# Demographic inference using summary statistics

Ben Peter

Max Planck Institute for Evolutionary Anthropology
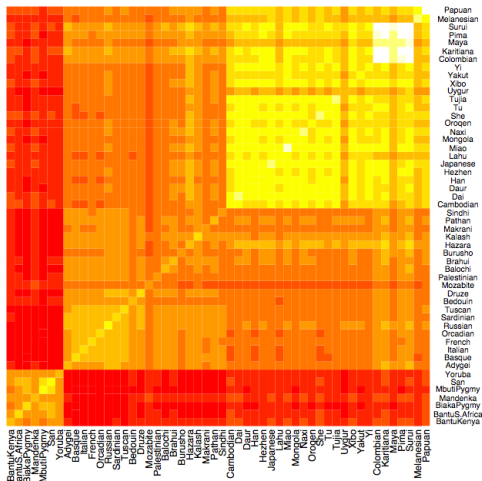
Leipzig, Germany

# Motivation

# Motivation



- Archaeological and linguistic sources of data give us alternative sources of data with which to confirm/contrast genetic inferences regarding population history
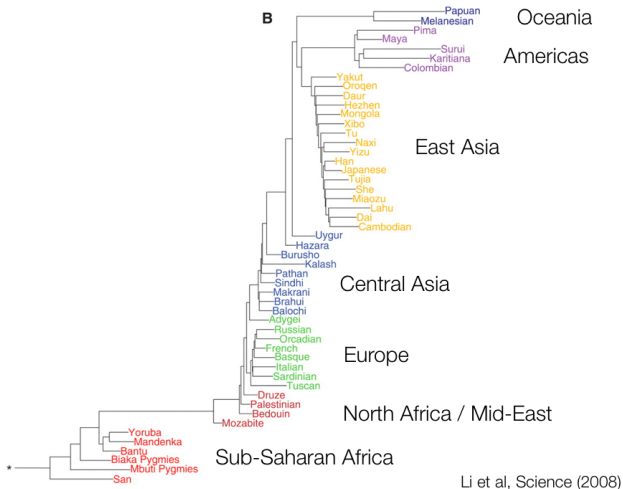- Observational studies only possible - so statistical methods are key for inference

# Similarity matrices

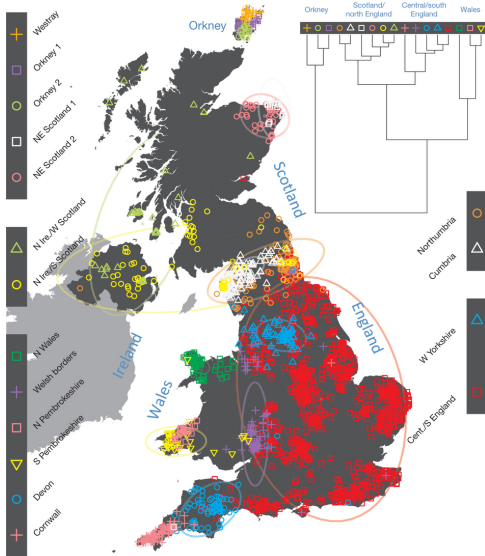Covariance matrix of allele frequencies across HGDP populations



Coop et al (2010) Genetics

# Phylogenetic trees

Neighbor-joining tree built with `PHYLIP` on the basis of similarity in allele frequencies:
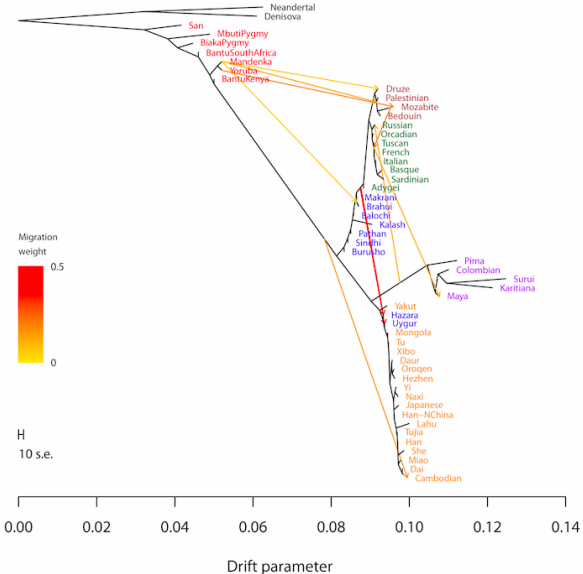


Li et al, Science (2008)

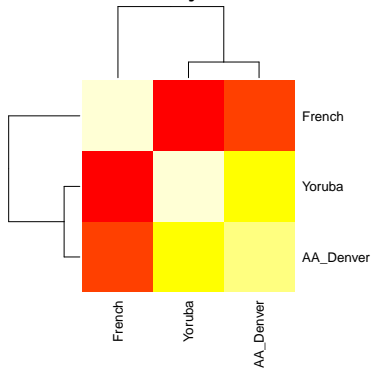# Finestructure algorithm (Leslie et al. 2015)

# Phylogenetic trees

Population tree with admixture events inferred using `TreeMix` software on the basis of allele frequencies:
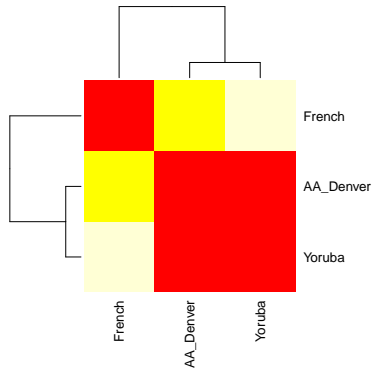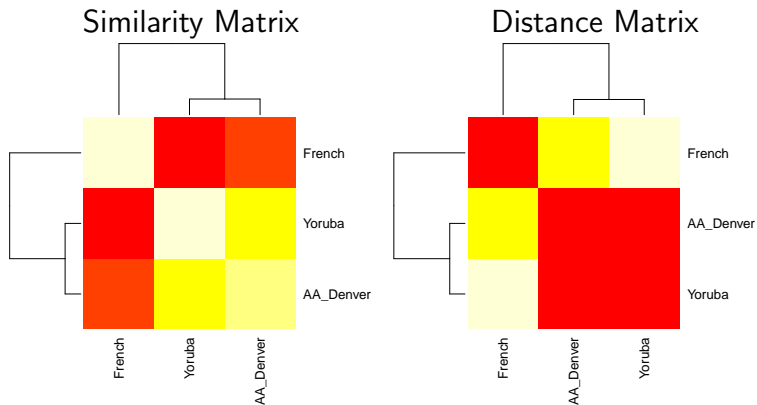
# Measuring Similarity vs measuring distance

# Measuring Similarity vs measuring distance



Similarity Matrix

Distance Matrix

Difference: $D_{i,i} = 0$

# How could we measure genetic similarity/dissimilarity in a population?

# How could we measure genetic similarity/dissimilarity in a population?



- change in allele frequency
- loss of heterozygosity
- probability of coalescence

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

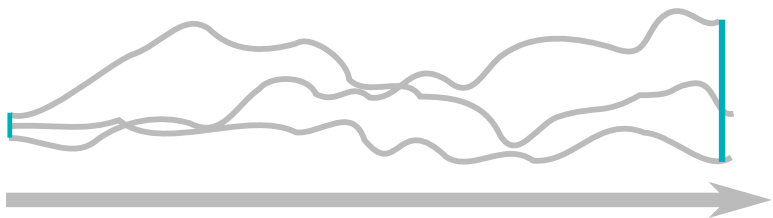# Introducing today's superhero

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

- change in allele frequency
- loss of heterozygosity
- probability of coalescence

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

C $\qquad F_2 = \frac{1}{2} f \, \mathbb{E} H_0$
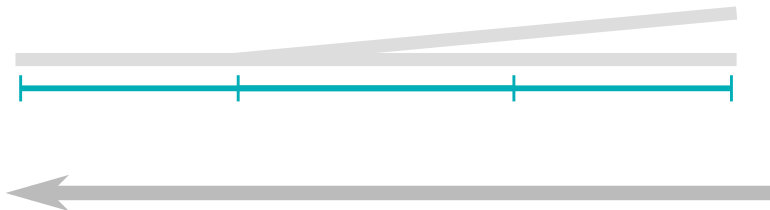
C

$$F_2 = \tfrac{1}{2} f\, \mathbb{E} H_0$$

# How could we measure genetic similarity/dissimilarity between populations?



sample 1

space

sample 2

# How could we measure genetic similarity/dissimilarity between populations?



sample 1

# space

sample 2

- difference in allele frequency
- Heterozygosity: $H_{\text{between}}$ vs $H_{\text{within}}$
- Coalescence: $T_{\text{between}}$ vs $T_{\text{within}}$

# How could we measure genetic similarity/dissimilarity between populations?



sample 1      space      sample 2

- difference in allele frequency
- Heterozygosity: $H_{\text{between}}$ vs $H_{\text{within}}$
- Coalescence: $T_{\text{between}}$ vs $T_{\text{within}}$

Conveniently, $F_2$, measures difference equivalently in this scenario

# From differences to trees



sample 1     space     sample 2

- difference in allele frequency
- Heterozygosity: $H_{\text{between}}$ vs $H_{\text{within}}$
- Coalescence: $T_{\text{between}}$ vs $T_{\text{within}}$

Conveniently, $F_2$, measures difference equivalently in this scenario

$$F_2(P_1, P_2) = 2\mathbb{E}\,T_{12} - \mathbb{E}\,T_{11} - \mathbb{E}\,T_{12}$$

$$F_{ST}(P_1, P_2) = \frac{2F_2(P_1, P_2)}{\mathbb{E}H}$$

# $F_2$ vs $F_{ST}$

$$F_{ST}(P_1, P_2) = \frac{2F_2(P_1, P_2)}{\mathbb{E}H}$$

Main difference is normalization:

- $F_{ST} = 0$ : no differentiation
- $F_{ST} = 1$ : maximum differentiation

$F_2$ vs $F_{ST}$

$$F_{ST}(P_1, P_2) = \frac{2F_2(P_1, P_2)}{\mathbb{E}H}$$

Main difference is normalization:

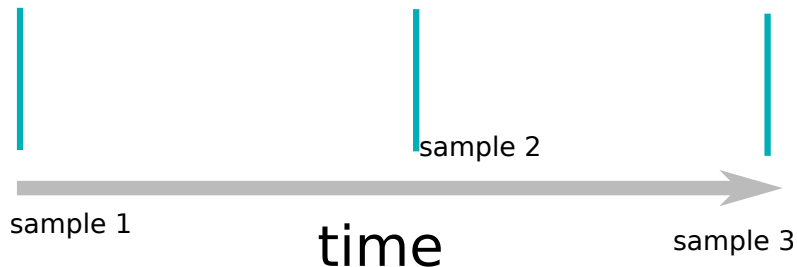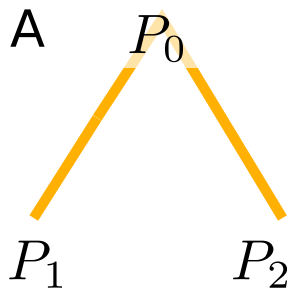- $F_{ST} = 0$ : no differentiation
- $F_{ST} = 1$ : maximum differentiation
- $F_2 = 0$ : no differentiation
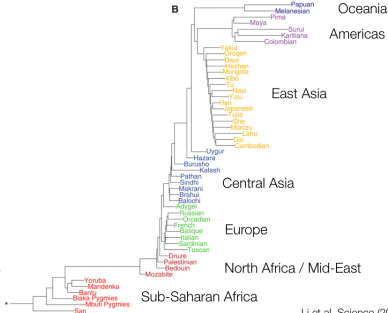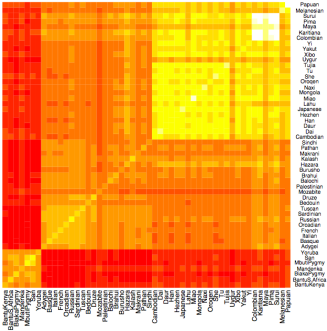- $F_2 = ???$ : maximum differentiation

# $F_2$ is additive



sample 2

sample 1

# time

sample 3

$$F_2(P_1, P_3) = F_2(P_1, P_2) + F_2(P_2, P_3)$$

$$F_2(P_1, P_2) = F_2(P_0, P_1) + F_2(P_0, P_2)$$
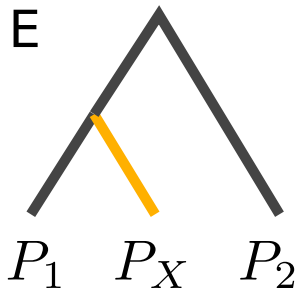
# Dissimilarity matrices vs Tree



Li et al, Science (2008)

$$F_2(P_1, P_2) = F_2(P_0, P_1) + F_2(P_0, P_2)$$

# testing treeness

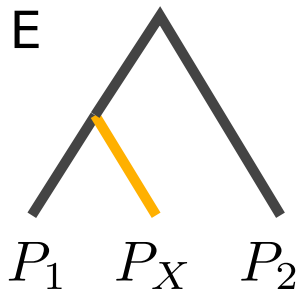$$2F_3(P_X; P_1, P_2) = F_2(P_X, P_1) + F_2(P_X, P_2) - F_2(P_1, P_2)$$

$$2F_3(P_X; P_1, P_2) = F_2(P_X, P_1) + F_2(P_X, P_2) - F_2(P_1, P_2)$$

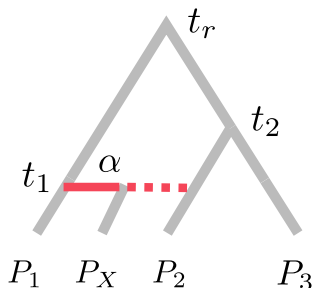$$2F_3(P_X; P_1, P_2) = F_2(P_X, P_1) + F_2(P_X, P_2) - F_2(P_1, P_2)$$



In a tree, $F_3 \geq 0$!

example when this is violated

$$2F_3(P_X; P_1, P_2) = F_2(P_X, P_1) + F_2(P_X, P_2) - F_2(P_1, P_2)$$

# example when this is violated

$$2F_3(P_X; P_1, P_2) = F_2(P_X, P_1) + F_2(P_X, P_2) - F_2(P_1, P_2)$$



$$F_3 = t_1 - 2\alpha(1 - \alpha)(1 - c_x)(t_r - t_1)$$
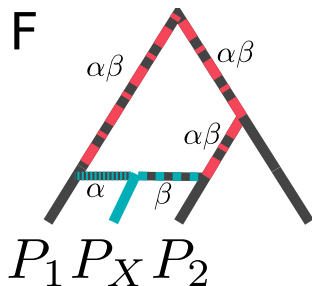
## alternative interpretation

overlap between paths:

$$F_2(P_1, P_2) = \mathbb{E}(p_1 - p_2)(p_1 - p_2)$$
$$F_3(P_X; P_1, P_2) = \mathbb{E}(p_x - p_1)(p_x - p_2)$$
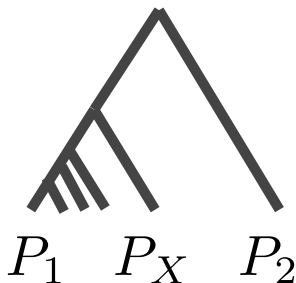
## alternative interpretation

overlap between paths:

$$F_2(P_1, P_2) = \mathbb{E}(p_1 - p_2)(p_1 - p_2)$$
$$F_3(P_X; P_1, P_2) = \mathbb{E}(p_x - p_1)(p_x - p_2)$$

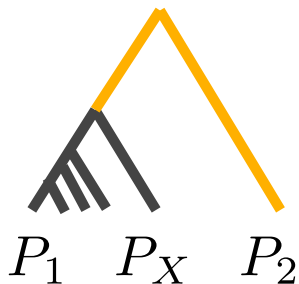Assume we have an unknown sample, and would like to know which potential population $P_1$ it is closest to:



$$P_1 \quad P_X \quad P_2$$

What statistic would you calculate?

Assume we have an unknown sample, and would like to know which potential population $P_1$ it is closest to:



$F_3(P_2; P_X, P_1)$ will be larger the closer $P_X$ and $P_1$ are!
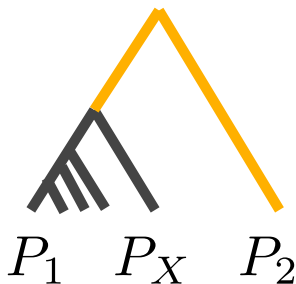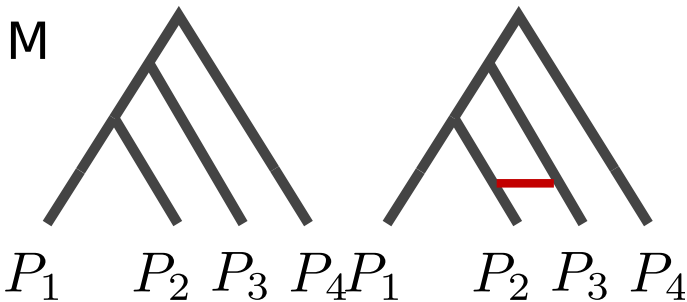
# outgroup-$F_3$

Assume we have an unknown sample, and would like to know which potential population $P_1$ it is closest to:



$F_3(P_2; P_X, P_1)$ will be larger the closer $P_X$ and $P_1$ are!
Advantage over direct measures of differentiation if sampling times of $P_1$ are different.

# $D$-statistic / ($F_4$-statistic)

Imagine you sequence a Neandertal for the first time. How do you test for gene flow?
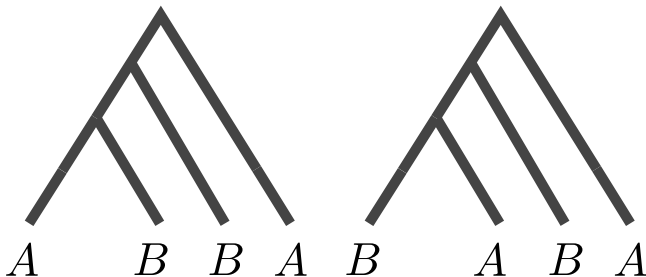
# D-statistic

Imagine you sequence a Neandertal for the first time. How do you test for gene flow?

# D-statistic

Imagine you sequence a Neandertal for the first time. How do you test for gene flow?



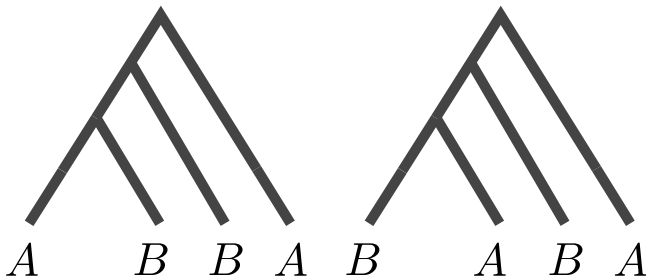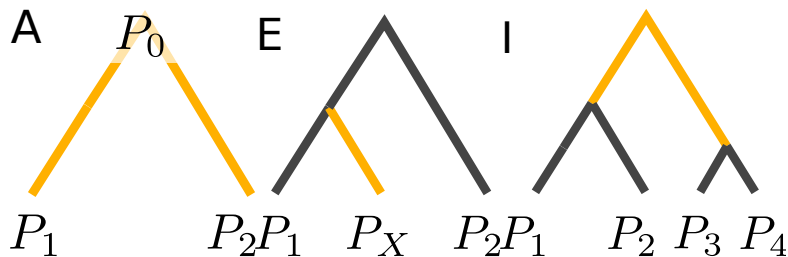$$D = \frac{\text{ABBA} - \text{BABA}}{\text{ABBA} + \text{BABA}}$$

# D-statistic

Imagine you sequence a Neandertal for the first time. How do you test for gene flow?



$A$ $B$ $B$ $A$ $B$ $A$ $B$ $A$

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

# What does $D/F_4$ actually measure?



A $\quad P_0 \quad$ E $\qquad$ I

$P_1 \qquad P_2 P_1 \quad P_X \quad P_2 P_1 \quad P_2 \; P_3 \; P_4$
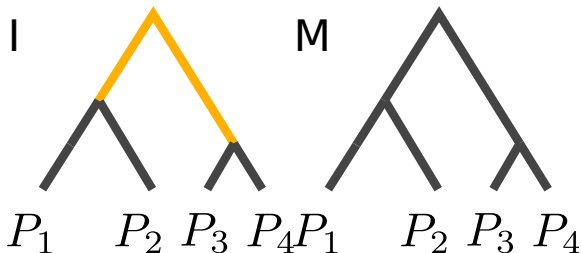
$$D = \frac{\text{ABBA} - \text{BABA}}{\text{ABBA} + \text{BABA}}$$

$$F4 = F_2(P_1, P_3) + F_2(P_2, P_4) - F_2(P_1, P_2) - F_2(P_3, P_4)$$
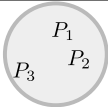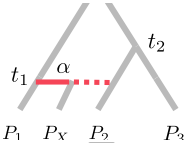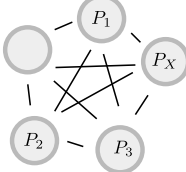
# What does $D/F_4$ actually measure?

Two possibilities:



$$2F_4 = F_2(P_1, P_3) + F_2(P_2, P_4) - F_2(P_1, P_2) - F_2(P_3, P_4)$$

$$2F_4 = F_2(P_1, P_3) + F_2(P_2, P_4) - F_2(P_1, P_4) - F_2(P_2, P_3)$$

# how do these statistics behave under other demographic models?

| Model | | $F_3(P_X; P_1, P_2)$ | $F_4(P_1; P_X; P_2, P_3)$ |
|---|---|---|---|
| Panmictic |  | 0 | 0 |
| Admixture Graph |  | $t_1 - 2\alpha(1-\alpha) \times (1 - c_x)t_r$ | $(1-\alpha)(t_2 - t_1)$ |
| Island Model |  | $\dfrac{1}{M}$ | 0 |

# how do these statistics behave under other demographic models?

| | | | |
|---|---|---|---|
| Stepping stone | $P_1 - P_X - P_2 - P_3$ | $\dfrac{2}{7M}$ | $-\dfrac{8}{7M}$ |
| Hierarchical stepping stone | $P_1\ P_1\ P_X\ P_X\ P_2\ P_2$ | $-\dfrac{\mathbf{0.06}}{\mathbf{M}}$ | $\dfrac{14}{55M}$ |
| Serial founder model | $P_1 \rightarrow P_X \rightarrow P_2 \rightarrow P_3$ | $t_x$ | $0$ |

# Recap

1. $F_3$ and $F_4$ are simple statistics that test for admixture
2. $F_3$ requires just 3 populations, and is most useful for recent admixture at approximately equal proportions
3. $F_4$ is suitable to more ancient admixture, but more sensitive