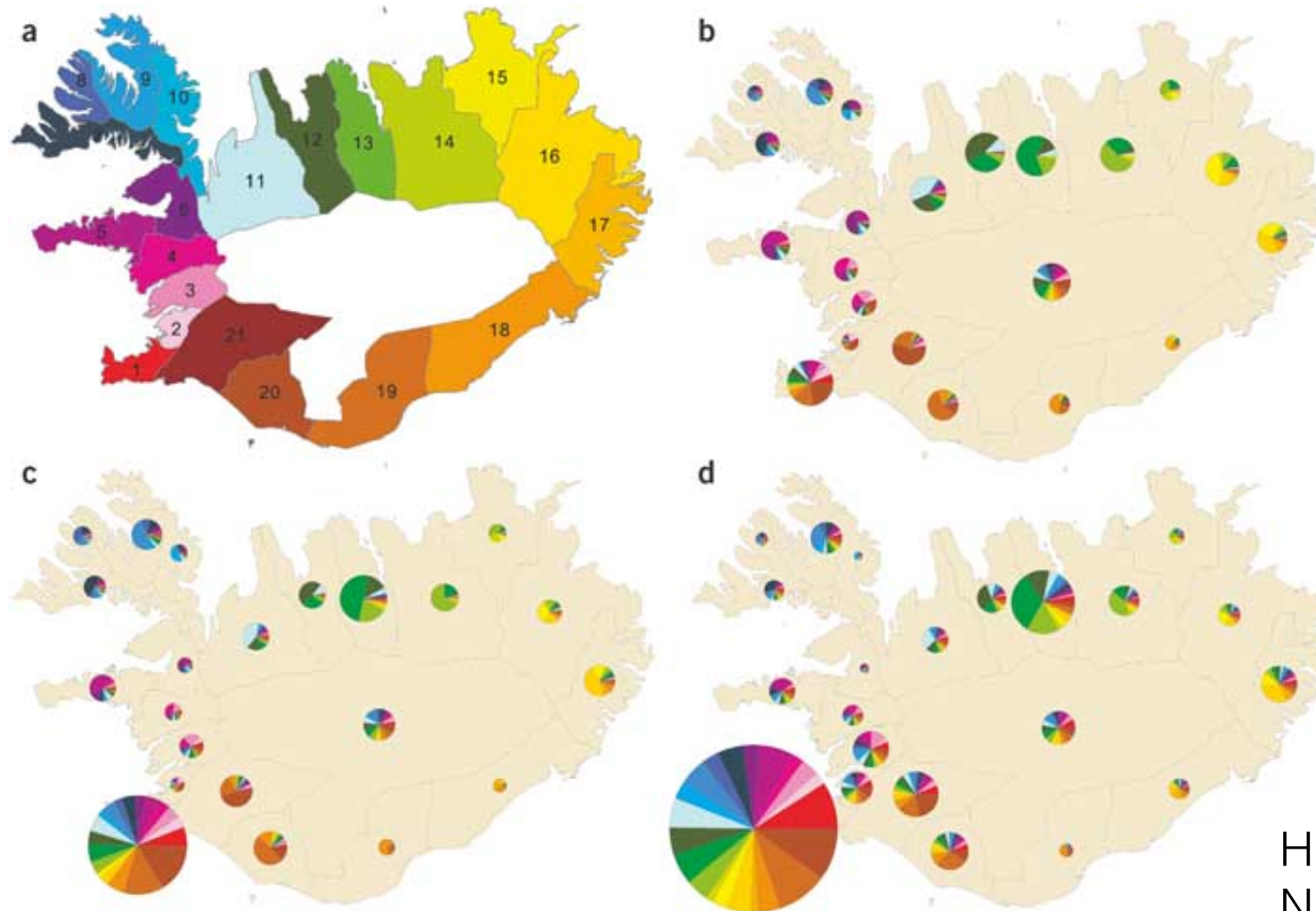# Cryptic Relatedness and fine scale population structure
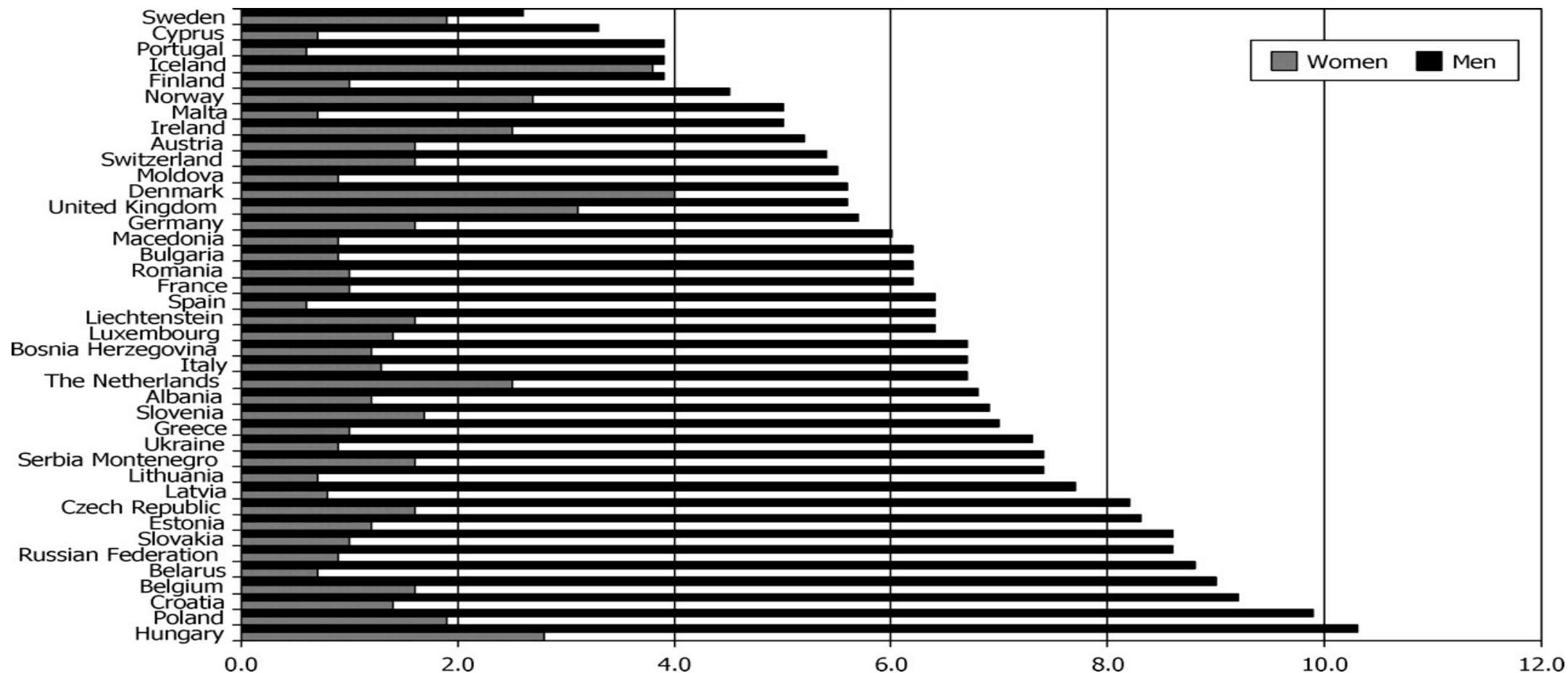
# Learning objectives

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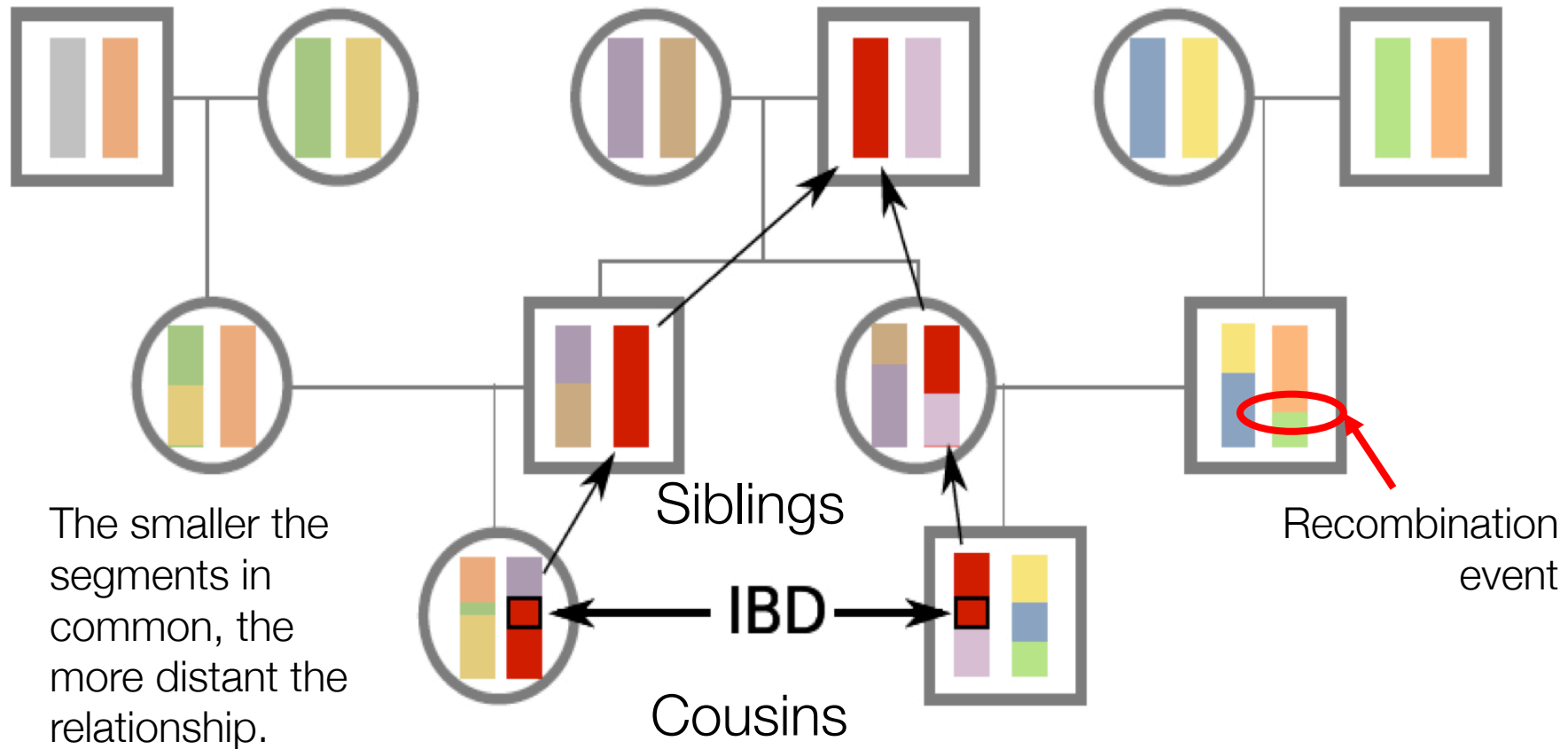
# Cryptic Population Structure



Helgason et al. (2004)
Nature Genet.
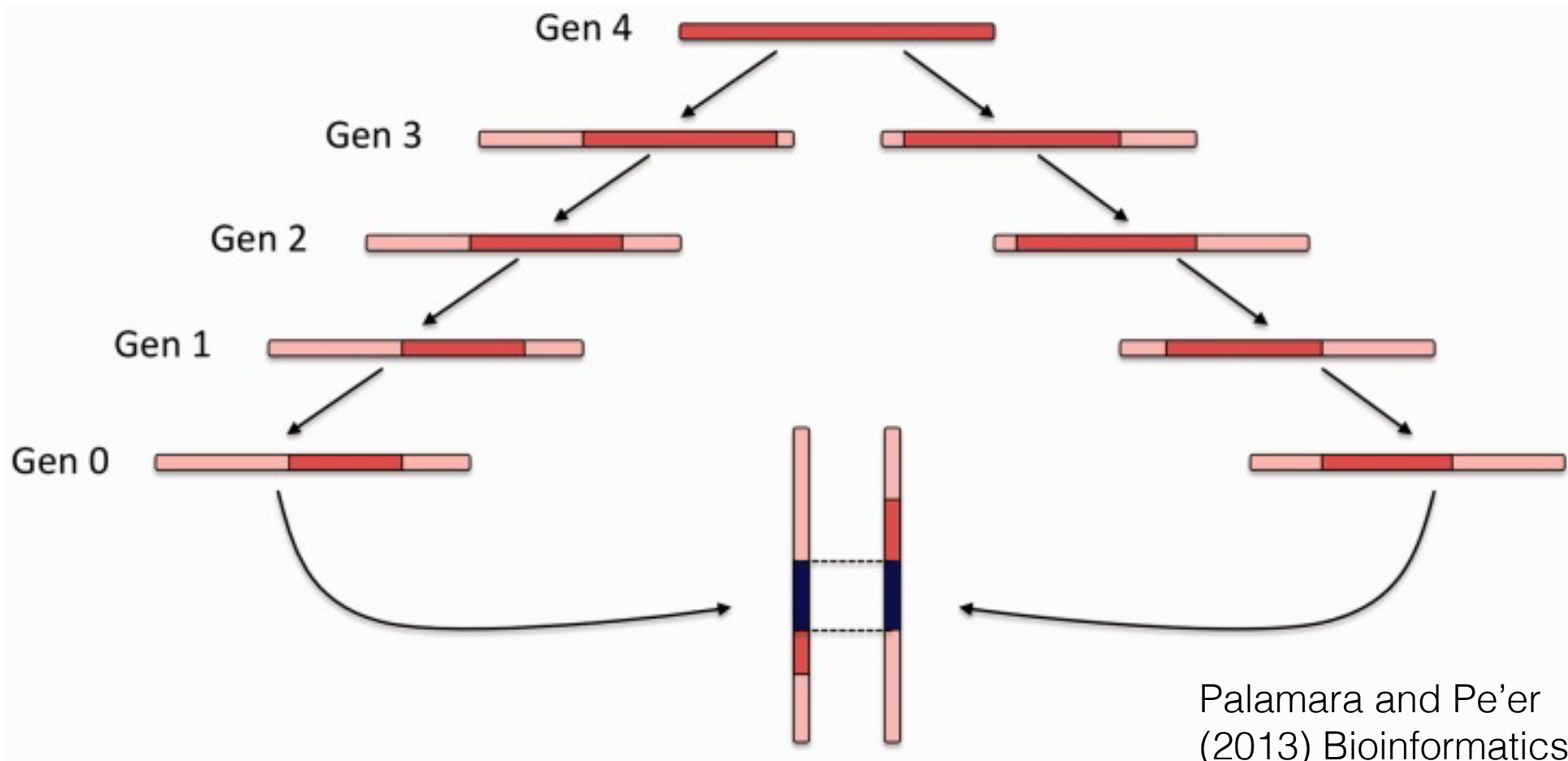
# Lung Cancer Prevalence in Europe



Boyle and Ferlay (2005) Annals of Oncology

# 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.
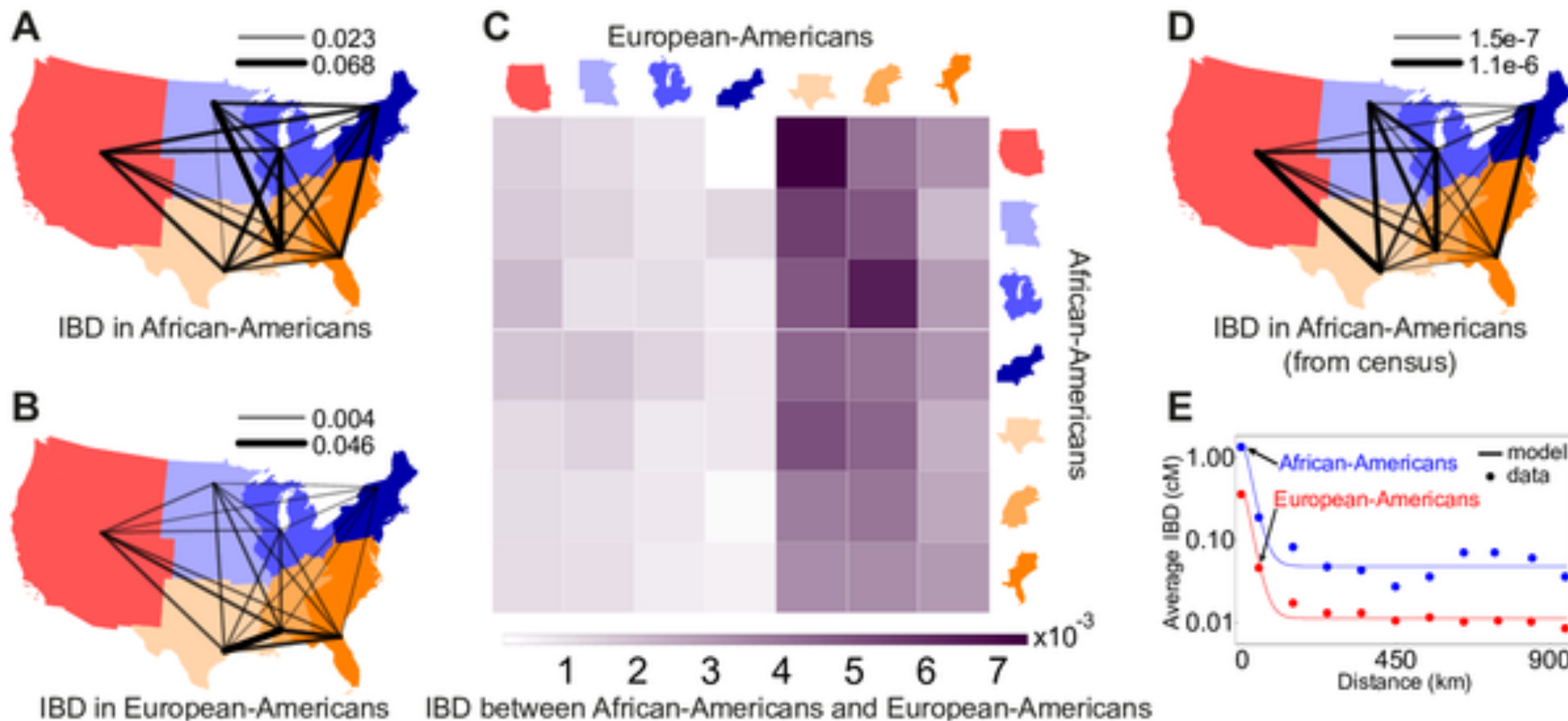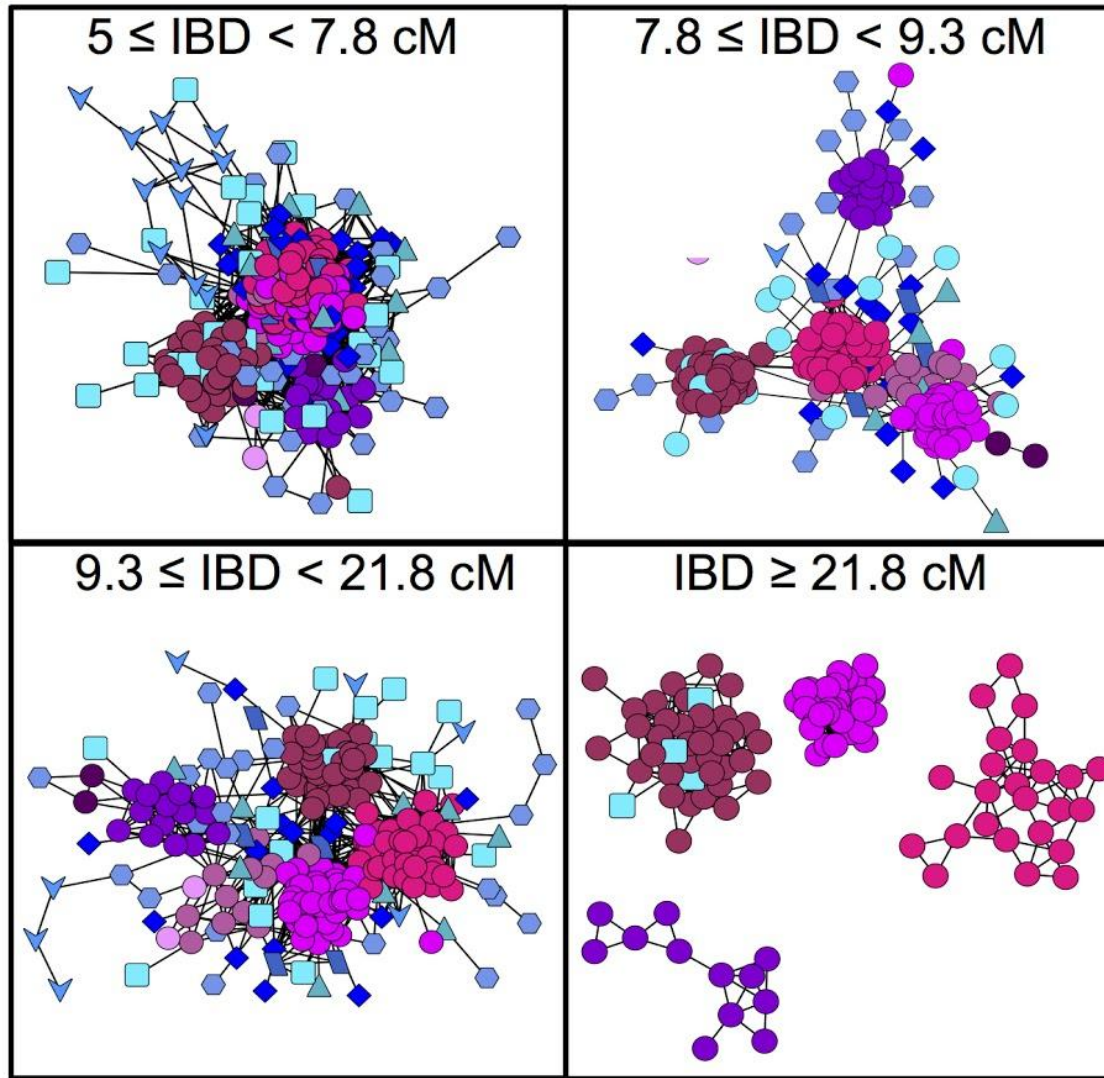


$$E[g|l] \cong \frac{3}{2 * l}$$

Baharian et al. (2016) PLoS Genet.

Palamara and Pe'er (2013) Bioinformatics
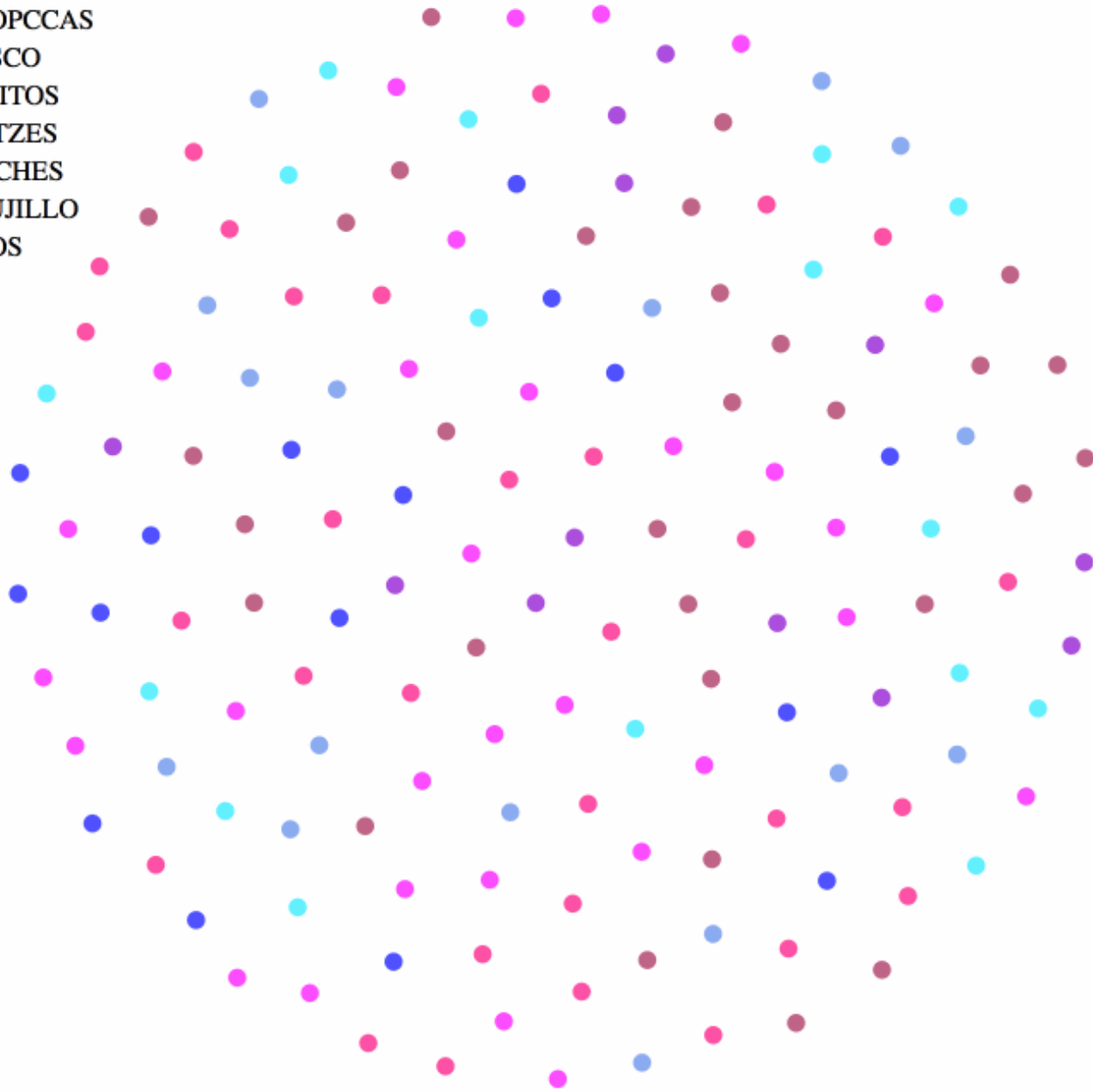
# Pairwise genetic relatedness across



A — 0.023 — 0.068
IBD in African-Americans

B — 0.004 — 0.046
IBD in European-Americans

C European-Americans
African-Americans
×10⁻³
1 2 3 4 5 6 7
IBD between African-Americans and European-Americans

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

E
Average IBD (cM)
1.00
African-Americans
European-Americans
0.10
0.01
— model
• data
0 450 900
Distance (km)

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
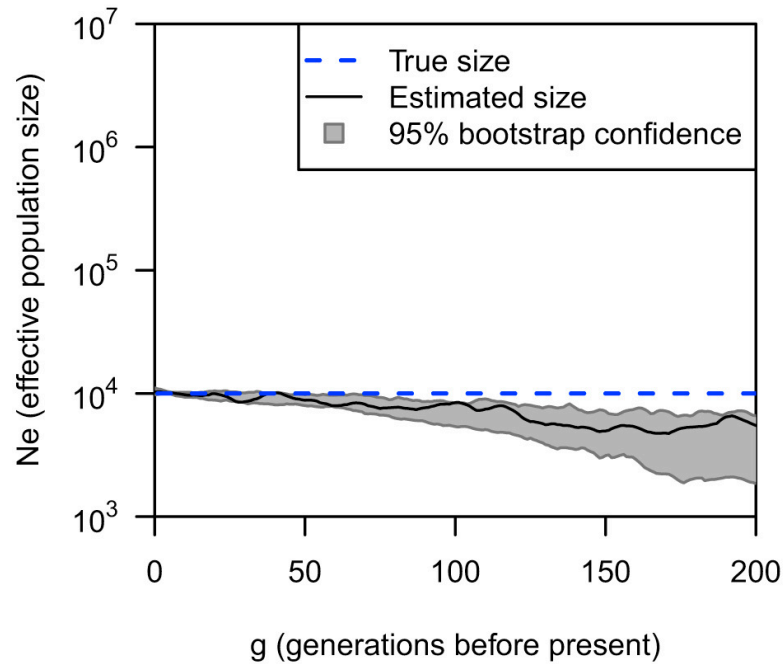
Harris et al. (2018) PNAS

**Legend:**
- 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.
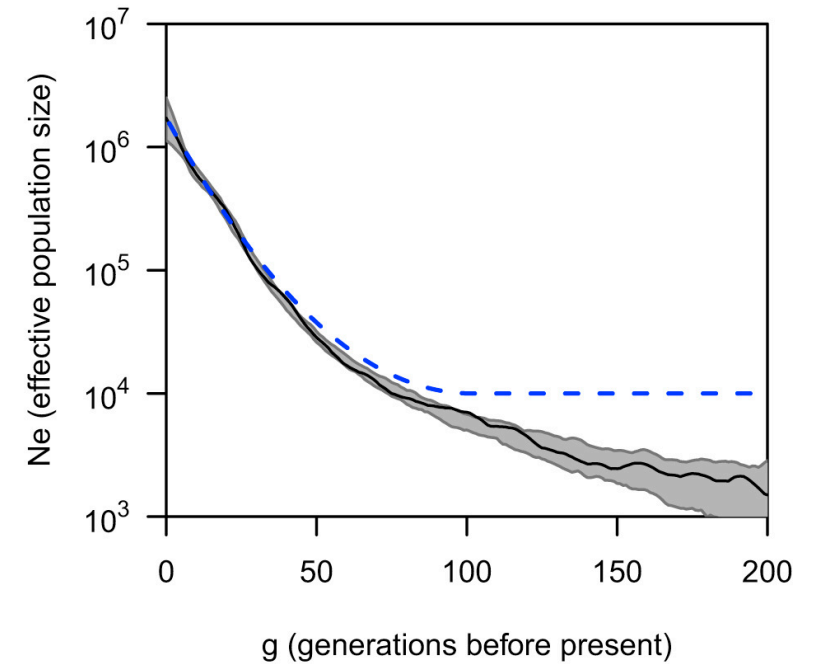


**Constant size: SNP array data**

**Exponential growth: SNP array data**

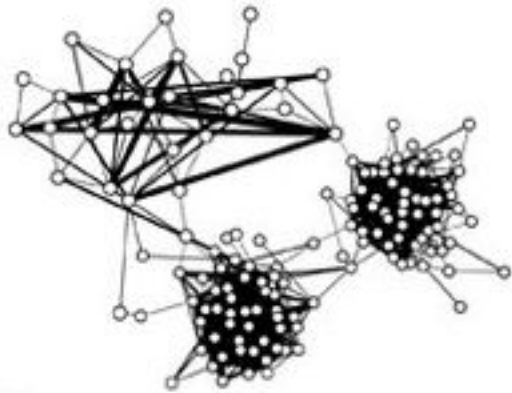**Super−exponential: SNP array data**

# IBD on a large scale



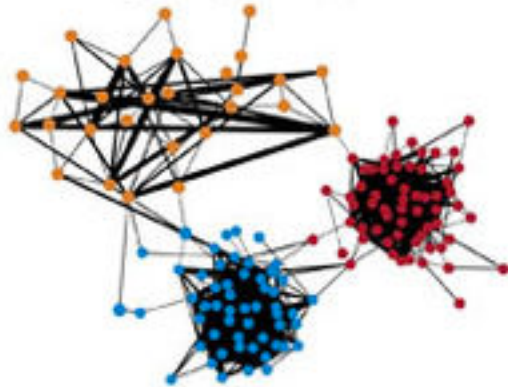**a** Construct network from IBD.
Join vertex pairs (genotyped samples) if IBD>12 cM. Edge weights are a function of total detected IBD.
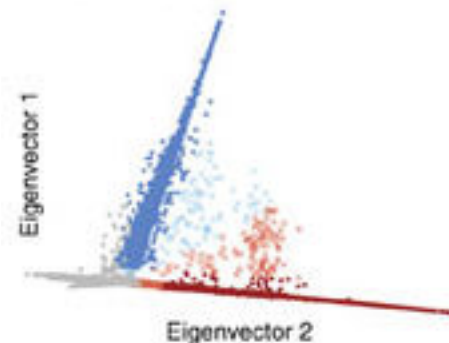
**b** Detect network clusters.
Recursively identify disjoint sets that maximize the modularity of the network. (Here one level of clustering hierarchy is shown.)
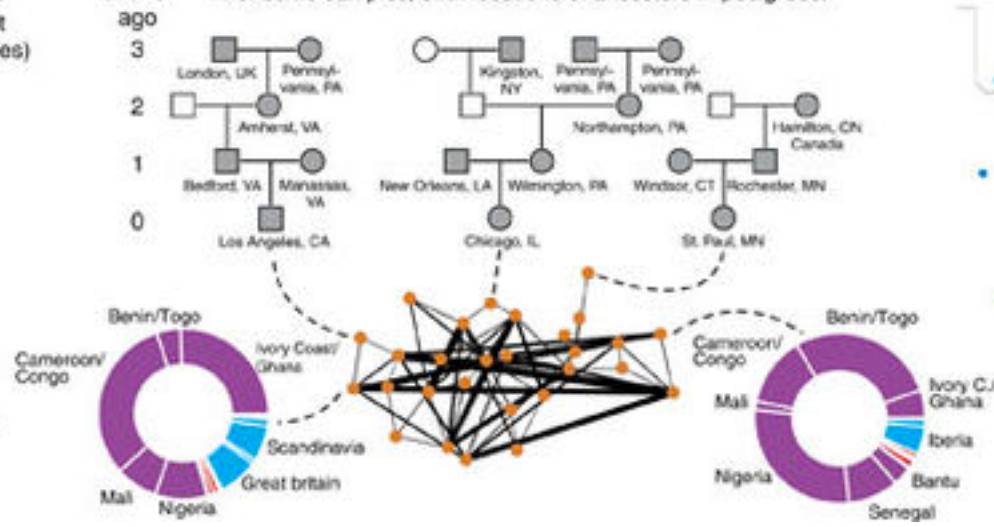
**c** Identify subsets of the clusters that separate in the spectral embedding.
Spectral embedding is computed from eigen-decomposition of Laplacian matrix. In the plot below, we identify "stable subsets" (filled circles) of the blue and red clusters.

Eigenvector 1
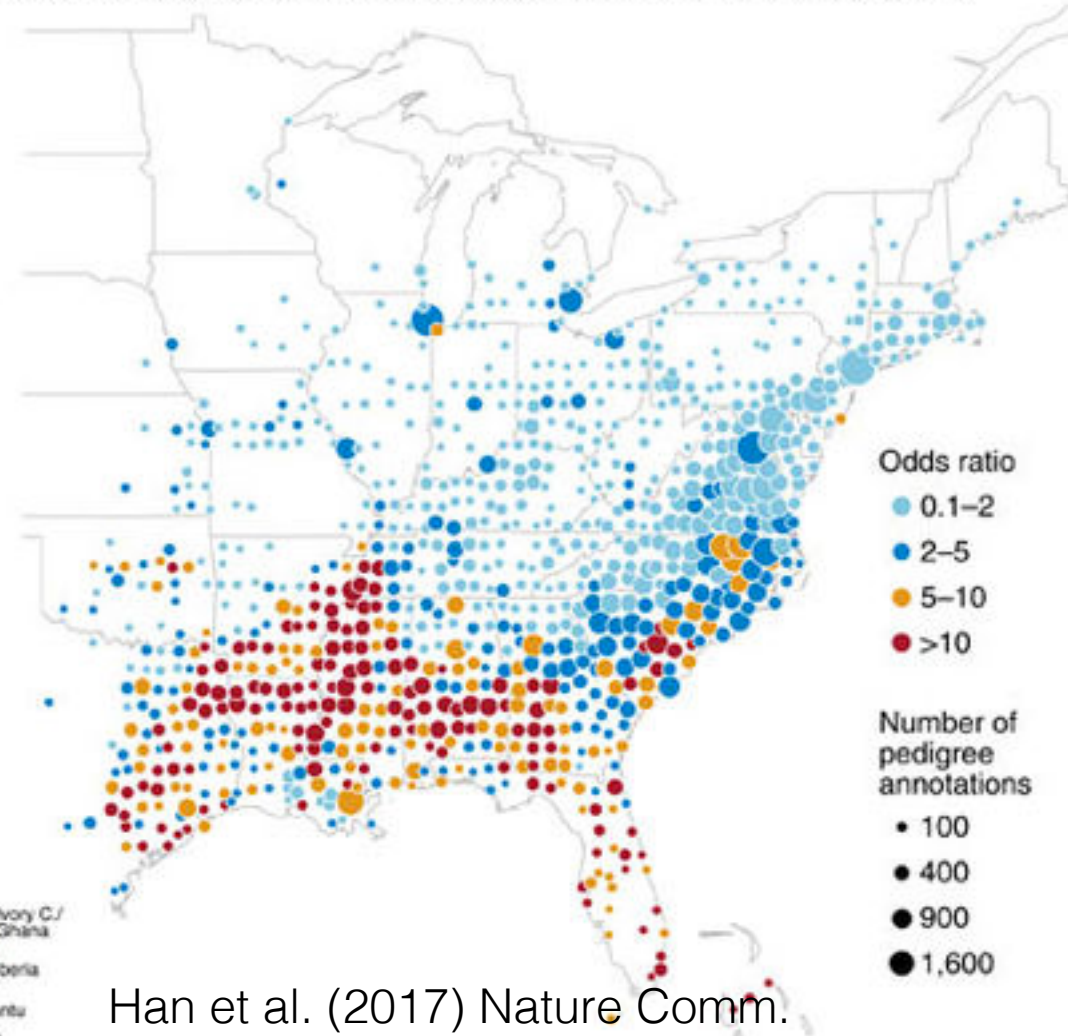Eigenvector 2

**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.

Generations ago
3
2
1
0

London, UK   Pennsylvania, PA
Amherst, VA
Bedford, VA   Manassas, VA
Los Angeles, CA

Kingston, NY   Pennsylvania, PA   Pennsylvania, PA
Northampton, PA
New Orleans, LA   Wilmington, PA
Chicago, IL

Hamilton, ON Canada
Windsor, CT   Rochester, MN
St. Paul, MN

Benin/Togo
Cameroon/Congo
Mali   Nigeria
Ivory Coast/Ghana
Scandinavia
Great britain

Benin/Togo
Cameroon/Congo
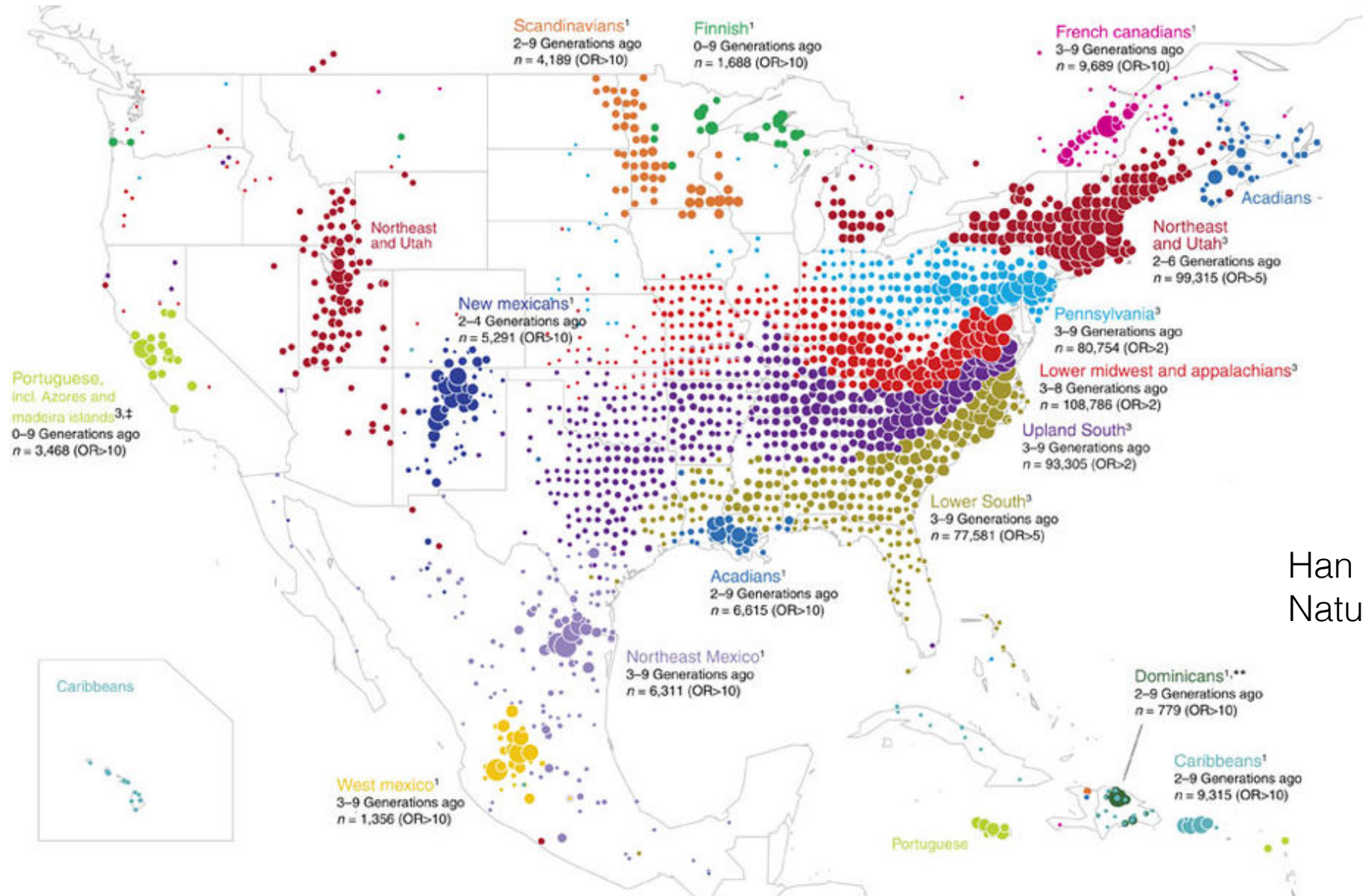Mali
Nigeria   Senegal
Ivory C./Ghana
Iberia
Bantu

**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.
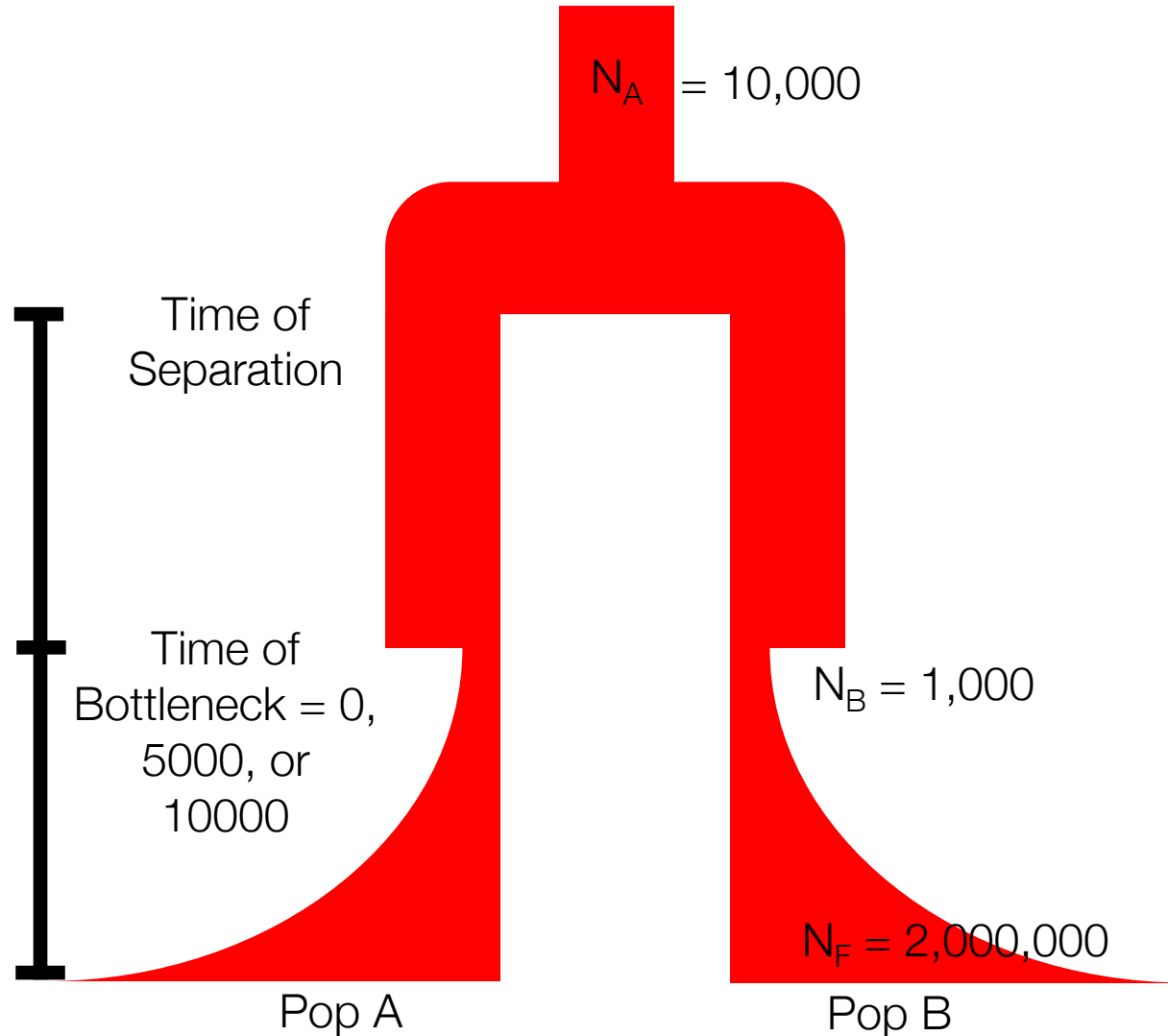
Odds ratio
● 0.1–2
● 2–5
● 5–10
● >10

Number of pedigree annotations
• 100
● 400
● 900
● 1,600

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)

Acadians

Northeast and Utah

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.

# Rare VS Common:
## Population Structure Simulations



$N_A = 10,000$

Time of Separation

Time of Bottleneck = 0, 5000, or 10000

$N_B = 1,000$

$N_F = 2,000,000$

Pop A

Pop B

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

# 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.

Chart: Expected Info. Gain vs Time of Separation, with Common (blue) and Rare (red) curves. Y-axis (Expected Info. Gain) ranges 0.000–0.008; X-axis (Time of Separation) ranges 2000–20000.

# 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?



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

# 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.

Legend:
- Unknown
- Ashkenazi
- Moroccan
- Sephardic
- Azerbaijan
- Bene Israel
- Cochin
- Ethiopian
- Georgia
- Iranian
- Iraq
- Uzbekistan
- Yemen

PC1: 0.00833 / −0.03529
PC2: 0.022089 / −0.000175
PC3: 0.0325 / −0.0530
PC4: 0.0114 / −0.0378
PC5: 0.0229 / −0.0023

# Population Average PCA with More Axes



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

# Trans-Omics for Precision Medicine (TOPMed) Cohorts

- N ≅ 18K
- This data freeze has 15 cohorts, each with 100s of samples
- Predominantly African, Latino, and European American
  - Samoa
  - Amish
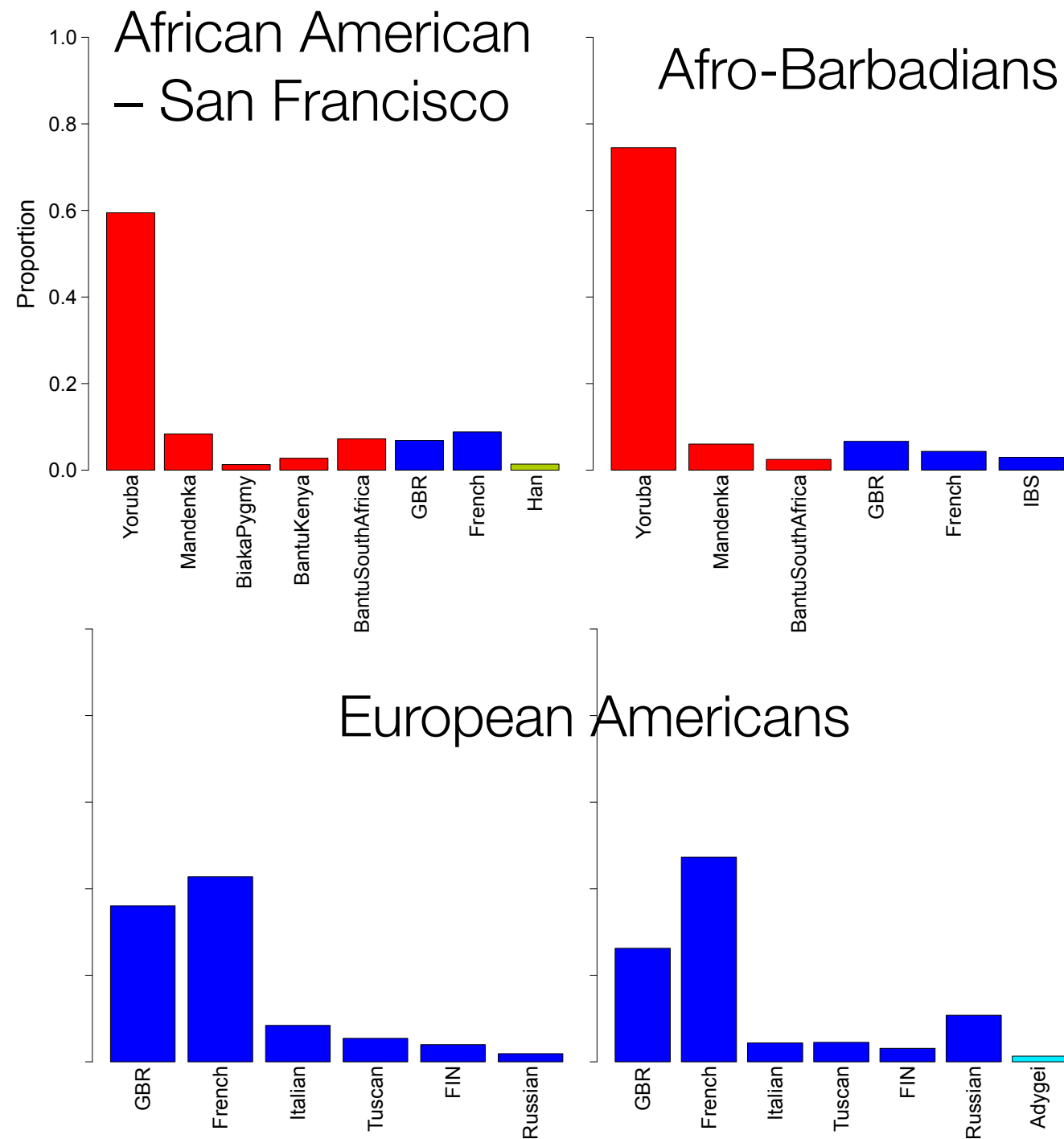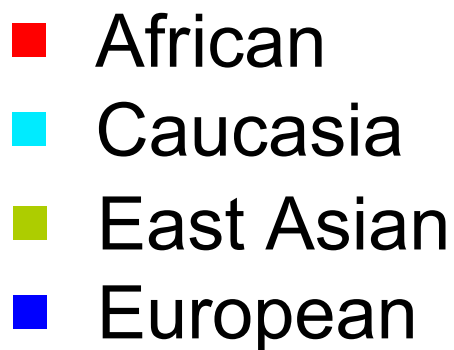- All are well characterized for heart, lung, blood, and sleep phenotypes

# Rare variant sharing across cohorts

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

# Rare variant sharing across cohorts

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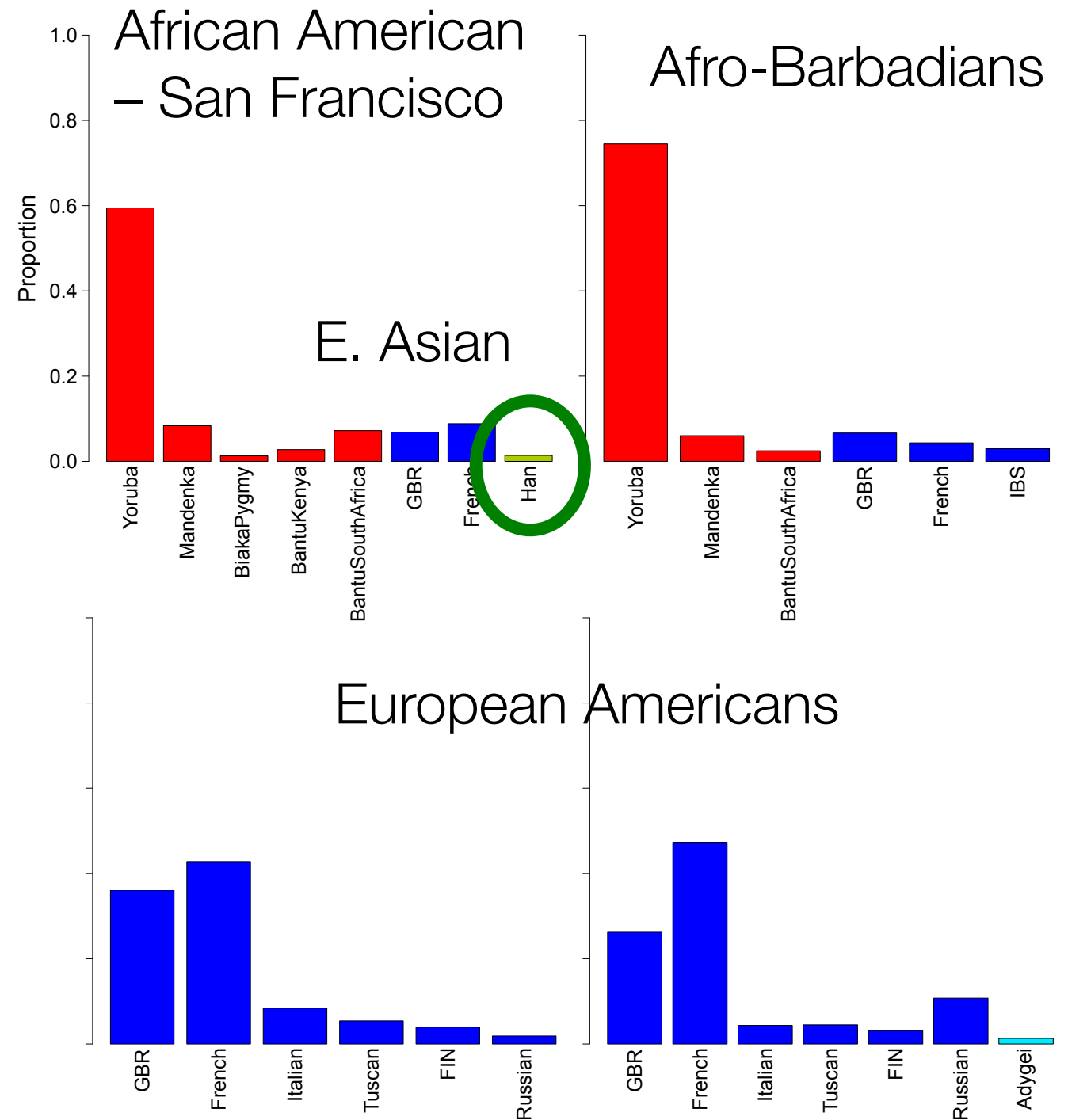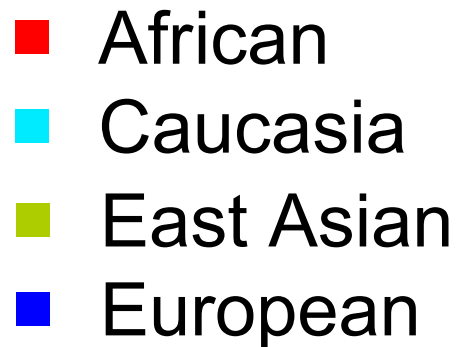- Corrected for:
  - sample size
  - Genome-wide heterozygosity

# Rare variant sharing across cohorts

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

# Rare variant sharing across cohorts

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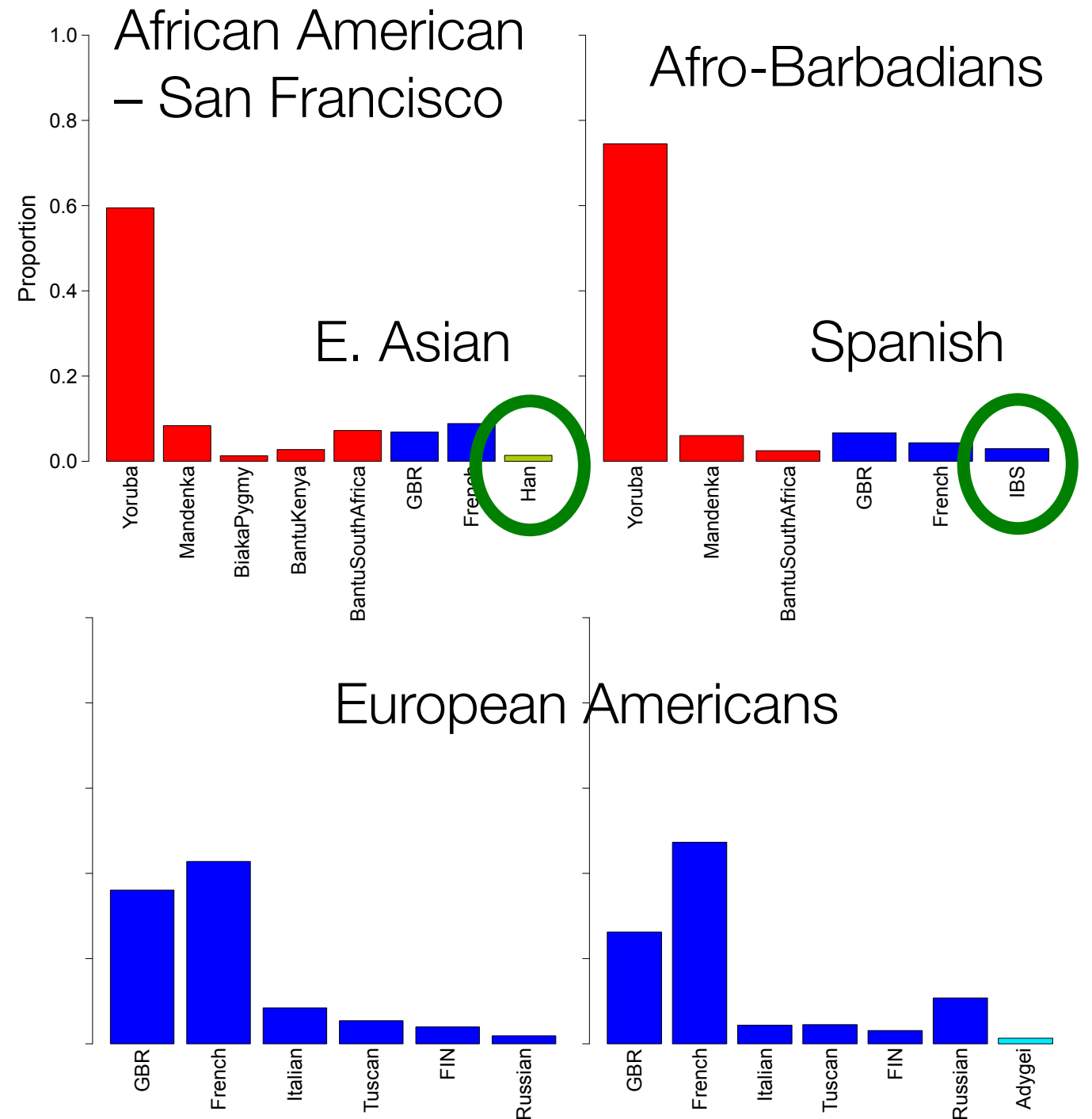  - sample size
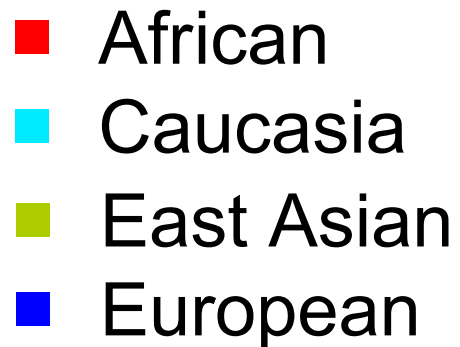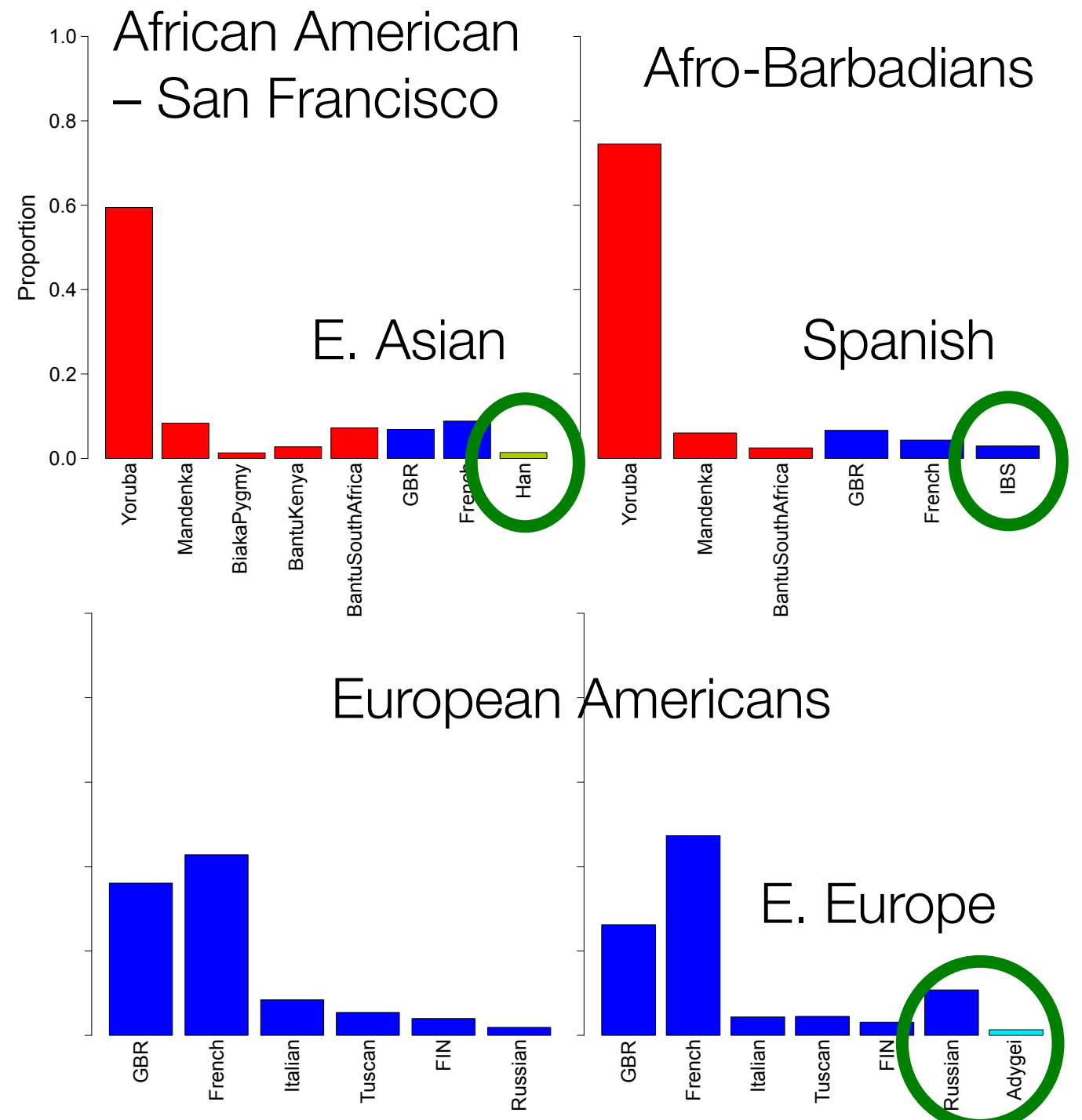  - Genome-wide heterozygosity

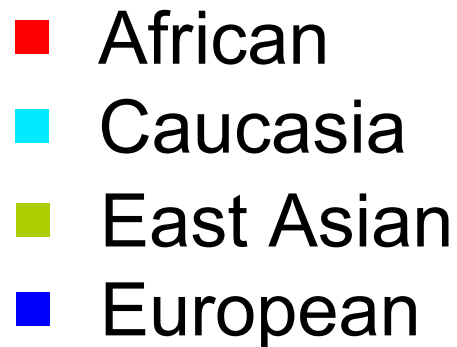fineStructure analysis of genome-wide ancestry

fineStructure analysis of genome-wide ancestry

fineStructure analysis of genome-wide ancestry

fineStructure analysis of genome-wide ancestry

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

African American – San Francisco
Afro-Barbadians
E. Asian
Spanish
European Americans
E. Europe

# 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
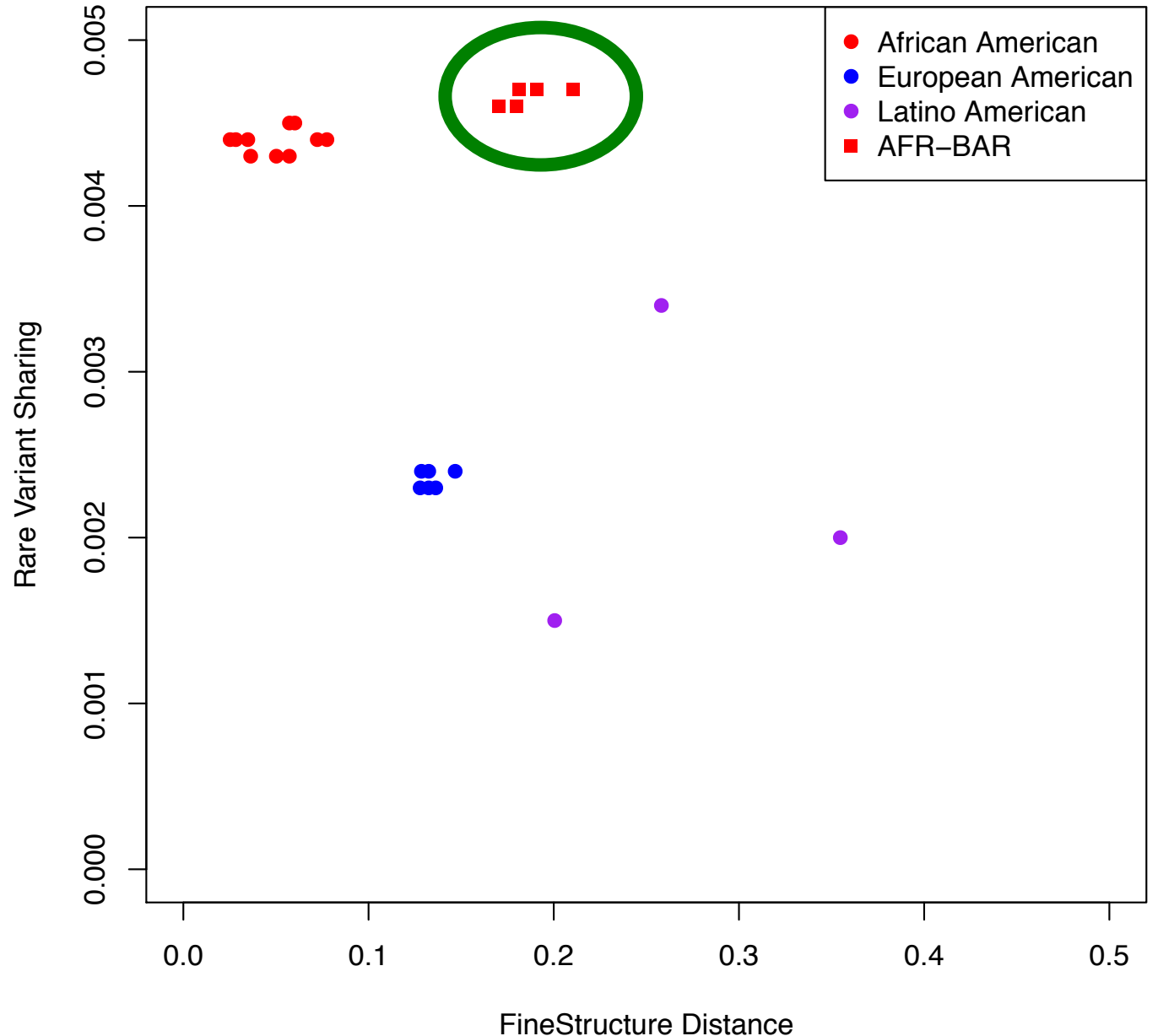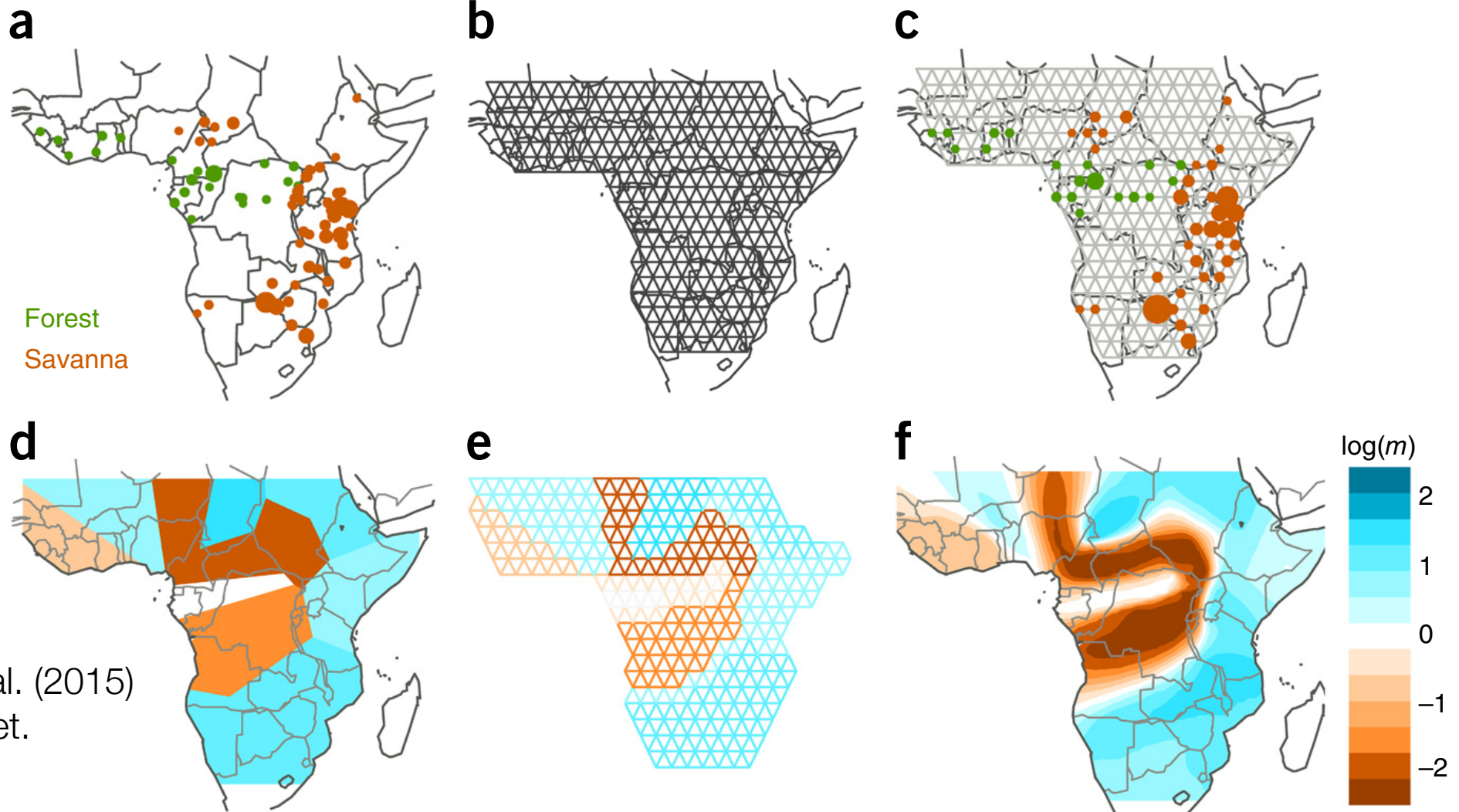
African American's have more homogeneous ancestral proportions

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

# Estimated Effective Migration Surfaces (EEMS)



Forest
Savanna

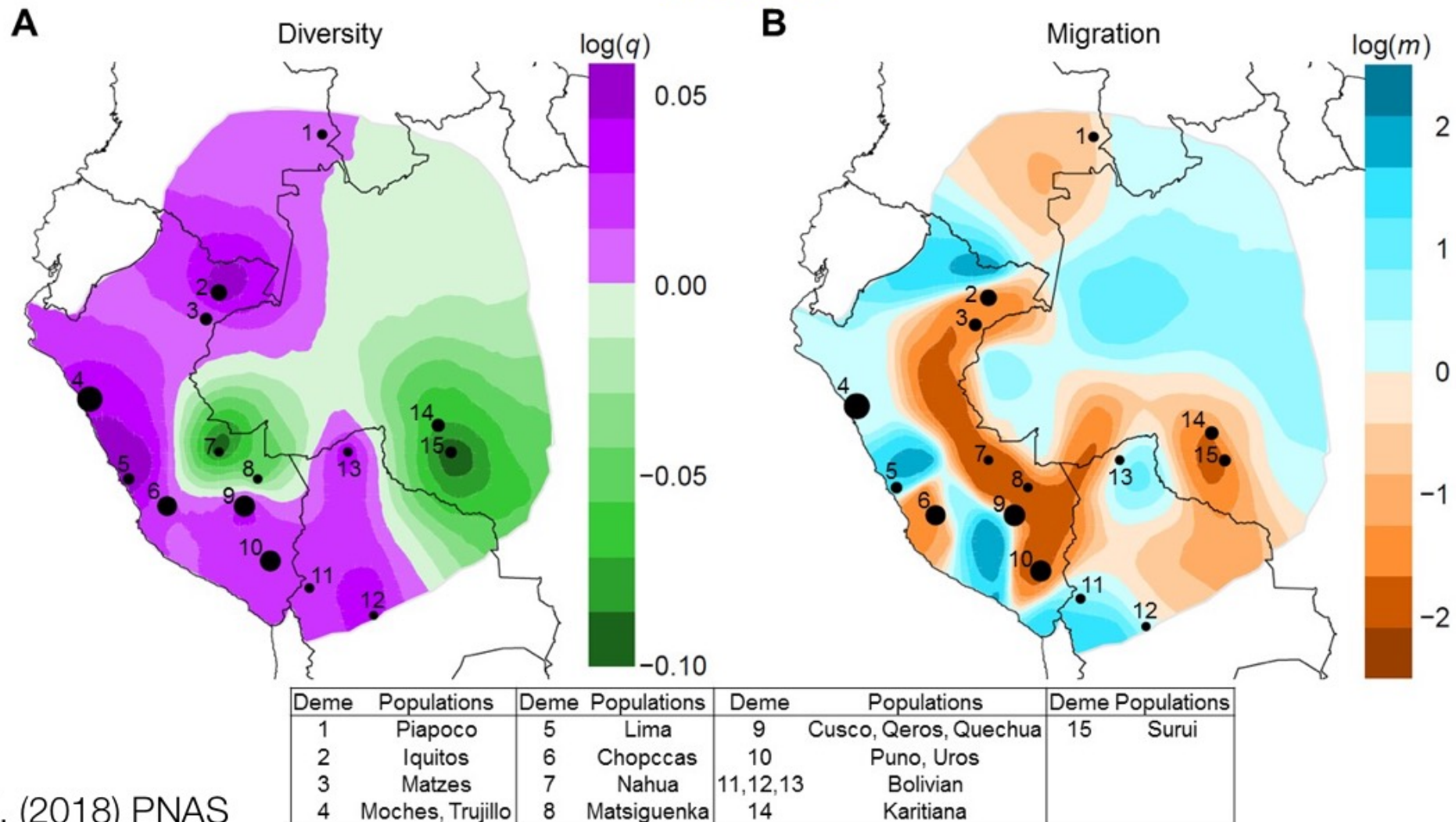Petkova et al. (2015)
Nature Genet.

$\log(m)$

# Assumptions: Stepping Stone Model

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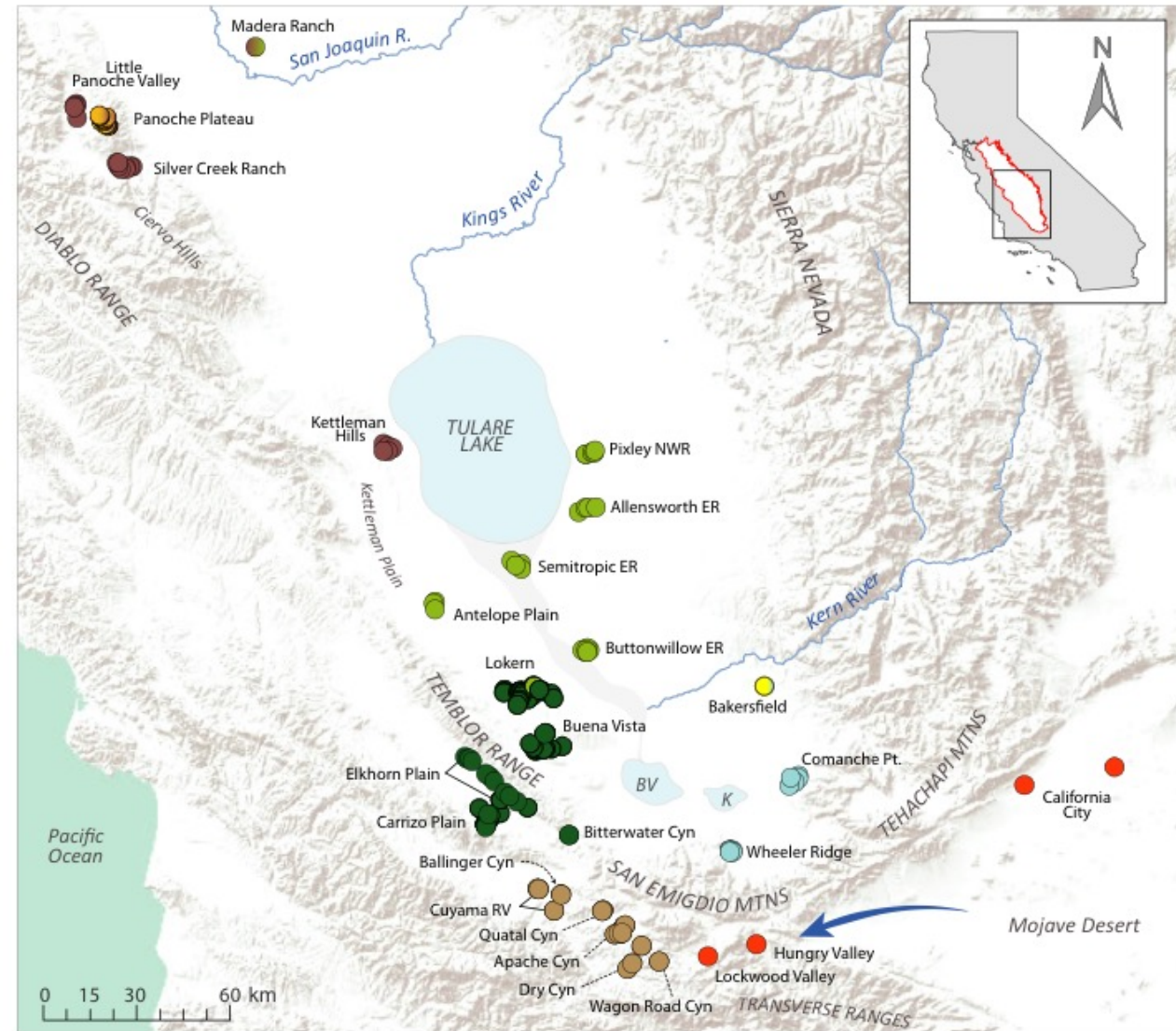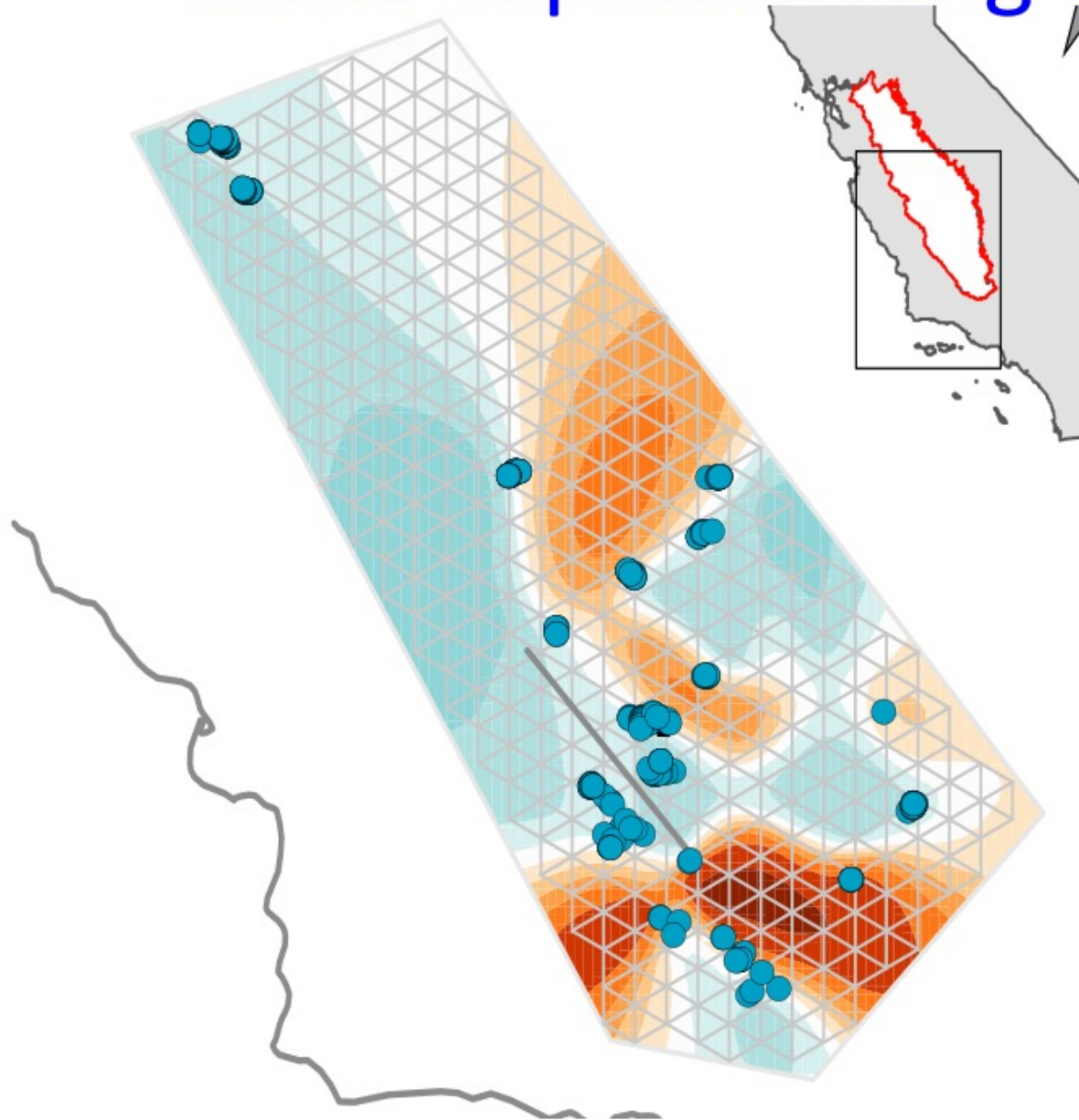- Migration rate between each deme is assumed to be equal



Kimura and Weiss (1964)

# EEMS: Migration and diversity within Peru



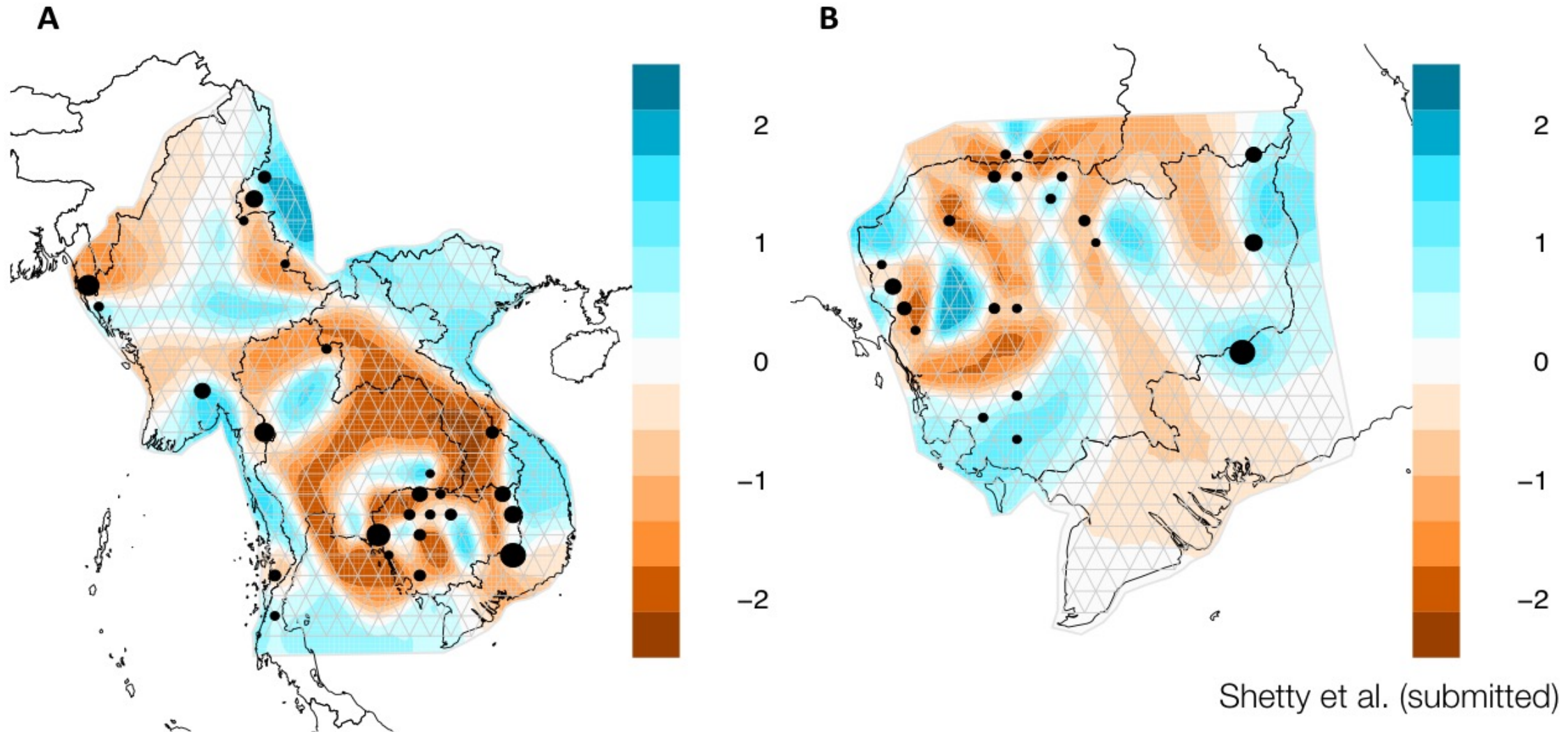| Deme | Populations | Deme | Populations | Deme | Populations | Deme | Populations |
|------|-------------|------|-------------|------|-------------|------|-------------|
| 1 | Piapoco | 5 | Lima | 9 | Cusco, Qeros, Quechua | 15 | Surui |
| 2 | Iquitos | 6 | Chopccas | 10 | Puno, Uros | | |
| 3 | Matzes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

Harris et al. (2018) PNAS

# EEMS captures long-term migration patterns



Richmond et al. (2015) Molecular Ecology

# EEMS in Malaria Parasites of South East Asia



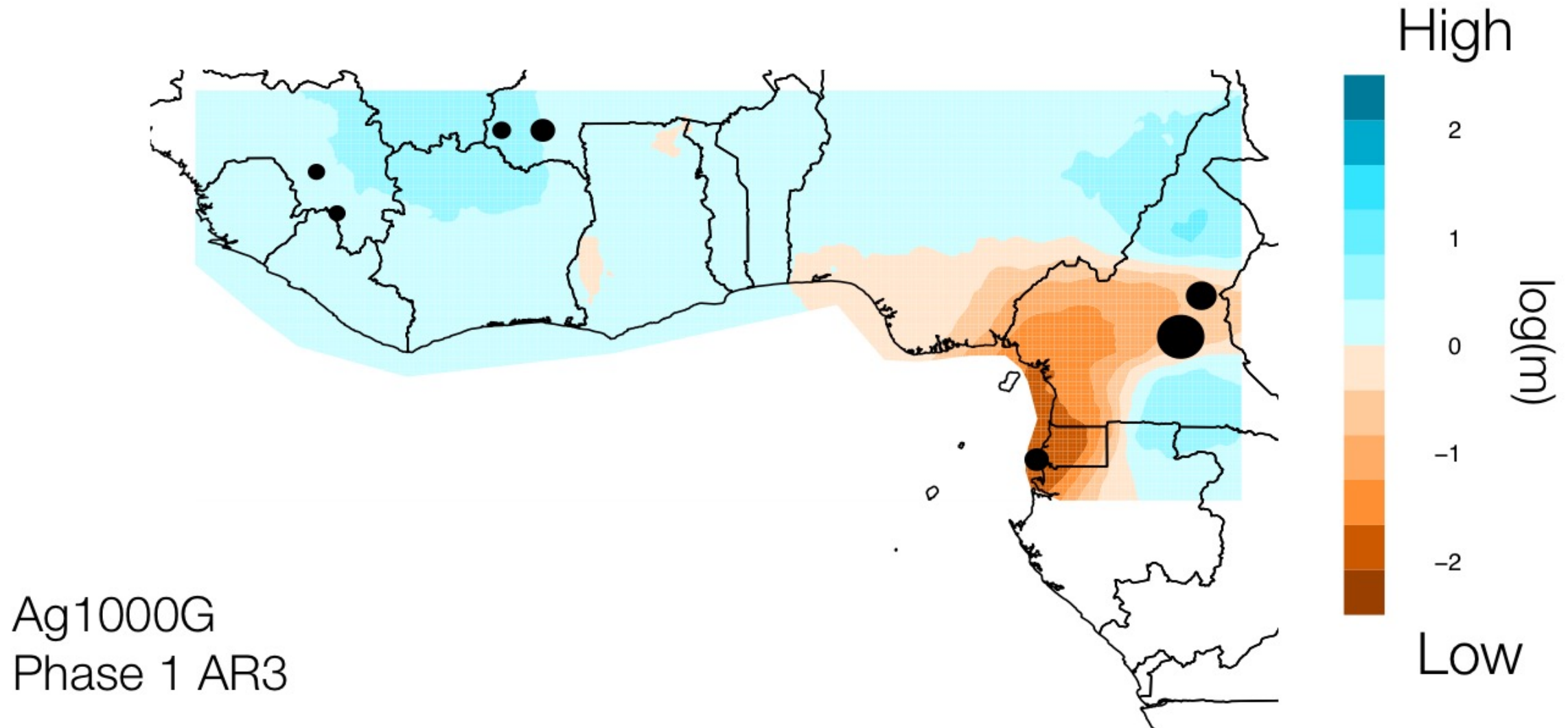Shetty et al. (submitted)

Application to Malaria Parasites in W. Africa

Pf3K Version 5.1

High

Low

log(m)

# Application to Mosquito in W. Africa

Ag1000G
Phase 1 AR3

log(m)

High

Low

# Robustness of Sampling on EEMS

# Concluding summary

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

- Identity-by-descent and sharing of rare variants are a powerful method 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