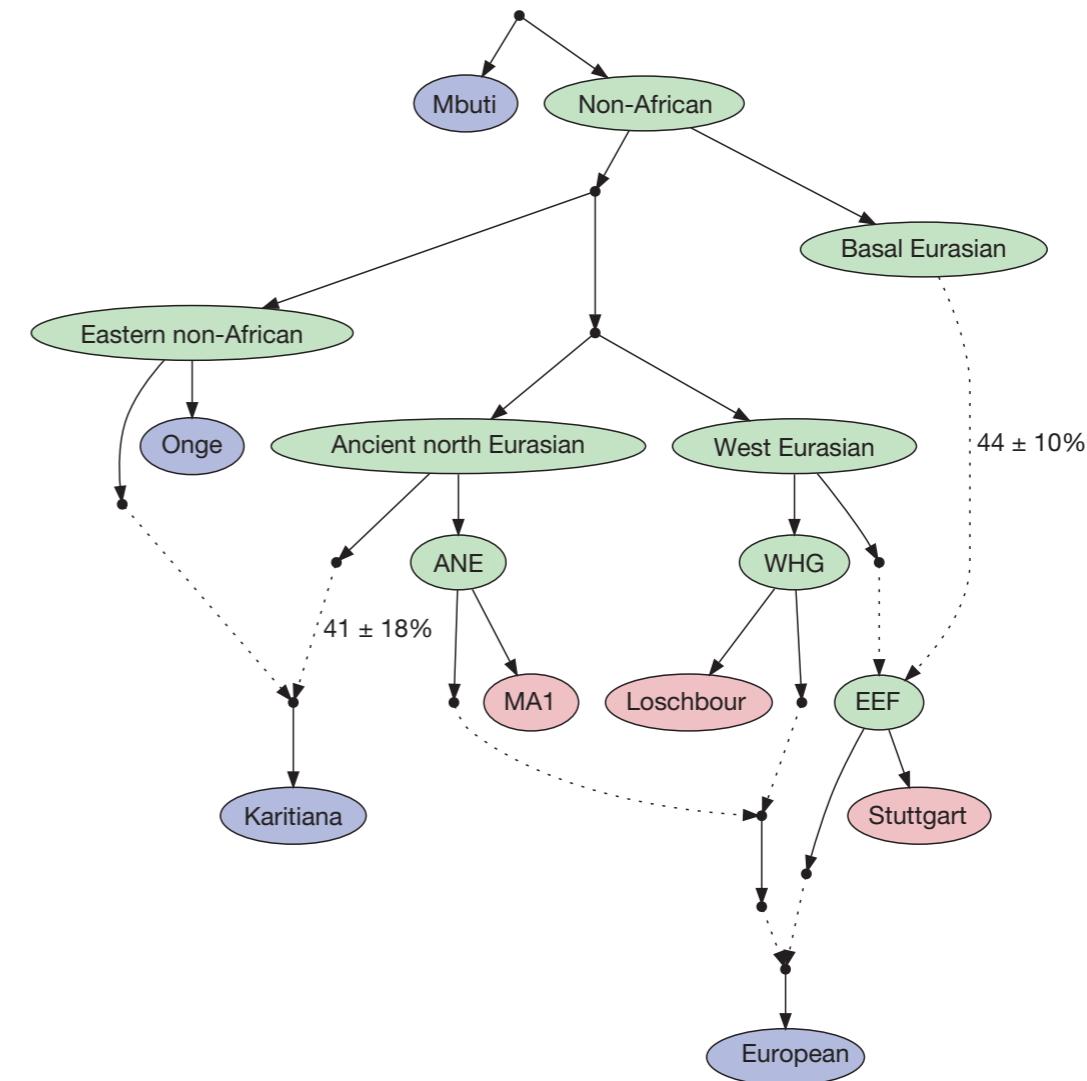
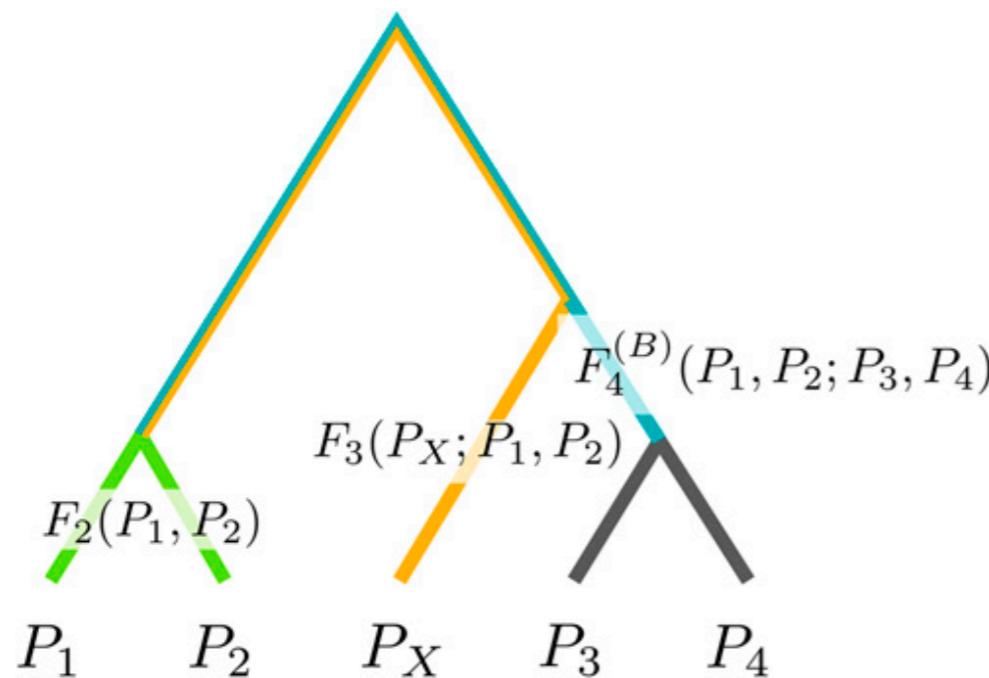


Inferring admixture from genomic data

EMBO Workshop Population Genomics, Napoli 2017



Martin Sikora, PhD

Centre for Geogenetics, Natural History Museum of Denmark, University of Copenhagen



Outline of today's lecture

- Motivation and Background
 - Admixture in human history
 - Overview of methods

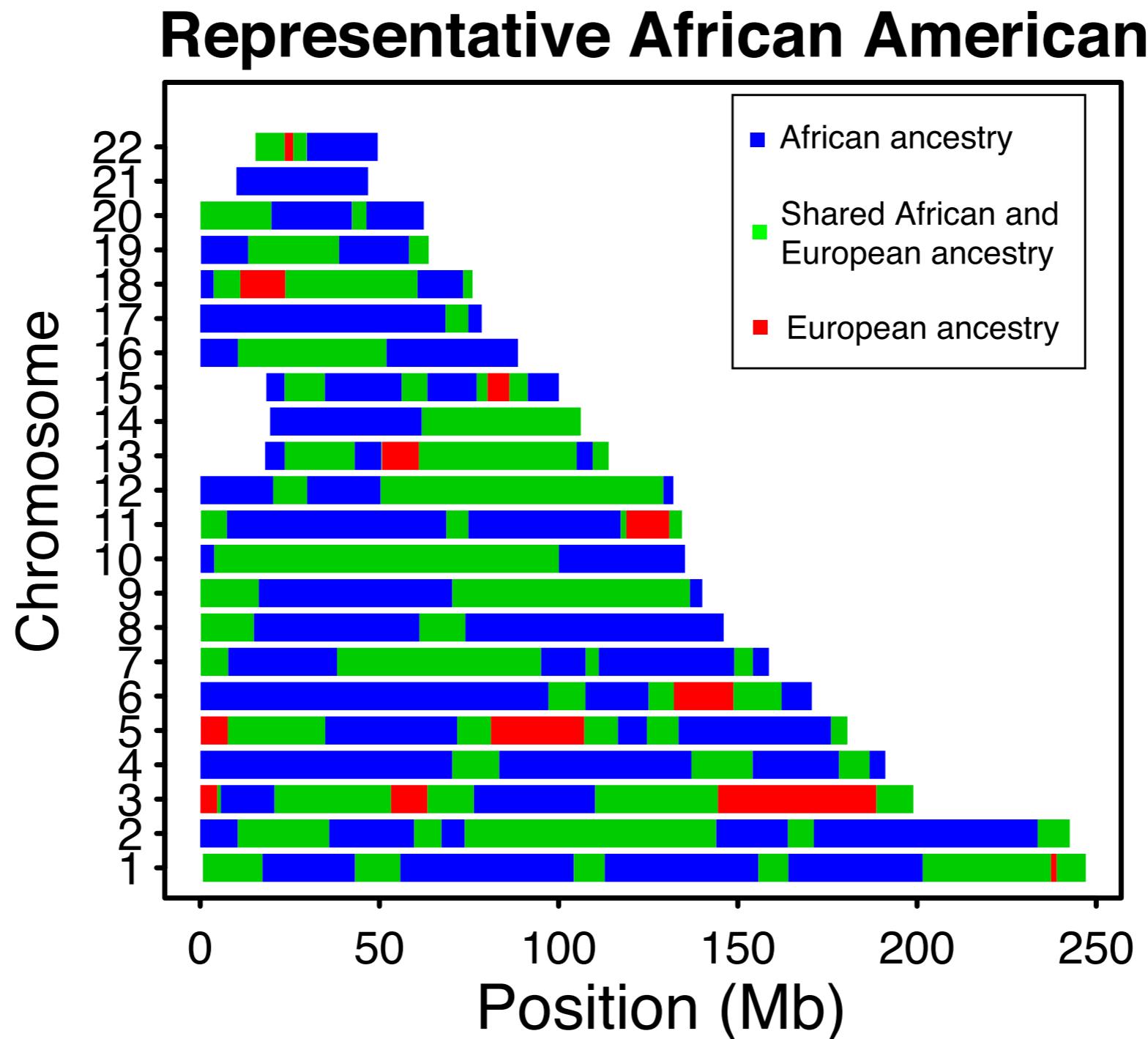
Outline of today's lecture

- Motivation and Background
 - Admixture in human history
 - Overview of methods
- Testing for admixture using f-statistics
 - Basic concepts
 - What can we test with f-statistics?
 - Estimation and interpretation
 - Admixture graph fitting

Outline of today's lecture

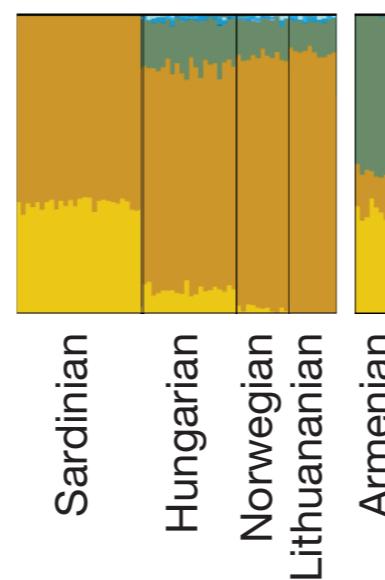
- Motivation and Background
 - Admixture in human history
 - Overview of methods
- Testing for admixture using f-statistics
 - Basic concepts
 - What can we test with f-statistics?
 - Estimation and interpretation
 - Admixture graph fitting
- Admixture dating
 - Basic concepts
 - ALDER

Pervasive admixture in human history

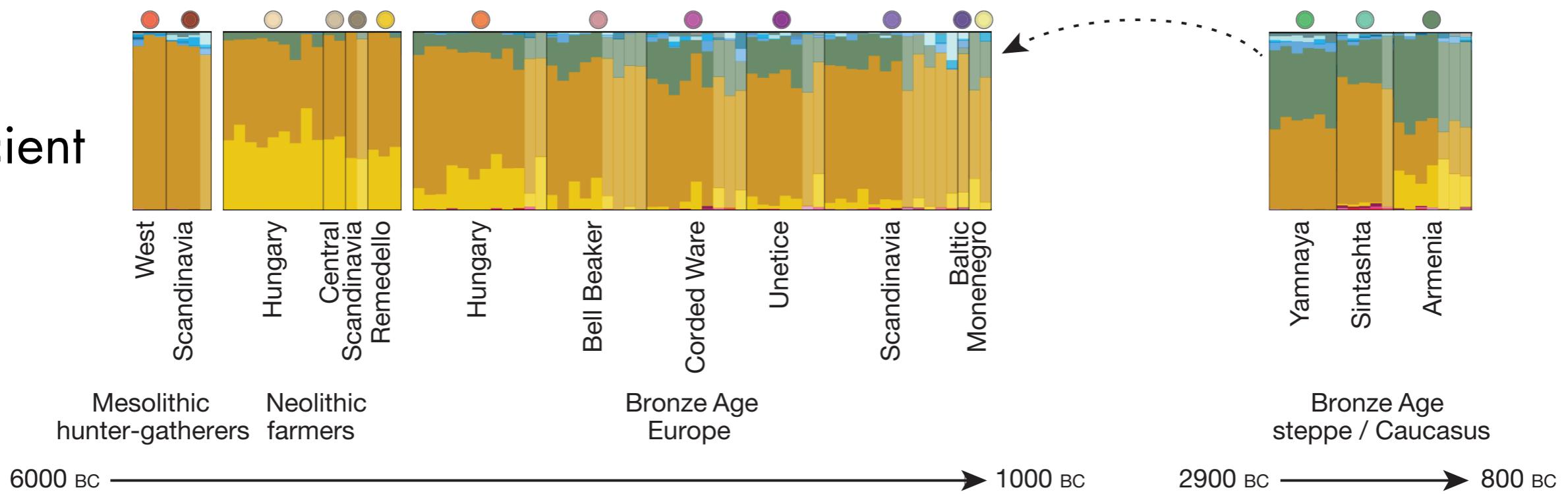


Pervasive admixture in human history

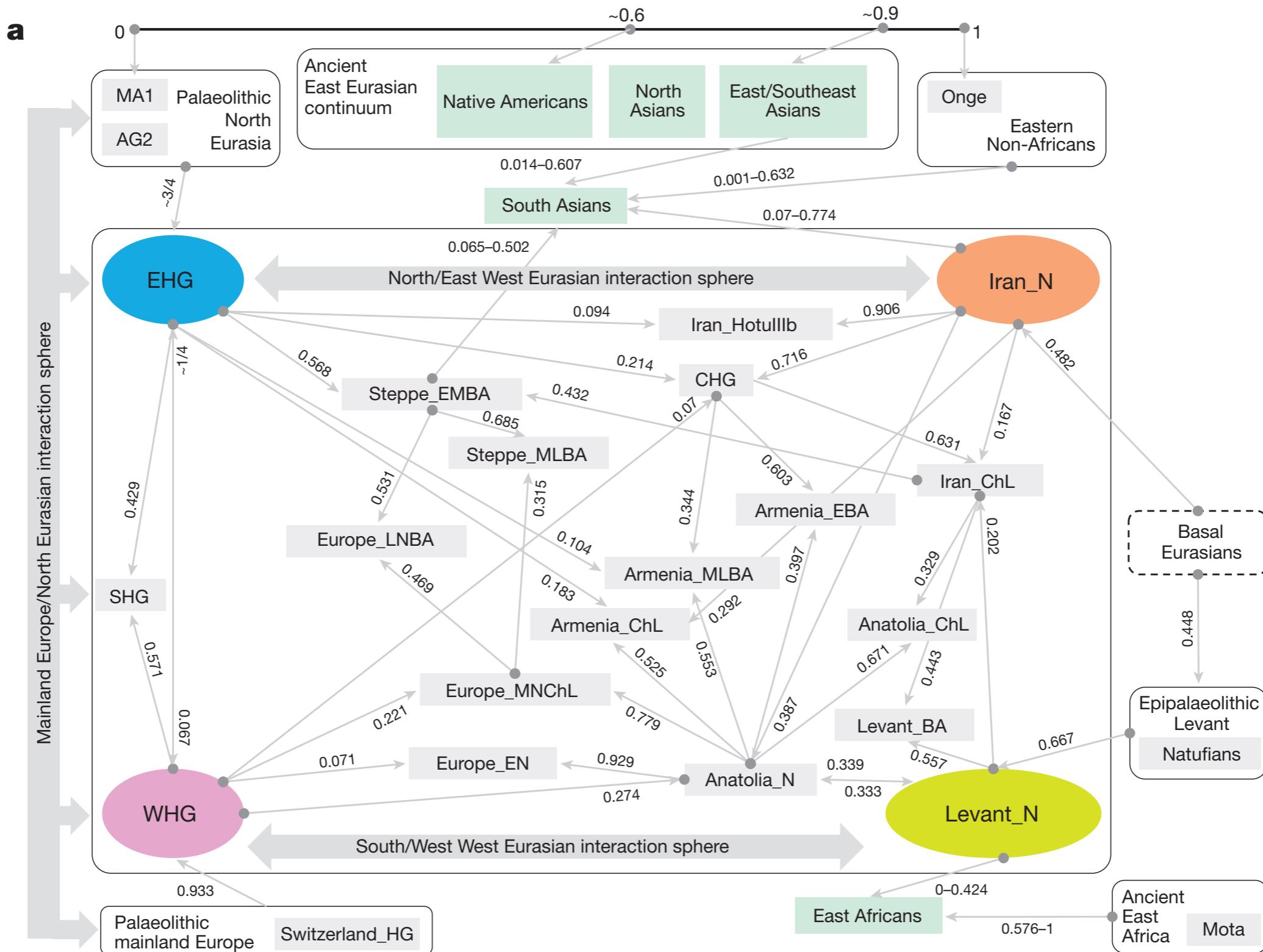
Modern



Ancient

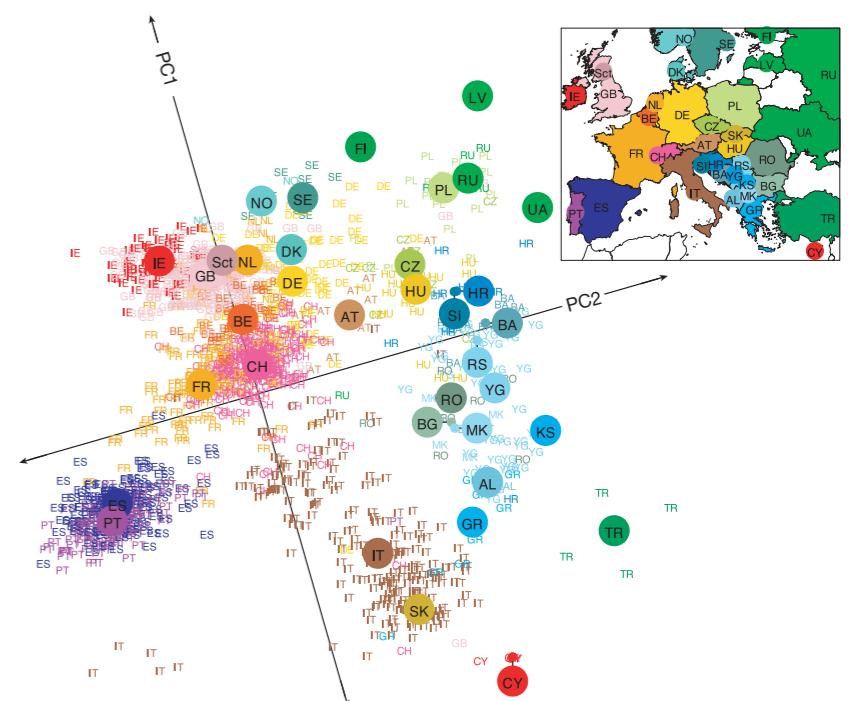
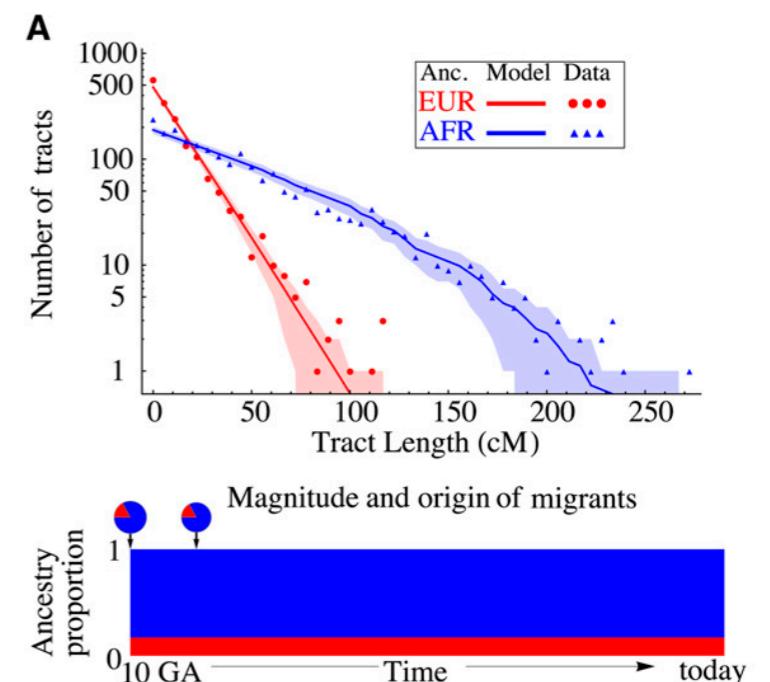


Pervasive admixture in human history



Methods to detect admixture

- Local methods
 - Local ancestry tracts (HAPMIX, TRACTS...)
 - Allow for complex admixture models
 - Need phased data
 - Less powerful for older admixture events
- Global methods
 - PCA, Model-based clustering (ADMIXTURE, STRUCTURE...)
 - Efficient and powerful method to detect structure in the data
 - Better at handling missing and lower-quality data
 - Difficult to interpret, or easy to over-interpret



Outline of today's lecture

- Motivation and Background
 - Admixture in human history
 - Overview of methods
- Testing for admixture using f-statistics
 - Basic concepts
 - What can we test with f-statistics?
 - Estimation and interpretation
 - Admixture graph fitting
- Admixture dating
 - Basic concepts
 - ALDER

Ancient Admixture in Human History

Nick Patterson,^{*1} Priya Moorjani,[†] Yontao Luo,[‡] Swapan Mallick,[†] Nadin Rohland,[†] Yiping Zhan,[‡] Teri Genschoreck,[‡] Teresa Webster,[‡] and David Reich^{*†}

^{*}Broad Institute of Harvard and Massachusetts Institute of Technology, Cambridge, Massachusetts 02142, [†]Department of Genetics, Harvard Medical School, Boston, Massachusetts 02115, and [‡]Affymetrix, Inc., Santa Clara, California 95051

Admixture, Population Structure, and F-Statistics

Benjamin M. Peter¹

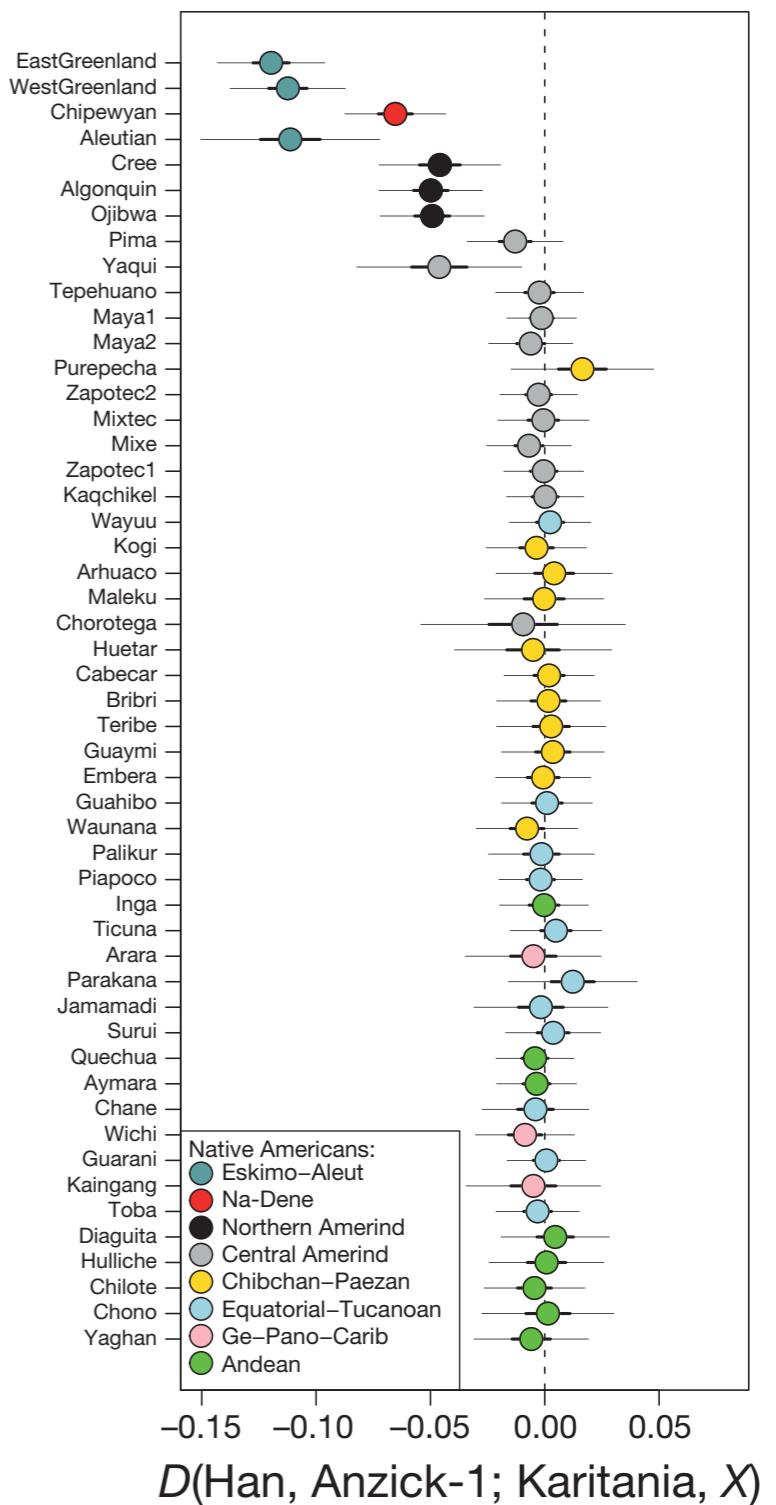
Department of Human Genetics, University of Chicago, Chicago, Illinois 60637

ORCID ID: 0000-0003-2526-8081 (B.M.P.)

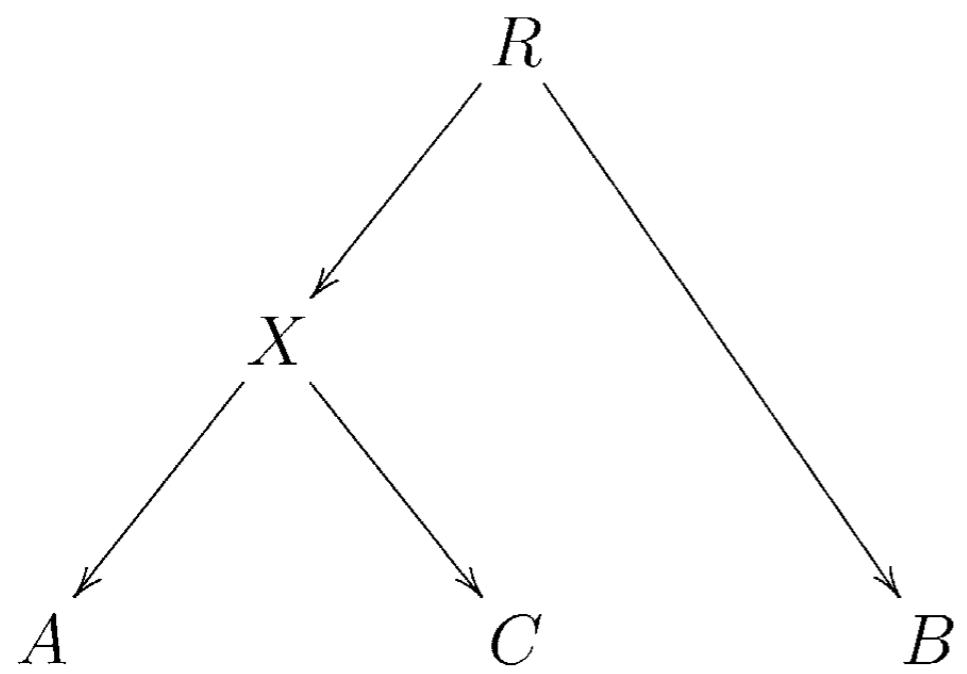
- Framework for admixture inference using allele frequency correlations / shared genetic drift among sets of populations
- Parts appeared earlier in various forms in the literature, but seminal paper summarising is Patterson et al. (2012) *Genetics*
- More recently additional theoretic work by Peter (2014) *Genetics*, interpreting f-statistics in the context of population genetics theory
- Has become a standard toolset to test hypotheses about population history and admixture

What can f-statistics be used for?

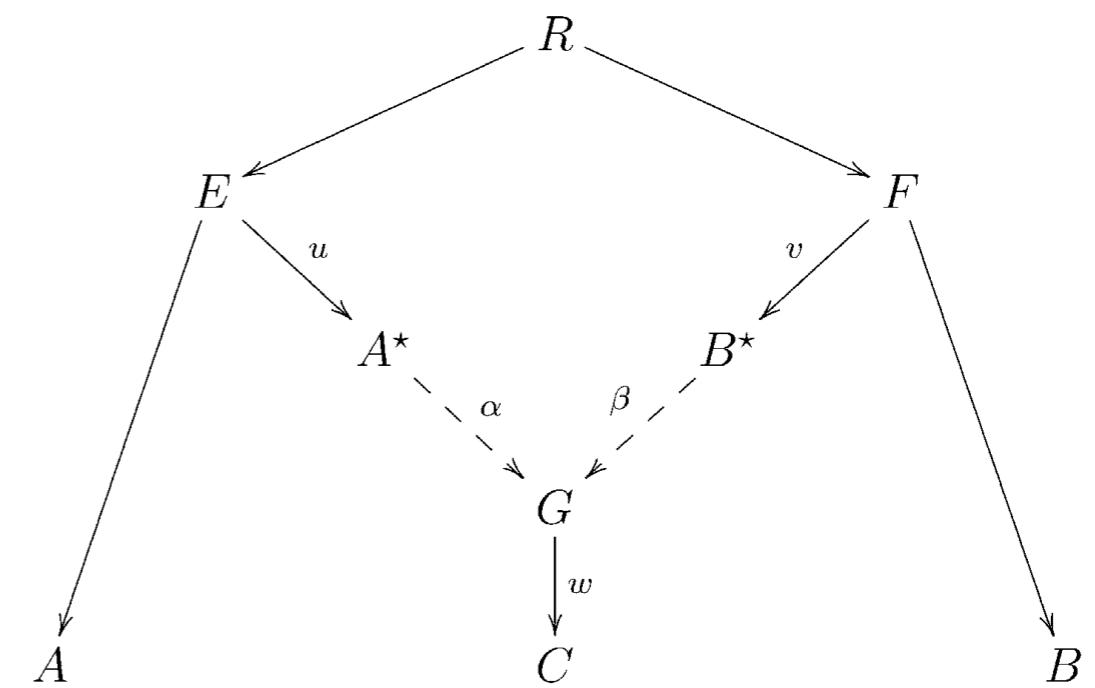
- Formal test for admixture (f_3)
- Test for treelessness (f_4)
- Estimation of admixture proportions (f_4 ratio)
- Admixture graph fitting
- Finding of closest related populations (outgroup f_3)
- Test for number of migration waves in a region (qpWave)
- Phylogeny-free admixture fitting (qpAdm)



Definitions

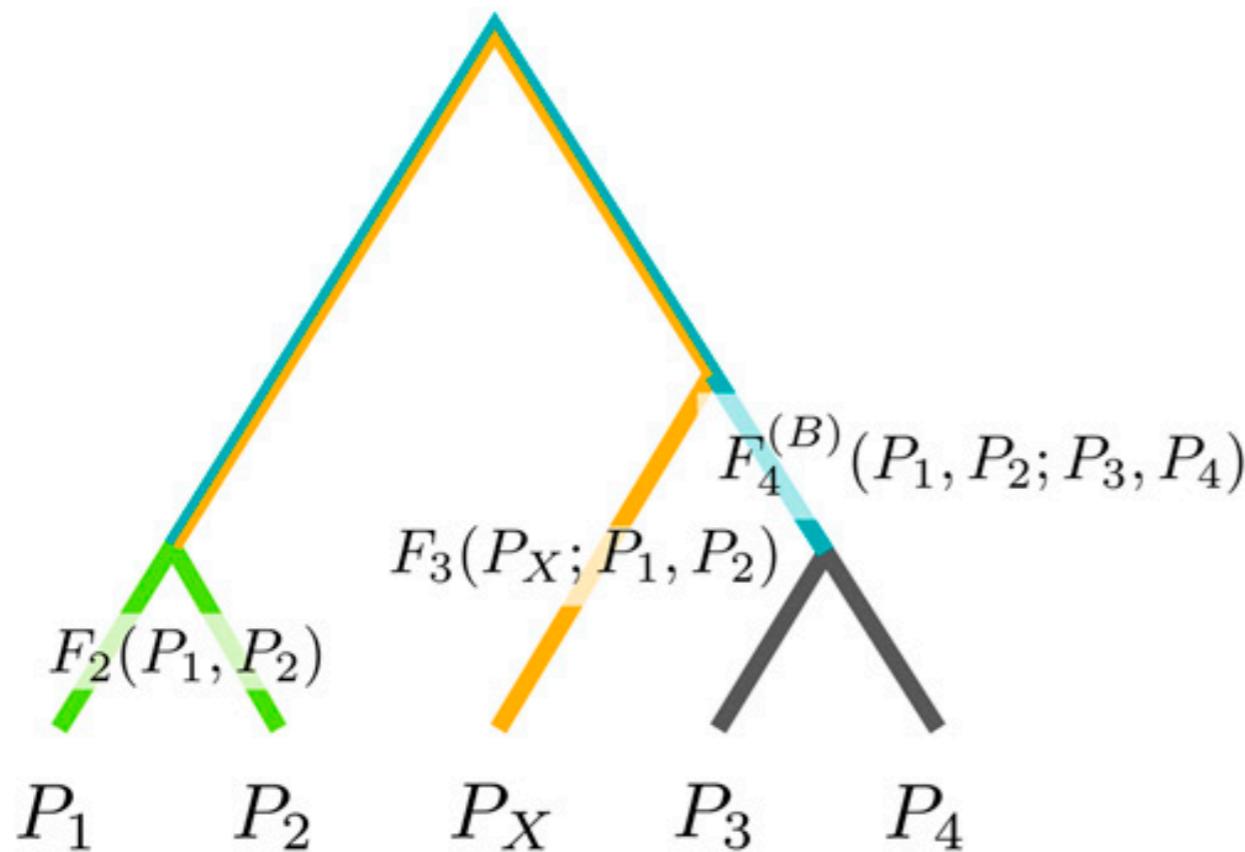


Population phylogeny



Admixture graph

Definitions



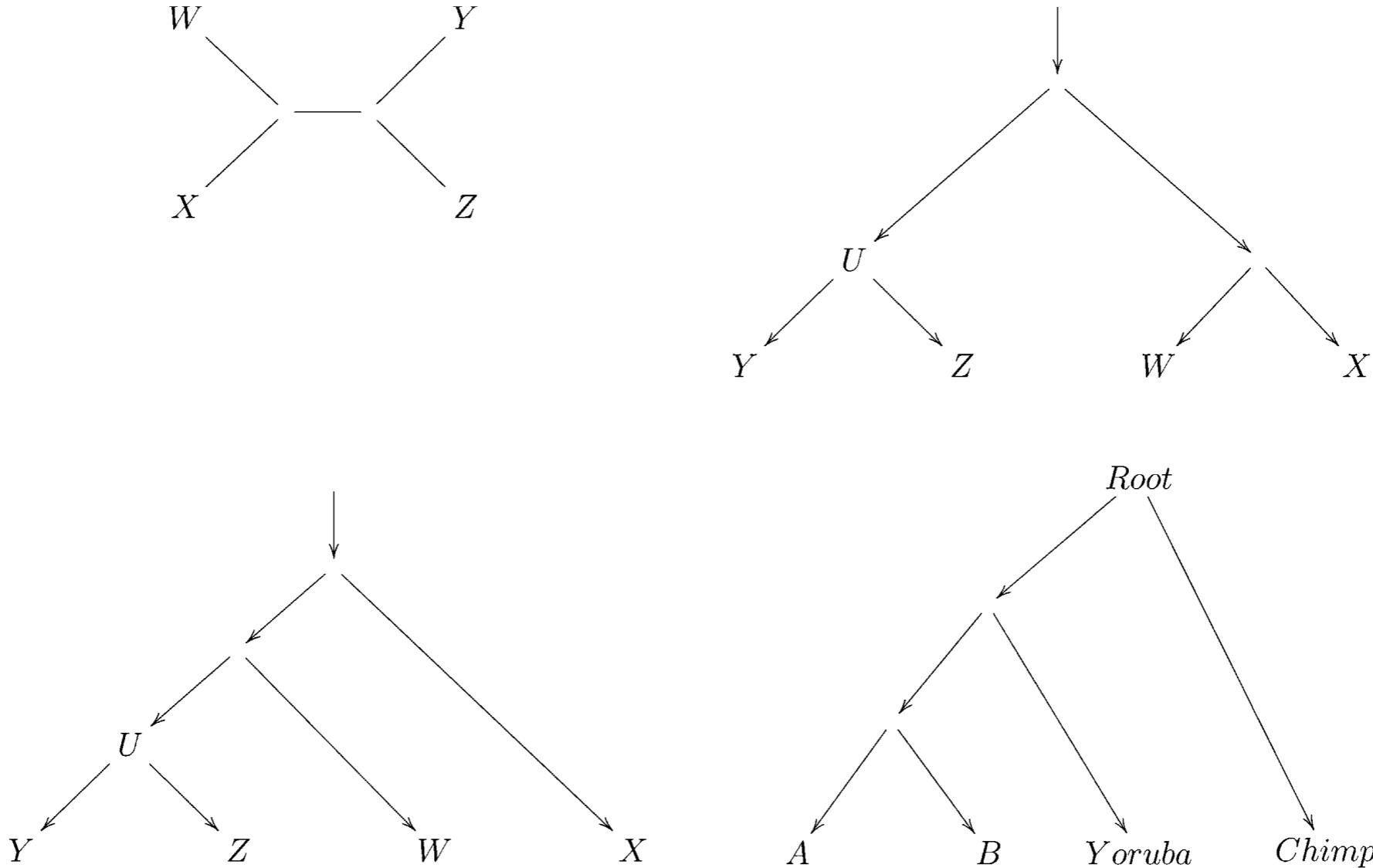
$$\frac{F_2(P_1, P_2)}{\mathbb{E}[(p_1 - p_2)^2]}$$

$$\frac{F_3(P_X; P_1, P_2)}{\mathbb{E}(p_X - p_1)(p_X - p_2)}$$

$$\frac{F_4(P_1, P_2, P_3, P_4)}{\mathbb{E}(p_1 - p_2)(p_3 - p_4)}$$

Interpretation in terms of branch lengths on a population phylogeny

Expectations and usage - f_4

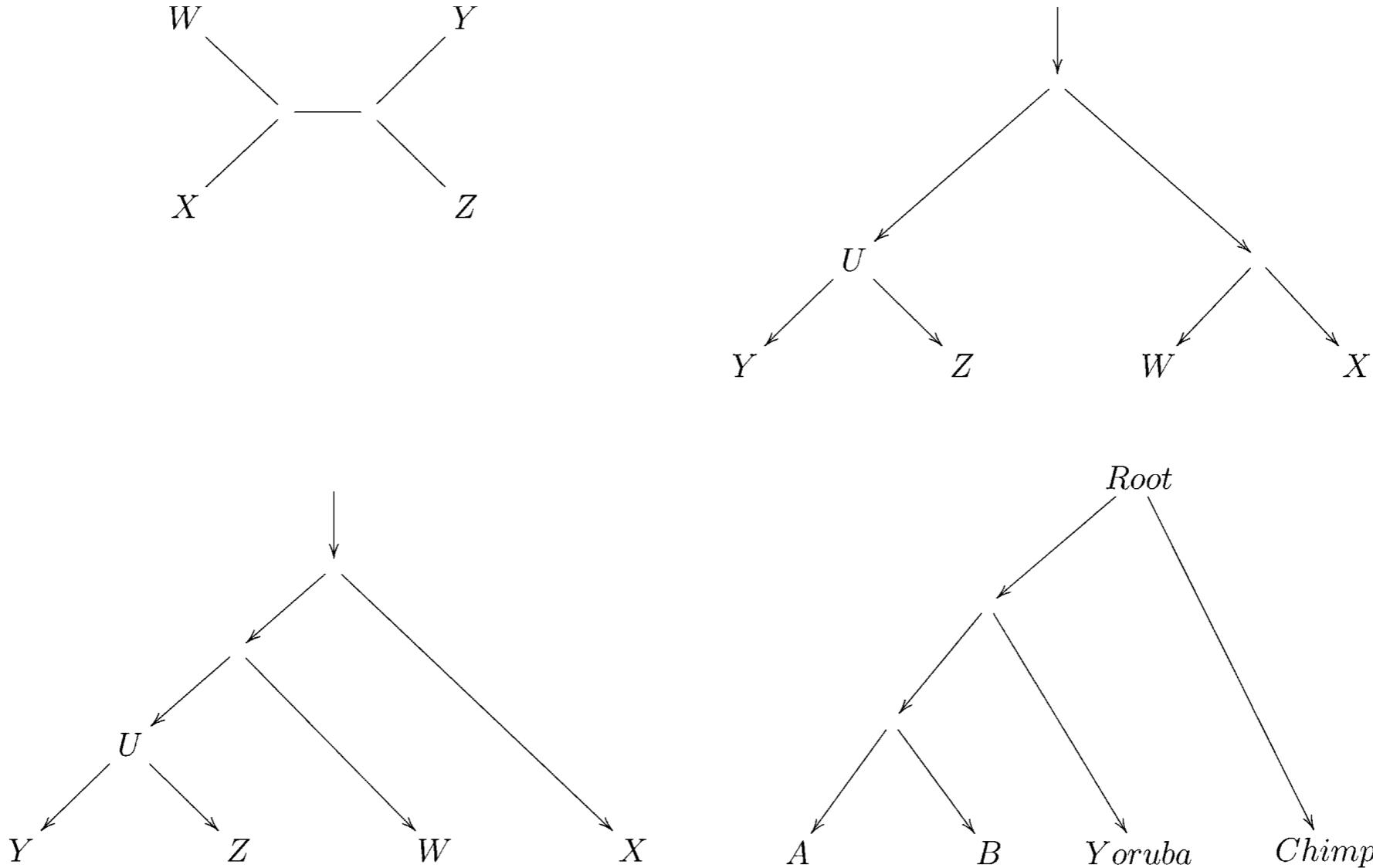


f_4 tests whether four populations are related through a simple tree
Test for “treeness”

Expectations and usage - f_4

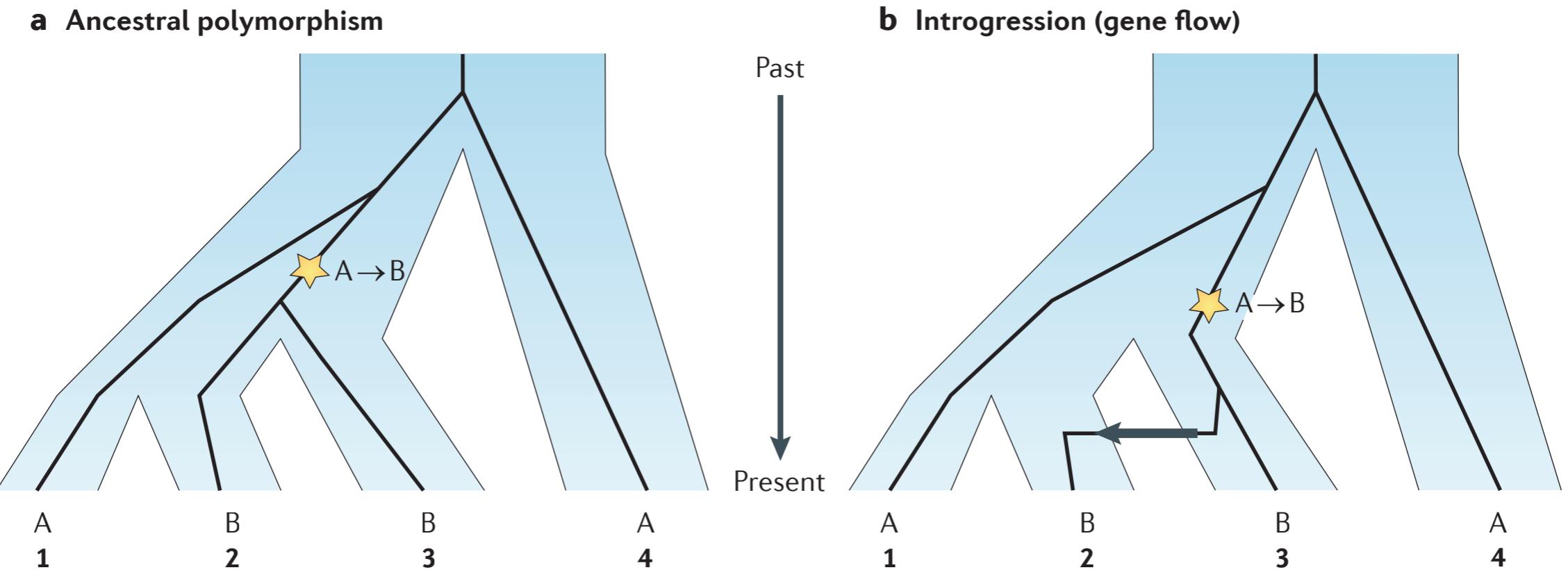
f_4 expectations on whiteboard

Expectations and usage - f_4



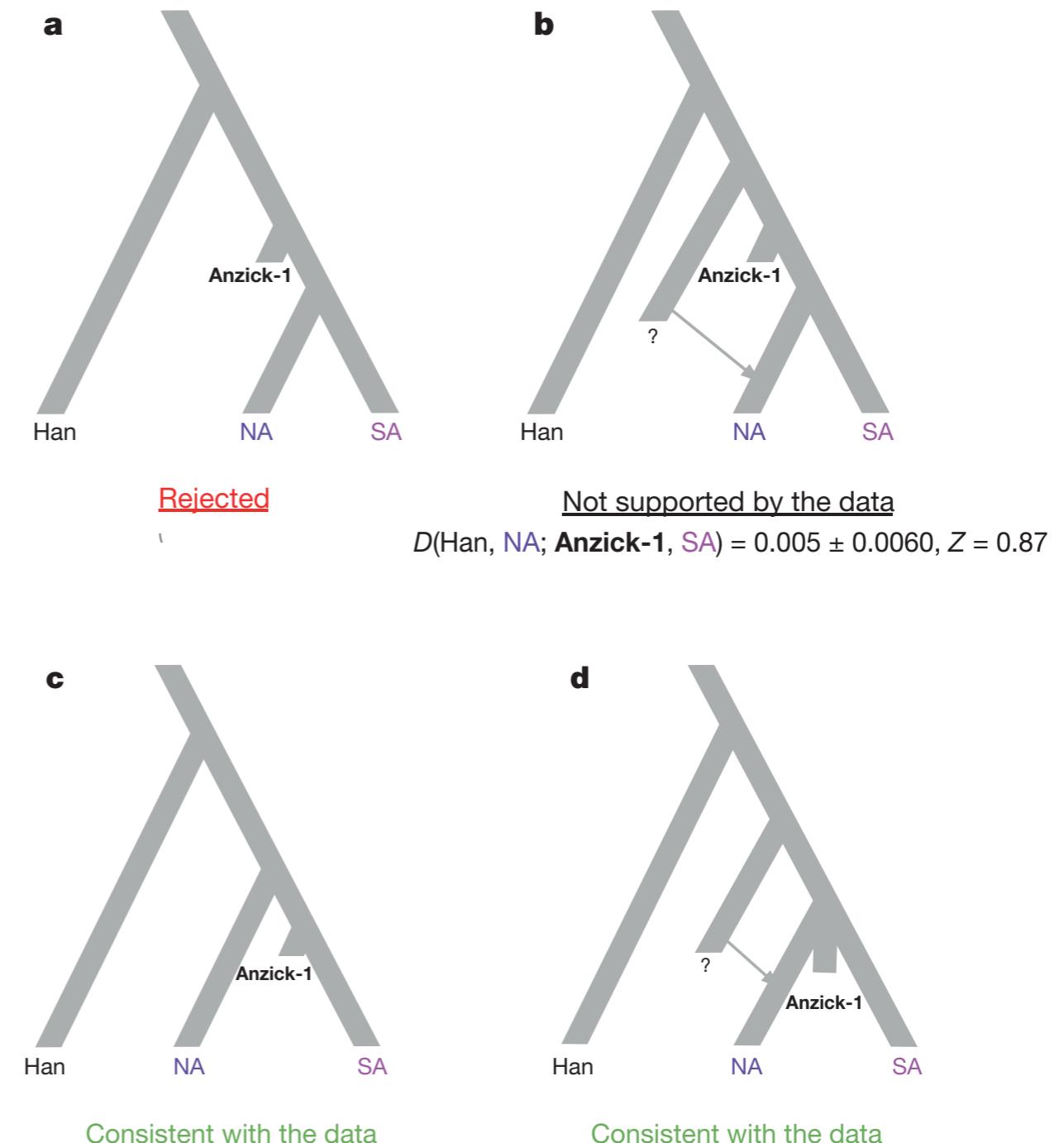
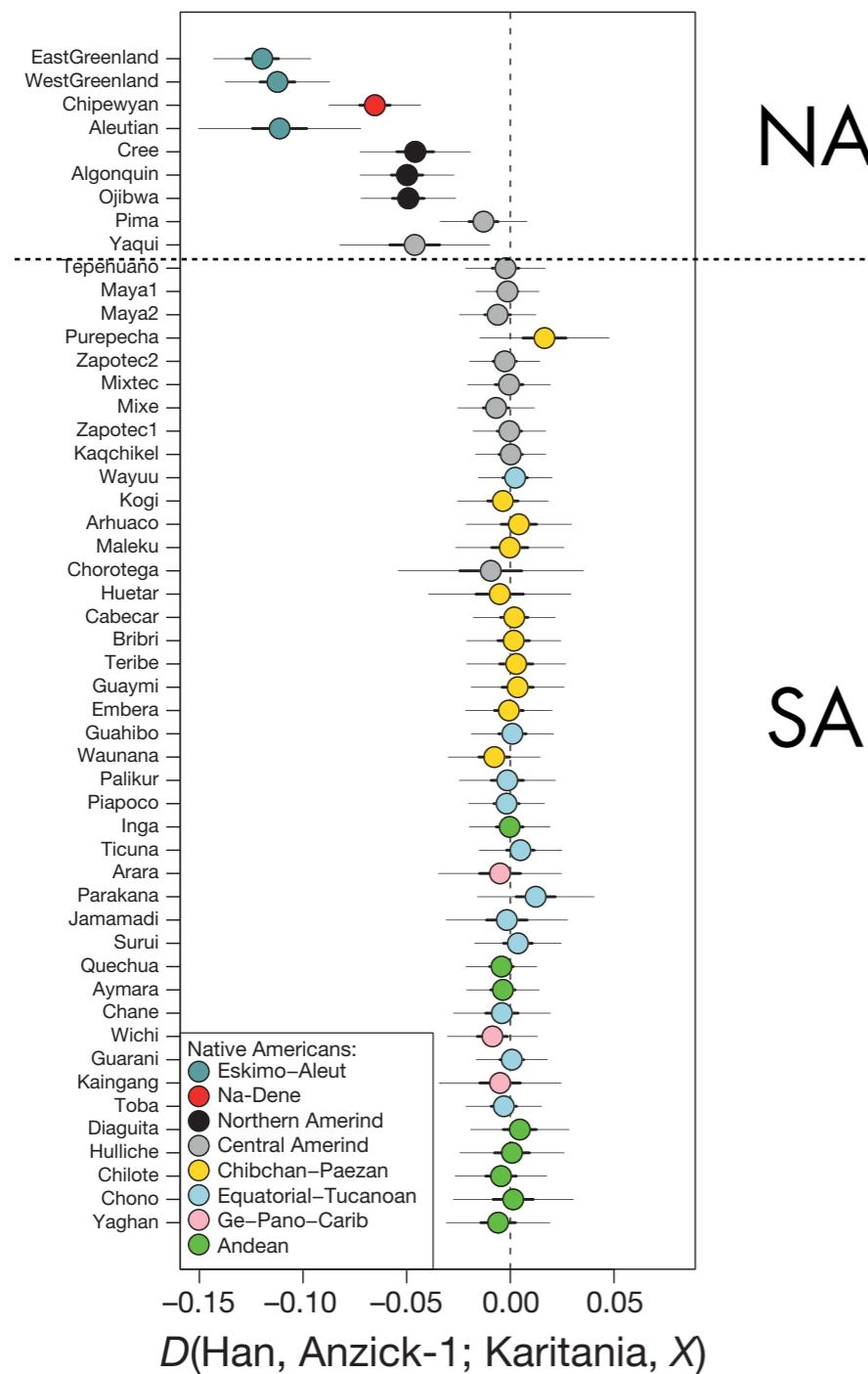
f_4 tests whether four populations are related through a simple tree
Test for “treeness”

f_4 is closely related to D-test / ABBA-BABA test



Test for Neanderthal admixture is test of treeness of modern humans

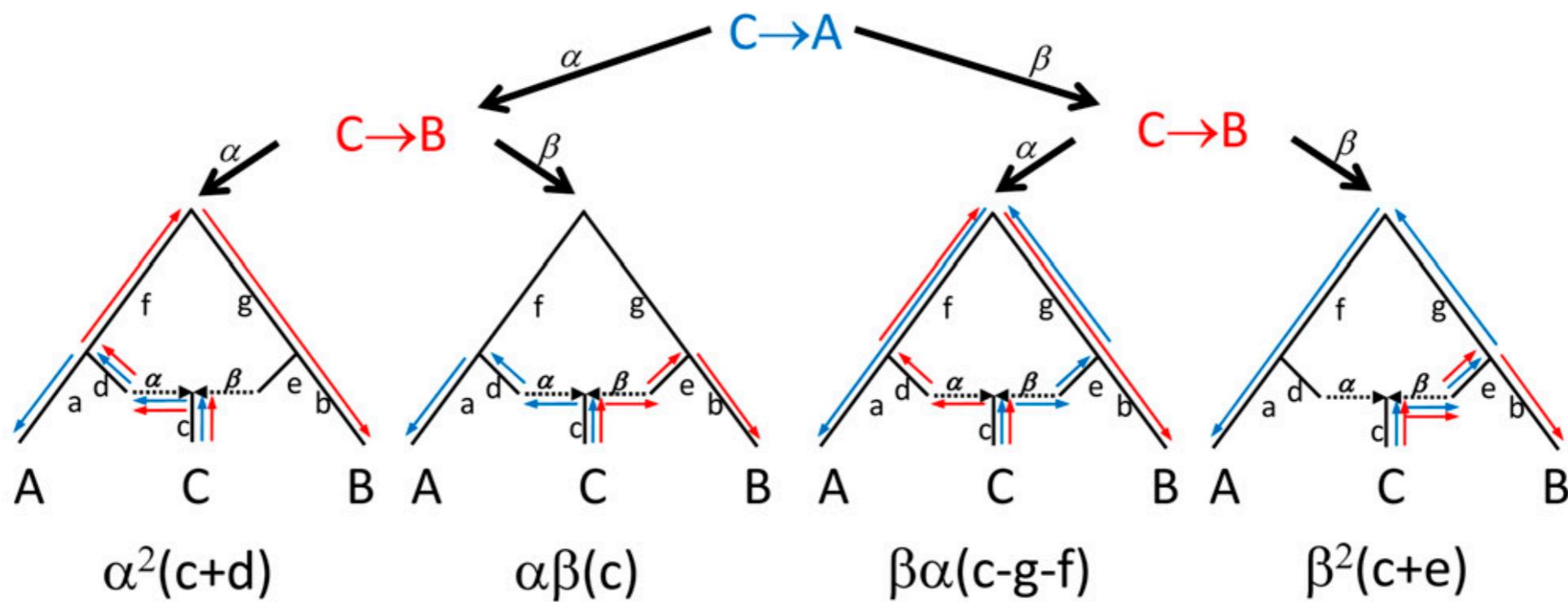
Interpretation of failed treeness tests



Failed treeness test can be interpreted in multiple ways!

Expectations and usage - f_3

$$F_3(C;A,B) = c + \alpha^2d + \beta^2e - \alpha\beta(g+f)$$



f_3 tests whether a target population C is the product of admixed from two source populations related (possibly distant) to sample populations A and B

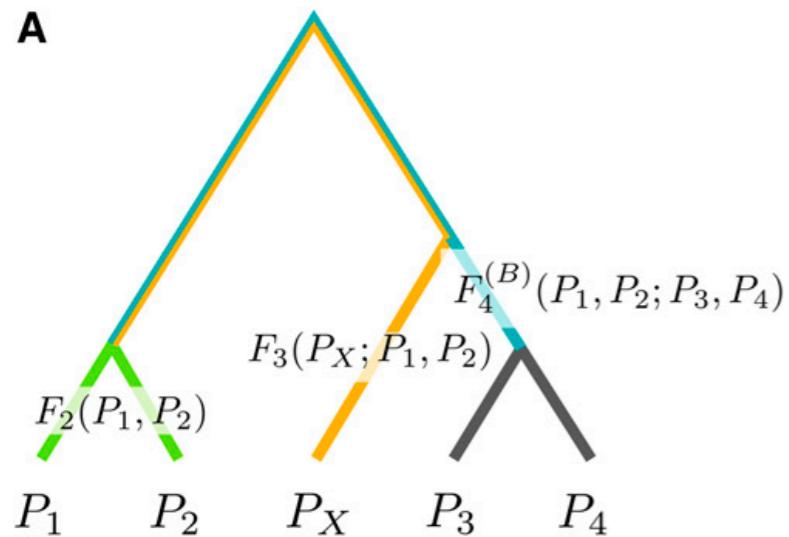
Statistical testing

F-statistic	Test	Interpretation
$f_3(X; A, B)$	$f_3 < 0$	X is admixed related to A, B
$f_4(A, B; C, D)$	$f_4 = 0$	(A, B) form a clade with respect to (C, D)
f_4 -ratio	$\alpha > 0$	Admixture proportion > 0

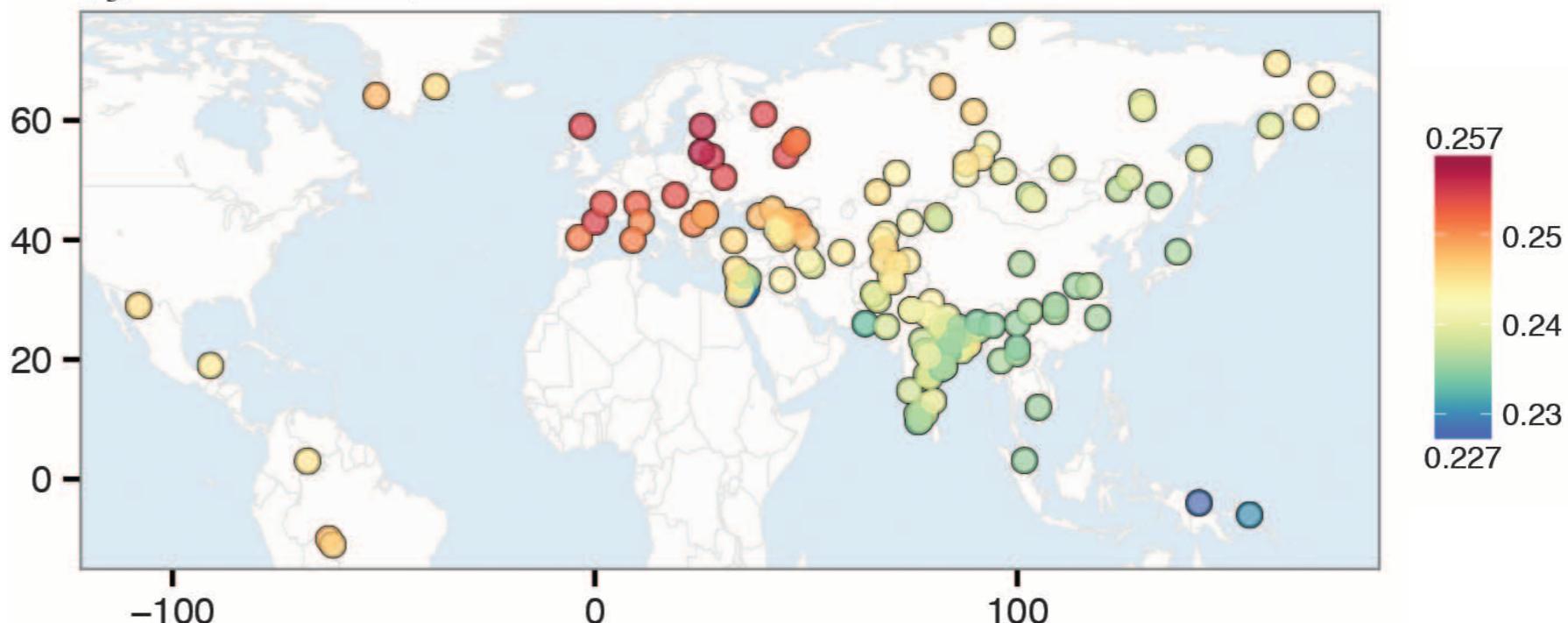
Standard errors are estimated using a weighted block jackknife and converted into a Z-score to determine significance

Outgroup f_3 statistics

A



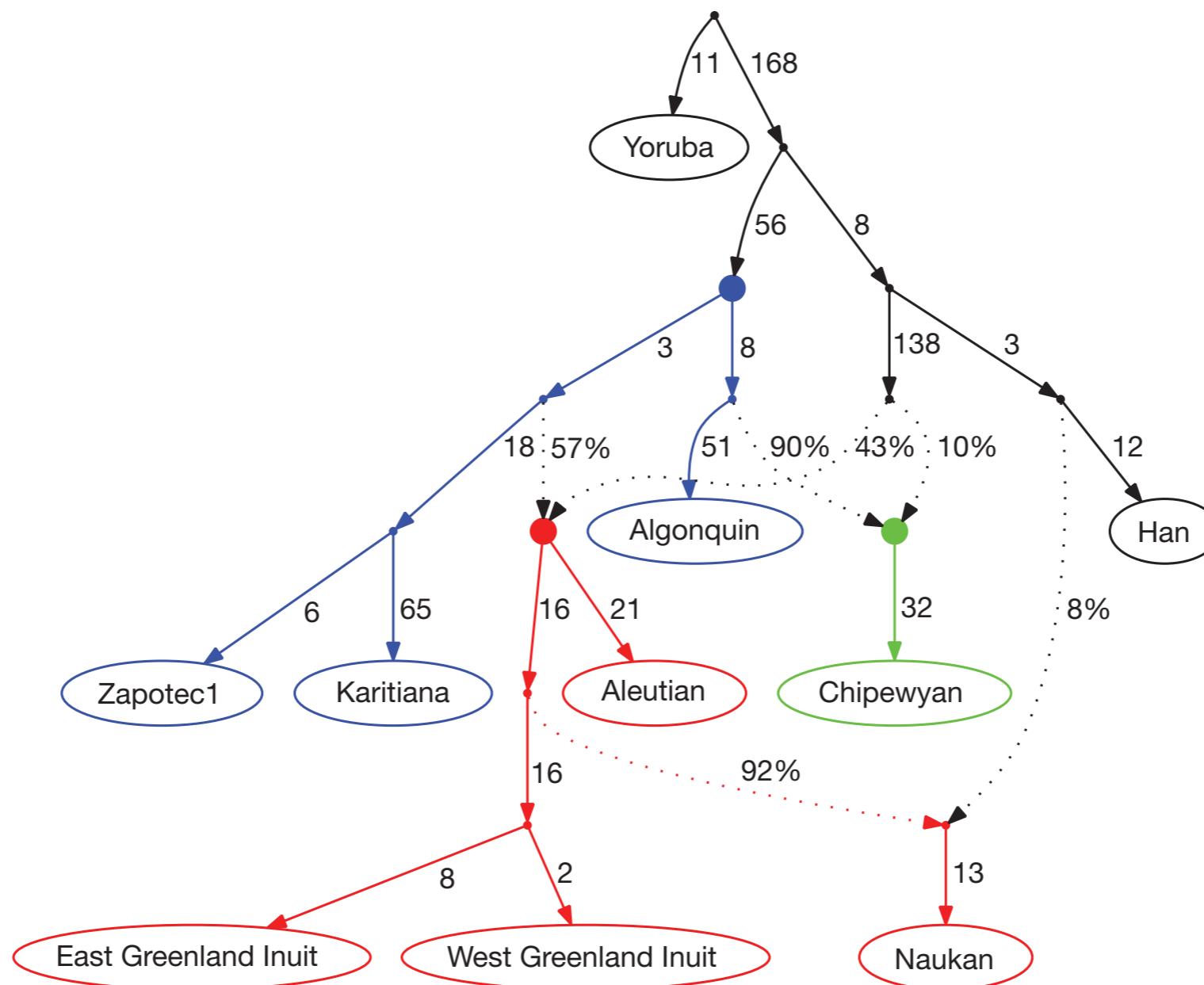
$f_3(\text{Mbuti}; \text{K14}, \text{X})$ - worldwide



In "outgroup" f_3 statistics, the target admixed population is replaced by an outgroup population

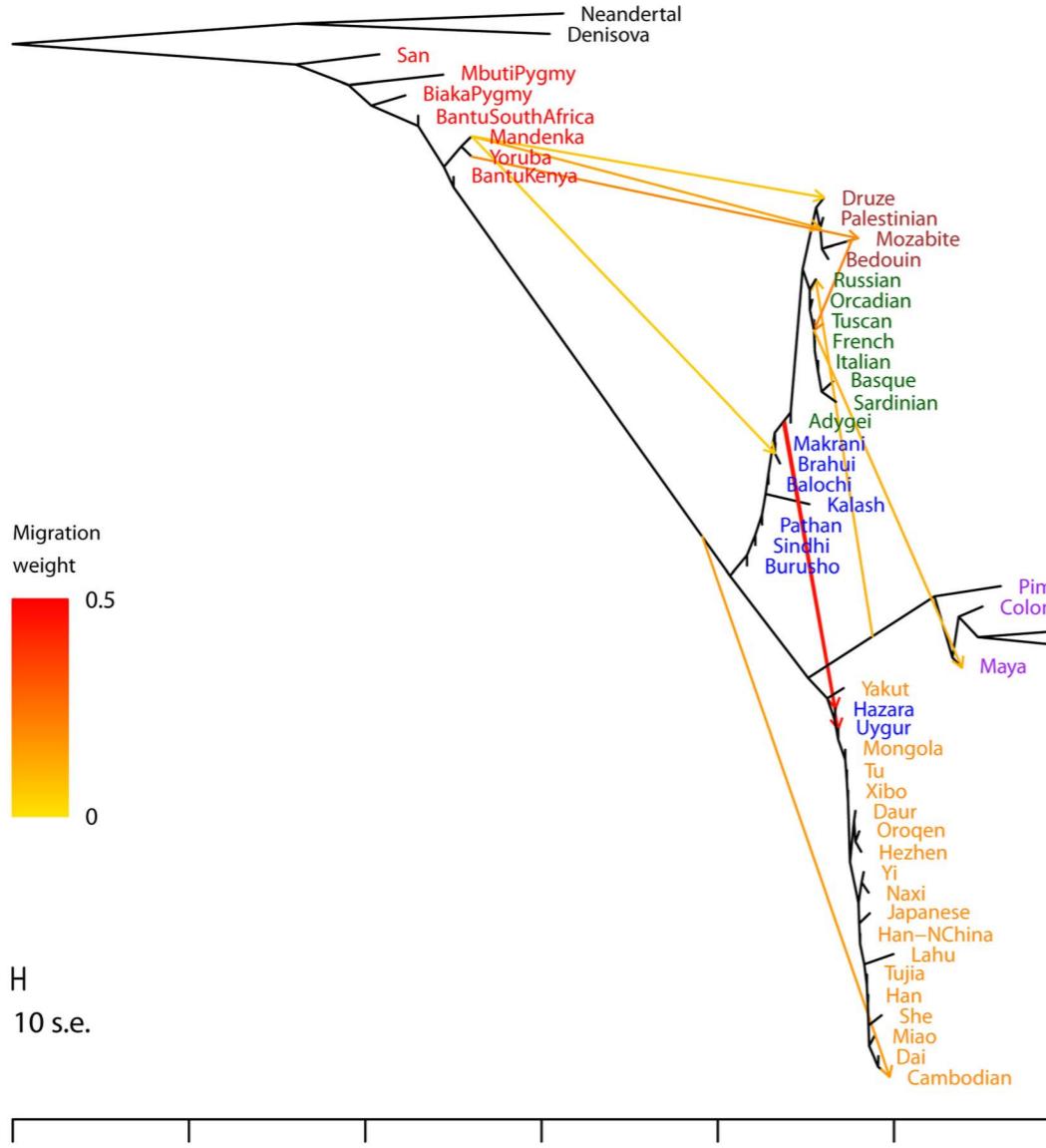
It measures the shared drift between the outgroup and the divergence of the two other populations

Admixture graph fitting

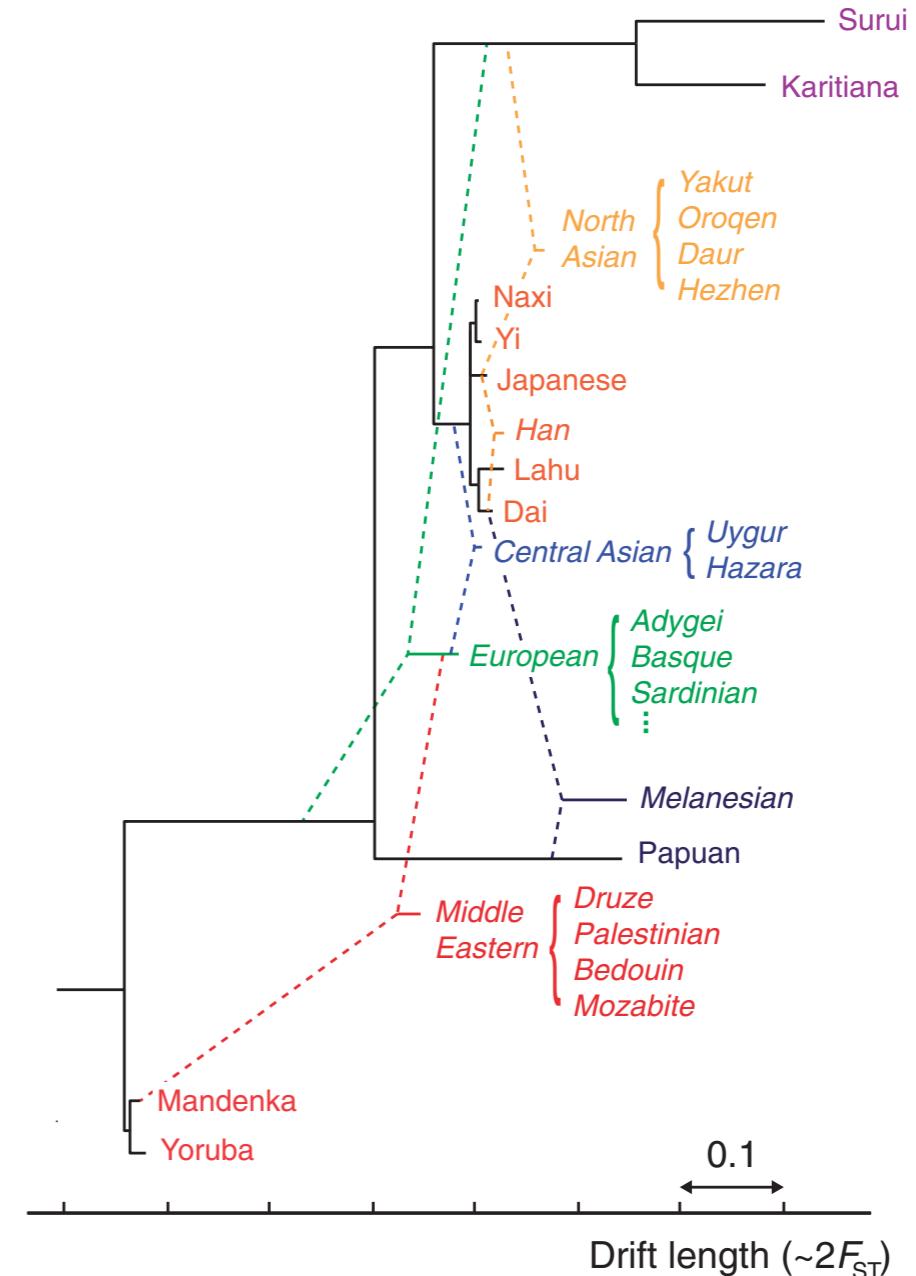


Given an admixture graph, calculate expected drift and admixture terms, and compare fit with observed values

Related methods

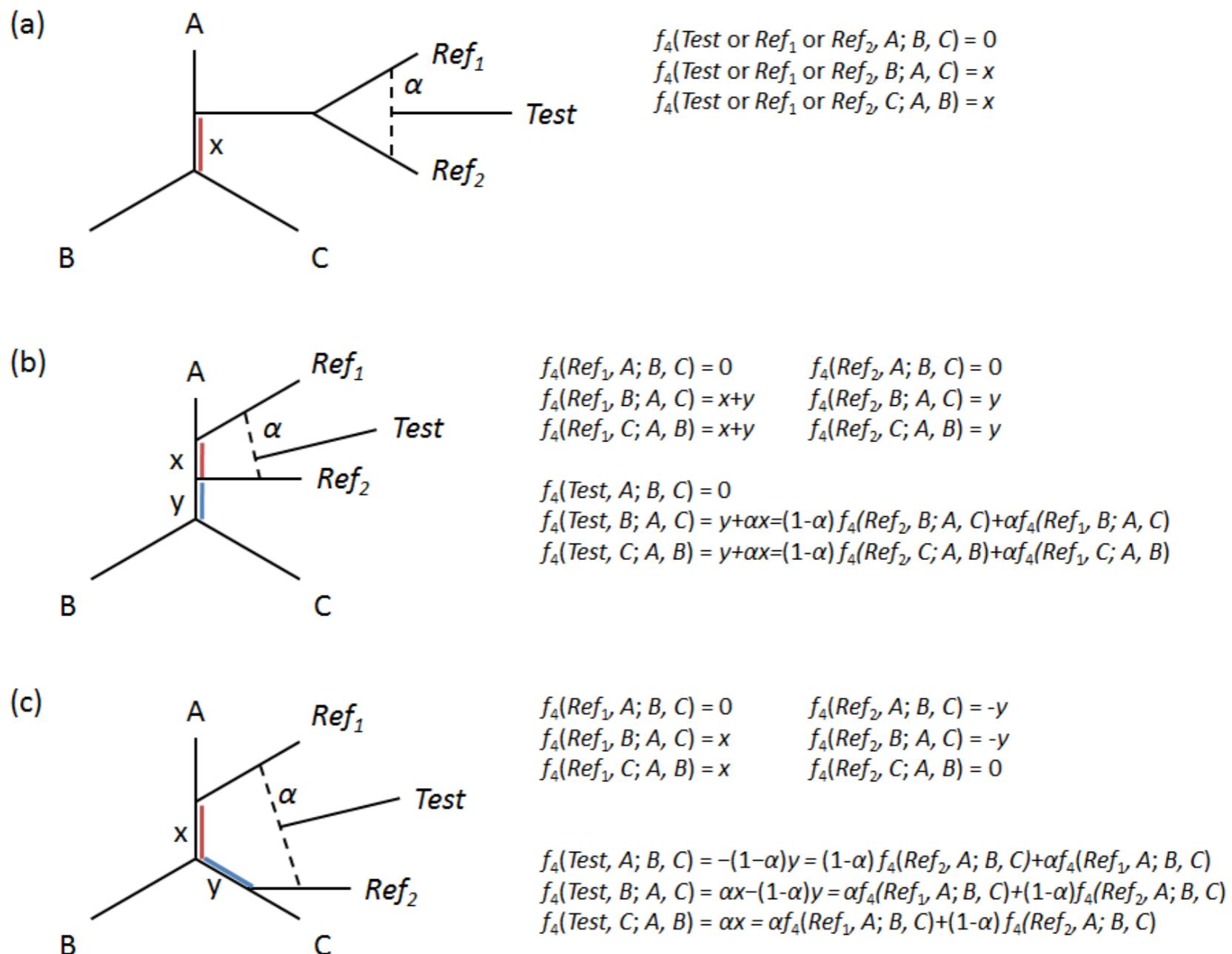


Treemix



MixMapper

Phylogeny-free admixture fitting (qpAdm)



$$f_4(\text{Test}; A; B, C) \approx \sum_{i=1}^N \alpha_i f_4(\text{Ref}_i; A; B, C)$$

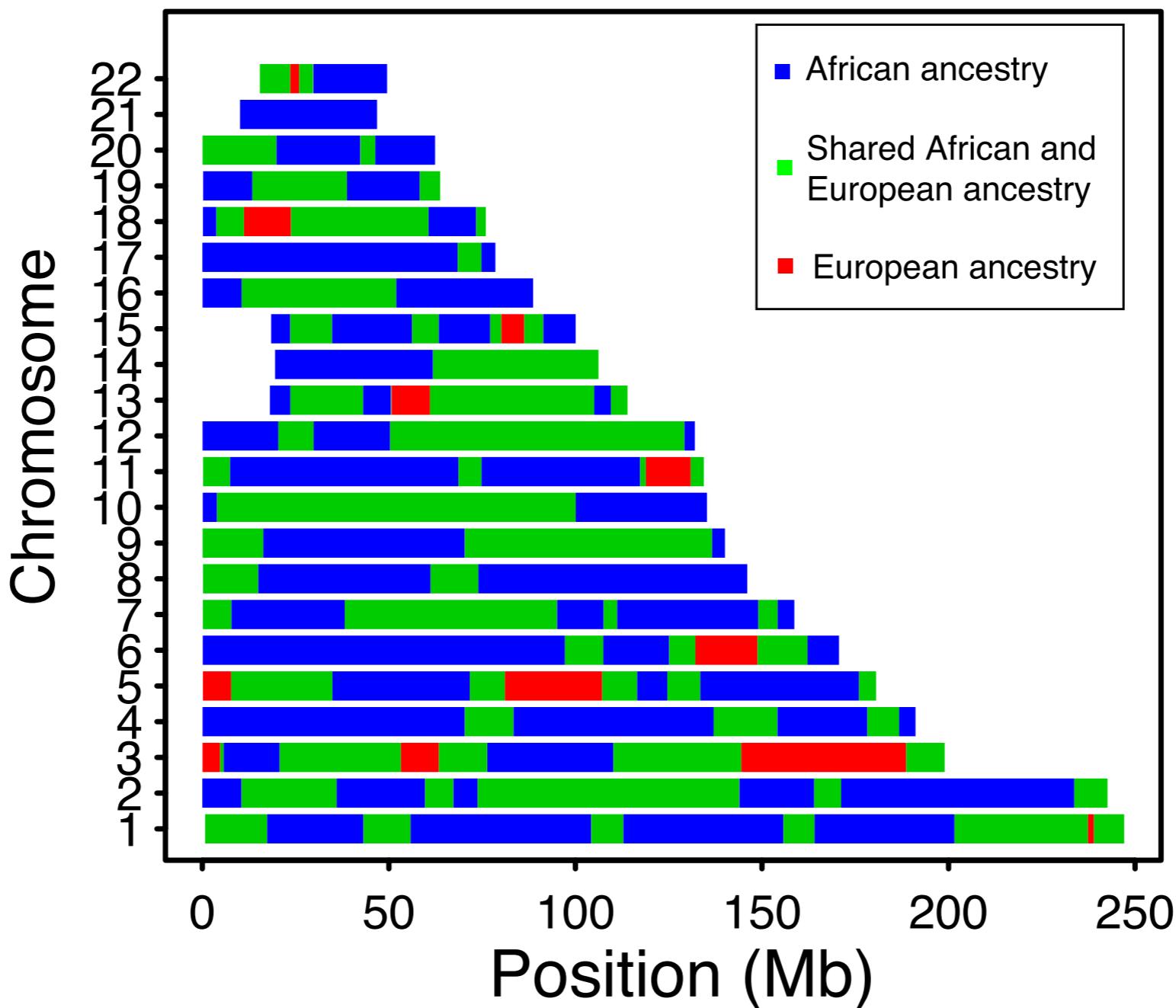
Differential drift sharing with sets of outgroup populations is leveraged to infer mixture proportions

Outline of today's lecture

- Motivation and Background
 - Admixture in human history
 - Overview of methods
- Testing for admixture using f-statistics
 - Basic concepts
 - What can we test with f-statistics?
 - Estimation and interpretation
 - Admixture graph fitting
- Admixture dating
 - Basic concepts
 - ALDER

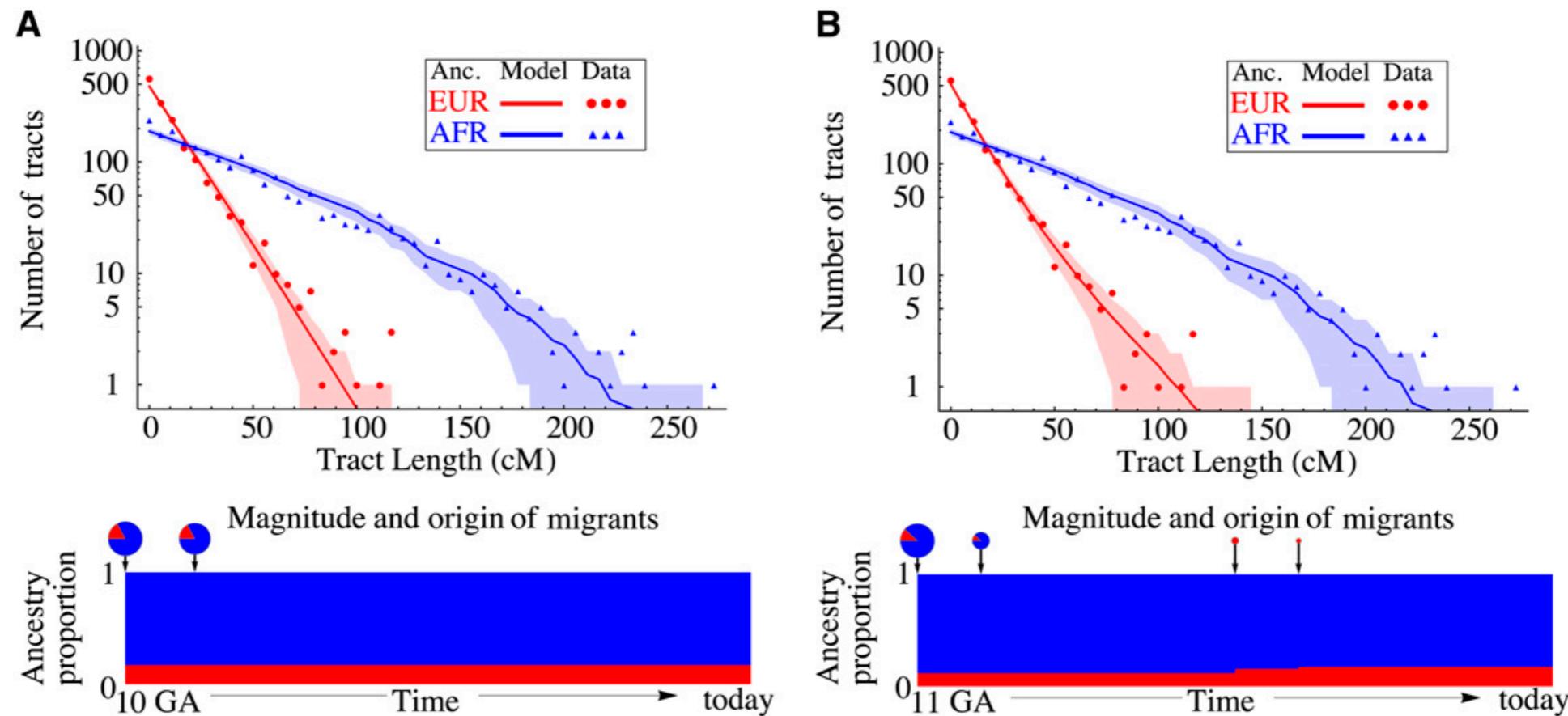
Admixture Linkage disequilibrium

Representative African American



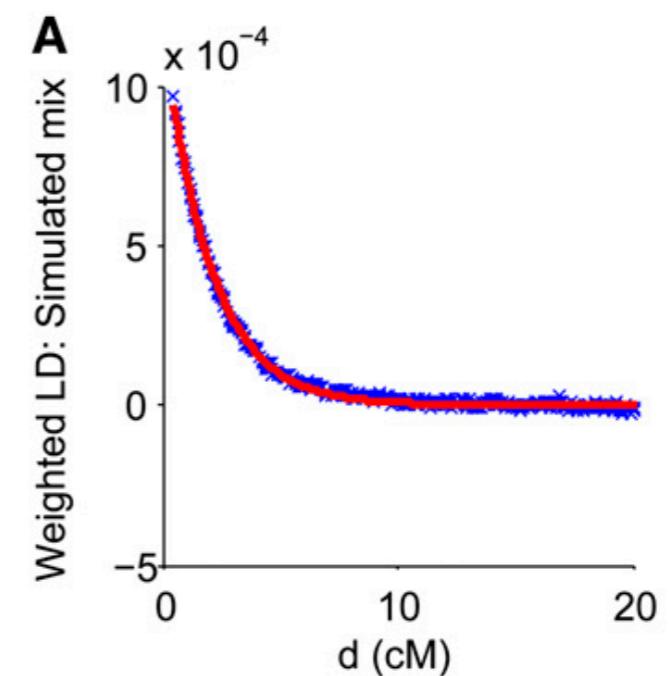
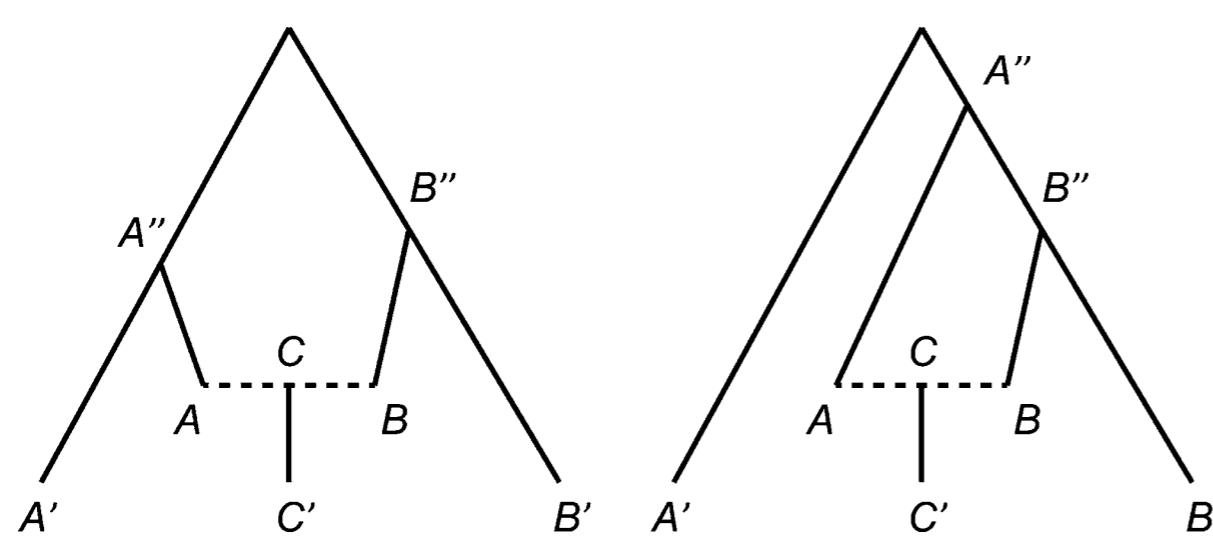
The length distribution of ancestry tracts is informative about time since admixture

Local methods



- Allow for complex admixture models
- Need phased data
- Less powerful for older admixture events

ALDER weighted LD model



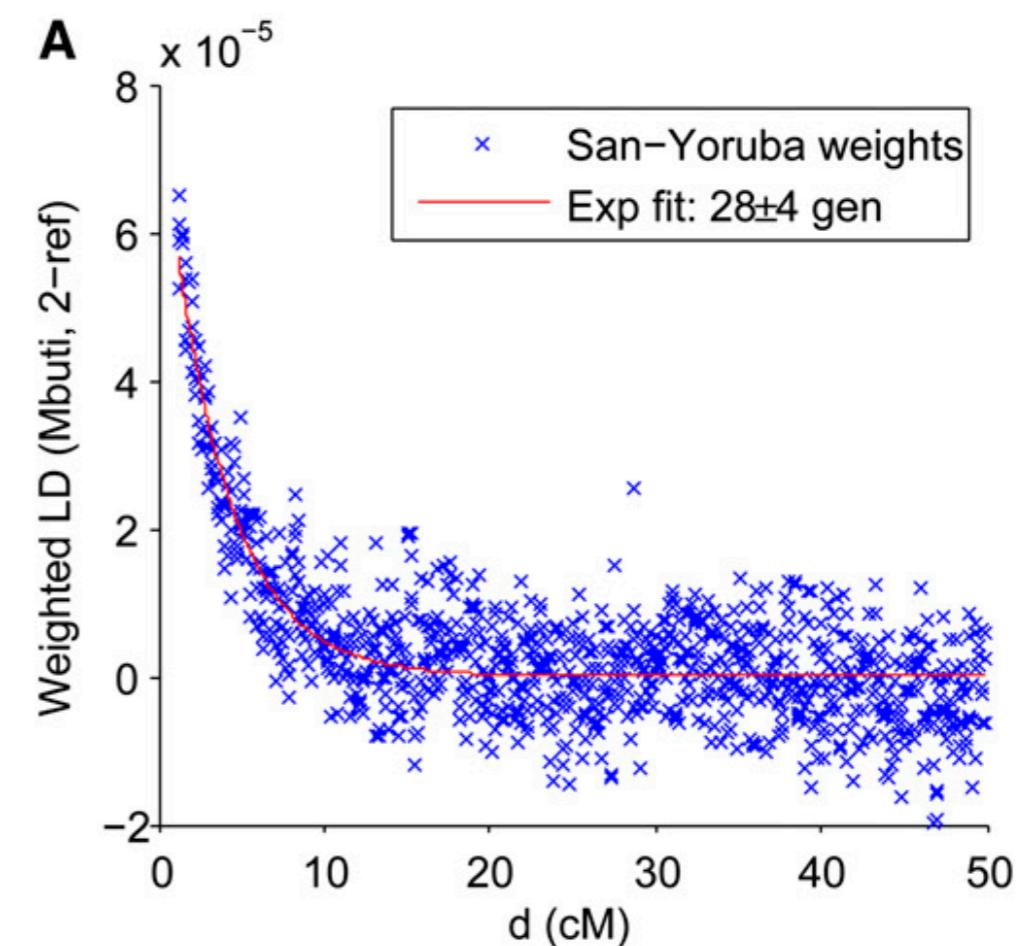
$$D_n = e^{-nd} D_0 = e^{-nd} \alpha \beta \delta(x) \delta(y)$$

$$a(d) := \frac{\sum_{\mathcal{S}(d)} z(x, y) w(x) w(y)}{|\mathcal{S}(d)|}. \quad E[a(d)] = 2\alpha\beta F_2(A'', B'')^2 e^{-nd},$$

LD between markers weighted by product of allele frequencies shows exponential decay with rate of generations since admixture

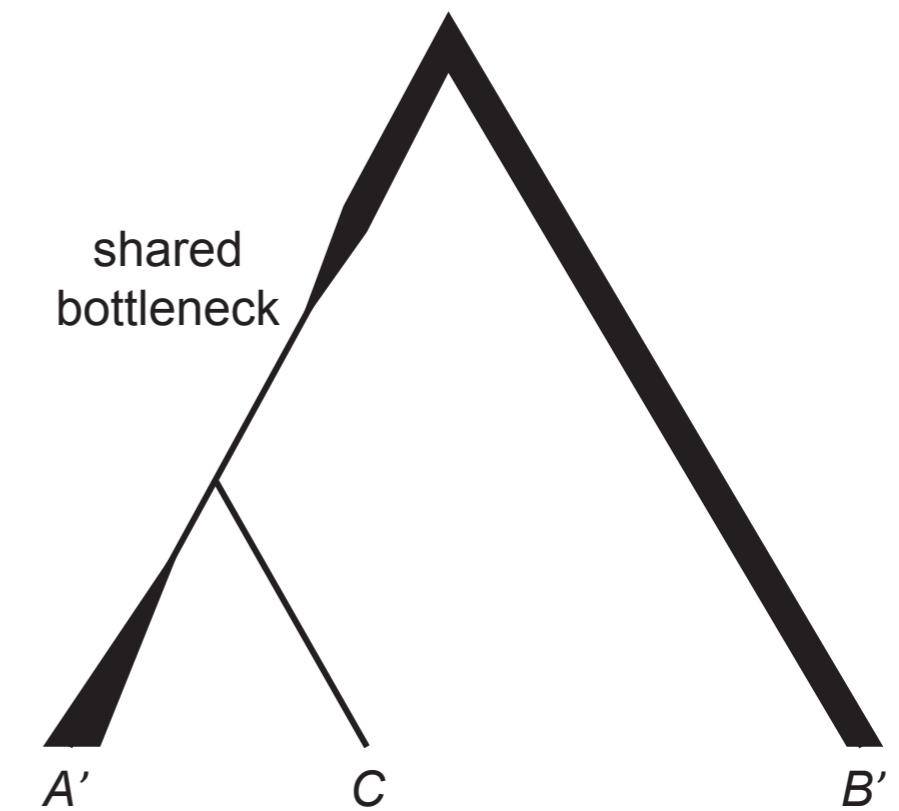
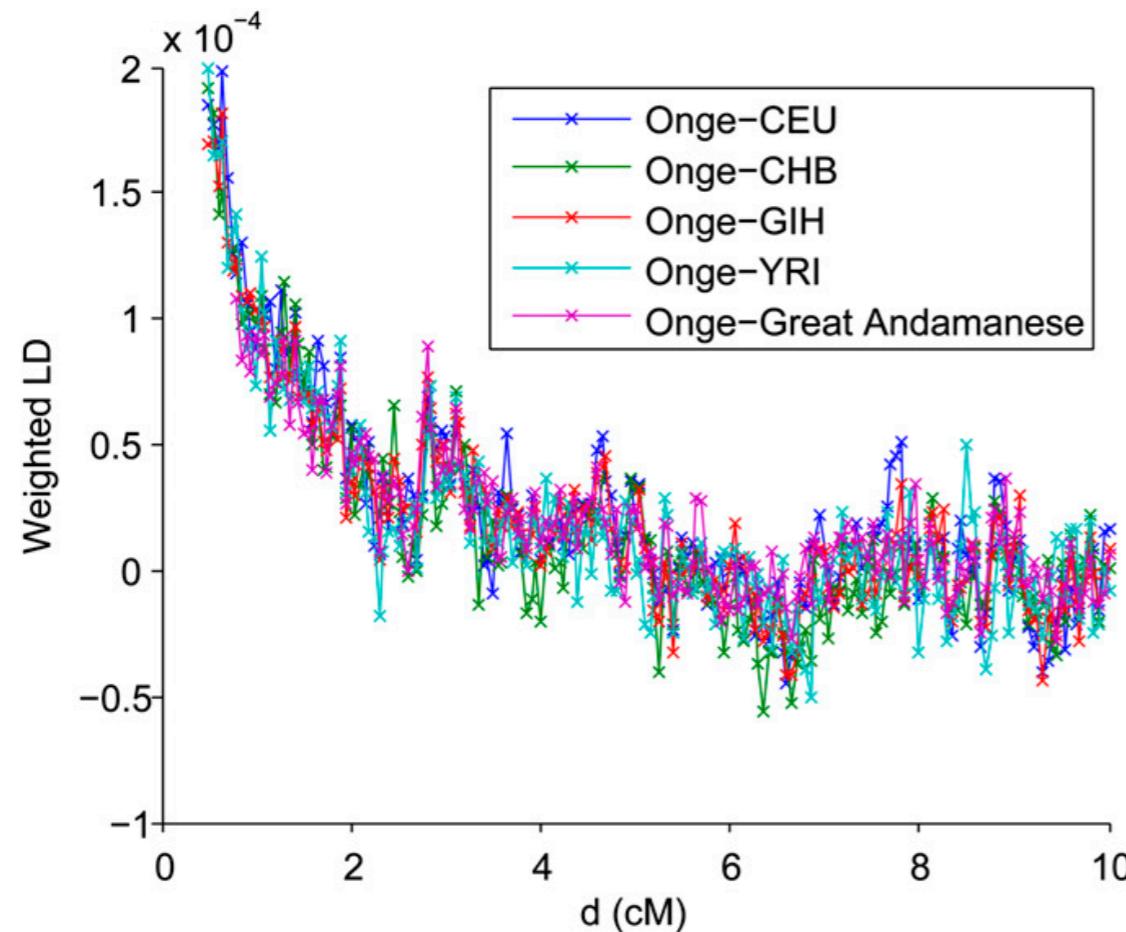
ALDER features

- Estimate time since admixture (decay of exponential)
- Infer phylogenetic placement of different source populations (amplitude)
- Can use admixed population itself as one reference
- Phasing not needed
- Single-pulse model, although implementation for multiple pulses exists (MALDER)



$$E[a(d)] = 2\alpha\beta F_2(A'', B'')^2 e^{-nd},$$

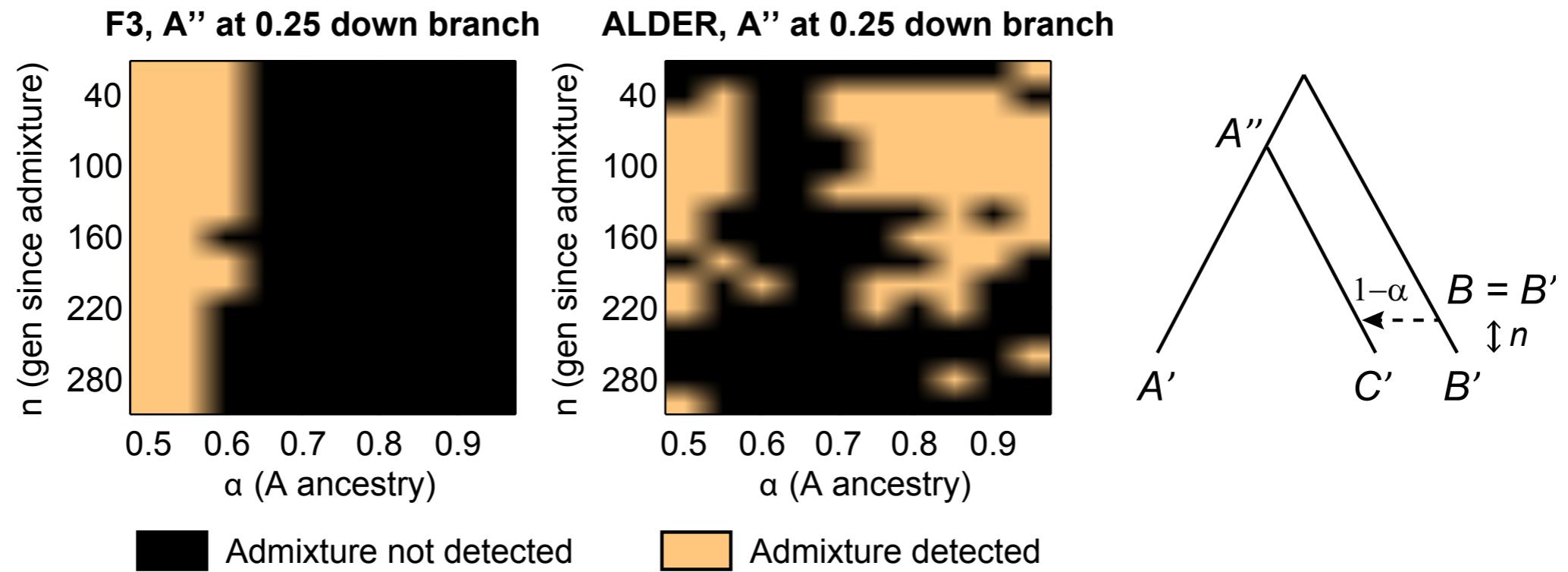
Weighted LD curves can occur without admixture



Single reference population curves have same amplitude

Two reference population curves more robust

ALDER vs f_3 -statistics



Complementary for detecting admixture

ALDER less affected by drift since admixture or very diverged source populations

f_3 more power for older events

Take - home messages

- f-statistics / ALDER are versatile frameworks for testing admixture hypotheses, including quite complex scenarios
- Suitable for low quality data (e.g. aDNA)
- Care has to be taken as f-statistics are susceptible to batch effects that induce spurious allele frequency correlations

Practicals

- Practical 1
 - f3 / f4 statistics
- Practical 2
 - qpGraph
 - ALDER