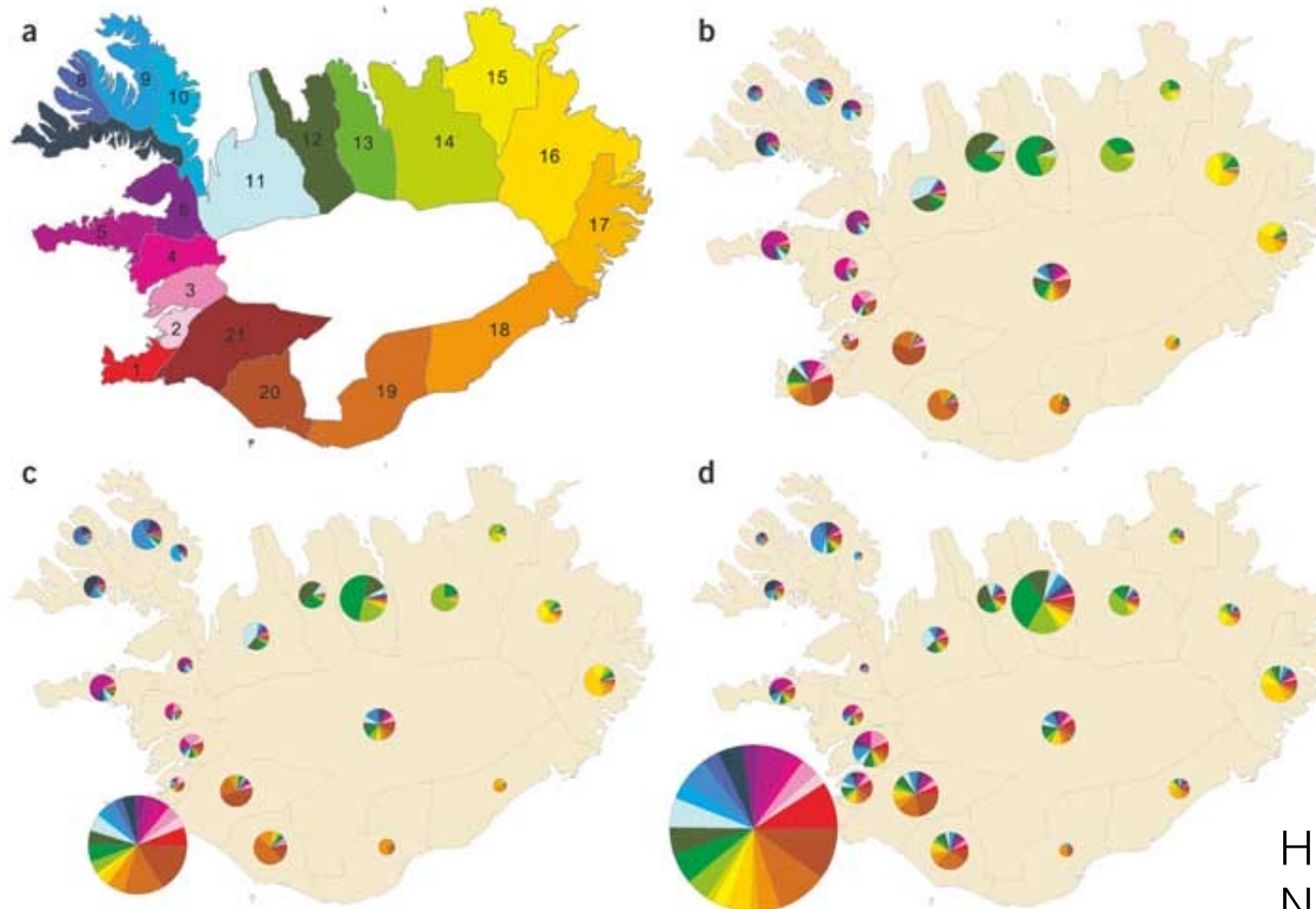# Cryptic Relatedness and fine-scale population structure
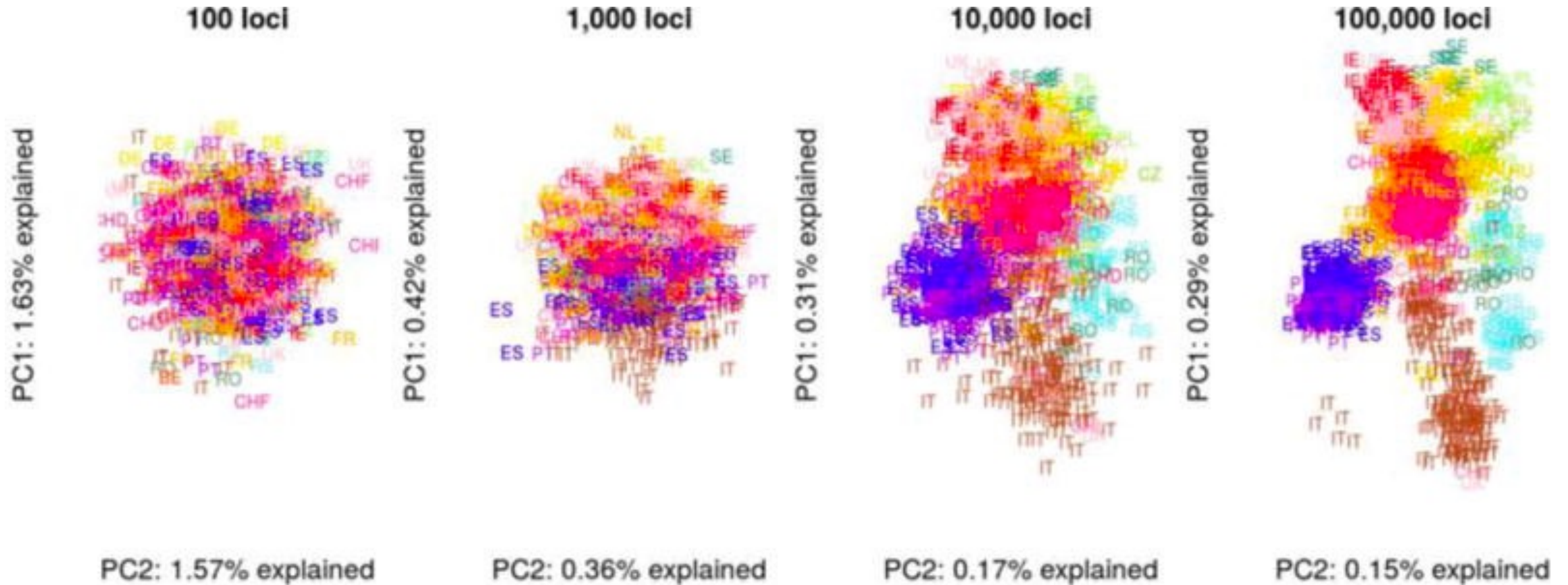
# Learning objectives

- Define fine-scale population structure and cryptic relatedness
- How is it identified
  - Identity-by-descent
  - Rare variation
  - Estimated Effective Migration Surfaces
- Why it can be important for association analyses, especially of rare variants.
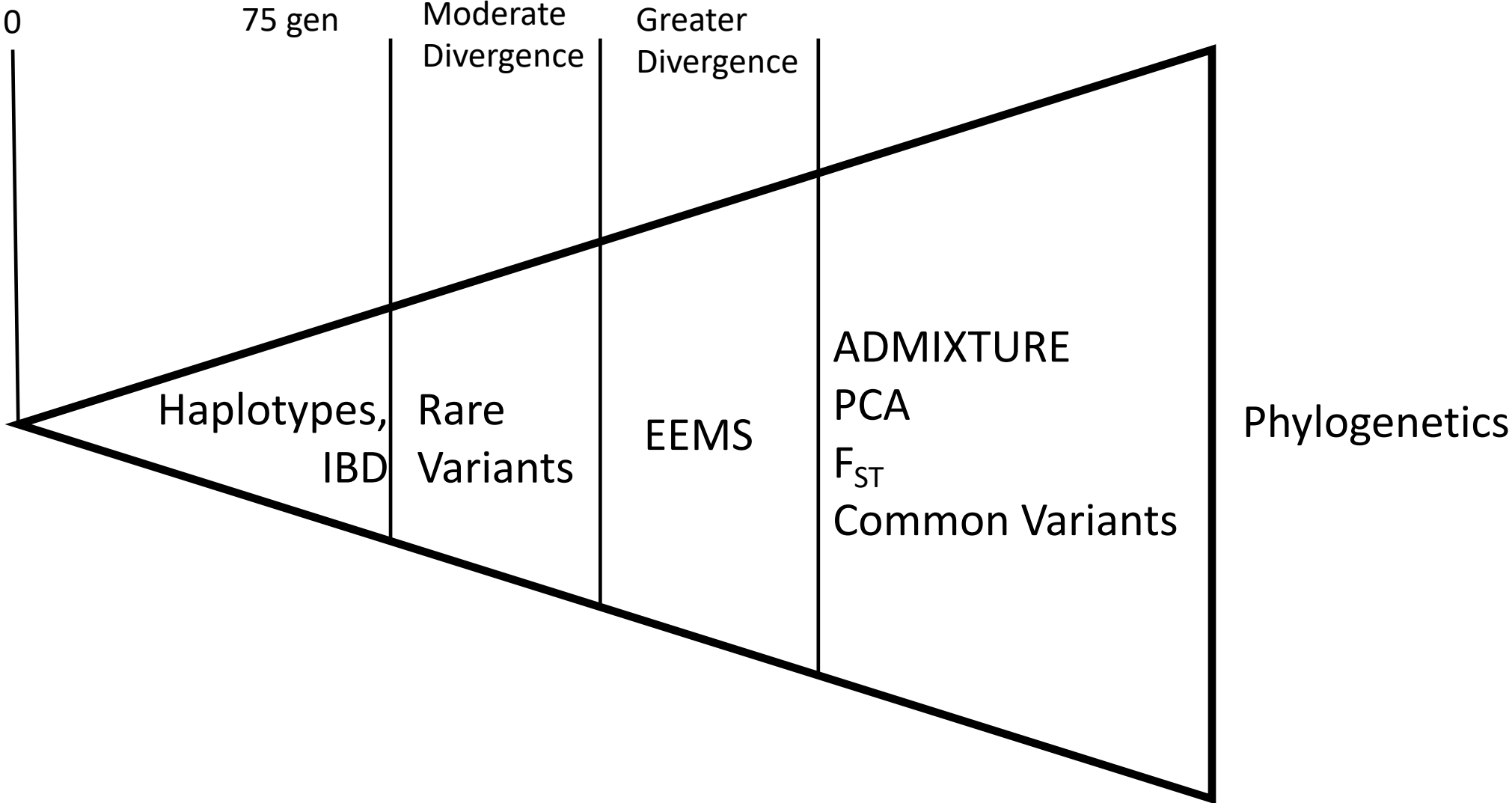
# Cryptic Population Structure


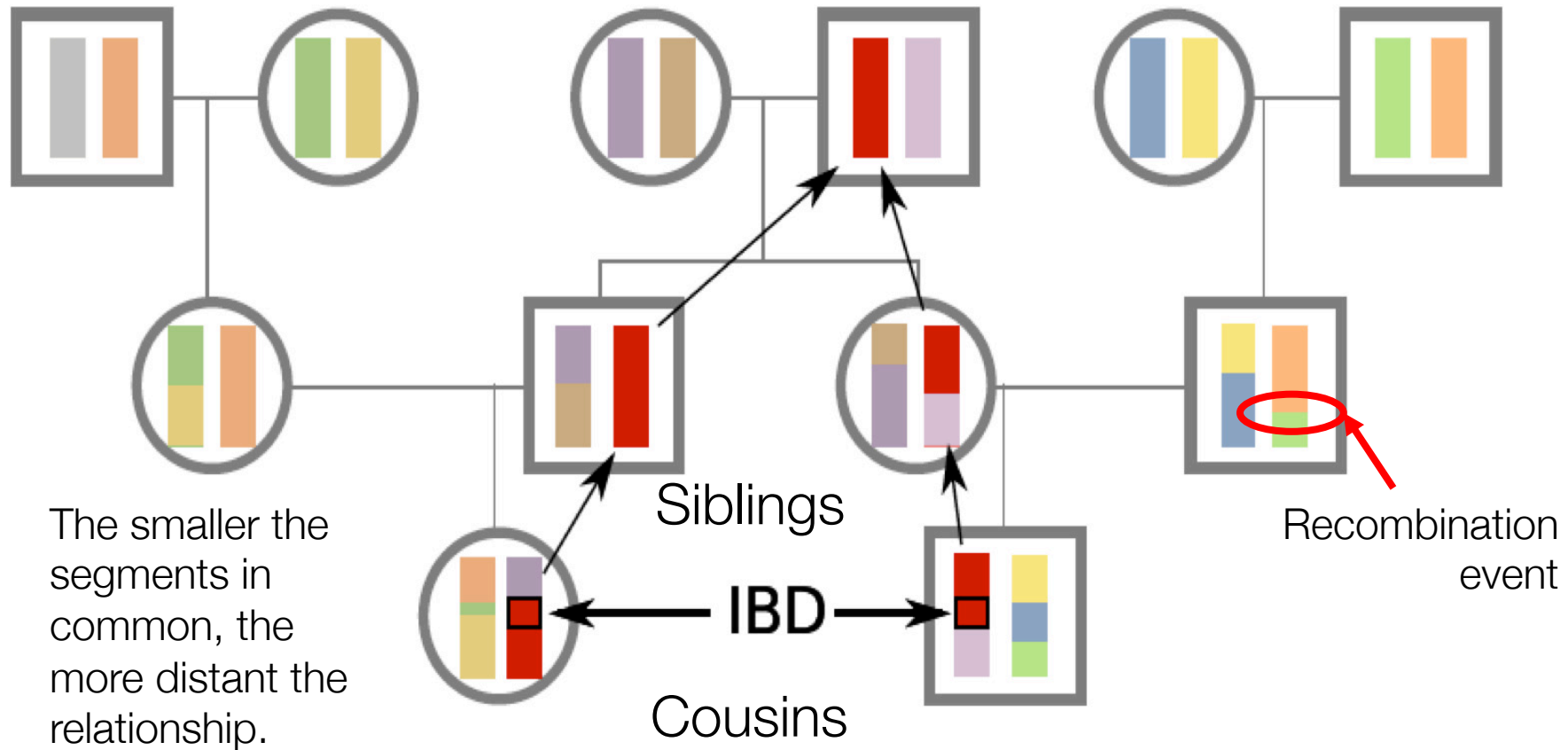
Helgason et al. (2004)
Nature Genet.

# "A large number of loci is required to reveal fine-scale population structure using PCA"
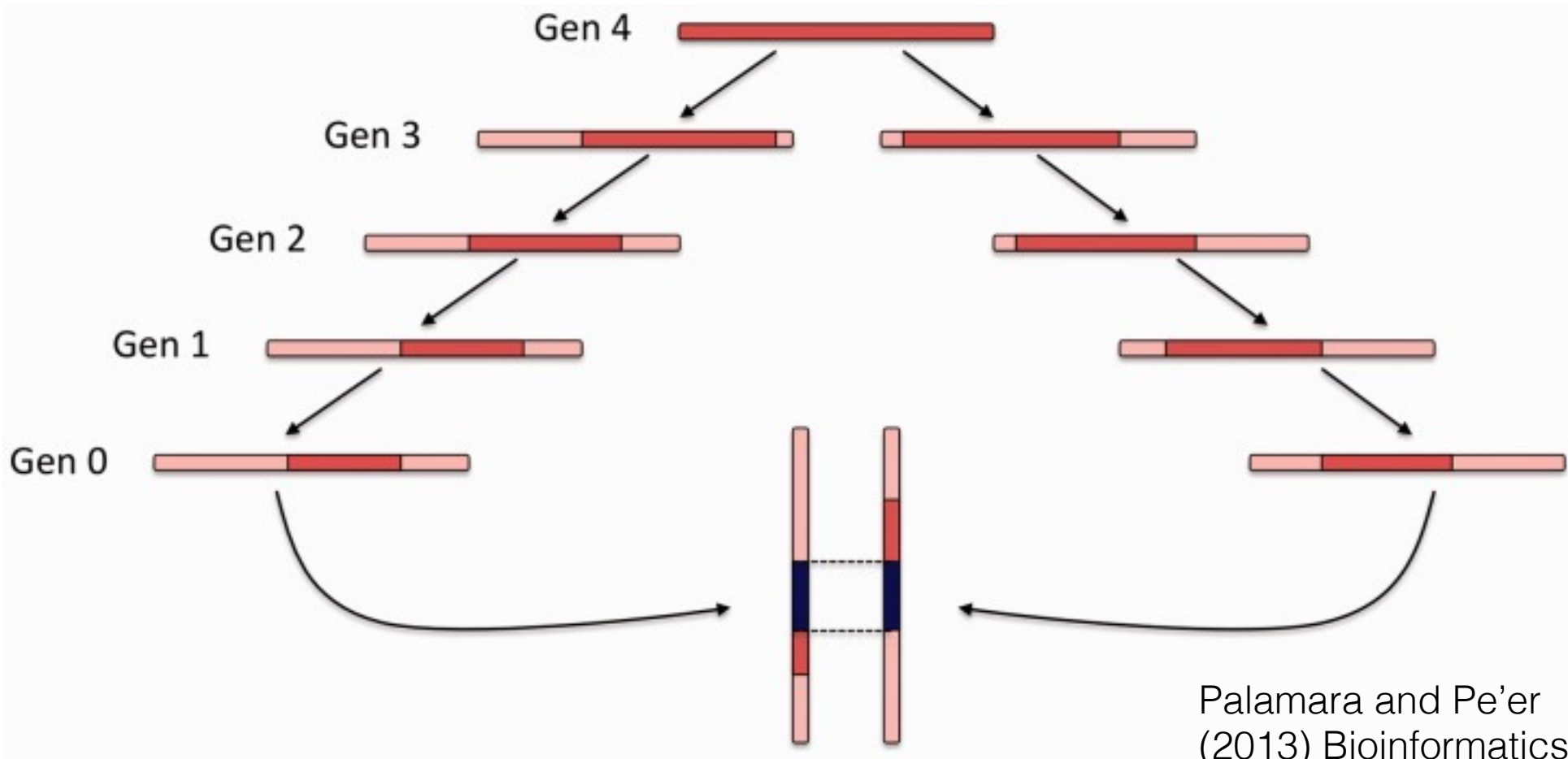


Novembre & Peter (2017) Curr Opin Genet Dev.

# Demographic Inference Time Frames

# Identity by Decent (IBD): A method to find both distant and recent relationships



The smaller the segments in common, the more distant the relationship.

Siblings

IBD

Cousins

Recombination event

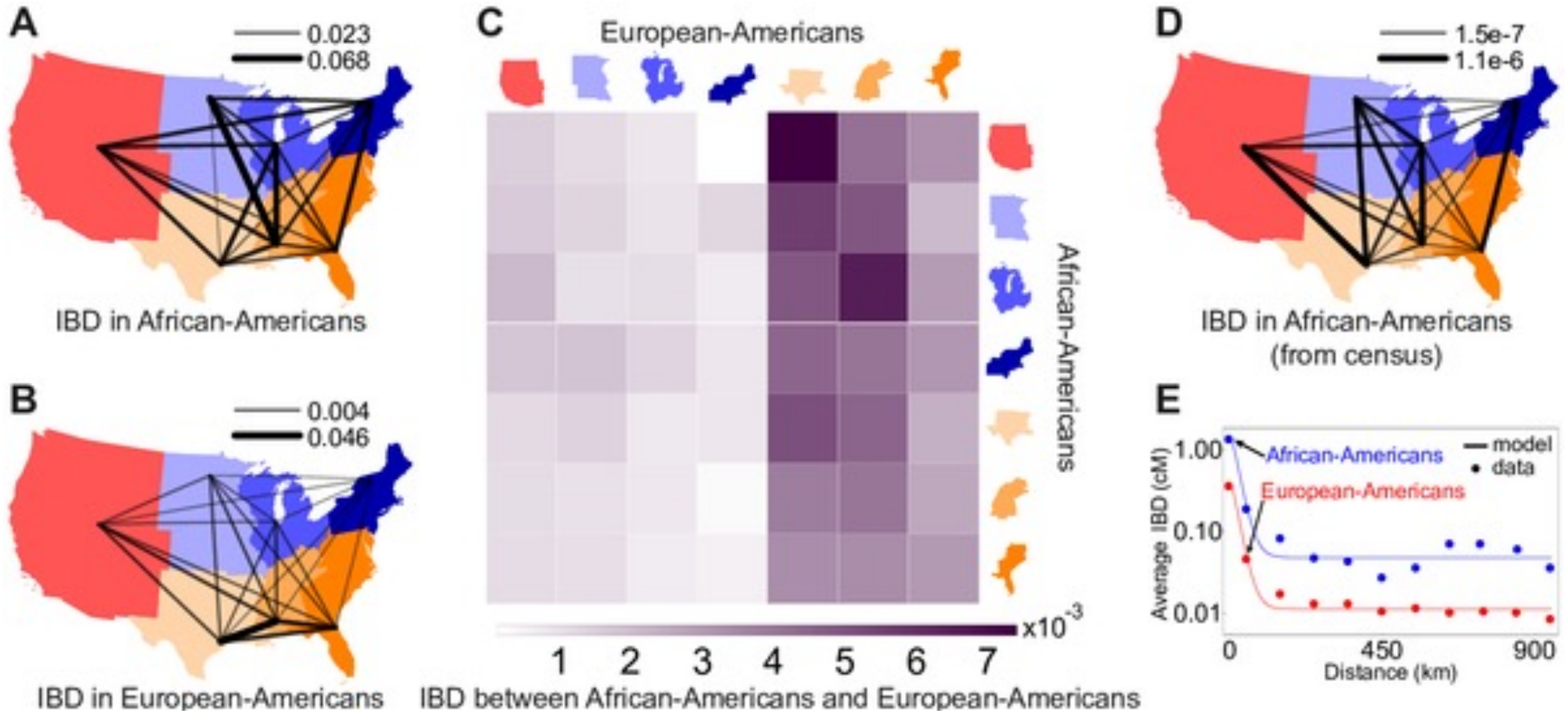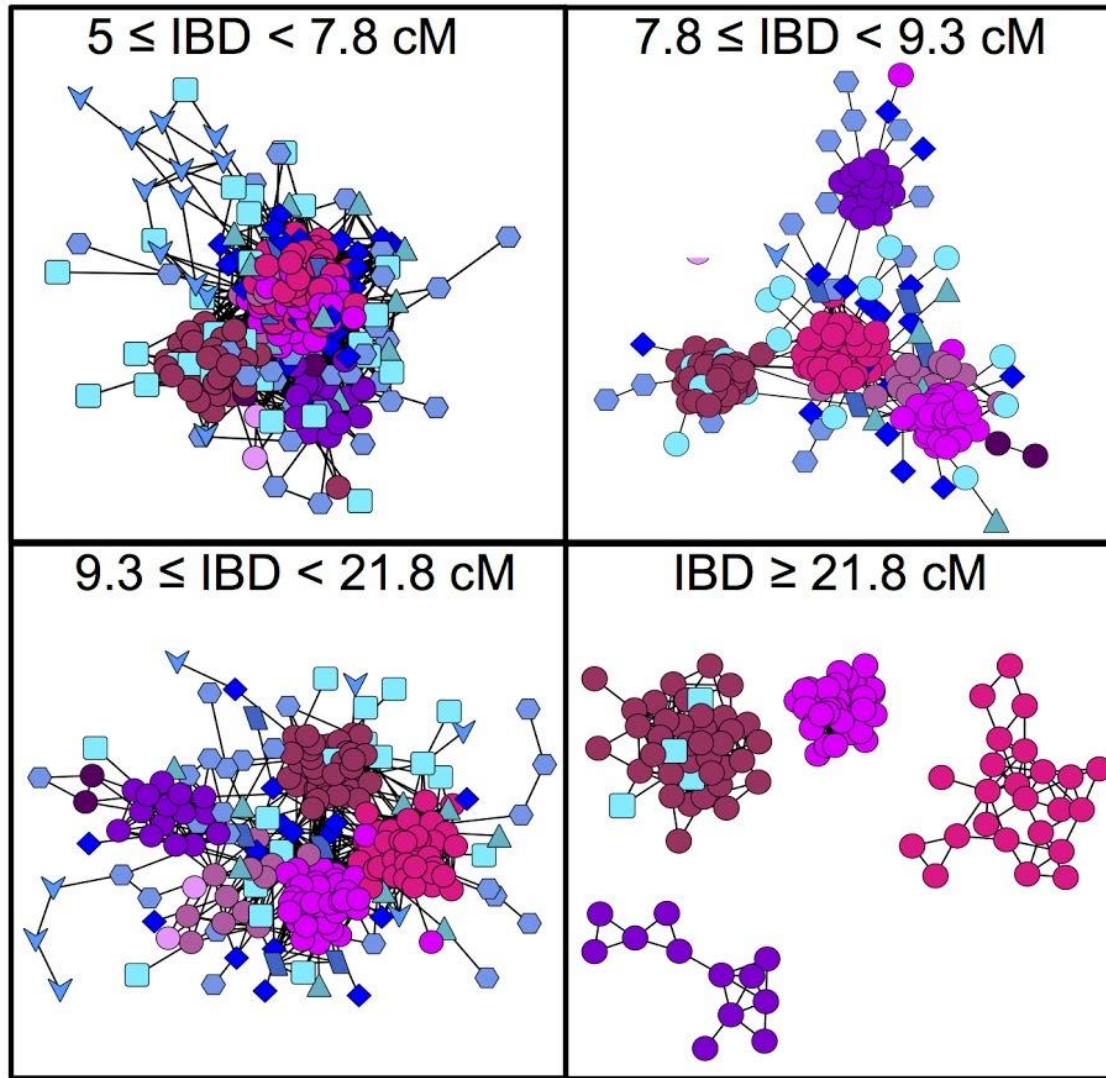# IBD length is correlated with historical relationships.



$$E[g|l] \cong \frac{3}{2 * l}$$

Baharian et al. (2016)
PLoS Genet.

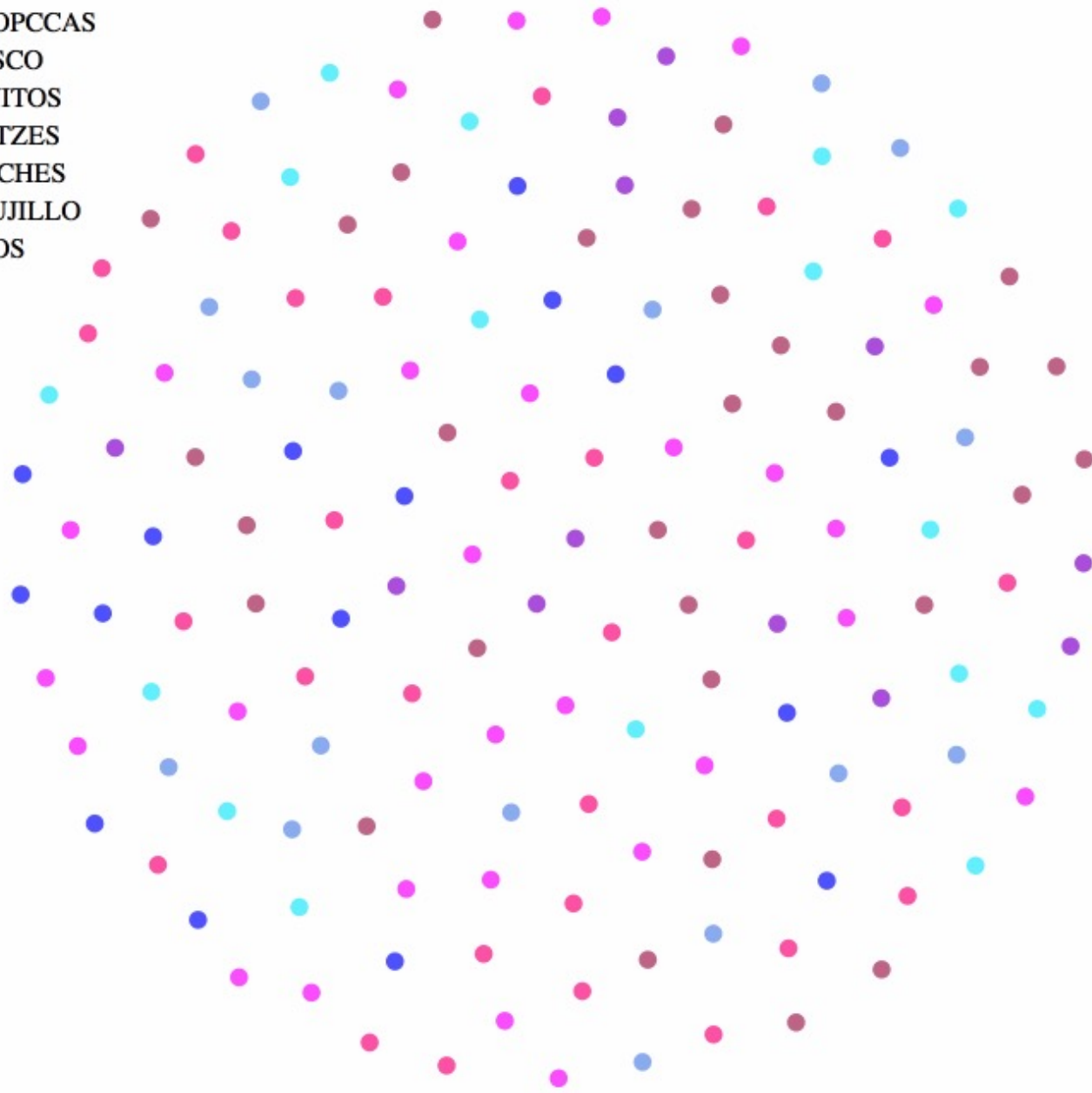Palamara and Pe'er
(2013) Bioinformatics

# Pairwise genetic relatedness across



**A** — IBD in African-Americans (0.023 / 0.068)

**B** — IBD in European-Americans (0.004 / 0.046)

**C** — European-Americans / African-Americans
IBD between African-Americans and European-Americans (1 2 3 4 5 6 7, ×10⁻³)

**D** — IBD in African-Americans (from census) (1.5e-7 / 1.1e-6)

**E** — Average IBD (cM) vs Distance (km); African-Americans, European-Americans; model, data

Baharian et al. (2016) PLoS Genet.

**C**

5 ≤ IBD < 7.8 cM

7.8 ≤ IBD < 9.3 cM

9.3 ≤ IBD < 21.8 cM

IBD ≥ 21.8 cM

Trujillo ∨ AP ● Chopccas ● Moches ● Qeros
▲ Lima ▮ Puno ● Matsig ● Nahua ● Uros
⬣ Iquitos ◆ Cusco ● Matzes

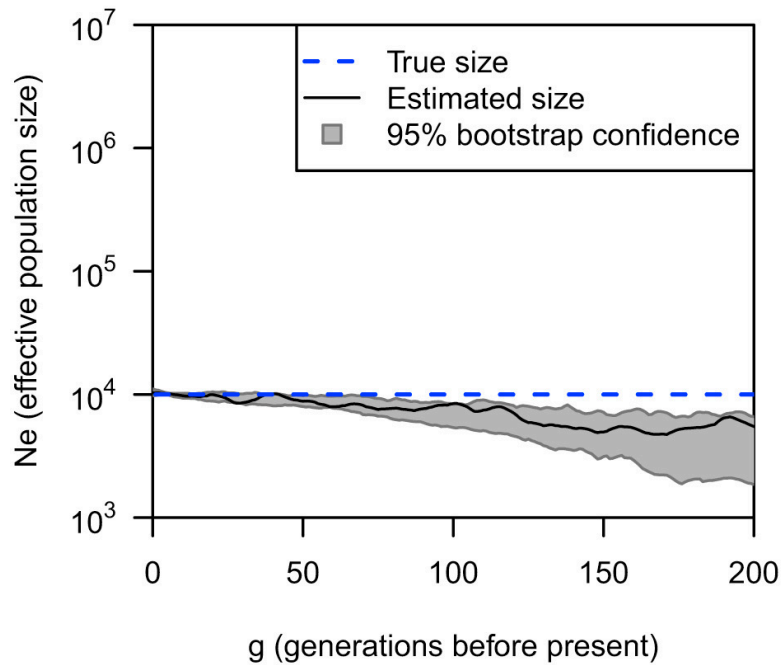Identity-by-descent as a means to look at fine-scale structure over time

Harris et al. (2018) PNAS

CHOPCCAS
CUSCO
IQUITOS
MATZES
MOCHES
TRUJILLO
UROS

Identity-by-descent as a means to look at fine-scale structure over time

Harris et al. (2018) PNAS

# IBD can estimate effective population size over time.



Constant size: SNP array data · Exponential growth: SNP array data · Super−exponential: SNP array data
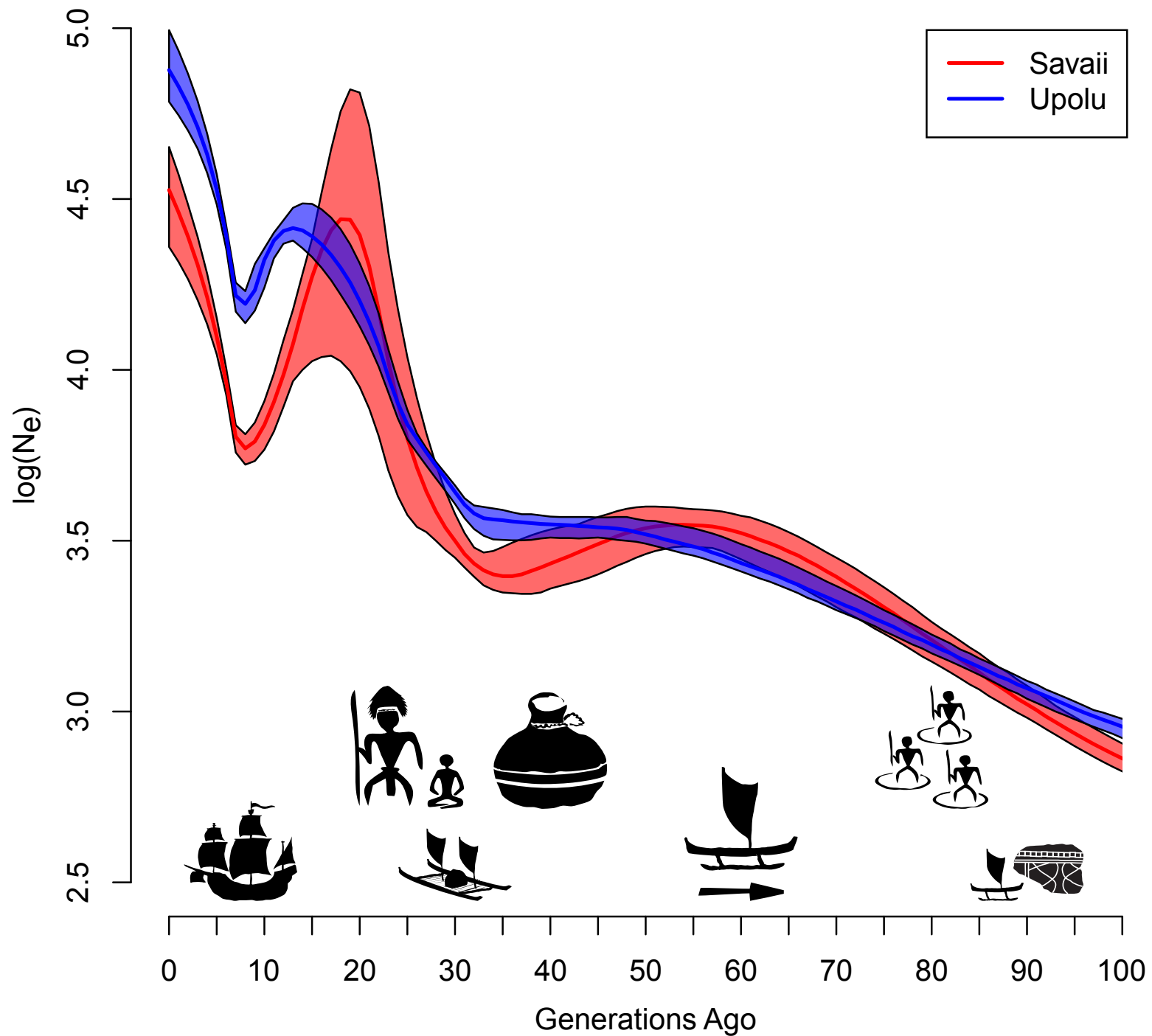
Browning &Browning (2016) Am. J. Human Genetics

IBDNe in Samoa!

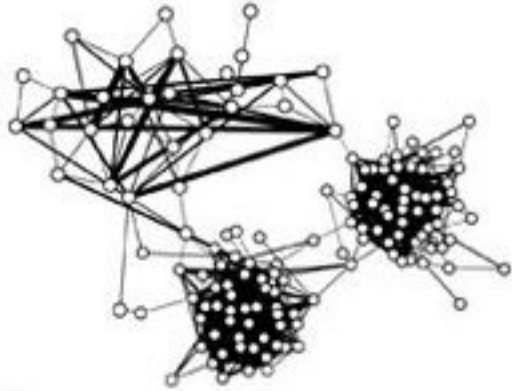Harris et al. (https://papers.ssrn.com/sol3/papers.cfm?abstract_id=3329885)

# IBD on a large scale
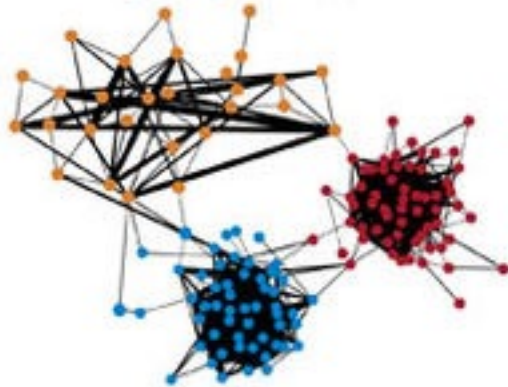


**a** Construct network from IBD.
Join vertex pairs (genotyped samples) if IBD>12 cM.
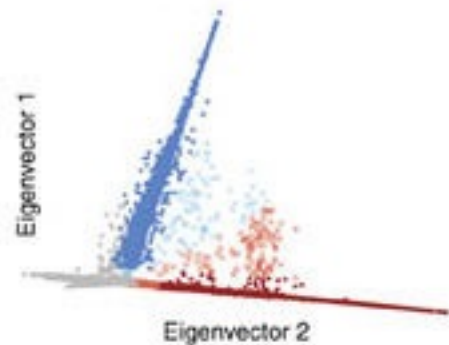Edge weights are a function of total detected IBD.

**b** Detect network clusters.
Recursively identify disjoint sets that maximize the modularity of the network. (Here one level of clustering hierarchy is shown.)
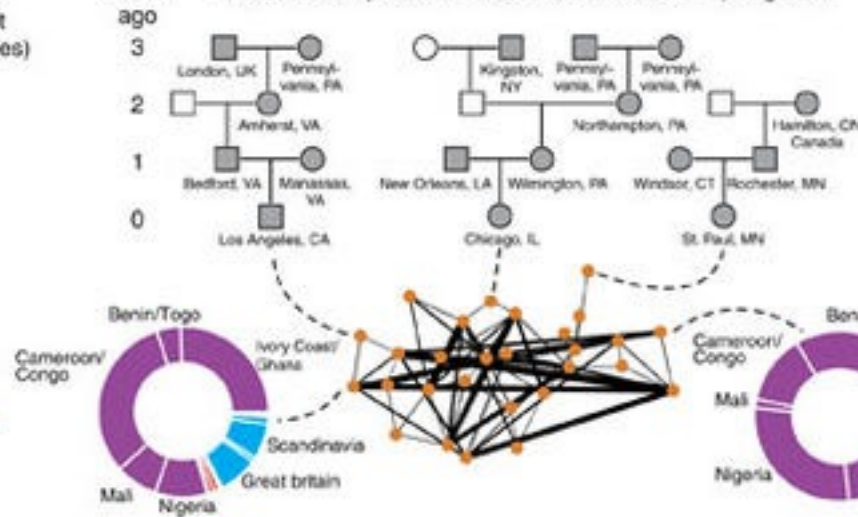
**c** Identify subsets of the clusters that separate in the spectral embedding.
Spectral embedding is computed from eigen-decomposition of Laplacian matrix. In the plot below, we identify "stable subsets" (filled circles) of the blue and red clusters.

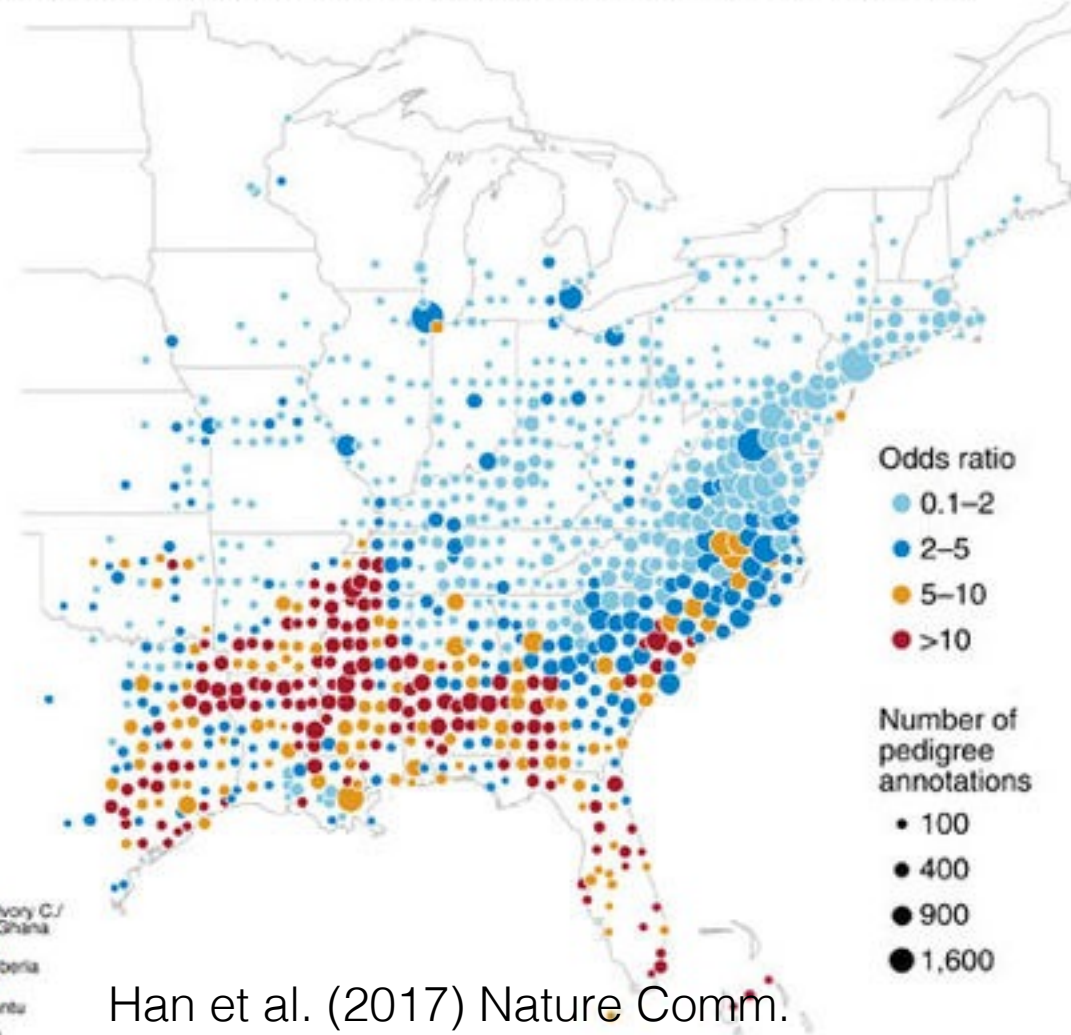**d** Annotate each cluster with two kinds of data:
• In all samples, global admixture of 20 populations (donut charts);
• For some samples, birth locations of ancestors in pedigrees.
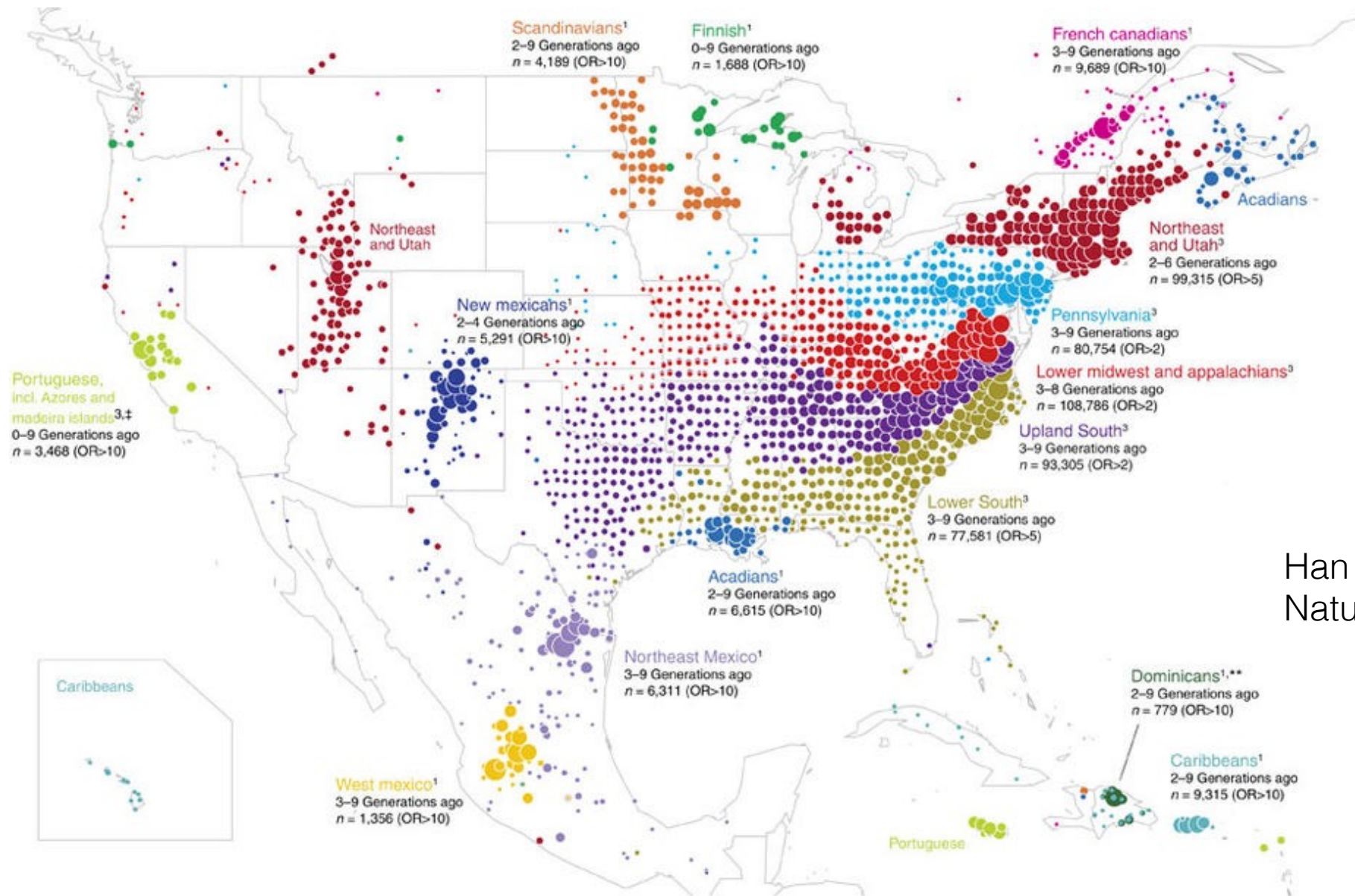
**e** Visualize geographic distribution of ancestral birth locations in each cluster.
Map below shows birth locations of ancestors in the African American cluster. Locations are colored by degree of over-representation (odds ratio), and scaled by number of birth location annotations.
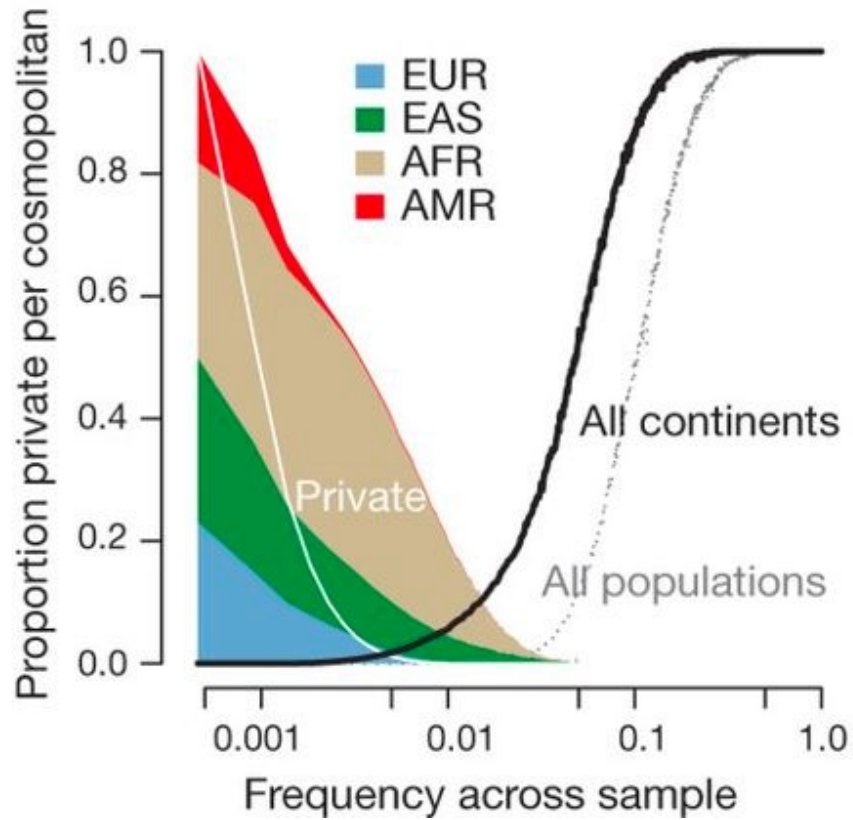
Han et al. (2017) Nature Comm.
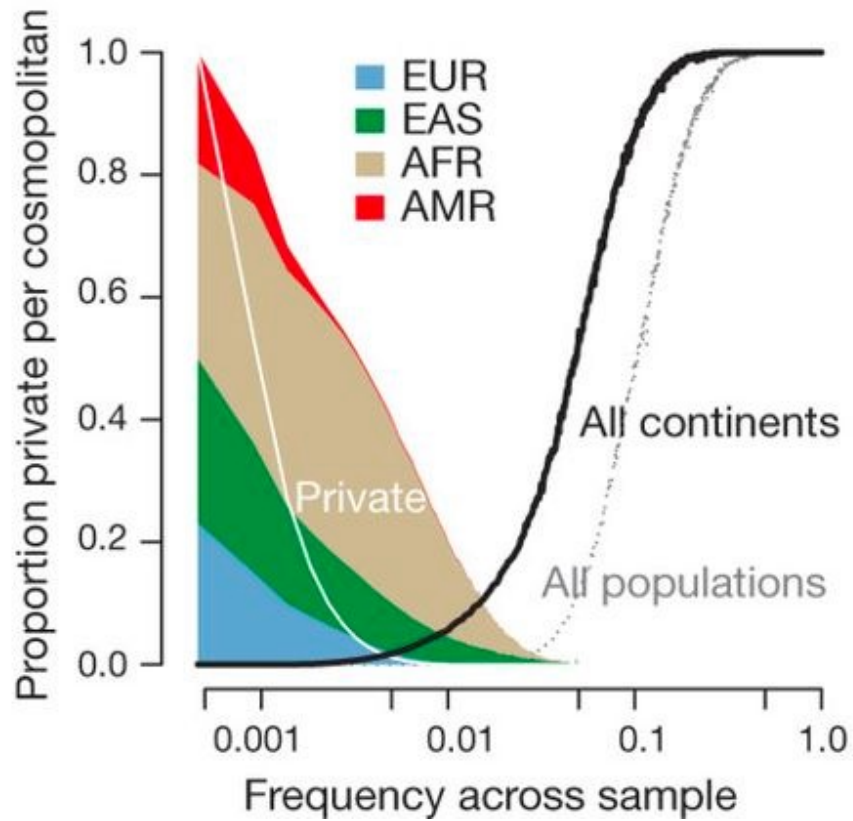
# IBD on a large scale



Scandinavians[1]
2–9 Generations ago
n = 4,189 (OR>10)

Finnish[1]
0–9 Generations ago
n = 1,688 (OR>10)

French canadians[1]
3–9 Generations ago
n = 9,689 (OR>10)

Northeast and Utah

Acadians

Northeast and Utah[3]
2–6 Generations ago
n = 99,315 (OR>5)

New mexicans[1]
2–4 Generations ago
n = 5,291 (OR>10)

Pennsylvania[3]
3–9 Generations ago
n = 80,754 (OR>2)

Lower midwest and appalachians[3]
3–8 Generations ago
n = 108,786 (OR>2)

Portuguese, incl. Azores and madeira islands[3,‡]
0–9 Generations ago
n = 3,468 (OR>10)

Upland South[3]
3–9 Generations ago
n = 93,305 (OR>2)

Lower South[3]
3–9 Generations ago
n = 77,581 (OR>5)

Acadians[1]
2–9 Generations ago
n = 6,615 (OR>10)

Caribbeans

Northeast Mexico[1]
3–9 Generations ago
n = 6,311 (OR>10)

Dominicans[1,**]
2–9 Generations ago
n = 779 (OR>10)

Caribbeans[1]
2–9 Generations ago
n = 9,315 (OR>10)

West mexico[1]
3–9 Generations ago
n = 1,356 (OR>10)

Portuguese

Han et al. (2017)
Nature Comm.

# Do rare variants help identify recent population structure?



1000 Genomes Project (2012) Nature

# Do rare variants help identify recent population structure?
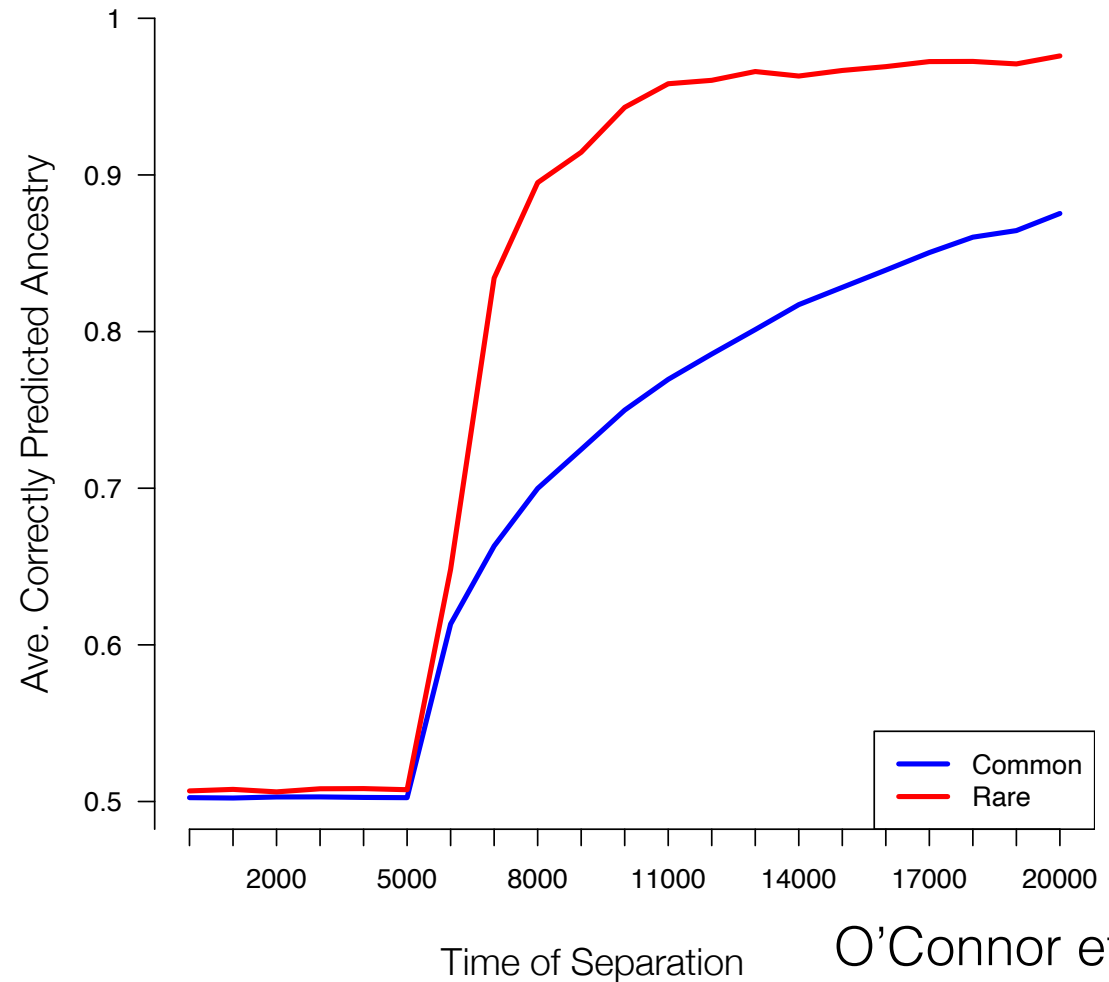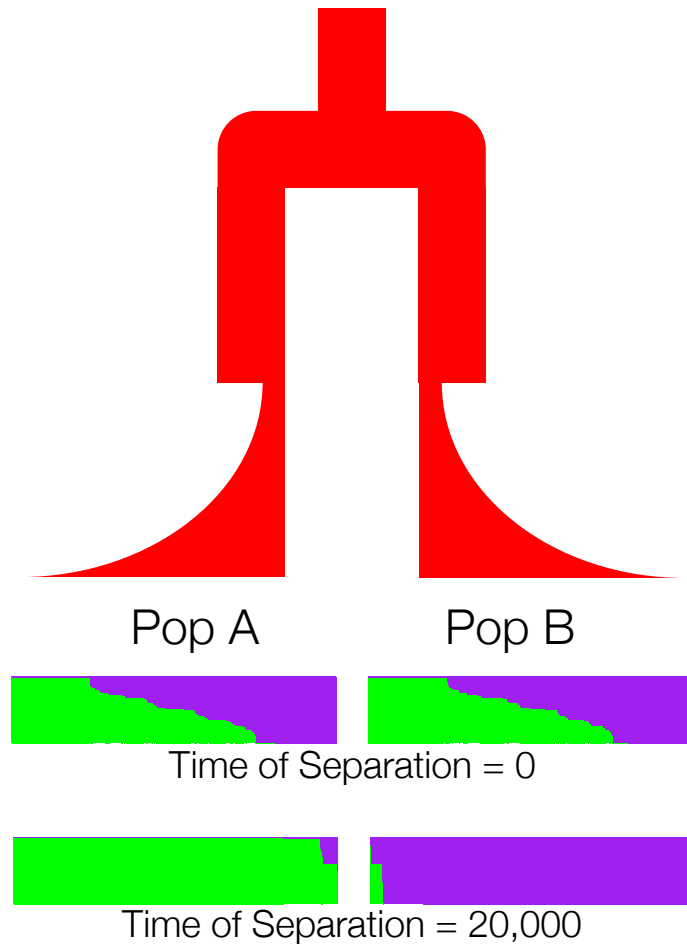


1000 Genomes Project (2012) Nature



Novembre et al. (2012) Science

# Rare VS Common:
## Assignment of Ancestry Proportions

Pop A   Pop B

Time of Separation = 0

Time of Separation = 20,000

O'Connor et al. (2014)
Mol. Biol. Evol.

# Rare VS Common: Which has Greater Information? And When?

Information Gain: how well a variant can distinguish between populations. (Rosenberg et al. 2003)

$$I_n(Q;J) = \sum_{j=1}^{N} \left( -p_j \ln p_j + \sum_{i=1}^{K} q_i p_{ij} \ln p_{ij} \right)$$

Expected Information Gain
- Calculate for a specific site count
- Correct for missing data
- Weighted average to calculate across a range of frequency (rare or common)

$$E(I_n \mid C, M) = \sum_{m \in M} \sum_{l=0}^{C} r_{lm} \times \sum_{j=1}^{N} \left( -p_{jlm} \ln p_{jlm} + \sum_{i=1}^{K} q_i p_{ijlm} \ln p_{ijlm} \right)$$



O'Connor et al. (2014)
Mol. Biol. Evol.

# Rare Variants Identify Cryptic Populations



Common (MAF > 10%)

O'Connor et al. (2014)
Mol. Biol. Evol.

# Rare Variants Identify Cryptic Populations



Common (MAF > 10%)

Rare (MAF < 0.5%)

O'Connor et al. (2014)
Mol. Biol. Evol.

# What is Their Geographic Ancestry?

# PCA of Global Diversity Including Cryptic Population



O'Connor et al. (2014)
Mol. Biol. Evol.
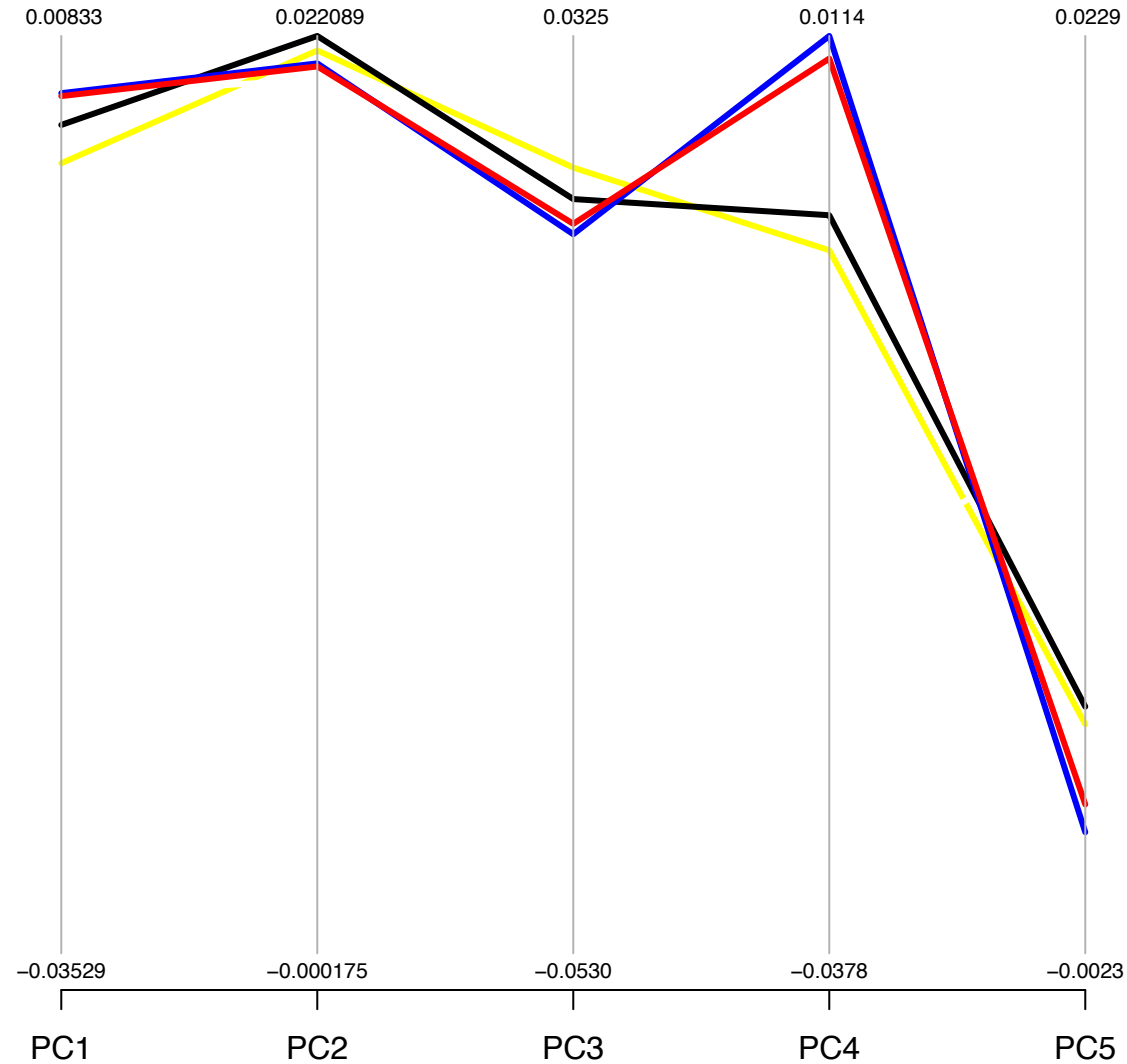
# PCA of Global Diversity Including Cryptic Population



O'Connor et al. (2014)
Mol. Biol. Evol.

# Population Average PCA with More Axes



O'Connor et al. (2014)
Mol. Biol. Evol.

# Population Average PCA with More Axes

0.00833    0.022089    0.0325    0.0114    0.0229

Unknown
Ashkenazi
Moroccan
Sephardic

−0.03529    −0.000175    −0.0530    −0.0378    −0.0023

PC1    PC2    PC3    PC4    PC5

O'Connor et al. (2014)
Mol. Biol. Evol.

# Trans-Omics for Precision Medicine (TOPMed) Cohorts

- N ≅ 55K
- Predominantly African, Latino, and European American
  - Samoa
  - Amish
- All are well characterized for heart, lung, blood, and sleep phenotypes

Taliun et al. (2019) Bioarxiv
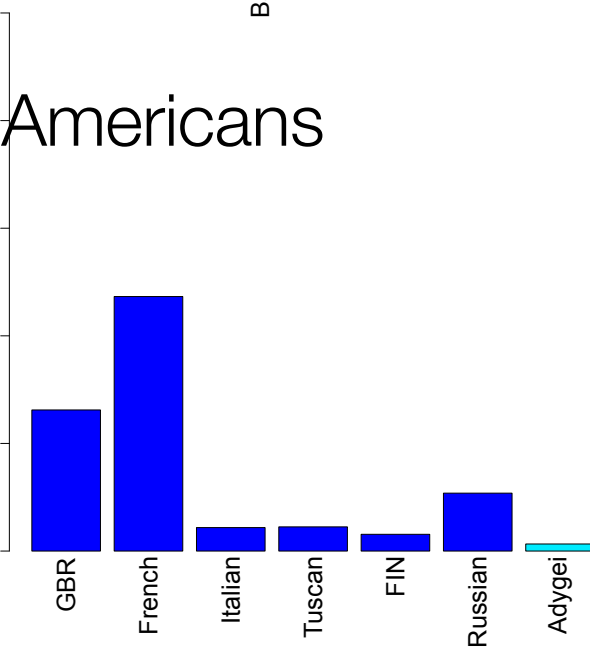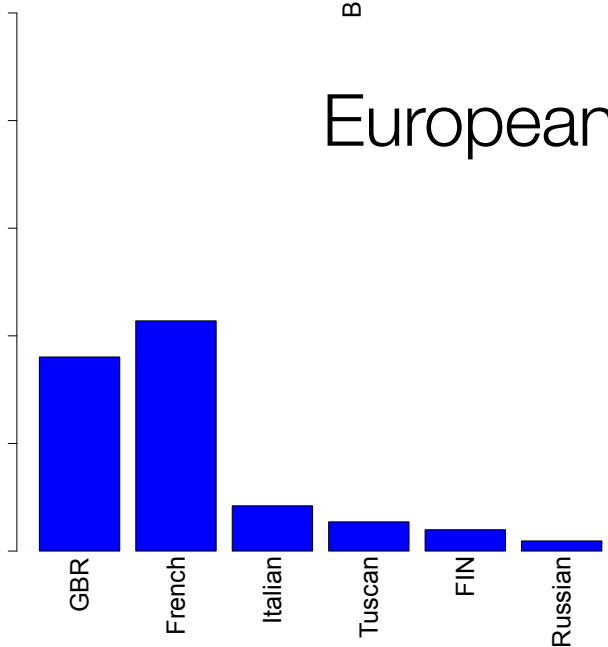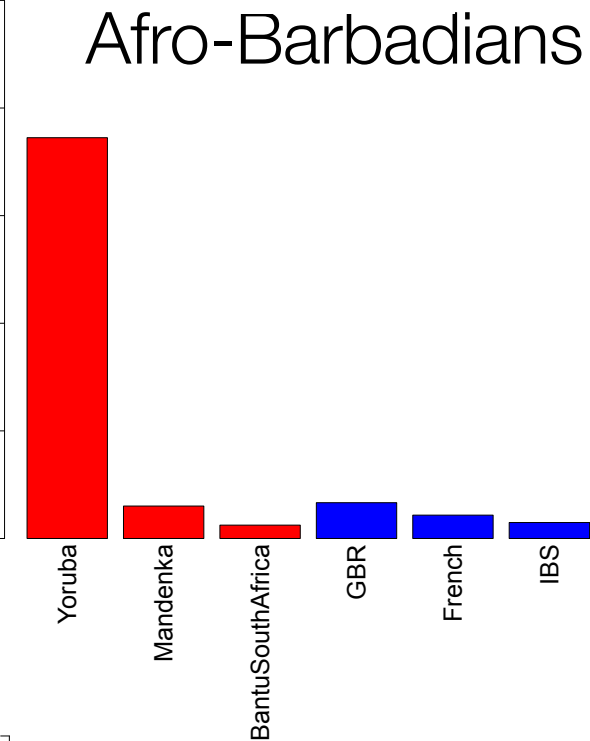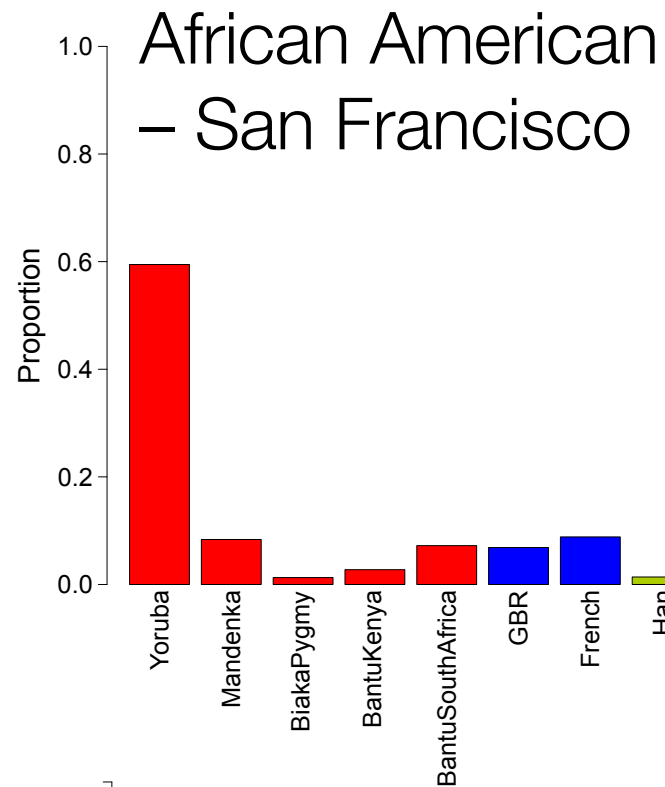
# Rare variant sharing across cohorts

- Allele Count 2 to 100
- Corrected for:
  - sample size
  - Genome-wide heterozygosity

Taliun et al. (2019) Bioarxiv

# Rare variant sharing across cohorts

- Allele Count 2 to 100
- Corrected for:
  - sample size
  - Genome-wide heterozygosity

Taliun et al. (2019) Bioarxiv

# Rare variant sharing across cohorts
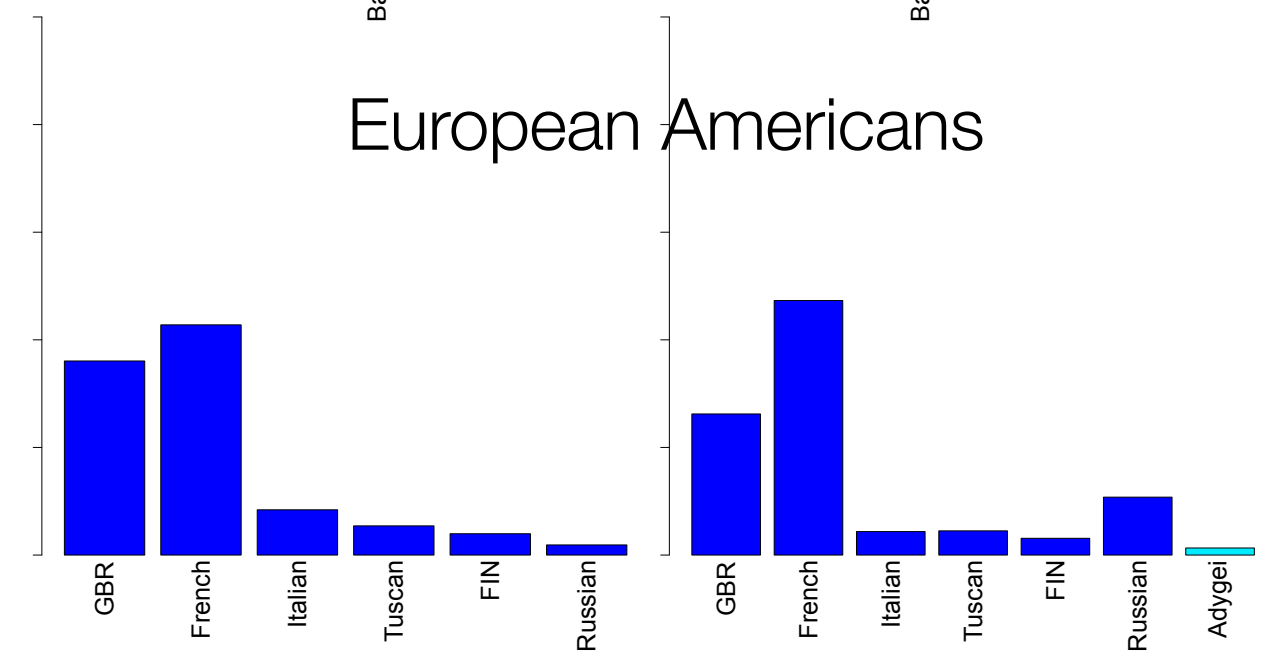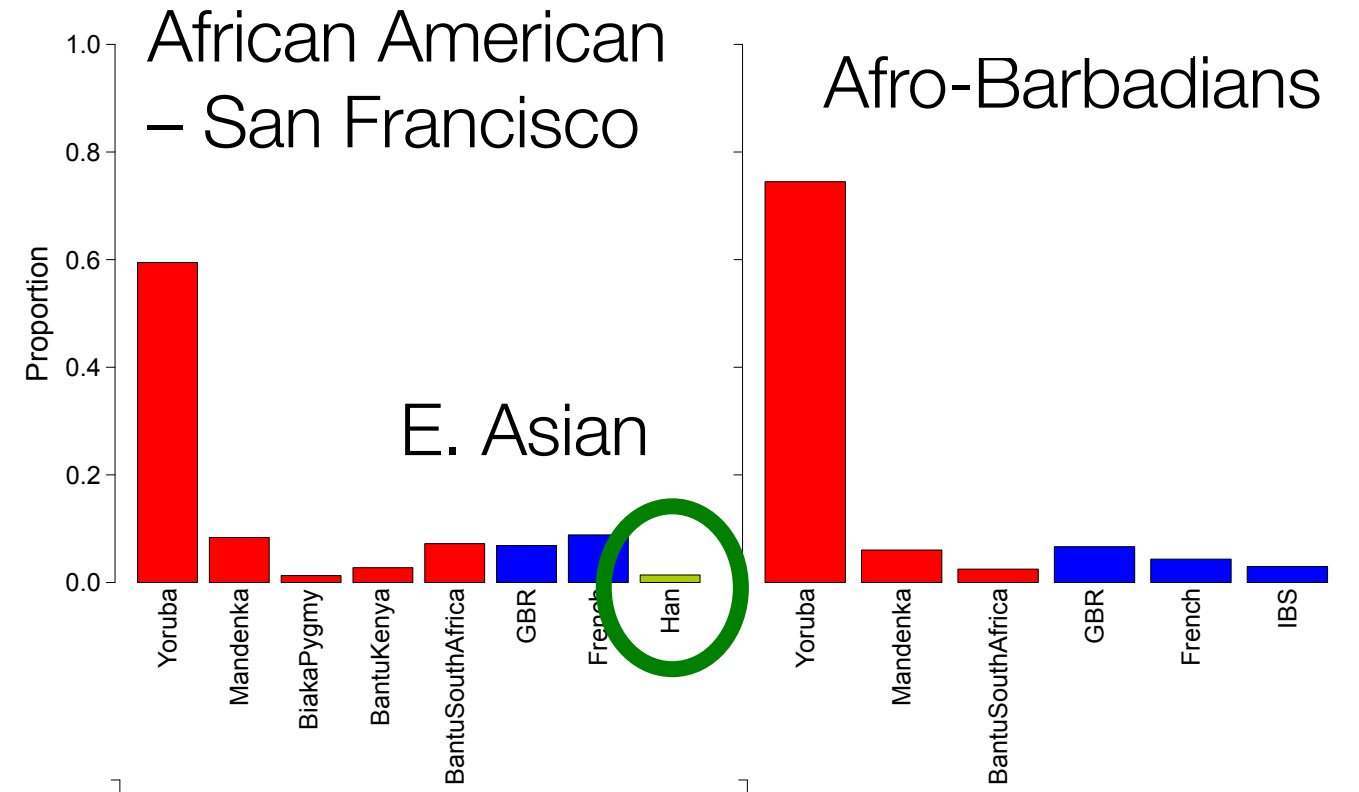
- Allele Count 2 to 100
- Corrected for:
  - sample size
  - Genome-wide heterozygosity

Taliun et al. (2019) Bioarxiv

# Rare variant sharing across cohorts

- Allele Count 2 to 100
- Corrected for:
  - sample size
  - Genome-wide heterozygosity

Taliun et al. (2019) Bioarxiv



Cohort Ancestry Types

- African American
- Amish
- Asian American
- European American
- Hispanic/Latino
- Samoan

fineStructure analysis of genome-wide ancestry

**Legend:**
- African (red)
- Caucasia (cyan)
- East Asian (olive)
- European (blue)

fineStructure analysis of genome-wide ancestry

**Legend:**
- African (red)
- Caucasia (cyan)
- East Asian (olive/yellow-green)
- European (blue)

fineStructure analysis of genome-wide ancestry

fineStructure analysis of genome-wide ancestry

Legend:
- African (red)
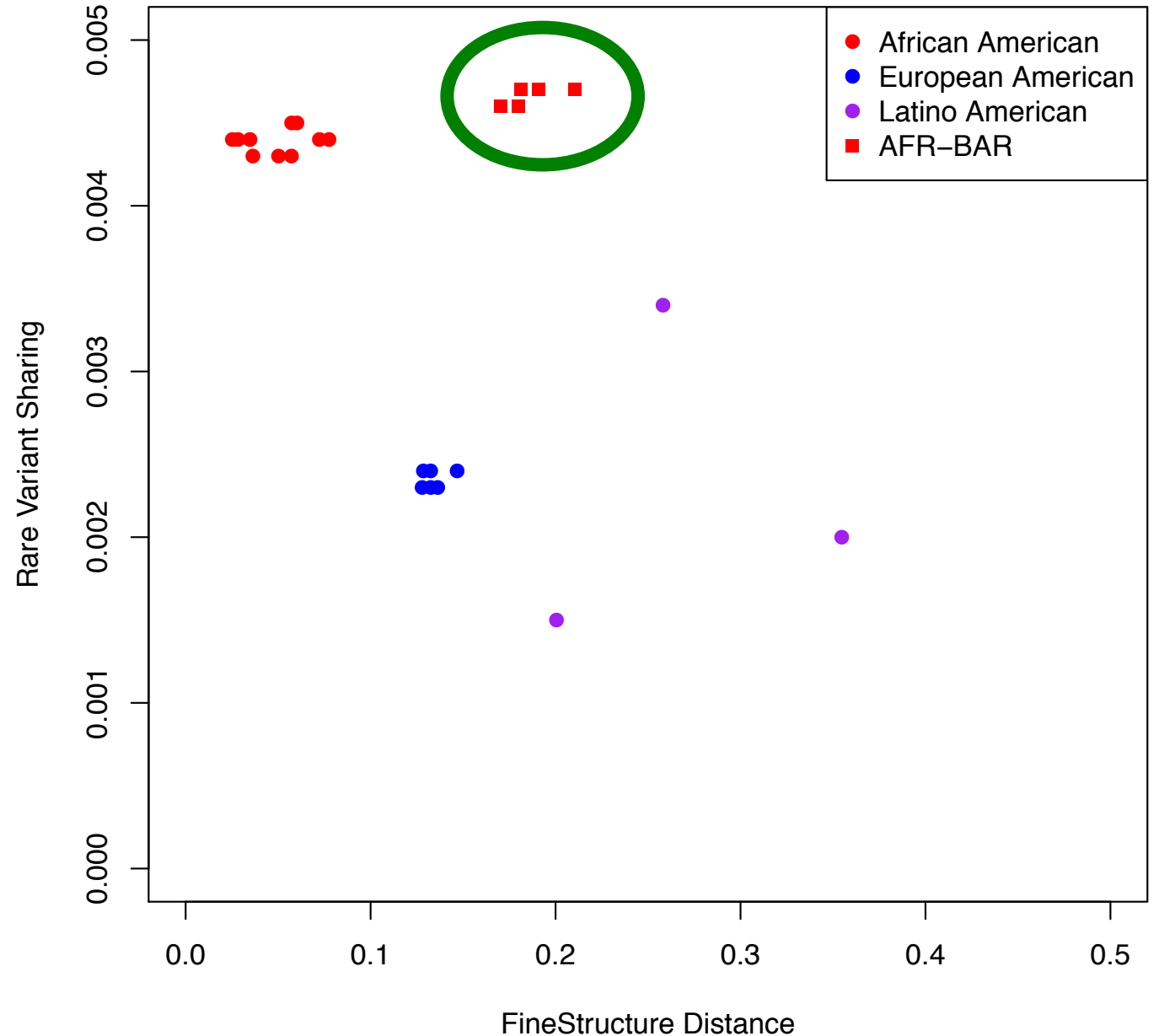- Caucasia (cyan)
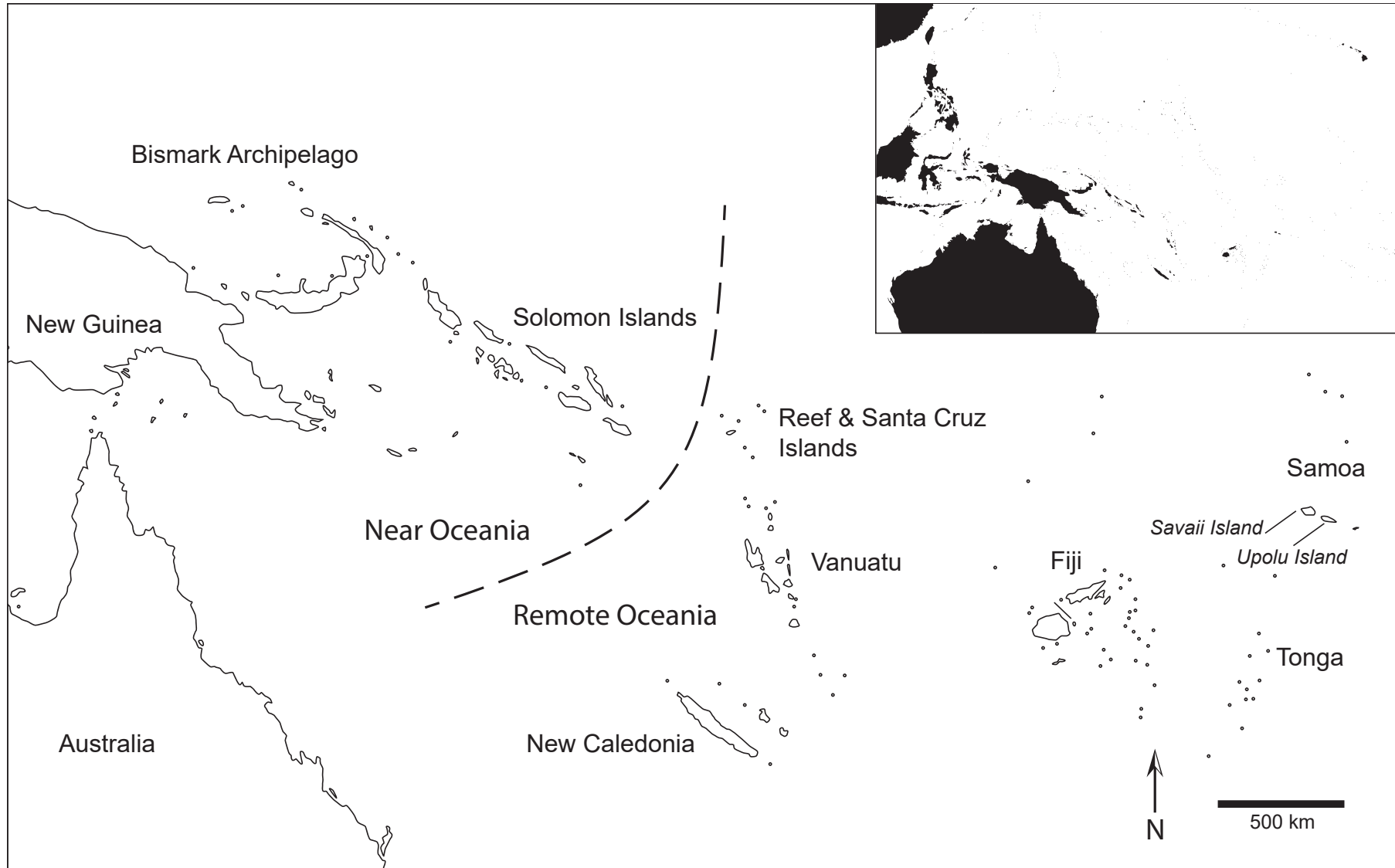- East Asian (olive)
- European (blue)

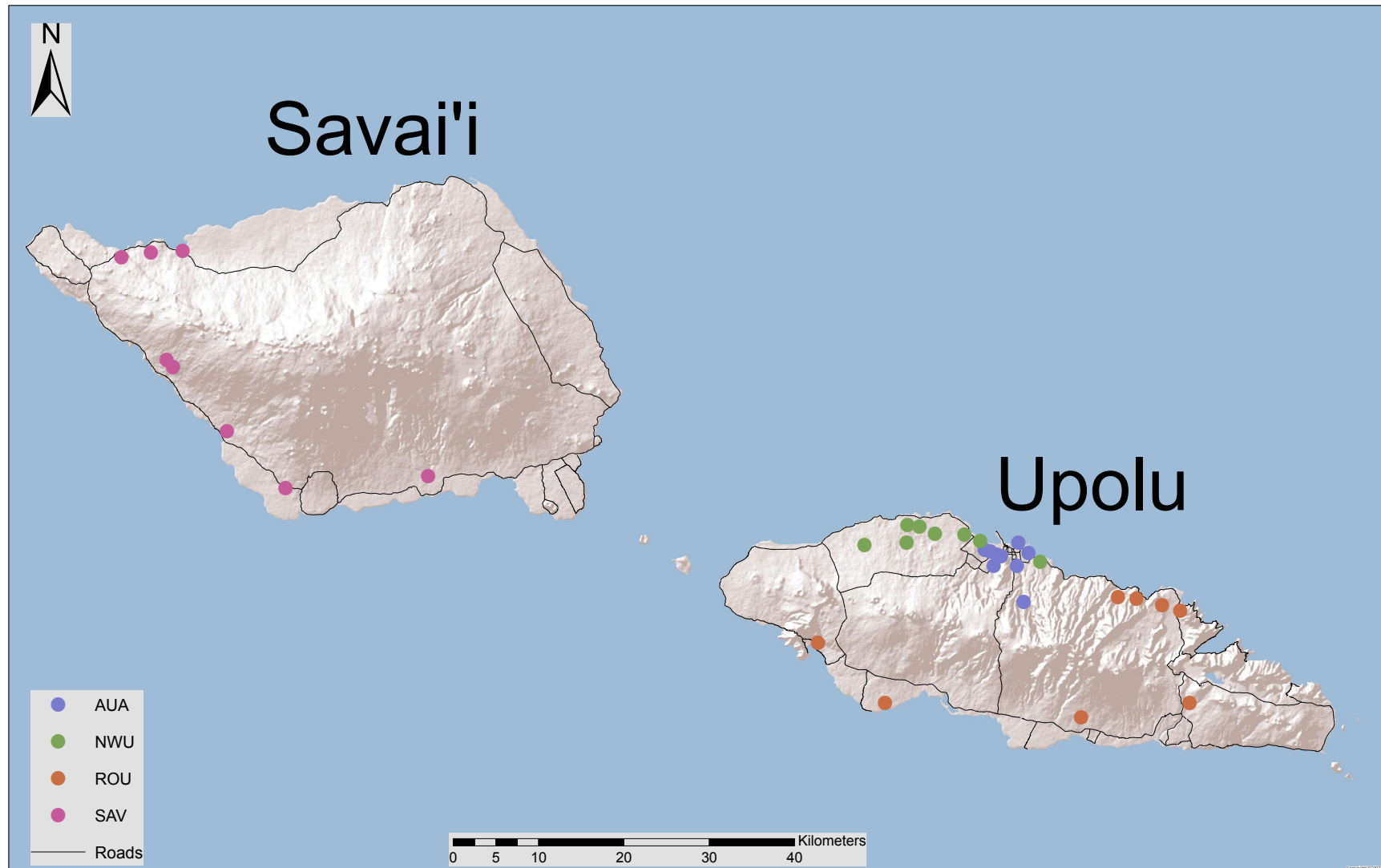# African American's have more homogeneous ancestral proportions

- Calculated Euclidian distance between fineSTRUCTURE proportions
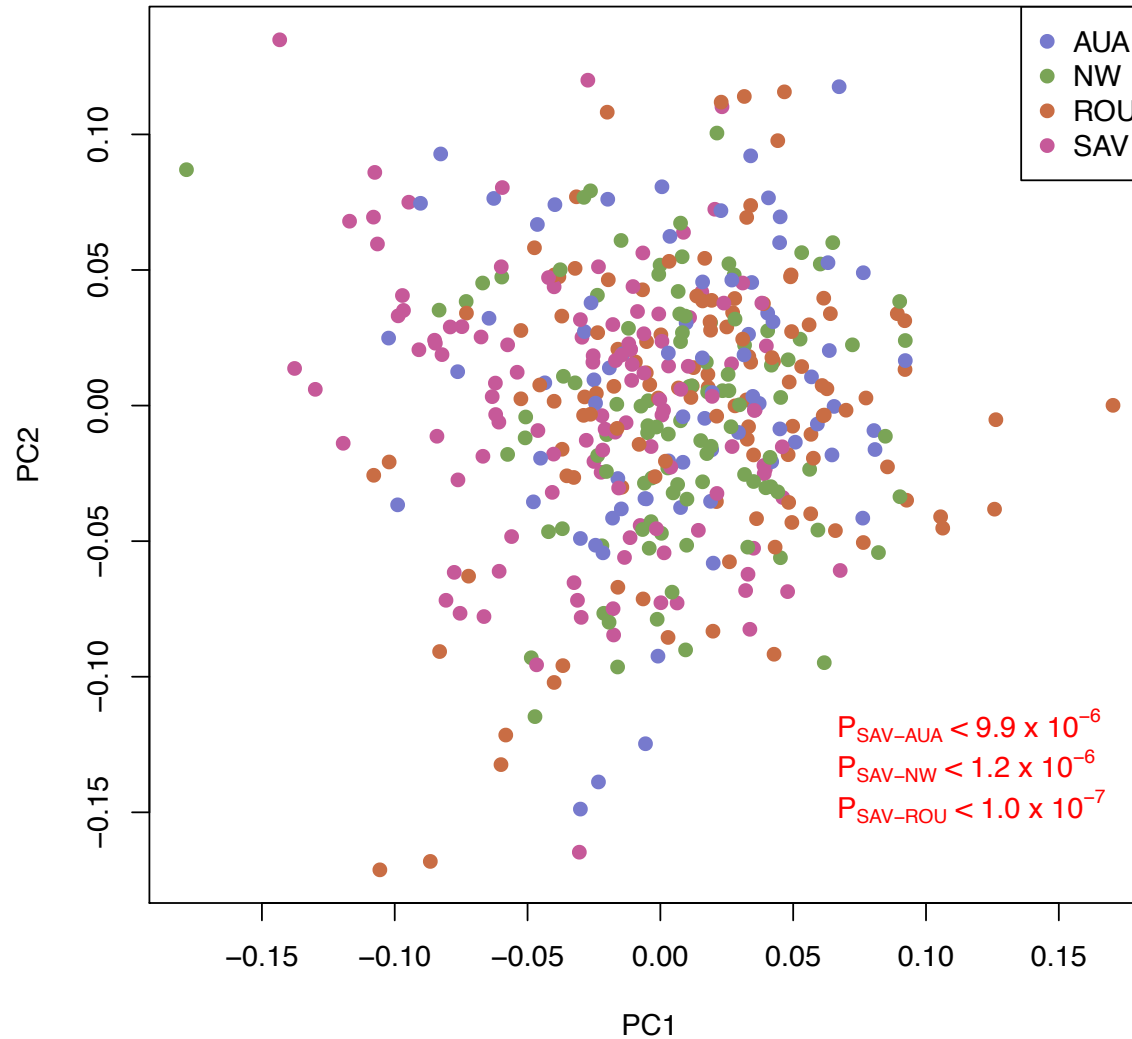- African American cohorts have the shortest distance and the greatest rare variant sharing

# African American's have more homogeneous ancestral proportions

- Calculated Euclidian distance between fineSTRUCTURE proportions
- African American cohorts have the shortest distance and the greatest rare variant sharing

# Quick background on Samoa



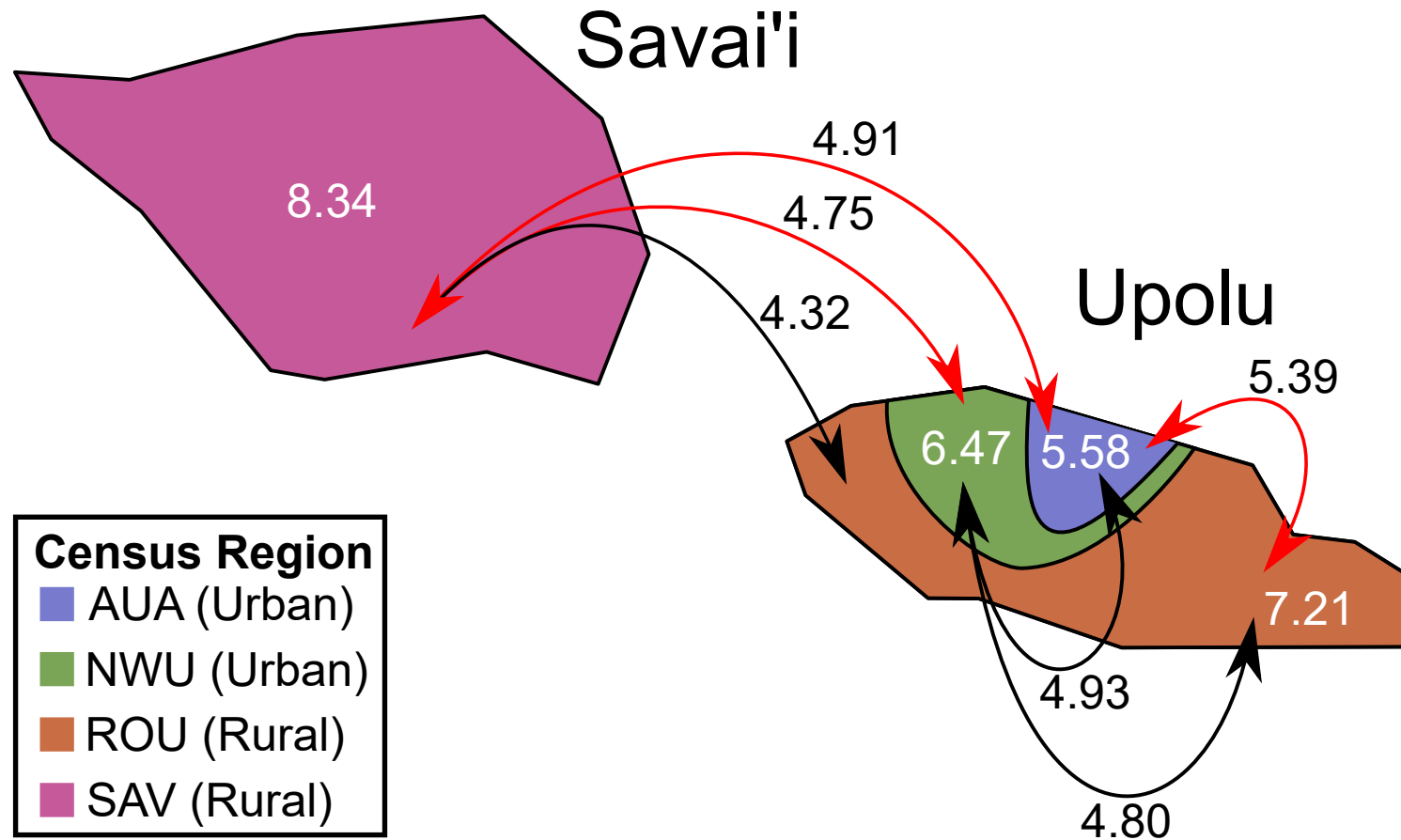Harris et al. (https://papers.ssrn.com/sol3/papers.cfm?abstract_id=3329885)

# Quick background on Samoa



Harris et al. (https://papers.ssrn.com/sol3/papers.cfm?abstract_id=3329885)

# PCA with all variants can't distinguish the two islands well.



Harris et al. (2020)
PNAS

Rare variant Sharing in Samoa

Savai'i

8.34

4.91

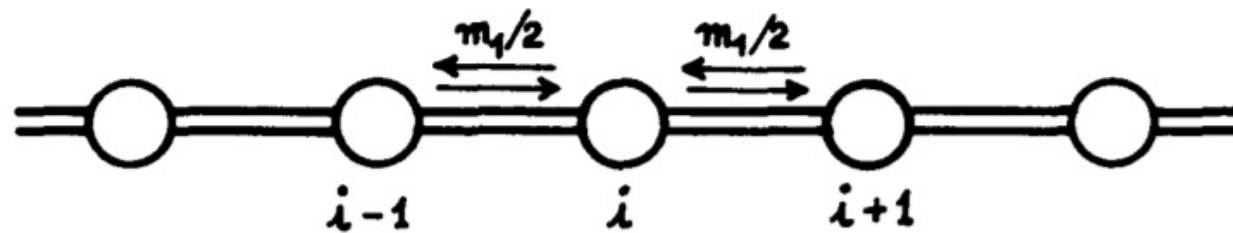4.75

Upolu

4.32

6.47    5.58    5.39

**Census Region**
- AUA (Urban)
- NWU (Urban)
- ROU (Rural)
- SAV (Rural)

7.21

4.93

4.80

Harris et al. (2020)
PNAS

# Estimated Effective Migration Surfaces (EEMS)



Forest
Savanna

Petkova et al. (2015)
Nature Genet.

log(*m*)

# Assumptions: Stepping Stone Model

- Migration can only occur between adjacent demes
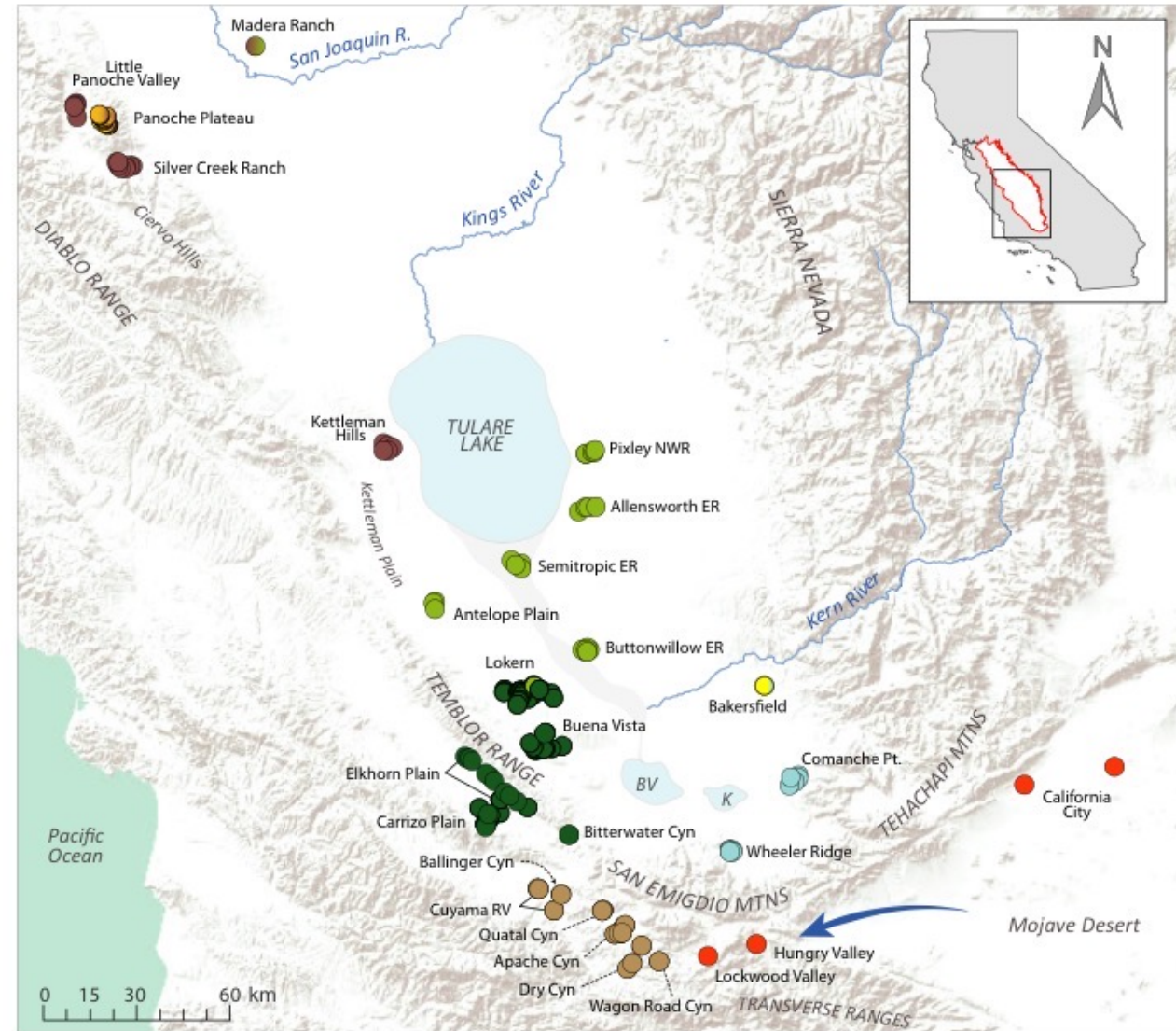- Migration rate between each deme is assumed to be equal
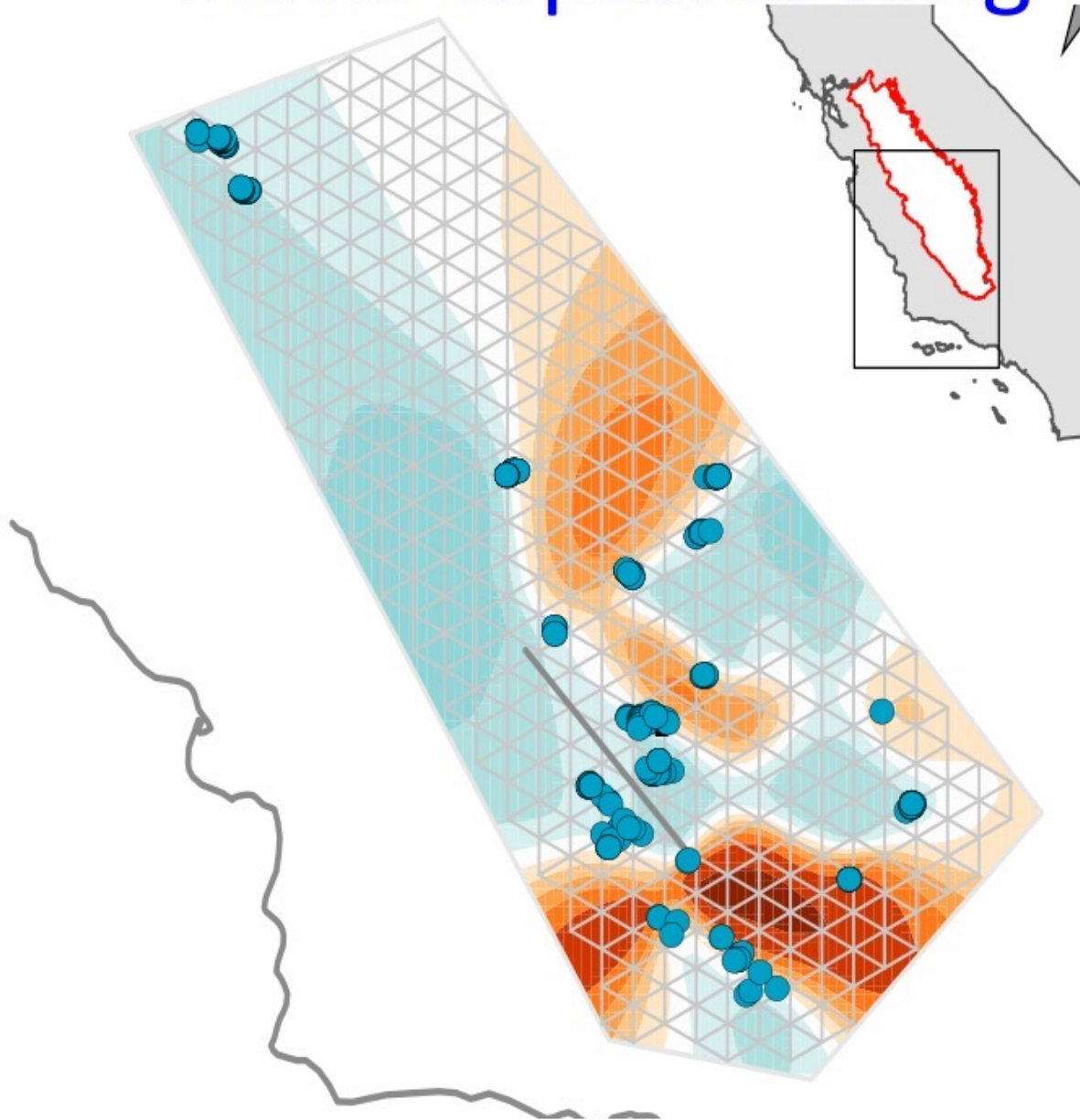


Kimura and Weiss (1964)

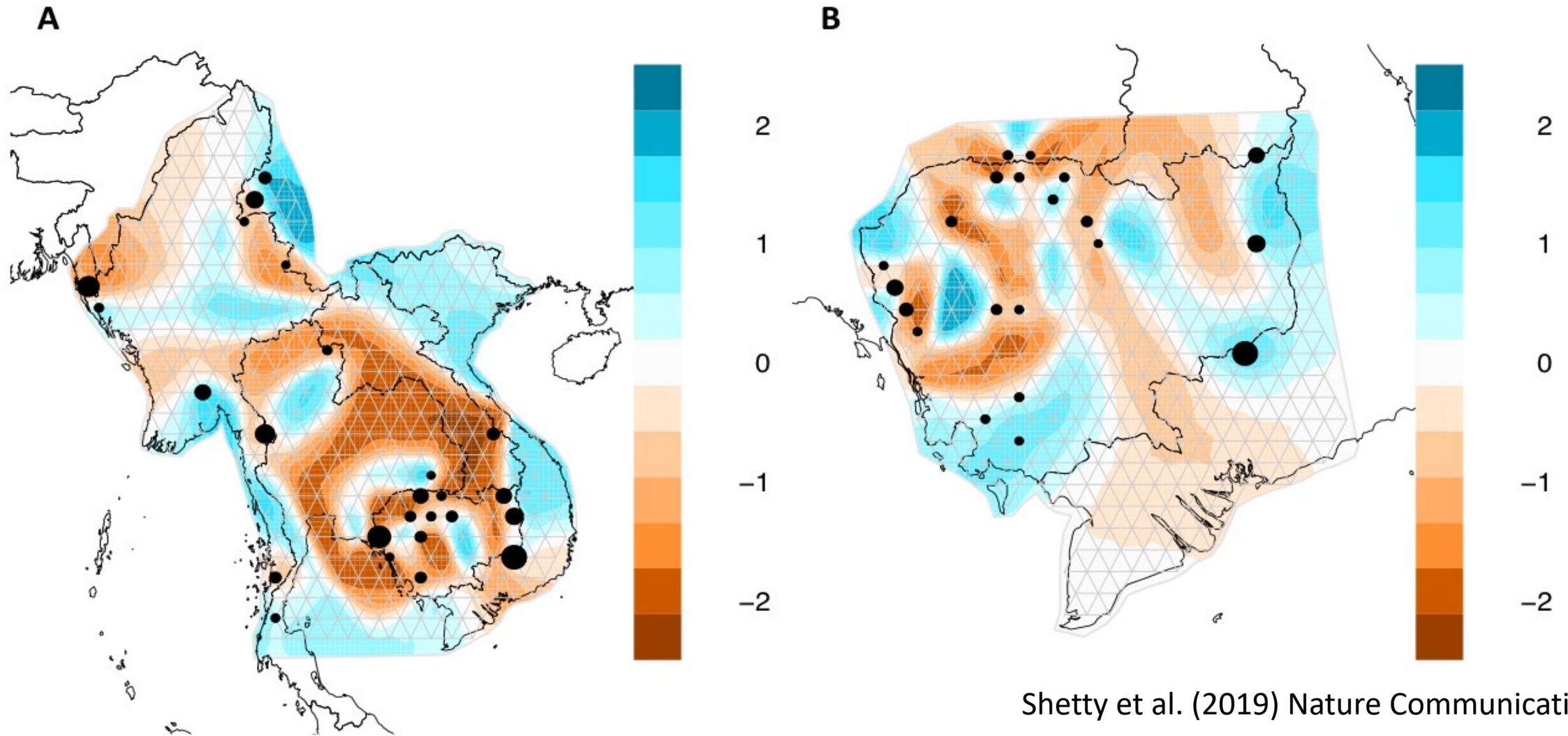# EEMS: Migration and diversity within Peru



Harris et al. (2018) PNAS

| Deme | Populations | Deme | Populations | Deme | Populations | Deme | Populations |
|------|-------------|------|-------------|------|-------------|------|-------------|
| 1 | Piapoco | 5 | Lima | 9 | Cusco, Qeros, Quechua | 15 | Surui |
| 2 | Iquitos | 6 | Chopccas | 10 | Puno, Uros | | |
| 3 | Matzes | 7 | Nahua | 11,12,13 | Bolivian | | |
| 4 | Moches, Trujillo | 8 | Matsiguenka | 14 | Karitiana | | |

# EEMS captures long-term migration patterns



Richmond et al. (2015) Molecular Ecology

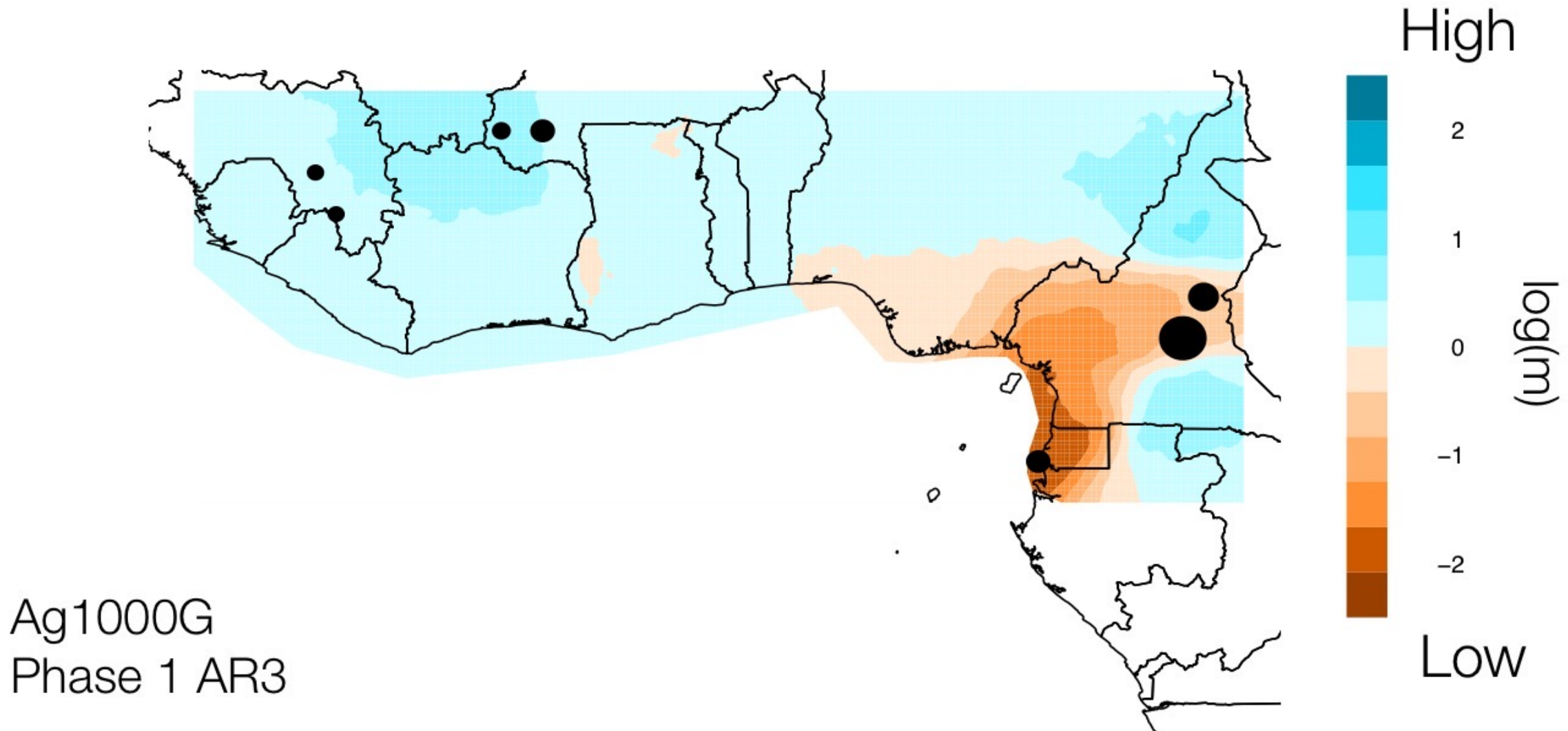# EEMS in Malaria Parasites of South East Asia



Shetty et al. (2019) Nature Communications

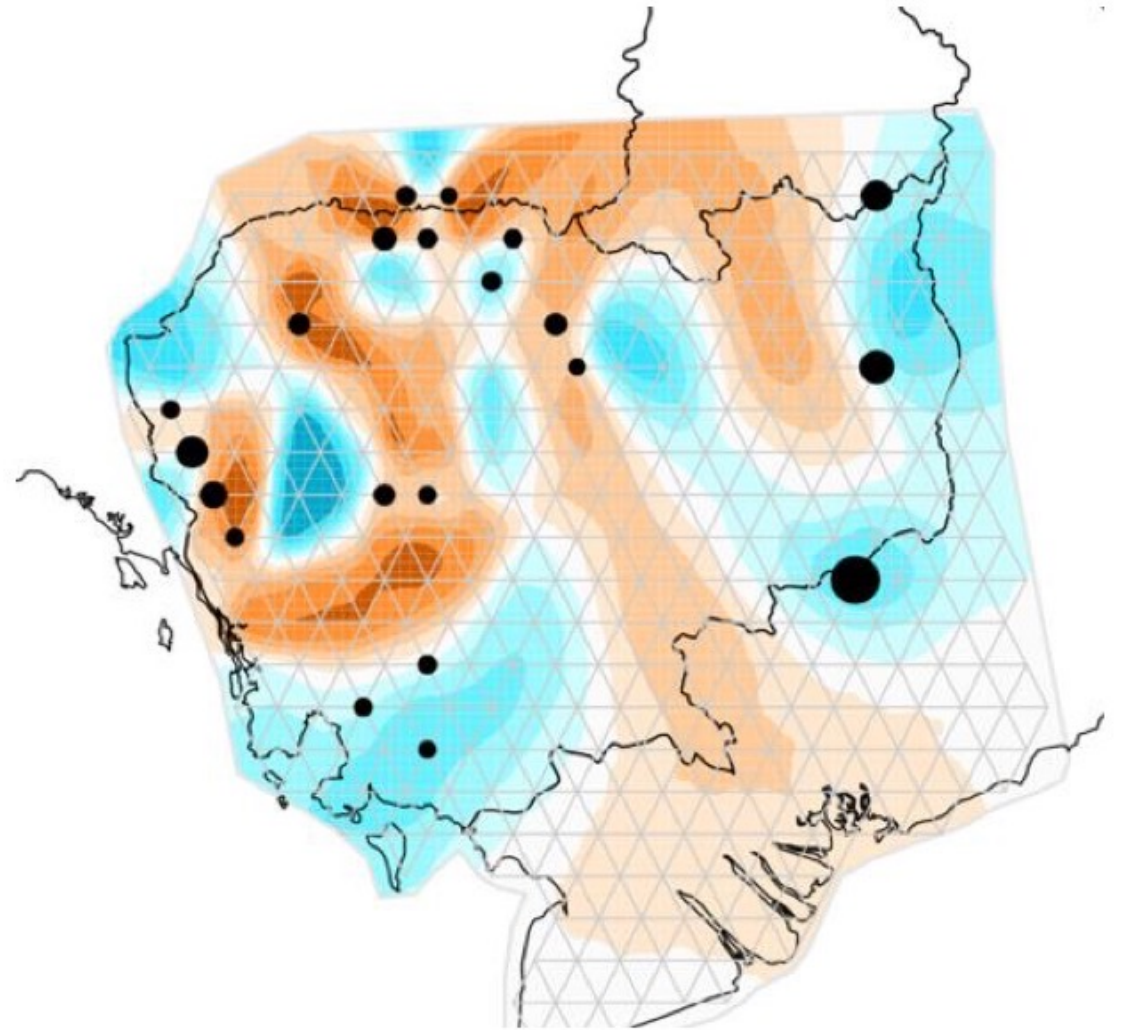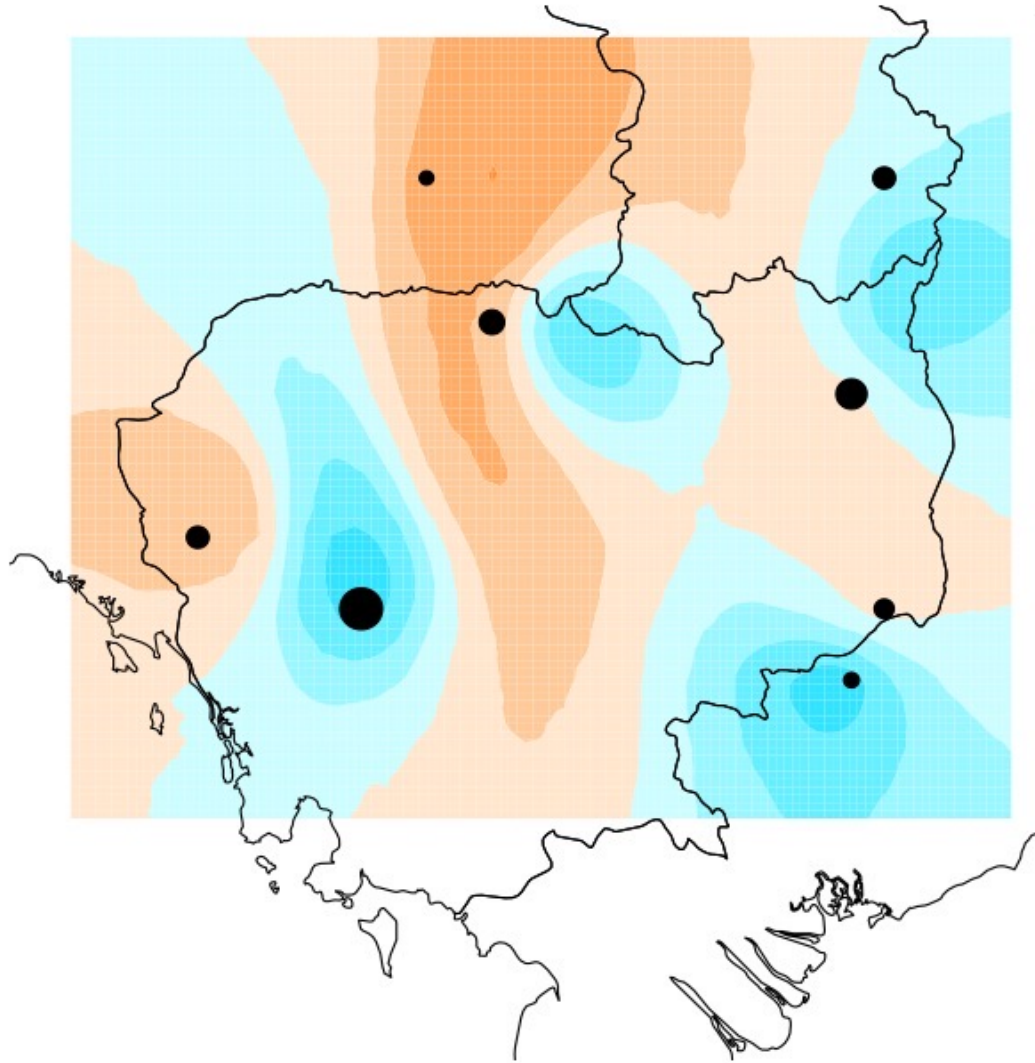# Application to Malaria Parasites in W. Africa

Pf3K Version 5.1

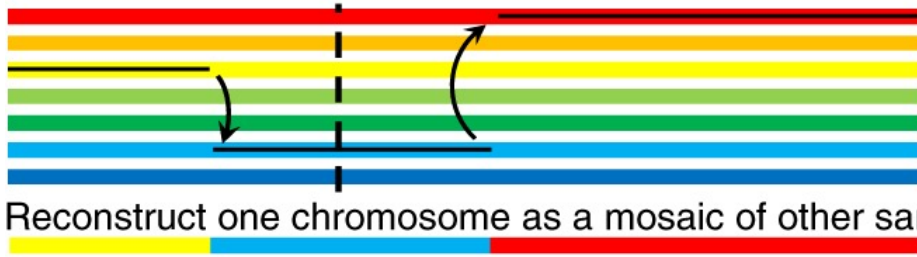Application to Mosquito in W. Africa

Ag1000G
Phase 1 AR3

# Robustness of Sampling on EEMS

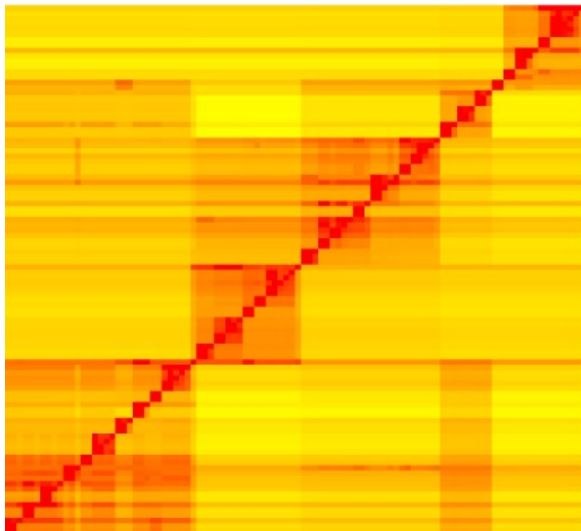# RELATE: a means of finding genealogical local genomic relationships
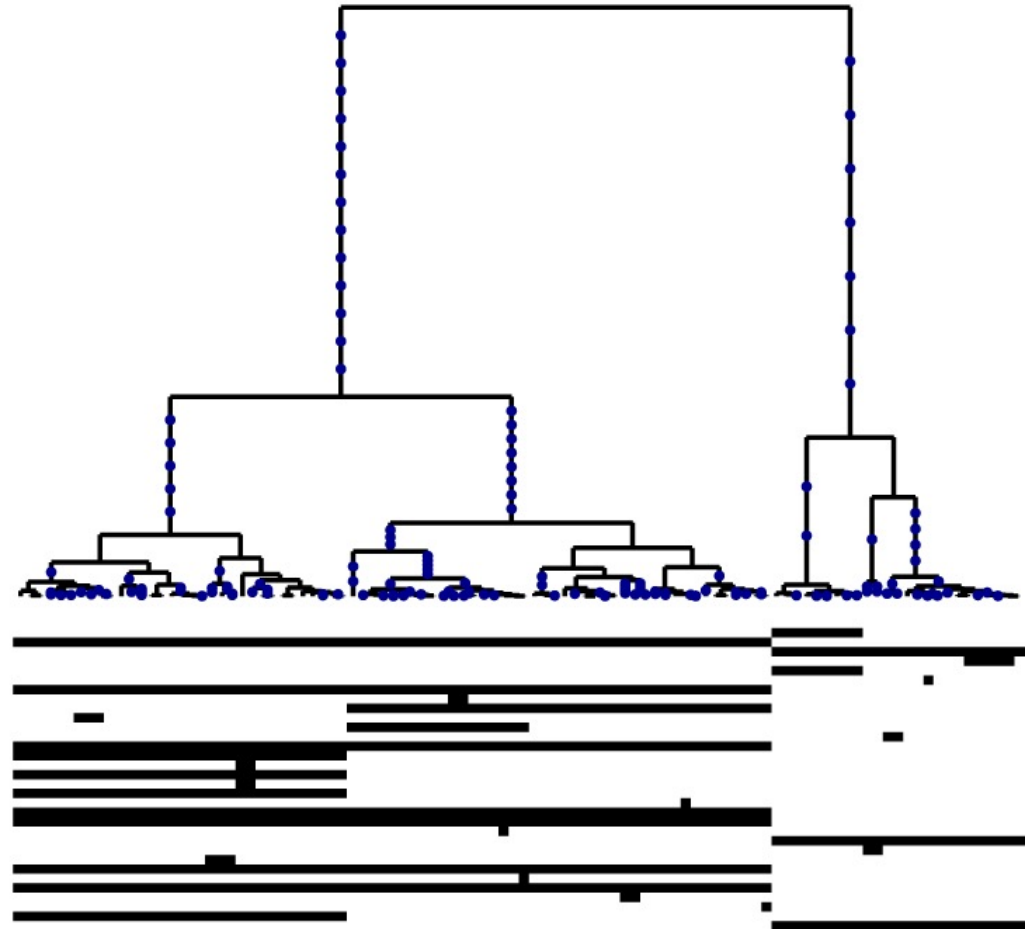


Modified Li and Stephens HMM

focal SNP

Reconstruct one chromosome as a mosaic of other samples

Store position specific distance matrix containing transformed probabilities of copying from each other sample
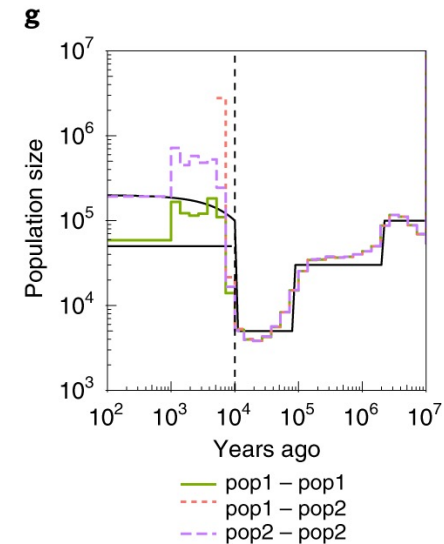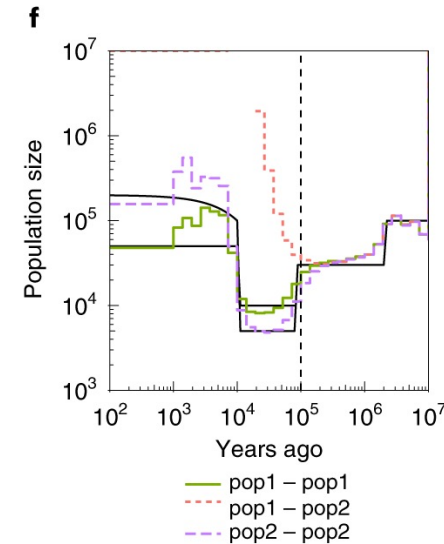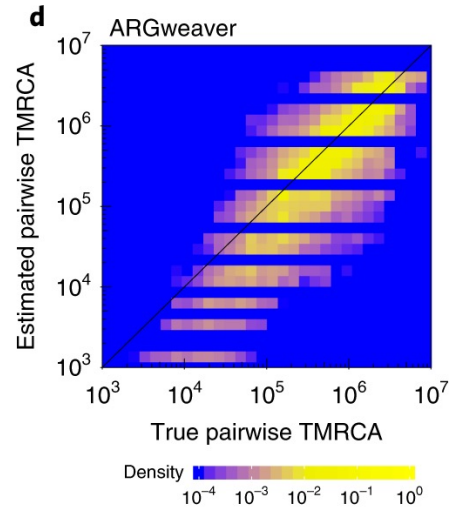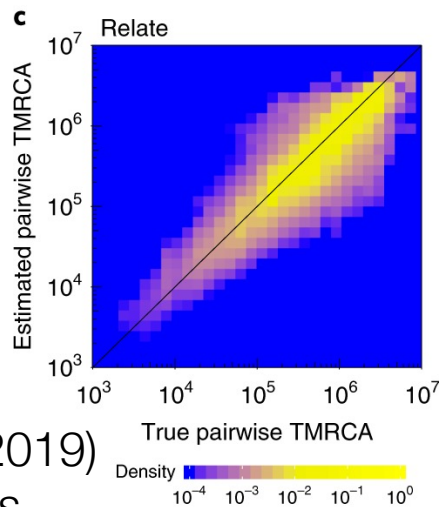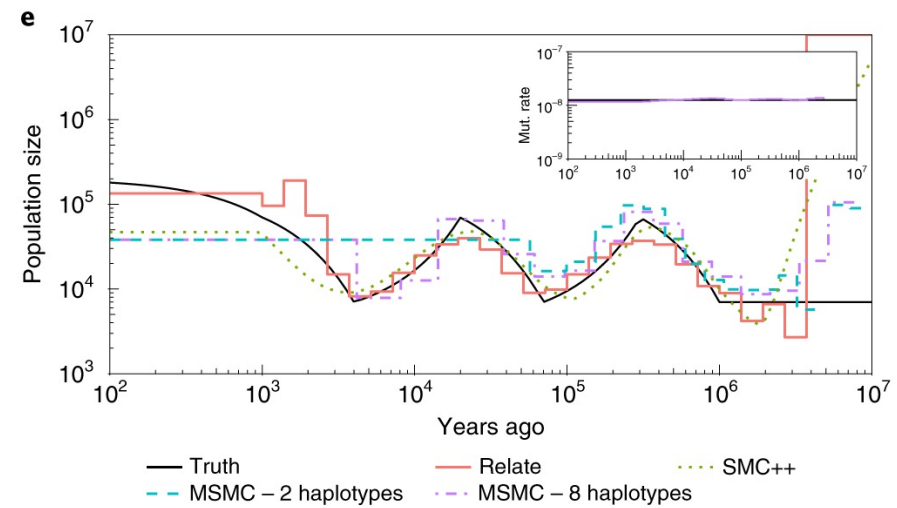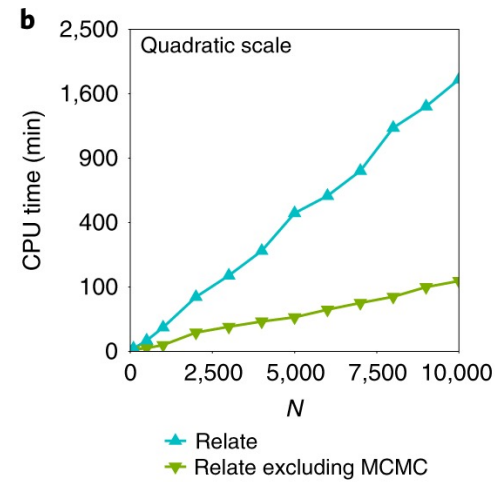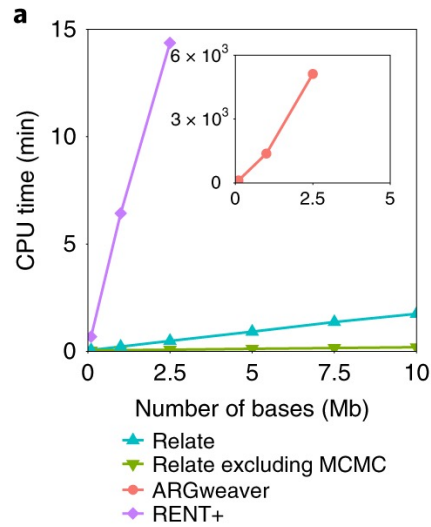
Hierarchical clustering and coalescent model-based branch length estimation produce local trees

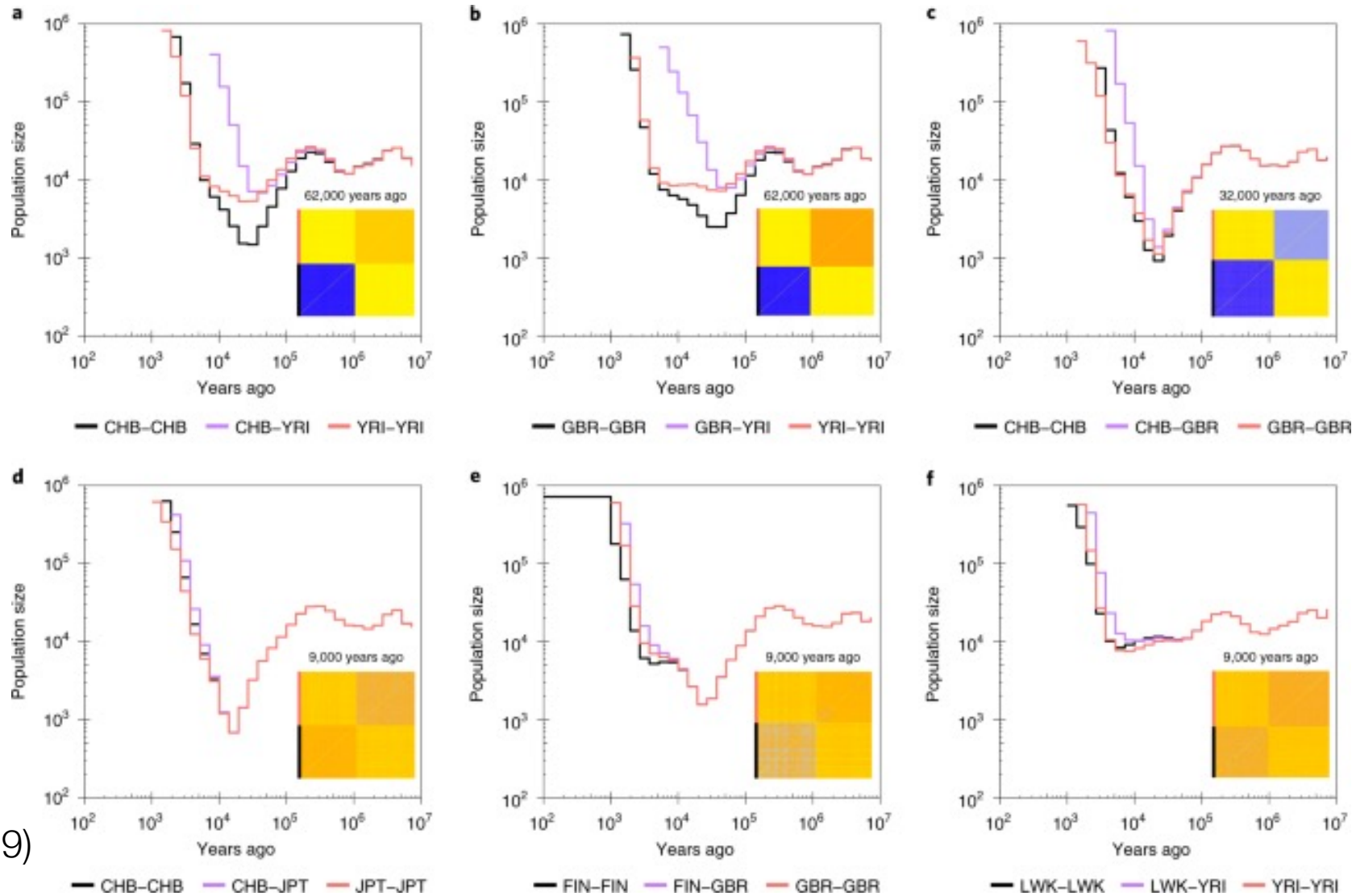Speidel et al. (2019) Nature Genetics

Haplotype data sorted using constructed tree

# RELATE of course was tested with simulation!



Speidel et al. (2019)
Nature Genetics

# RELATE tested on 1000 Genomes Data



Speidel et al. (2019) Nature Genetics

# Concluding summary

- Fine-scale population structure is subdivisions of individuals on an ever increasingly granular scale

- Identity-by-descent and rare variant sharing are a powerful methods of identifying recent relationships and can be scaled by time.

- Cryptic population structure arises with extended relationships within a cohort, unknown to the investigators.

- EEMS can visualize migration patterns on a fine-scale illustrating cryptic structure not observed with other methods

# Questions?



CHOPCCAS
CUSCO
IQUITOS
MATZES
MOCHES
TRUJILLO
UROS