted by individual j = Y_{ij}

Haplotype clustering



Assign individuals to a random set of clusters k

$(Y_{i1}, \dots, Y_{ik}) \sim \text{Mult}(P_{A1}, \dots, P_{Ak})$
for ind I assigned to cluster A

Infer P_{Ak} – number of chunks individuals in cluster A are painted with individuals in cluster k

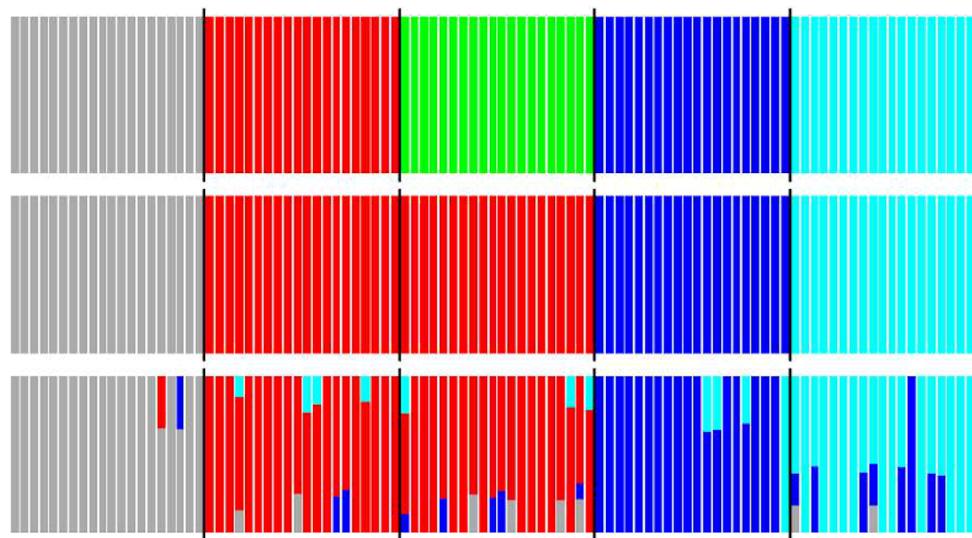
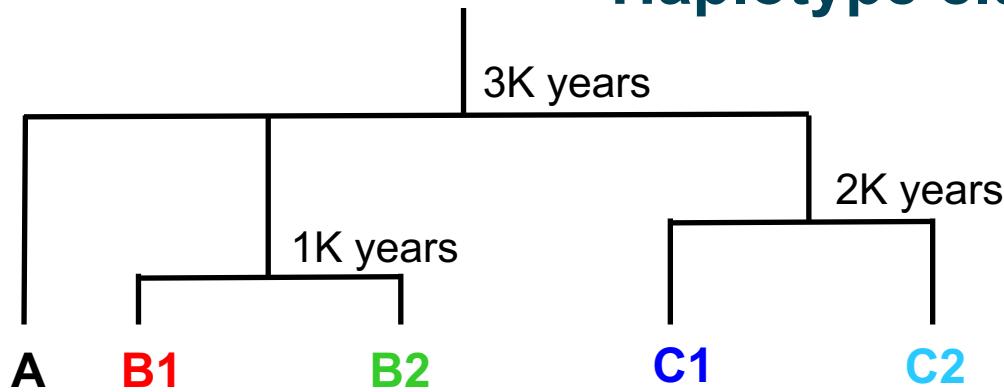
Test if $\Pr(Y_{i1}, \dots, Y_{ik})$ is low – change cluster and repeat!

Donor individual = j

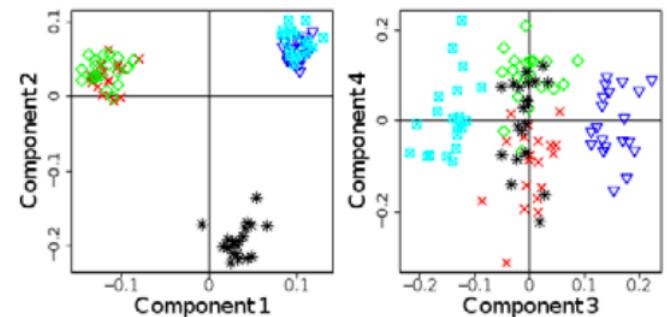
Recipient individual = i

Number of chunks of DNA by which individual I is painted by individual j = Y_{ij}

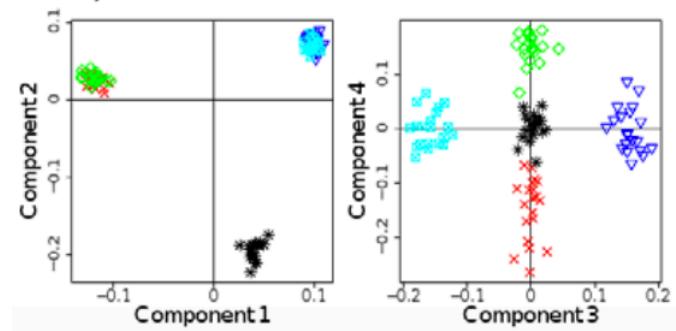
Haplotype clustering



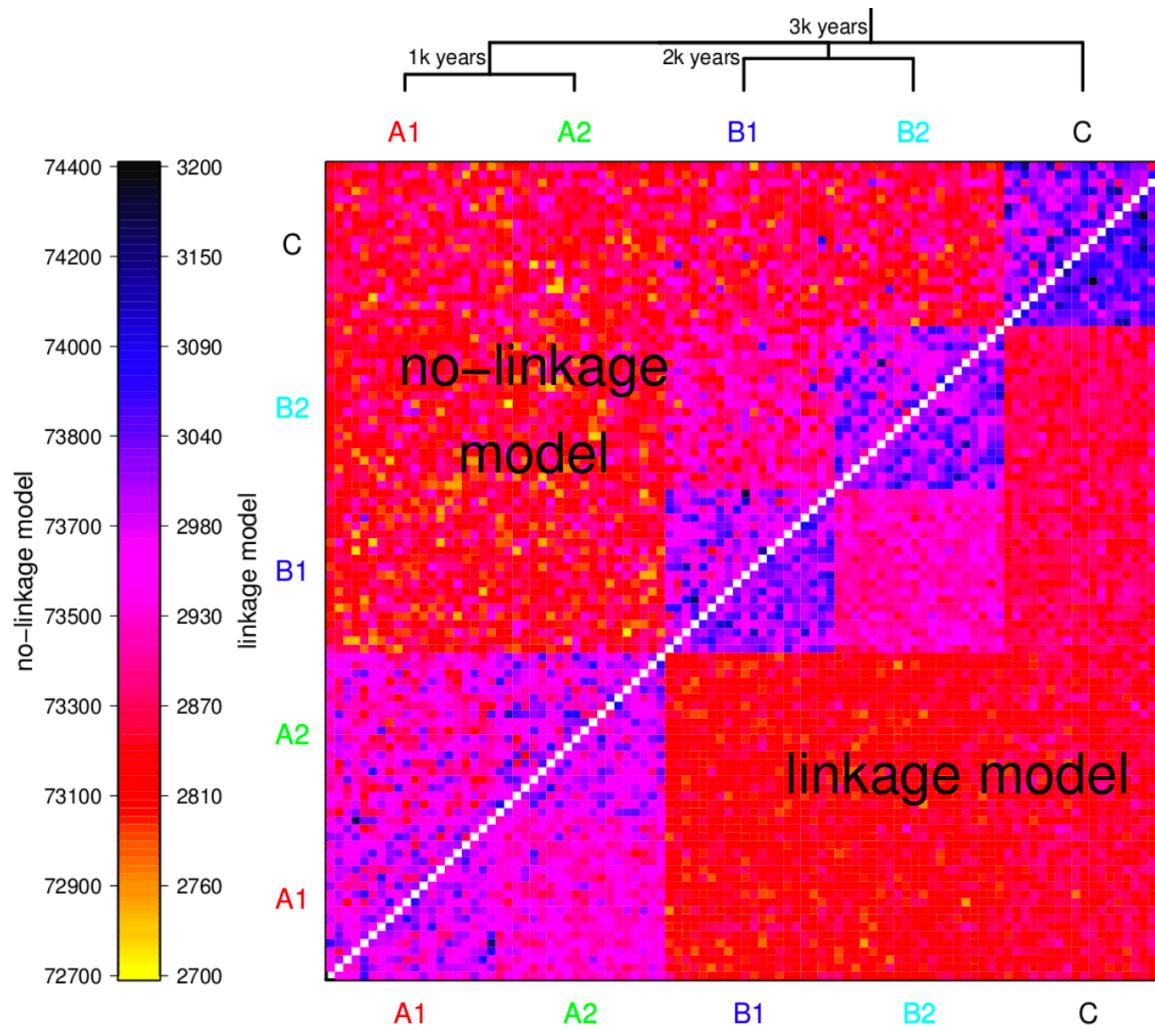
E) Unlinked model PCA



F) Linked model PCA

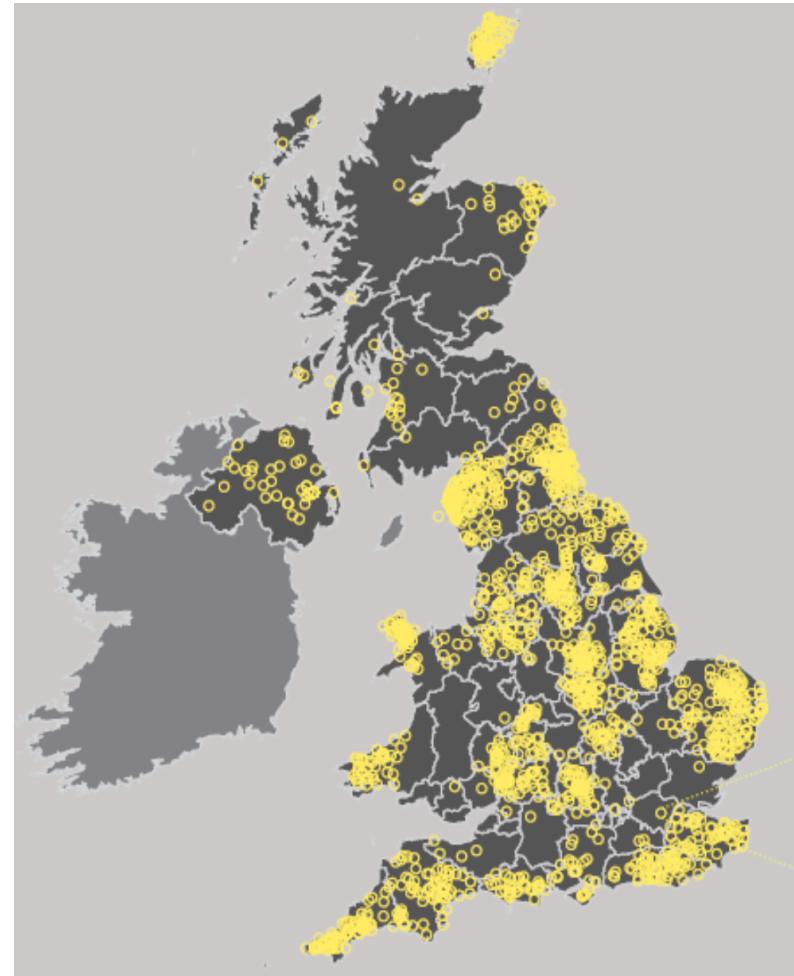


Haplotype clustering

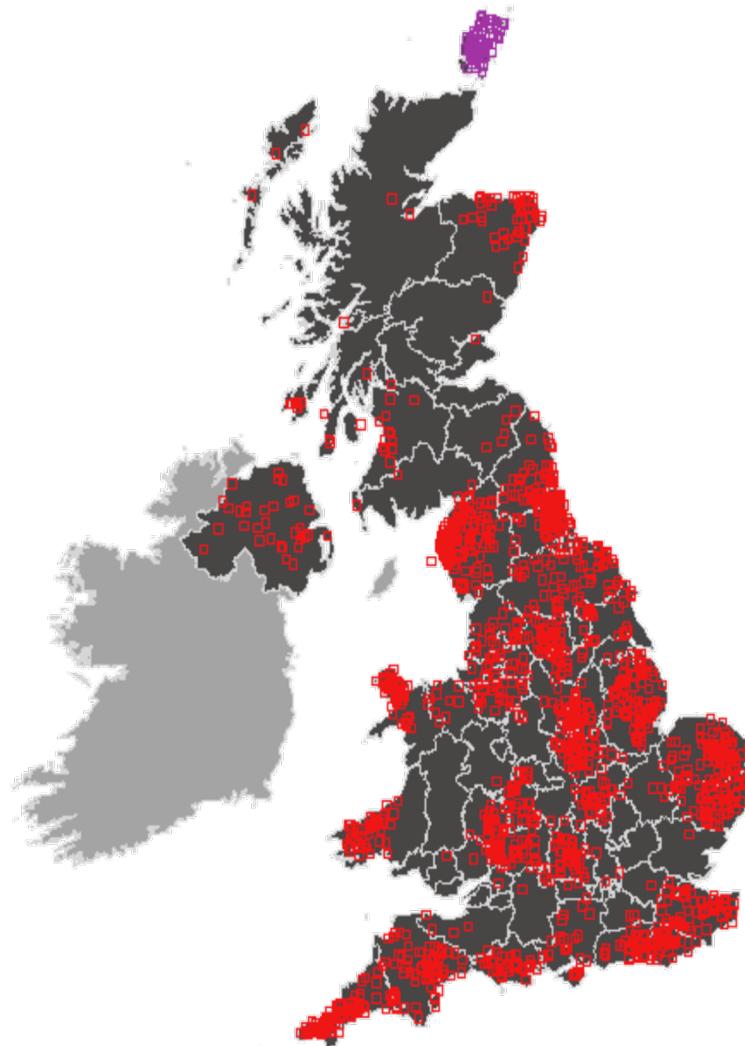


Haplotype clustering in the British Isles

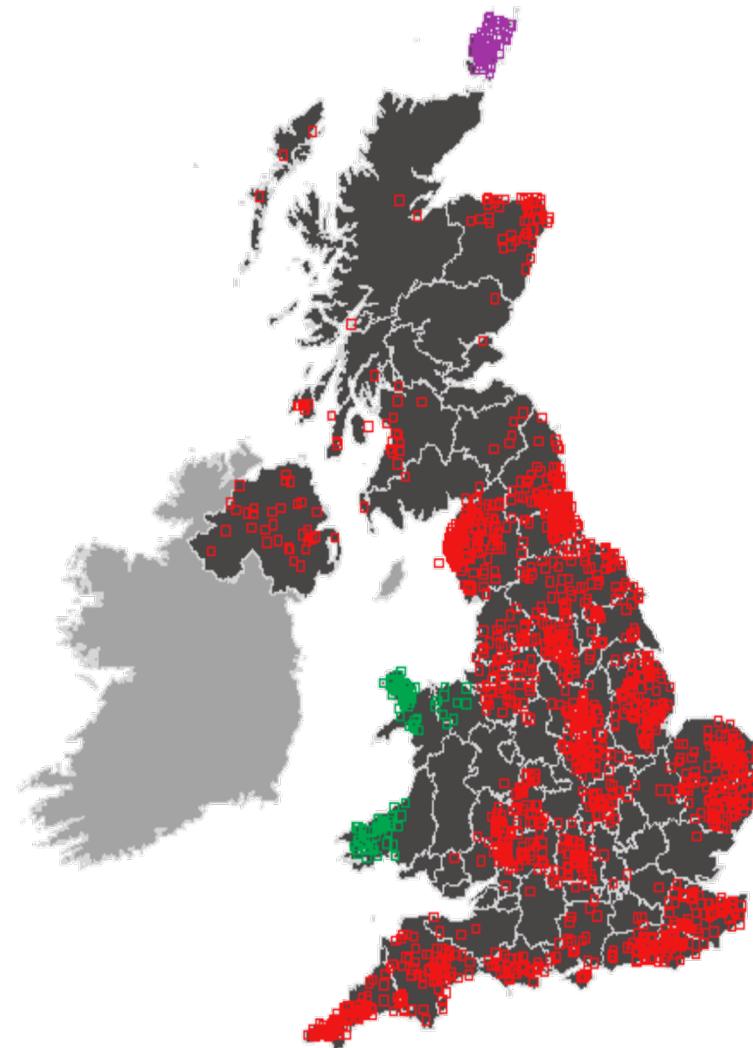
- 2,039 individuals collected across rural areas of UK
- All four-grandparents born within 80km of one another
- Inferred DNA matching patterns amongst all individuals
- Grouped individuals into hierarchical groups based on these DNA patterns



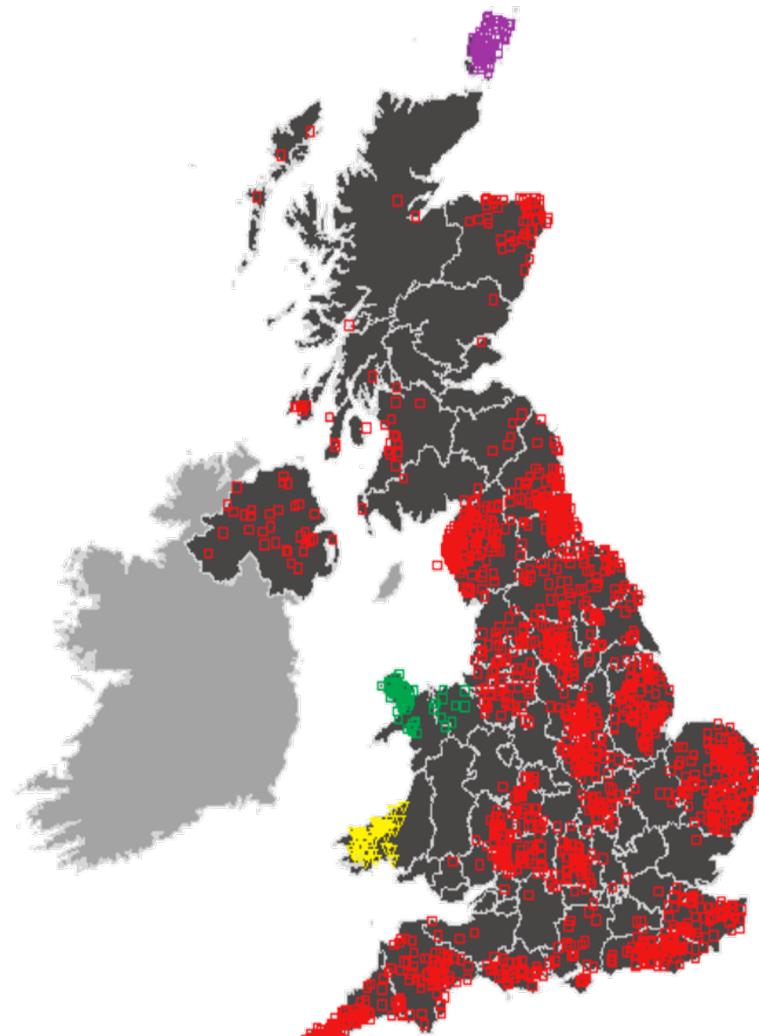
Haplotype clustering in the British Isles



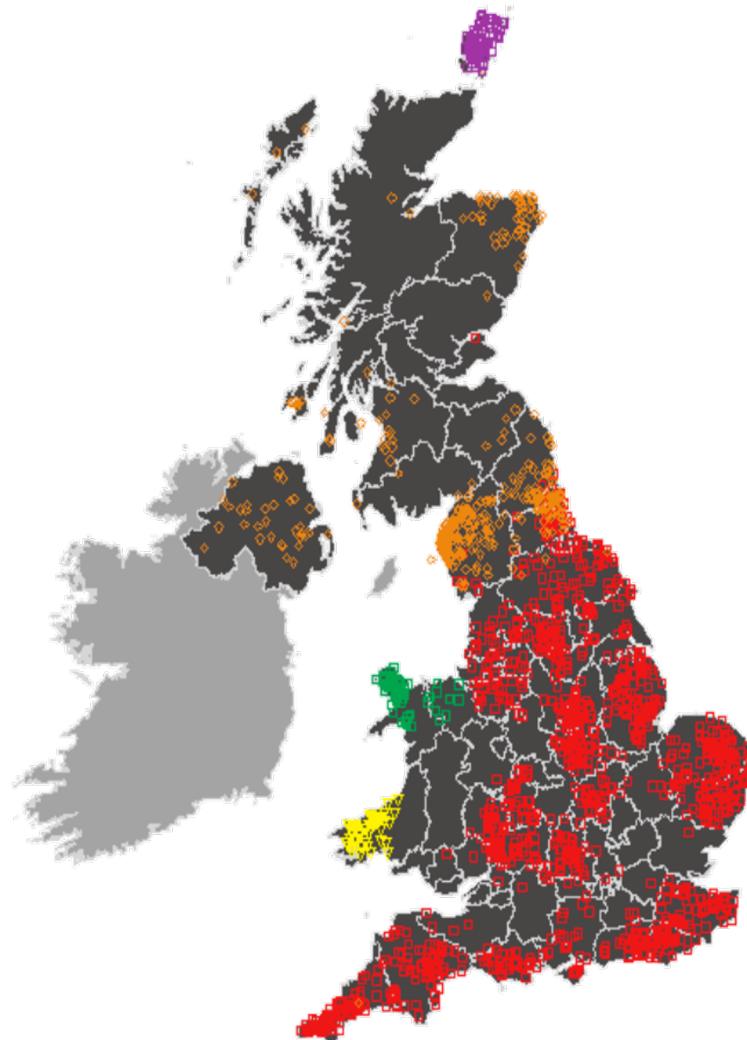
Haplotype clustering in the British Isles



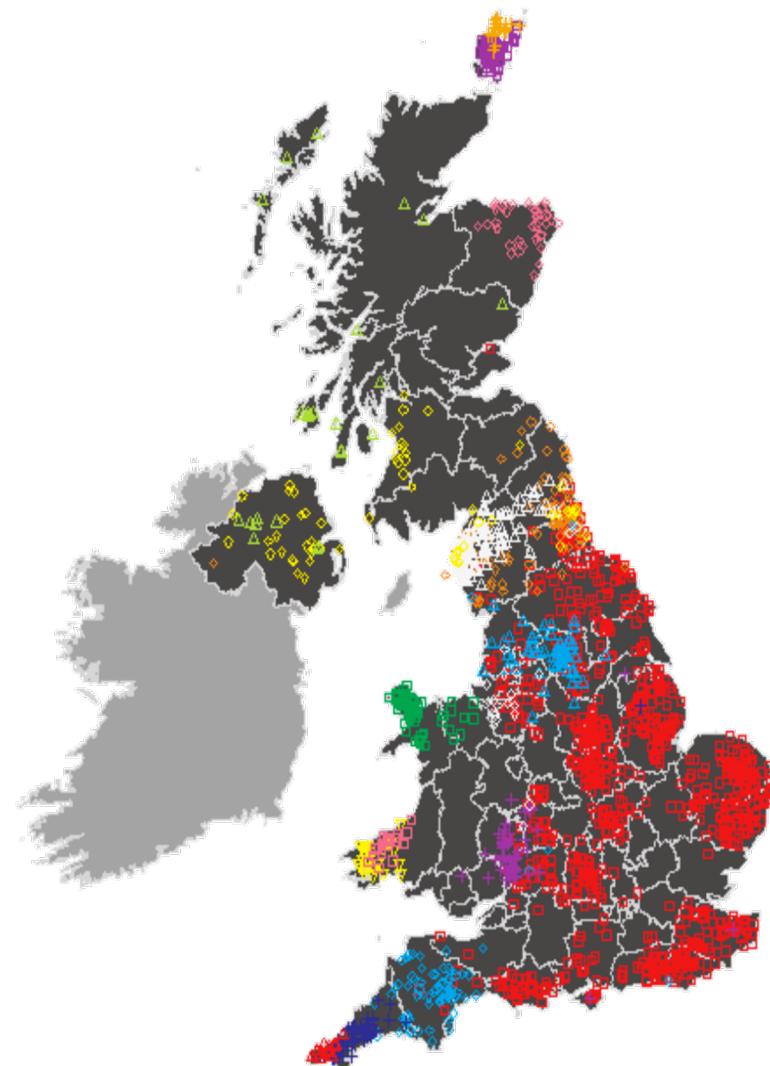
Haplotype clustering in the British Isles



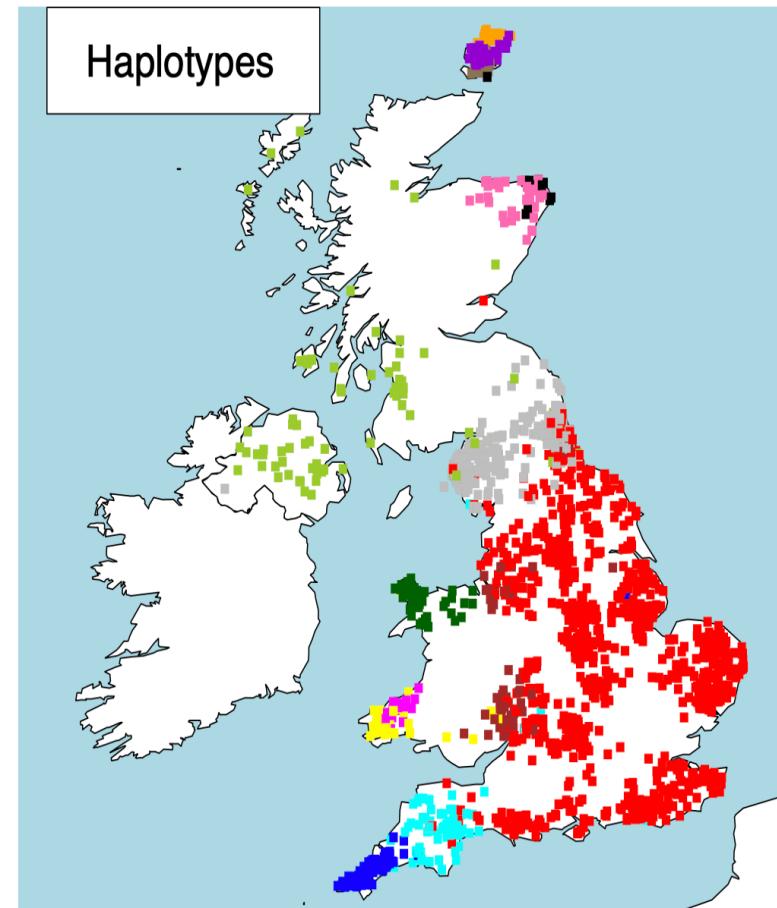
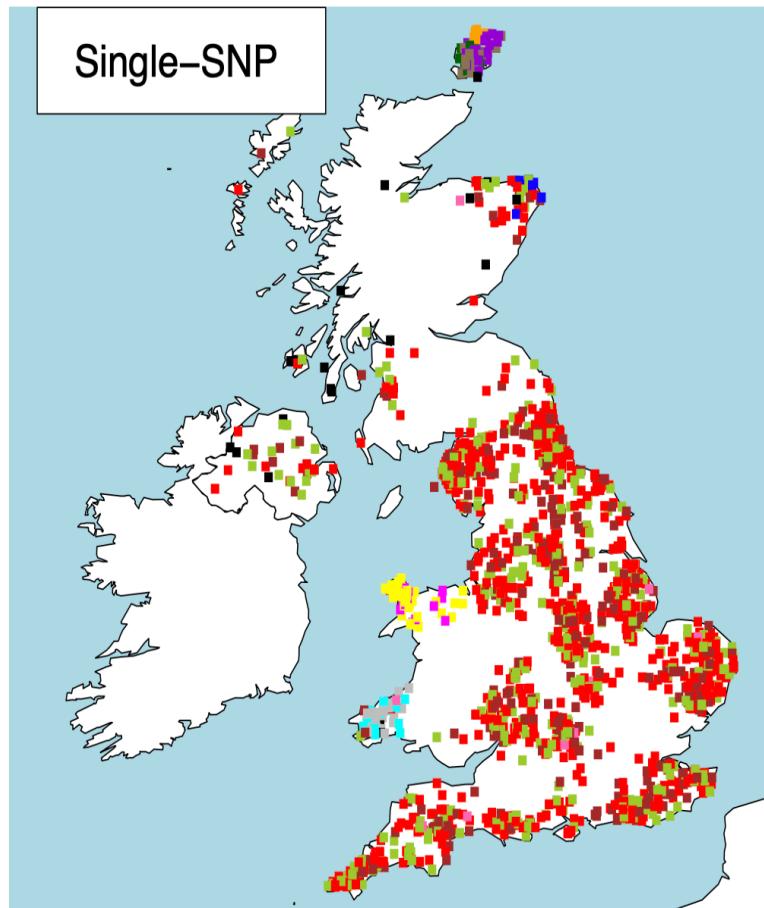
Haplotype clustering in the British Isles



Haplotype clustering in the British Isles



Haplotype clustering in the British Isles



Haplotype clustering - fineSTRUCTURE

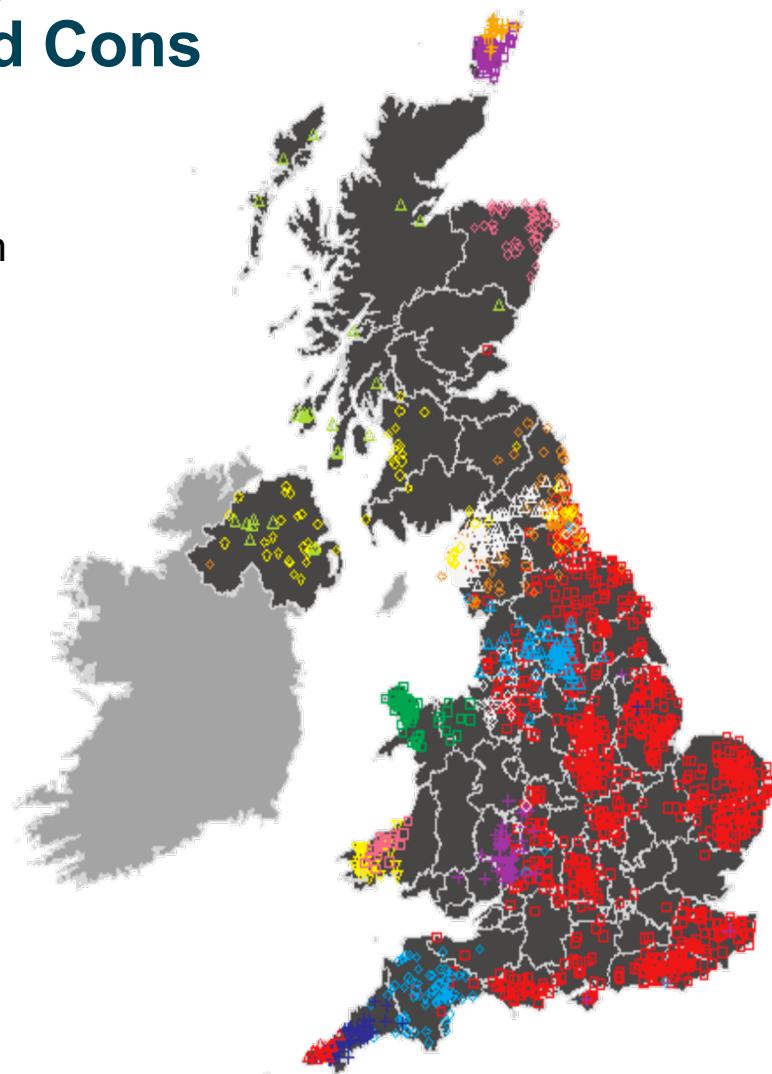
Pros and Cons

Advantages:

- Increased power to fine more subtle population structure
- Heatmaps demonstrate presence of structure/admixture
- Current implementation infers number of clusters K automatically and builds a tree

Disadvantages

- Challenges of interpretation – drift/admixture/other
- Requires phased data
- Computationally slow compared to ADMIXTURE

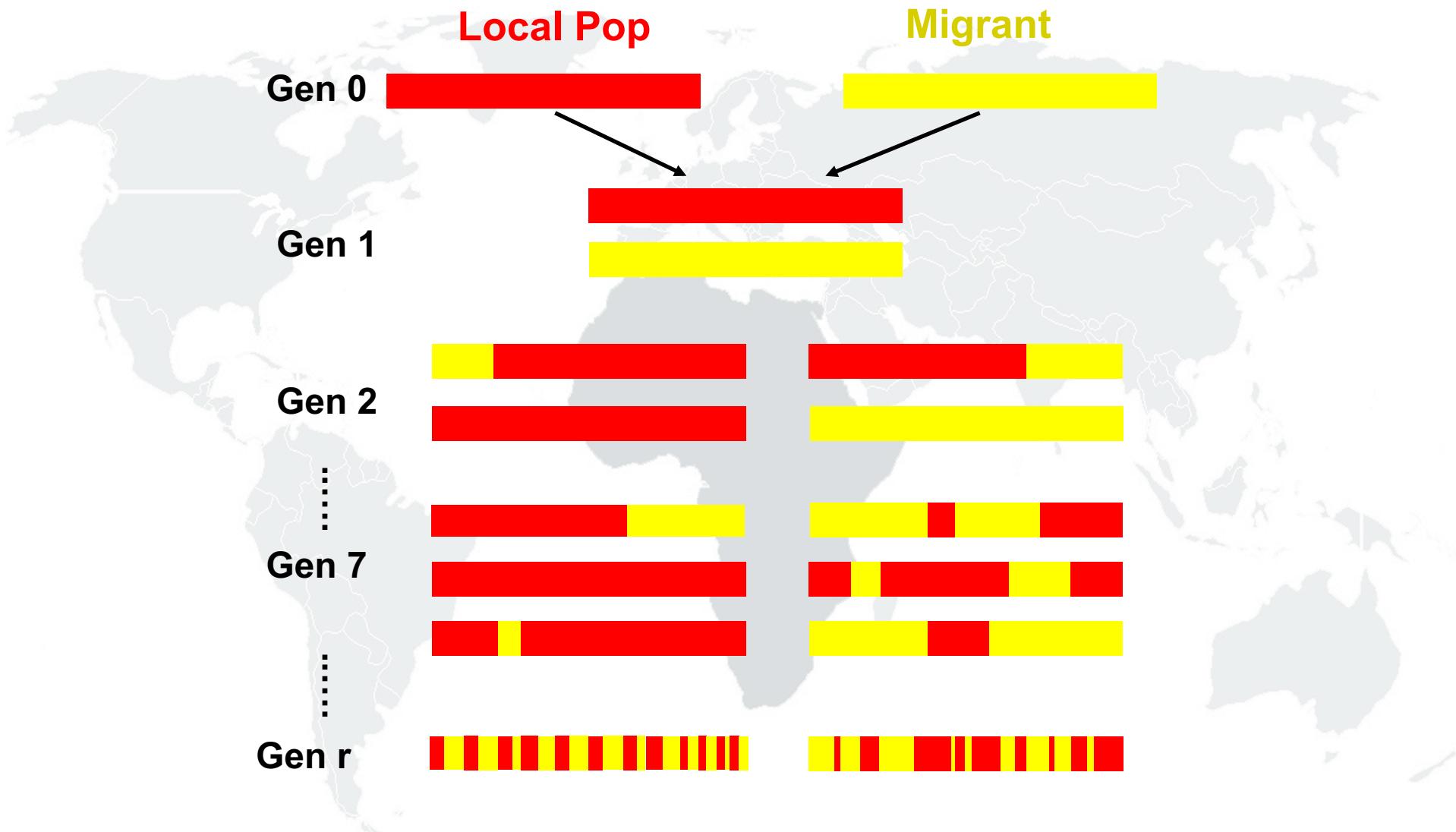


Haplotype sharing patterns can also be used to infer admixture events in human populations

- Admixture = mixing/interbreeding of two genetically distinguishable groups
- We can detect admixture in the genome by exploring patterns of decay of LD over genetic distance using chunks inferred from chromosome painting



Haplotype sharing patterns can also be used to infer admixture events in human populations



LD decays with genetic distance and this can be used for admixture inference

- Genetic recombinations are well modeled by a Poisson random variable (Falush et al, *Genetics*, 2003) → where the time between ‘arrivals’ of a Poisson process follows an exponential distribution.

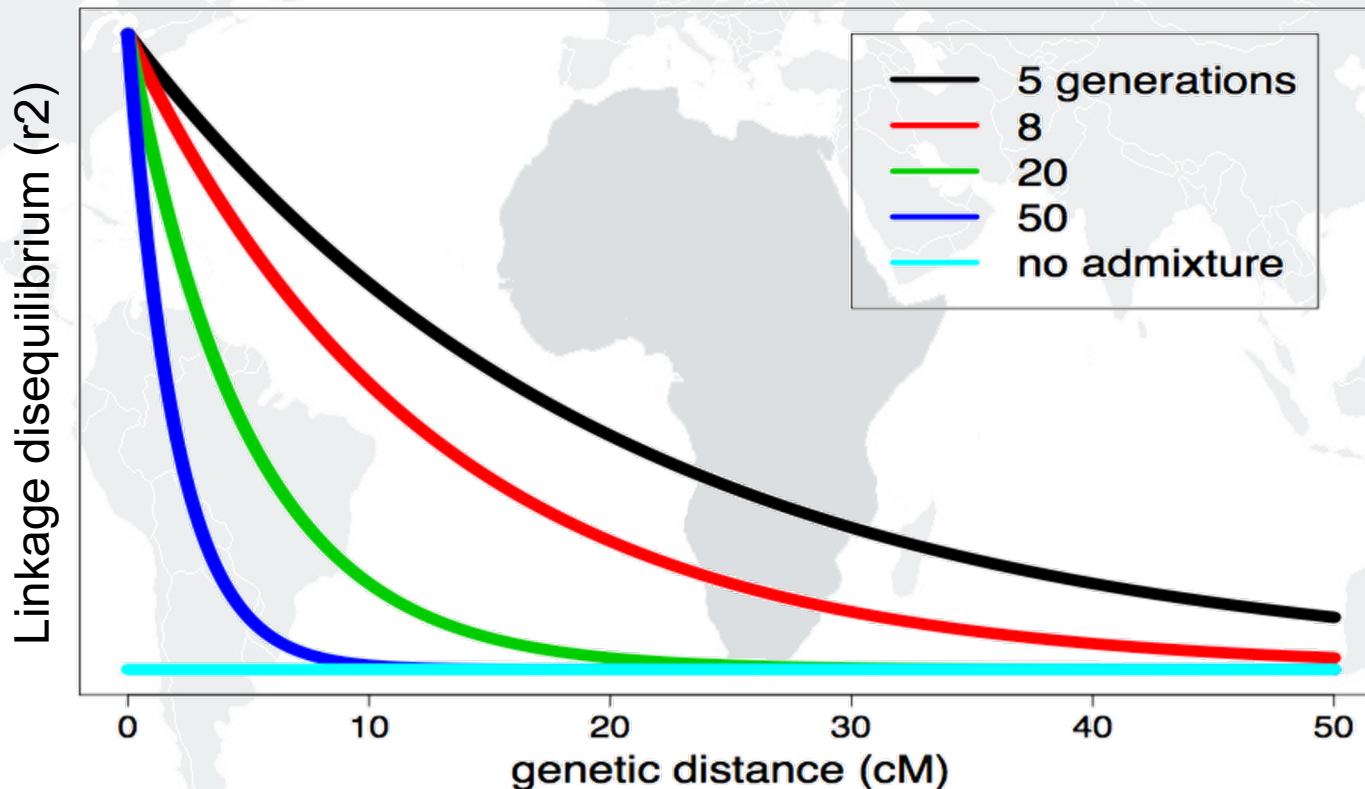


- The length of red and yellow depends on the number of recombination's that have occurred.
 - Eg. a continuous tract is one in which no ancestry-breaking recombinations have occurred.
 - Thus, tract length distribution should be exponential.
- Can we use this distribution to approximate the number of generations since admixture?

LD decays with genetic distance and this can be used for admixture inference

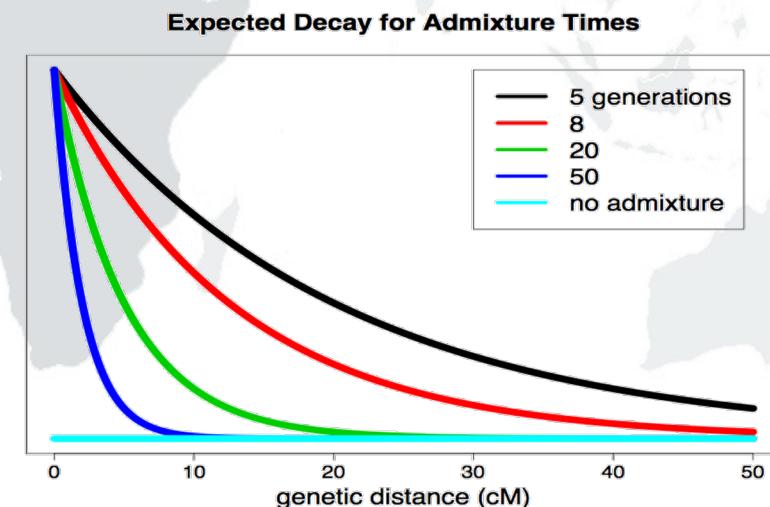
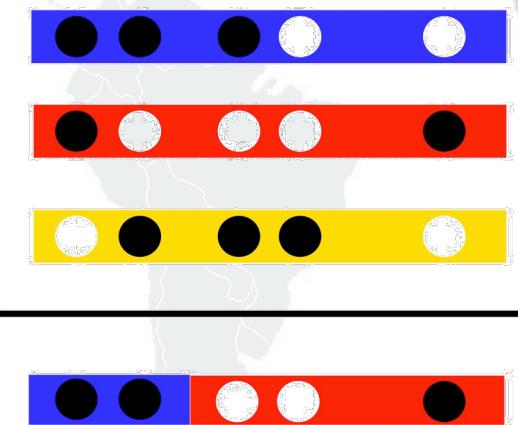


Expected Decay for Admixture Times

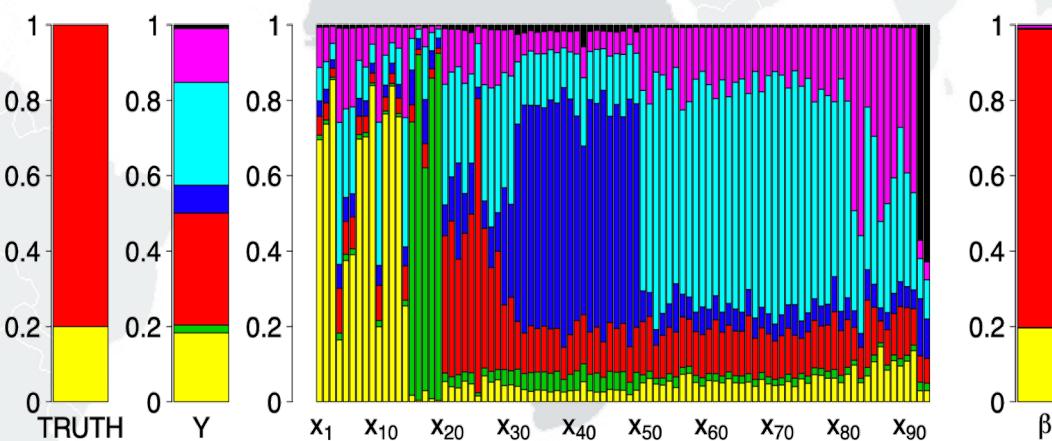
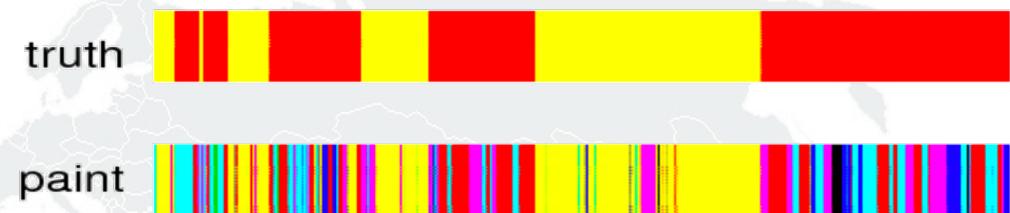
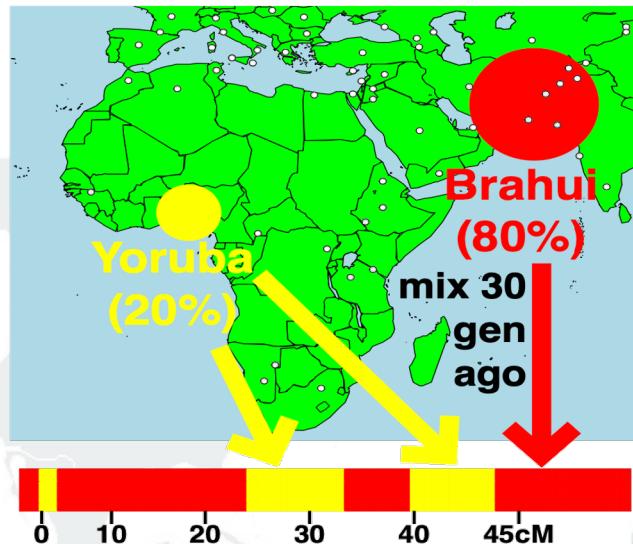


Inferring admixture events in human populations

- Several methods that use admixture LD to infer, identify and date admixture events eg. ROLLOFF (Moorjani et al, AJHG, 2013), ALDER (Loh et al, PLoS Genetics, 2013), and GLOBETROTTER (Hellenthal et al, Science, 2014).
- DNA → chromosome painting (infer which groups) → size of segments → fit exponentials



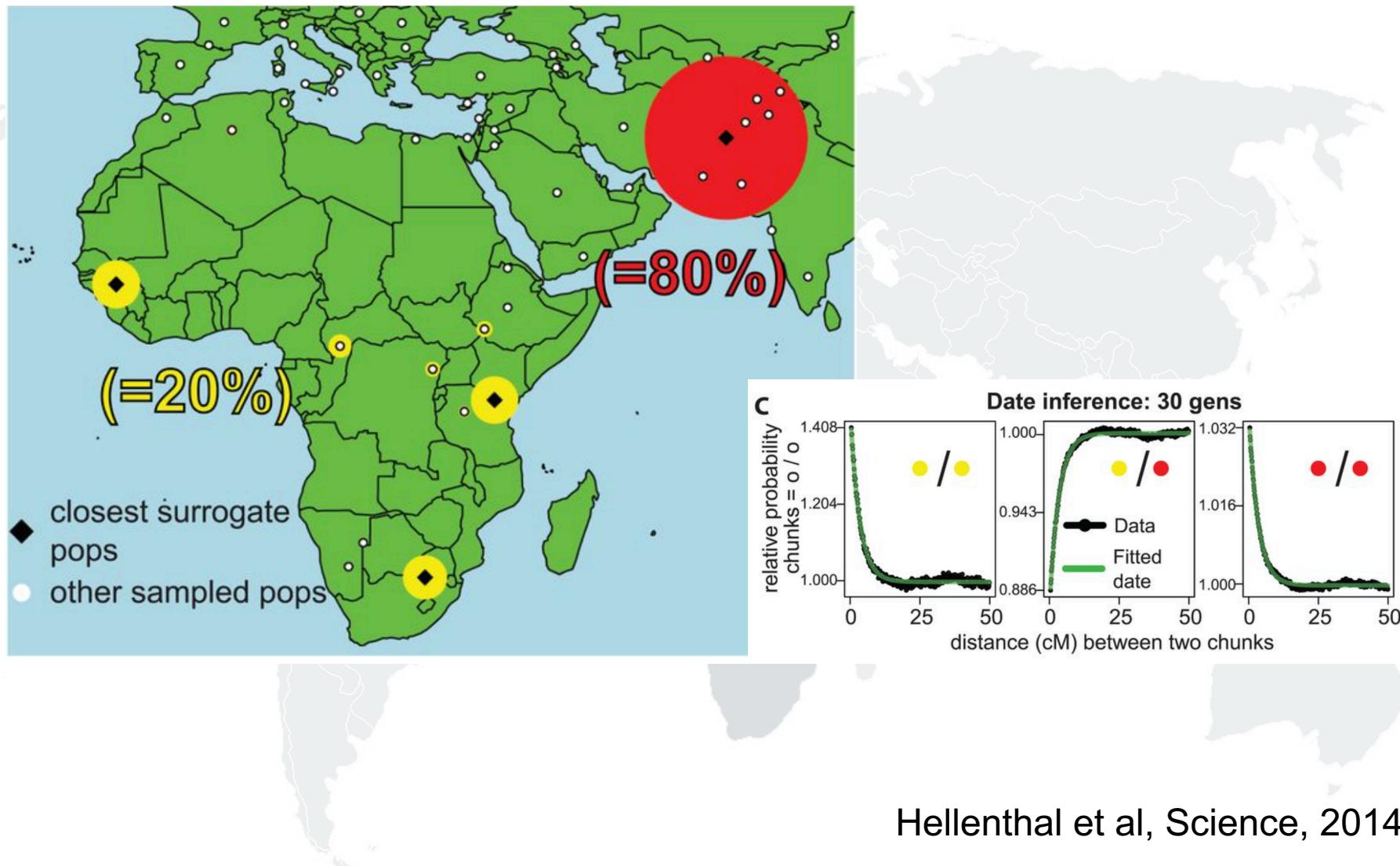
Inferring admixture events - Globetrotter



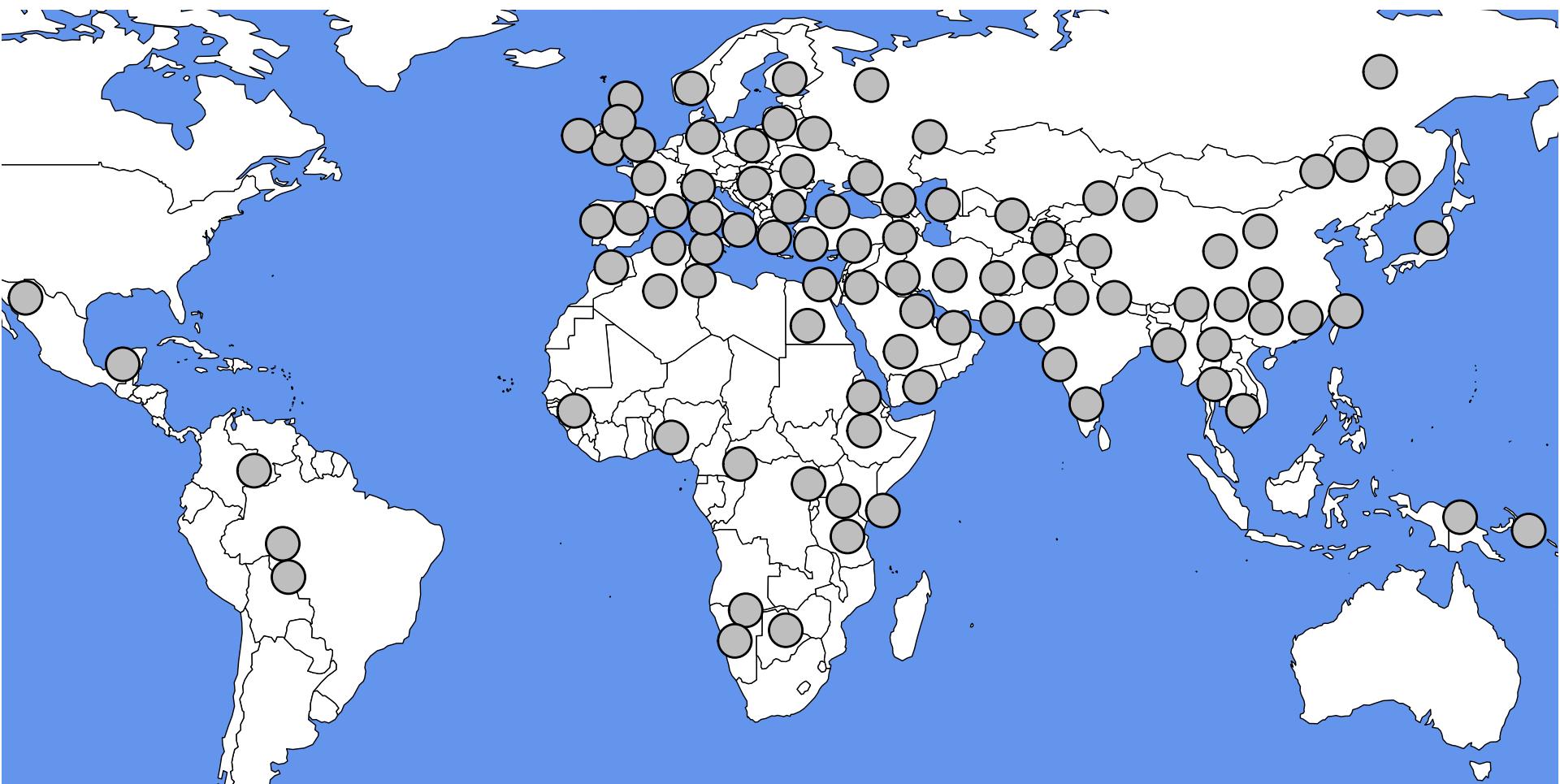
$$E[Y] = \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_{93} X_{93}$$

Hellenthal et al, Science, 2014

Inferring admixture events - Globetrotter

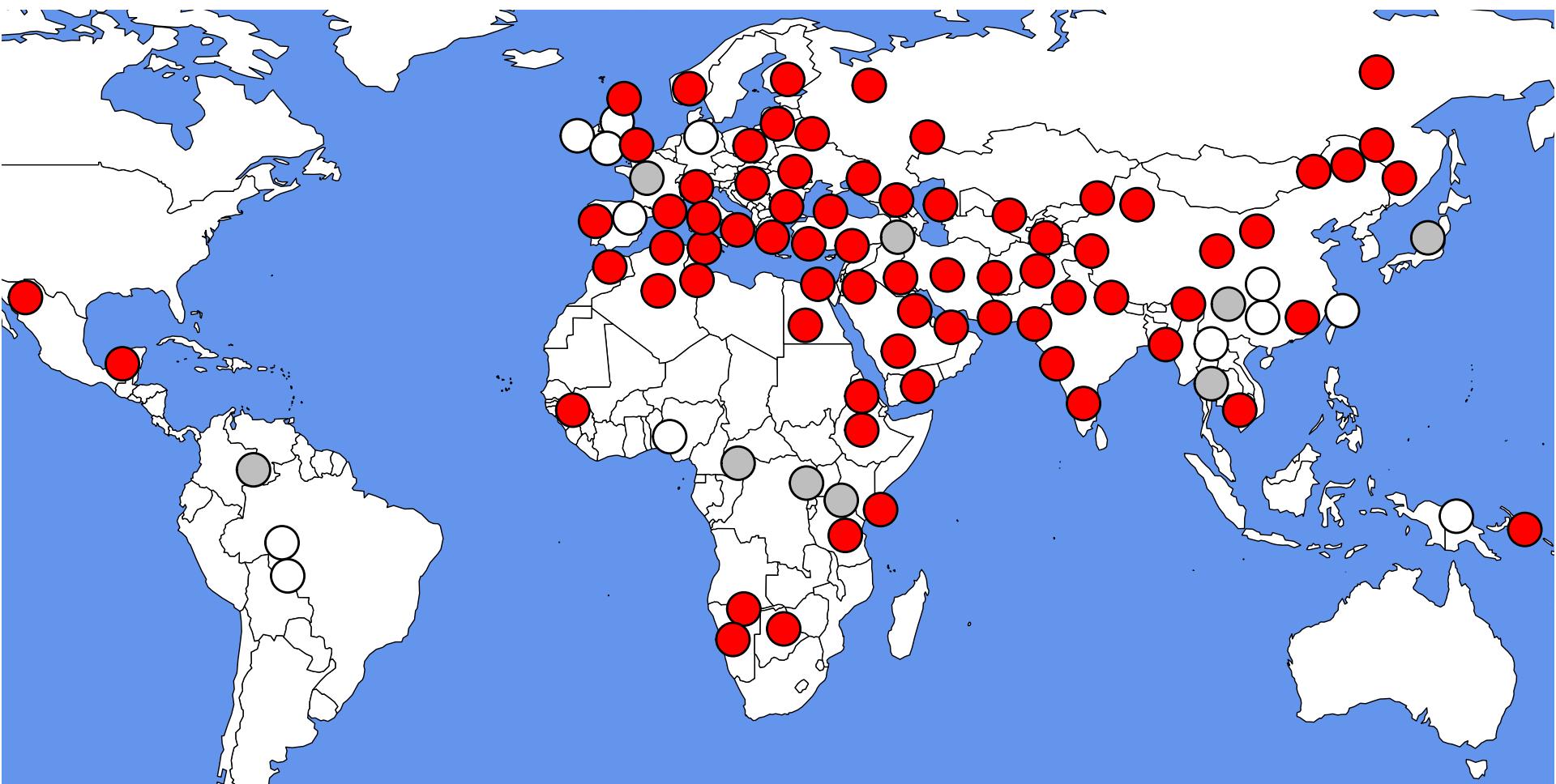


Genetic mixing over the past 4000 years



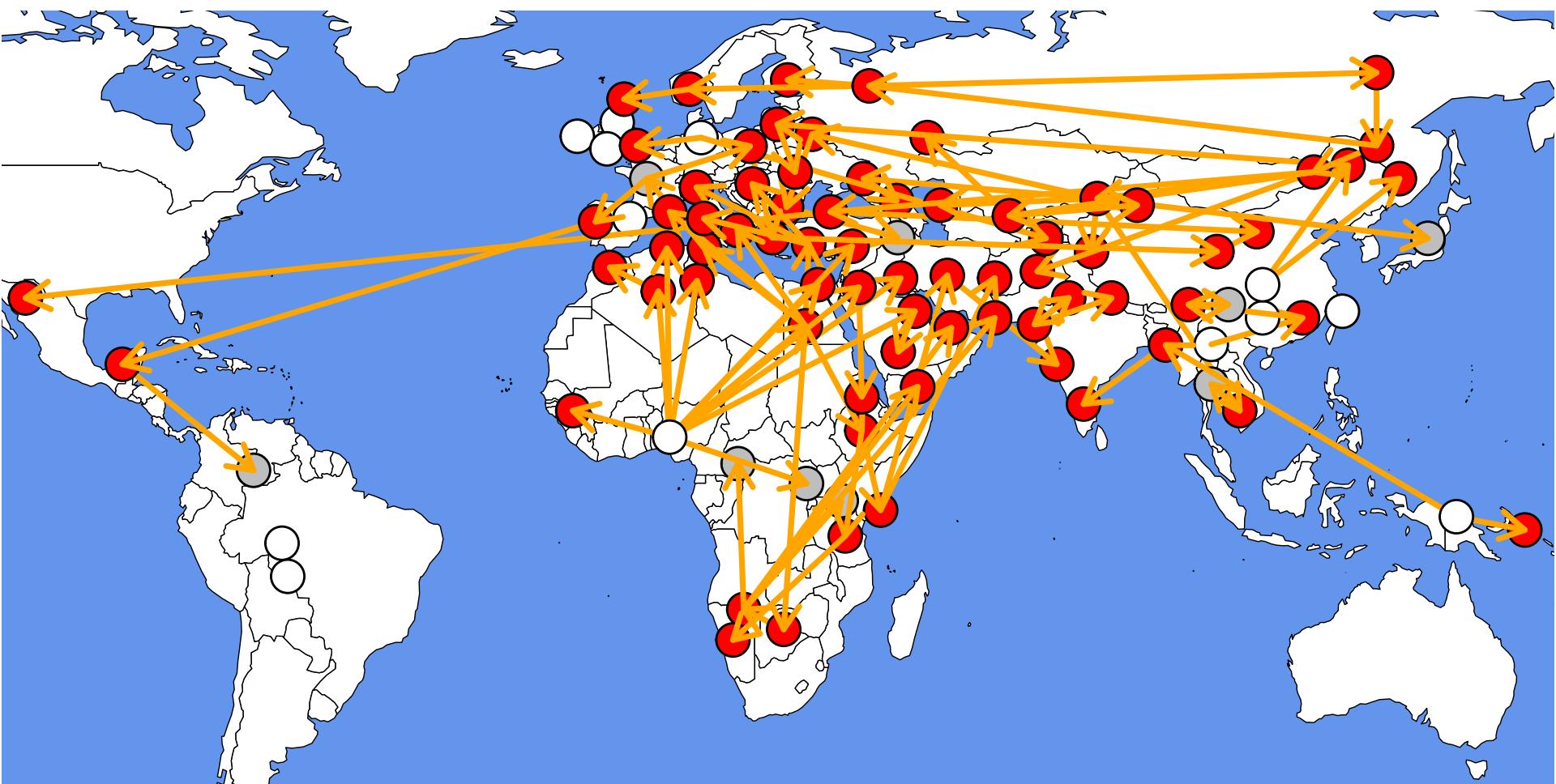
Hellenthal et al, *Science*, 2014
Image credit G Hellenthal

Genetic mixing over the past 4000 years



Hellenthal et al, *Science*, 2014
Image credit G Hellenthal

Genetic mixing over the past 4000 years



Hellenthal et al, *Science*, 2014
Image credit G Hellenthal

Summary

Many programs for inferring and determining population structure – central to understanding demographic history

PCA and model-based clustering make use of unlinked allele frequencies.
Fast and tractable, though interpretation can be challenging.

Haplotype-based methods – a toolkit for exploring local ancestry, LD, recombination and subsequently fine-scale patterns of population structure at different scales.

Chromopainter and fineSTRUCTURE as tools for resolving subtle differences between populations.

Admixture LD can be utilized together with chromosome painting to identify when admixture has occurred, its sources and its timings.

Such analyses are highlighting that admixture is common place within human populations over the past 4000 years.

Suggested Reading/Software

Overview of methods in PopGen:

Schraiber & Akey. Methods and models for unravelling human evolutionary history. *Nature Reviews Genetics*. (2015) 6(12):727-40.

How not to overinterpret ADMIXTURE plots:

Lawson, van Dorp & Falush. A tutorial on how not to over-interpret STRUCTURE and ADMIXTURE bar plots. *Nature Communications*. (2018) 9,3258.

Chromopainter/fineSTRUCTURE:

Lawson, Hellenthal, Myers & Falush. Inference of Population Structure using Dense Haplotype Data. *PLoS Genetics*. (2012) 8(1): e1002453.

https://people.maths.bris.ac.uk/~madjl/finestructure-old/chromopainter_info.html

Fine-scale structure of the British Isles:

Lesie, Winney & Hellenthal et al. The fine-scale genetic structure of the British population. *Nature*. (201) 519,309-314.

GLOBETROTTER:

Hellenthal et al. A genetic atlas of human admixture history. *Science*. 14; 343(6712):747-751.

<https://people.maths.bris.ac.uk/~madjl/finestructure/globetrotter.html>

Thank you to Garrett Hellenthal!