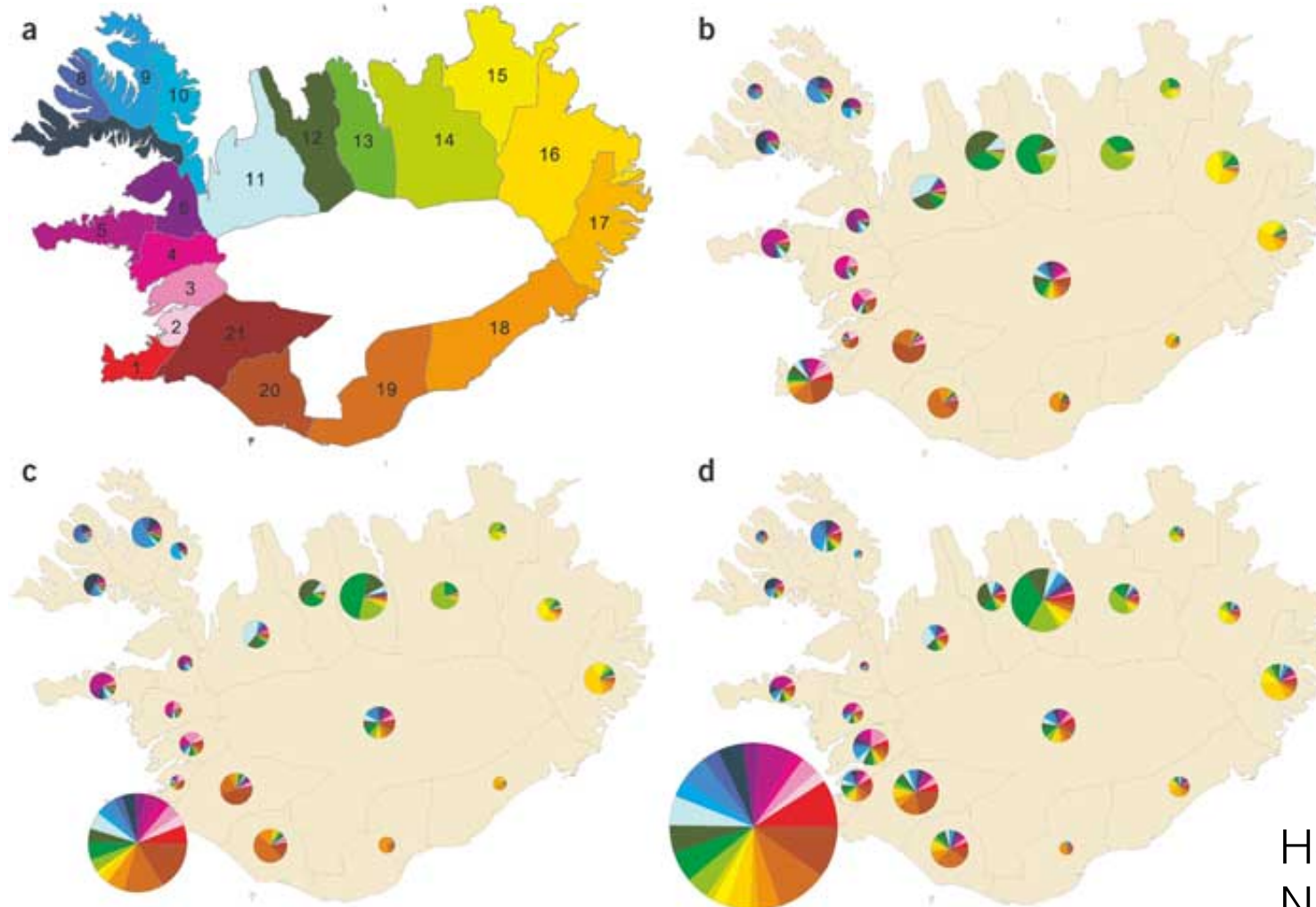# Cryptic Relatedness and fine scale population structure
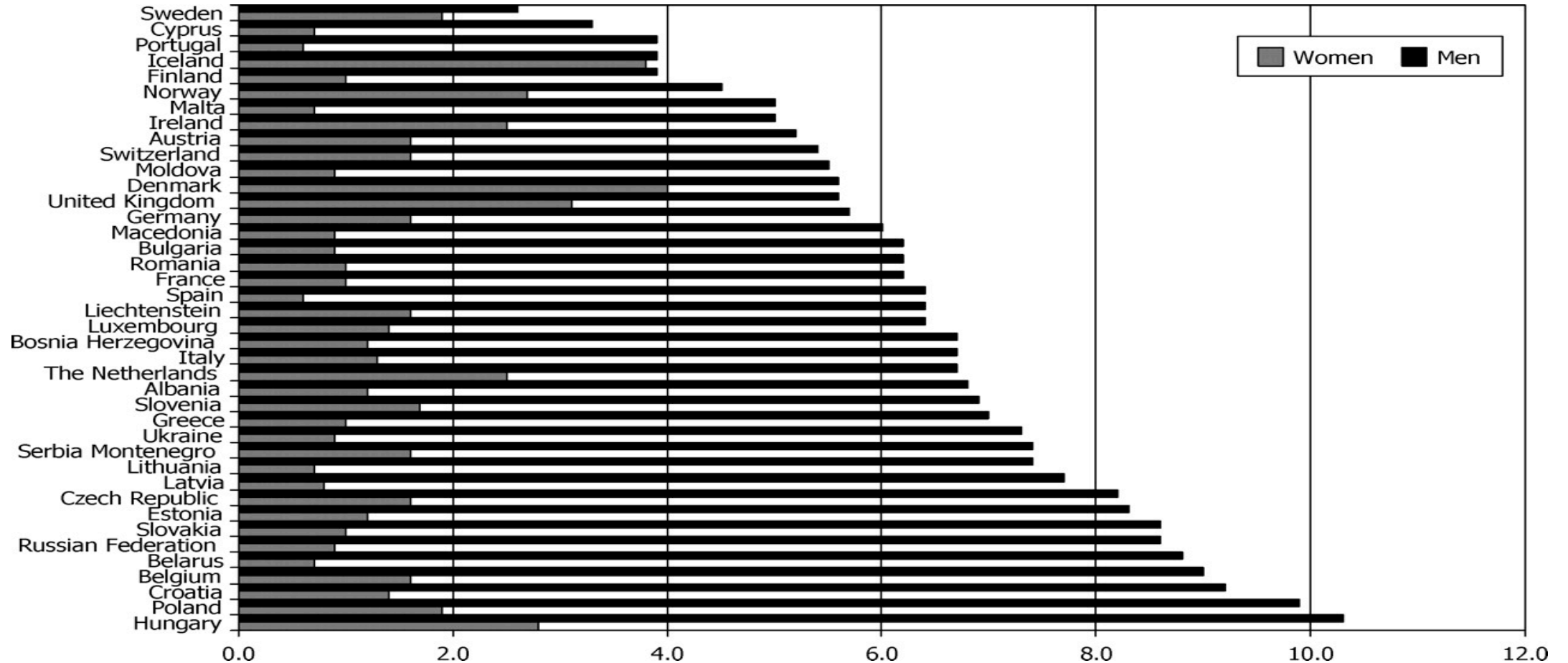
# Learning objectives

- 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.
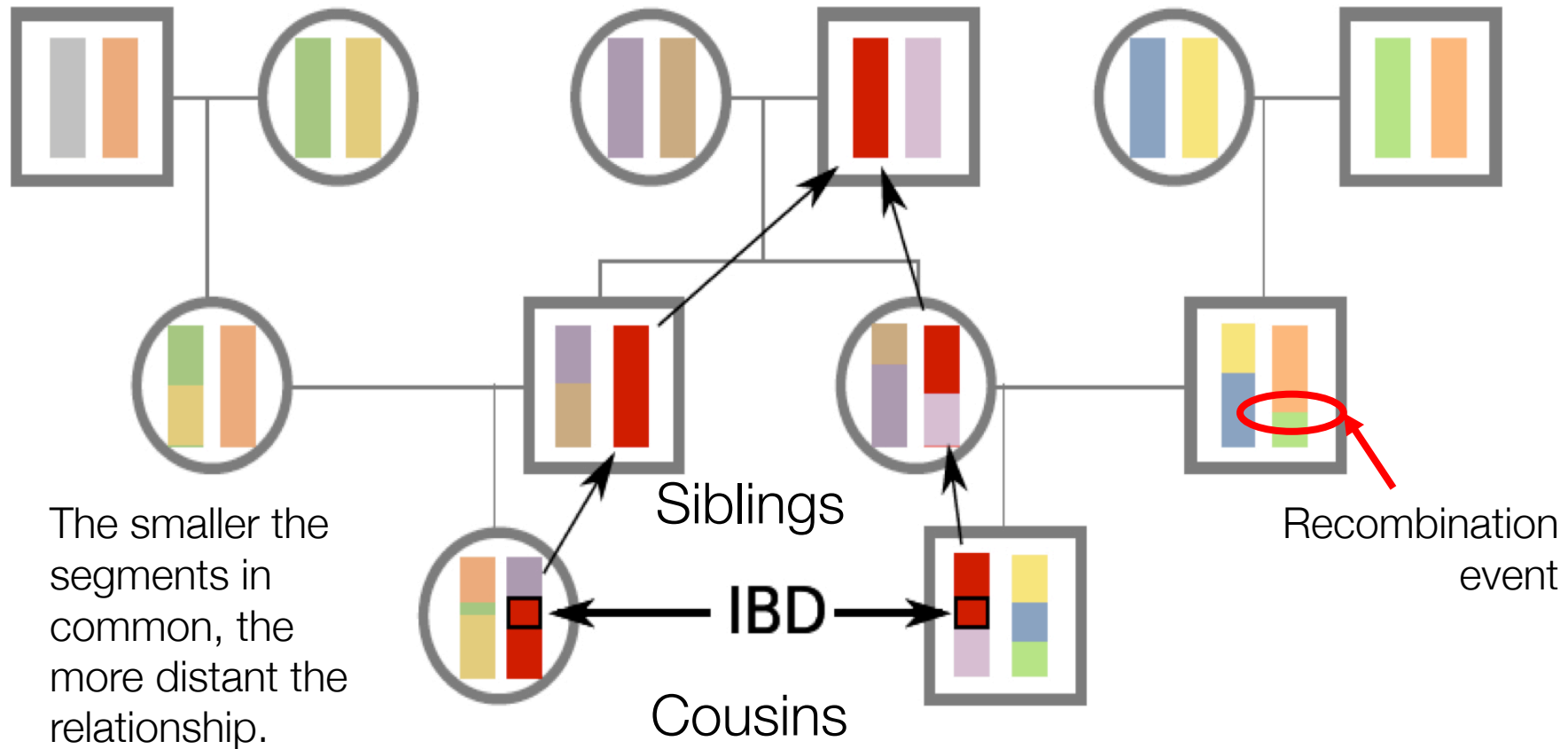
# Cryptic Population Structure
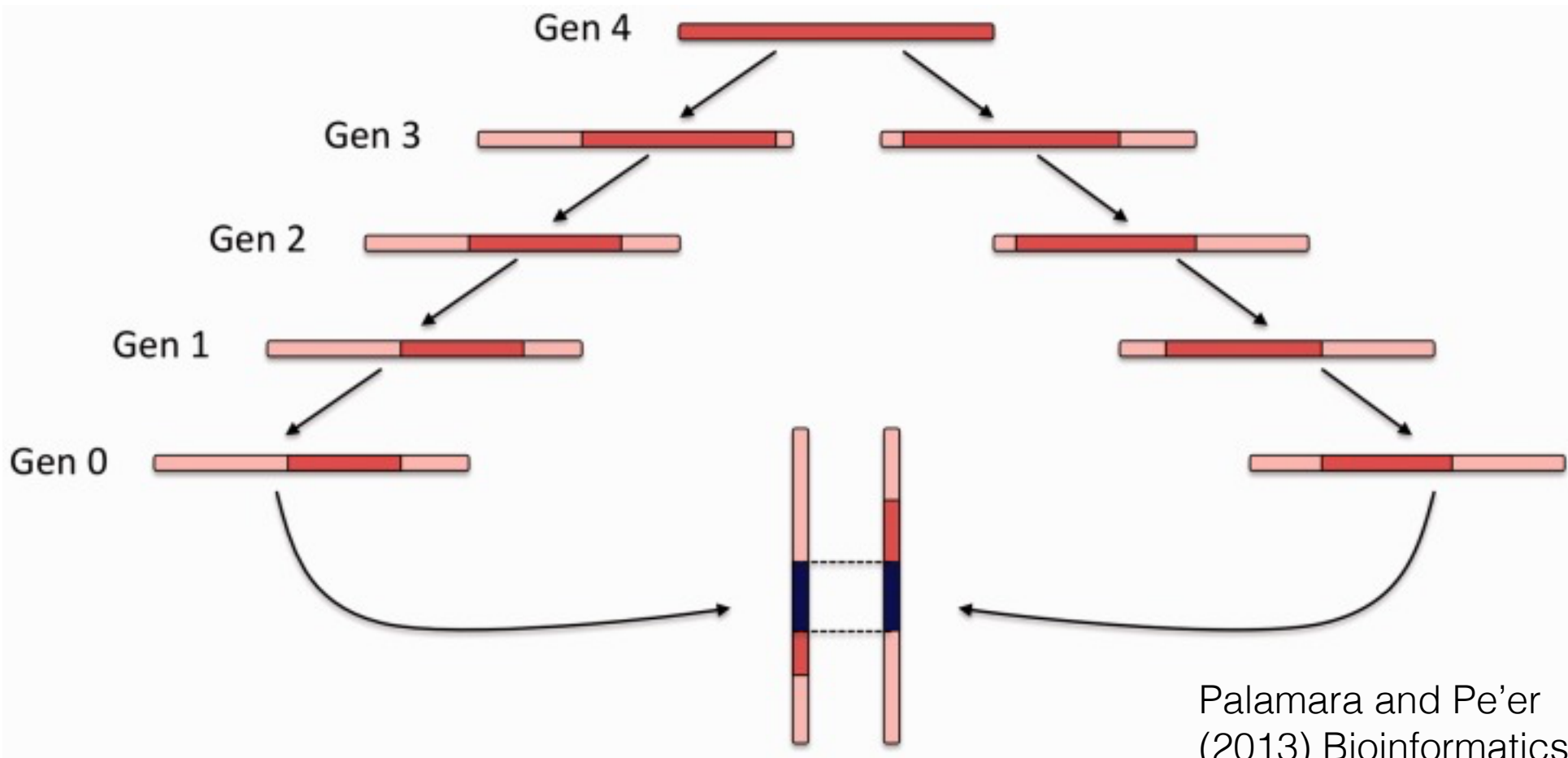


Helgason et al. (2004)
Nature Genet.

# 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

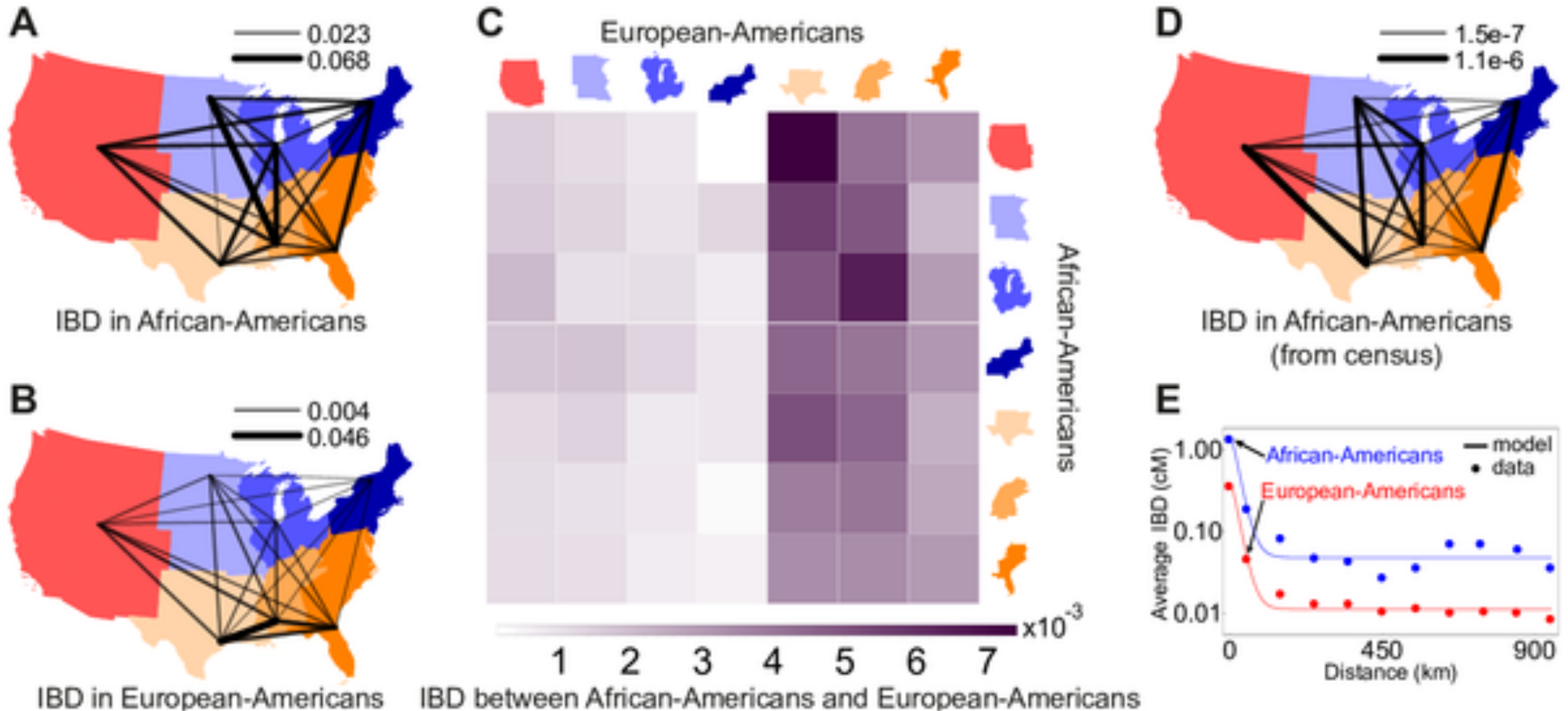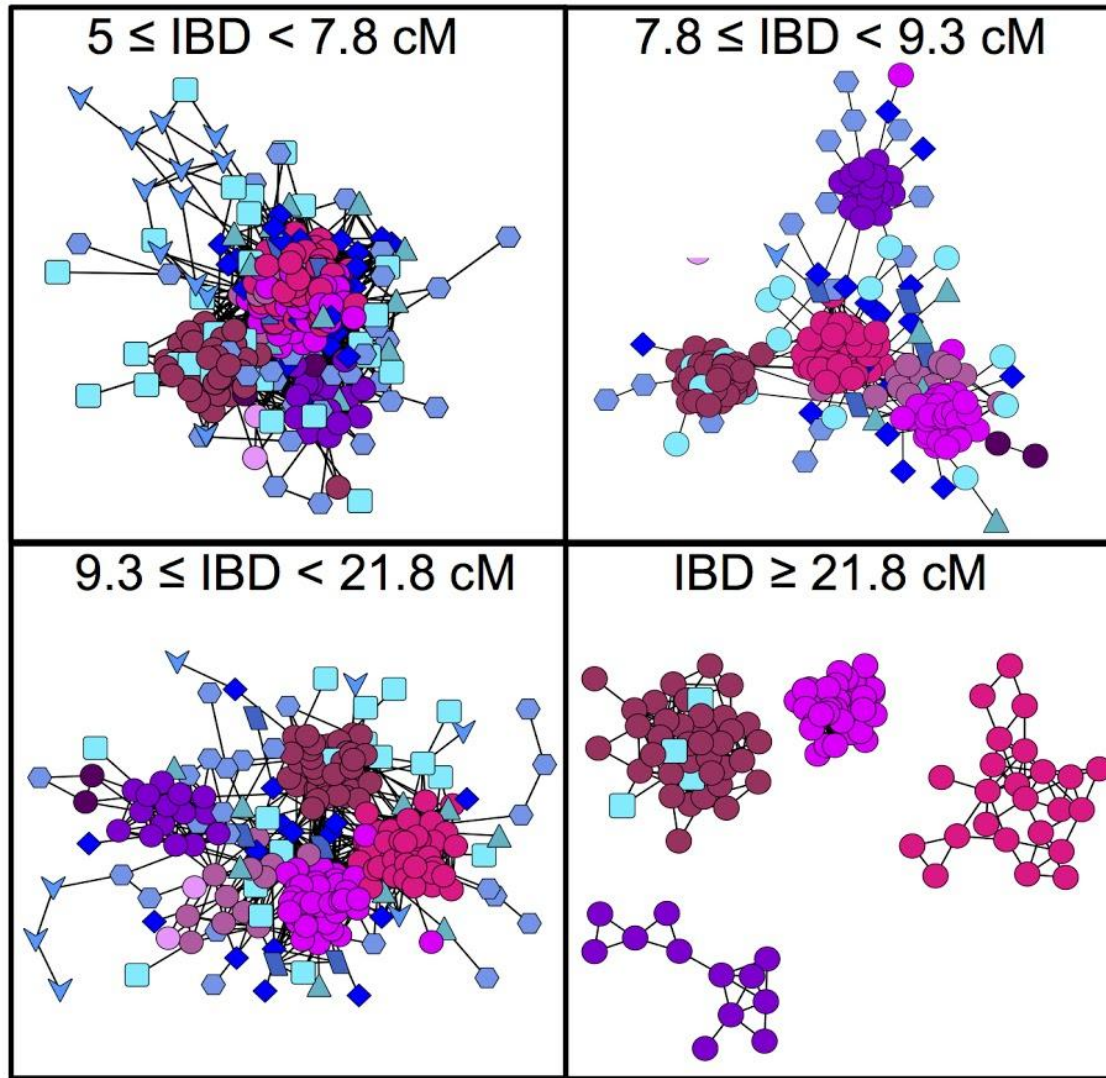# IBD length is correlated with historical relationships.



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

Baharian et al. (2016) PLoS Genet.

Palamara and Pe'er (2013) Bioinformatics

# Pairwise genetic relatedness across



**A** IBD in African-Americans — 0.023, 0.068

**B** IBD in European-Americans — 0.004, 0.046

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

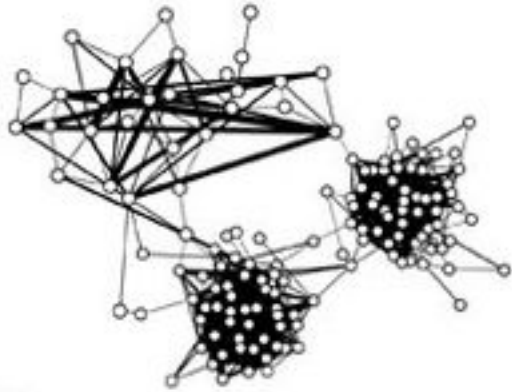Identity-by-descent as a means to look at fine-scale structure over time

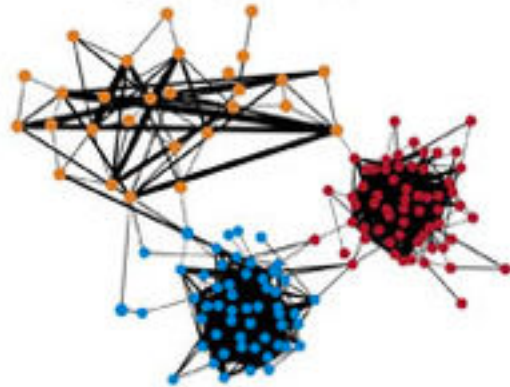Harris et al. (submitted)

# IBD on a large scale



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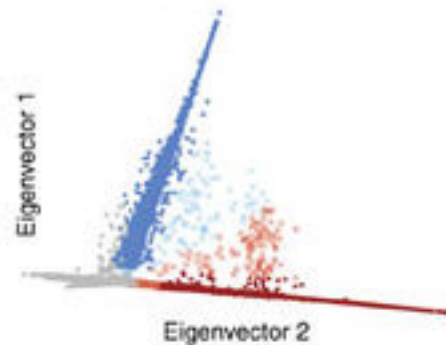Join vertex pairs (genotyped samples) if IBD>12 cM. Edge weights are a function of total detected IBD.

**b** Detect network clusters.
Recursively identify disjoint sets that maximize the modularity of the network. (Here one level of clustering hierarchy is shown.)
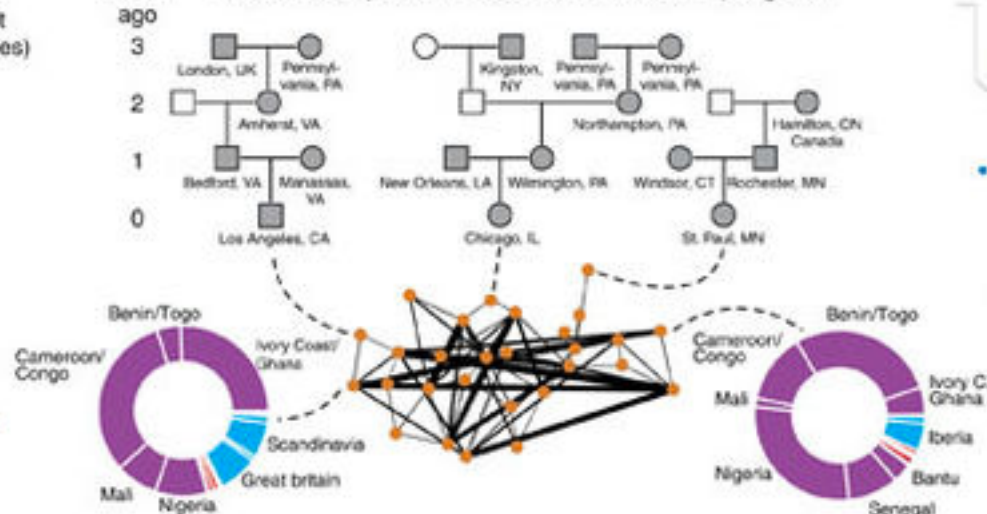
**c** Identify subsets of the clusters that separate in the spectral embedding.
Spectral embedding is computed from eigen-decomposition of Laplacian matrix. In the plot below, we identify "stable subsets" (filled circles) of the blue and red clusters.
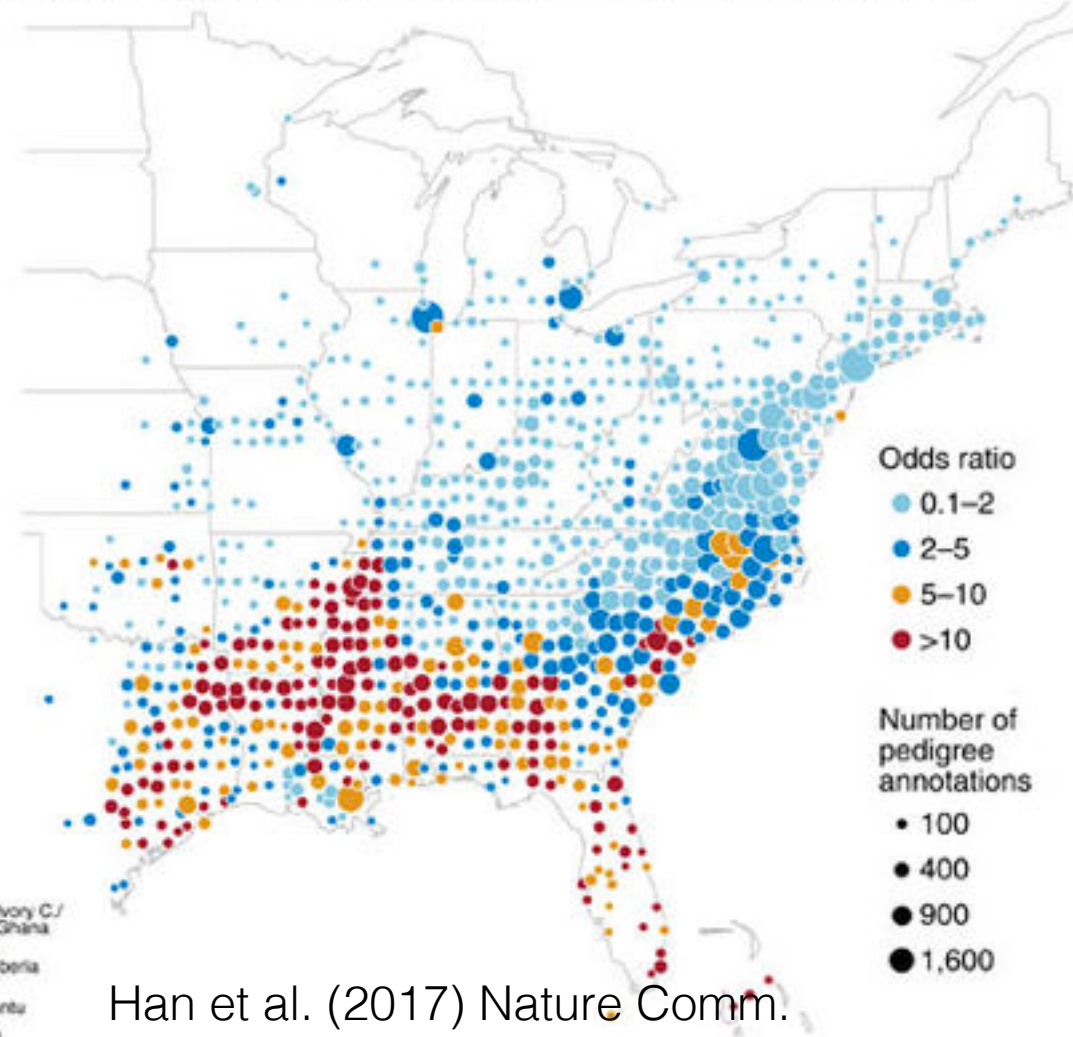
**d** Annotate each cluster with two kinds of data:
• In all samples, global admixture of 20 populations (donut charts);
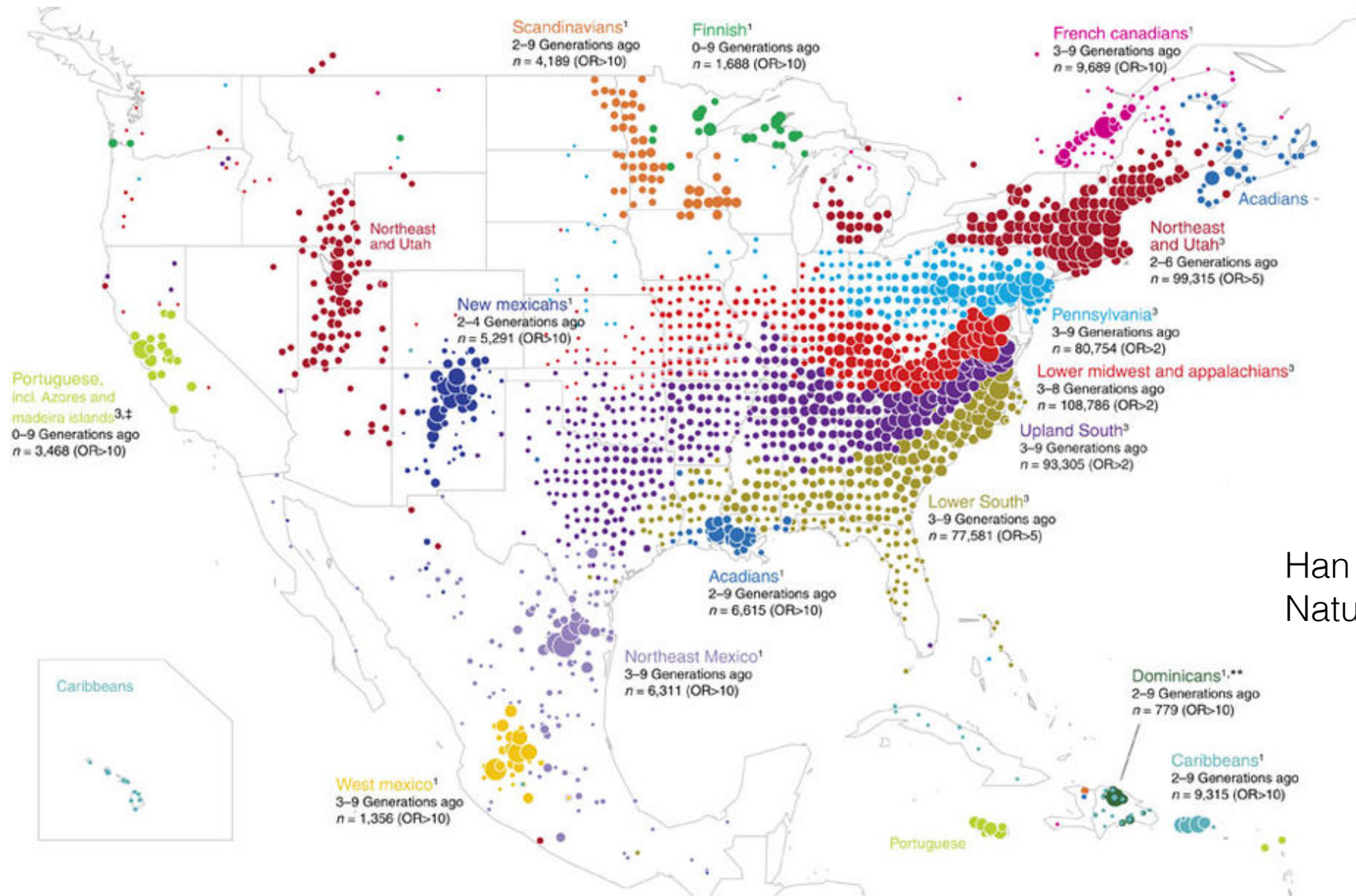• For some samples, birth locations of ancestors in pedigrees.

**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.
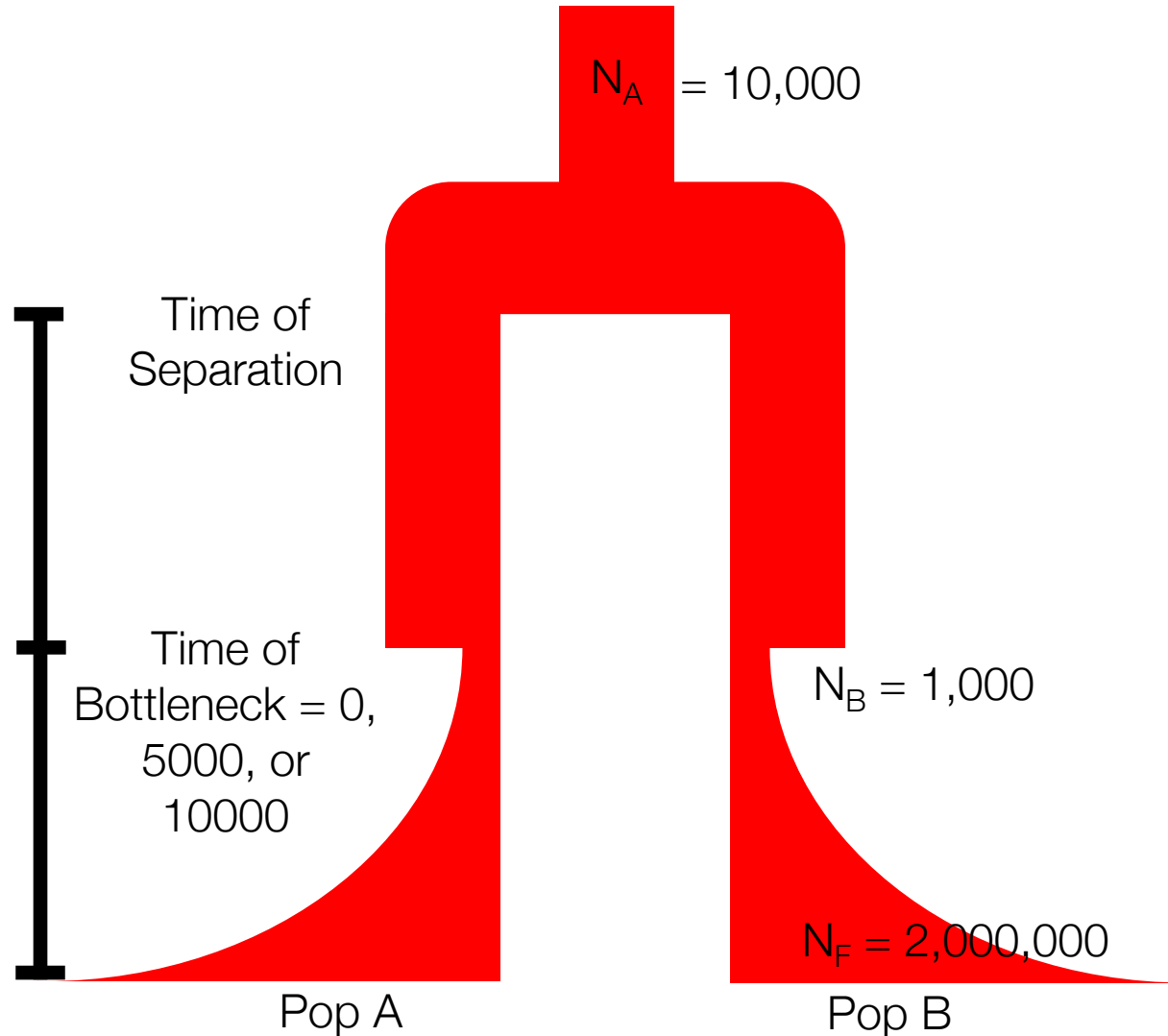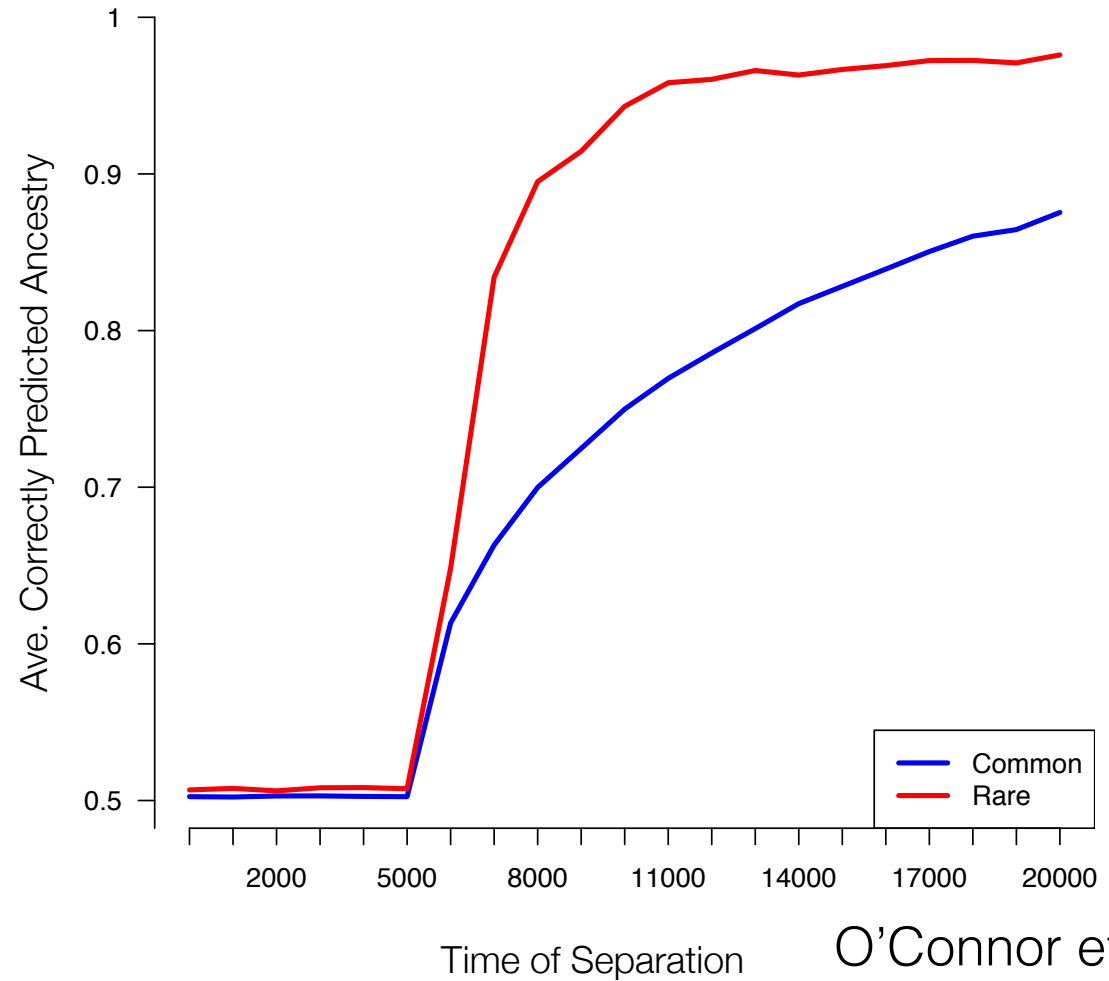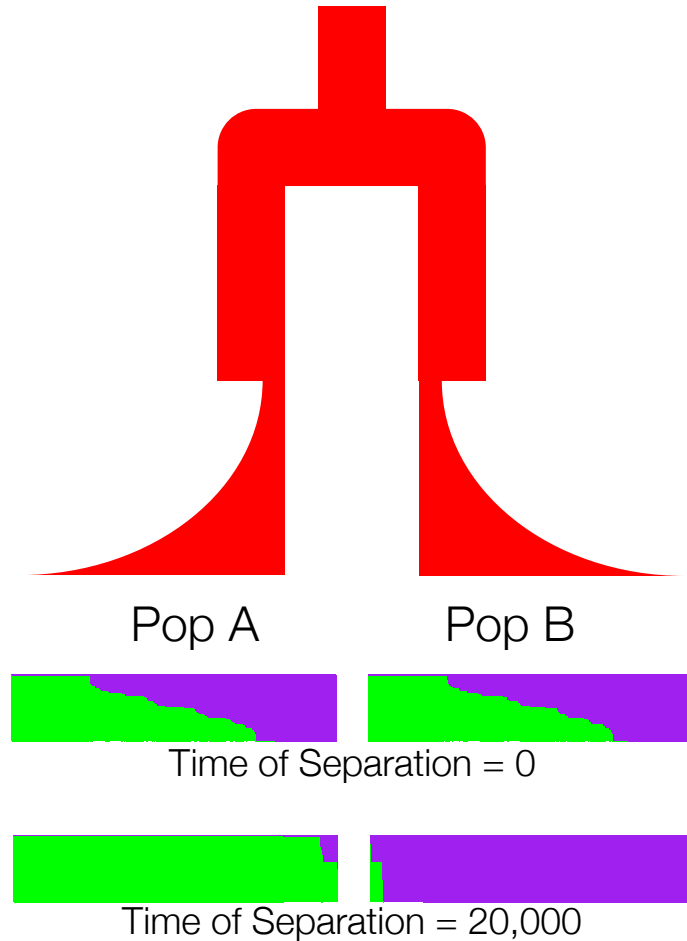
Han et al. (2017) Nature Comm.

# 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)

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)

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

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

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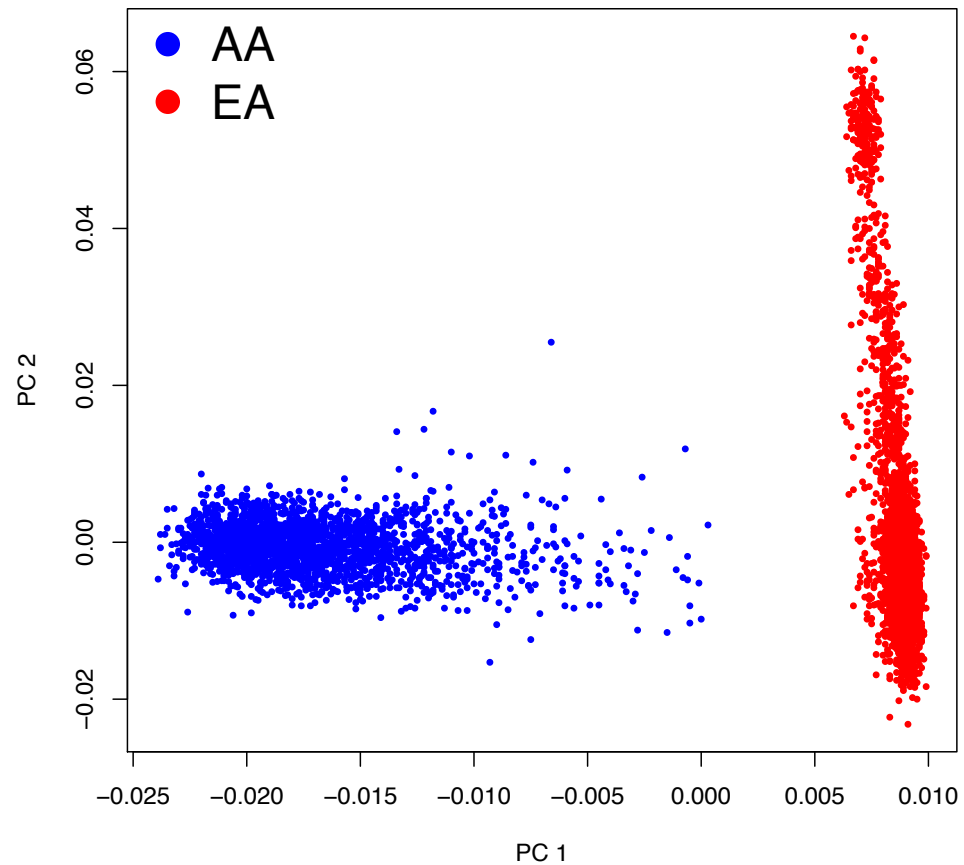
$$E(I_n \mid C,M) = \sum_{m\in M}\sum_{l=0}^{C} r_{lm} \times \sum_{j=1}^{N}\left(-p_{jlm} \ln p_{jlm} + \sum_{i=1}^{K} q_i p_{ijlm} \ln p_{ijlm}\right)$$
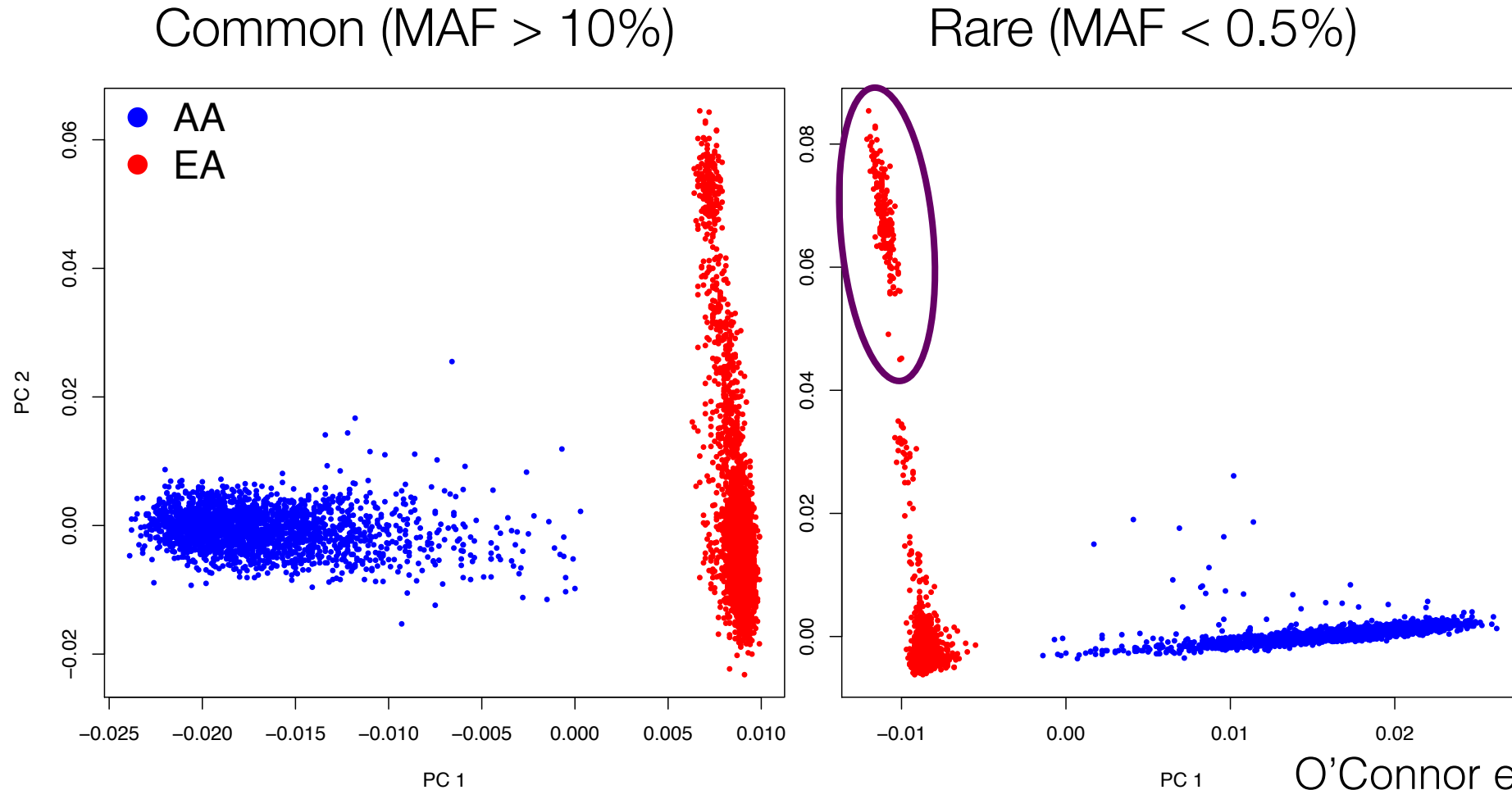


O'Connor et al. (2014)
Mol. Biol. Evol.

# Rare Variants Identify Cryptic Populations
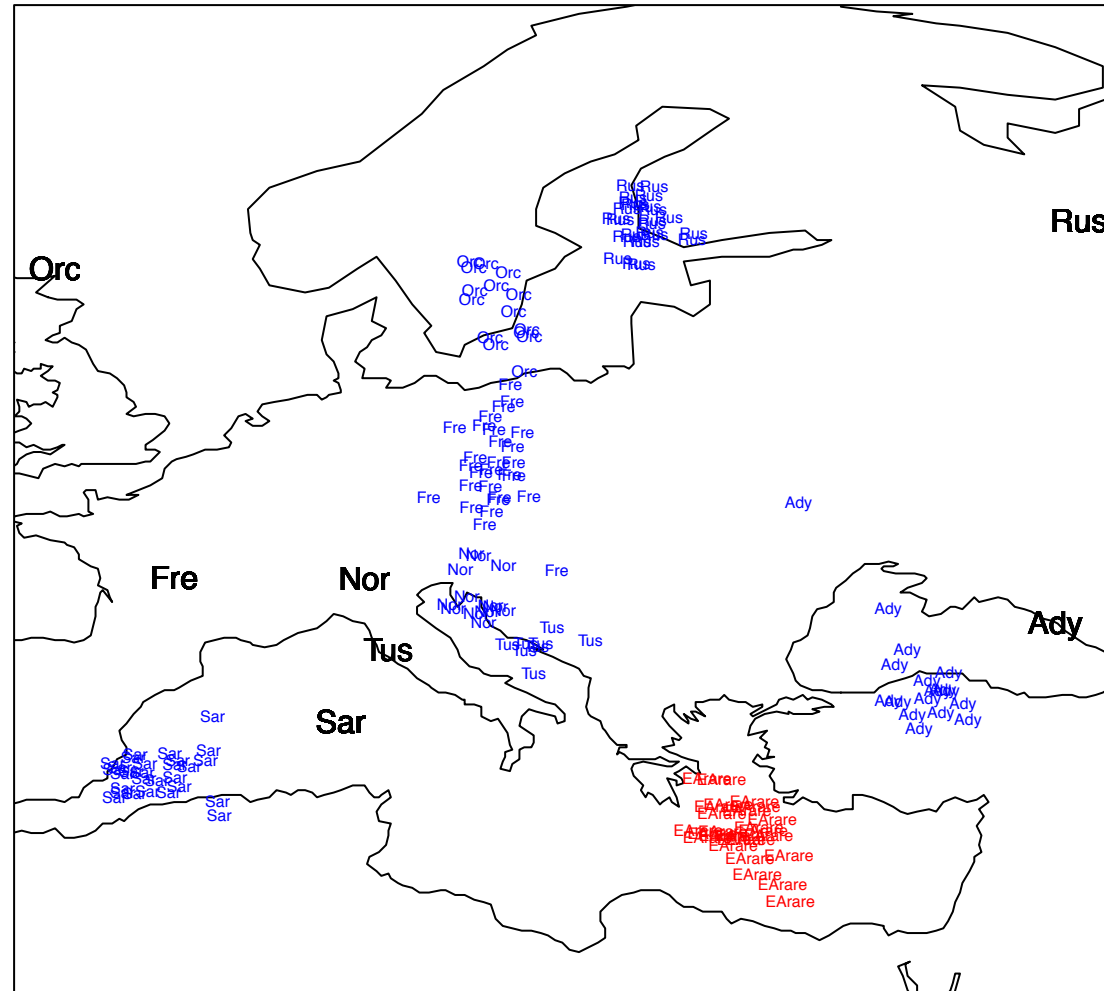


Common (MAF > 10%)

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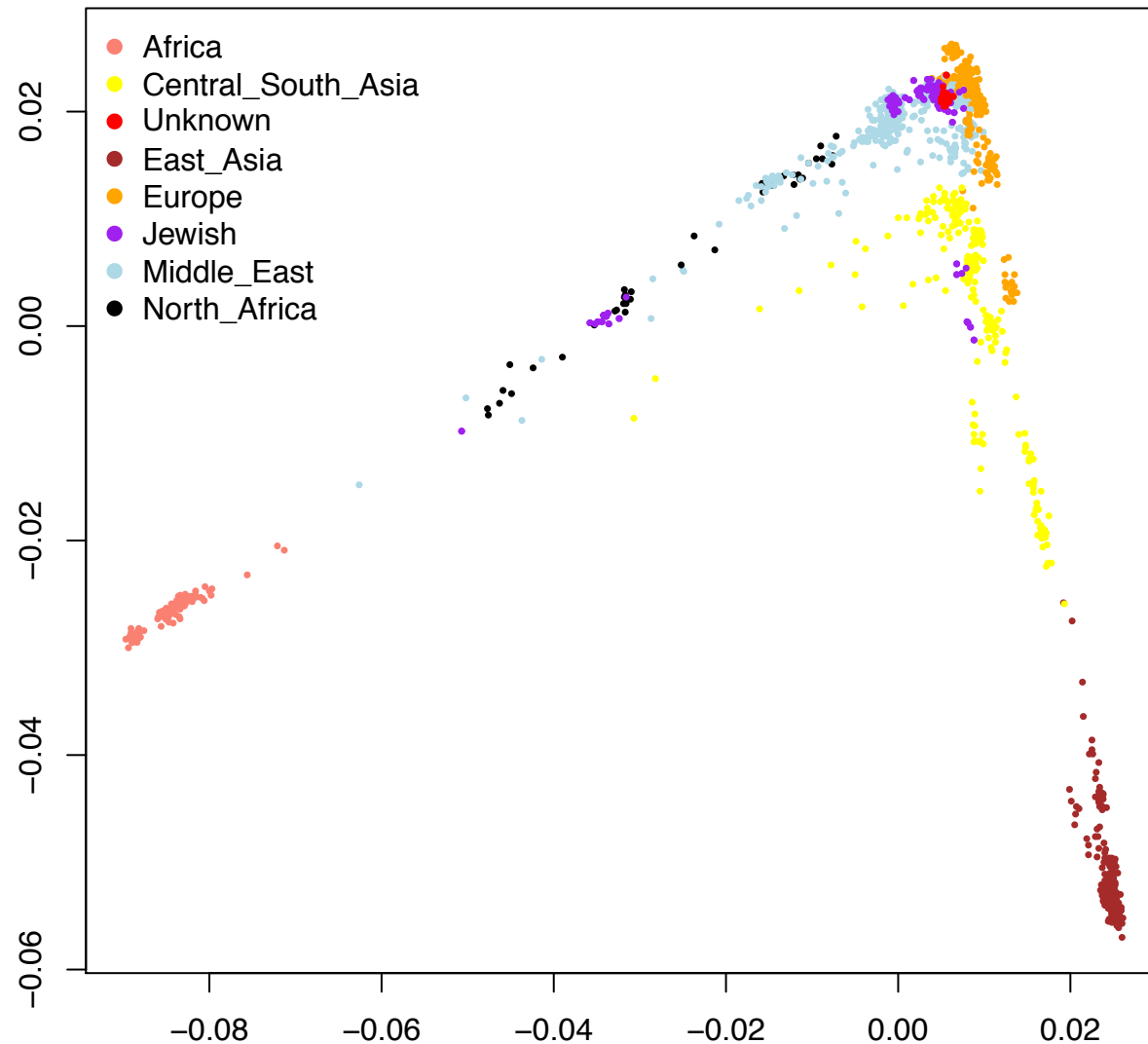Mol. Biol. Evol.

# Rare Variants Identify Cryptic Populations



Common (MAF > 10%)

Rare (MAF < 0.5%)

AA
EA

PC 2

PC 1

PC 1

O'Connor et al. (2014)
Mol. Biol. Evol.
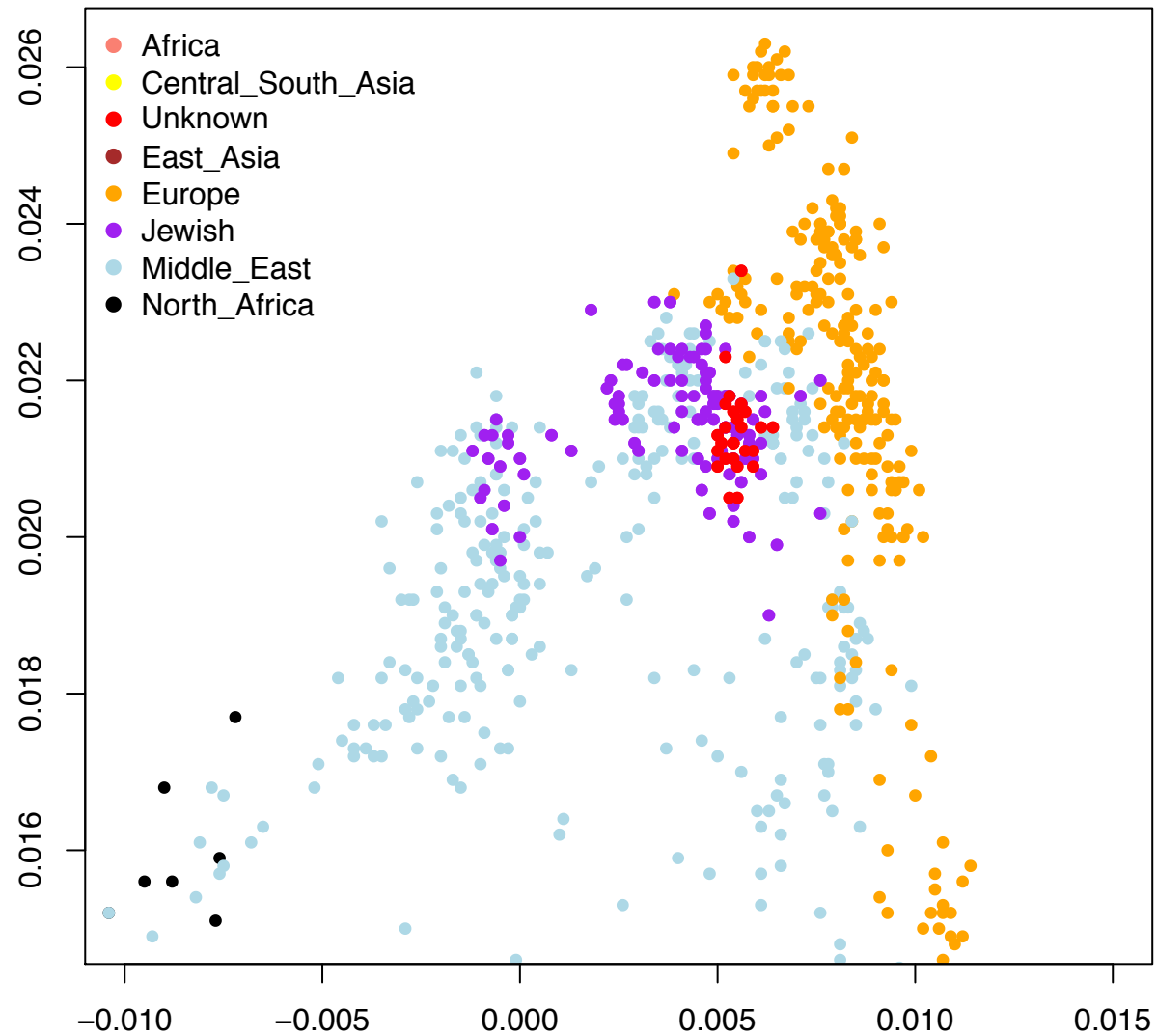
# What is Their Geographic Ancestry?



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

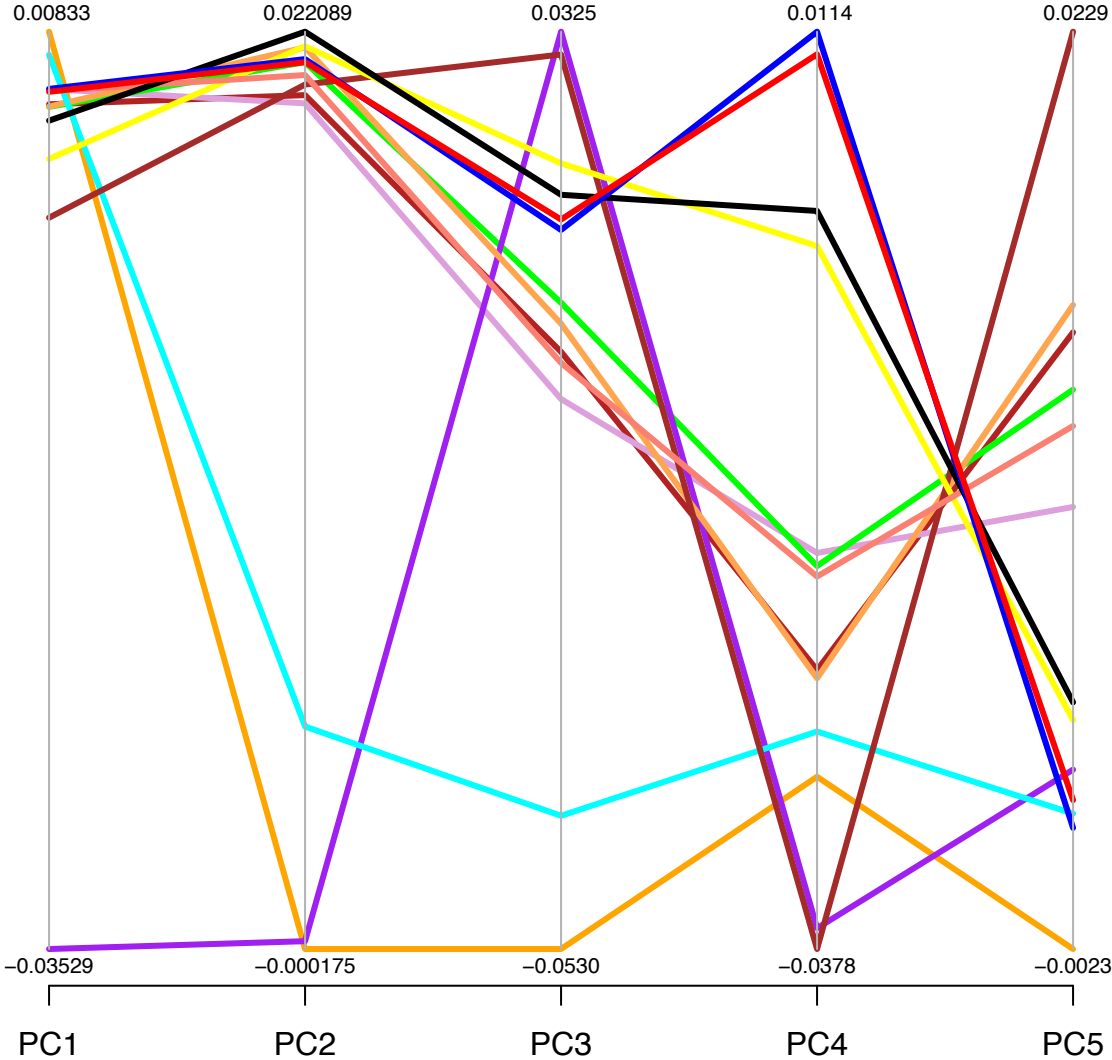# PCA of Global Diversity Including Cryptic Population



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

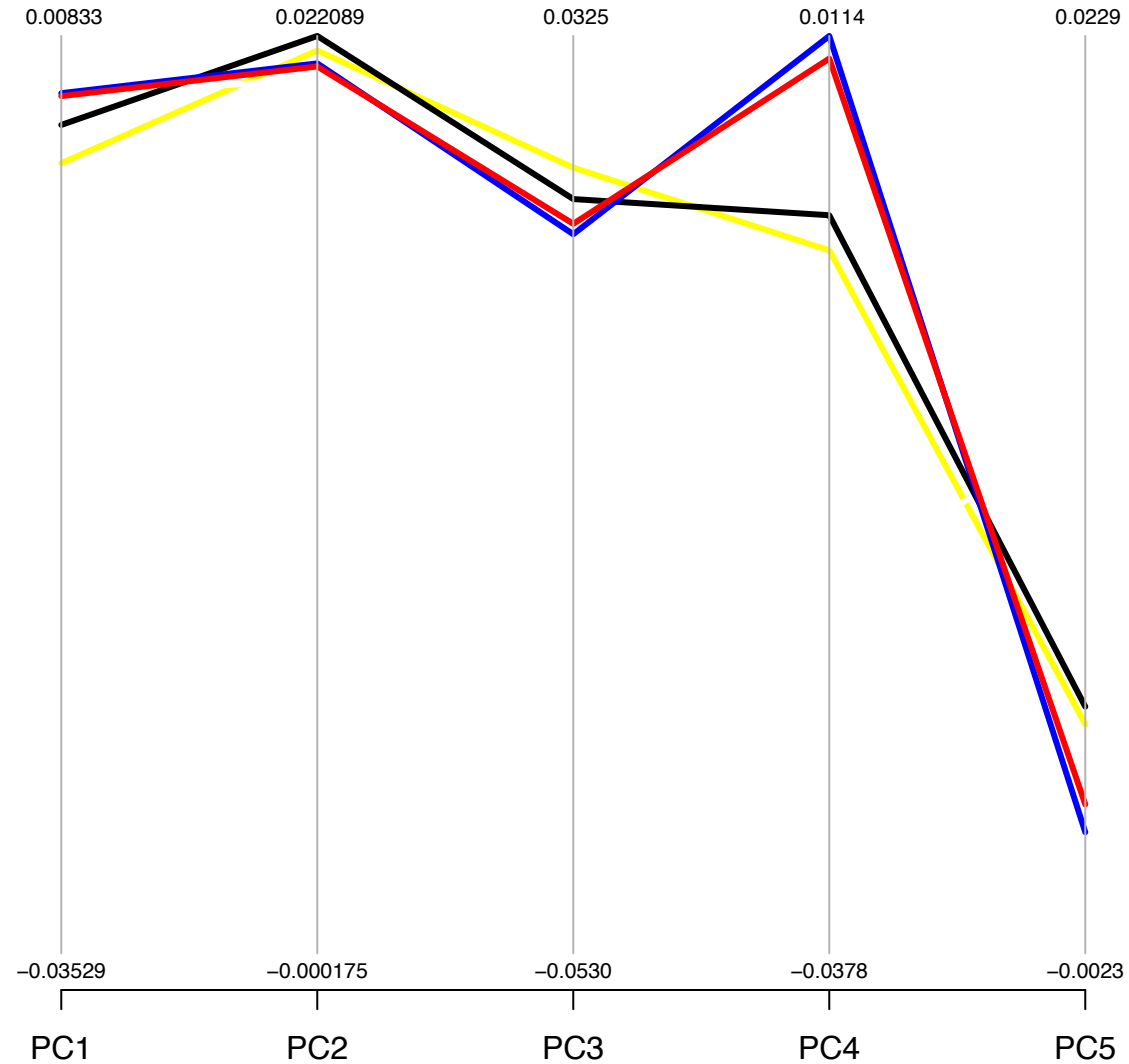# PCA of Global Diversity Including Cryptic Population



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

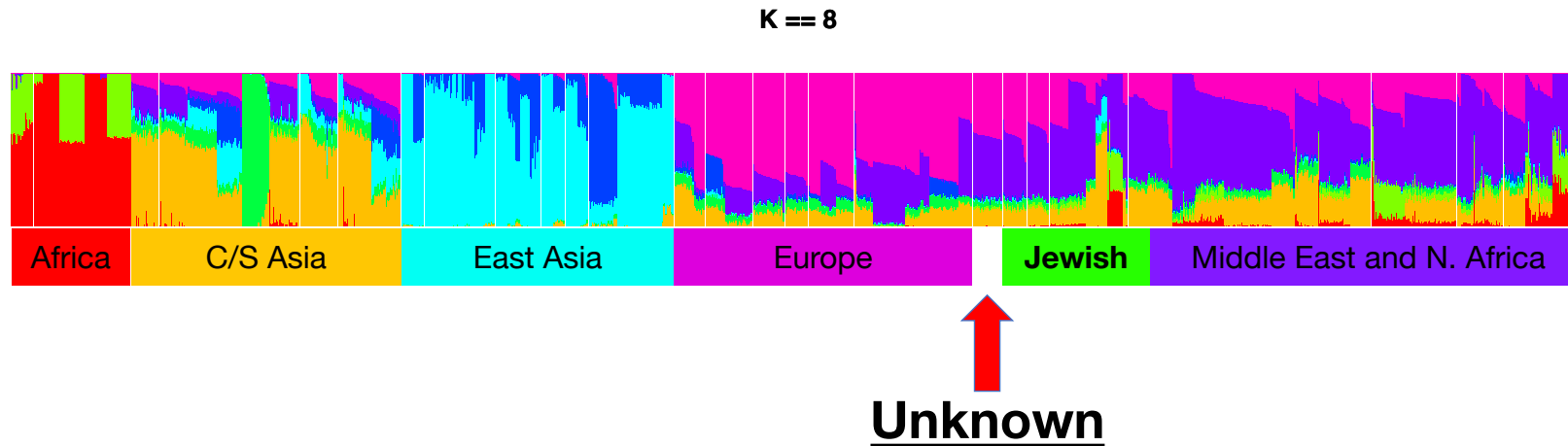# Population Average PCA with More Axes



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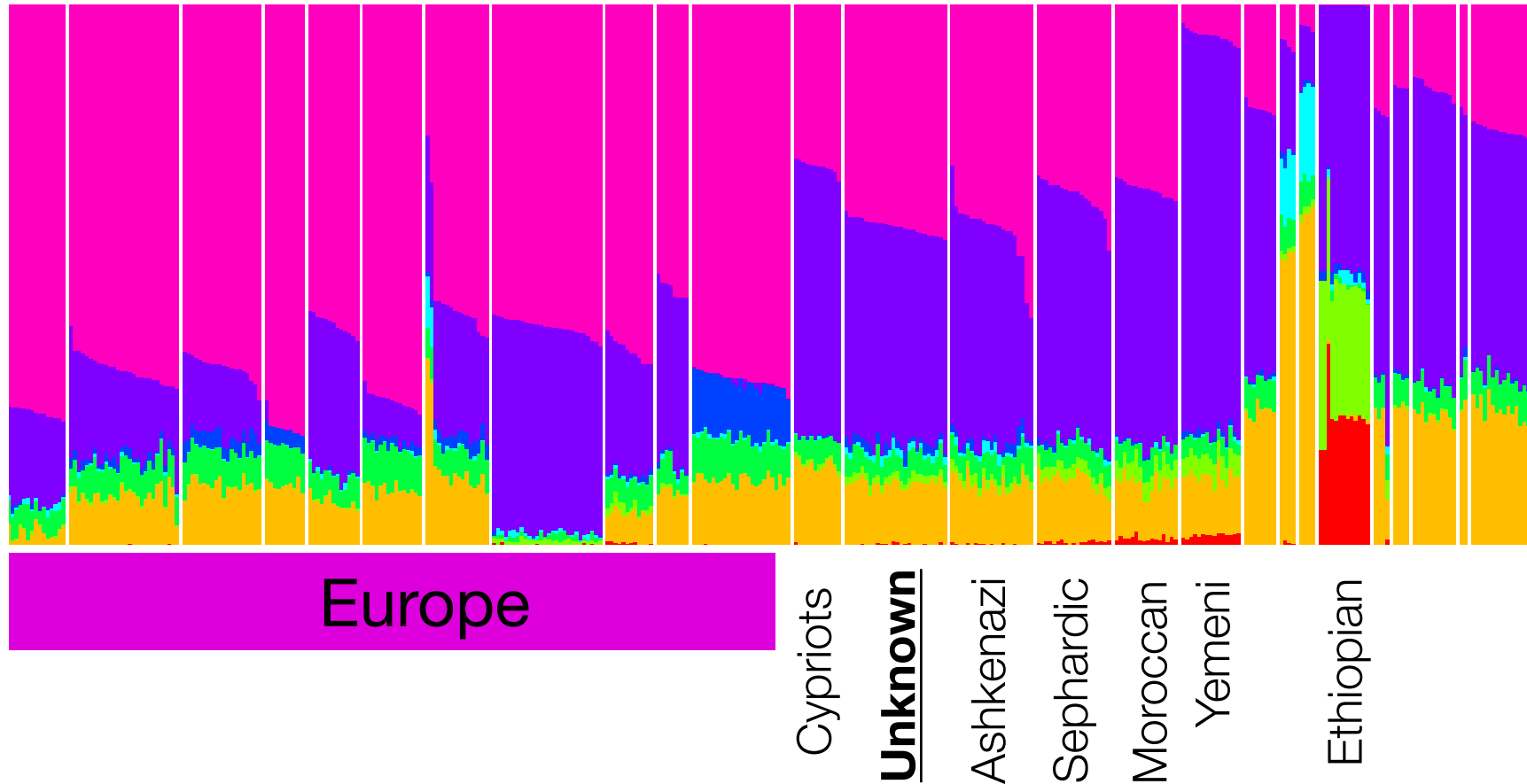Mol. Biol. Evol.

Population Average PCA with More Axes

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

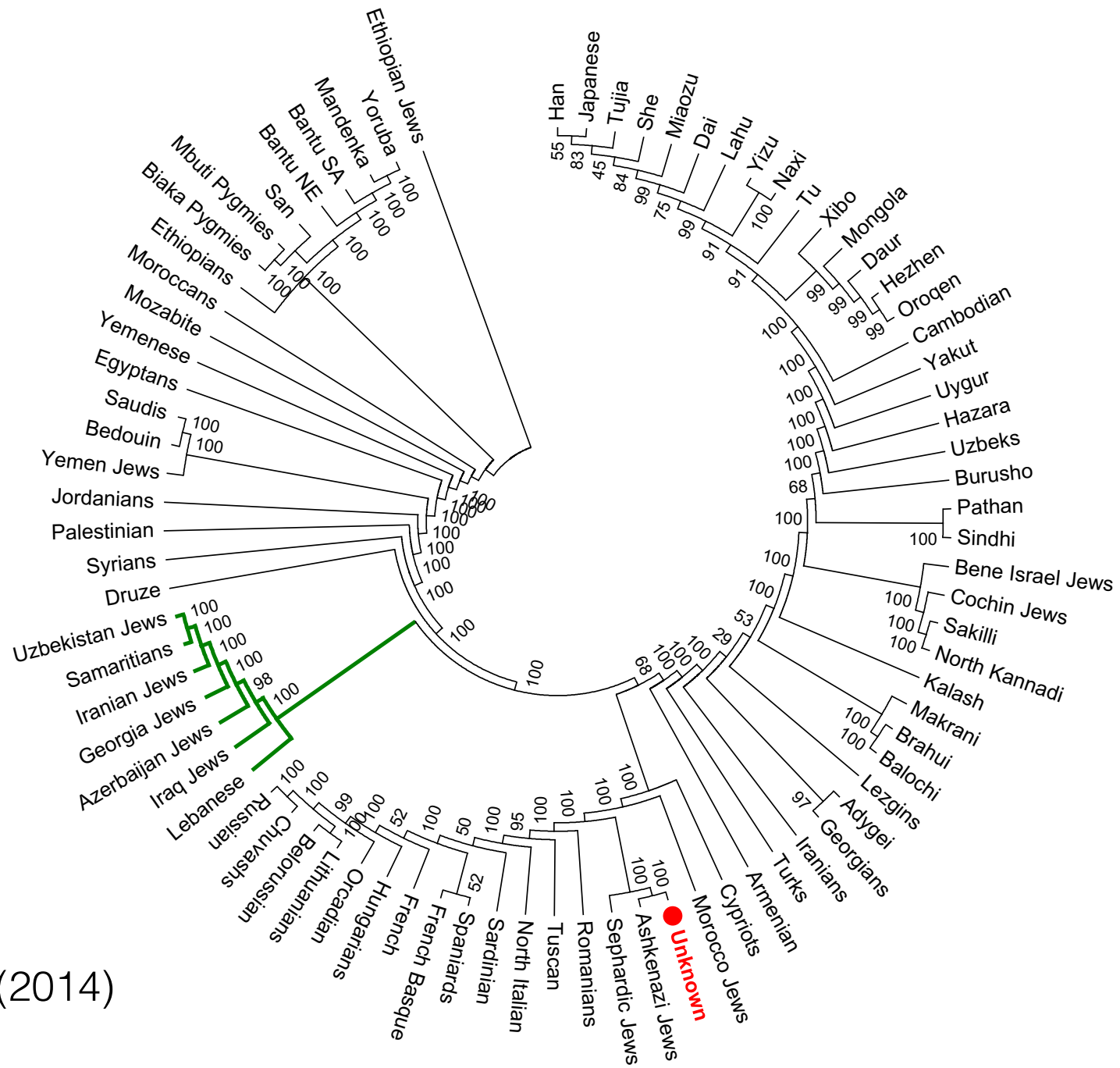# Cryptic Group has Similar Admixture Proportions to Jewish Groups.

**K == 8**



Africa | C/S Asia | East Asia | Europe | **Jewish** | Middle East and N. Africa

**Unknown**

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

# Cryptic Group has Similar Admixture Proportions to Jewish Groups.

Europe

Cypriots

**Unknown**

Ashkenazi

Sephardic

Moroccan

Yemeni

Ethiopian

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

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

# Population Stratification - Concept



Hirschhorn (2002) Genetics in Medicine

# Population Stratification – Example of spurious association

## Population 1

% with disease: 10%

% with variant allele (A*): 20%

|     | Case | Cont. |     |
|-----|------|-------|-----|
| A*  | 2    | 18    | 20  |
| G   | 8    | 72    | 80  |
|     | 10   | 90    | 100 |

OR = (2 *72)/(8*18) = 1

## Population 2

% with disease: 40%

% with variant allele (A*): 50%

|     | Case | Cont. |     |
|-----|------|-------|-----|
| A*  | 20   | 30    | 50  |
| G   | 20   | 30    | 50  |
|     | 40   | 60    | 100 |

OR = (20 *30)/(20*30) = 1

|     | Case | Cont. |     |
|-----|------|-------|-----|
| A*  | 22   | 48    | 70  |
| G   | 28   | 102   | 130 |
|     | 50   | 150   | 200 |

OR = (22 *102)/(28*48)
   = 1.67

# Population Stratification – thought question

The problem – poor sample matching.
Cases and controls are not selected from
the same source populations.

**Population 1**

% with variant allele (A*):  20%

**Population 2**

% with variant allele (A*):  50%

Is a scenario like this an issue for
continuous traits vs. case-control analysis?
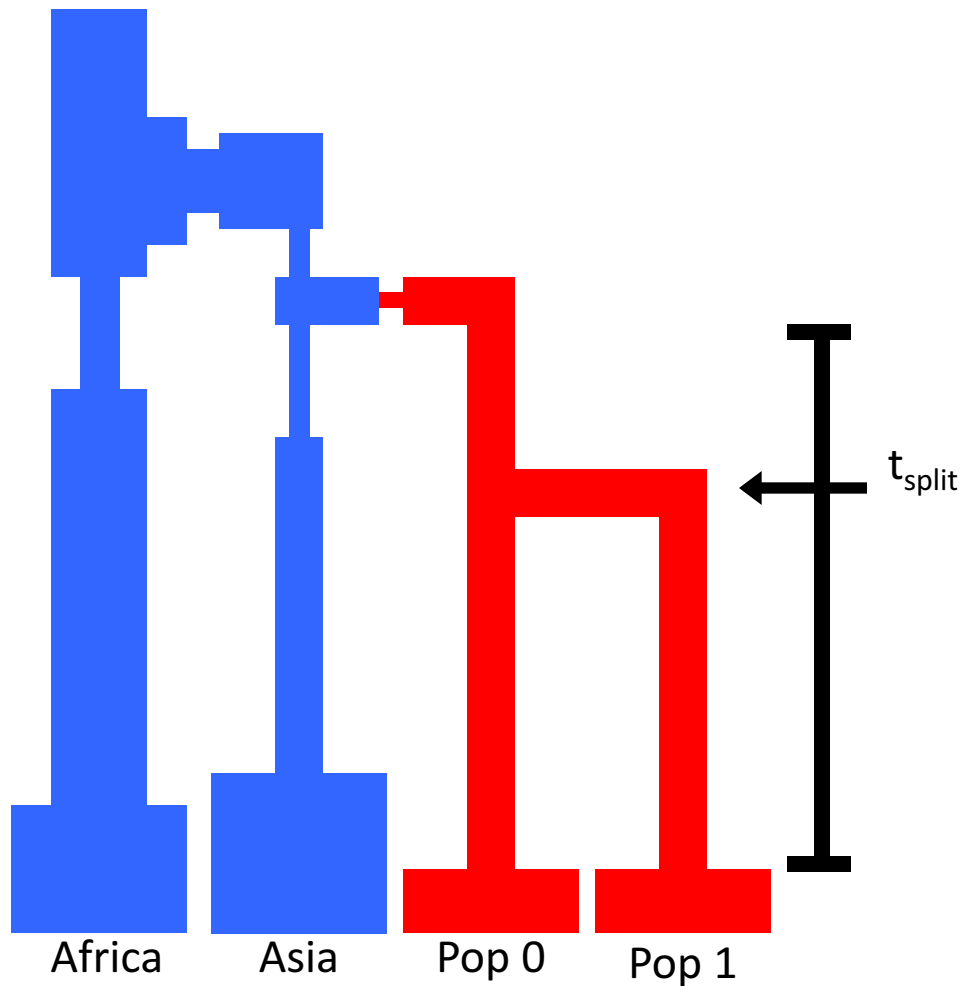
# Determining Proportions of Case and Control

$$P(i = X) = \text{The proportion of the subpopulation } i \text{ in the full population.}$$

$$P(d = c \mid i = X) = \text{The probability of subpopulation } i \text{ being a case (ie disease risk).}$$

$$P(d = c) = \sum_{i=1}^{N} P(d = c \mid i = X) \times P(i = X)$$

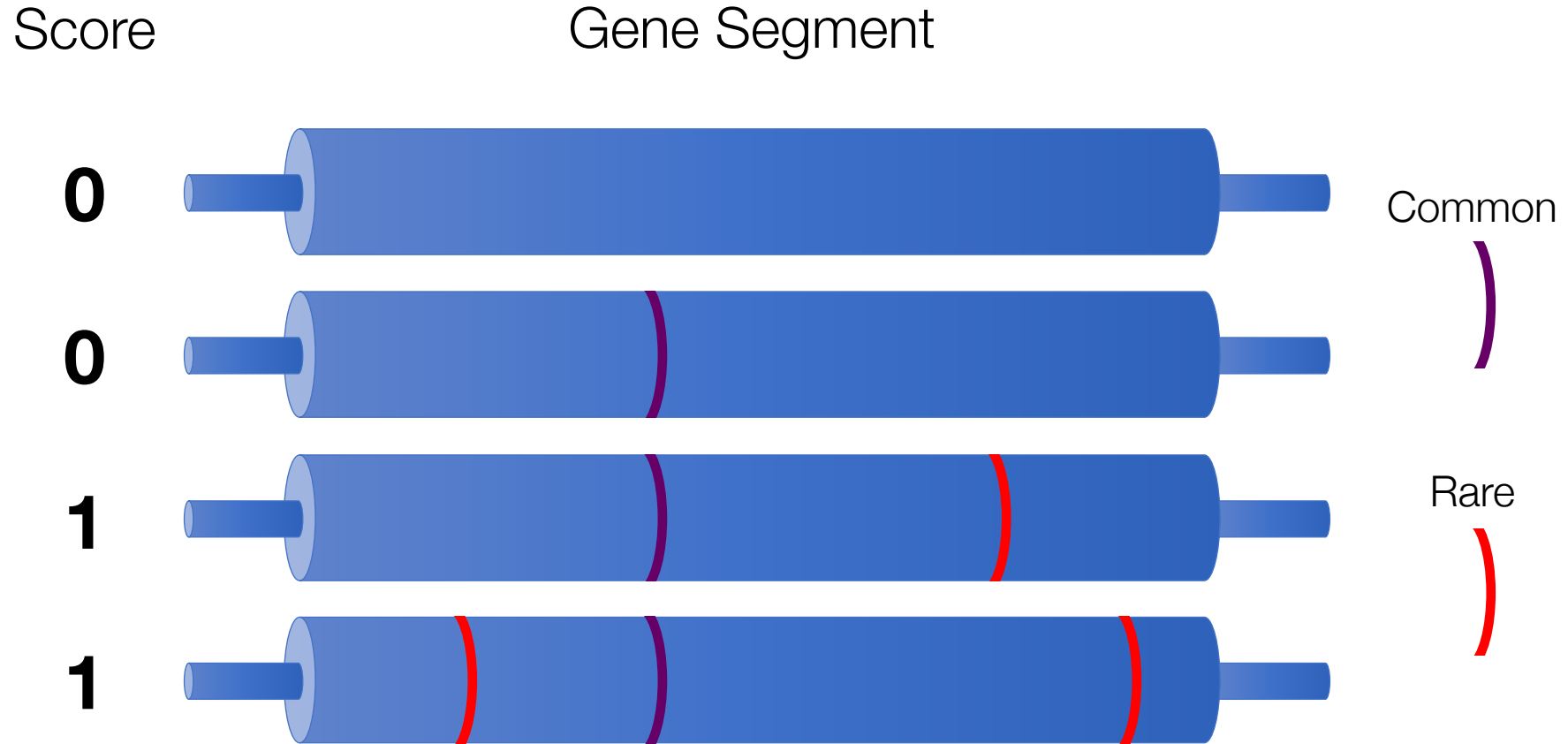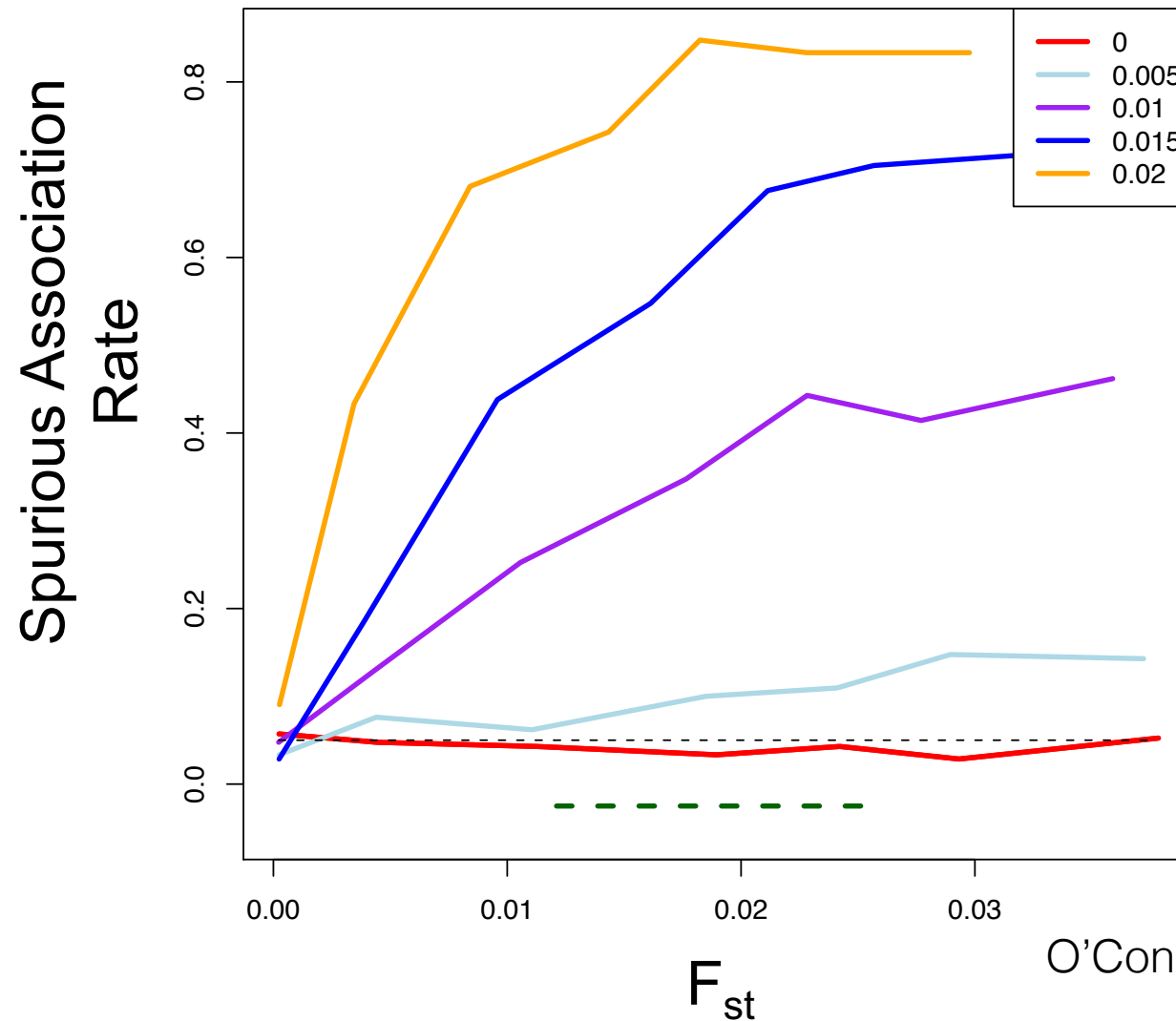$$P(i = X \mid d = c) = \frac{P(d = c \mid i = X) \times P(i = X)}{P(d = c)}$$

O'Connor et al. (2013) PLoS One

# Two Population Case



$$P(i = X) = 0.5$$

$$P(d = c \mid i = 0) = 0.02 + x$$
$$P(d = c \mid i = 1) = 0.02 - x$$

O'Connor et al. (2013) PLoS One

# Spurious associations as a function of confounding
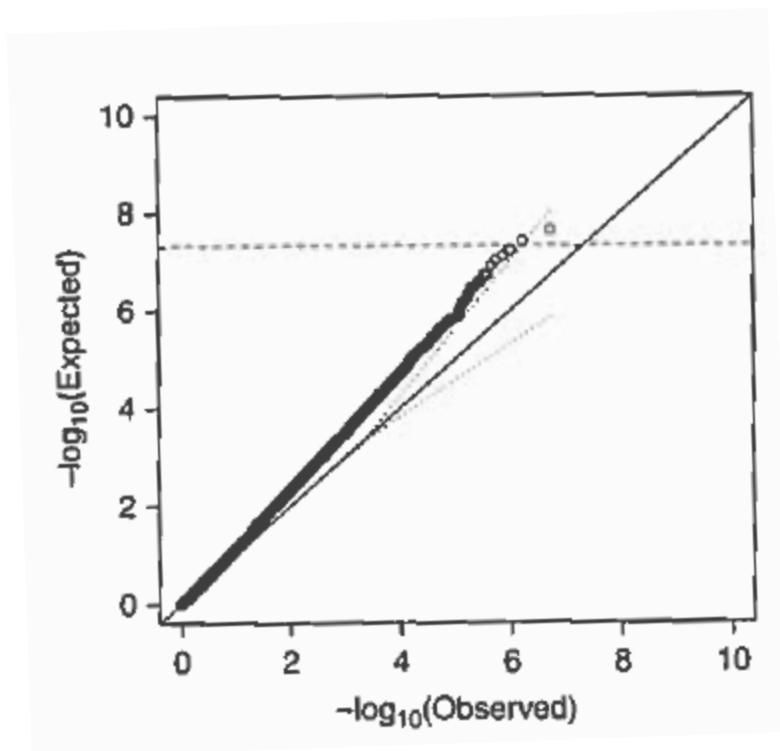


O'Connor et al. (2013) PLoS One

# General Approaches to handle structure in association analyses

- Stratify by race – carefully match cases and controls.
- Control with family data
- Genomic control
- Control using genetic markers –
  - Ancestry informative markers (AIMs)
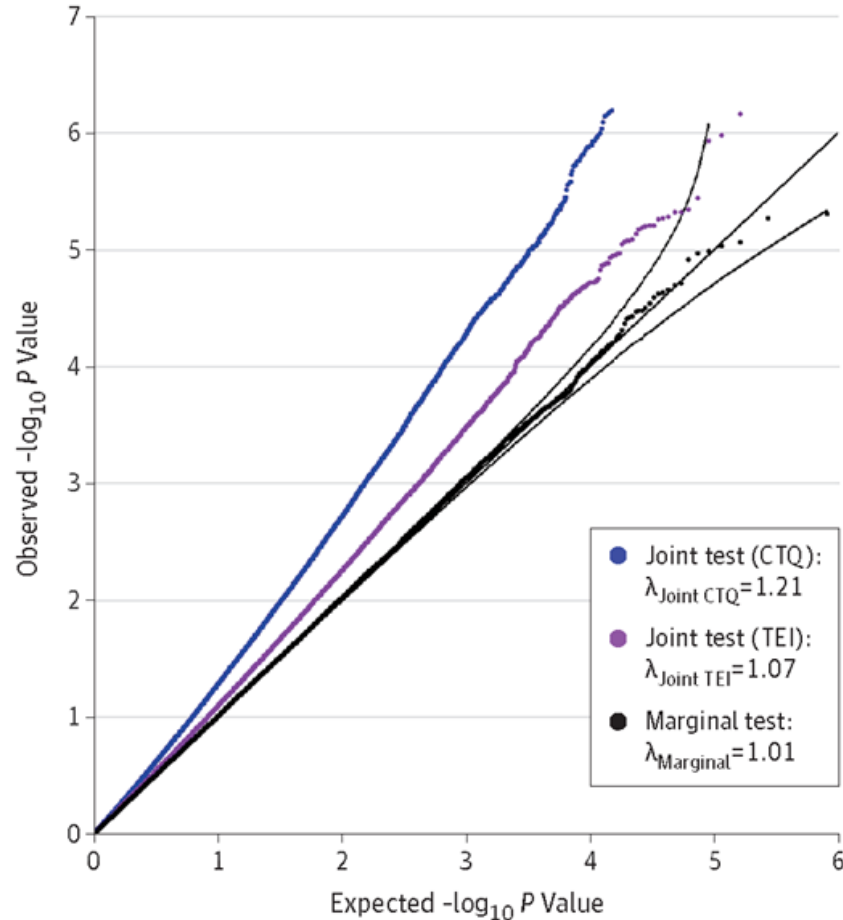  - Principle components analysis

# Genomic Control

• <u>How inflated is this?</u>



<u>Genomic Control</u>
• Assumption 1 - inflation affects the entire distribution of test statistics
• Assumption 2 – Underlying distribution of test statistics is $\chi^2$ distribution
• Take median of observed test statistics
• Divide by expected =0.4549

$$\lambda = \chi^2_{observed\ median} / 0.4549$$

• Divide all test statistics by $\lambda$

Austin. Genetic Epidemiology Methods & Applications 2013.

# Genomic Control Cont.
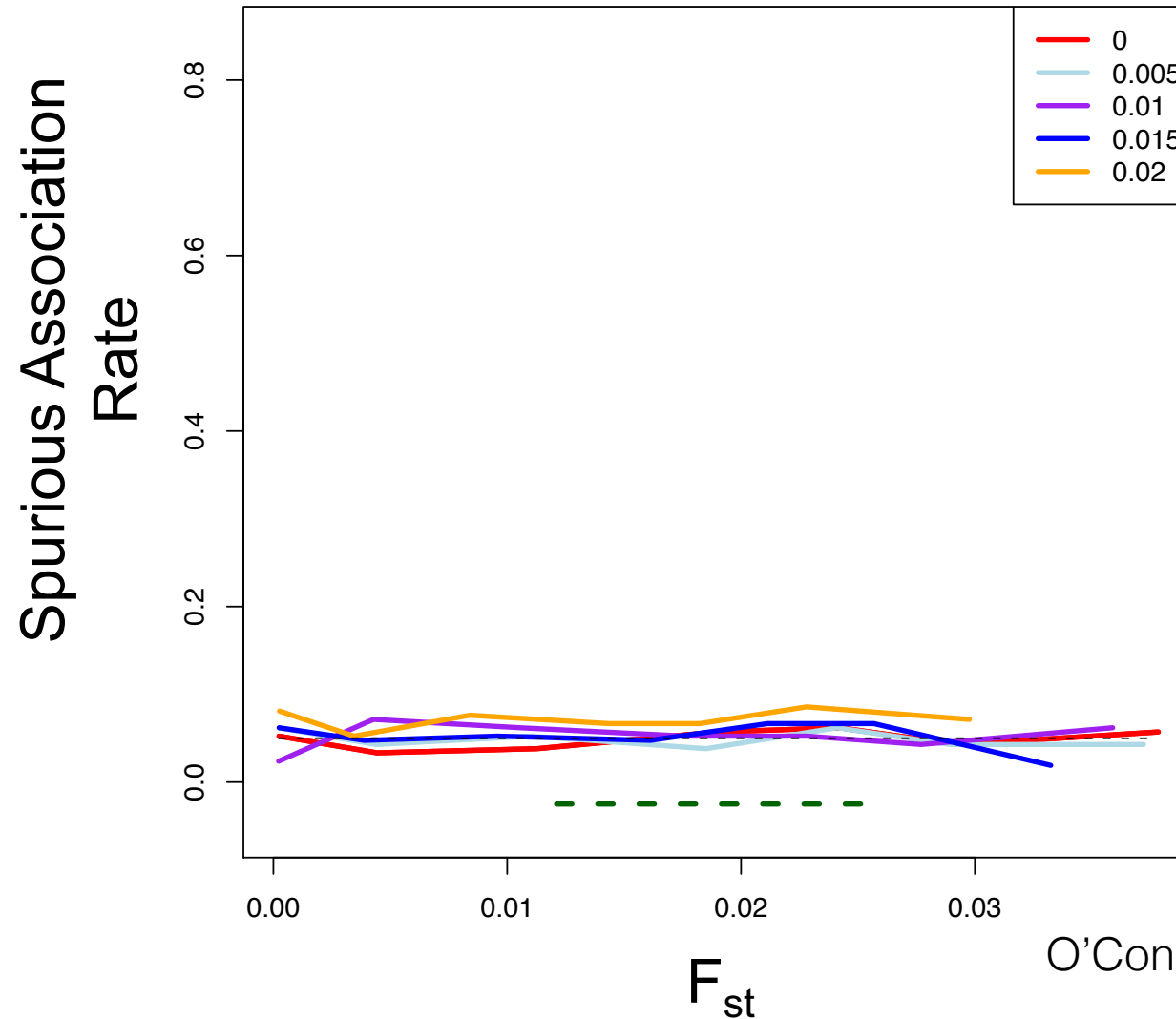


- Larger lambda → more systematic bias
- High lambda → consequence not cause
- Not only admixture (eg. cryptic relatedness)

# AIMs and PCs

- General idea – Use genetic markers to measure and account for ancestry differences
- AIMs – fixed markers known to detect population substructure
- PCs – use large scale data (e.g. GWAS) to conduct principle components analysis
  - Dimension reduction
  - The first principal components summarize most of the variation
  - PCs can be used as variable in regression models

# Spurious associations as a function of confounding with PC correction



O'Connor et al. (2013) PLoS One

# Concluding summary

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

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

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

- Cryptic relatedness can increase spurious associations for phenotype studies, but should be handled within the models.