

# Introduction to population genetics

leo.planche@universite-paris-saclay.fr

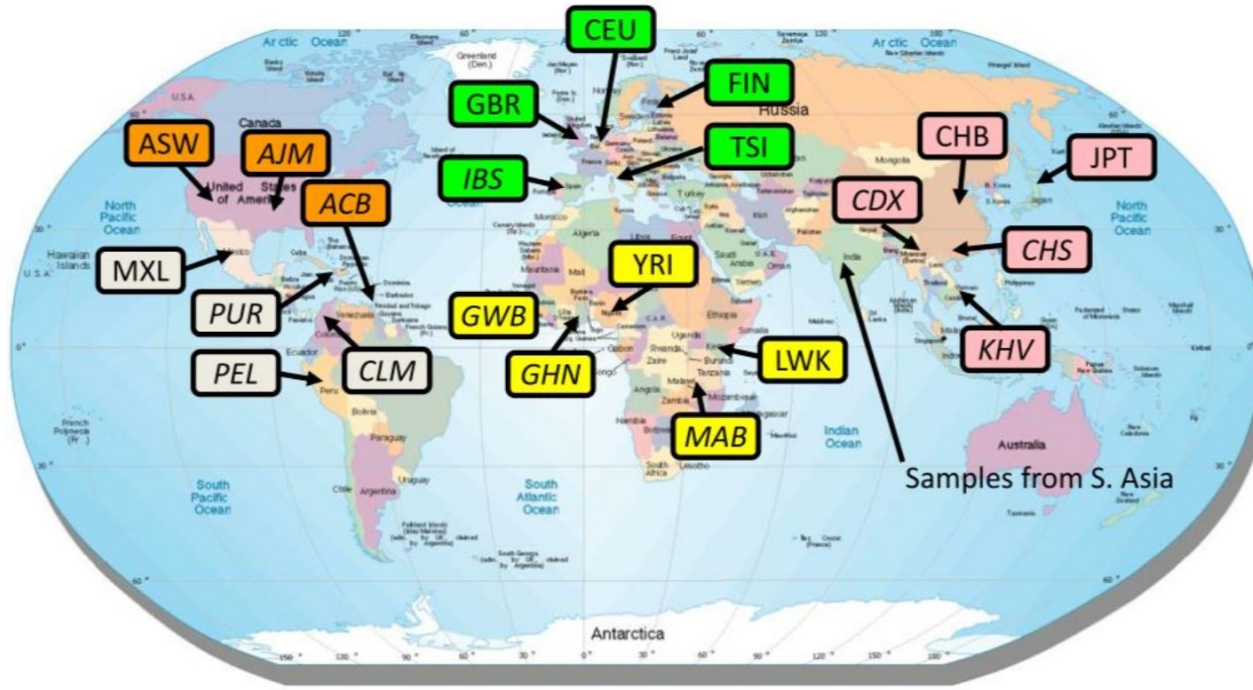
# Today's Overview

- i. What is population genomics
- ii. Variation between populations
- iii. Population size
- iv. Sapiens meets Neandertal and Denisovan

# What is population genetics?

Population genetics is a subfield of genetics that deals with genetic differences within and among populations.

What is a populations? A group of freely interbreeding individuals.



How many labeled populations are displayed on this map?

# Why do population genetics?

- For the love of history (human or non human)
- For medical research
- For research in ecology (conservation genetics)

# Genetic variations

# Mutations

Human mutation rate is around  $1.5 \times 10^{-8}$  per base pair (bp). Meaning that around 30 new variants will be introduced in each gamete. The vast majority of mutations are **neutral**.

# Genetic variations

Chromosome 1			Population 1				Population 2			
Position	Ref	Alt	Hap1	Hap2	Hap3	Hap4	Hap5	Hap6	Hap7	Hap8
1	A	.	A	A	A	A	A	A	A	A
2 3	T	.	T	T	T	T	T	T	T	T
4 5	G	.	G	G	G	G	G	G	G	G
6 7	G	C	G	G	G	G	G	C	G	C
8 9	A	.	A	A	A	A	A	A	A	A
	T	.	T	T	T	T	T	T	T	T
	G	A	A	A	G	A	G	G	G	A
	G	.	G	G	G	G	G	G	G	G
	T	A	A	T	T	T	T	T	T	A

# Genetic variations

Chromosome 1			Population 1				Population 2			
Position	Ref	Alt	Hap1	Hap2	Hap3	Hap4	Hap5	Hap6	Hap7	Hap8
1	A	.	A	A	A	A	A	A	A	A
2 3	T	.	T	T	T	T	T	T	T	T
4 5	G	.	G	G	G	G	G	G	G	G
6 7	G	C	G	G	G	G	G	C	G	C
8 9	A	.	A	A	A	A	A	A	A	A
	T	.	T	T	T	T	T	T	T	T
	G	A	A	A	G	A	G	G	G	A
	G	.	G	G	G	G	G	G	G	G
	T	A	A	T	T	T	T	T	T	A

← Vocab recap!  
Positions that  
are not identical  
for all humans  
are called  
**SNPs** or  
“variants” or  
“segregating  
sites”...



# Genetic variations: allele frequency

Chromosome 1			Population 1				Population 2			
Position	Ref	Alt	Hap1	Hap2	Hap3	Hap4	Hap5	Hap6	Hap7	Hap8
4 7	G	C	G	G	G	G	G	C	G	C
9	G	A	A	A	G	A	G	G	G	A
	T	A	A	T	T	T	T	T	T	A

# Genetic variations: allele frequency

Chromosome 1				Population 1				Population 2			
Position	Ref	Alt		Hap1	Hap2	Hap3	Hap4	Hap5	Hap6	Hap7	Hap8
4	7	G	C	0	0	0	0	0	1	0	1
	9	G	A	1	1	0	1	0	0	0	1
		T	A	1	0	0	0	0	0	0	1

The allele frequency is the frequency of an allele in a population

Frequency of the reference allele at position 4:

In population 1 :  $4/4 = 1$

In population 2 :  $2/4 = 0.5$

Frequency of the reference allele at position 7:

In population 1 :  $1/4 = 0.25$

In population 2 :  $3/4 = 0.75$

Allele frequencies:

Position	Pop1	Pop2
4	1	0.5
7	0.25	0.75
9	0.75	0.75

# Genetic variations: allele frequency

## Fixation index (FST)

The FST is a statistic to measure differentiation between two population.

$H_T = 2 \cdot \sum f_i \cdot (1 - f_i)$  with  $f_i$  the allele frequency at position  $i$  in both populations. This measure the allele diversity taking both populations together.

$H_S = \sum (f_{i1}^1 \cdot (1 - f_{i1}^1) + f_{i1}^2 \cdot (1 - f_{i1}^2))$  with  $f_{i1}^1$  (resp.  $f_{i1}^2$ ) the allele frequency at position  $i$  in populations 1 and (resp. 2). This measure the allele diversity in both populations separately.

Allele frequencies:

Position	Pop1	Pop2	Pop1&Pop2
4	1	0.5	0.75
7	0.25	0.75	0.5
9	0.75	0.75	0.75

$$FST = \frac{H_T - H_S}{H_T}$$

# Genetic variations: allele frequency

## Fixation index (FST)

The FST is a statistic to measure differentiation between two population.

$H_T = 2 * \sum f_{i.} (1 - f_{i.})$  with  $f_{i.}$  the frequency of allele  $i$  in both populations. This measure the allele diversity taking both populations together.

$H_S = \sum (f_{i1}^1 (1 - f_{i1}^1) + f_{i1}^2 (1 - f_{i1}^2))$  with  $f_{i1}^1$  (resp.  $f_{i1}^2$ ) the frequency of allele  $i$  in populations 1 and (resp. 2).

This measure the allele diversity in both populations separately.

Allele frequencies:

Position	Pop1	Pop2	Pop1&Pop2
4	1	0.5	0.75
7	0.25	0.75	0.5
9	0.75	0.75	0.75

$$H_T = 2 * (0.75 * 0.25 + 0.5 * 0.5 + 0.75 * 0.25) = 1.25$$

$$FST = \frac{H_T - H_S}{H_T}$$

# Genetic variations: allele frequency

## Fixation index (FST)

The FST is a statistic to measure differentiation between two population.

$H_T = 2 * \sum f_{i.} (1 - f_{i.})$  with  $f_{i.}$  the frequency of allele  $i$  in both populations. This measure the allele diversity taking both populations together.

$H_S = \sum (f_{i1}^1 (1 - f_{i1}^1) + f_{i1}^2 (1 - f_{i1}^2))$  with  $f_{i1}^1$  (resp.  $f_{i1}^2$ ) the frequency of allele  $i$  in populations 1 and (resp. 2).

This measure the allele diversity in both populations separately.

Allele frequencies:

Position	Pop1	Pop2	Pop1&Pop2
4	1	0.5	0.75
7	0.25	0.75	0.5
9	0.75	0.75	0.75

$$H_T = 2 * (0.75 * 0.25 + 0.5 * 0.5 + 0.75 * 0.25) = 1.25$$

$$H_S = 1 * 0 + 0.25 * 0.75 + 0.75 * 0.25 + \dots$$

$$FST = \frac{H_T - H_S}{H_T}$$

# Genetic variations: allele frequency

## Fixation index (FST)

The FST is a statistic to measure differentiation between two population.

$H_T = 2 \sum f_{i.}(1-f_{i.})$  with  $f_{i.}$  the frequency of allele  $i$  in both populations. This measure the allele diversity taking both populations together.

$H_S = \sum (f_{i.}^1(1-f_{i.}^1) + f_{i.}^2(1-f_{i.}^2))$  with  $f_{i.}^1$  (resp.  $f_{i.}^2$ ) the frequency of allele  $i$  in populations 1 and (resp. 2).

This measure the allele diversity in both populations separately.

Allele frequencies:

Position	Pop1	Pop2	Pop1&Pop2
4	1	0.5	0.75
7	0.25	0.75	0.5
9	0.75	0.75	0.75

$$H_T = 2 * (0.75 * 0.25 + 0.5 * 0.5 + 0.75 * 0.25) = 1.25$$

$$H_S = 1 * 0 + 0.25 * 0.75 + 0.75 * 0.25 + 0.5 * 0.5 + 0.75 * 0.25 + 0.75 * 0.25 = 1$$

$$FST = \frac{H_T - H_S}{H_T}$$

# Genetic variations: allele frequency

## Fixation index (FST)

The FST is a statistic to measure differentiation between two population.

$H_T = 2 \cdot \sum f_{i.} \cdot (1 - f_{i.})$  with  $f_{i.}$  the frequency of allele  $i$  in both populations. This measure the allele diversity taking both populations together.

$H_S = \sum (f_{i.}^1 \cdot (1 - f_{i.}^1) + f_{i.}^2 \cdot (1 - f_{i.}^2))$  with  $f_{i.}^1$  (resp.  $f_{i.}^2$ ) the frequency of allele  $i$  in populations 1 and (resp. 2).

This measure the allele diversity in both populations separately.

Allele frequencies:

Position	Pop1	Pop2	Pop1&Pop2
4	1	0.5	0.75
7	0.25	0.75	0.5
9	0.75	0.75	0.75

$$H_T = 2 \cdot (0.75 \cdot 0.25 + 0.5 \cdot 0.5 + 0.75 \cdot 0.25) = 1.25$$

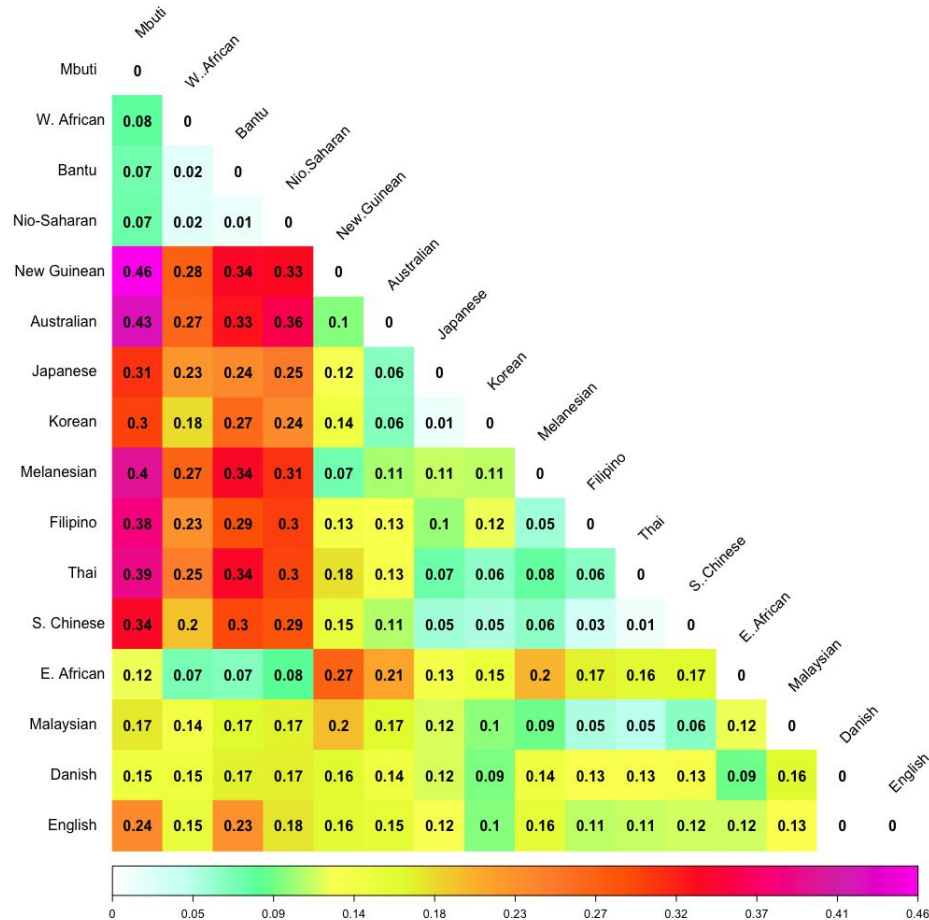
$$H_S = 1 \cdot 0 + 0.25 \cdot 0.75 + 0.75 \cdot 0.25 + 0.5 \cdot 0.5 + 0.75 \cdot 0.25 + 0.75 \cdot 0.25 = 1$$

$$FST = (1.25 - 1) / 1.25 = 0.2$$

$$FST = \frac{H_T - H_S}{H_T}$$

Conclusion: There is 20% more differences between the two populations than within each one.

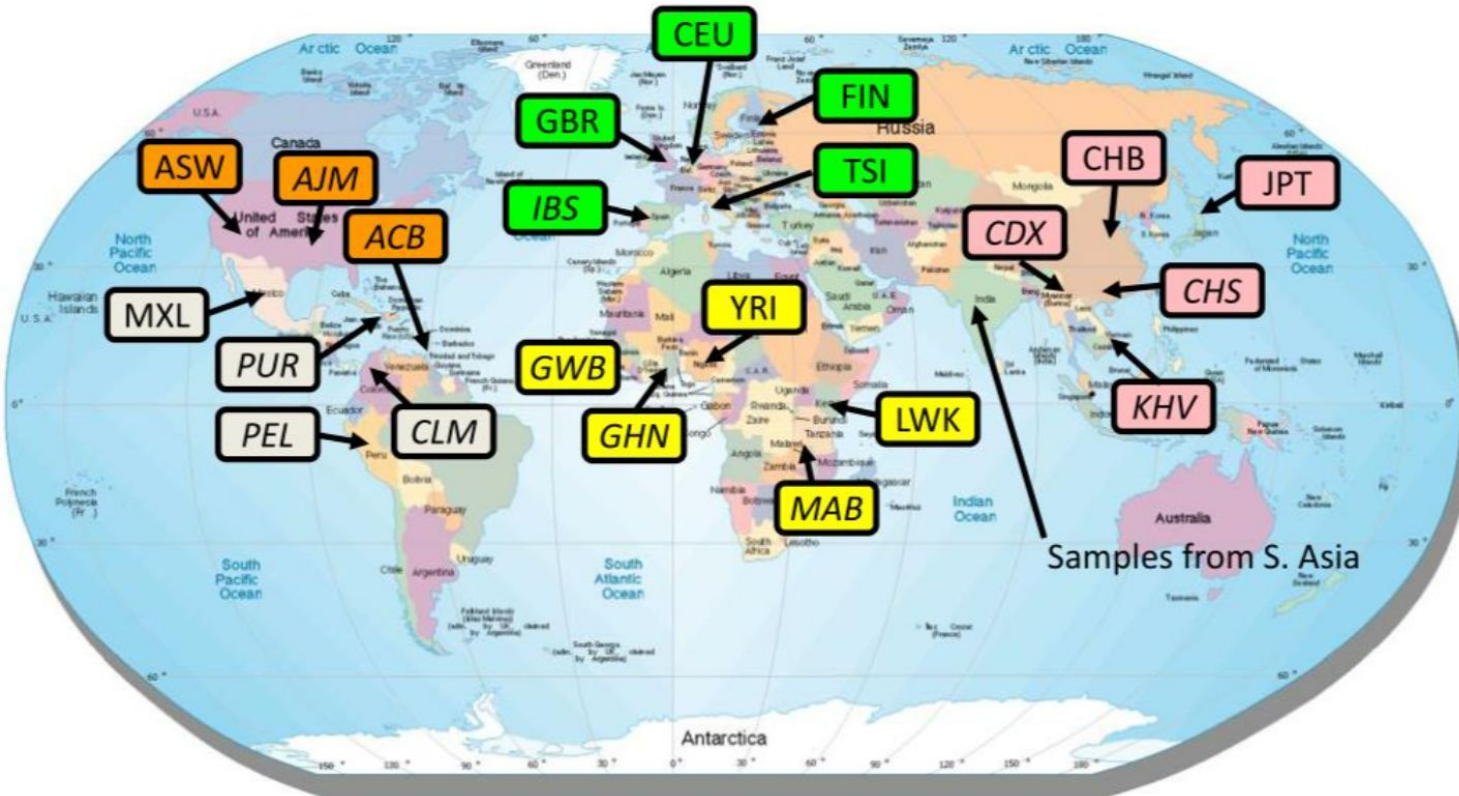
# Genetic variations: allele frequency - (FST)



- The average FST between continent is at around 10%.
- If we were to do an FST-like statistic on skin color, that number would be around 90%.
- The reason is strong selection on genes responsible for the skin color.



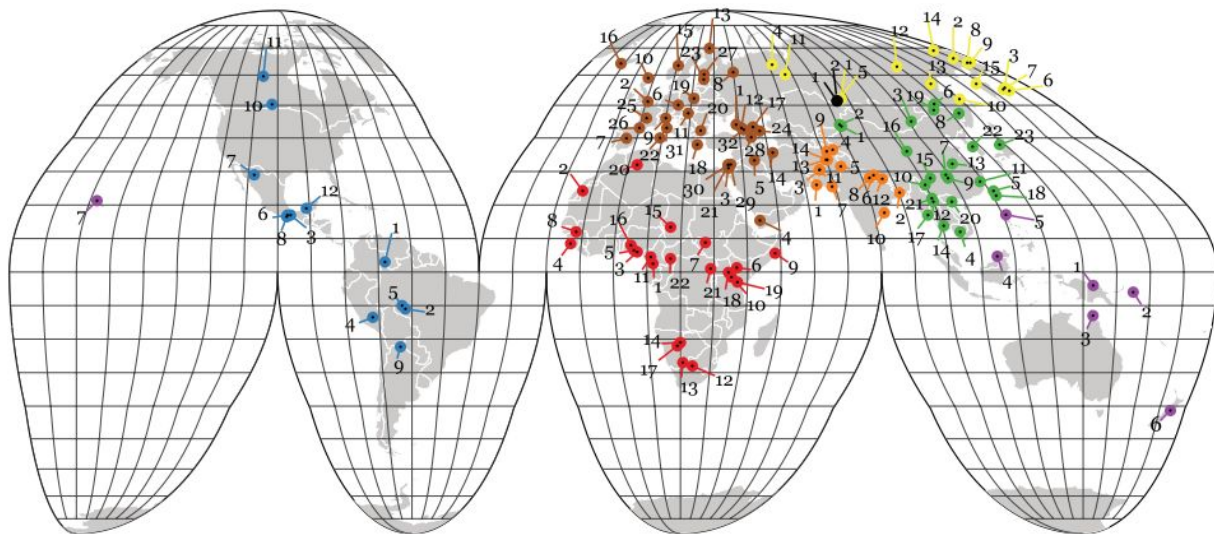
# Data availability - 1000 Genome project



Sequencing of around 100 genomes per 12 different populations around the globe.

Data is publicly available.

# Data availability - SGDP



236 individual  
genomes (total) from  
125 distinct human  
populations

- |                   |                      |                       |                      |                      |                               |
|-------------------|----------------------|-----------------------|----------------------|----------------------|-------------------------------|
| ● 1 - Neanderthal | ● 9 - Miao           | ● 1 - Adygei          | ● 23 - Estonian      | ● 13 - KhomaniSan    | ● 13 - Hazara                 |
| ● 2 - Denisova    | ● 10 - Naxi          | ● 2 - EnglandGBR      | ● 24 - Lezgin        | ● 14 - San           | ● 14 - Pathan                 |
| ● 1 - Piapoco     | ● 11 - She           | ● 3 - Bedouin         | ● 25 - French        | ● 15 - GambiaGWD     | ● 1 - Tubalar                 |
| ● 2 - Surui       | ● 12 - Dai           | ● 4 - YemeniteJew     | ● 26 - Basque        | ● 16 - Yoruba        | ● 2 - EskimoNaukan            |
| ● 3 - Mixe        | ● 13 - Han           | ● 5 - IraqiJew        | ● 27 - FinlandFIN    | ● 17 - BantuSEHerero | ● 3 - AleutBeringIsland       |
| ● 4 - Quechua     | ● 14 - Thai          | ● 6 - Czechoslovakian | ● 28 - Armenian      | ● 18 - Luo           | ● 4 - MansiSosvaRiver         |
| ● 5 - Karitiana   | ● 15 - Yi            | ● 7 - SpainIBS        | ● 29 - Jordanian     | ● 19 - BantuKenya    | ● 5 - Altaian                 |
| ● 6 - Mixtec      | ● 16 - Tu            | ● 8 - Russian         | ● 30 - Druze         | ● 20 - Mozabite      | ● 6 - TlingitPreobrazhenskoye |
| ● 7 - Pima        | ● 17 - Burmese       | ● 9 - Bergamo         | ● 31 - Tuscan        | ● 21 - Mbuti         | ● 7 - TlingitNikolskoye       |
| ● 8 - Zapotec     | ● 18 - Ami           | ● 10 - Orcadian       | ● 32 - Georgian      | ● 22 - Biaka         | ● 8 - EskimoChaplino          |
| ● 9 - Chane       | ● 19 - Daur          | ● 11 - Hungarian      | ● 1 - Kongo          | ● 1 - Makrani        | ● 9 - EskimoSireniki          |
| ● 10 - Cree       | ● 20 - KinhKVVH      | ● 12 - Abkhasian      | ● 2 - Saharawi       | ● 2 - BengaliBEB     | ● 10 - Ulchi                  |
| ● 11 - Chipewyan  | ● 21 - Lahu          | ● 13 - Saami          | ● 3 - Igbo           | ● 3 - Balochi        | ● 11 - MansiKondaRiver        |
| ● 12 - Maya       | ● 22 - Korean        | ● 14 - Iranian        | ● 4 - MendeMSL       | ● 4 - Burusho        | ● 12 - Yakut                  |
| ● 1 - Xibo        | ● 23 - Japanese      | ● 15 - Norwegian      | ● 5 - EsanESN        | ● 5 - PunjabiPJL     | ● 13 - Even                   |
| ● 2 - Uyghur      | ● 1 - Papuan         | ● 16 - Icelandic      | ● 6 - LuhyaLWK       | ● 6 - Tibetan        | ● 14 - Chukchi                |
| ● 3 - Mongola     | ● 2 - Bougainville   | ● 17 - Chechen        | ● 7 - Dinka          | ● 7 - Sindhi         | ● 15 - Itelman                |
| ● 4 - Cambodian   | ● 3 - AustraliaECCAC | ● 18 - Greek          | ● 8 - Mandenka       | ● 8 - Kusunda        |                               |
| ● 5 - Atayal      | ● 4 - Dusun          | ● 19 - Polish         | ● 9 - Somali         | ● 9 - Kalash         |                               |
| ● 6 - Oroqen      | ● 5 - Igorot         | ● 20 - Bulgarian      | ● 10 - MasaiMKK      | ● 10 - Madiga        |                               |
| ● 7 - Tujia       | ● 6 - Maori          | ● 21 - Palestinian    | ● 11 - Lemande       | ● 11 - Brahui        |                               |
| ● 8 - Hezhen      | ● 7 - Hawaiian       | ● 22 - Sardinian      | ● 12 - BantuSETswana | ● 12 - Sherpa        |                               |

# Genetic distances

Allele frequencies:

Position	Pop1	Pop2	Pop3
4	1	0.5	0.85
7	0.25	0.75	0.35
9	0.75	0.75	0.8

How to compute the distance  
between population i and j at  
some position:

$$d_{ij} = |p_i - p_j|$$

With  $p_i$  and  $p_j$  allele frequencies.

We then compute mean by taking  
the average value for all  
positions.

# Genetic distances - Population trees

Allele frequencies:

Position	Pop1	Pop2	Pop3
4	1	0.5	0.85
7	0.25	0.75	0.35
9	0.75	0.75	0.8

$$d(\text{pop1}, \text{pop2}) = (|1-0.5| + |0.25-0.75| + |0.75-0.75|)/3 = 0.33$$

$$d(\text{pop1}, \text{pop3}) = 0.1$$

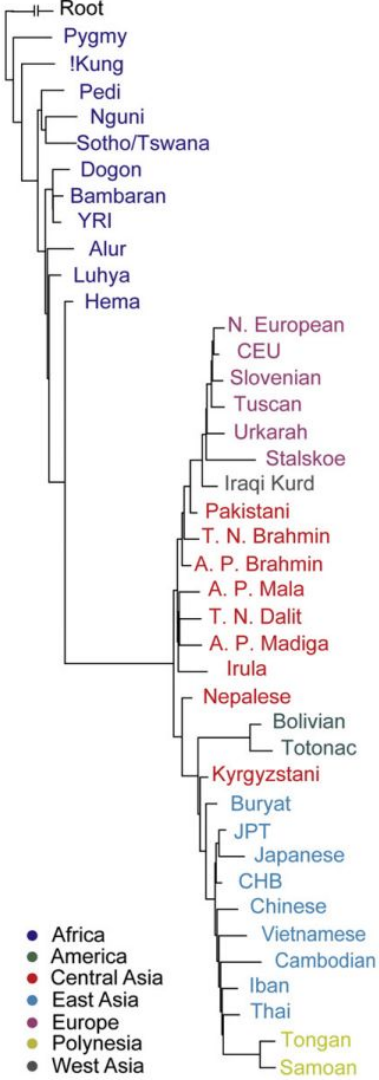
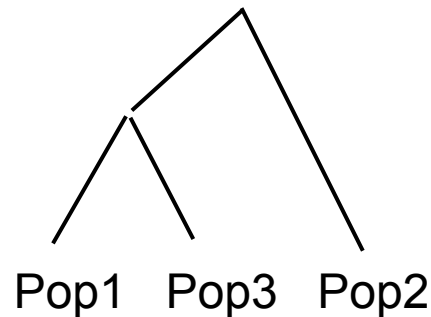
$$d(\text{pop2}, \text{pop3}) = 0.26$$

How to compute the distance between population i and j at some position:

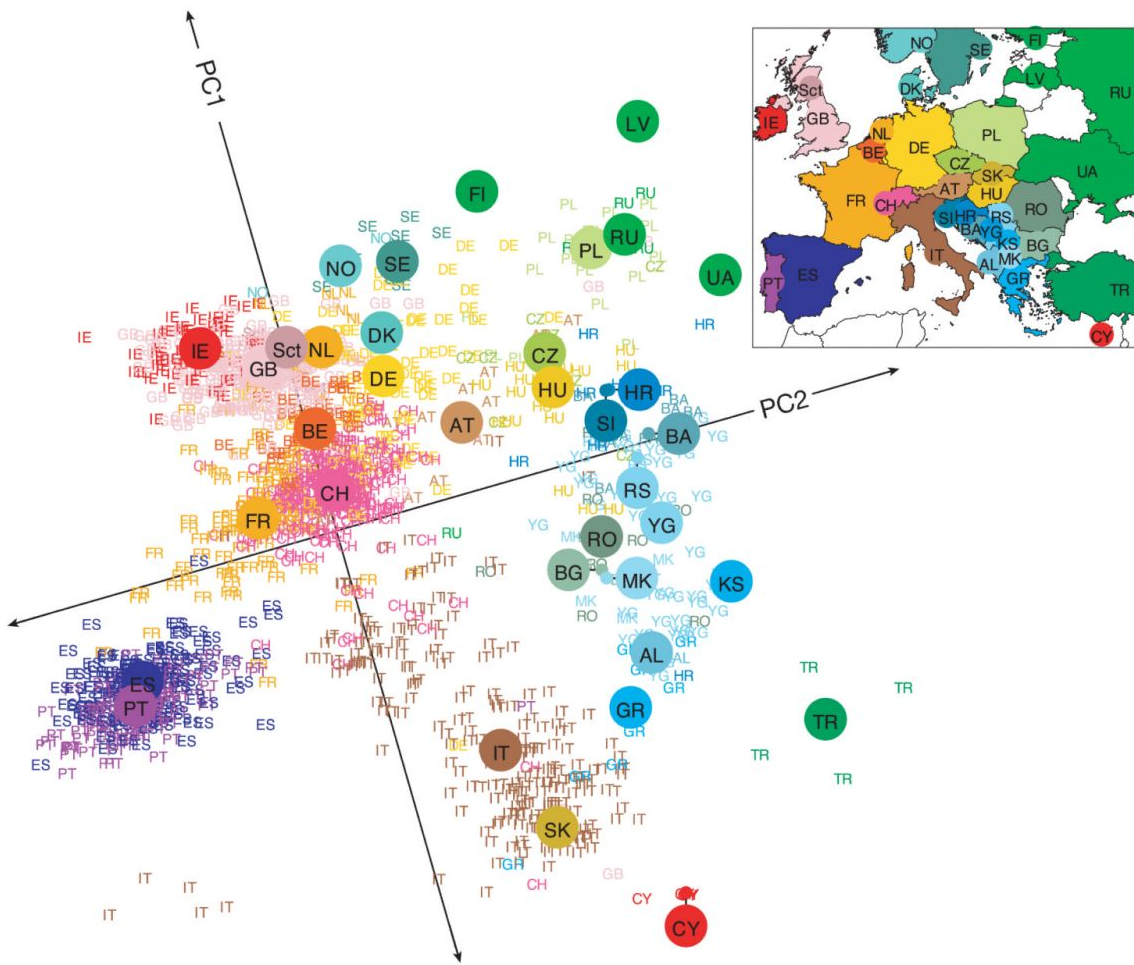
$$d_{ij} = |p_i - p_j|$$

With  $p_i$  and  $p_j$  allele frequencies.

We then compute mean by taking the average value for all positions.



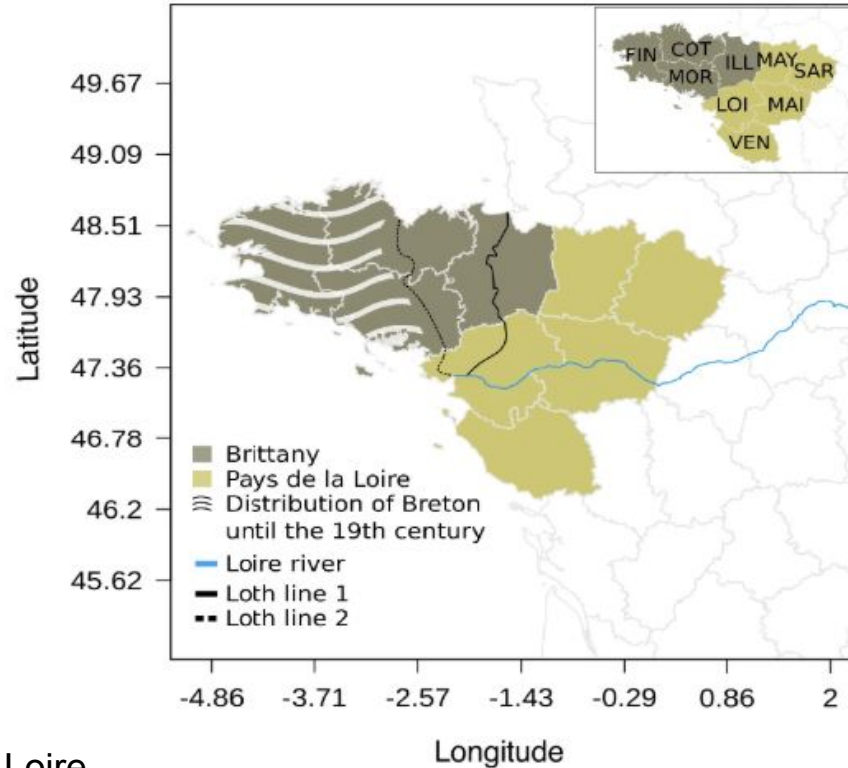
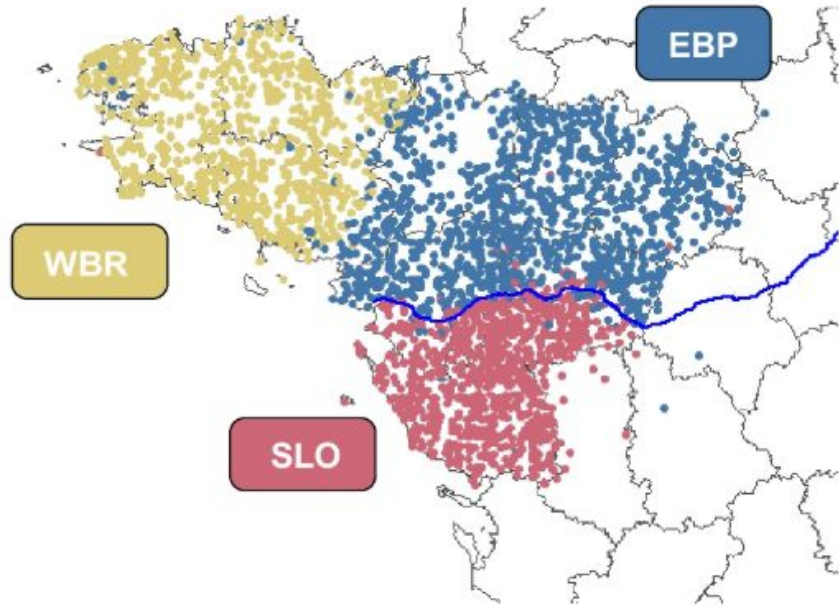
# Principal Component Analysis (PCA)



Dataset of 197,146 SNPs per individual projected on 2-d space.

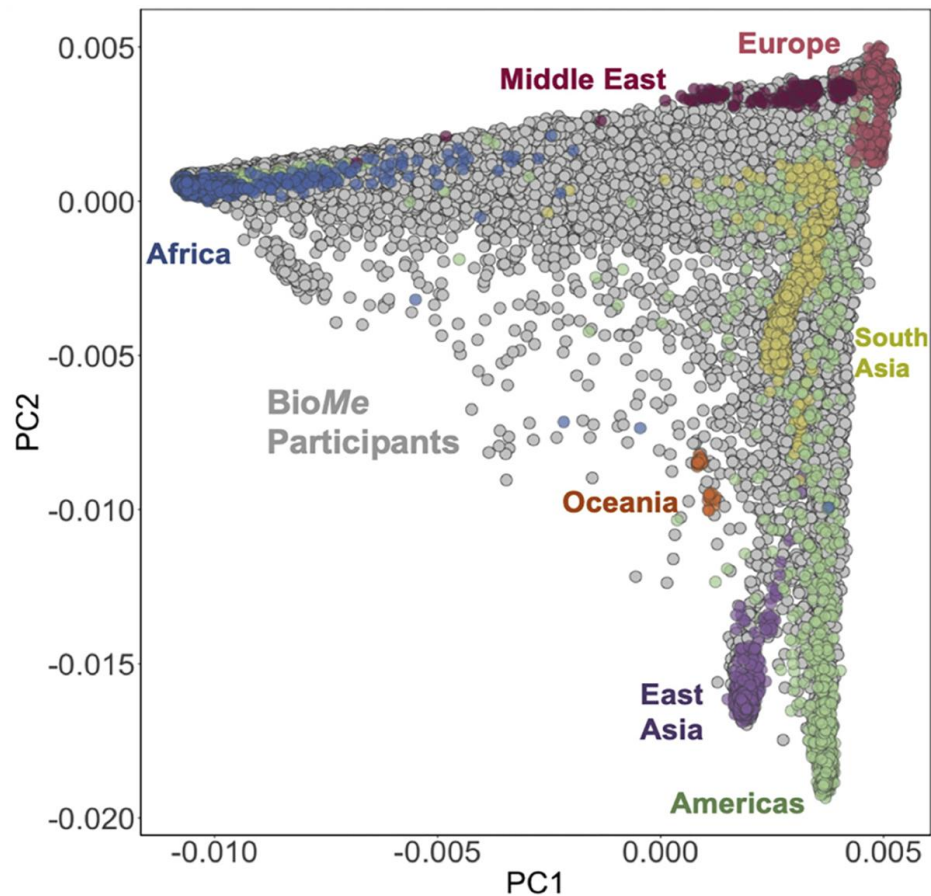


# k-clustering



WBR: Western Brittany, EBP: Eastern Brittany and Pays-de-la-Loire,  
SLO:South Loire. 3234 present day individuals.  
Alves et al. (2024)

# Population as strict categories?

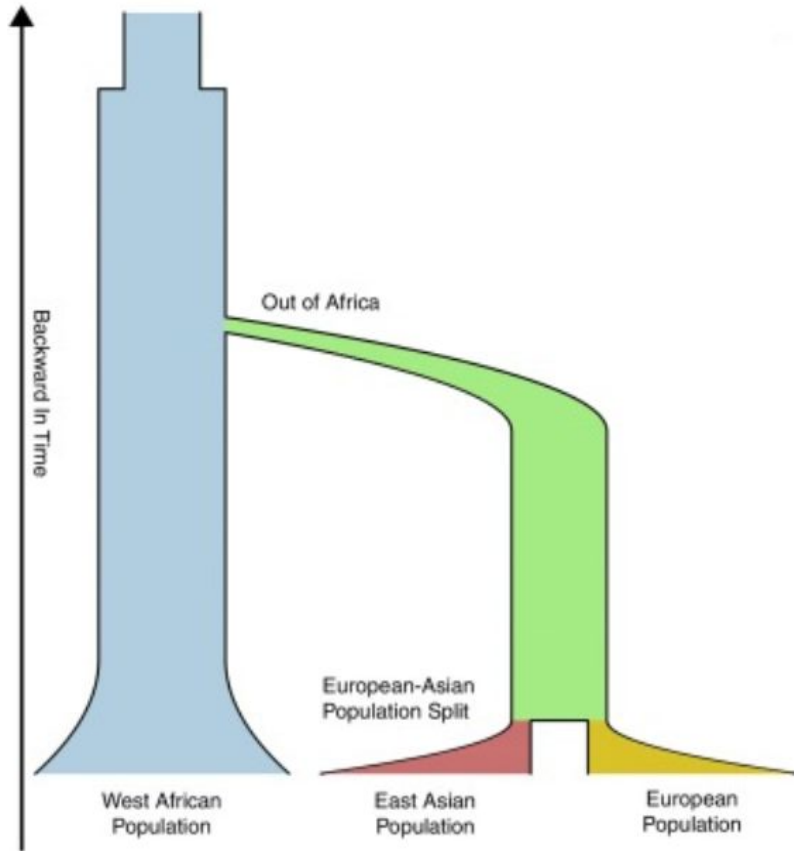


Lewis et al. (2022)

Population size

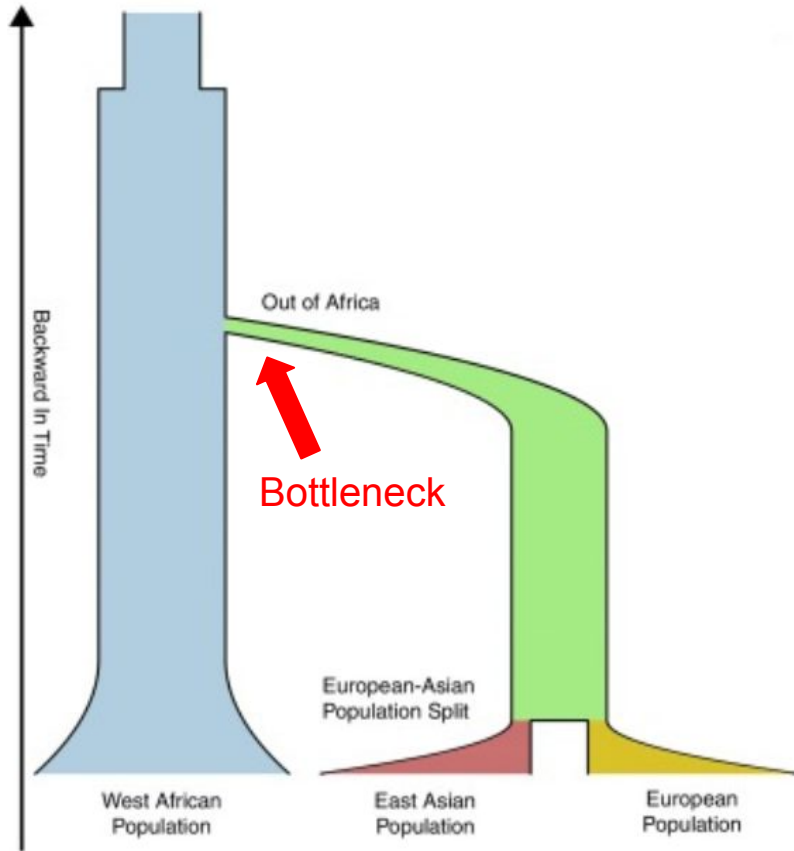


# Population size



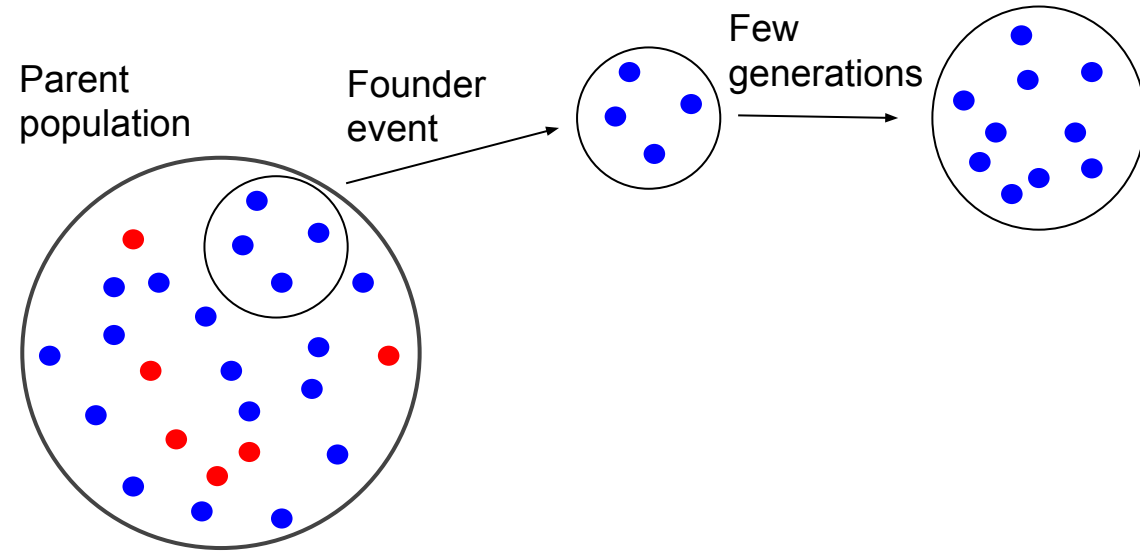
- Population sizes vary through time
- This can be inferred from present day genomes
- First idea:  
low genetic diversity  $\Leftrightarrow$  low population size

# Population size: bottlenecks

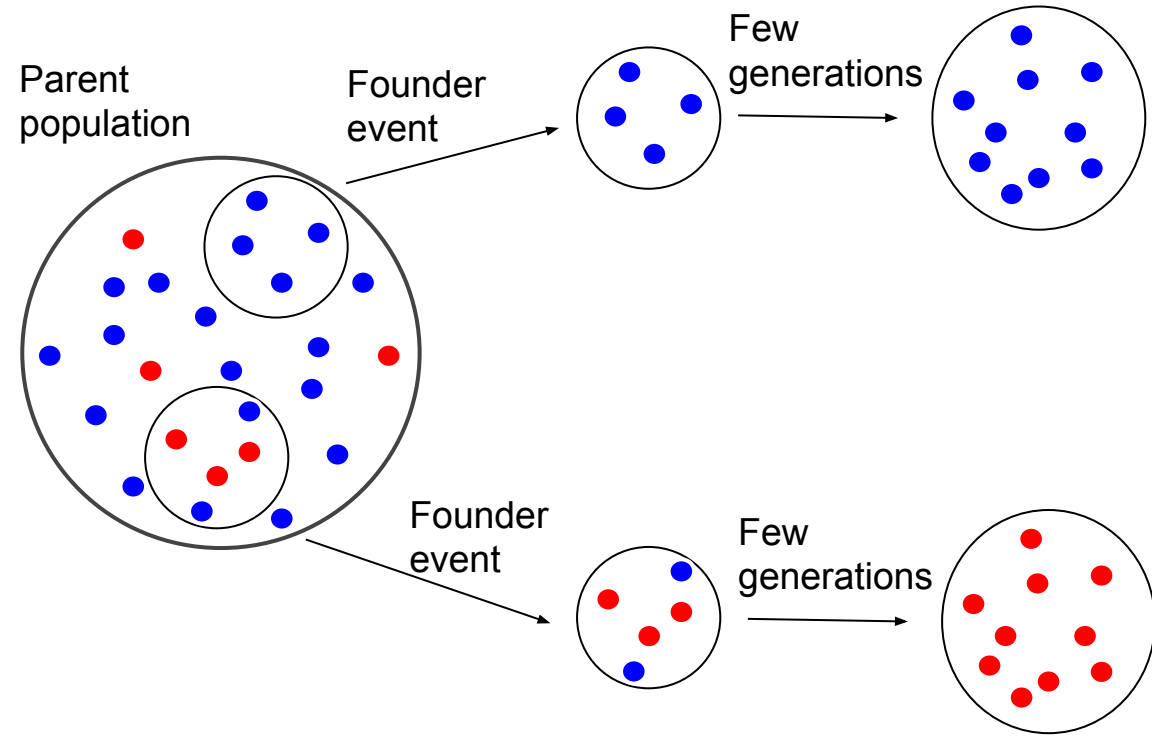


- A bottleneck is a large reduction in population size
- During such events, lots of variants are lost and genetic diversity suddenly decreases

# Population size - bottlenecks and founder effect

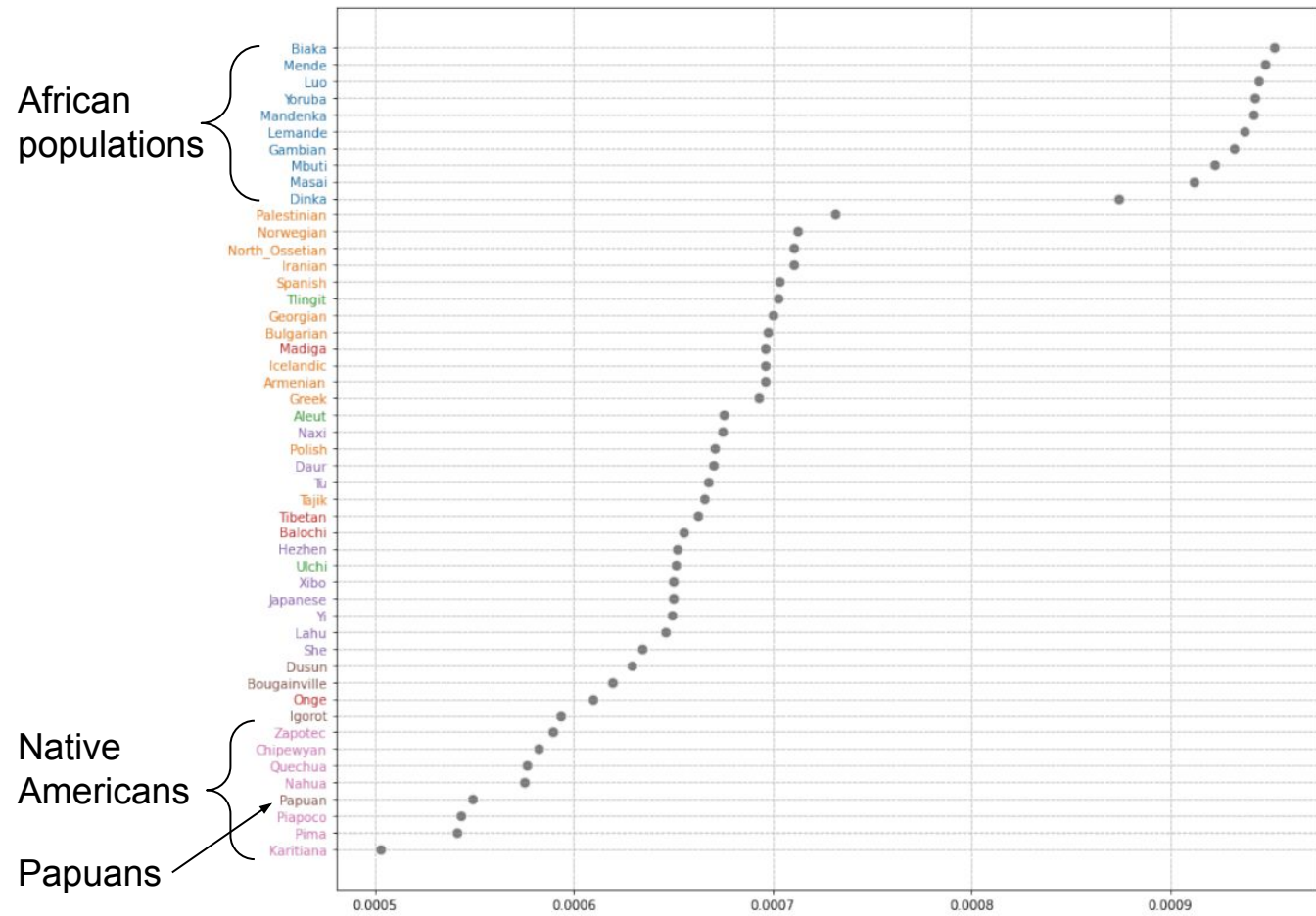


# Population size - bottlenecks and founder effect



- A bottleneck can be caused by migrations (e.g. the Out of Africa event or the peopling of the Americas) or by environmental events.
- Due to the **founder effect**, a bottleneck will cause a loss of genetic diversity that cannot be recovered even if the population grows afterward.

# Population size - bottlenecks consequences



- African populations have the highest amount of genetic diversity.
- Native Americans have the least genetic diversity due to multiple consecutive bottlenecks.

x-axis: heterozygosity proportion (higher -> more genetic diversity)

# Inferring population size history

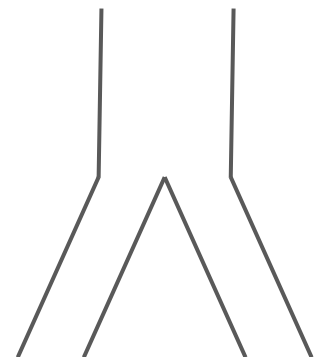
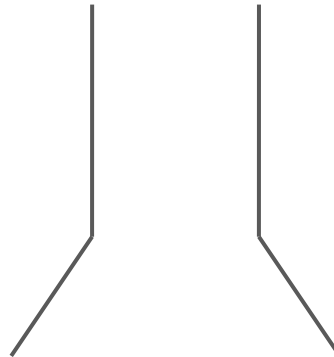
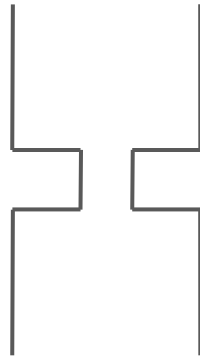
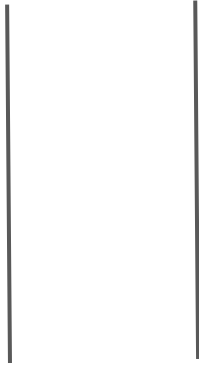
Constant

Bottleneck

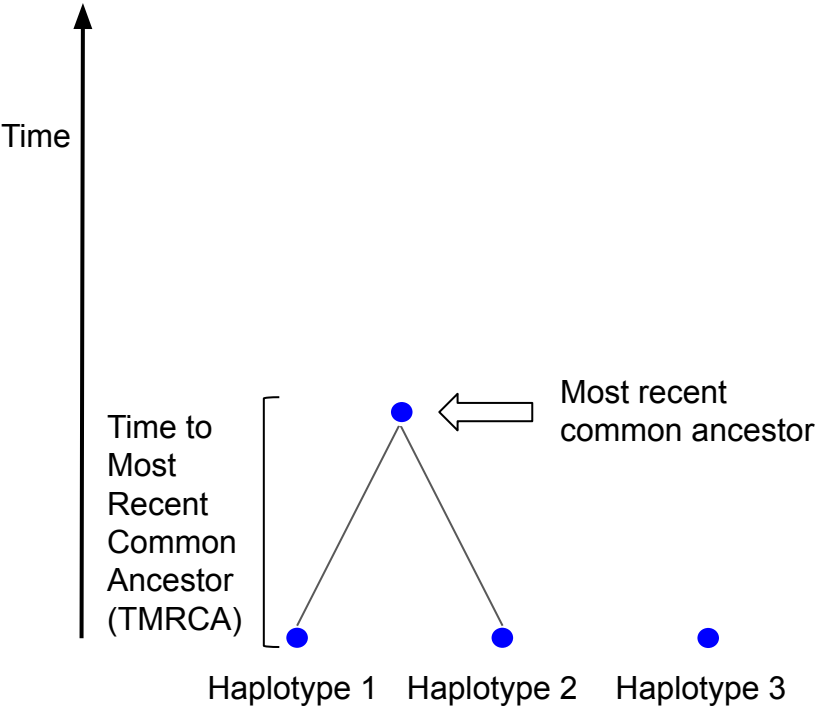
Expansion

Structured  
population

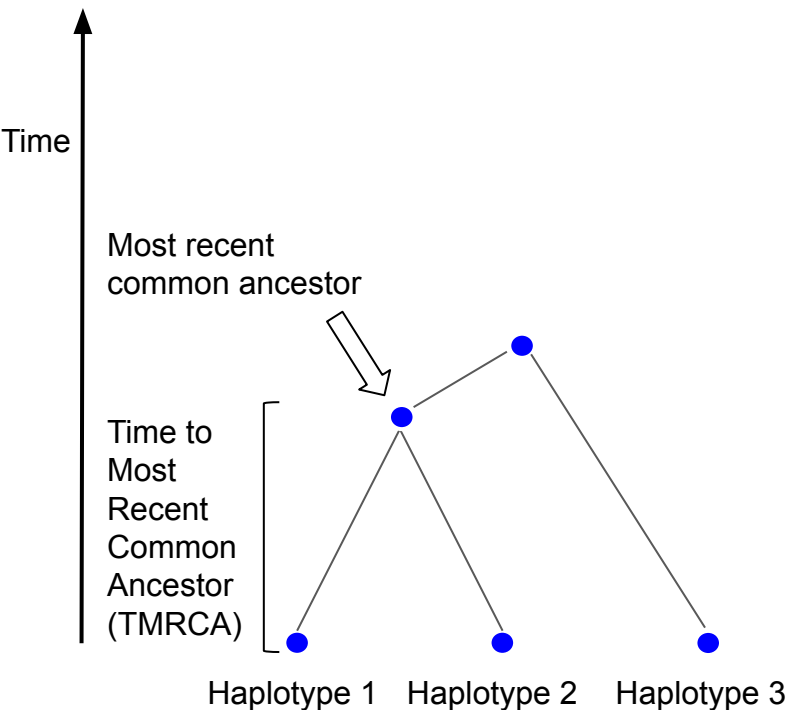
Models



# Inferring population size - Coalescence trees



# Inferring population size - Coalescence trees

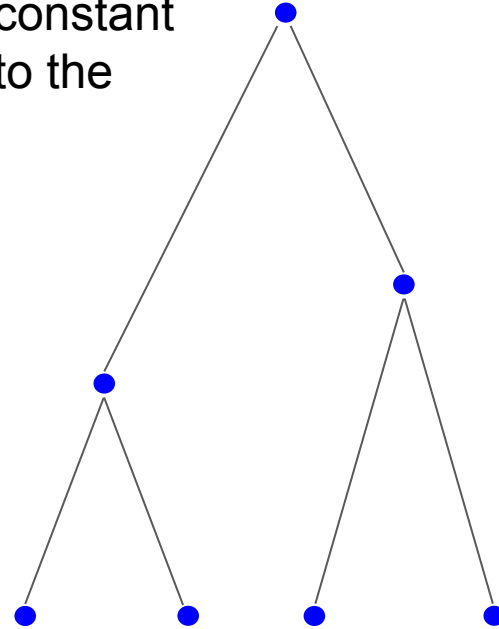
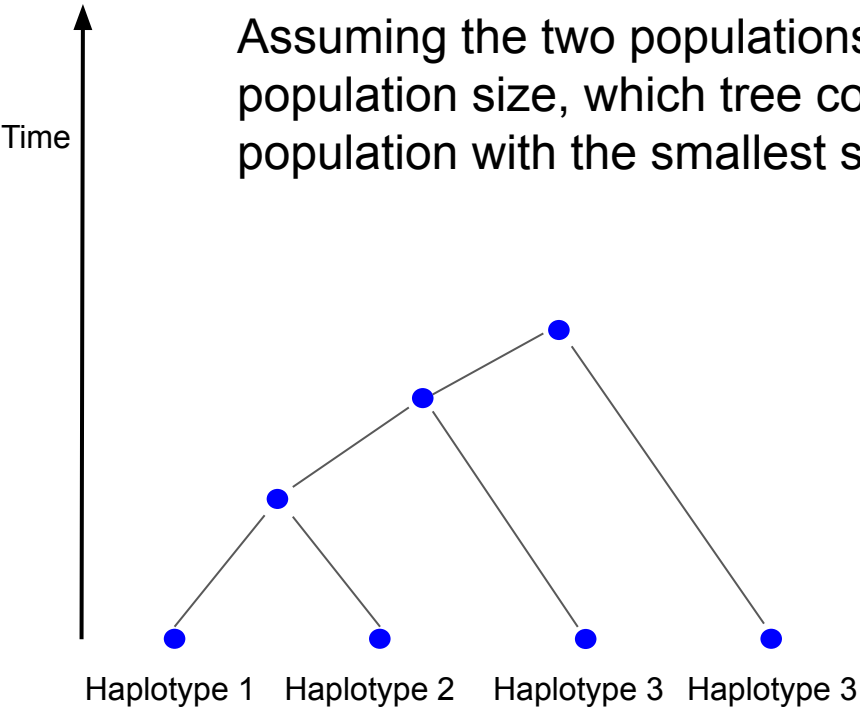


Some vocabulary:

- Any two haplotypes must come from the same ancestral haplotype at some point in the past.
- We say that the haplotypes **coalesce** (merge) together at time **TMRCA**.
- If we take multiple haplotypes, we can build a tree from this, it is called a **coalescence tree**.

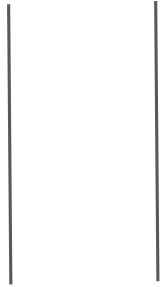


# Inferring population size - Coalescence trees

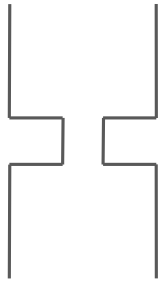


# Inferring population size history

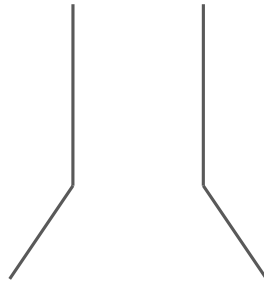
Constant



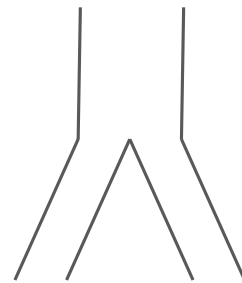
Bottleneck



Expansion

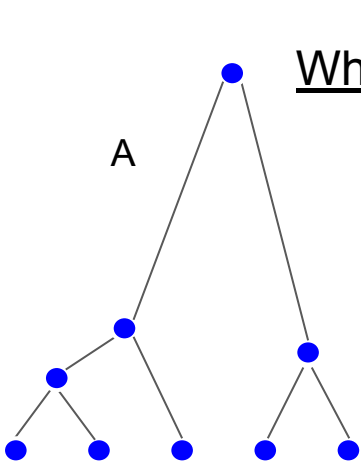


Structured  
population

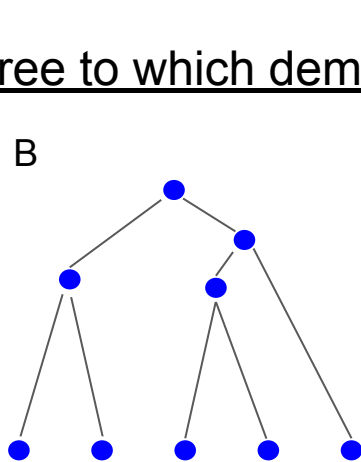


Which tree to which demography?

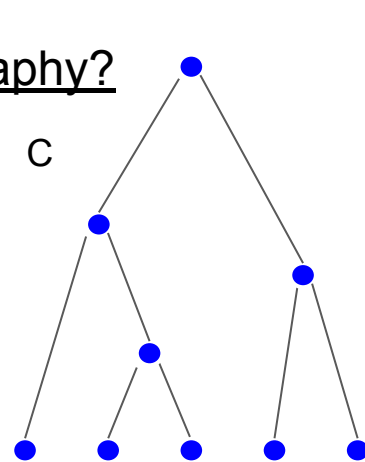
A



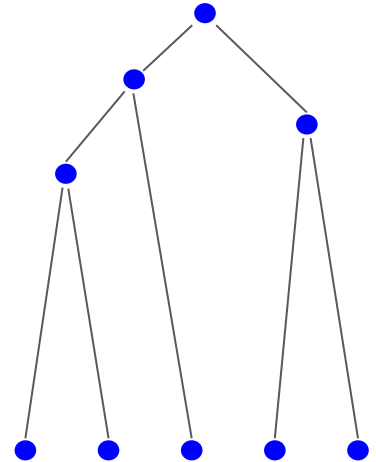
B



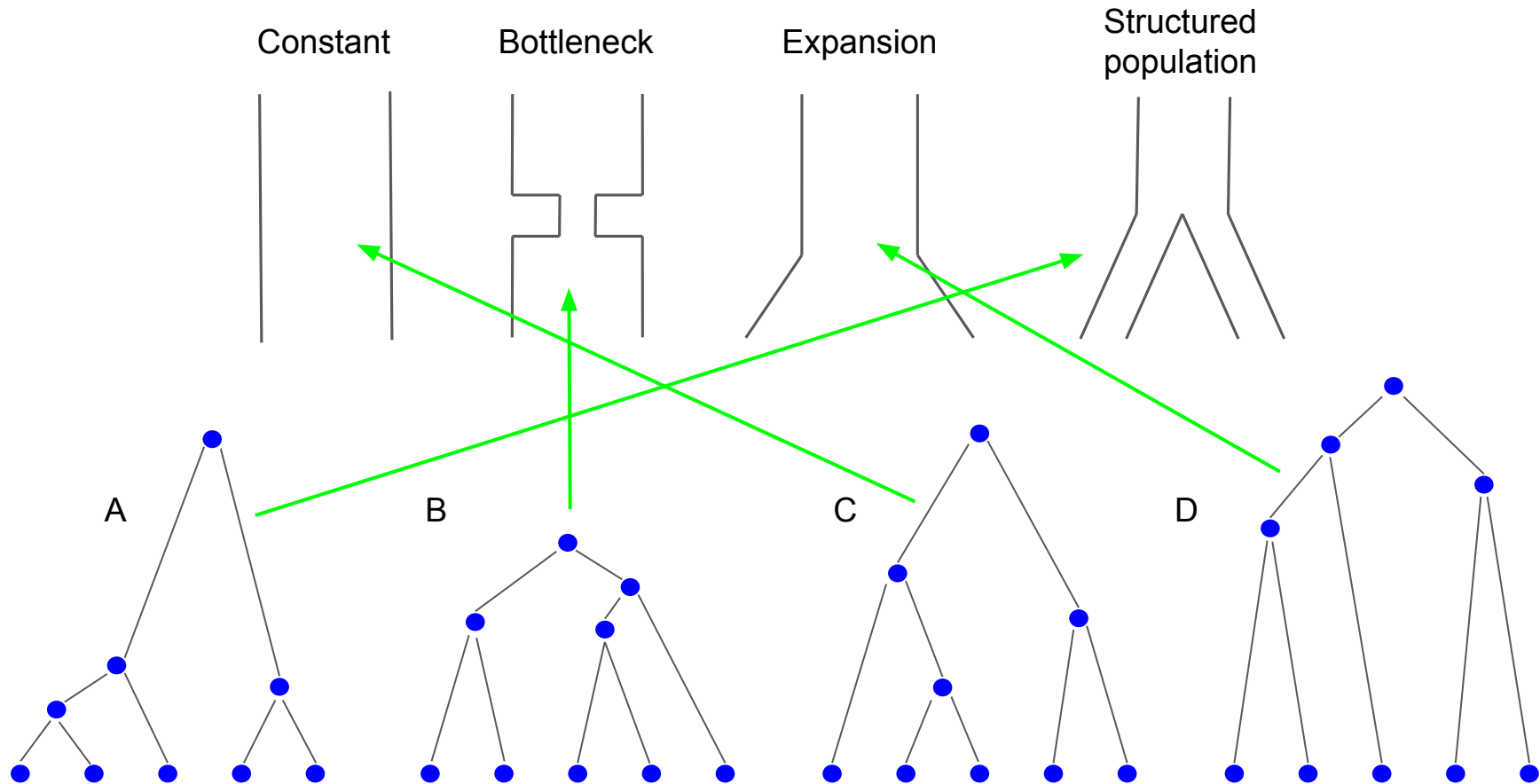
C



D



# Inferring population size history



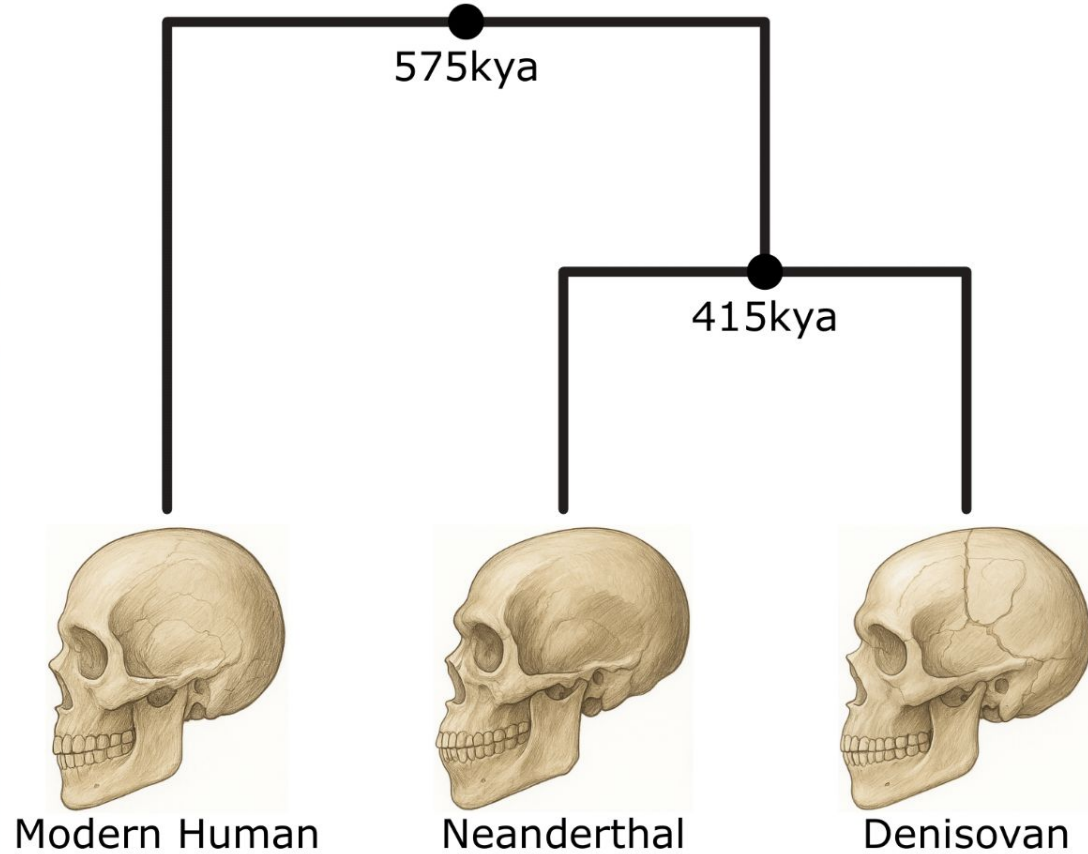
# Inferring population size history

## Conclusion:

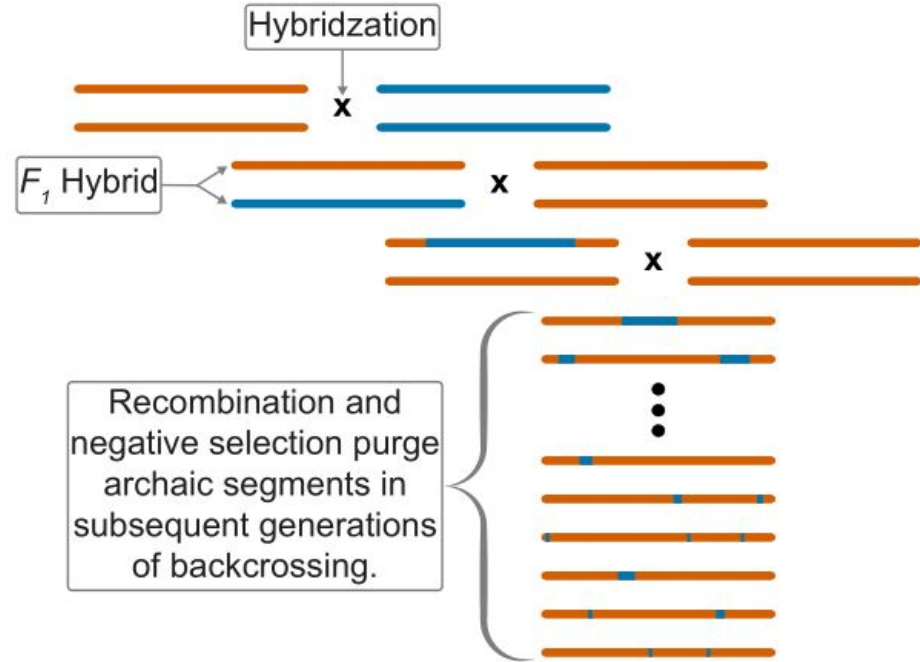
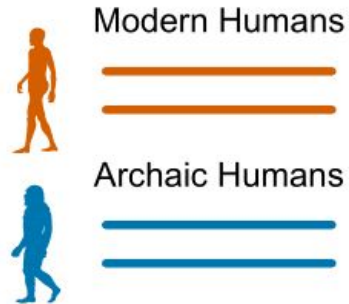
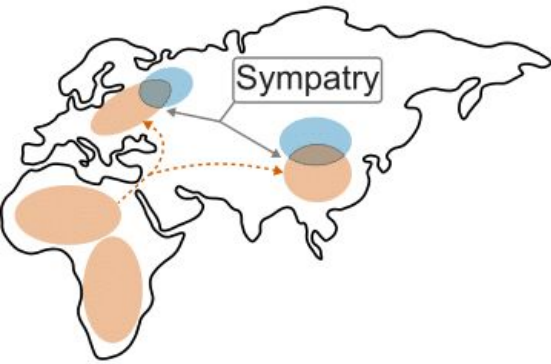
- If we are able to construct the coalescence tree from a population, we can know its whole population size history.
- Population size history can notably inform us about past migrations or ecological events.
- Still, how do we construct this coalescence tree from genomic data?
- The full answer is not for today! Key idea: If two individuals are closely genetically related, then they have a recent common ancestor and conversely.

Archaic ancestry

# Archaic Humans



# Archaic Introgression



# Archaic Introgression

Possible questions:

- Do we have archaic DNA, if so, how much and in which populations?
- Can we locate where archaic DNA is present in our genomes?
- Was there positive or negative selection for archaic DNA in Sapiens?
- If Archaic DNA is found in multiple contemporary populations (i.e Europe and Asia), does it have the same source?
- Can we date the introgression(s) event(s)?
- How did this archaic DNA spread through time?



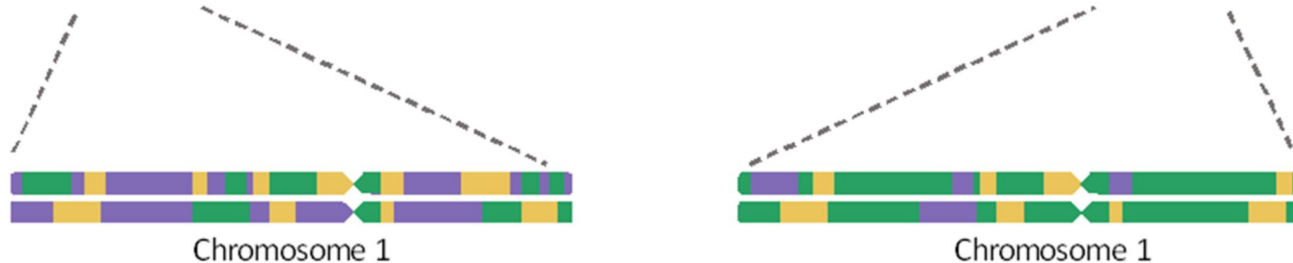
# Global and Local ancestry inference

A.

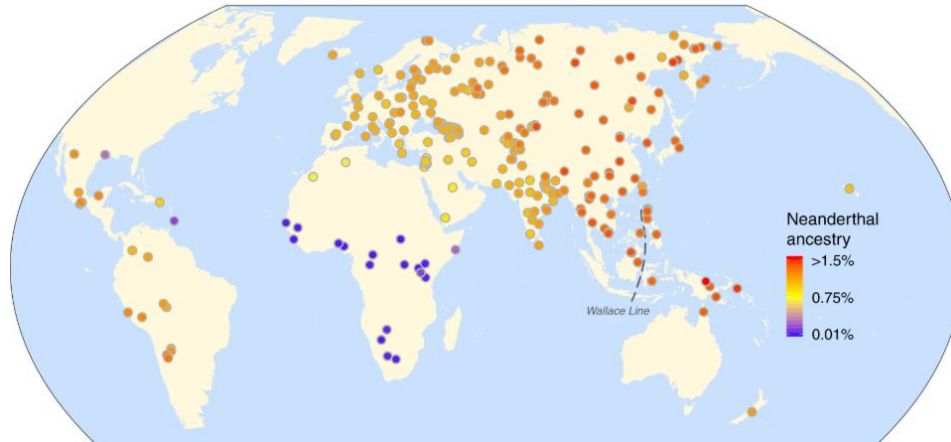
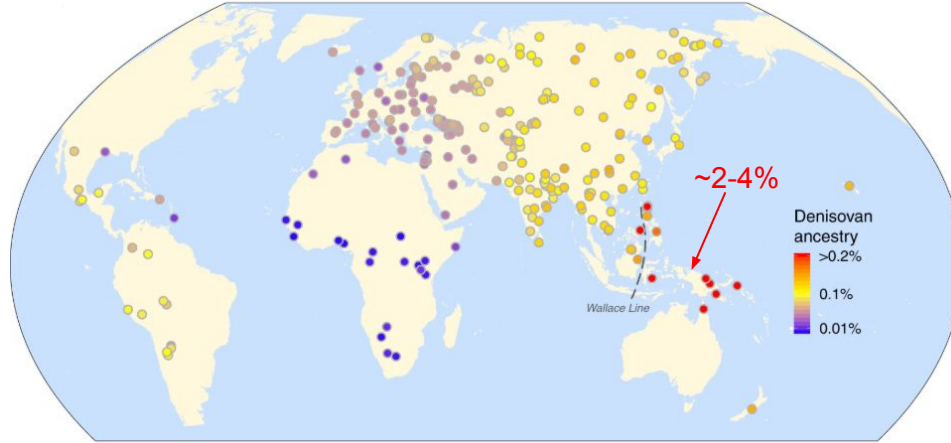


- A. Global ancestry inference.
- B. Local ancestry inference (LAI)

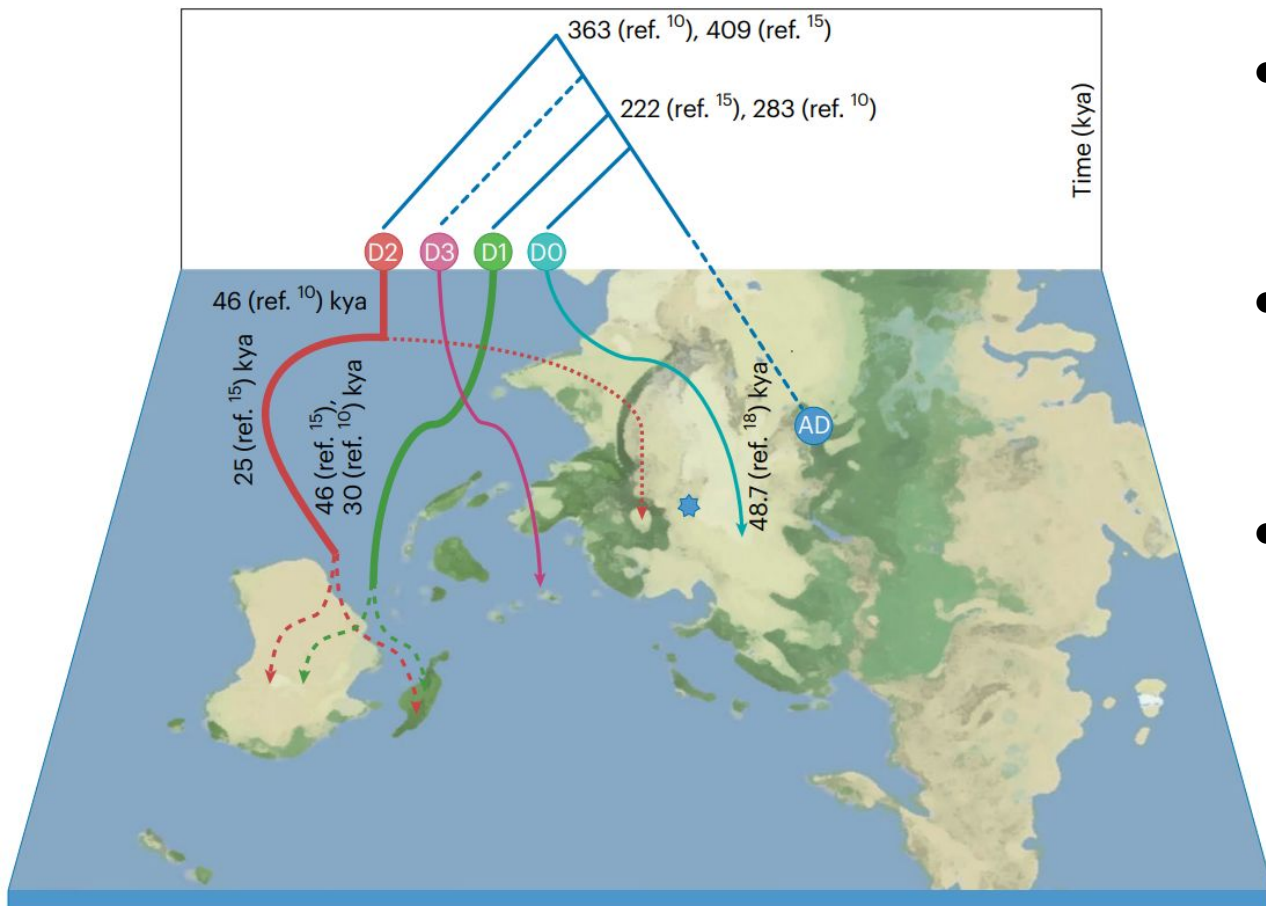
B.



# Archaic Introgression - Global ancestry



# Archaic Introgression - Multiple sources



- Most likely a single Neanderthal introgression event
- Multiple Denisovan introgression events, the number is still debated
- This can be deduced by comparing Denisovan segments in modern humans to the Denisovan reference genome

# Archaic Introgression - Adaptively introgressed genes

