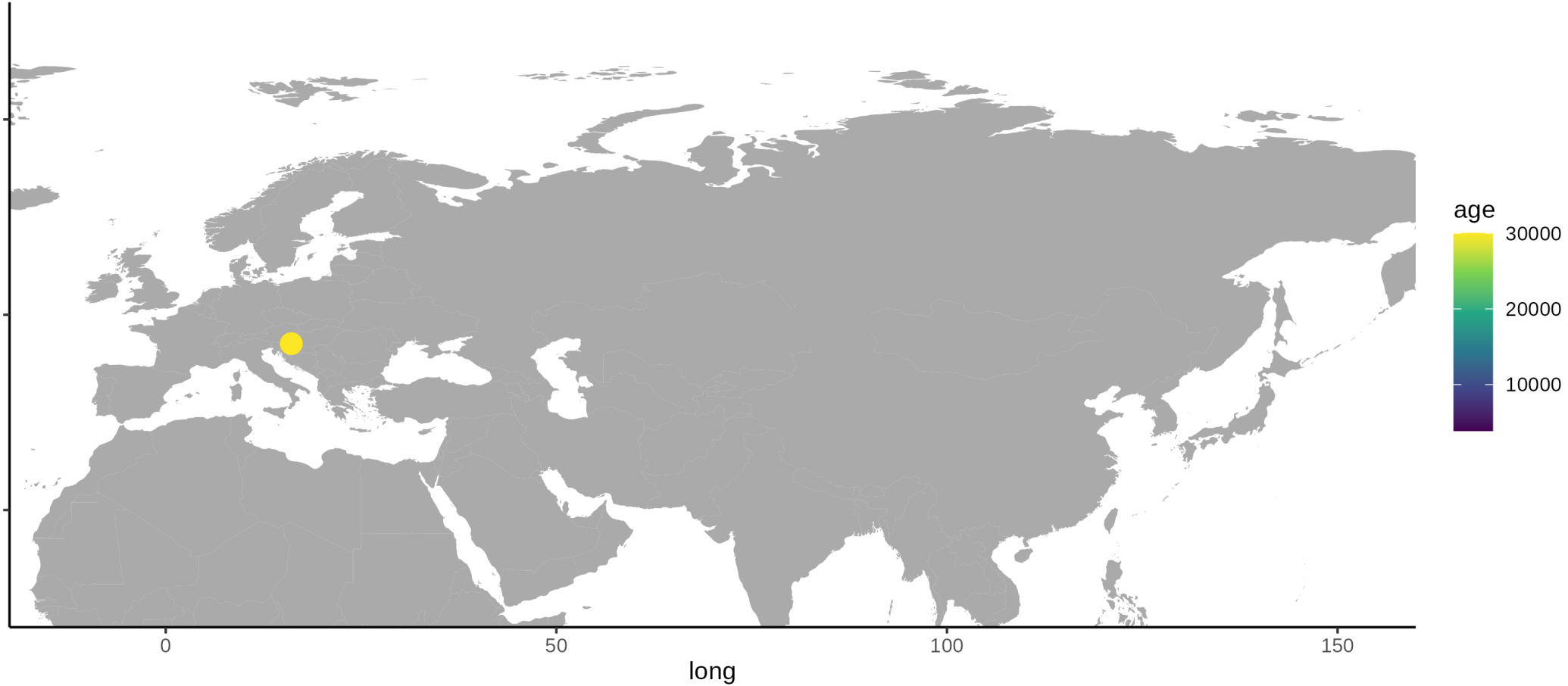# F-statistics and Population Structure
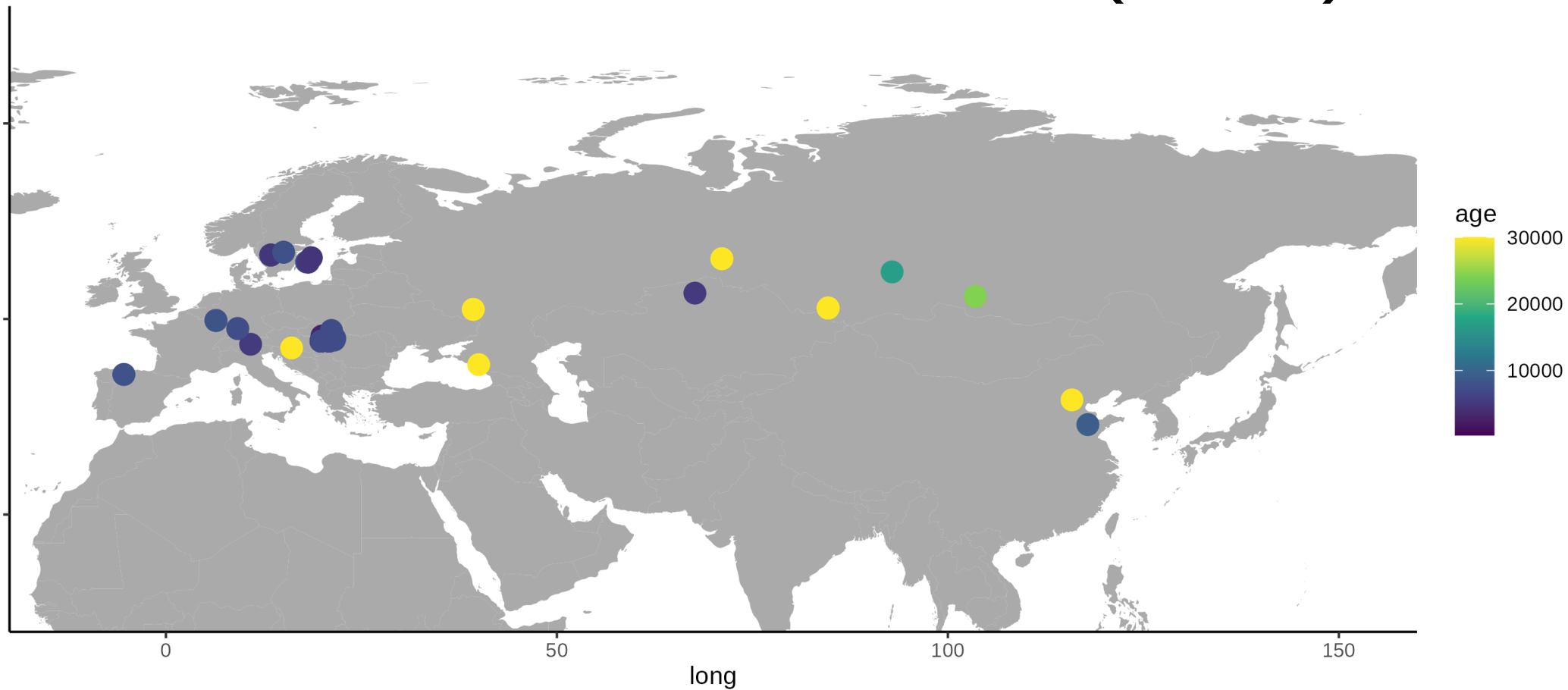
Benjamin Peter, MPI for Evolutionary Anthropology

# Hominin Ancient DNA (2010)
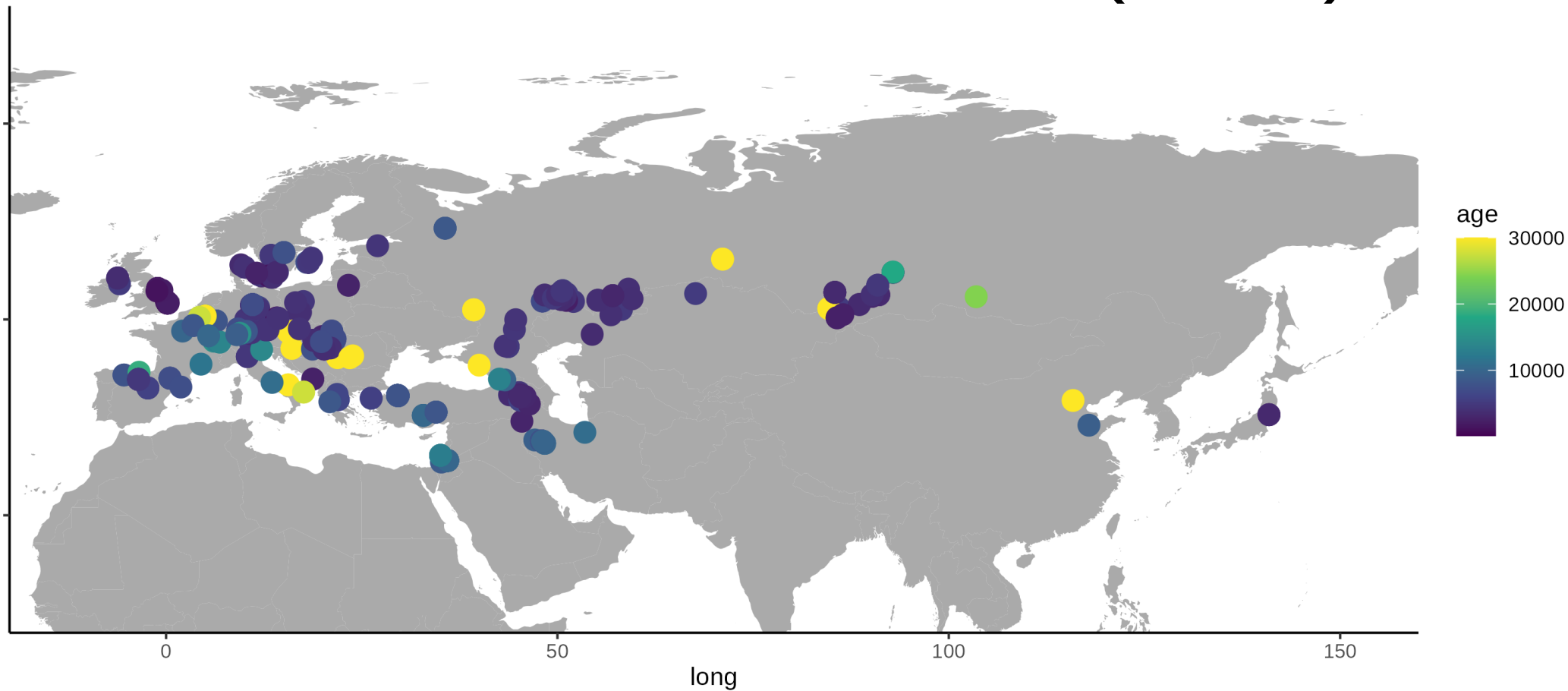
# Hominin Ancient DNA (2014)

# Hominin Ancient DNA (2016)
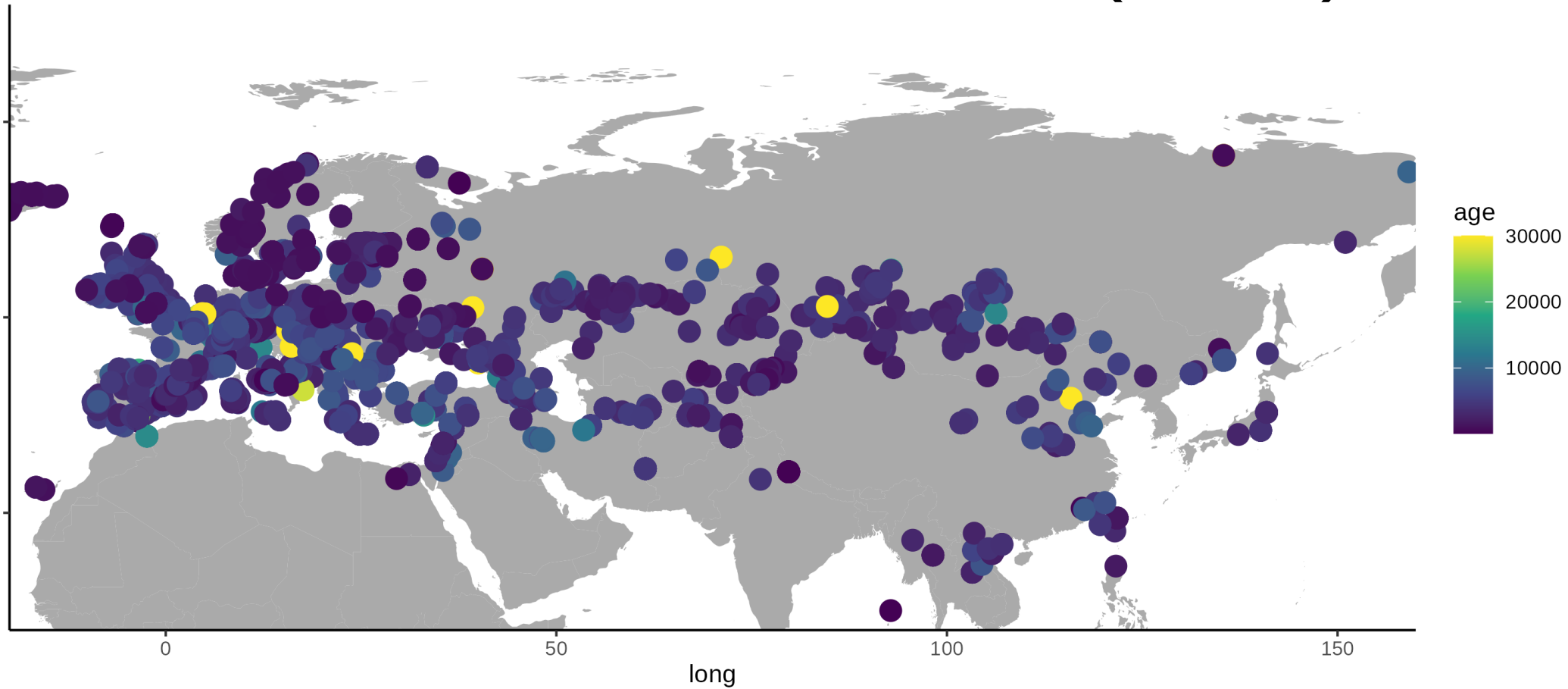


age

30000

20000

10000

long

# Hominin Ancient DNA (2020)

# Main Reference

## Admixture, Population Structure, and *F*-Statistics

**Benjamin M. Peter[1]**
Department of Human Genetics, University of Chicago, Chicago, Illinois 60637
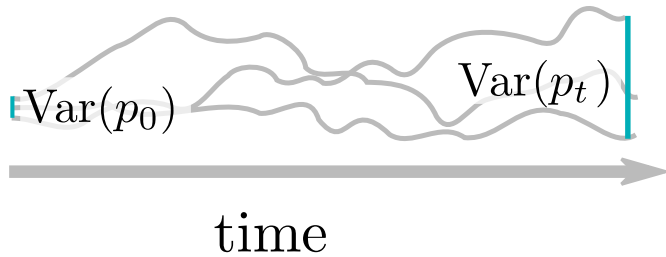ORCID ID: 0000-0003-2526-8081 (B.M.P.)

# Setup

- Today: Theory of F-statistics and Computations

- Tomorrow: Using F-statistics to build more complex models
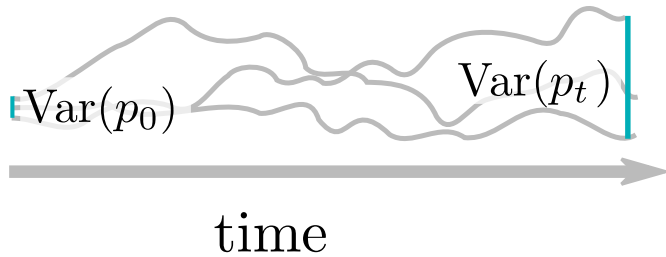
# Measuring Genetic Drift

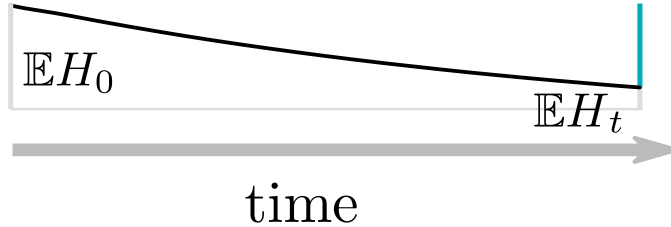# Measuring Genetic Drift

Change in Allele Frequency

# Measuring Genetic Drift
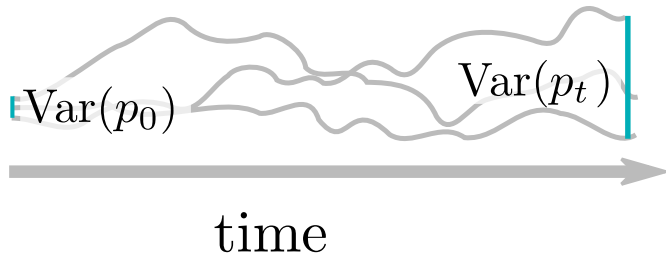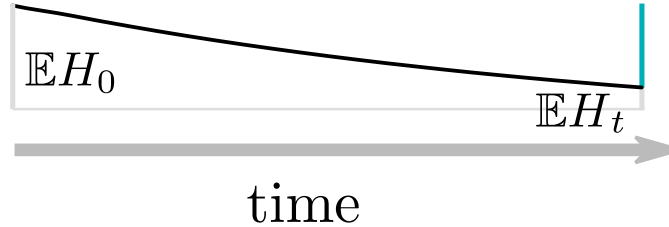
Change in Allele Frequency

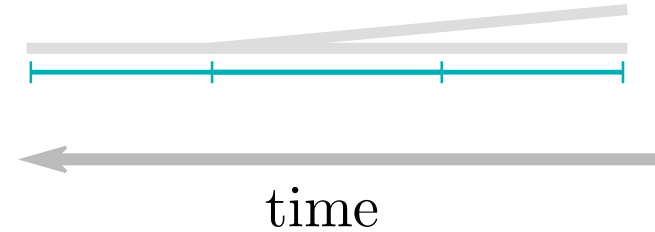Decay of Heterozygosity

# Measuring Genetic Drift

Change in Allele Frequency
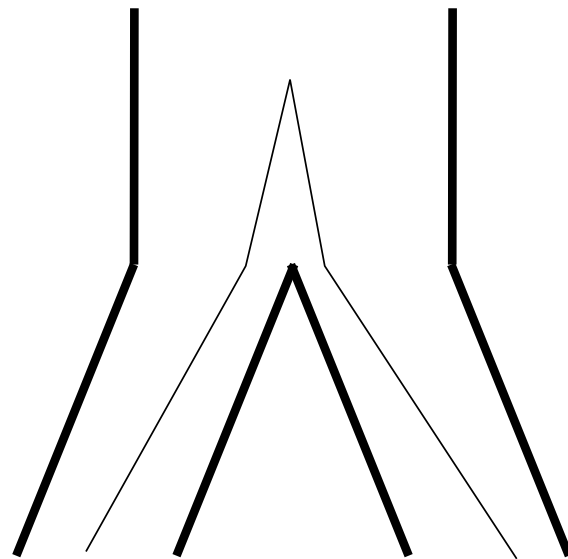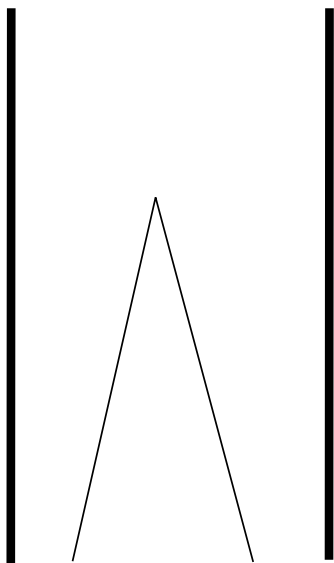
Decay of Heterozygosity

Coalescence rates

# Pairwise differences

$$\mathbb{E}[\pi] = 4N\mu = \theta \qquad \mathbb{E}[\pi_{12}] = t_{12} + 4N_{anc}\mu = t_{12} + \theta$$
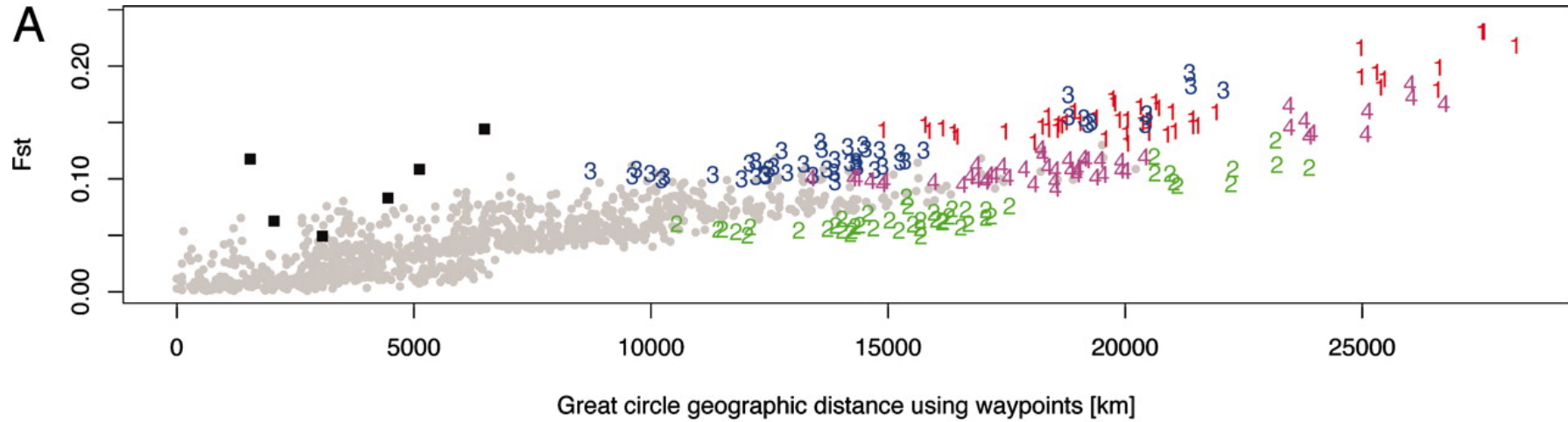
# Fixation Index $F_{ST}$

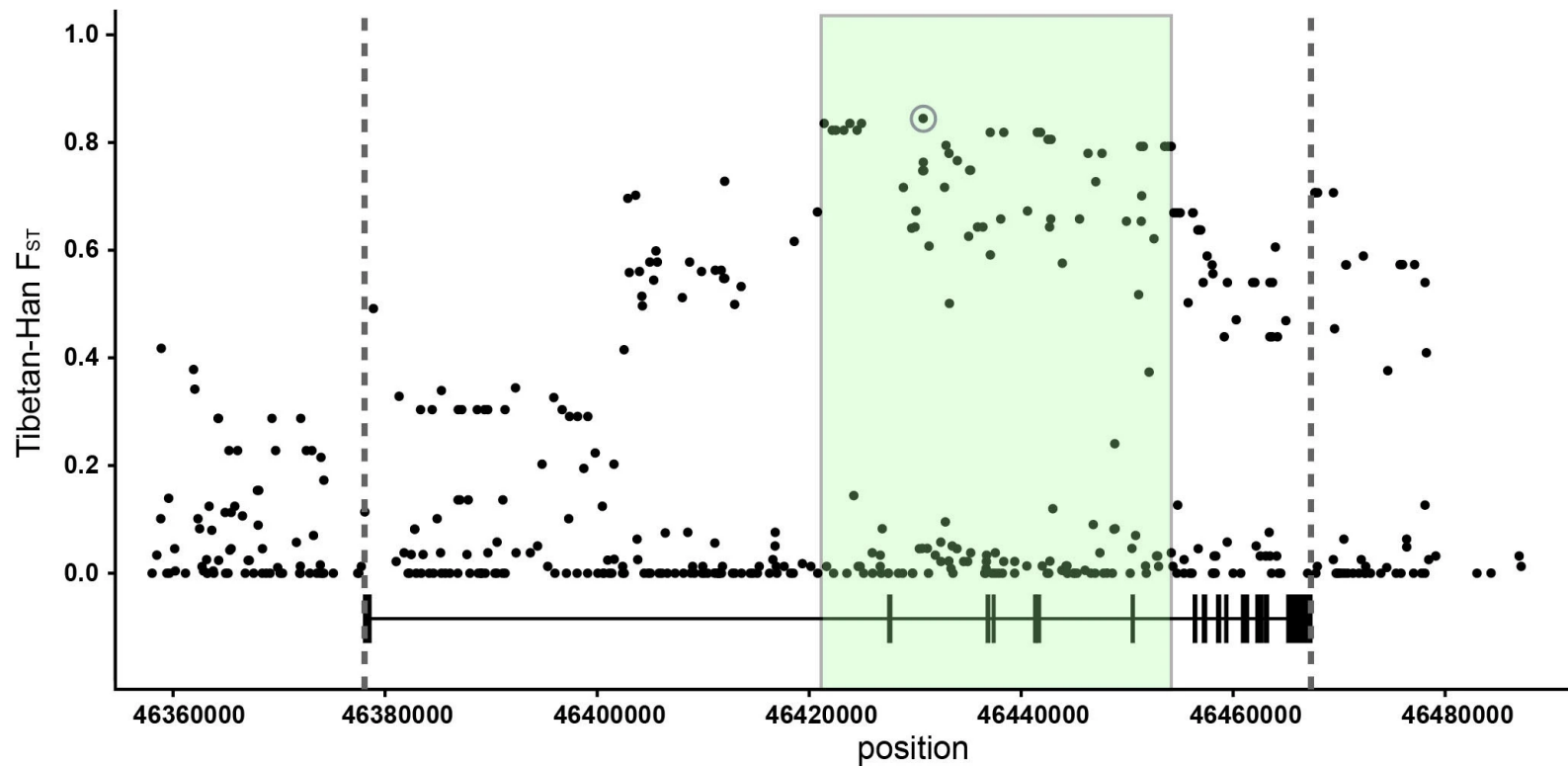$$F_{ST}(P_1, P_2) = \frac{\pi_{12} - \frac{\pi_1 + \pi_2}{2}}{\pi_{12}}$$

- $F_{ST}$ is a correlation coefficient
- Between 0 and 1
- Hierarchical partitioning (AMOVA)
- Many estimators exist
  - Hudson (1991)
  - Weir & Cockerham (1984)

# Fixation Index F$_{ST}$



Great circle geographic distance using waypoints [km]

Ramachandran et al. 2005

# F$_{ST}$ Outliers



Huerta-Sanchez et al. 2014

# $F_2$-statistic
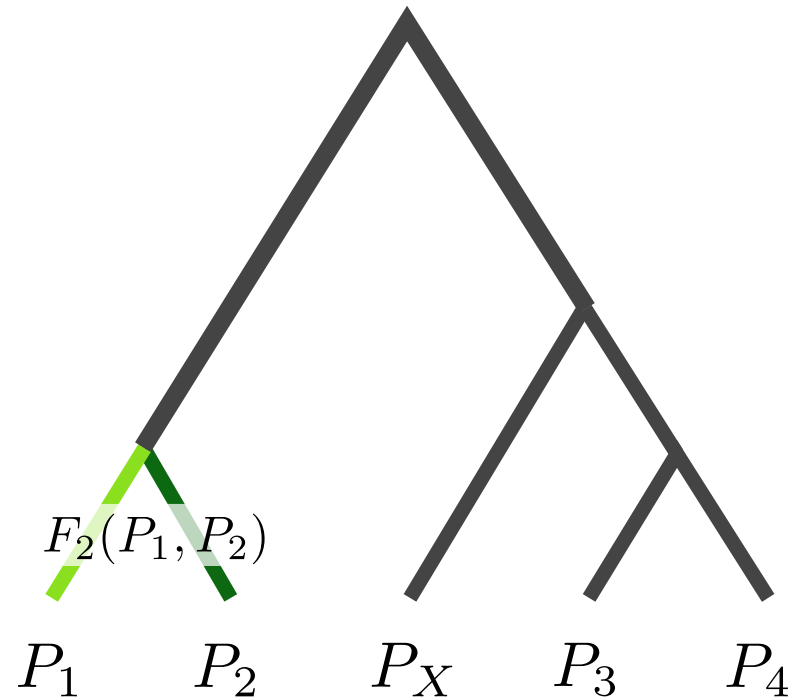
$$F_{ST}(P_1, P_2) = \frac{\pi_{12} - \frac{\pi_1 + \pi_2}{2}}{\pi_{12}}$$

$$F_2(P_1, P_2) = 2\pi_{12} - \pi_1 - \pi_2$$
$$= \sum_l (p_{1l} - p_{2l})^2$$

- $F_{ST}$ is a correlation coefficient
- Between 0 and 1
- Hierarchical partitioning (AMOVA)
- Many estimators exist
  - Hudson (1991)
  - Weir & Cockerham (1984)

- $F_2$ is a covariance
- Bigger than 0
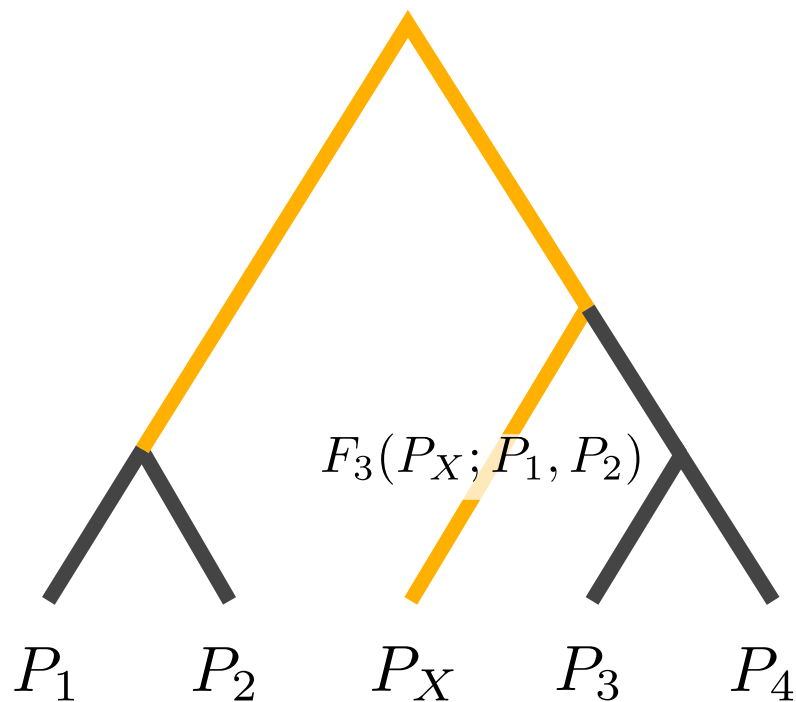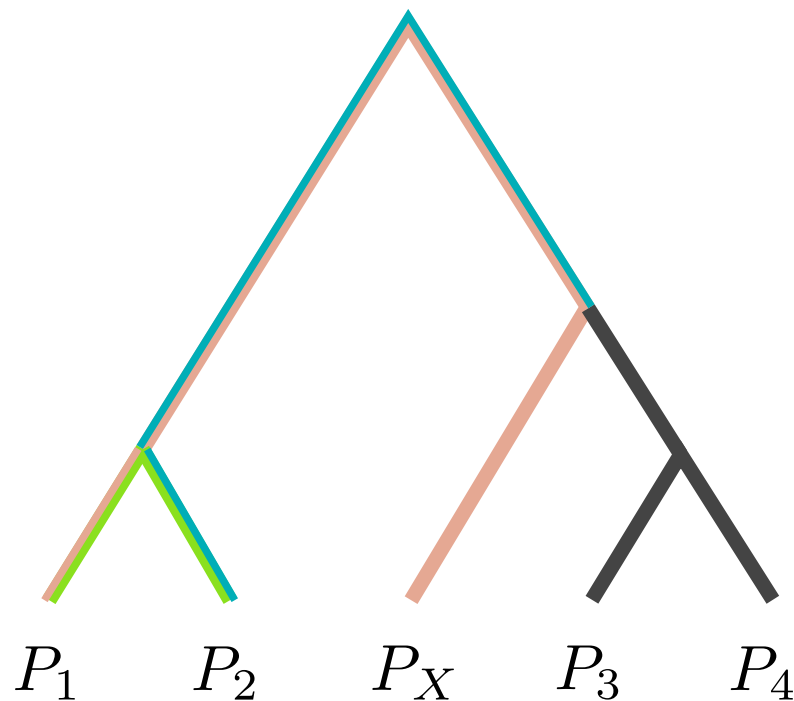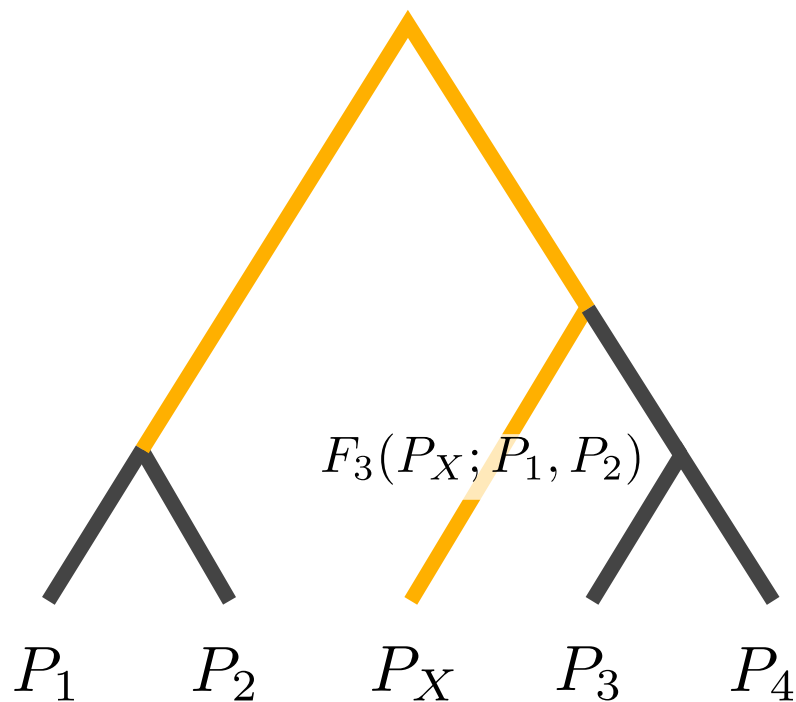- Tree-additive
- Testing for treeness

Tree-additive

$F_2(P_1, P_2)$

$P_1$　$P_2$　$P_X$　$P_3$　$P_4$

# F$_3$-statistic

Given all F2-values, how can we calculate the yellow branch length?



$$F_3(P_X; P_1, P_2)$$

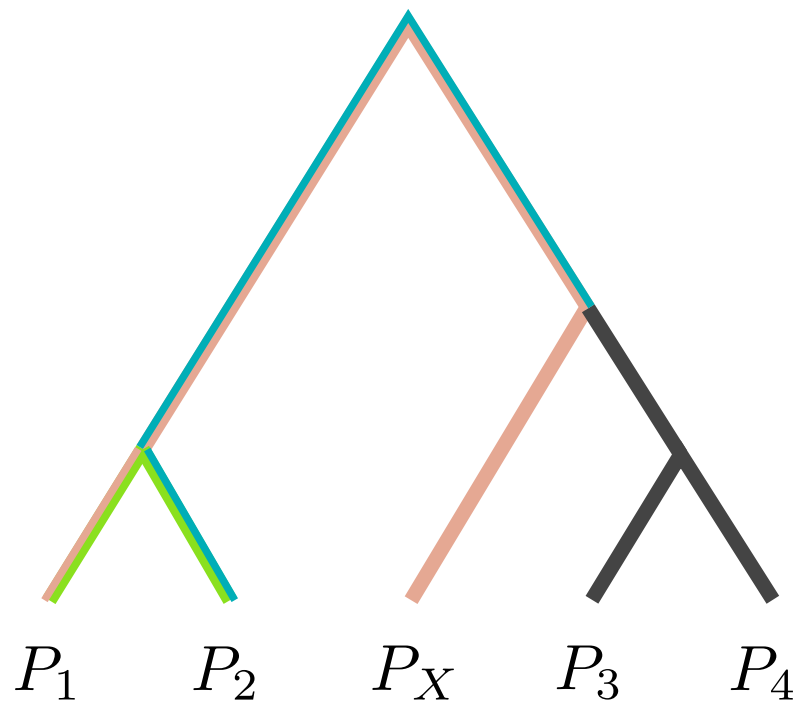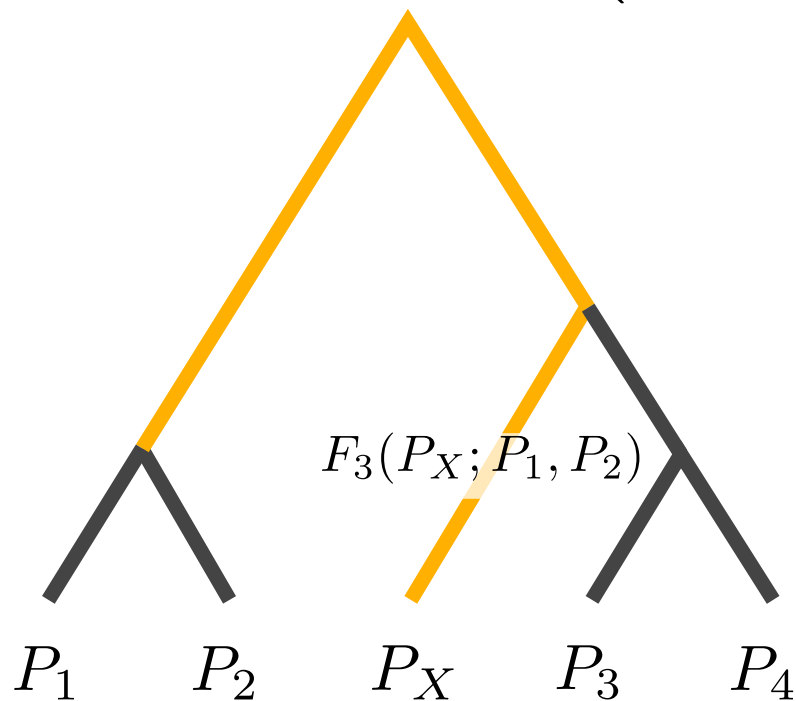$P_1 \quad P_2 \quad P_X \quad P_3 \quad P_4$

# F₃-statistic

Given all F2-values, how can we calculate the yellow branch length?



$F_3(P_X; P_1, P_2)$

# F₃-statistic

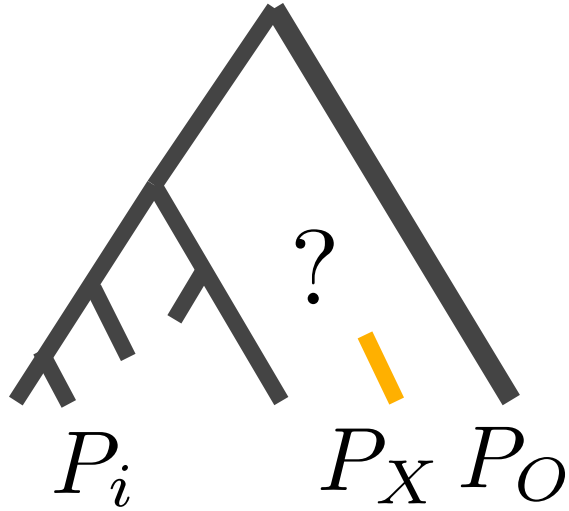$$F_3(P_X; P_1, P_2) = \frac{1}{2}\Big(F_2(P_X, P_1) + F_2(P_X, P_2) - F_2(P_1, P_2)\Big)$$

# $F_3$ -statistic equations

$$F_3(P_X; P_1, P_2) = \frac{1}{2}\left( F_2(P_X, P_1) + F_2(P_X, P_2) - F_2(P_1, P_2) \right)$$
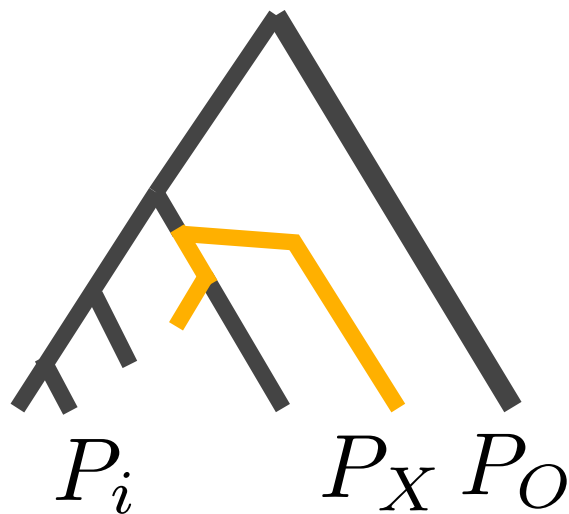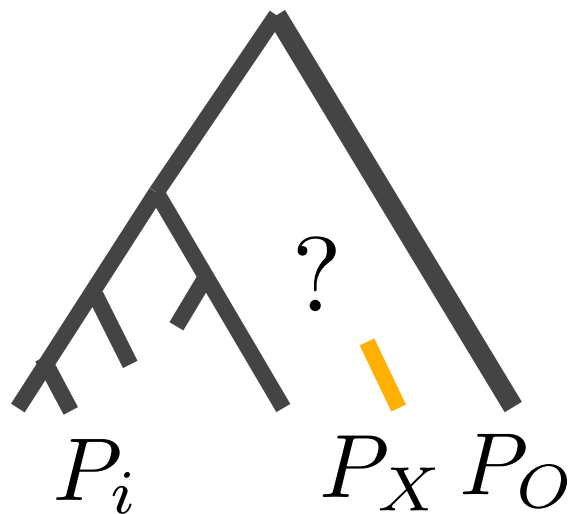
$$F_3(P_X; P_1, P_2) = \sum_l (p_{xl} - p_{x1})(p_{xl} - p_{x2})$$

$$F_3(P_X; P_1, P_2) = \pi_{1x} + \pi_{2x} - \pi_{12} - \pi_x$$
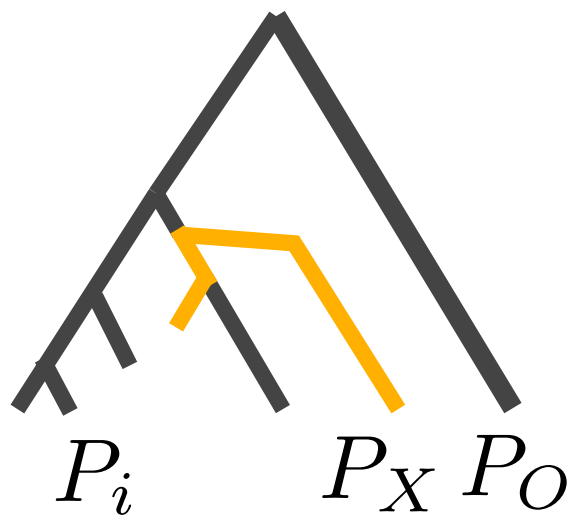
# Outgroup-F₃-statistic
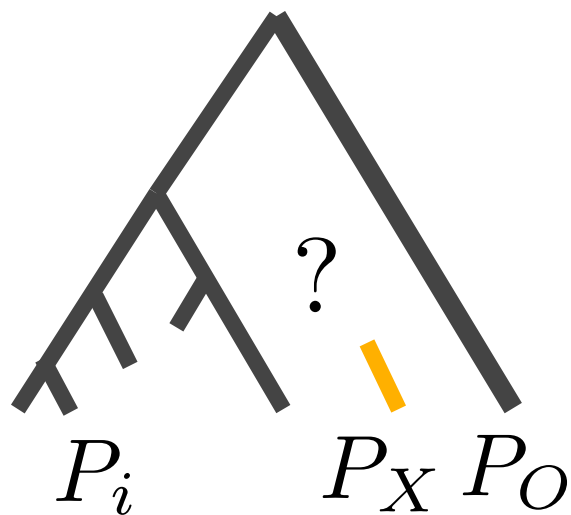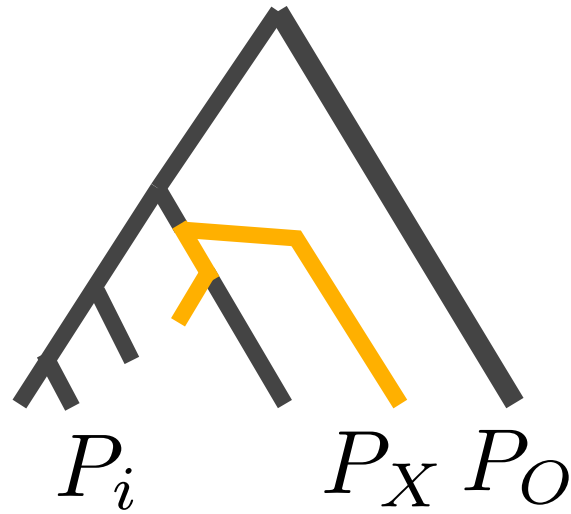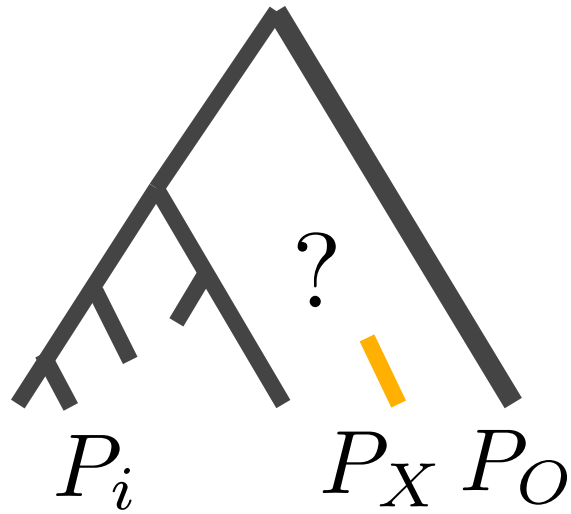
# Outgroup-F₃-statistic

# Outgroup-F₃-statistic

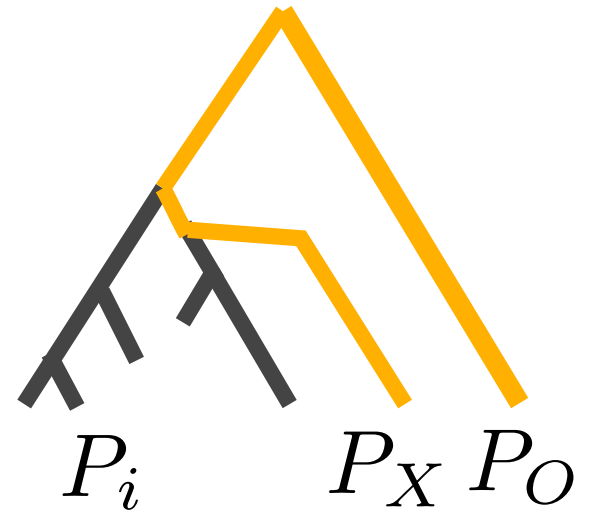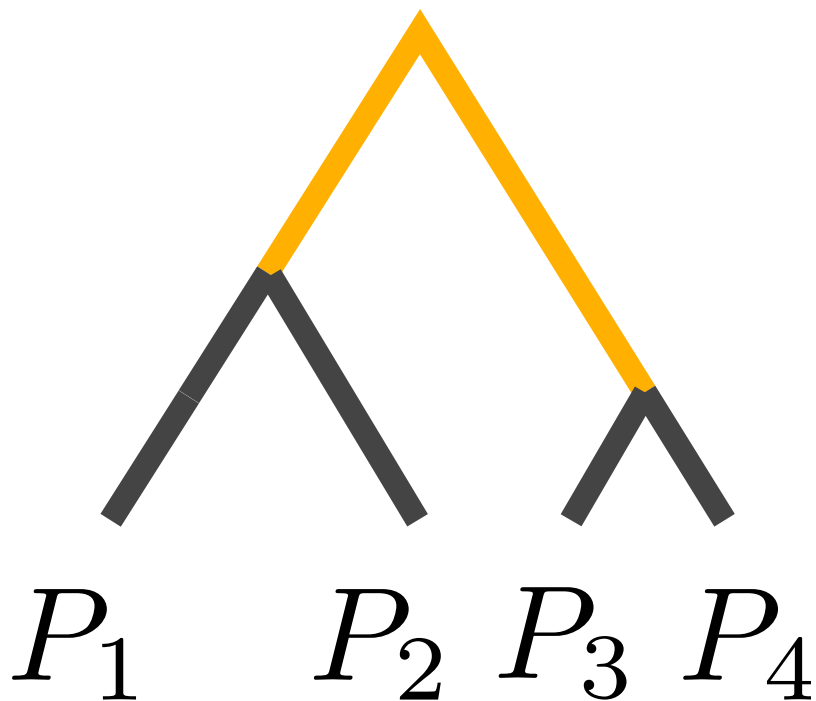

$$F_2(P_X, P_i)$$

# Outgroup-F₃-statistic



$$F_2(P_X, P_i)$$

$$F_3(P_O; P_X, P_i)$$

# (Branch)-F$_4$-statistic

$$F_4^{(B)}(P_1, P_2; P_3, P_4) = \frac{1}{2}\Big(F_2(P_1, P_3) + F_2(P_2, P_4) - F_2(P_1, P_2) - F_2(P_3, P_4)\Big)$$
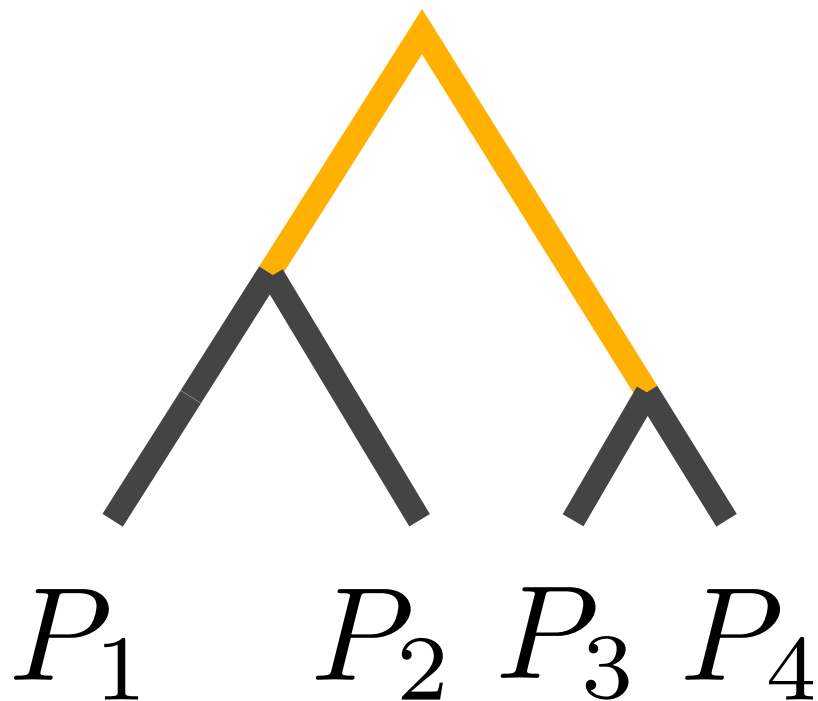


$P_1 \qquad P_2 \quad P_3 \quad P_4$

# (Branch)-F$_4$-statistic

$$F_4^{(B)}(P_1, P_2; P_3, P_4) = \frac{1}{2}\left( F_2(P_1, P_3) + F_2(P_2, P_4) - F_2(P_1, P_2) - F_2(P_3, P_4) \right)$$
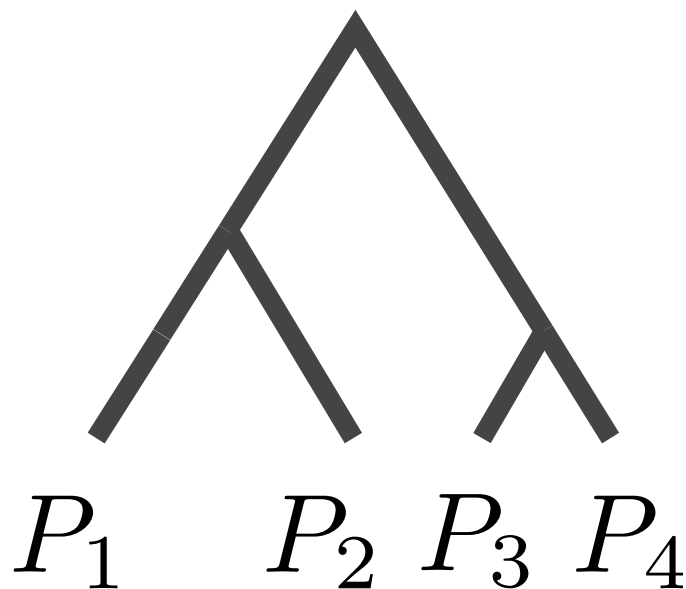


What if we reorder the arguments?

$$F_4^{(T)}(P_1, P_2; P_3, P_4) = F_4^{(B)}(P_1, P_4; P_3, P_2)$$

$P_1 \qquad P_2 \quad P_3 \quad P_4$

# (Treeness)-F$_4$-statistic

$$F_4^{(T)}(P_1, P_2; P_3, P_4) = \frac{1}{2}\Big( F_2(P_1, P_3) + F_2(P_2, P_4) - F_2(P_1, P_4) - F_2(P_2, P_3)\Big)$$



$P_1 \quad P_2 \; P_3 \; P_4$

# F$_4$-statistic-equations

$$F_4^{(T)}(P_1, P_2; P_3, P_4) = \frac{1}{2}\bigg(F_2(P_1, P_3) + F_2(P_2, P_4) - F_2(P_1, P_4) - F_2(P_2, P_3)\bigg)$$

$$F_4^{(T)}(P_1, P_2; P_3, P_4) = \sum_l (p_{l1} - p_{l2})(p_{l3} - p_{l4})$$
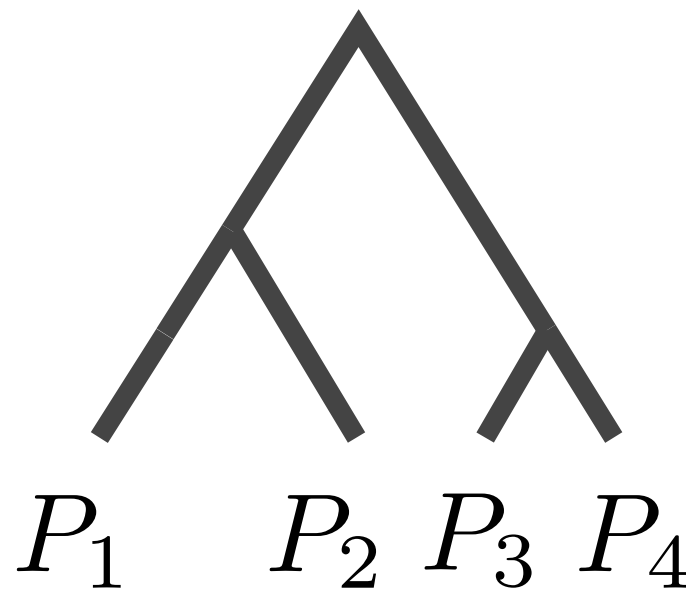
$$F_4^{(T)}(P_1, P_2; P_3, P_4) = \pi_{13} + \pi_{24} - \pi_{14} - \pi_{23}$$
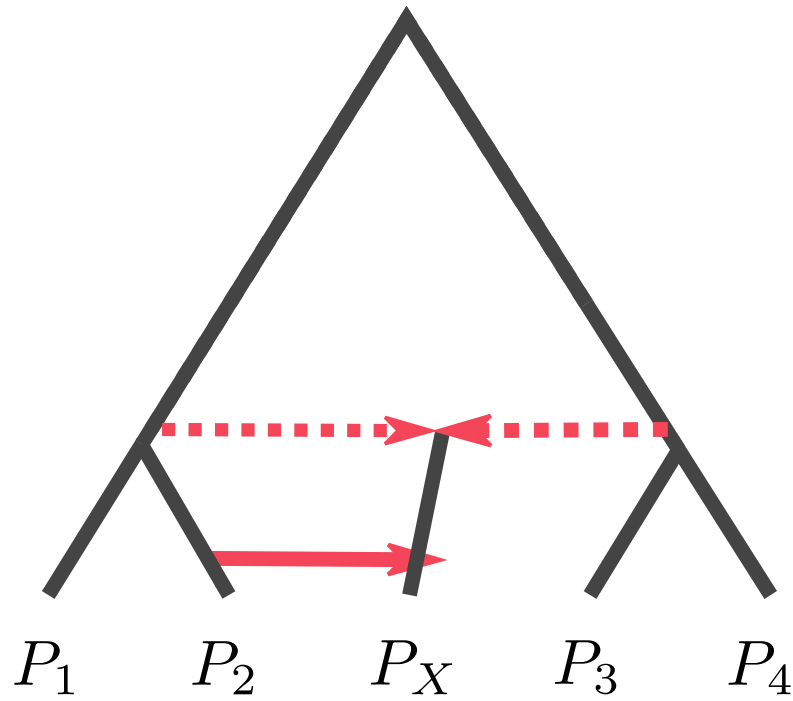
# Testing Treeness

If data is generated from a tree:

$$F_3(P_3; P_1; P_2) \geq 0$$

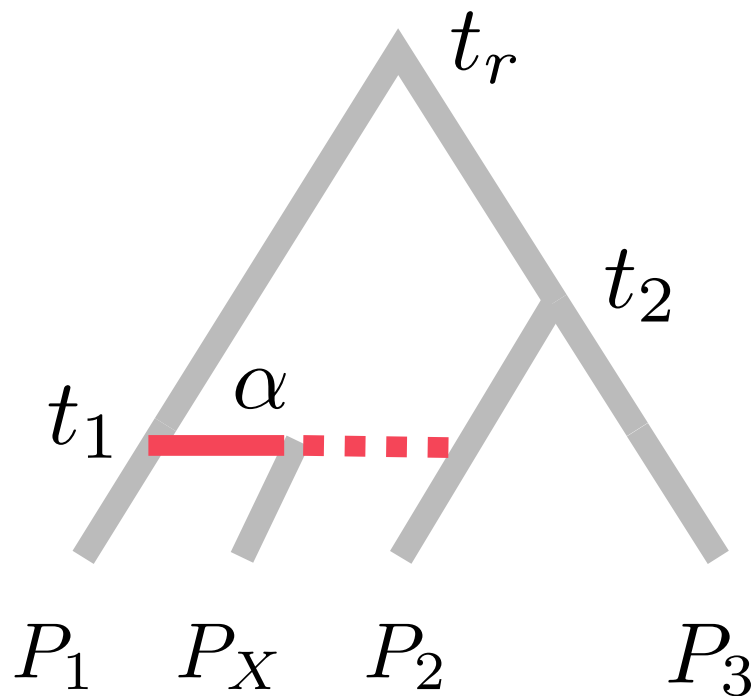$$F_4^{(T)}(P_1, P_2; P_3, P_4) = 0$$



Buneman 1974

# Admixture Graphs

# F3 in an admixture graph

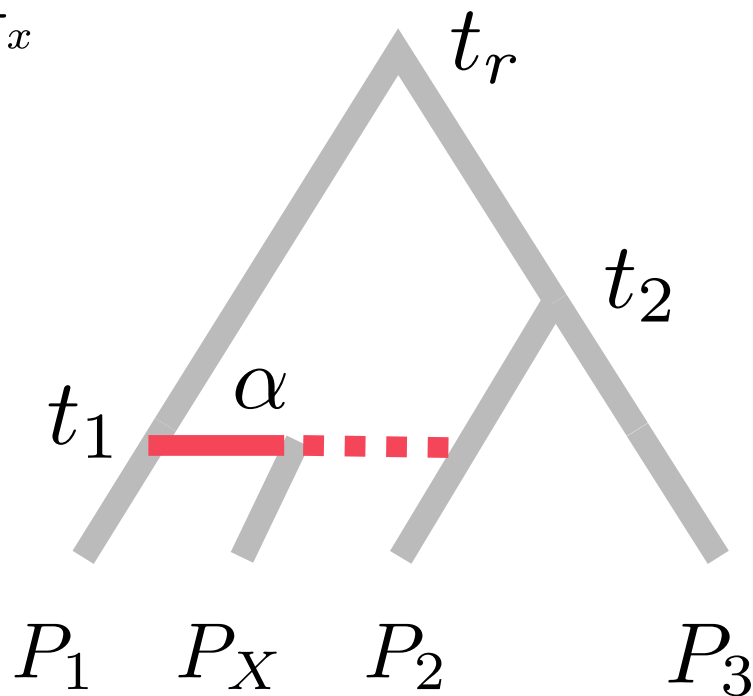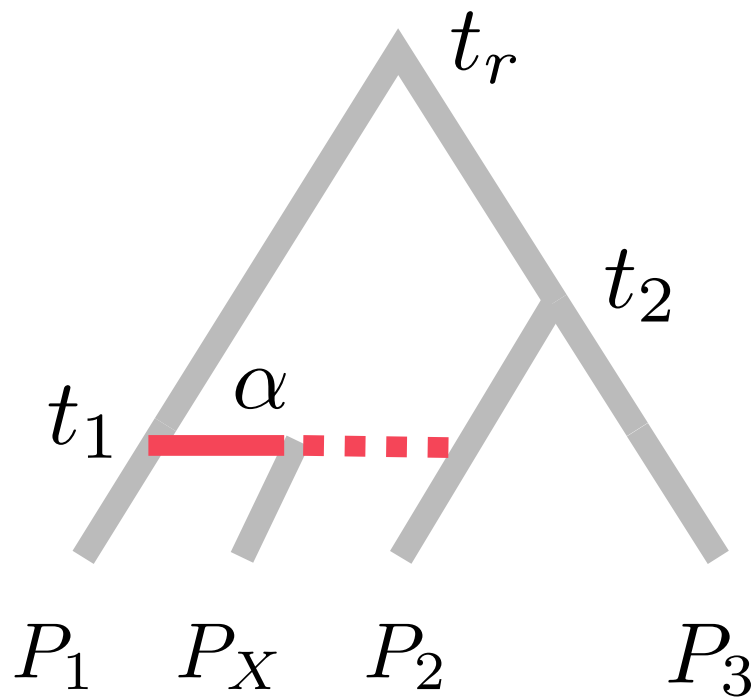$$F_3(P_X; P_1, P_2) \approx \theta\big[t_1 - 2\alpha(1-\alpha)t_r\big]$$

# F3 in an admixture graph

$$F_3(P_X; P_1, P_2) = \pi_{1x} + \pi_{2x} - \pi_{12} - \pi_x$$

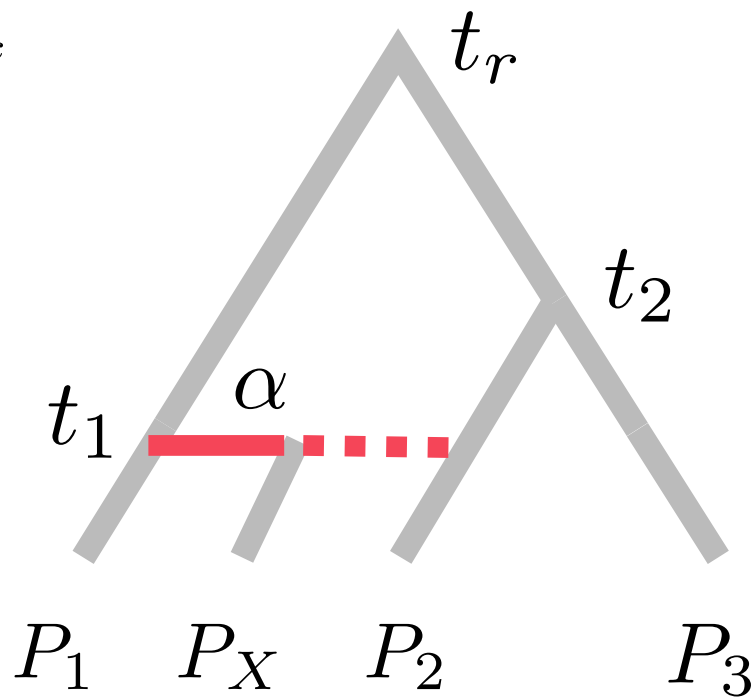$$F_3(P_X; P_1, P_2) \approx \theta\big[t_1 - 2\alpha(1-\alpha)t_r\big]$$

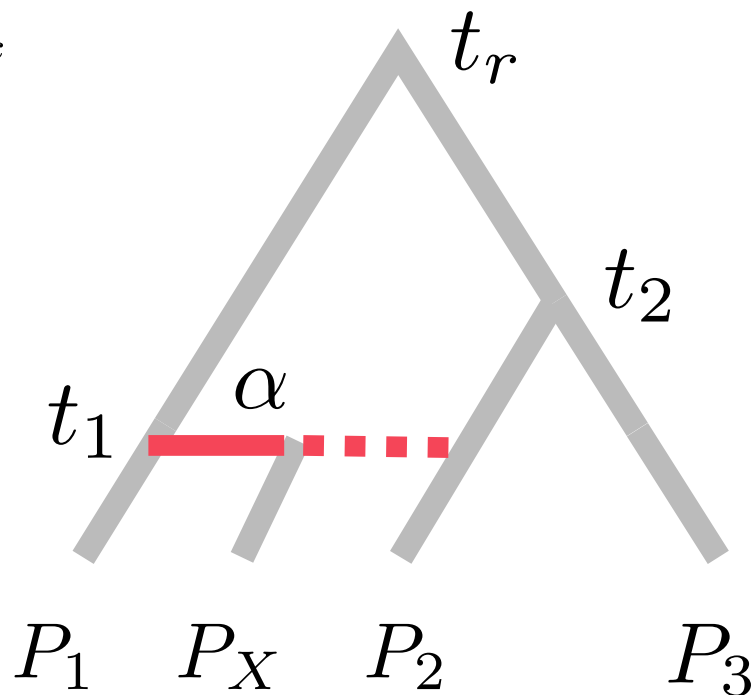# F4 in an admixture graph

# F4 in an admixture graph

$$F_4^{(T)}(P_1, P_X; P_2, P_3) = \pi_{12} + \pi_{3x} - \pi_{13} - \pi_{2x}$$

# F4 in an admixture graph

$$F_4^{(T)}(P_1, P_X; P_2, P_3) = \pi_{12} + \pi_{3x} - \pi_{13} - \pi_{2x}$$

$$F_4^{(T)}(P_1, P_X; P_2, P_3) = (1 - \alpha)(t_2 - t_1) \neq 0$$
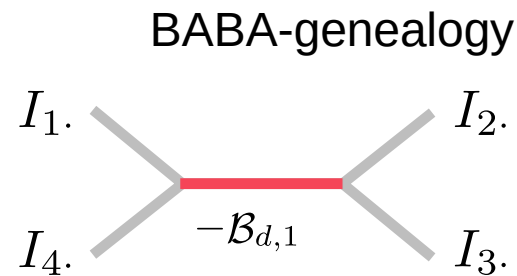


$t_r$

$t_2$

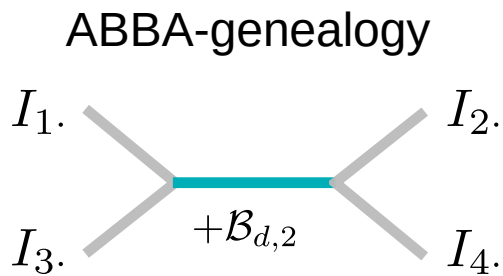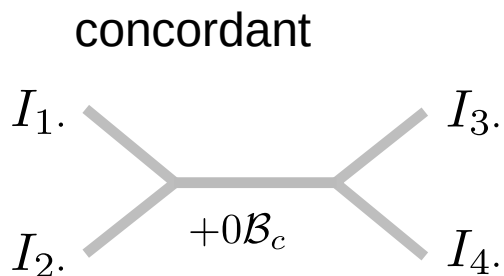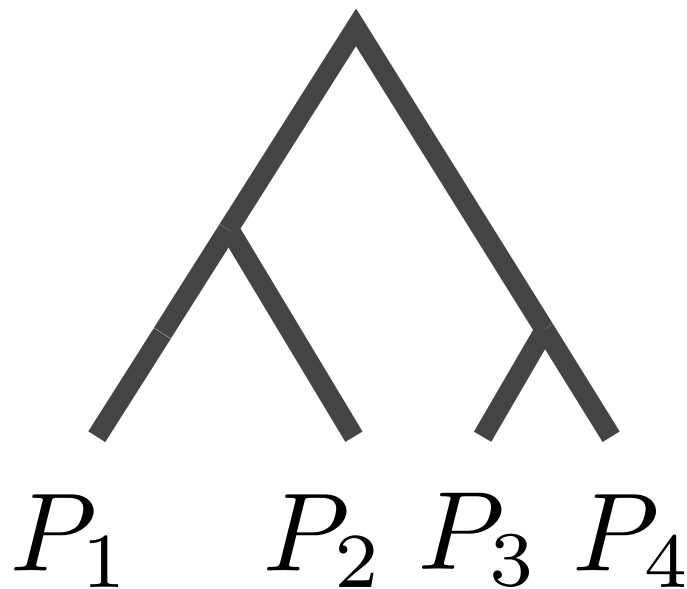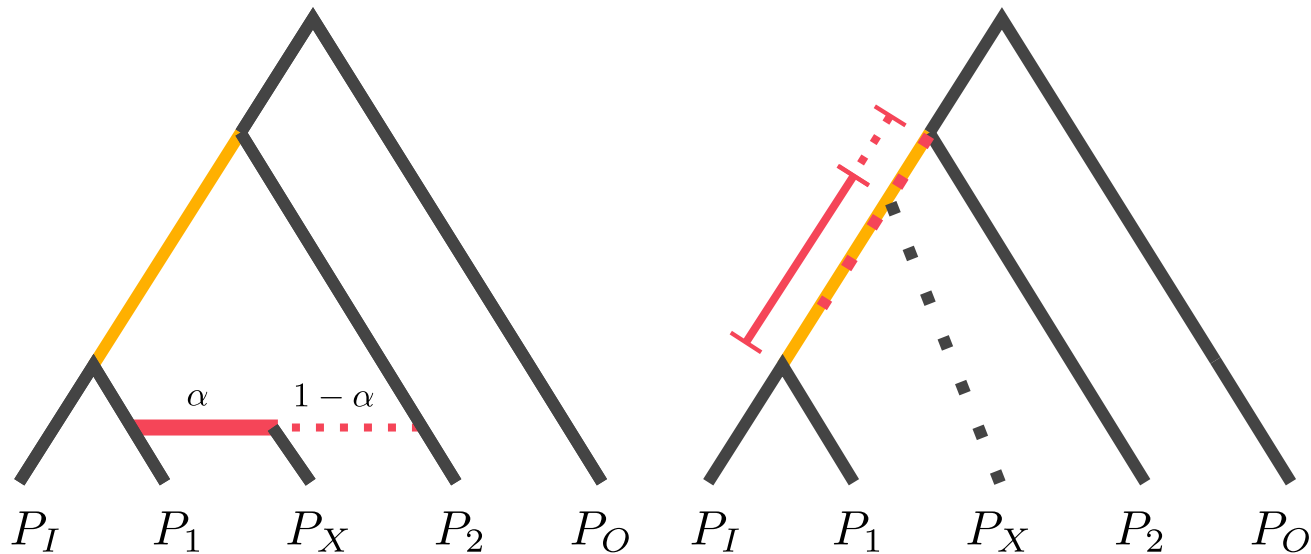$t_1$  $\alpha$

$P_1$  $P_X$  $P_2$  $P_3$

# D-statistic

$$D = \frac{\text{ABBA} - \text{BABA}}{\text{BABA} + \text{ABBA}}$$

- D-statistic and F4 are closely related

$$F_4^{(T)}(P_1, P_2; P_3, P_4) = \pi_{13} + \pi_{24} - \pi_{14} - \pi_{23}$$
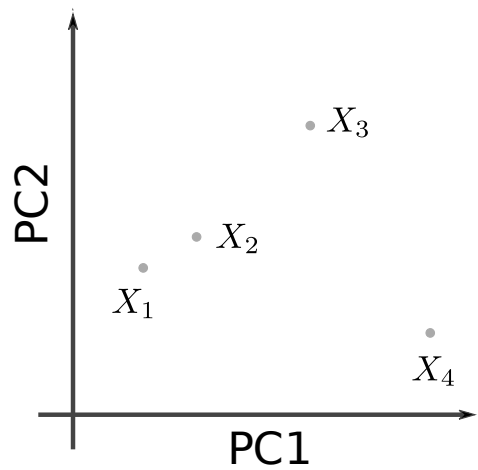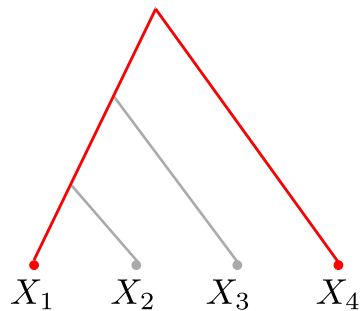


concordant

$I_1.$   $I_3.$

$+0\mathcal{B}_c$

$I_2.$   $I_4.$

ABBA-genealogy

$I_1.$   $I_2.$

$+\mathcal{B}_{d,2}$

$I_3.$   $I_4.$

BABA-genealogy
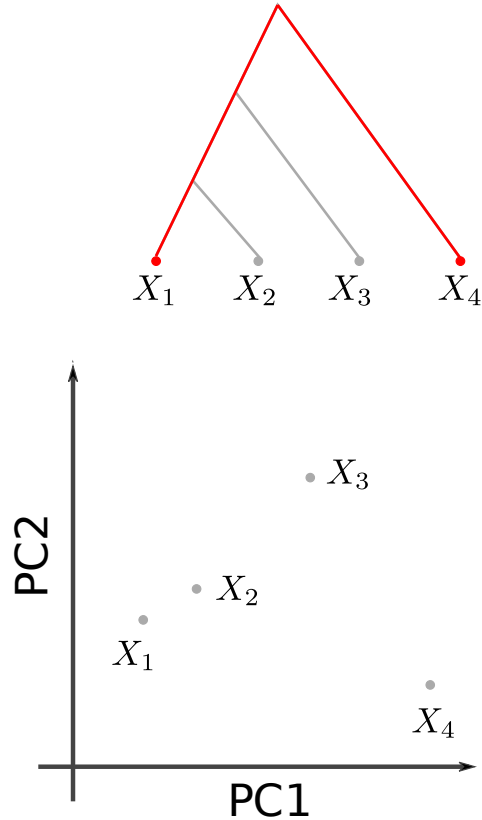
$I_1.$   $I_2.$

$-\mathcal{B}_{d,1}$

$I_4.$   $I_3.$

# F4-ratio



$$\alpha = 1 - \frac{F_4^{(B)}(P_I, P_1; P_X, P_O)}{F_4^{(B)}(P_I, P_1; P_2, P_O)}$$

**A** $F_2(X_1; X_4)$

A $F_2(X_1; X_4)$

# A geometric relationship of $F_2$, $F_3$ and $F_4$-statistics with principal component analysis
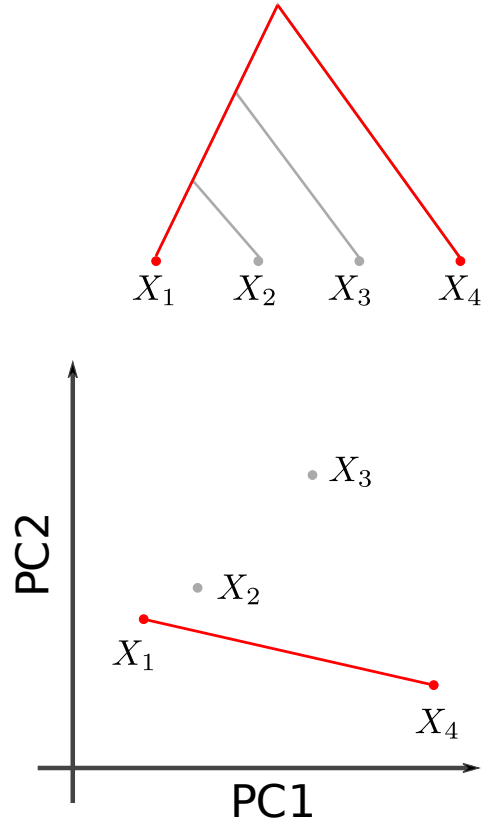
Benjamin M. Peter

Max-Planck-Institute for Evolutionary Anthropology, Leipzig 04103, Germany

BMP, 0000-0003-2526-8081

Principal component analysis (PCA) and *F*-statistics *sensu* Patterson are two of the most widely used population genetic tools to study human genetic variation. Here, I derive explicit connections between the two approaches and show that these two methods are closely related. *F*-statistics have a simple geometrical interpretation in the context of PCA, and orthogonal projections are a key concept to establish this link. I show that for any pair of populations, any population that is admixed as determined by an $F_3$-statistic will lie inside a circle on a PCA plot. Furthermore, the $F_4$-statistic is closely related to an angle measurement, and will be zero if the differences between pairs of populations intersect at a right angle in PCA space. I illustrate my results on two examples, one of Western Eurasian, and one of global human diversity. In both examples, I find that the first few PCs are sufficient to approximate most *F*-statistics, and that PCA plots are effective at predict

A $F_2(X_1; X_4)$

# A geometric relationship of $F_2$, $F_3$ and $F_4$-statistics with principal component analysis
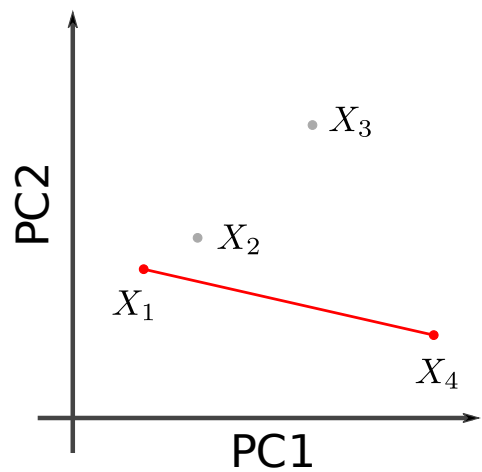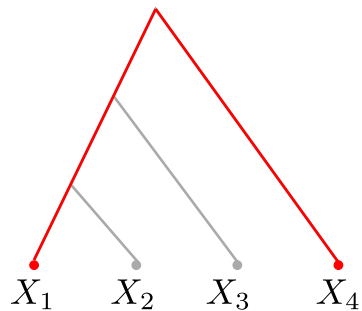
Benjamin M. Peter

Max-Planck-Institute for Evolutionary Anthropology, Leipzig 04103, Germany

BMP, 0000-0003-2526-8081

Principal component analysis (PCA) and $F$-statistics *sensu* Patterson are two of the most widely used population genetic tools to study human genetic variation. Here, I derive explicit connections between the two approaches and show that these two methods are closely related. $F$-statistics have a simple geometrical interpretation in the context of PCA, and orthogonal projections are a key concept to establish this link. I show that for any pair of populations, any population that is admixed as determined by an $F_3$-statistic will lie inside a circle on a PCA plot. Furthermore, the $F_4$-statistic is closely related to an angle measurement, and will be zero if the differences between pairs of populations intersect at a right angle in PCA space. I illustrate my results on two examples, one of Western Eurasian, and one of global human diversity. In both examples, I find that the first few PCs are sufficient to approximate most $F$-statistics, and that PCA plots are effective at predict
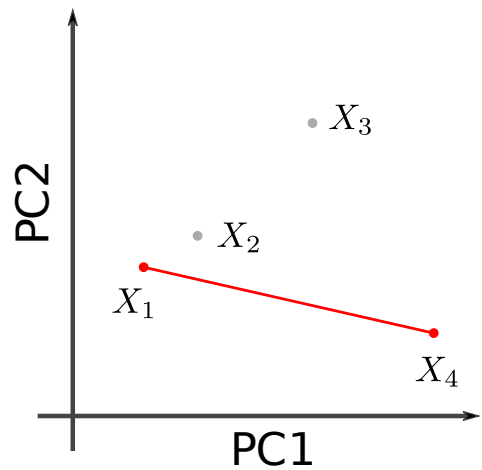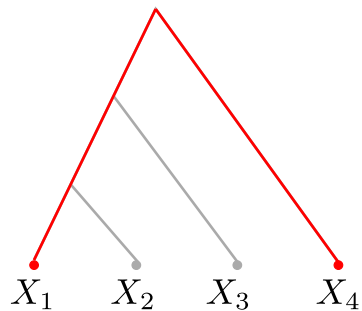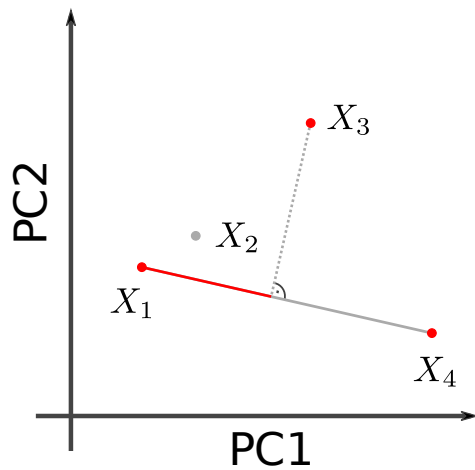
* true if all PCs are used, optimal approximation otherwise

Peter (2022): doi 10.1098/rstb.2020.0413

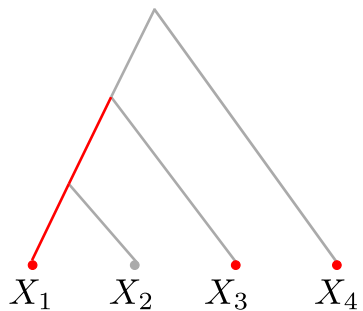A $F_2(X_1; X_4)$

A $F_2(X_1; X_4)$  B $F_3(X_1; X_3, X_4)$

* true if all PCs are used, optimal approximation otherwise

Peter (2022): doi 10.1098/rstb.2020.0413

A $F_2(X_1; X_4)$

B $F_3(X_1; X_3, X_4)$

C $F_4(X_1, X_4; X_2, X_3)$

* true if all PCs are used, optimal approximation otherwise
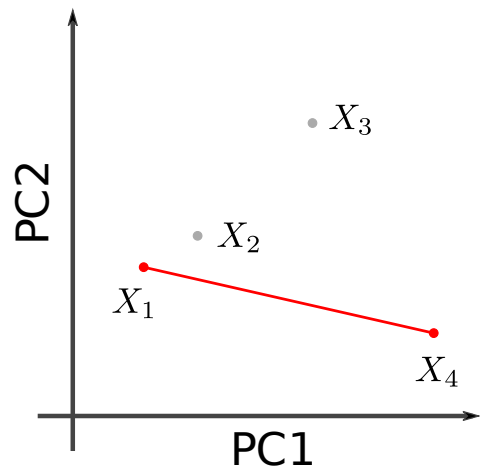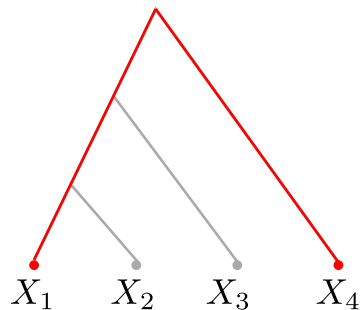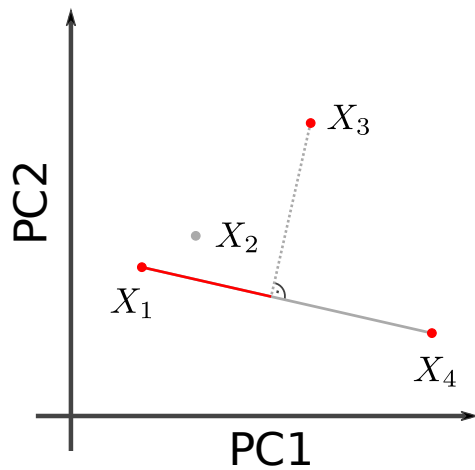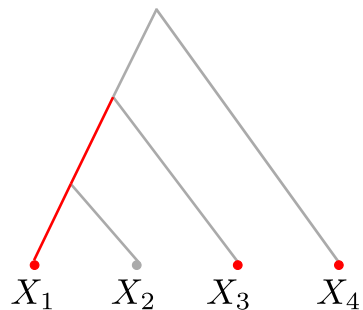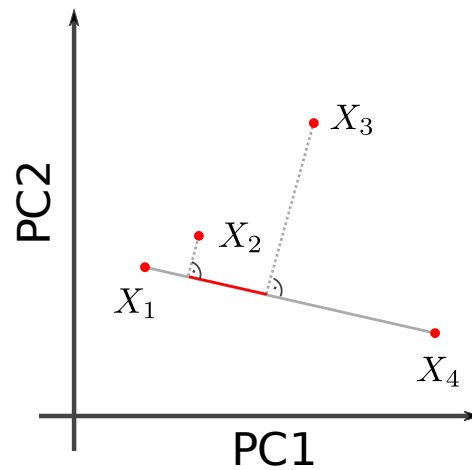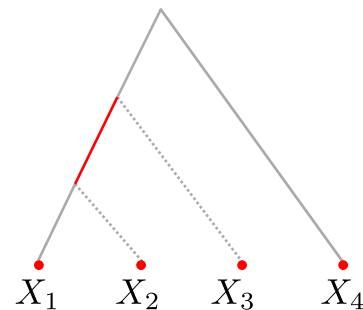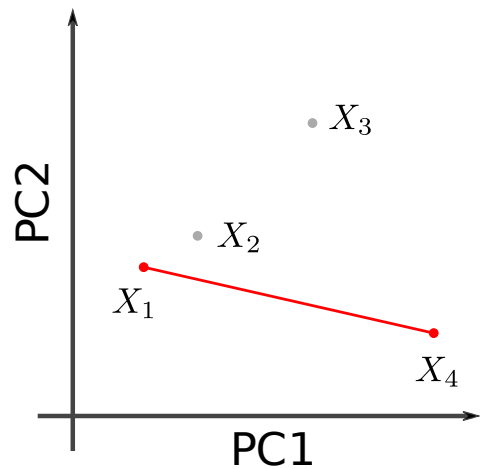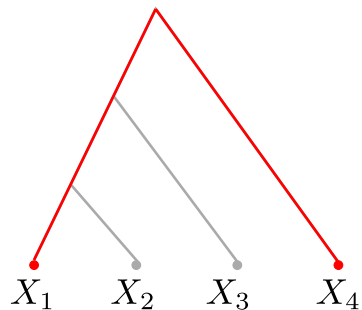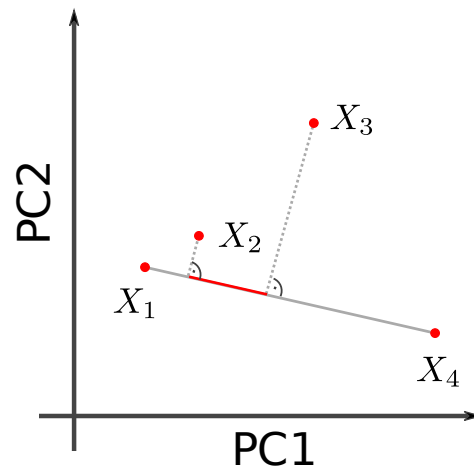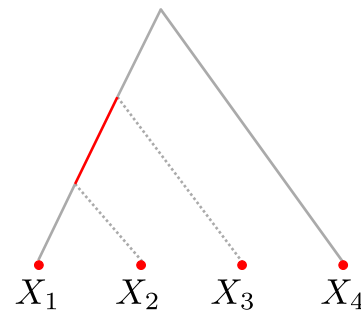
A $F_2(X_1; X_4)$

B $F_3(X_1; X_3, X_4)$

C $F_4(X_1, X_4; X_2, X_3)$

D $F_4(X_1, X_2; X_3, X_4)$

* true if all PCs are used, optimal approximation otherwise