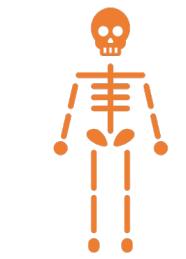




aDNA: Methods and Applications

aDNA analysis



Samples



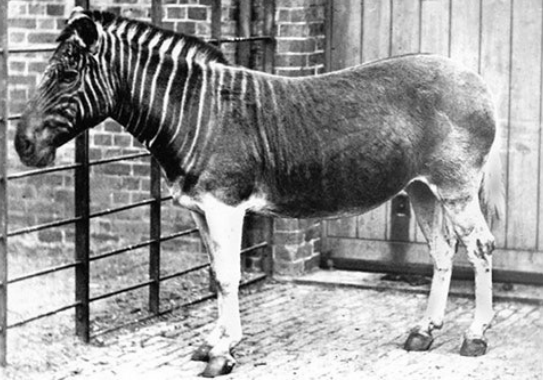
Genetic
information



Bioinformatic
analysis



Evolutionary and
Historical
reconstruction



Quagga

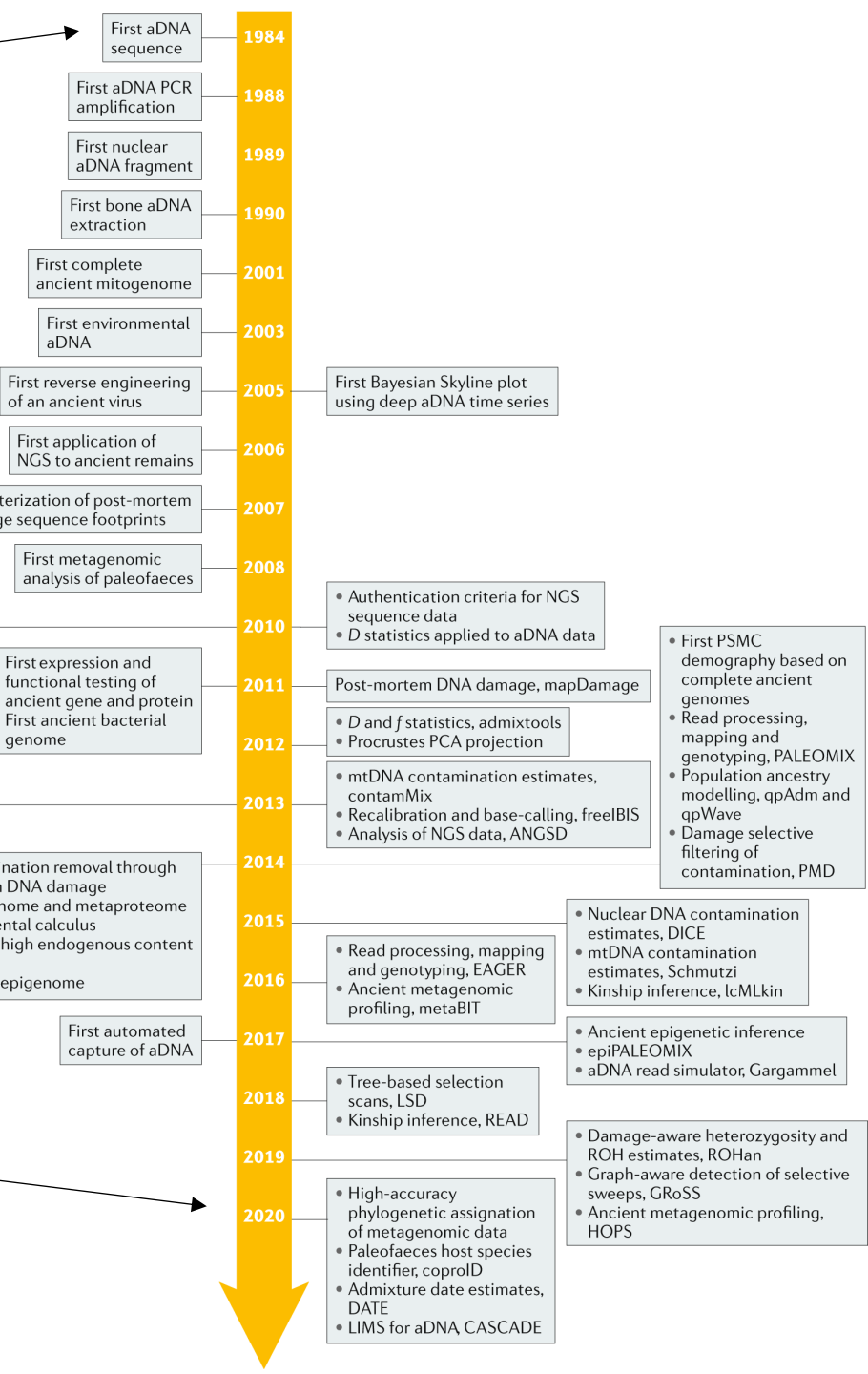


Neanderthal

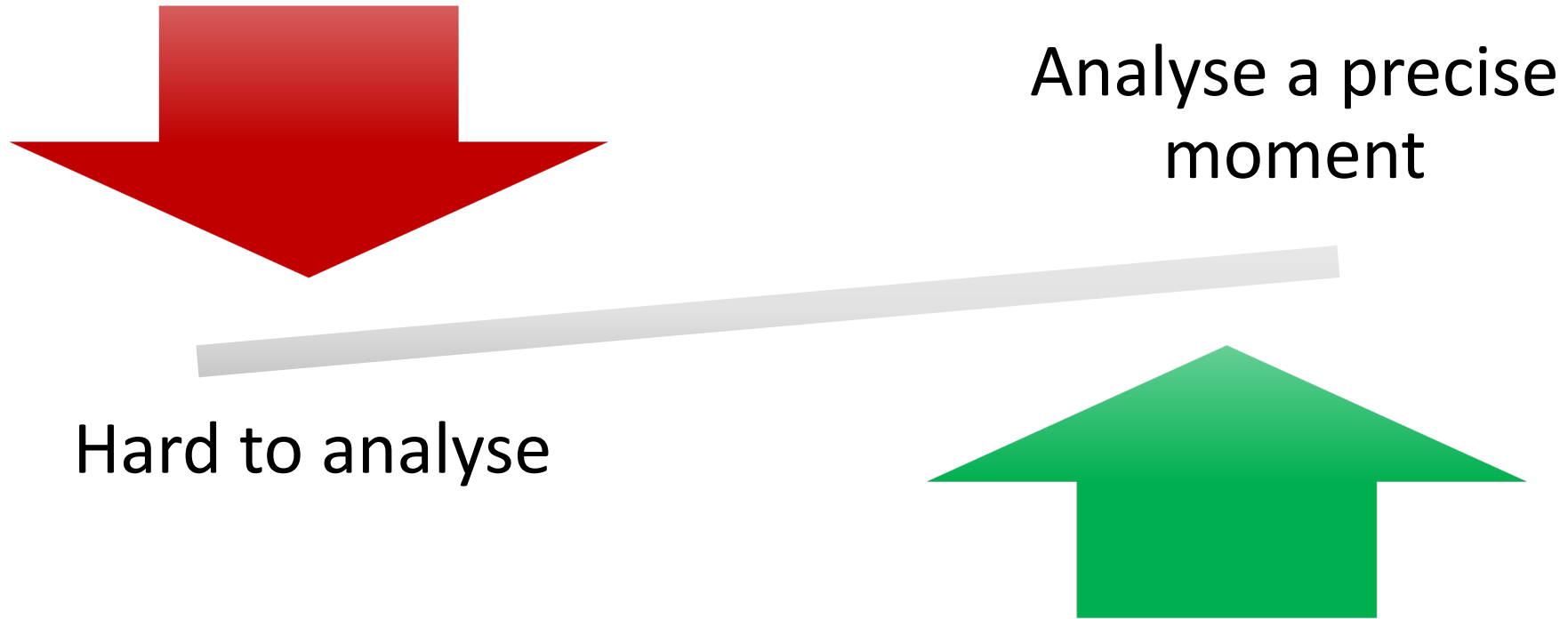


Pleistocene horse

More than 5,000 ancient humans analysed
https://umap.openstreetmap.fr/en/map/ancient-human-dna_41837#5/45.106/17.534



Pros and Cons of aDNA analysis



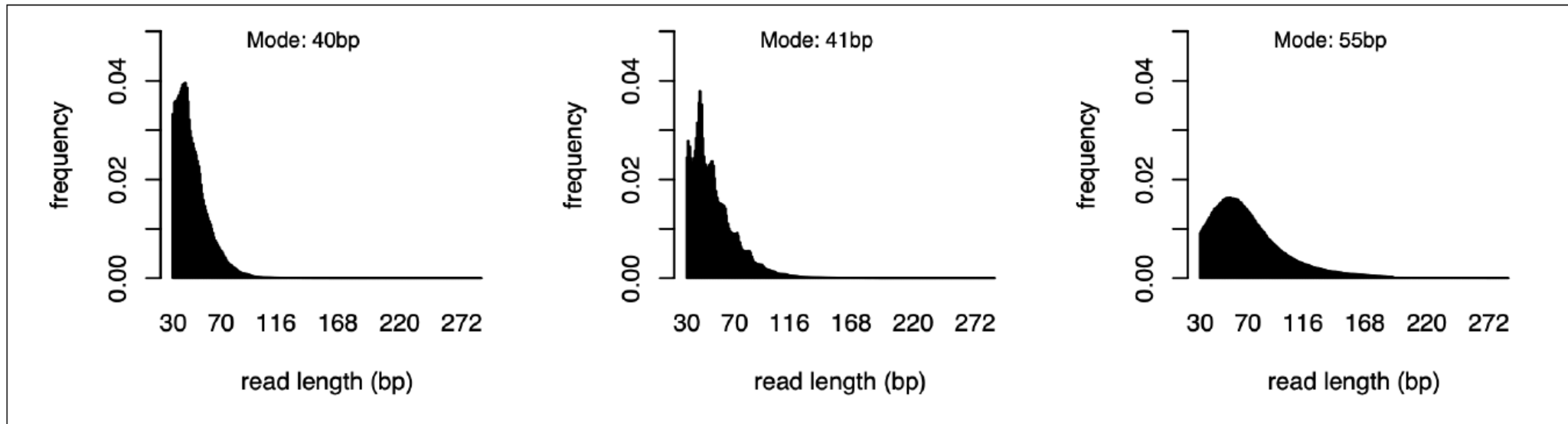
Characteristics of aDNA

Degradation: Fragmentation and post-mortem damage



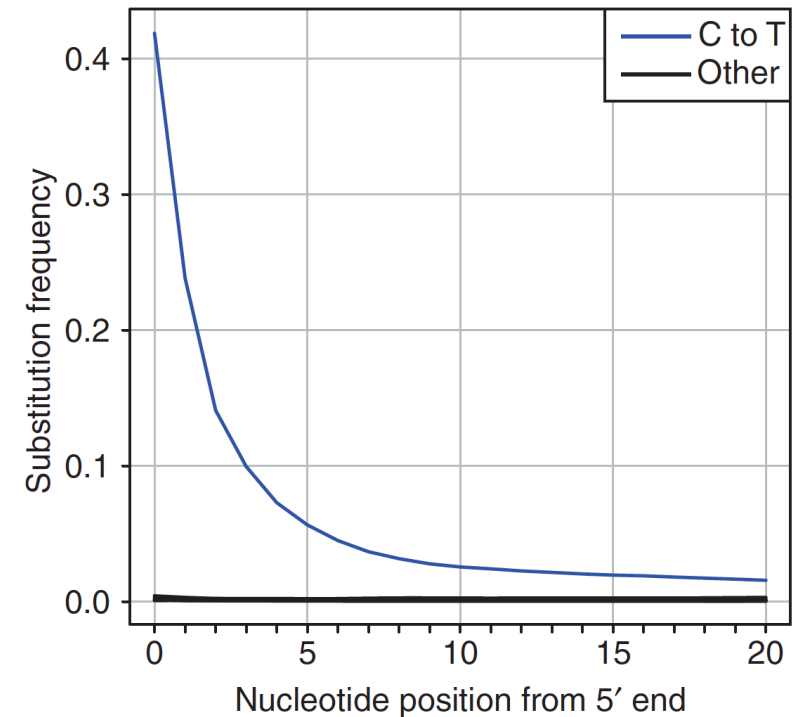
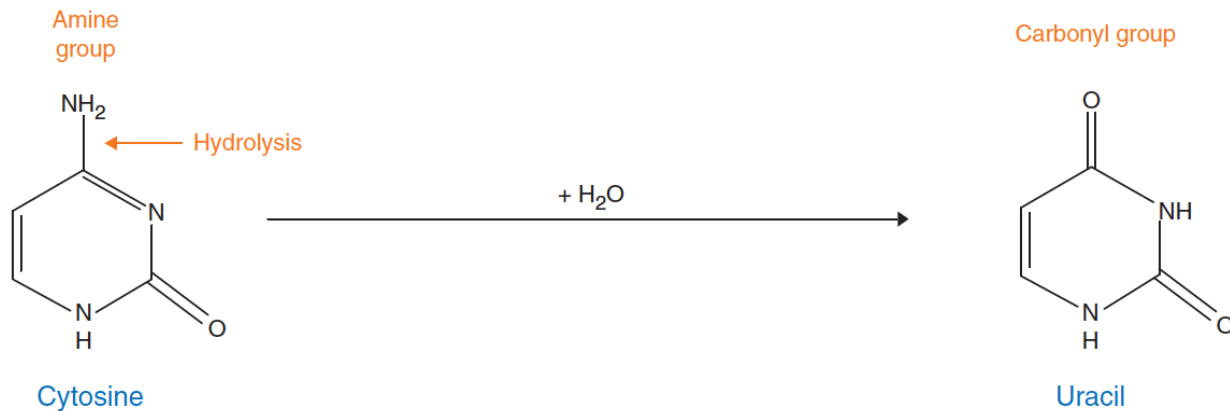
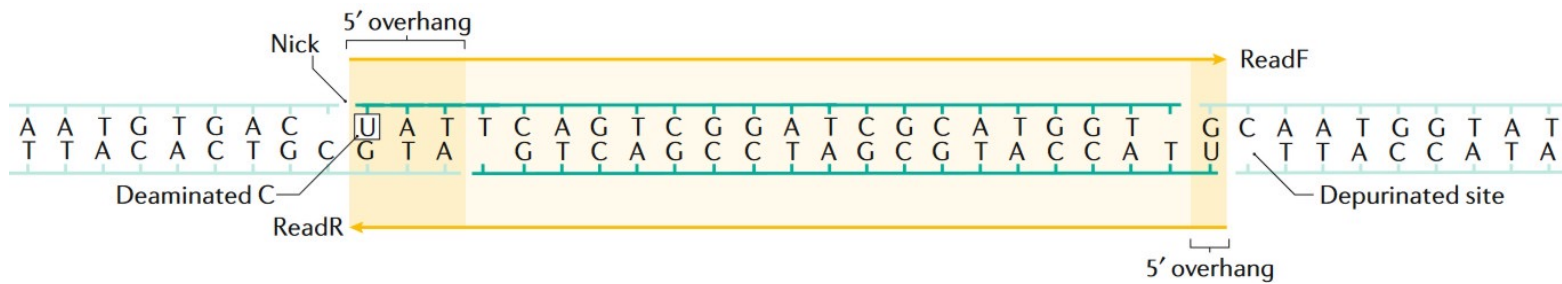
Characteristics of aDNA

Degradation: Fragmentation and post-mortem damage



Characteristics of aDNA

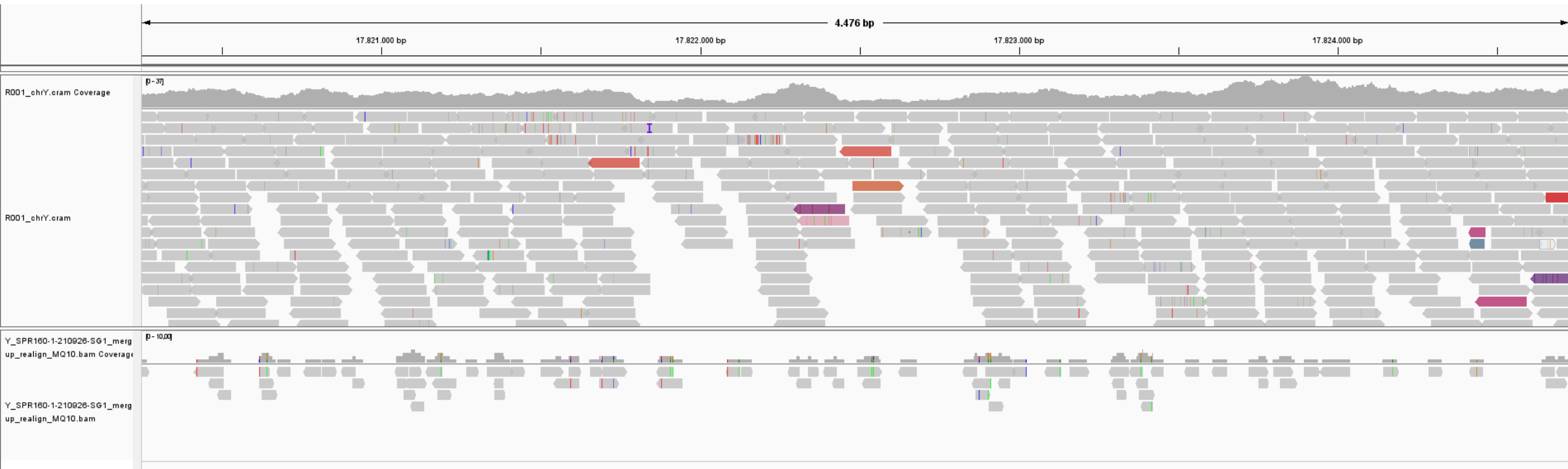
Degradation: Fragmentation and **post-mortem damage**



Characteristics of aDNA

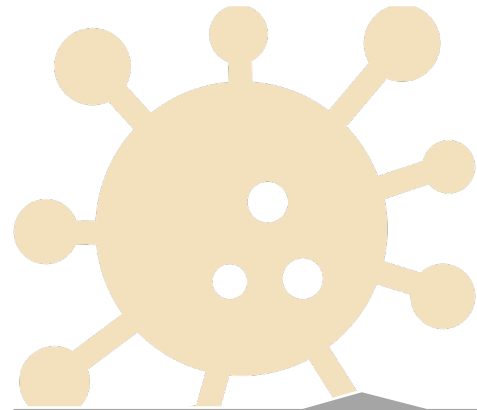
Usually found in low quantities

Resulting in low coverage sequences



Characteristics of aDNA

Potentially contaminated



Environmental DNA

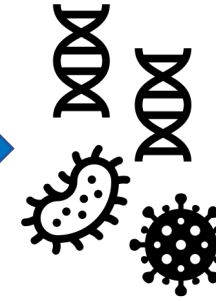
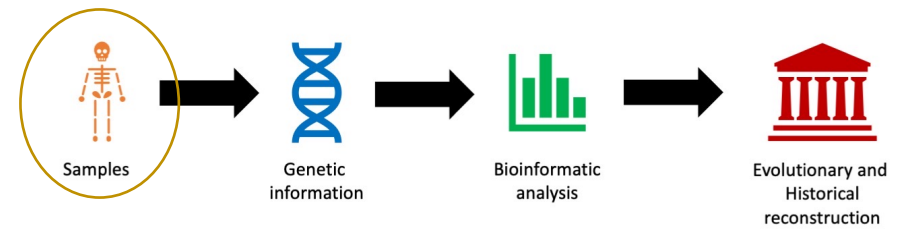
Not mapped to the human
reference sequence



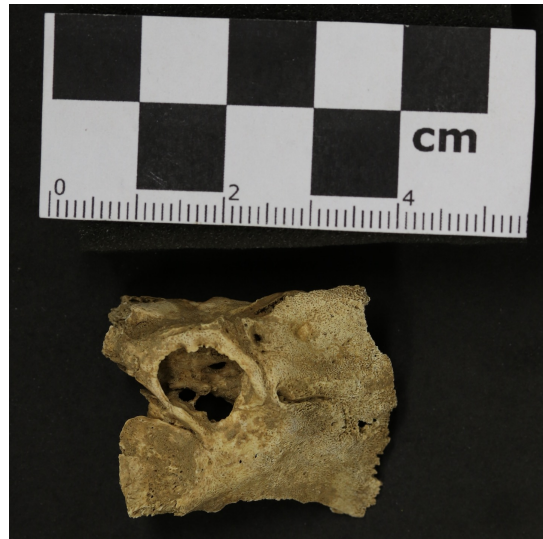
Human DNA

Controlled environment to
prevent contamination

Sample collection



Tooth: relatively less DNA molecules, but greater chance to find ancient pathogens.



Petrous bones: relatively more DNA molecules. Not optimal for ancient pathogen search.

Sample collection

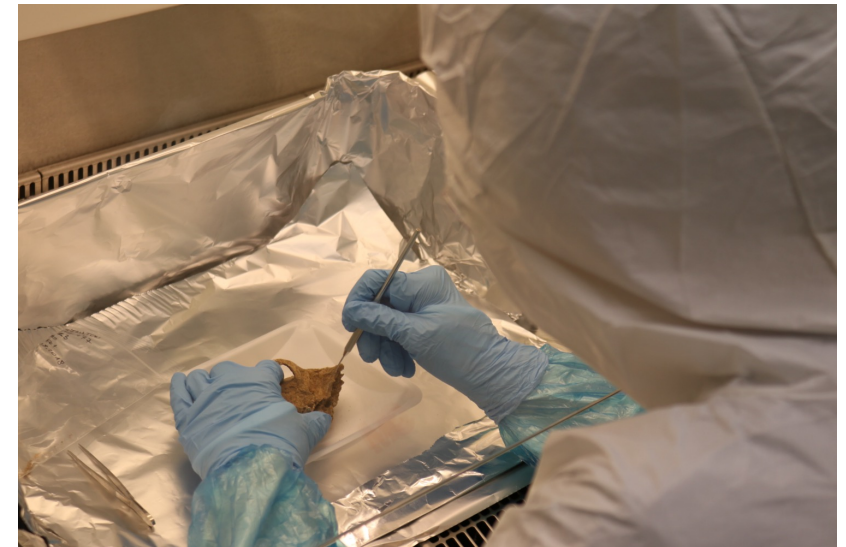
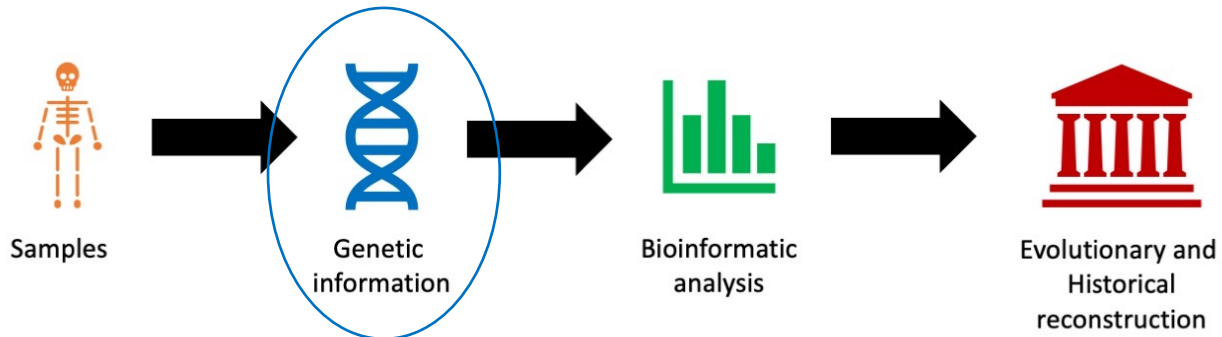
- More endogenous DNA in the petrous bone
- Possibility to recover ancient pathogens from teeth
- How destructive is the method?
- Samples may be used for other analysis



aDNA clean lab

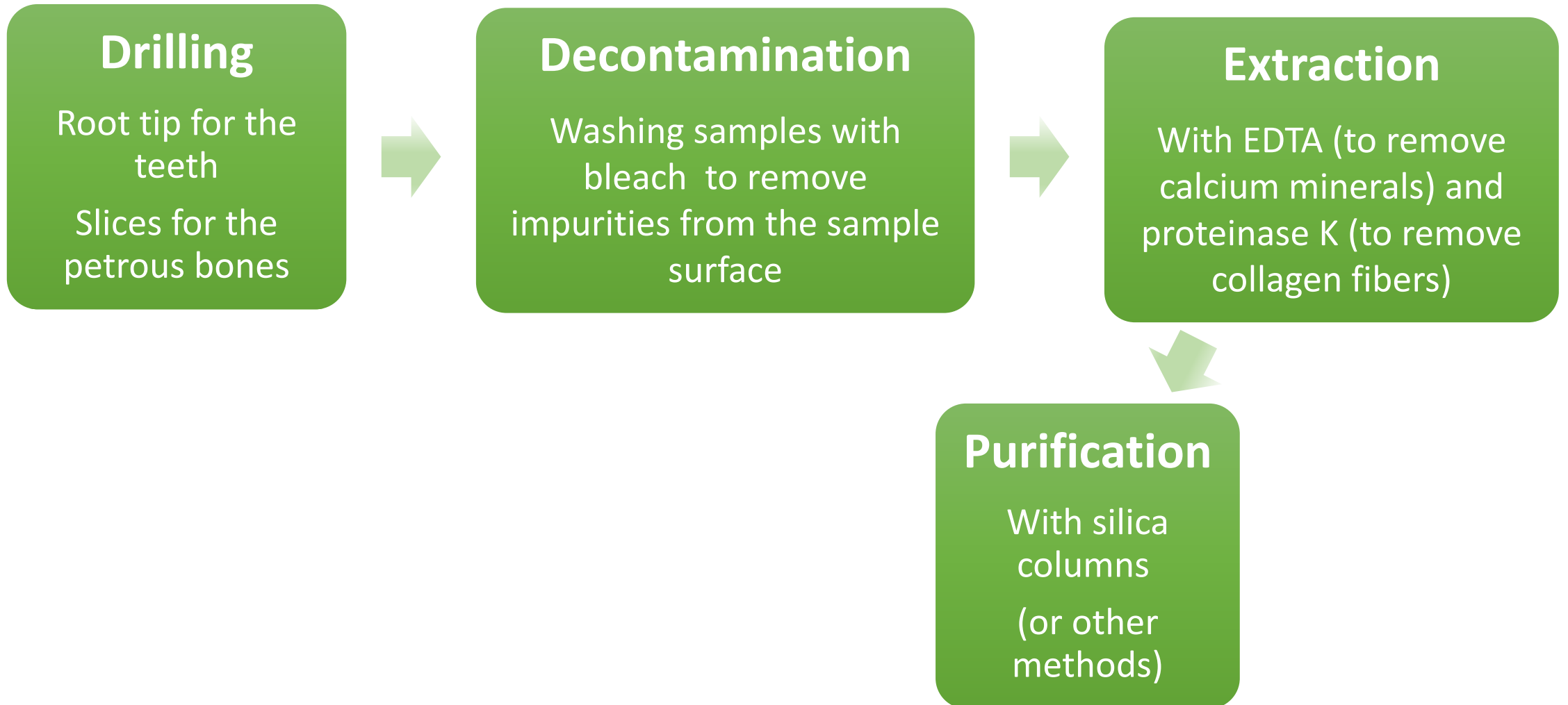
The laboratory is designed to prevent contamination:

- Controlled environment
- Positive pressure
- Filtered air
- UV light (optional)
- No entry without security devices (suits, masks, gloves etc.)
- Compartmentalized laboratory (one room for each operation)
- All objects brought from outside must be cleaned with appropriate products (or bleached)
- Daily cleaning

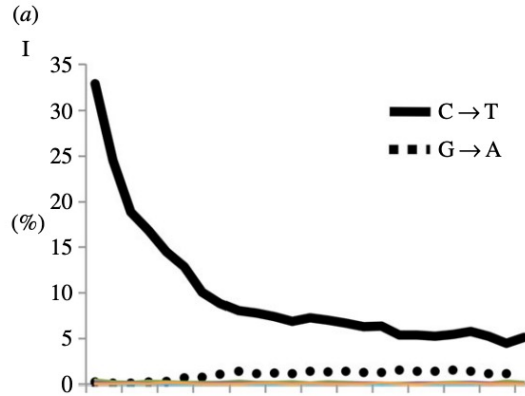


DNA extraction

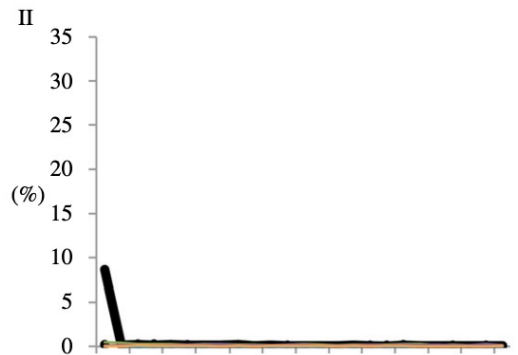
There are several protocols that can be used to extract DNA from bones and teeth



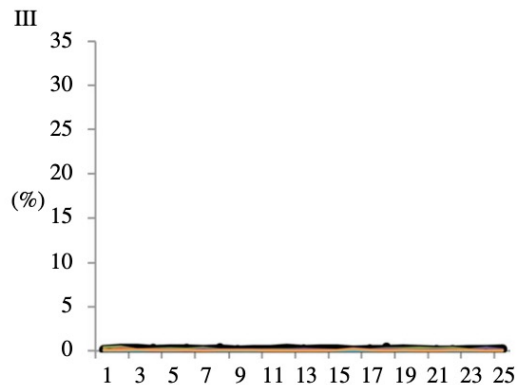
UDG treatment (optional)



No UDG treatment



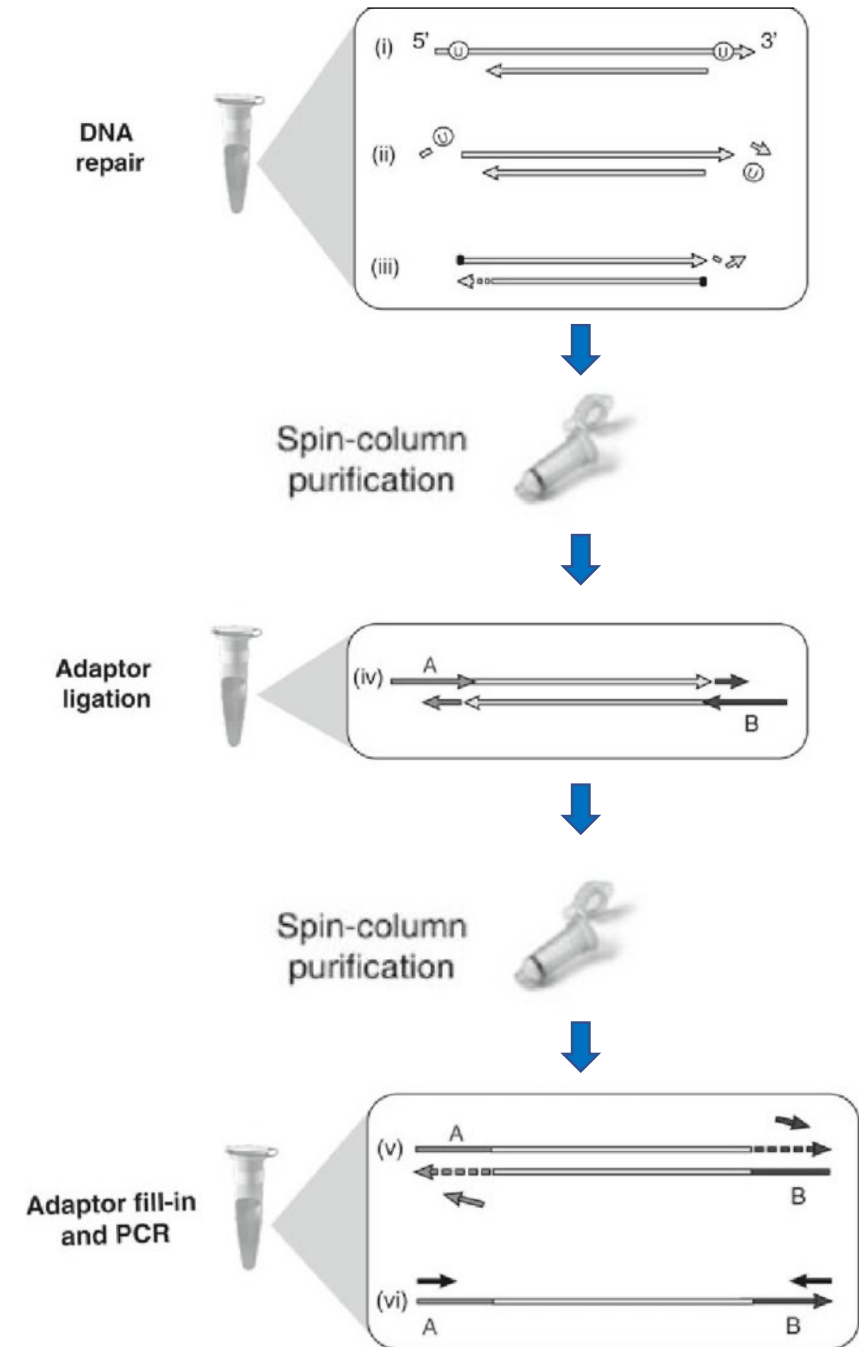
Partial UDG treatment



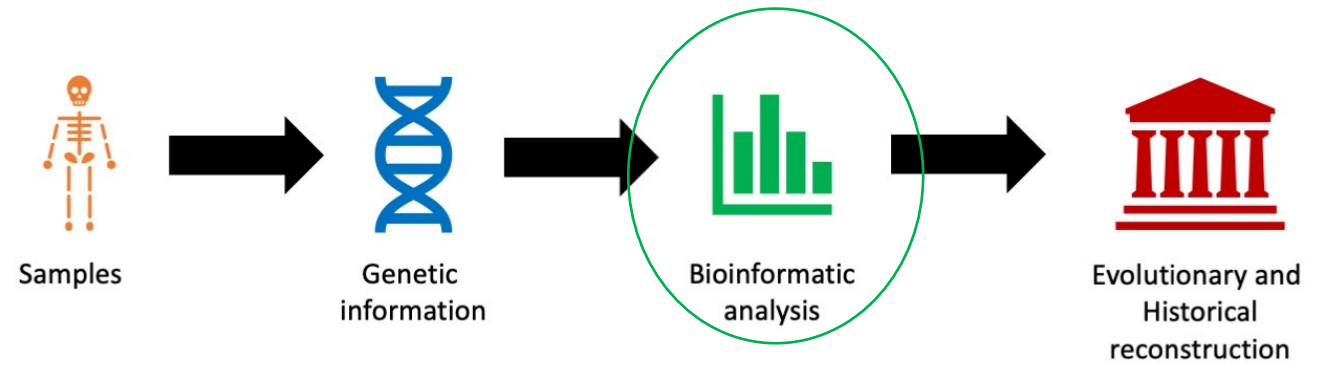
Full UDG treatment

Library preparation

- Fragmentation (not needed in aDNA)
- DNA molecule end repair
- Adaptor ligation (with indexes for the sequencing)
- Adaptor fill-in
- PCR (outside the clean lab)



Sequencing



NGS Sequencing:

- Whole Genome Sequencing
- SNP capture
- Output: fastq files

A sequence identifier with information about the read

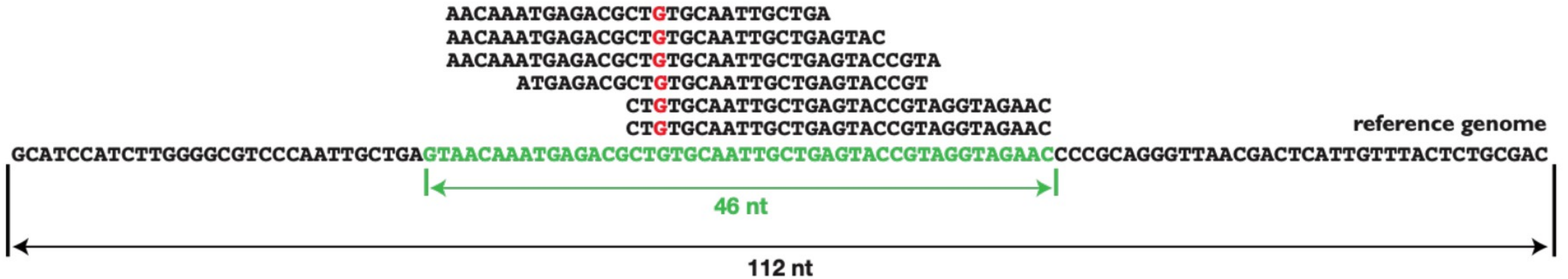
The sequence

```
@NB501163:37:HGK7WBGX7:1:11101:20397:1053 1:N:0:CGTACTAG+TACTCCTT  
AGATCNGAAGAGCACACGTCTGAACTCCAGTCACCGTACTAGATCTCGTATGCCGTCTTCTGCTTGAAAAAAAAA  
+  
AAAAA#/EEE/6EEEEEEAEE6AEEAEE/EEA/EEEEEEAEEEA//EAEEEEEEAEEEA/EEEE/AEEEE/E//
```

Base call quality scores

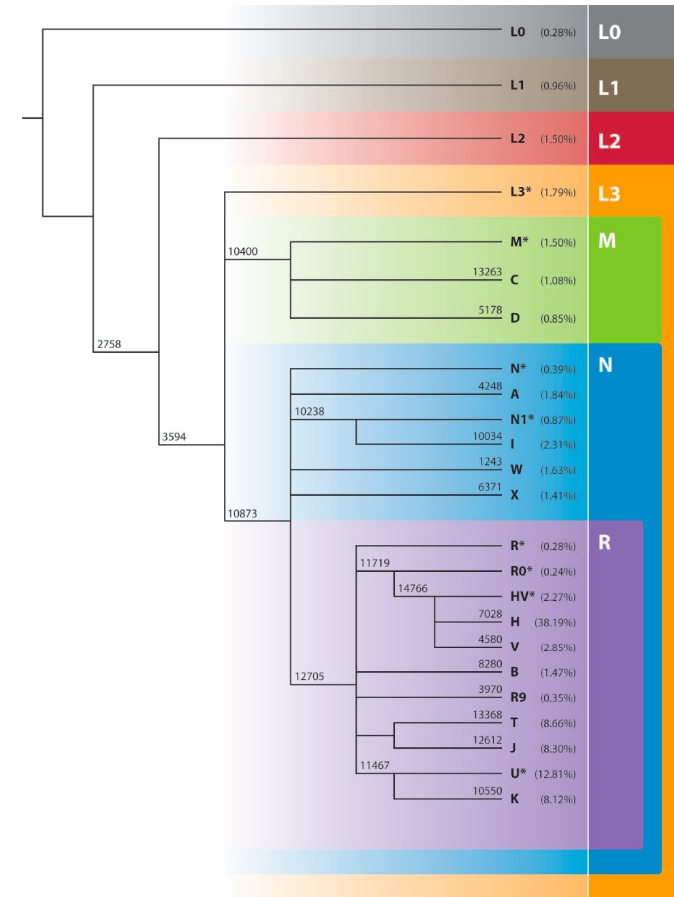
Mapping

Mapping sequencing reads (from fastq files) to the reference genome



Authentication

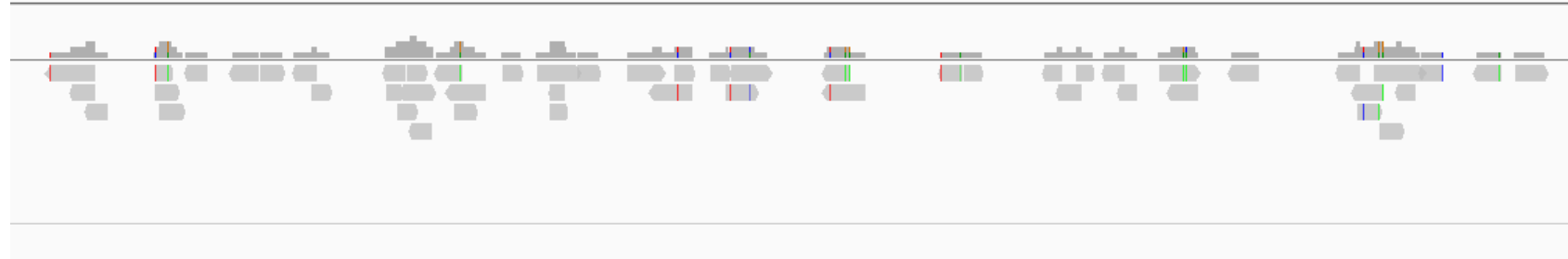
- **Amount of endogenous DNA** (mapped/unmapped reads ratio)
- **Ancient or modern DNA**
 - Read length
 - aDNA damage
- **Contamination**
 - X-based method (only for male samples)
 - mtDNA method (Calculating the percentage of non-consensus bases at haplogroup-defining positions)



Variant calling

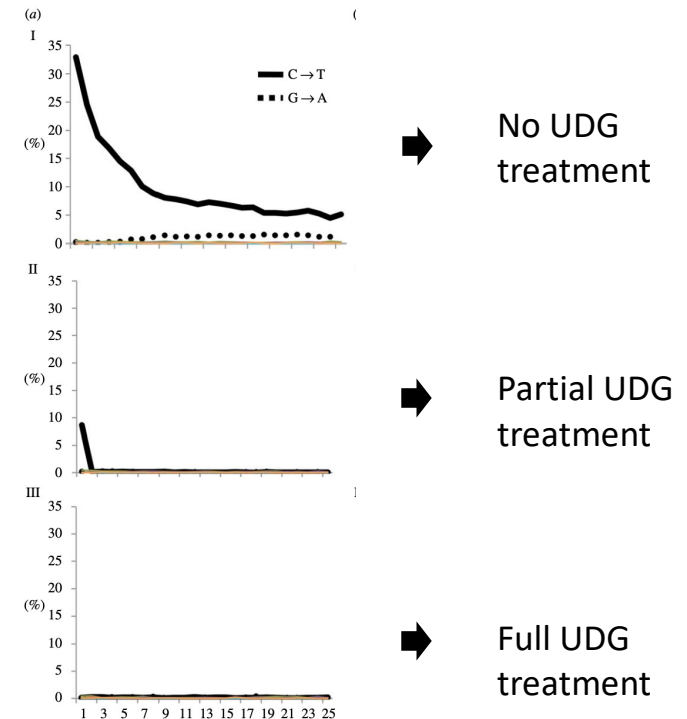
Variant type:

- Genotypes
- Pseudo-haploid genotype
- Genotype likelihoods

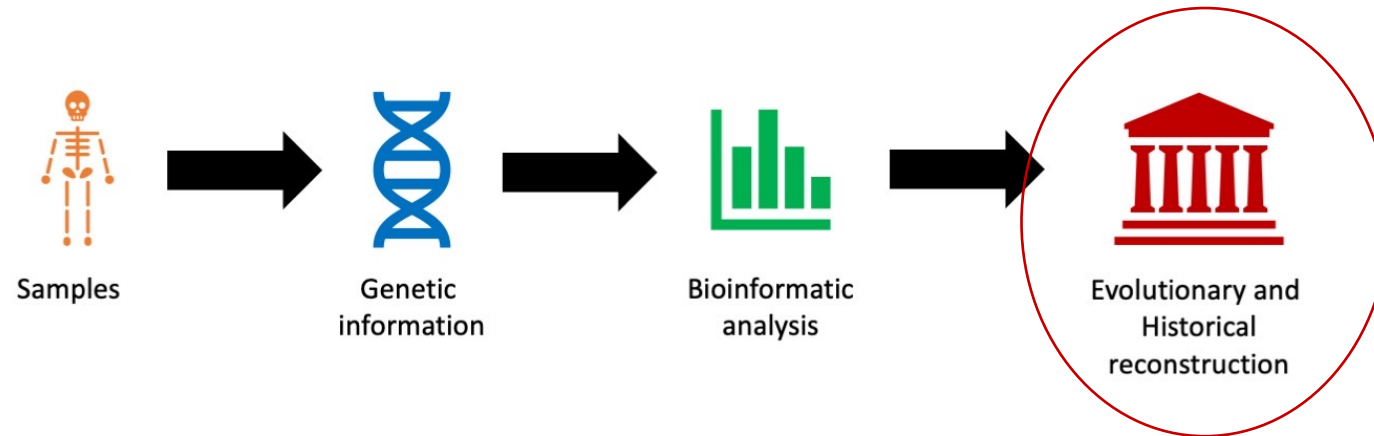


Deal with post-mortem damage:

- Trim reads for partially UDG-treated samples
- Remove transitions (C \leftrightarrow T, G \leftrightarrow A)
- Likelihood methods

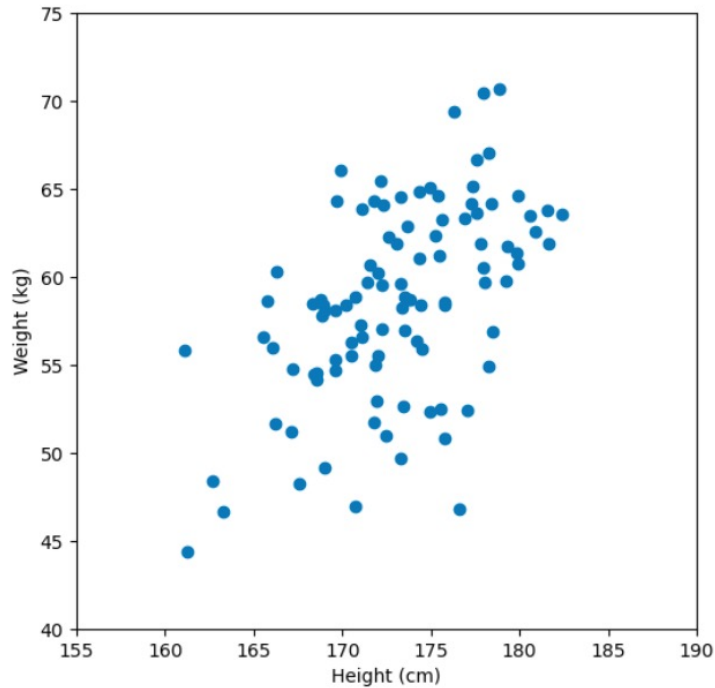


Population genetics analysis for aDNA data

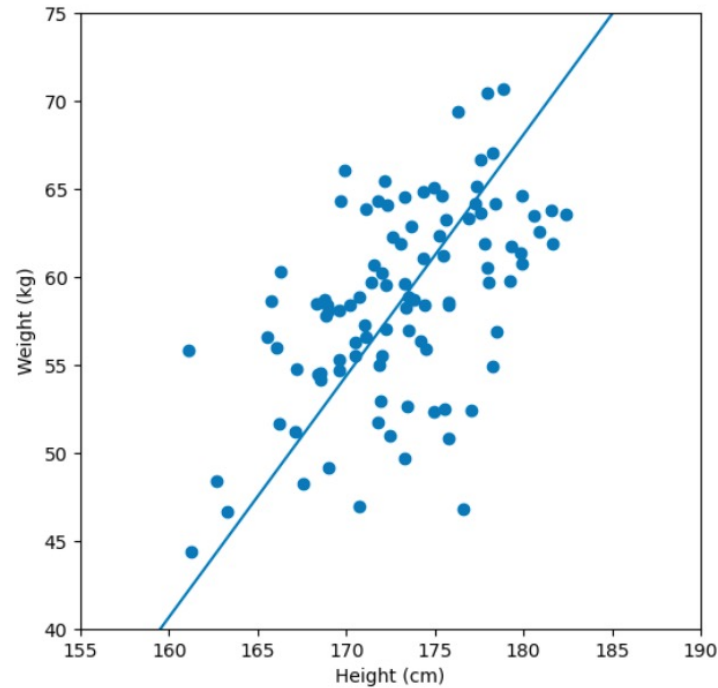


Principal Component Analysis (PCA)

- PCA is a linear transformation to a new coordinate system
- **Reduction of dimensions:** the genetic information contained in 1M SNPs can be summarized by a few new variables



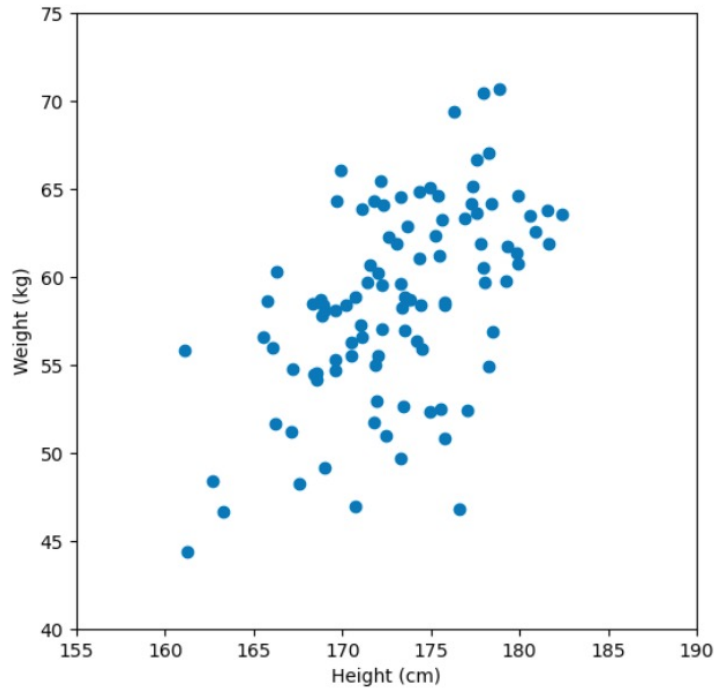
Each individual (point) is represented by two variables.



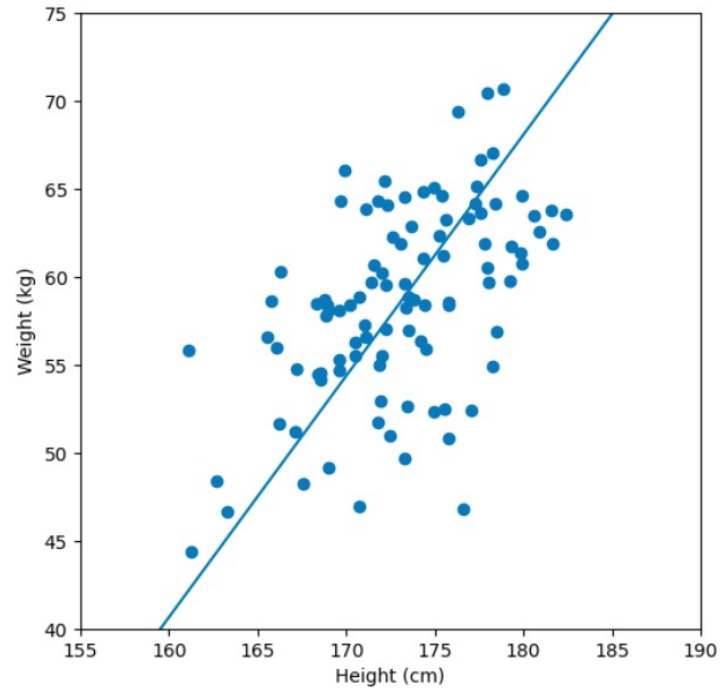
Find the axis of greatest variation (fit line) → The principal component.

Principal Component Analysis (PCA)

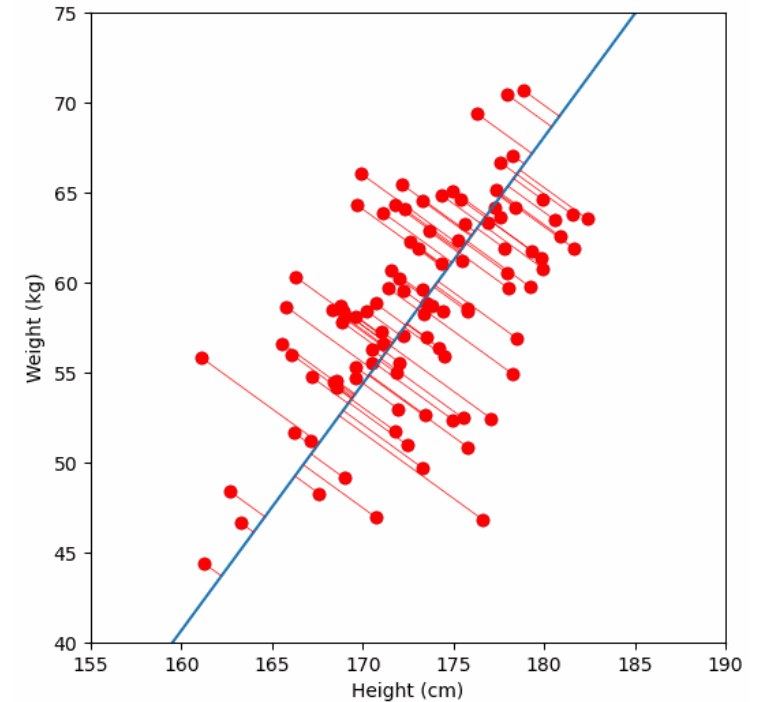
- PCA is a linear transformation to a new coordinate system
- **Reduction of dimensions:** the genetic information contained in 1M SNPs can be summarized by a few new variables



Each individual (point) is represented by two variables.



Find the axis of greatest variation (fit line) → The principal component.



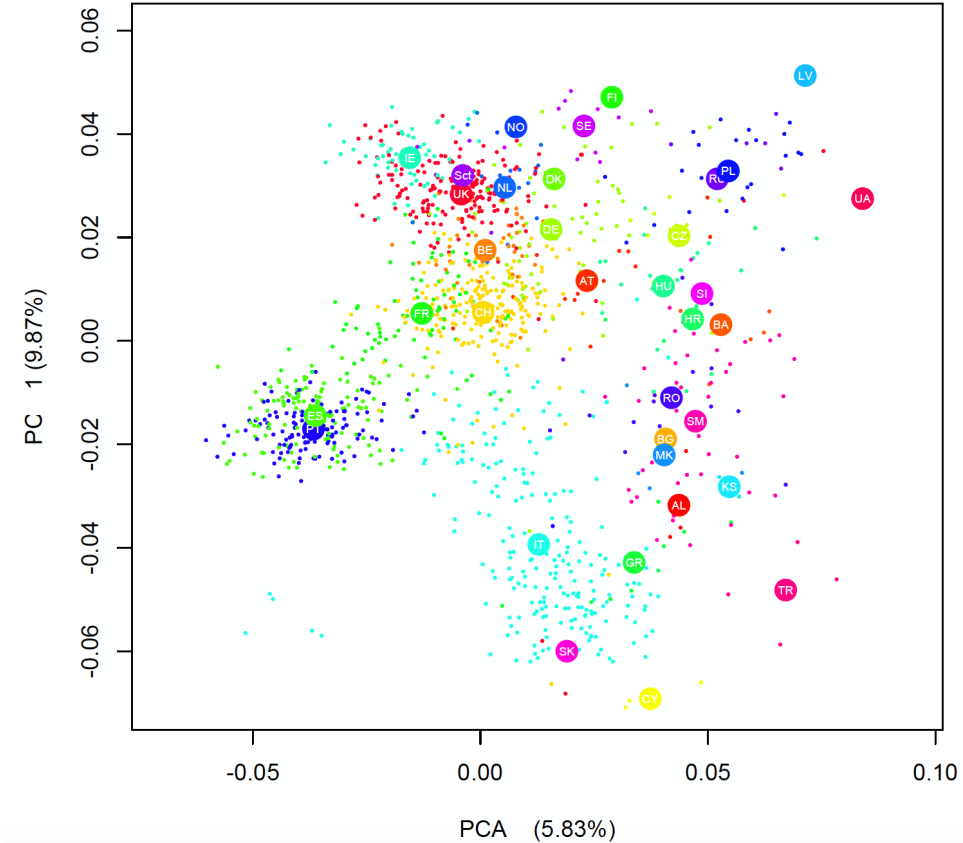
“Project” each point onto the line. Now each individual is represented by one variable.

Principal Component Analysis (PCA)

Ind(1): 0101110110101110
Ind(2): 0111110110101111
Ind(3): 0100110110101011
Ind(4): 0111111111101111
Ind(5): 0101110110100001
.
.
.
Ind(n): 0101110110101111



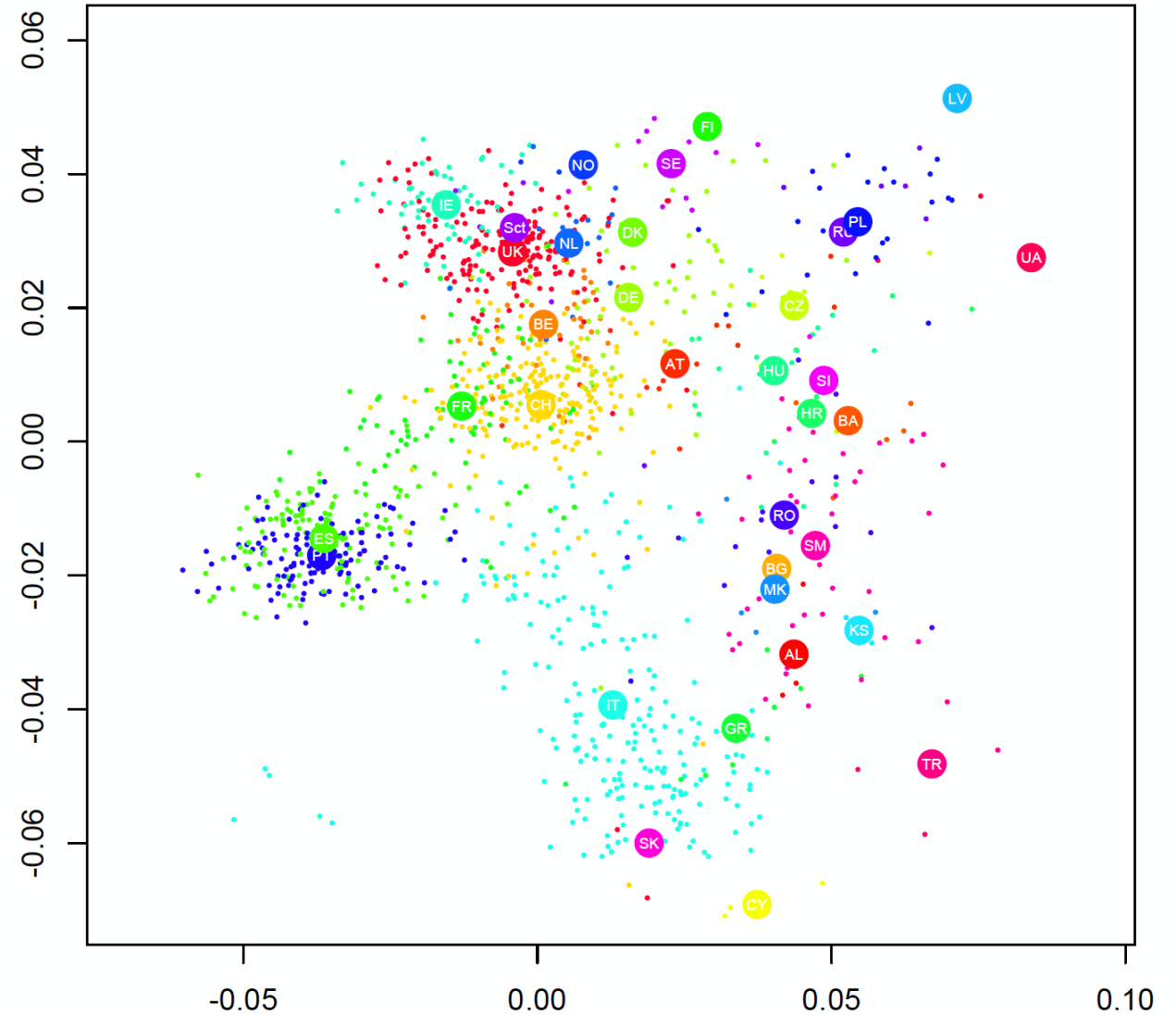
	PC1	PC2
	0.01	-0.02
	0.50	0.03
	0.07	-0.13
	0.02	-0.04
	0.01	-0.05
.	.	.
.	.	.
.	.	.
	-0.03	0.03



Principal Component Analysis (PCA)

Novembre et al. 2009

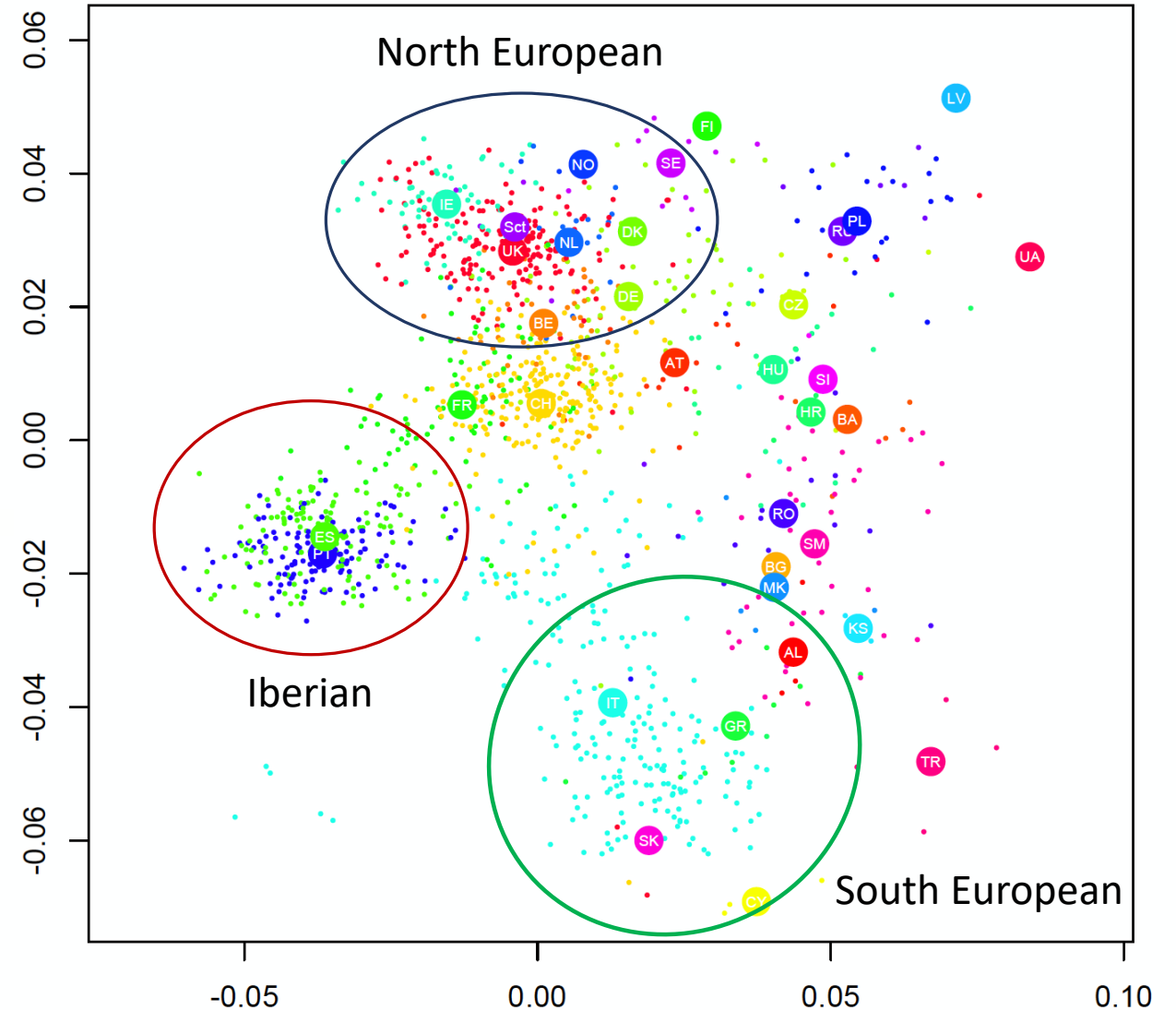
- PCA reveal population structure
- Genetic Distance \approx Physical distance
- Easily identify genetic outliers and isolated populations



Principal Component Analysis (PCA)

Novembre et al. 2009

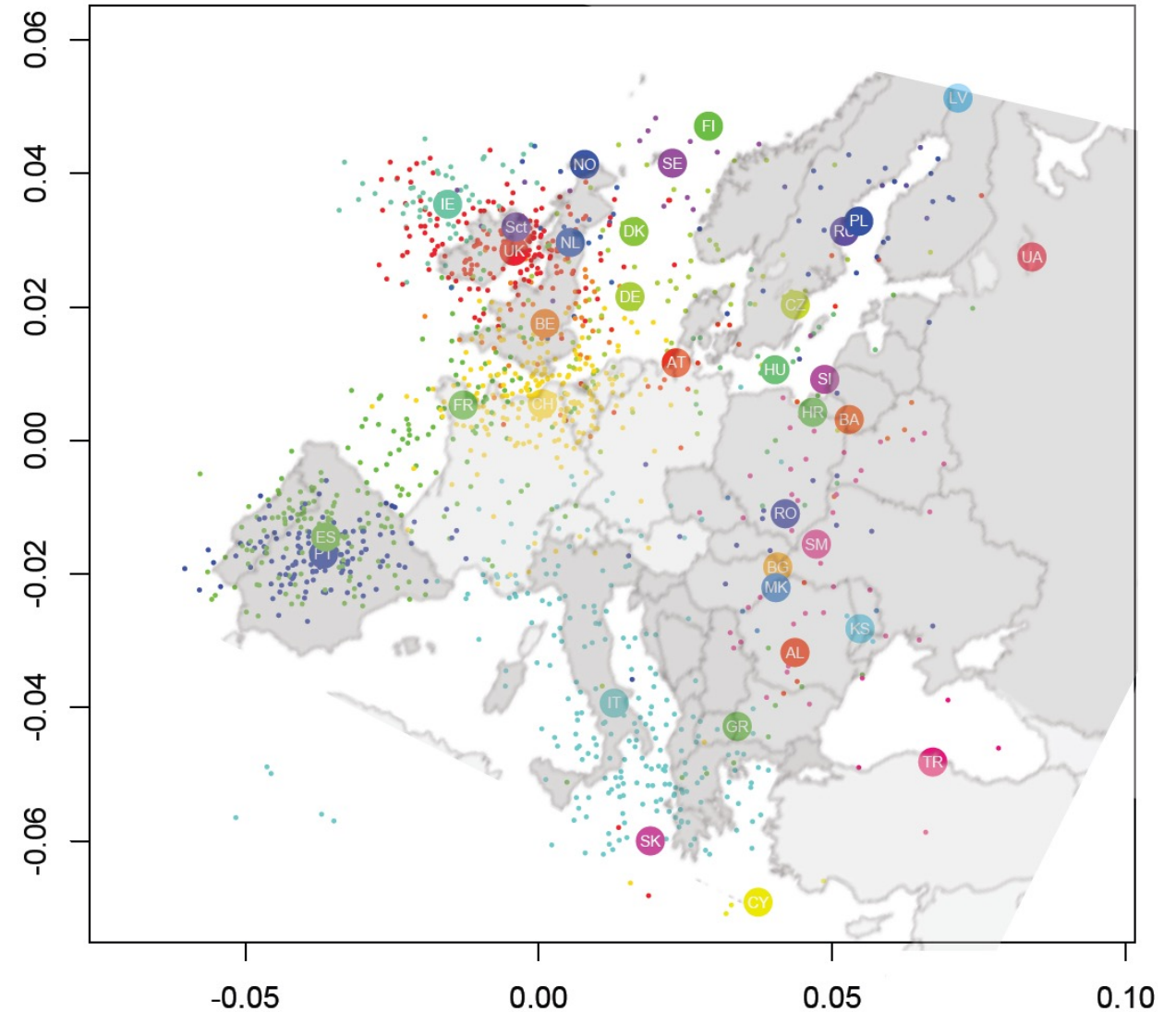
- PCA reveal population structure
- Genetic Distance \approx Physical distance
- Easily identify genetic outliers and isolated populations



Principal Component Analysis (PCA)

Novembre et al. 2009

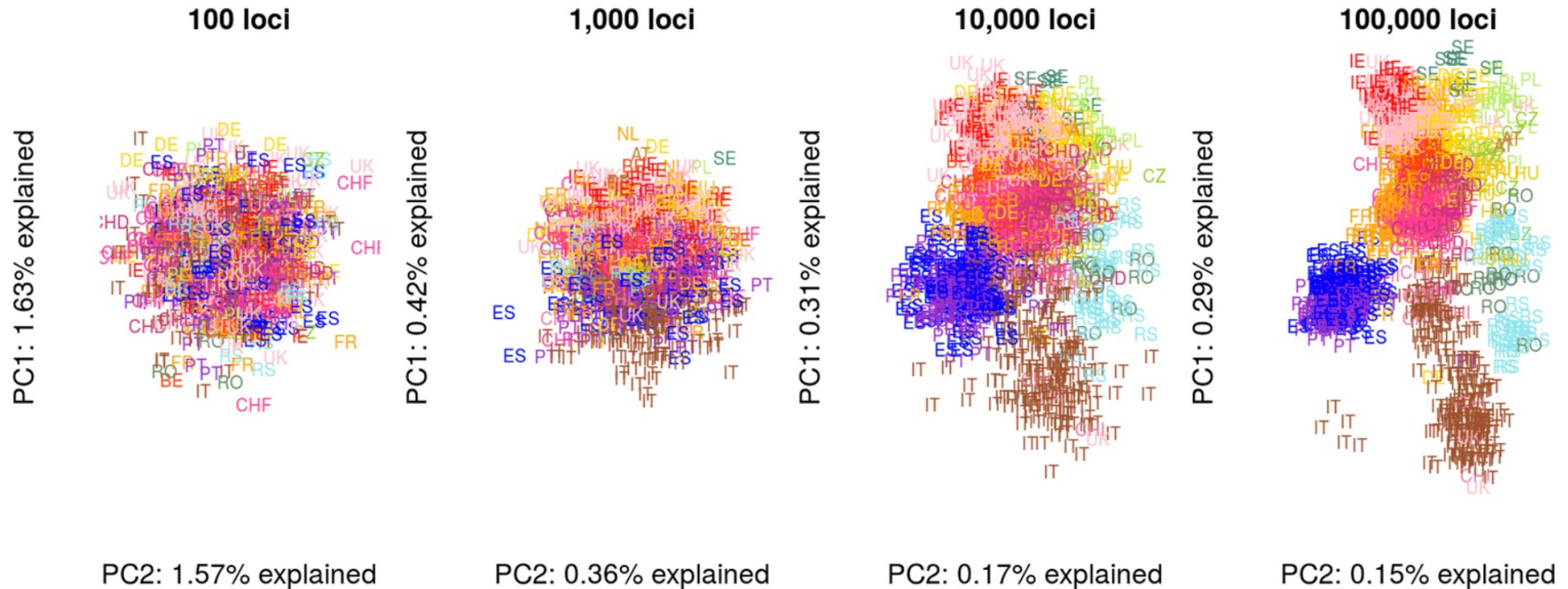
- PCA reveal population structure
- Genetic Distance \approx Physical distance
- Easily identify genetic outliers and isolated populations



Principal Component Analysis (PCA)

Produce good results even when the information is low

Novembre et al. 2009

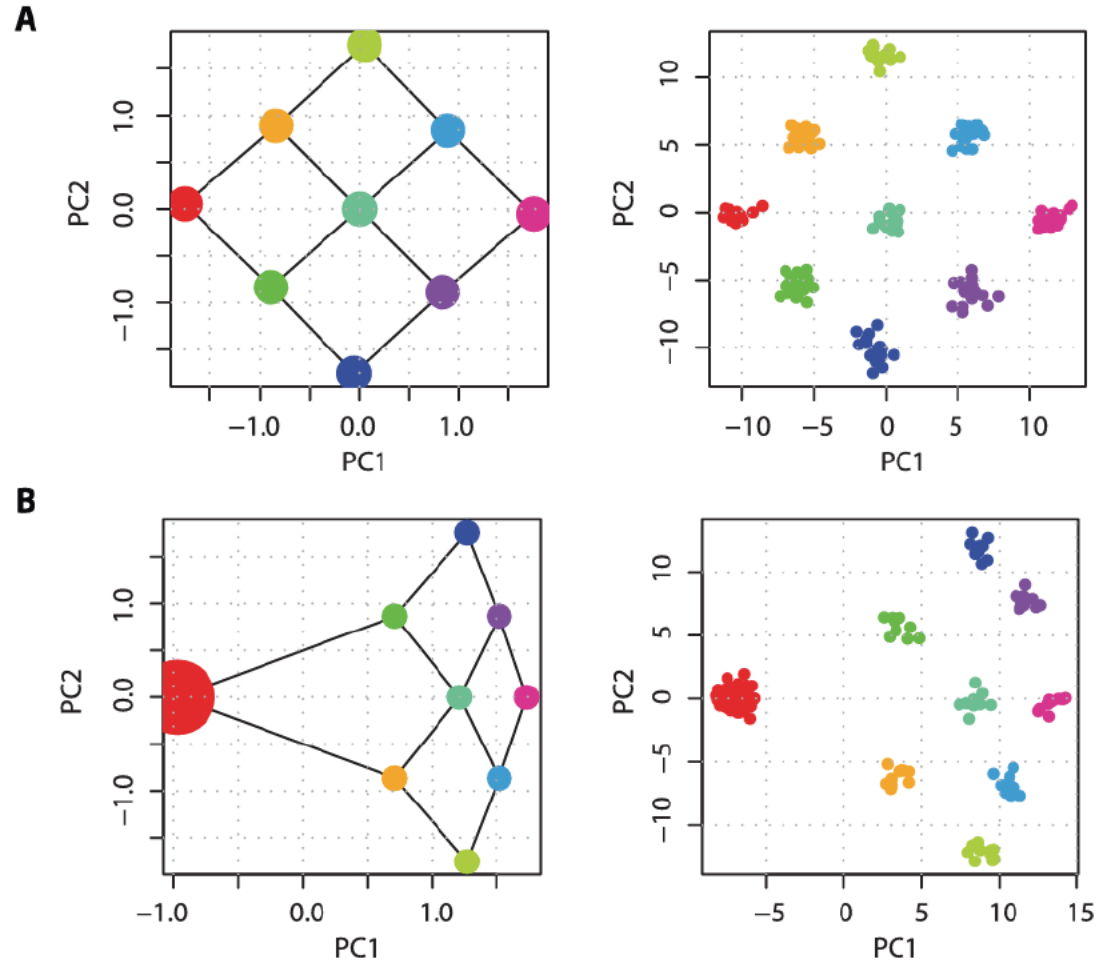


Principal Component Analysis (PCA)

Factors that influence PCA:

- Migration
- Genetic drift
- Admixture

- Population size
- SNP selection



PCA with ancient samples

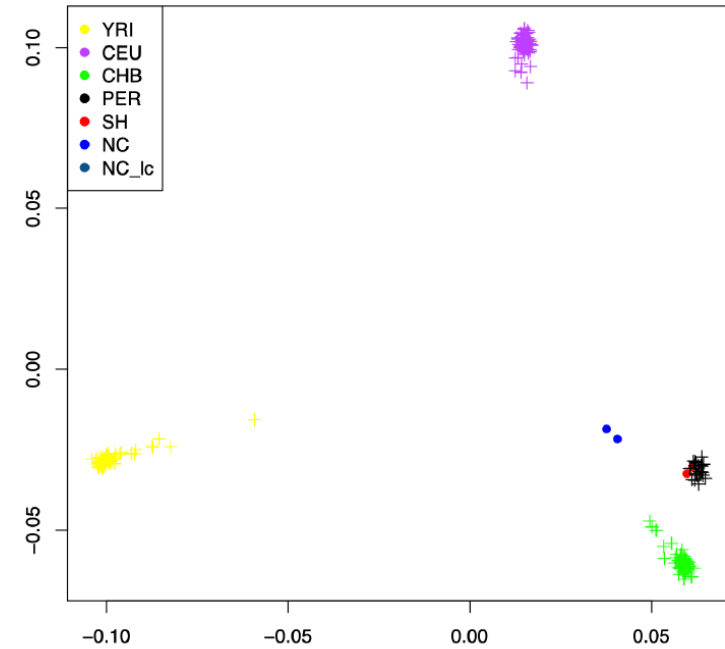
Low coverage individuals result in many SNPs with missing data

Usually, PCA methods will fill in all missing data. This results in PCA plots that have ancient individuals near/at the origin (0,0 coordinate).

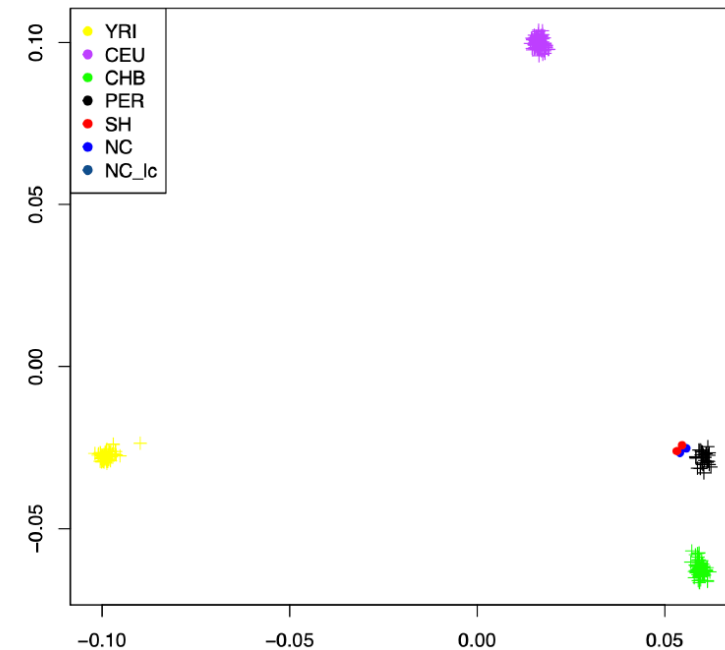
Solution: Projection of ancient individuals.

We can infer eigenvectors using the reference set and then project ancient individuals onto those eigenvectors.

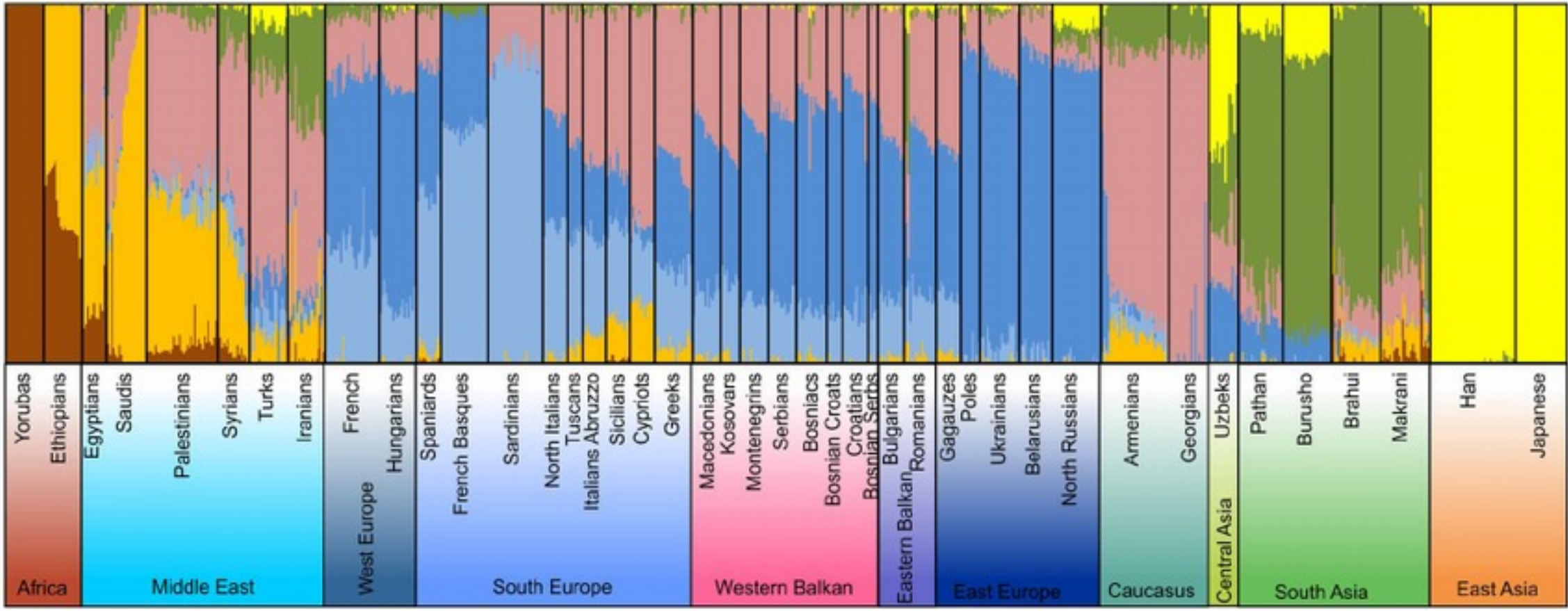
Not projected



Projected



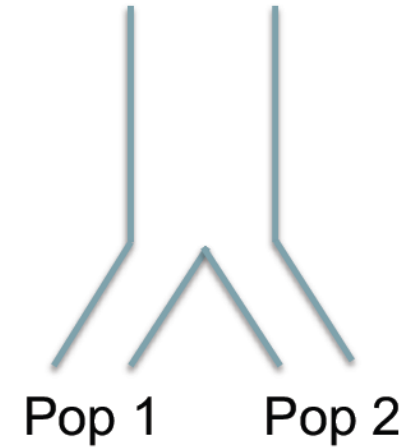
Ancestry proportion inference (ADMIXTURE)



Allele frequency-based clustering

Thought experiment:

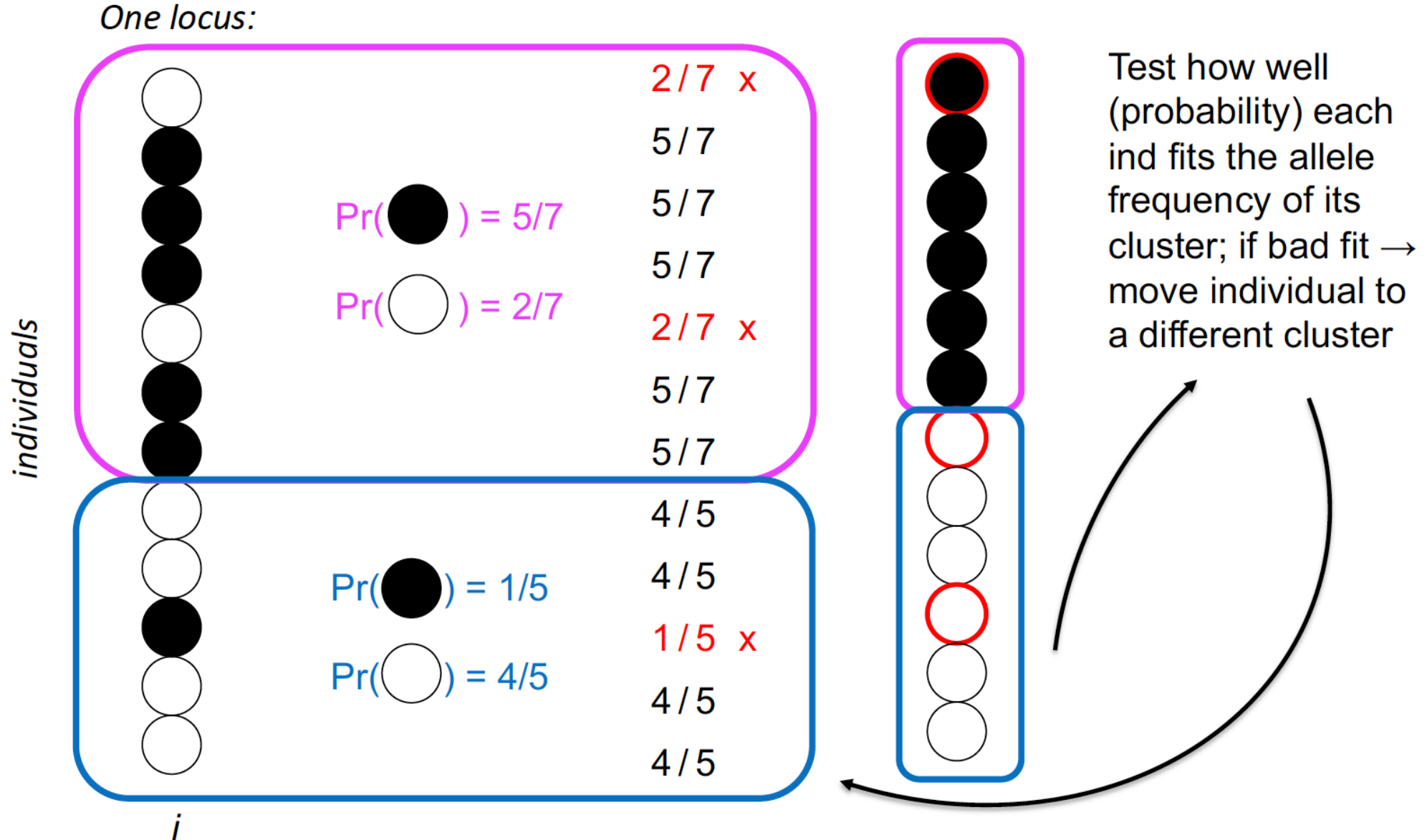
- Assume we sequenced individuals we knew are from Pop 1 and Pop 2
- We now sequence another individual, where we are unsure whether they are from Pop1 or Pop 2
- How could we try to assign this individual to Pop1 or Pop 2?



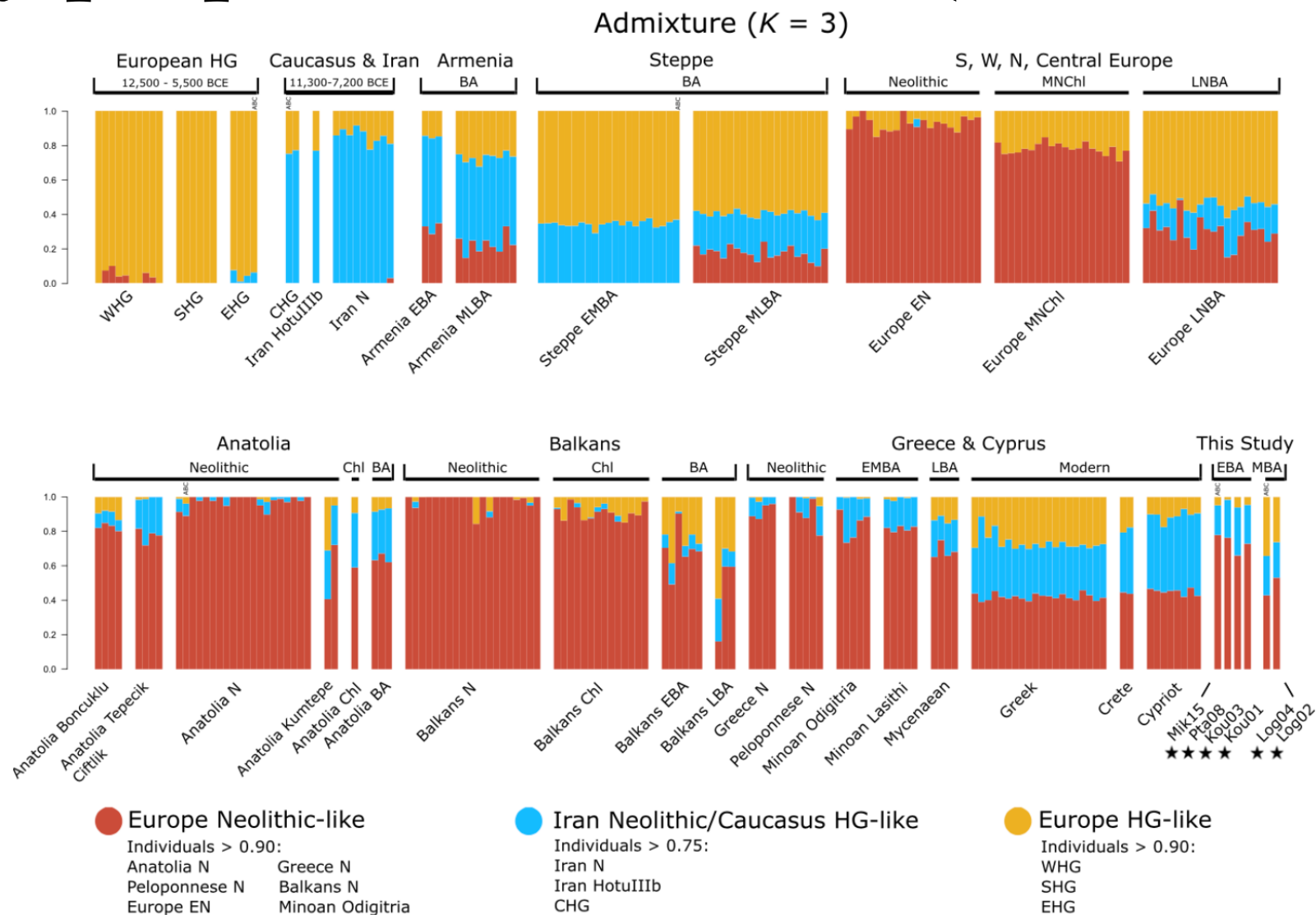
	Allele frequency in Pop1	Pop2	Genotype of individual	
SNP1	0.8	0.4	1	→ Pop1?
SNP2	0.3	0.7	0	→ Pop1?
SNP3	0.4	0.6	1	→ Pop2? (Pop1?)
SNP4	0.9	0.1	1	→ Pop1?

Allele frequency-based clustering

We can do the same **without** knowing allele frequencies in Pop1 and Pop2 by clustering

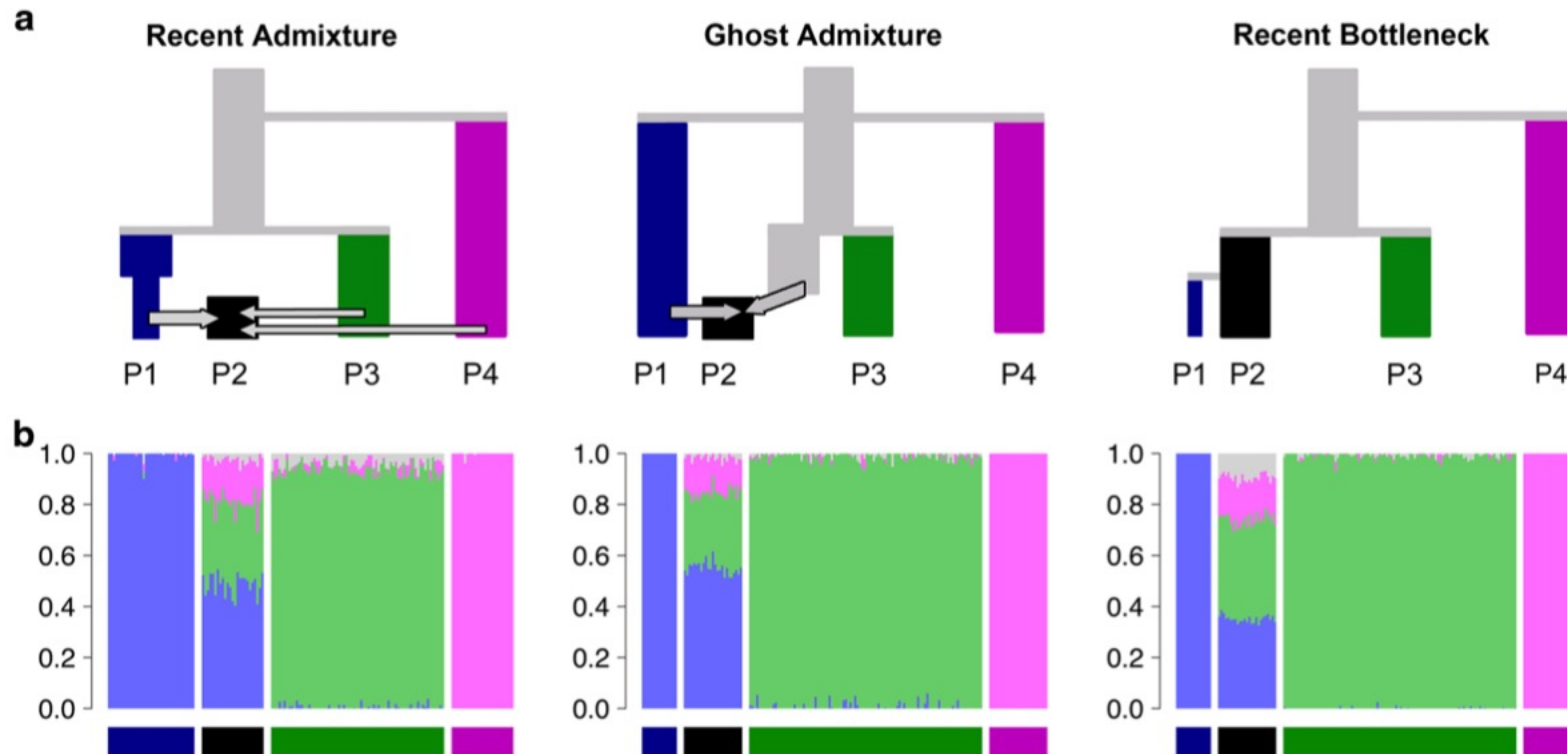


Ancestry proportion inference (ADMIXTURE)



Ancestry proportion inference (ADMIXTURE)

Be careful when interpreting ADMIXTURE results!



Ancestry proportion inference (ADMIXTURE)

Be careful when interpreting ADMIXTURE results!



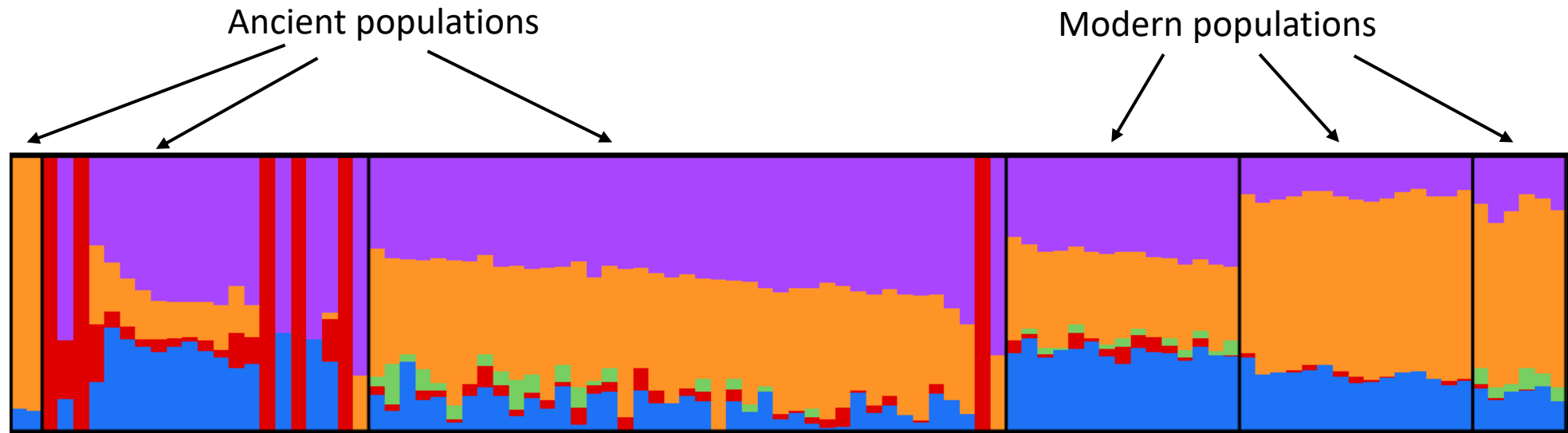
In this case, clustering will be the same as that for discrete populations



Ancestry proportion inference (ADMIXTURE)

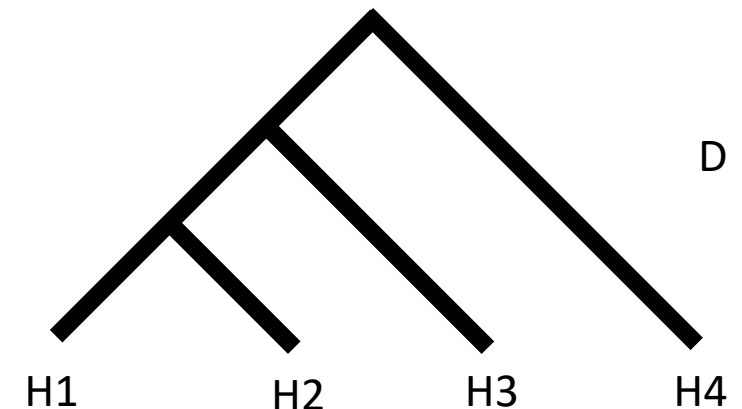
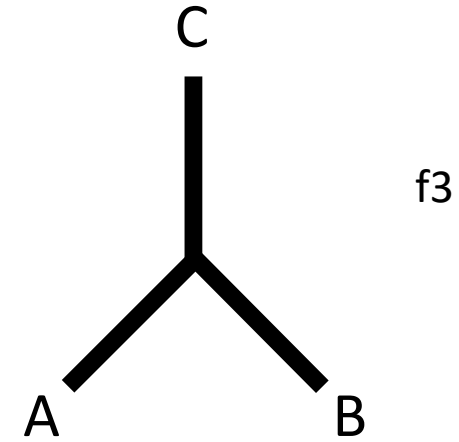
Be careful when interpreting ADMIXTURE results!

Possible problem with low coverage samples



Tests of “treeness” – f and Patterson’s D statistics

- Testing if a tree of population is correct
- Identify admixture and gene flow
- Simple to analyse
- Results (relatively) easy to interpret
- Statistically robust even with a small number of loci
- Ideal for aDNA data!

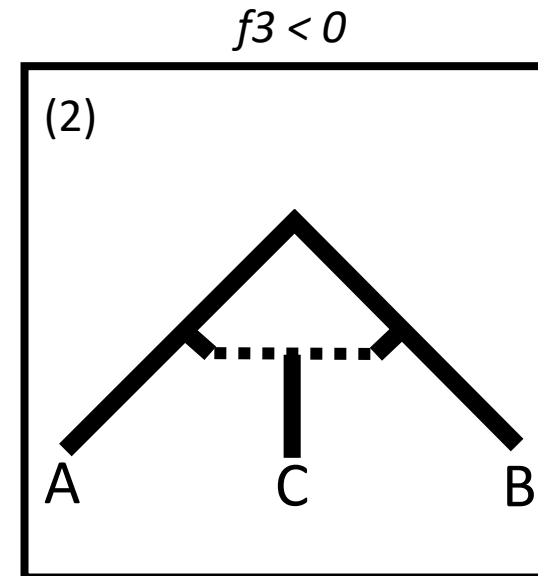
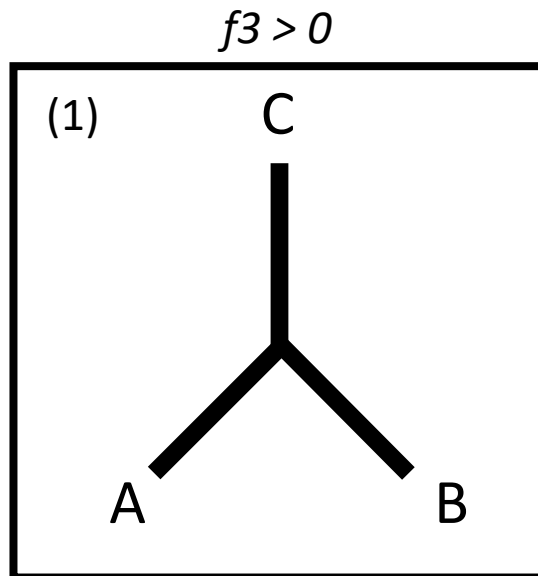


f_3 statistic

$$f_3(C; A, B) = \frac{1}{J} \sum_{j=1}^J (c_j - a_j)(c_j - b_j)$$

Two main purposes:

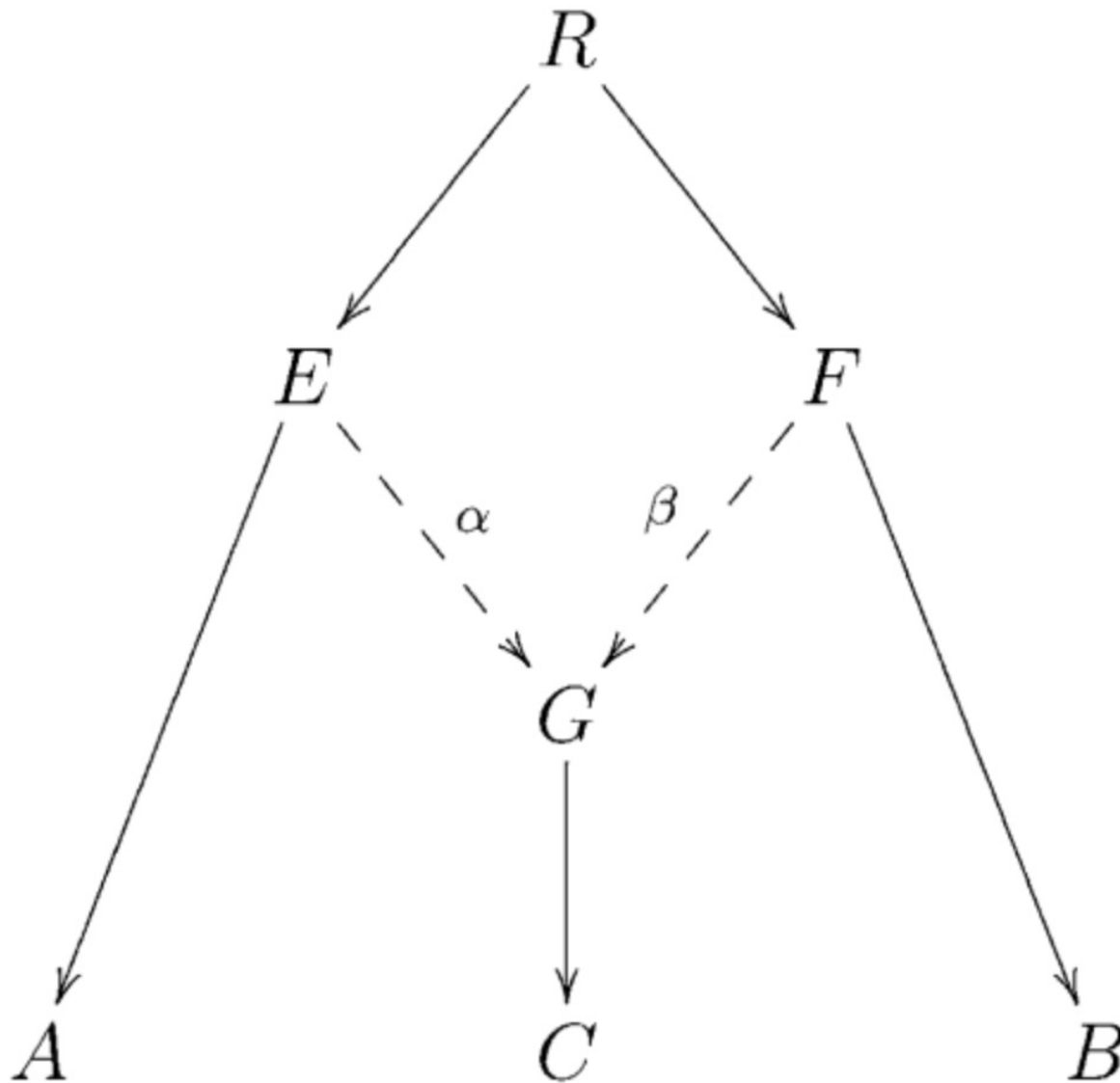
- Measuring how much two populations are similar with respect to an outgroup (1)
- Testing if a population is the result of an admixture between the other two populations (2)



f_3 statistic

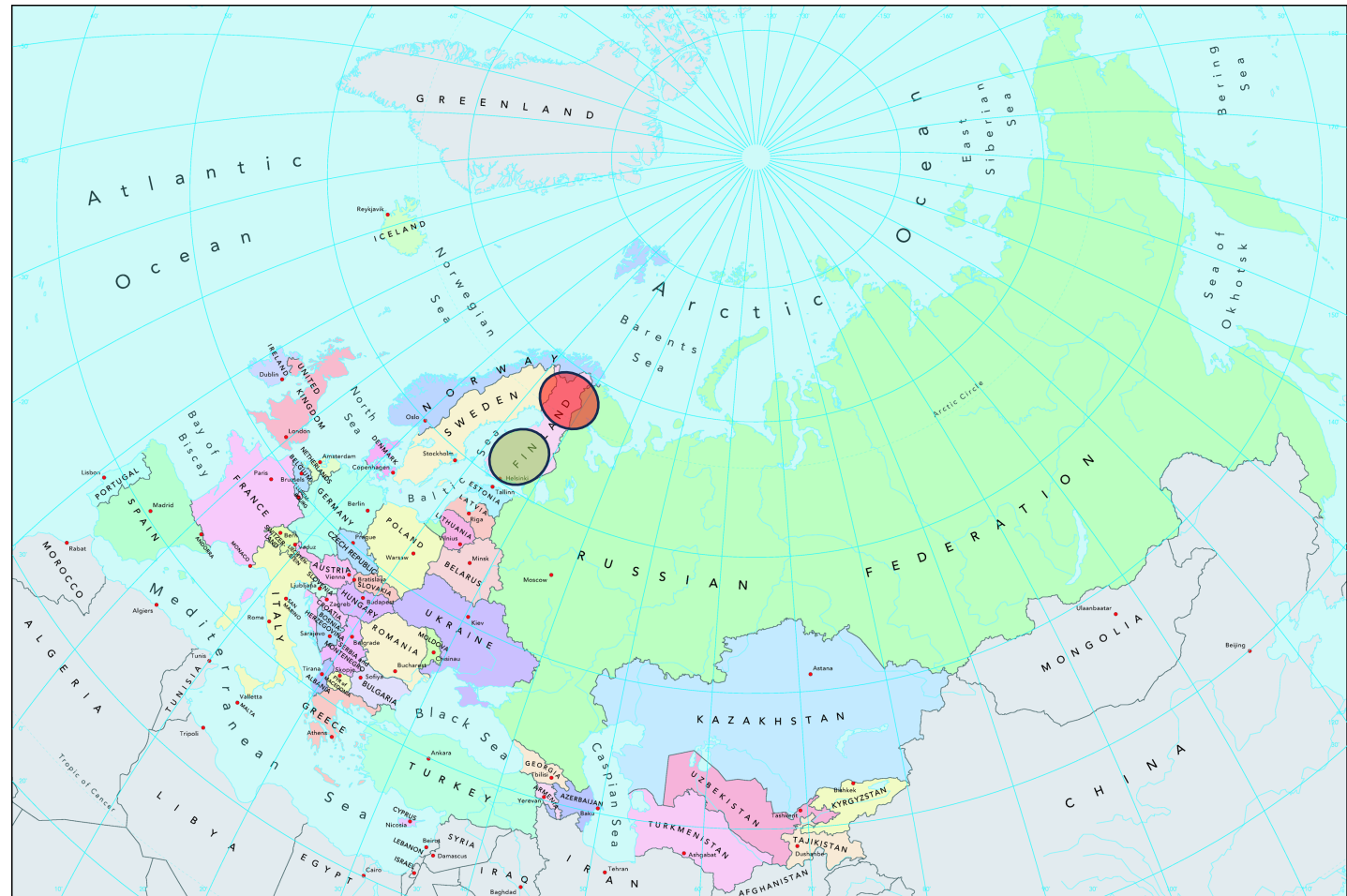
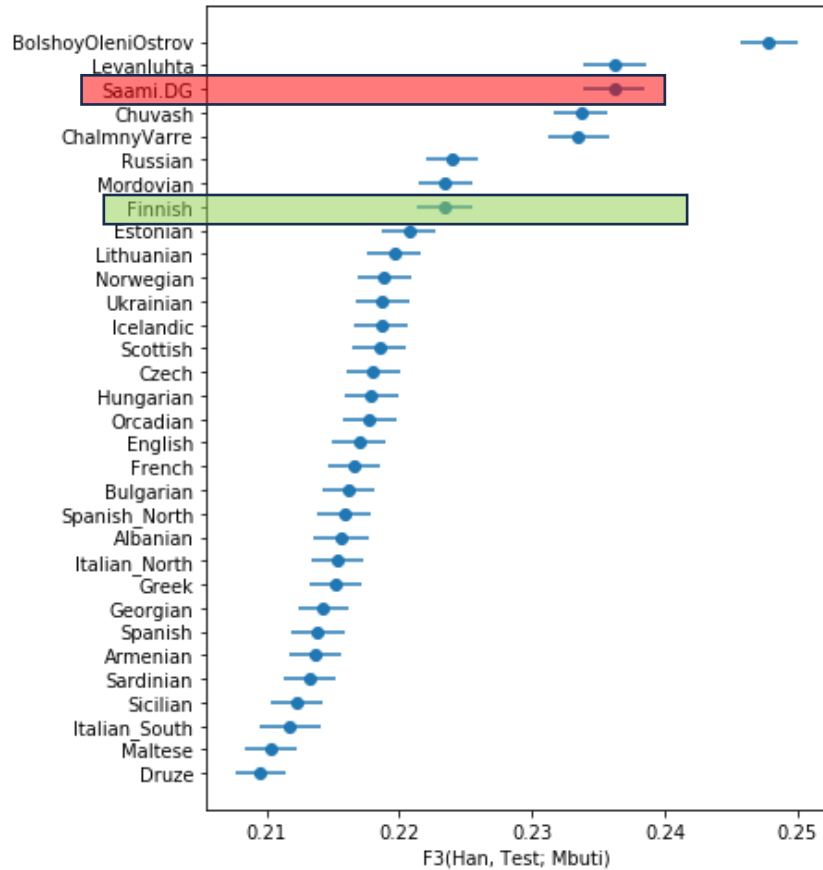
$$f_3(C; A, B) = \frac{1}{J} \sum_{j=1}^J (c_j - a_j)(c_j - b_j)$$

$$f_3 < 0$$



Outgroup f_3 statistic – Example

Goal: We want to test the genetic affinity of European populations to East Asia, by performing the statistic $f_3(\text{Han}, X; \text{Mbuti})$, where Mbuti is a distant African population and acts as outgroup here, Han denotes Han Chinese, and X denotes various European populations



Target f_3 statistic – Example

We can use target f_3 to better understand what is the genetic relationship between East Asia and Europe (and the Americas)

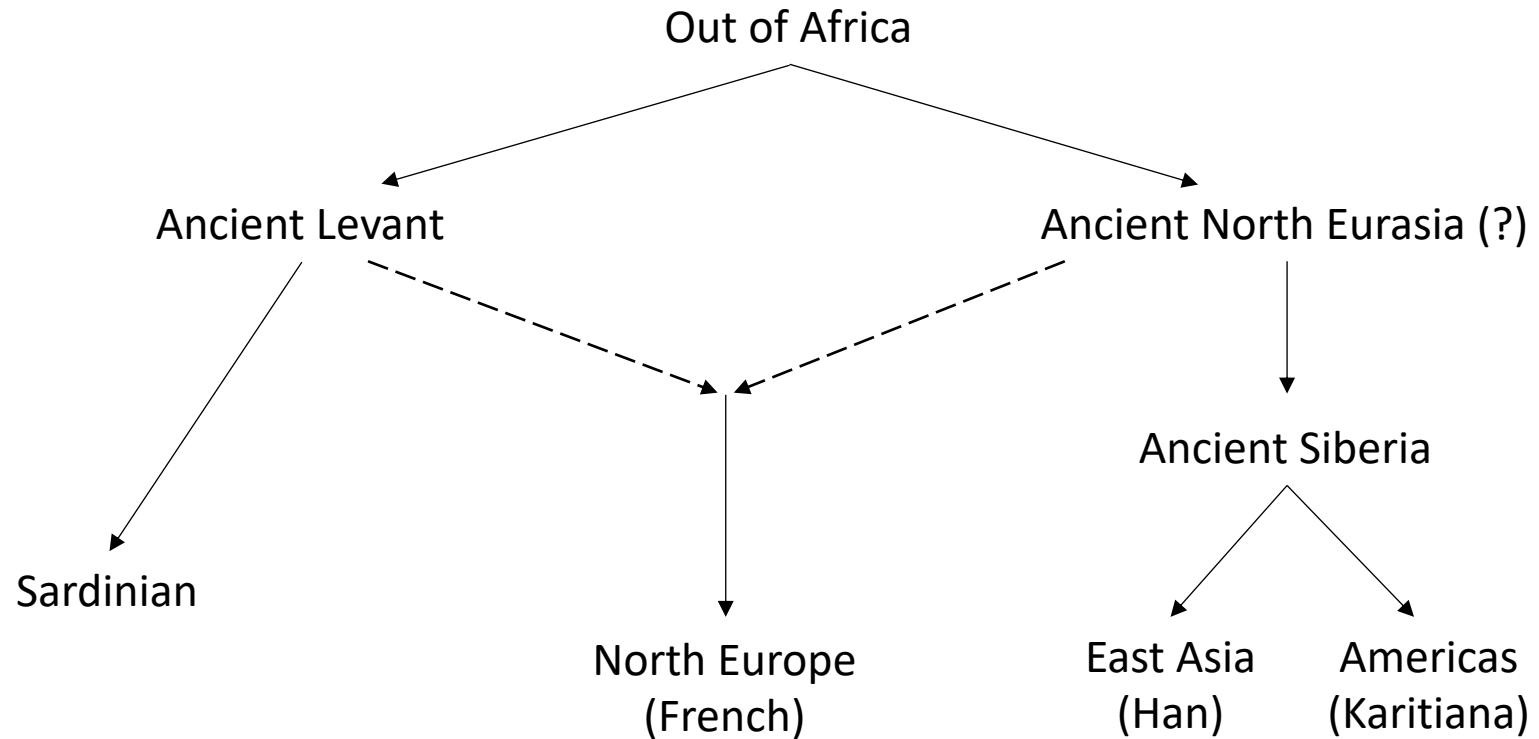
Source1	Source2	Target	f_3	Z-score
Japanese	Italian	Uygur	-0.0259	-74.79
Japanese	Italian	Hazara	-0.0230	-74.05
Yoruba	Sardinian	Mozabite	-0.0211	-56.95
Mozabite	Surui	Maya	-0.0149	-19.67
Yoruba	San	Bantu-SA	-0.0107	-31.39
Yoruba	Sardinian	Palestinian	-0.0107	-36.70
Yoruba	Sardinian	Bedouin	-0.0104	-33.73
Druze	Yi	Burusho	-0.0090	-27.62
Sardinian	Karitiana	Russian	-0.0086	-20.68
Druze	Karitiana	Pathan	-0.0084	-22.25
Han	Orcadian	Tu	-0.0076	-20.64
Mbuti	Orcadian	Makrani	-0.0076	-19.56
Han	Orcadian	Mongola	-0.0075	-19.21
Han	French	Xibo	-0.0069	-16.92
Druze	Dai	Sindhi	-0.0067	-21.99
Sardinian	Karitiana	French	-0.0060	-18.36
Dai	Italian	Cambodian	-0.0060	-13.16
Sardinian	Karitiana	Adygei	-0.0057	-13.03
Biaka	Sardinian	Bantu-Kenya	-0.0054	-13.42
Sardinian	Karitiana	Tuscan	-0.0052	-11.26
Sardinian	Pima	Italian	-0.0045	-12.48
Druze	Karitiana	Balochi	-0.0044	-11.58
Daur	Dai	Han	-0.0026	-13.20
Han	Orcadian	Han-NChina	-0.0025	-7.09
Han	Yakut	Daur	-0.0025	-9.05
Druze	Karitiana	Brahui	-0.0025	-6.43
Hezhen	Dai	Tujia	-0.0021	-6.97
Sardinian	Karitiana	Orcadian	-0.0019	-4.31
She	Yakut	Oroqen	-0.0017	-5.13



Target f_3 statistic – Example

We can use target f_3 to better understand what is the genetic relationship between East Asia and Europe (and the Americas)

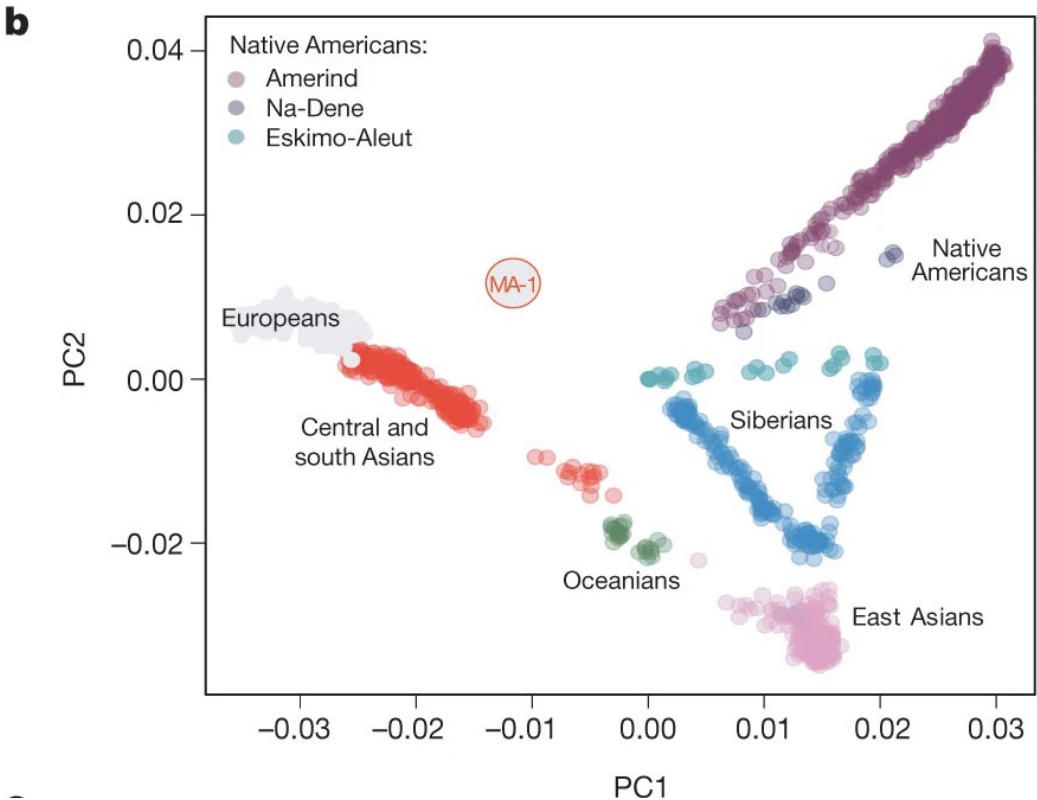
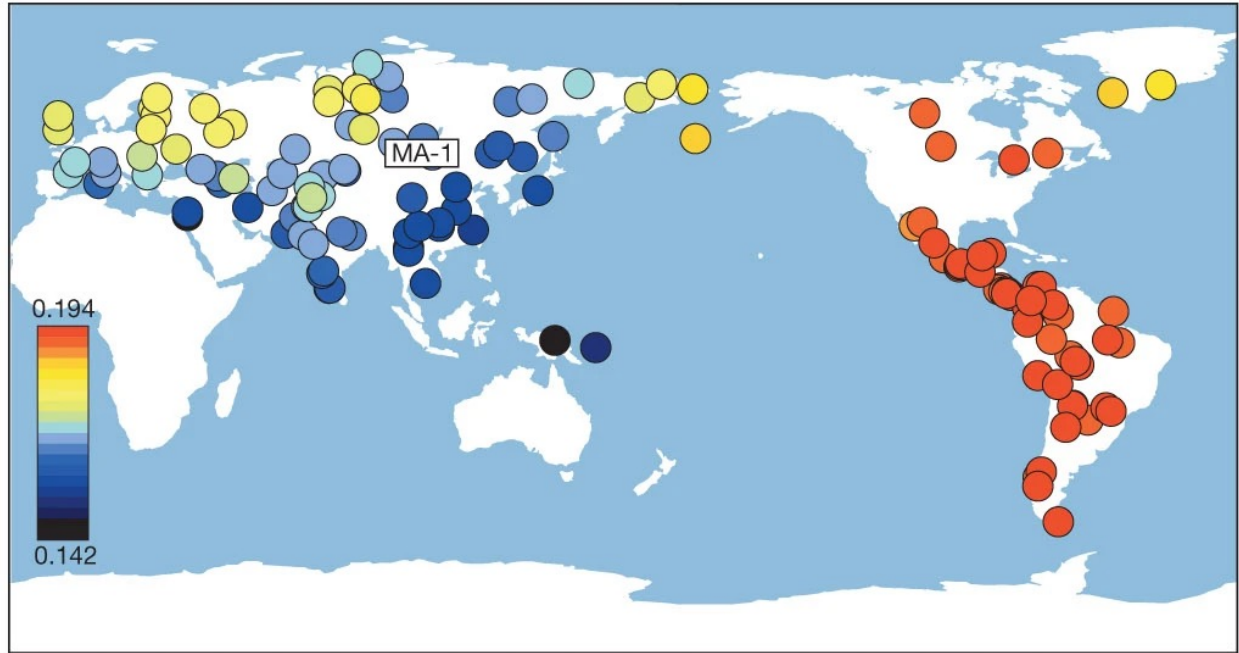
Source1	Source2	Target	f_3	Z-score
Japanese	Italian	Uygur	-0.0259	-74.79
Japanese	Italian	Hazara	-0.0230	-74.05
Yoruba	Sardinian	Mozabite	-0.0211	-56.95
Mozabite	Surui	Maya	-0.0149	-19.67
Yoruba	San	Bantu-SA	-0.0107	-31.39
Yoruba	Sardinian	Palestinian	-0.0107	-36.70
Yoruba	Sardinian	Bedouin	-0.0104	-33.73
Druze	Yi	Burusho	-0.0090	-27.62
Sardinian	Karitiana	Russian	-0.0086	-20.68
Druze	Karitiana	Pathan	-0.0084	-22.25
Han	Orcadian	Tu	-0.0076	-20.64
Mbuti	Orcadian	Makrani	-0.0076	-19.56
Han	Orcadian	Mongola	-0.0075	-19.21
Han	French	Xibo	-0.0069	-16.92
Druze	Dai	Sindhi	-0.0067	-21.99
Sardinian	Karitiana	French	-0.0060	-18.36
Dai	Italian	Cambodian	-0.0060	-13.16
Sardinian	Karitiana	Adygei	-0.0057	-13.03
Biaka	Sardinian	Bantu-Kenya	-0.0054	-13.42
Sardinian	Karitiana	Tuscan	-0.0052	-11.26
Sardinian	Pima	Italian	-0.0045	-12.48
Druze	Karitiana	Balochi	-0.0044	-11.58
Daur	Dai	Han	-0.0026	-13.20
Han	Orcadian	Han-NChina	-0.0025	-7.09
Han	Yakut	Daur	-0.0025	-9.05
Druze	Karitiana	Brahui	-0.0025	-6.43
Hezhen	Dai	Tujia	-0.0021	-6.97
Sardinian	Karitiana	Orcadian	-0.0019	-4.31
She	Yakut	Oroqen	-0.0017	-5.13

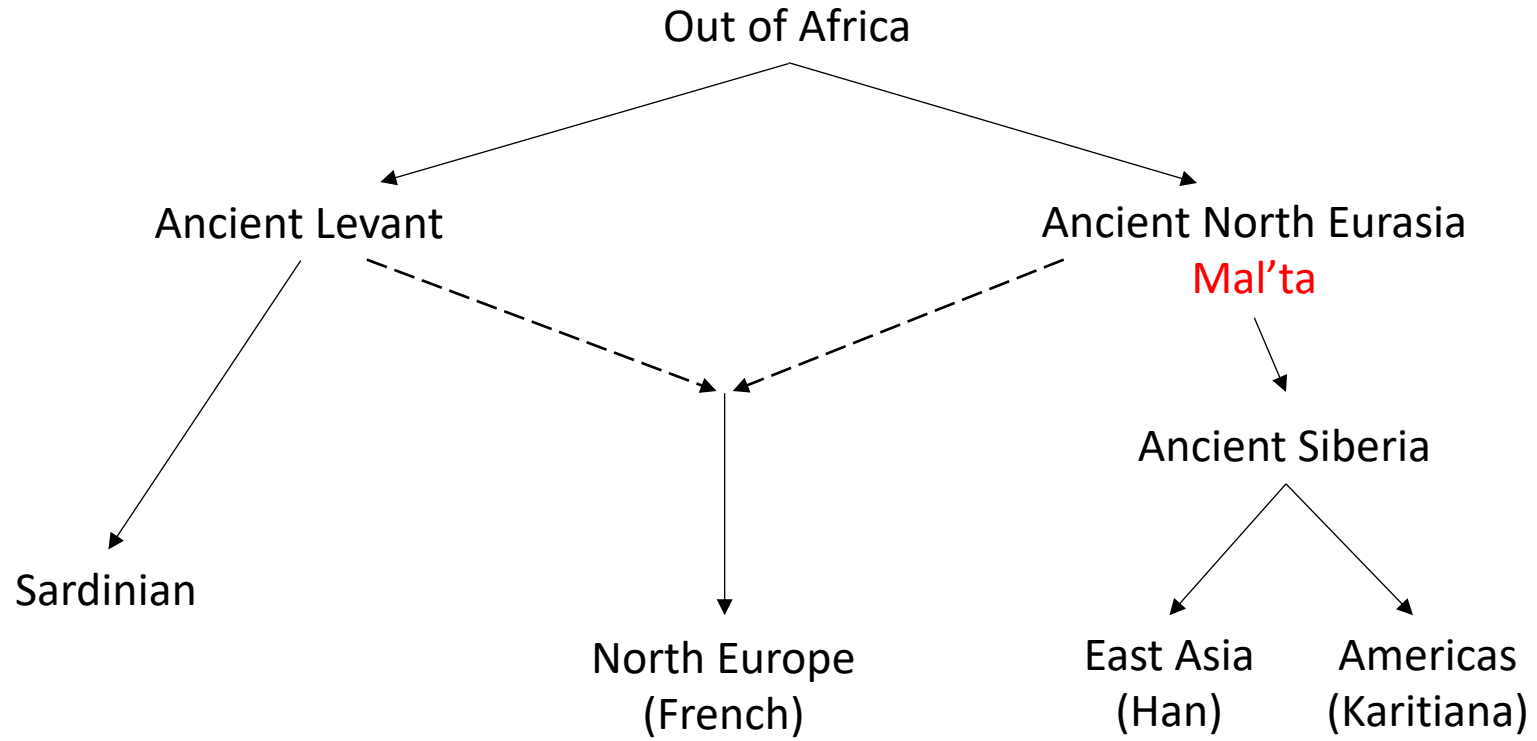




Ancient North Eurasia

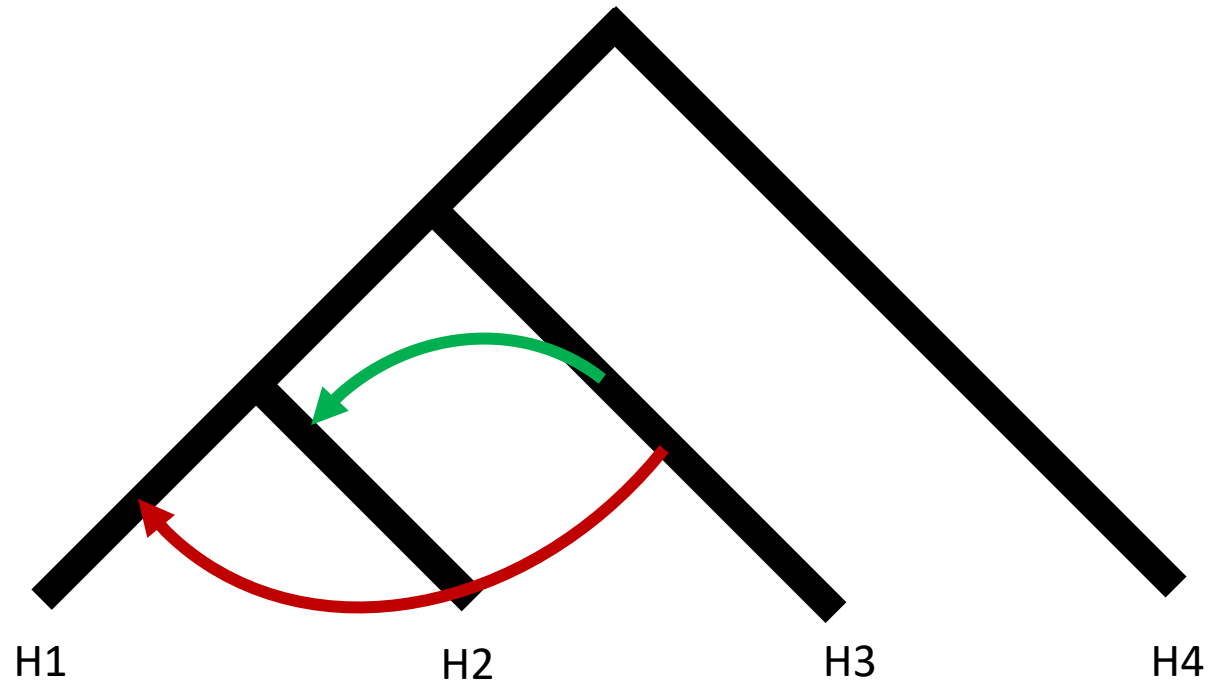
24,000-year-old individual (MA-1) from Mal'ta





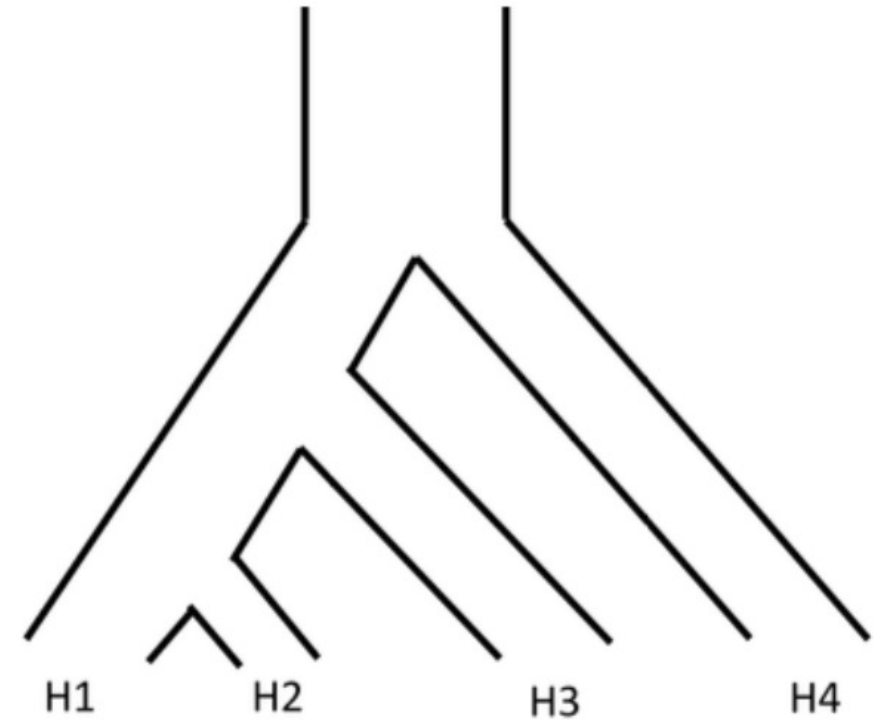
D statistic

Detect signature of admixture between populations



D statistic

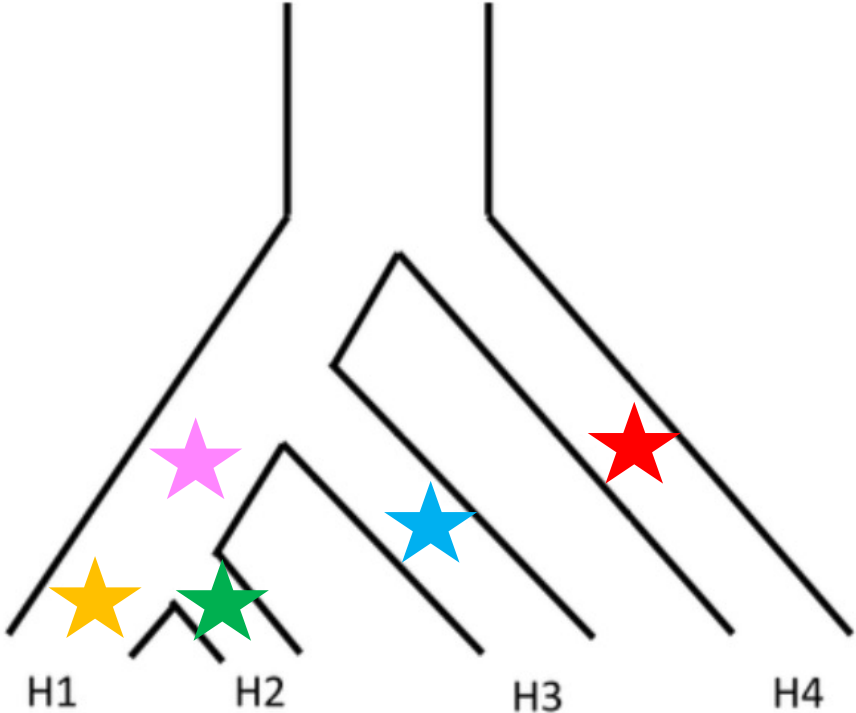
- Analyse a tree with four population
- Pick one individual for each population (it can be performed also with the whole population)
- Look at a polymorphic site – “A” is the ancestral state and “B” is the derived one
- Possible observable pattern of allele sharing



B	A	A	A
A	B	A	A
A	A	B	A
A	A	A	B
A	B	B	A
B	A	B	A
B	B	A	A

D statistic

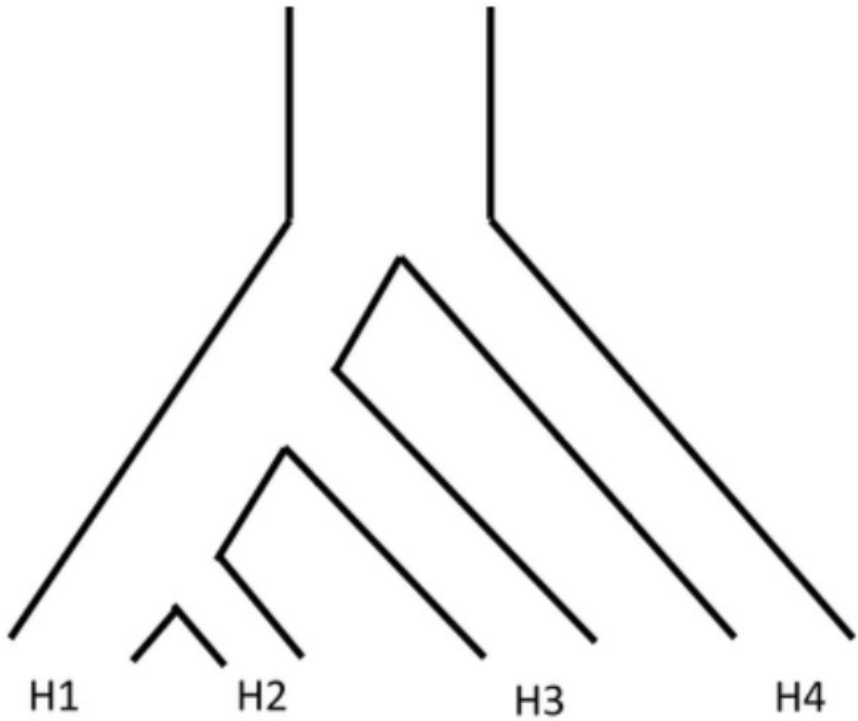
How to explain the patterns?



B	A	A	A
A	B	A	A
A	A	B	A
A	A	A	B
A	B	B	A
B	A	B	A
B	B	A	A

D statistic

How to explain the patterns?

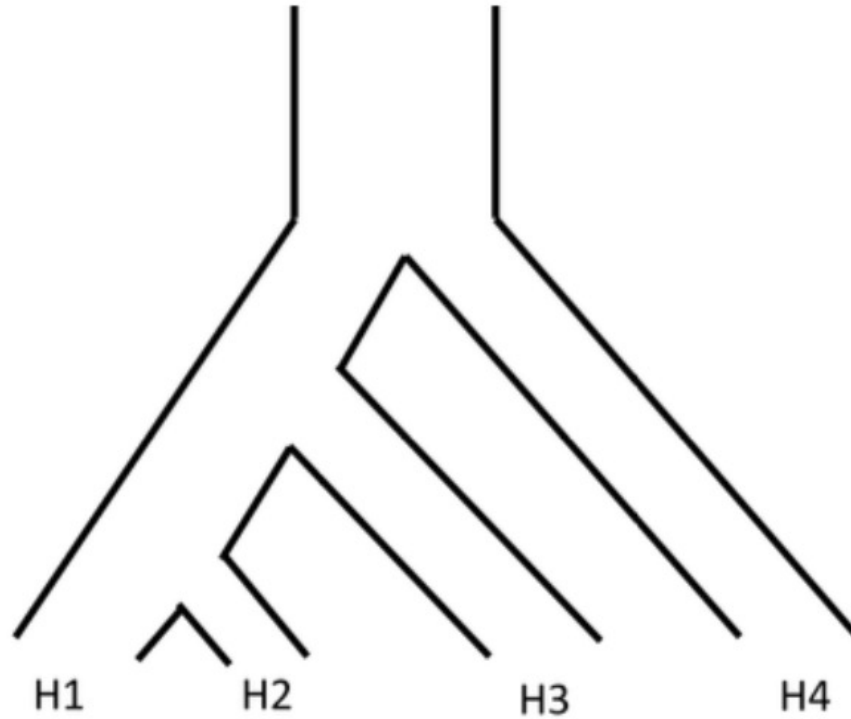


B	A	A	A
A	B	A	A
A	A	B	A
A	A	A	B
A	B	B	A
B	A	B	A
B	B	A	A

D statistic

How to explain the patterns?

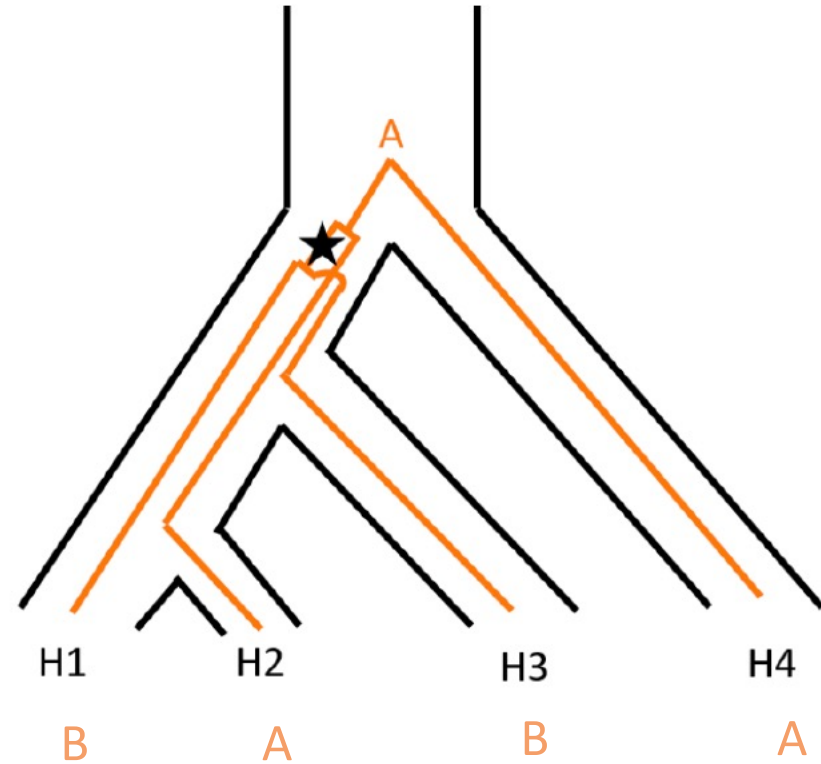
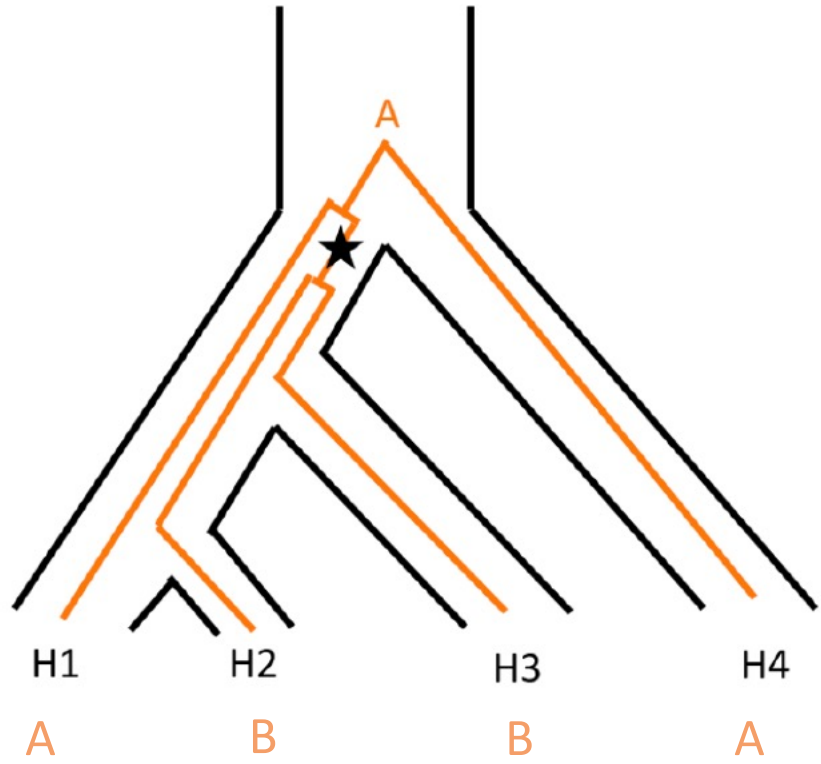
gene genealogies not necessarily follow the **population tree**



B	A	A	A
A	B	A	A
A	A	B	A
A	A	A	B
A	B	B	A
B	A	B	A
B	B	A	A

D statistic

ABBA and BABA sites



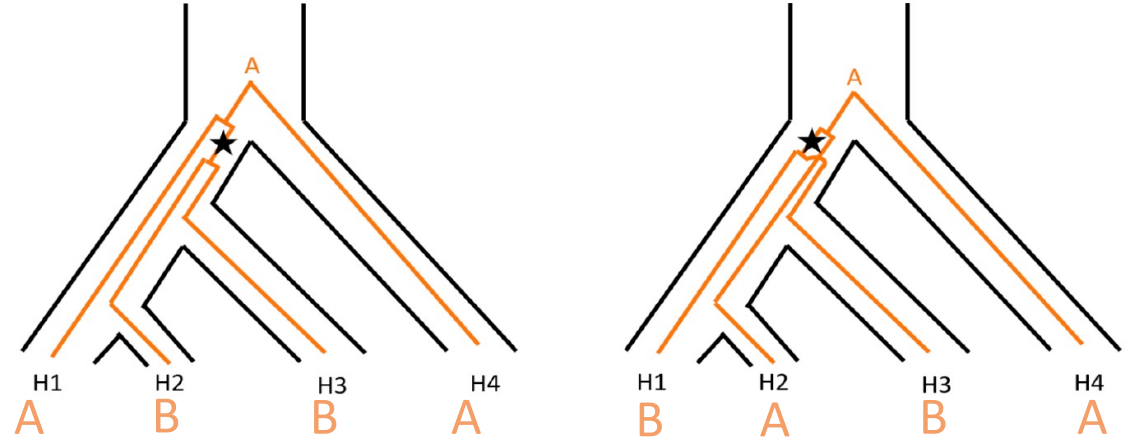
D statistic

D statistic is calculated in this way:

$$D(H_1, H_2; H_3, H_4) = \frac{(n_{ABBA} - n_{BABA})}{(n_{ABBA} + n_{BABA})}$$

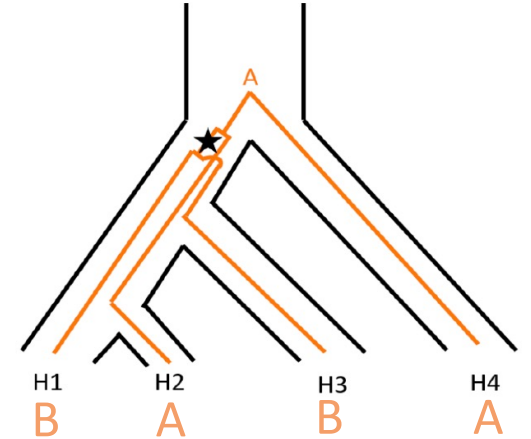
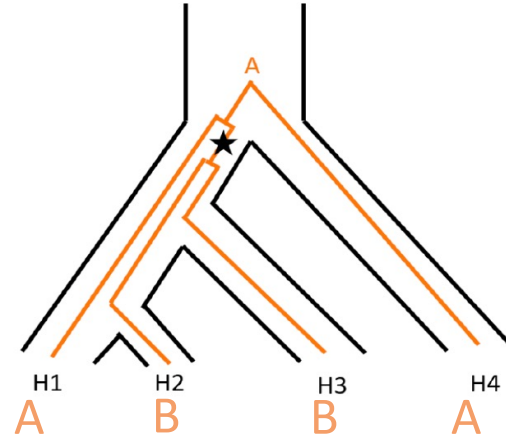
Using several (all) the loci in the genome

We are observing which pattern is the most frequent, ABBA or BABA



D statistic

$$D(H_1, H_2; H_3, H_4) = \frac{(n_{ABBA} - n_{BABA})}{(n_{ABBA} + n_{BABA})}$$



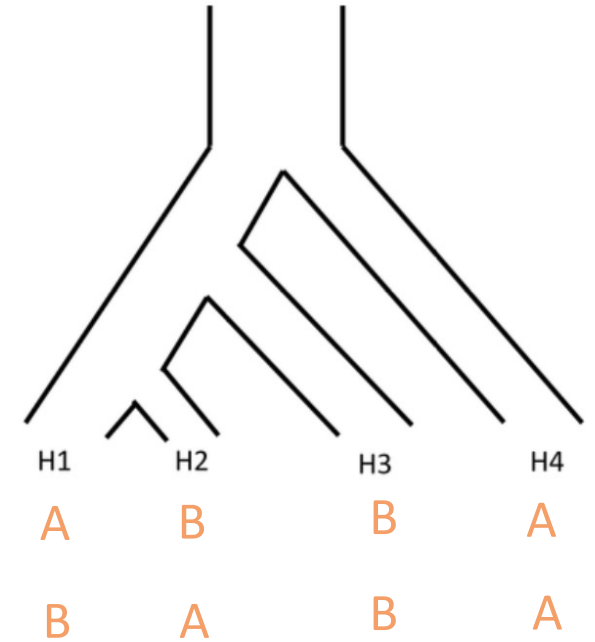
$D = (1000 - 500) / (1000 + 500) = 0.33$ $D > 0$ if ABBA is more common

$D = (500 - 1000) / (500 + 1000) = -0.33$ $D < 0$ if BABA is more common

Interpreting D statistic

ABBA and BABA sites should be equally represented

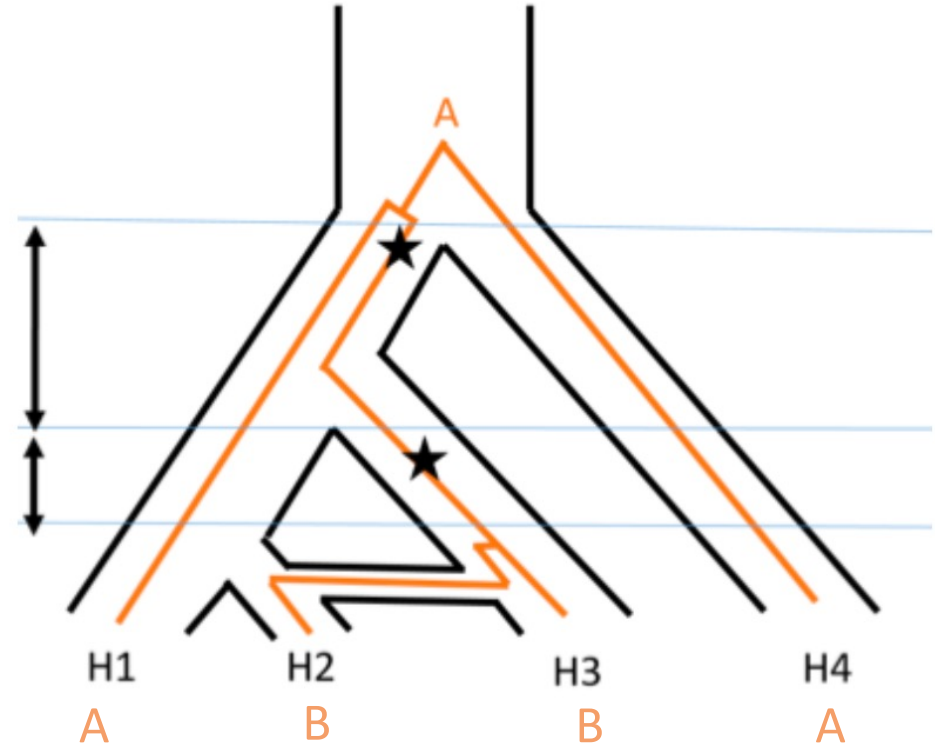
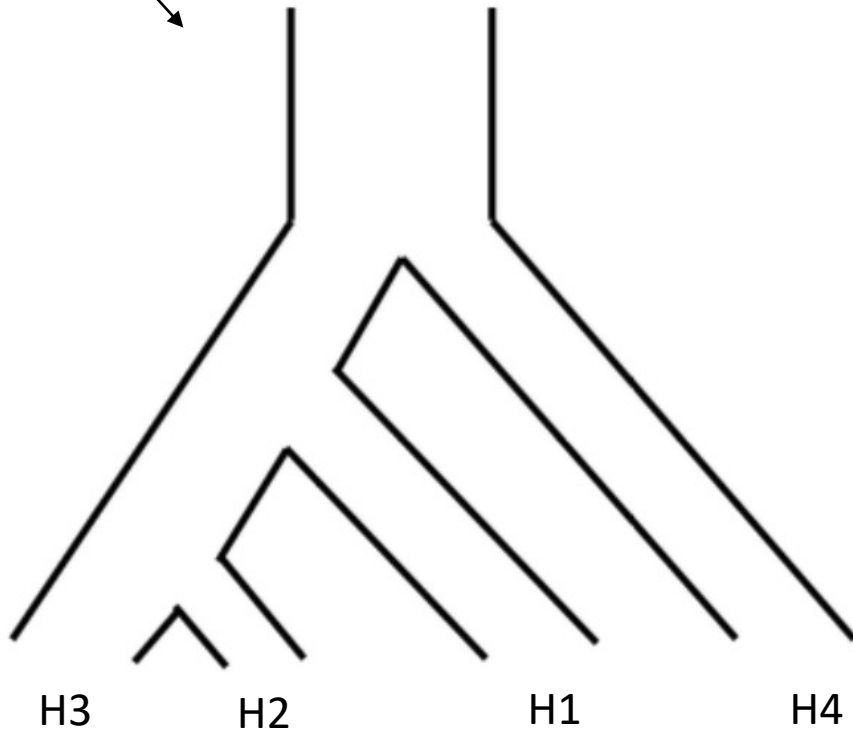
What is happening if they are not?



Interpreting D statistic

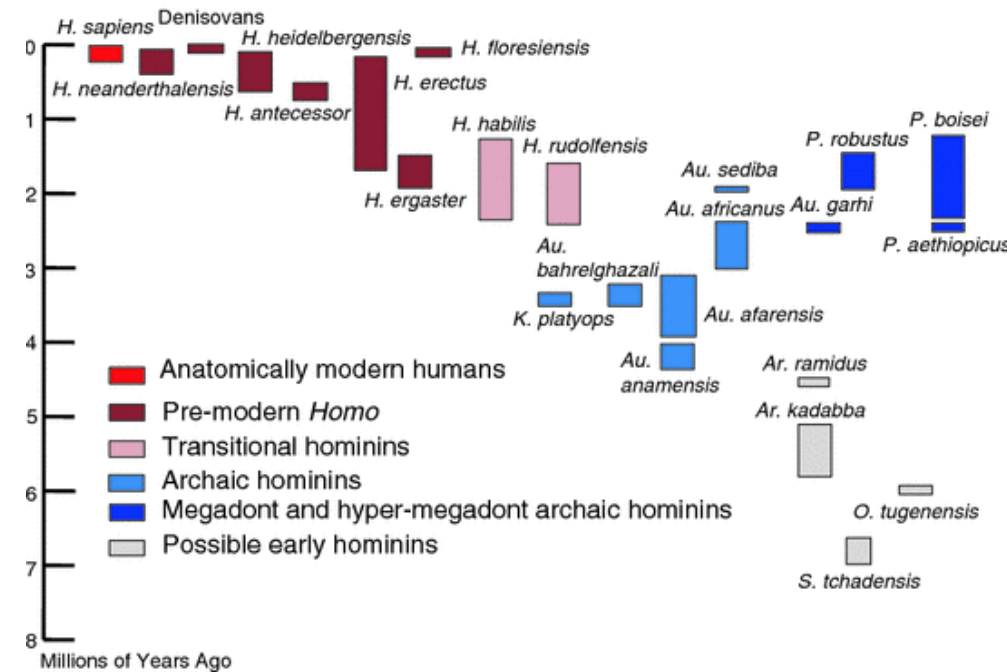
What if $D \neq 0$?

- Gene flow
- The tree is not correct

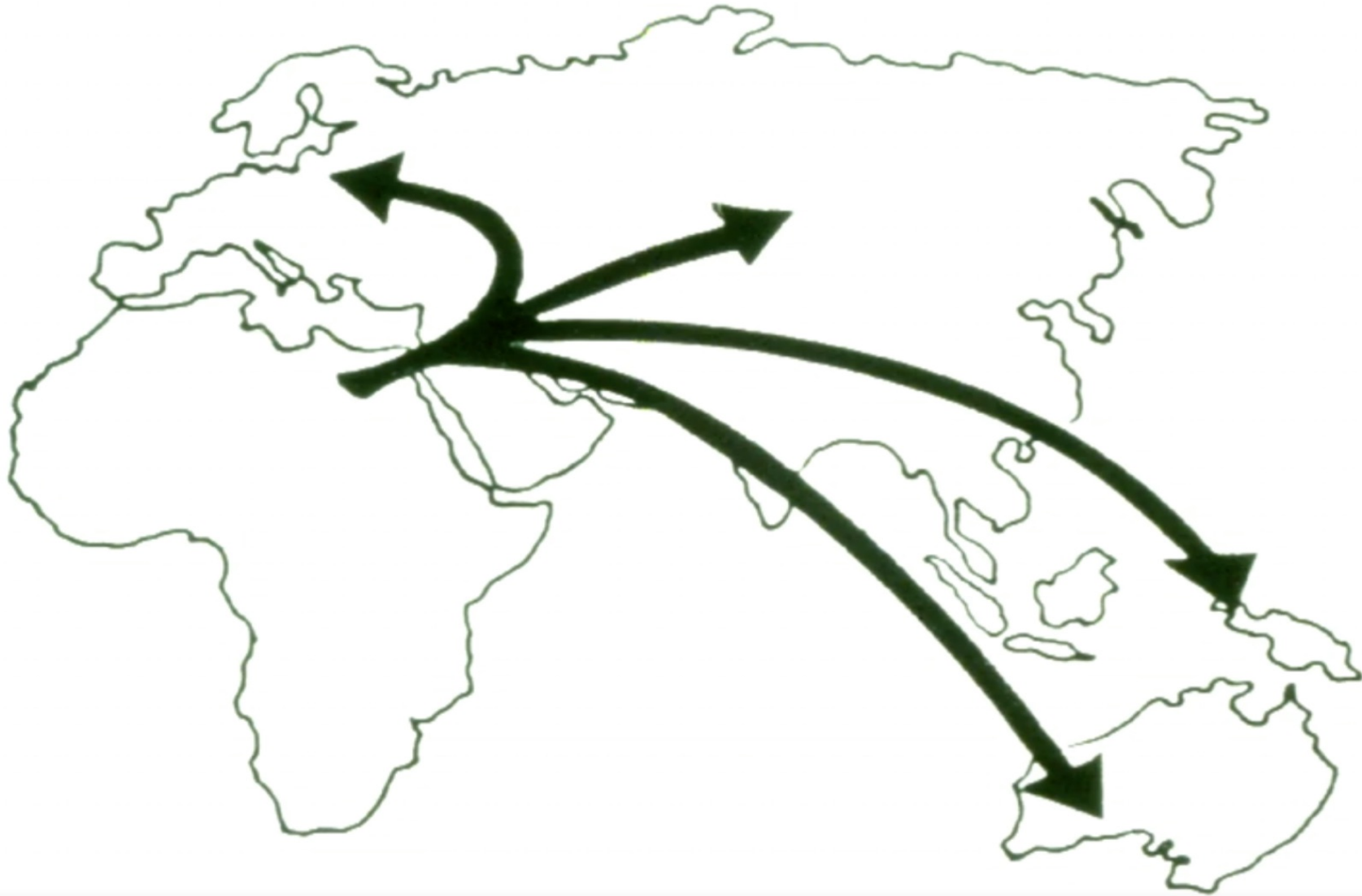


Neanderthal

- First ancient hominin discovered
- Modern humans closest relative
- Lived between $\approx 400,000$ and $40,000$ years ago



Out of Africa



D statistic – Human/Neanderthal admixture

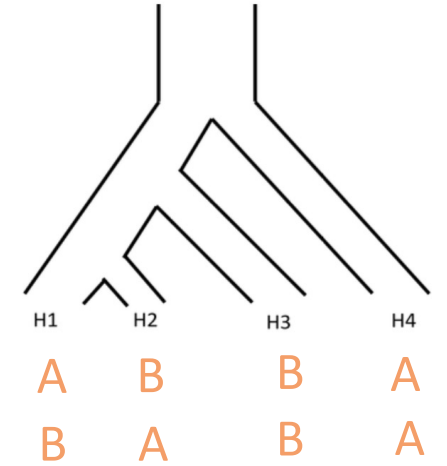
Whole genome sequences for one individual (or more) from each of the six following populations:

- Neanderthal
- Yoruba (Africa)
- Dinka (Africa)
- French (Europe)
- Han Chinese (East Asia)
- Chimpanzee (Outgroup)

We can compare their genomes and calculate the number of ABBA and BABA sites.

D statistic – Human/Neanderthal admixture

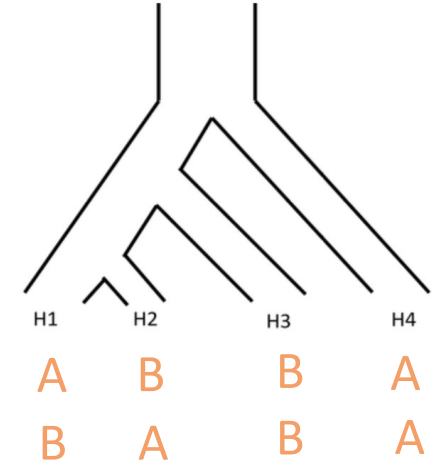
H1	H2	H3	H4	N° ABBA	N° BABA
Yoruba	Dinka	Neanderthal	Chimpanzee	44,161	44,221
Yoruba	French	Neanderthal	Chimpanzee	46,449	44,347
Yoruba	Han	Neanderthal	Chimpanzee	48,227	43,863



$$D(H_1, H_2; H_3, H_4) = \frac{(n_{ABBA} - n_{BABA})}{(n_{ABBA} + n_{BABA})}$$

D statistic – Human/Neanderthal admixture

H1	H2	H3	H4	N° ABBA	N° BABA
Yoruba	Dinka	Neanderthal	Chimpanzee	44,161	44,221
Yoruba	French	Neanderthal	Chimpanzee	46,449	44,347
Yoruba	Han	Neanderthal	Chimpanzee	48,227	43,863



$$D(H_1, H_2; H_3, H_4) = \frac{(n_{ABBA} - n_{BABA})}{(n_{ABBA} + n_{BABA})}$$

	Test	D-stat	Standard error	Z-score
Scenario 1	(Yoruba, Dinka; Neanderthal, Chimp)	-0.000678	0.00336	-0.201
Scenario 2	(Yoruba, French; Neanderthal, Chimp)	0.02315	0.00473	4.894
Scenario 3	(Yoruba, Han; Neanderthal, Chimp)	0.04738	0.00543	8.725

D statistic – Human/Neanderthal admixture

	Test	D-stat	Standard error	Z-score
Scenario 1	(Yoruba, Dinka; Neanderthal, Chimp)	-0.000678	0.00336	-0.201

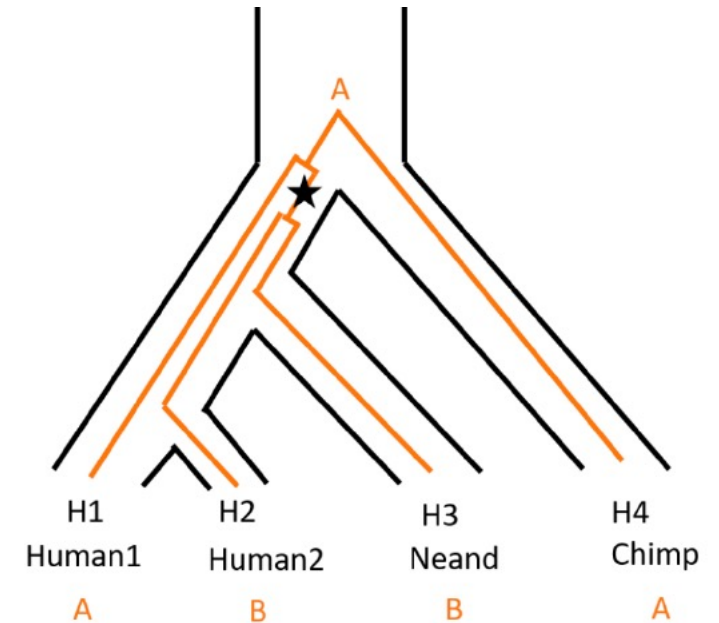
This result suggest that the pair of African genomes are symmetrically related to the Neanderthal and the chimp. Therefore, we infer that these two Africans form a clade to the exclusion of the Neanderthal and the chimp.

Moreover, we observe **no statistically significant evidence of gene flow between the African individuals and the Neanderthal.**

D statistic – Human/Neanderthal admixture

	Test	D-stat	Standard error	Z-score
Scenario 2	(Yoruba, French; Neanderthal, Chimp)	0.02315	0.00473	4.894

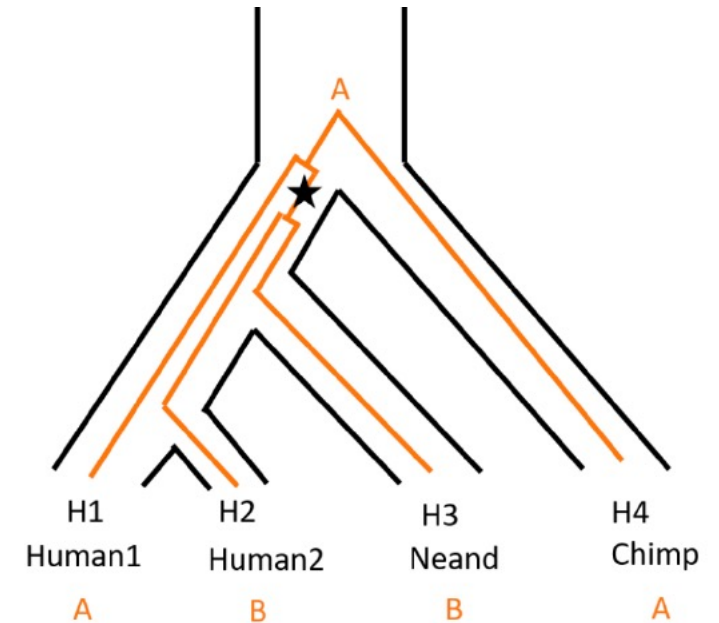
This result suggests that the **French genome shares a statistically significant larger proportion of derived alleles with the Neanderthal genome** (excess of ABBA sites), than the Yoruba does.



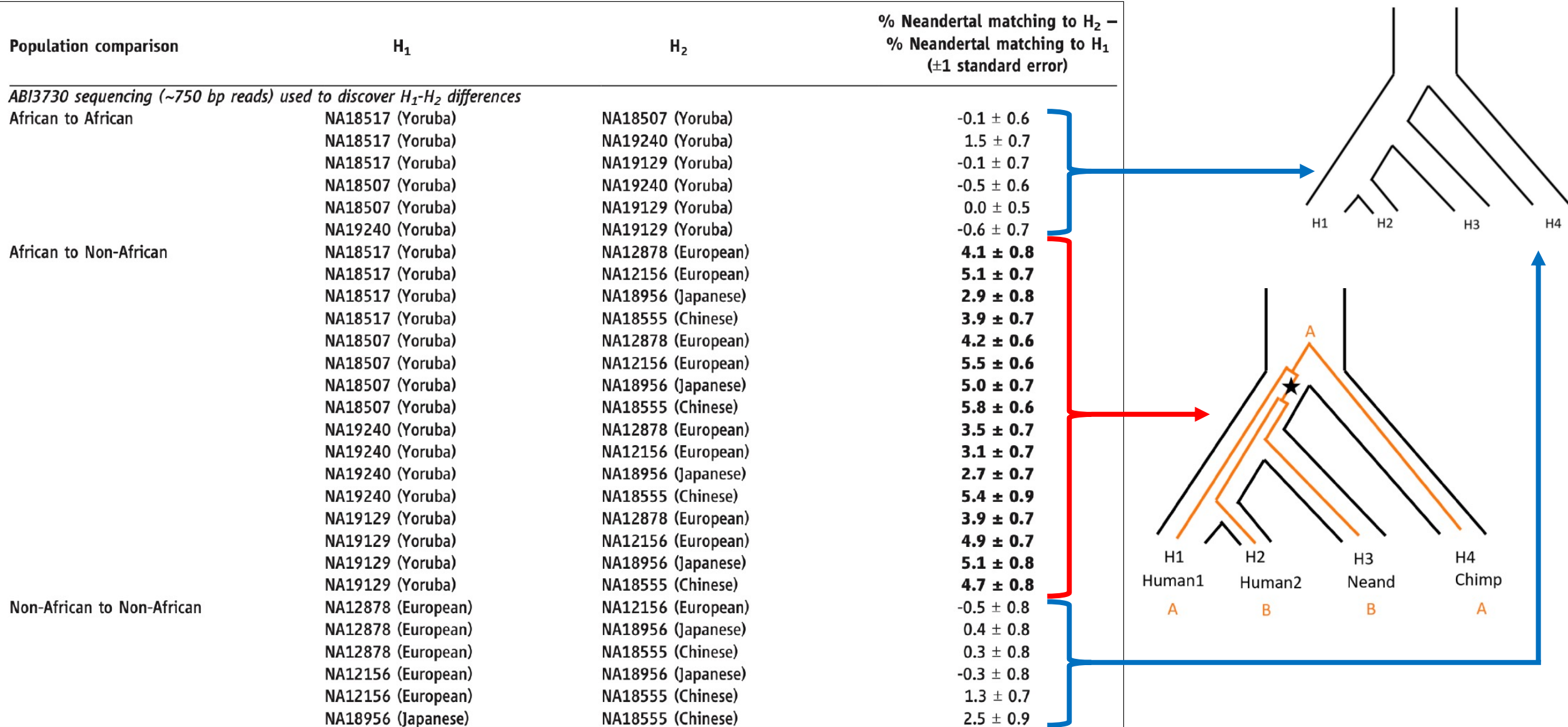
D statistic – Human/Neanderthal admixture

	Test	D-stat	Standard error	Z-score
Scenario 3	(Yoruba, Han; Neanderthal, Chimp)	0.04738	0.00543	8.725

Similar to what we observed for Scenario 2, this suggests that the **Han genome shares a statistically significant larger proportion of derived alleles with the Neanderthal genome** (excess of ABBA sites), than the Yoruba does.



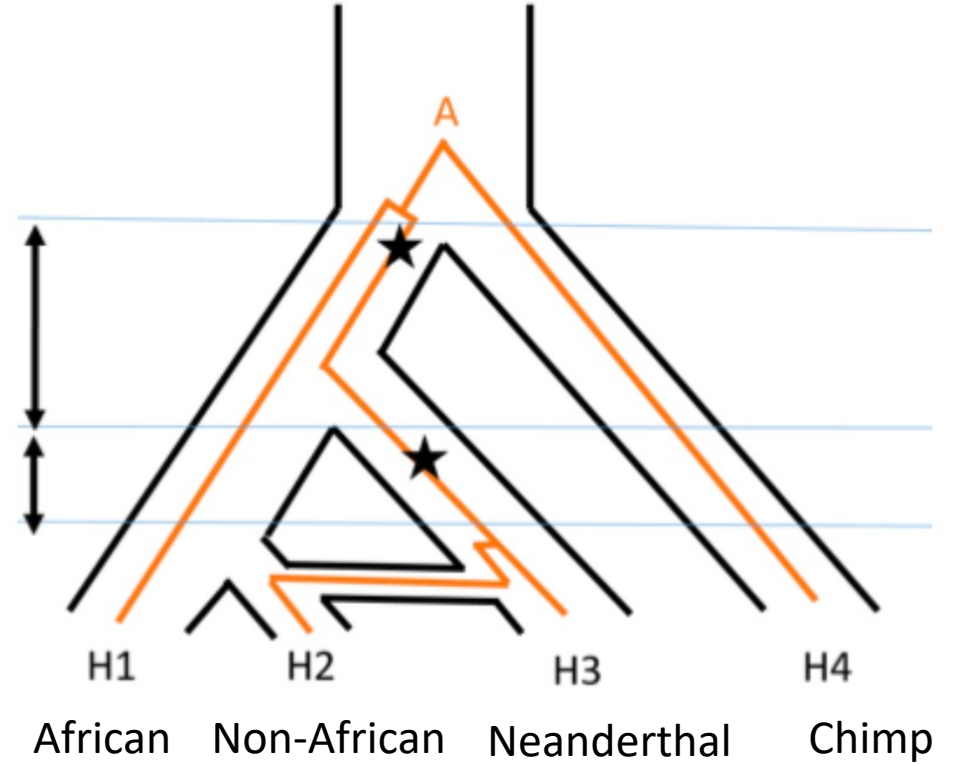
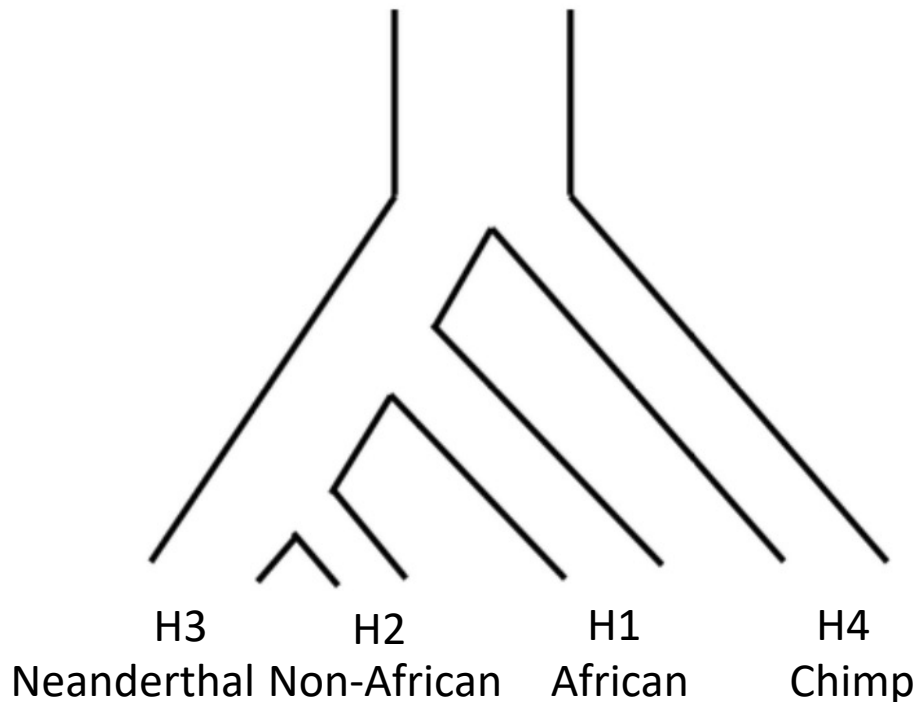
D statistic – Human/Neanderthal admixture



D statistic – Human/Neanderthal admixture

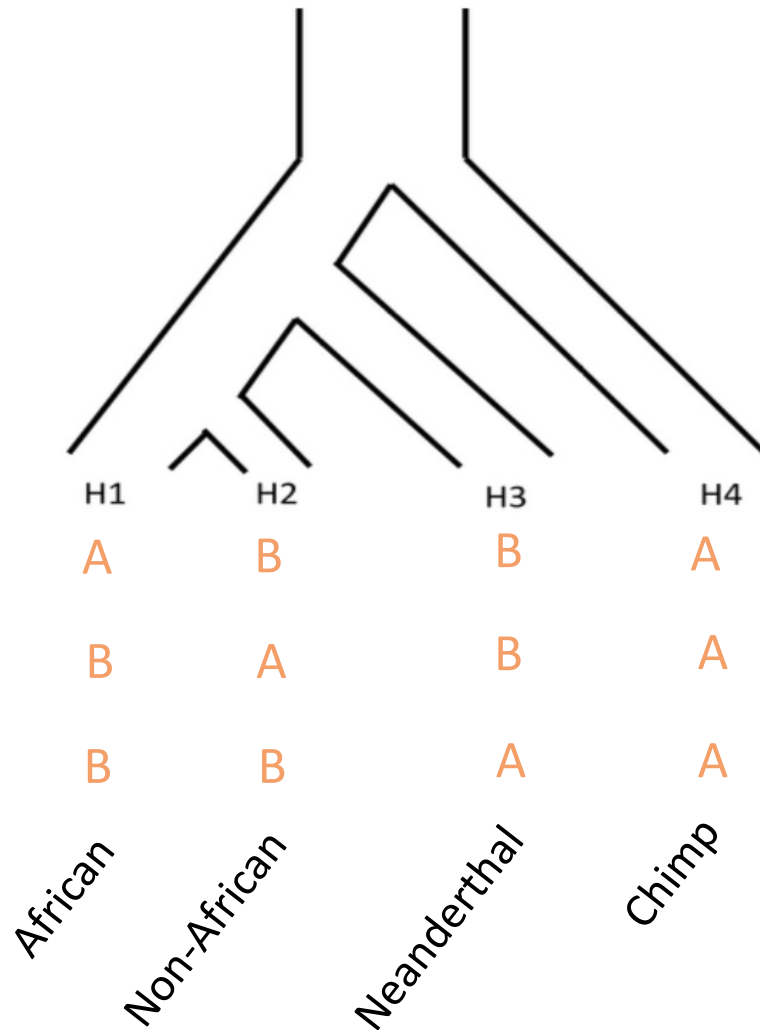
What if $D \neq 0$?

- Gene flow
- The tree is not correct



What is the right model?

D statistic – Human/Neanderthal admixture



How we can discriminate between the two model:

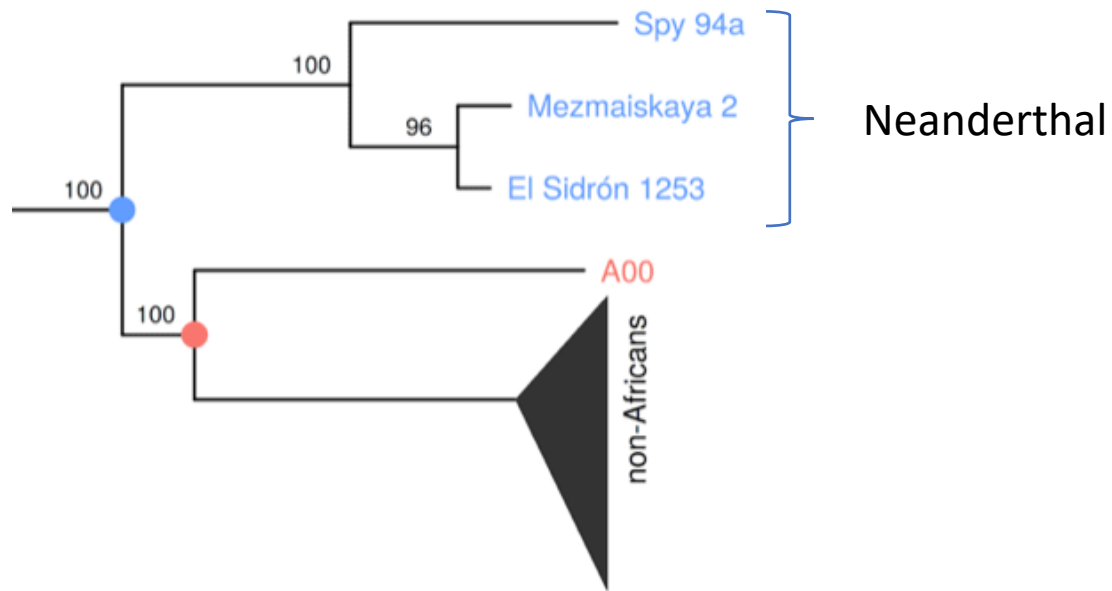
- Look for BBAA sites

D statistic – Human/Neanderthal admixture

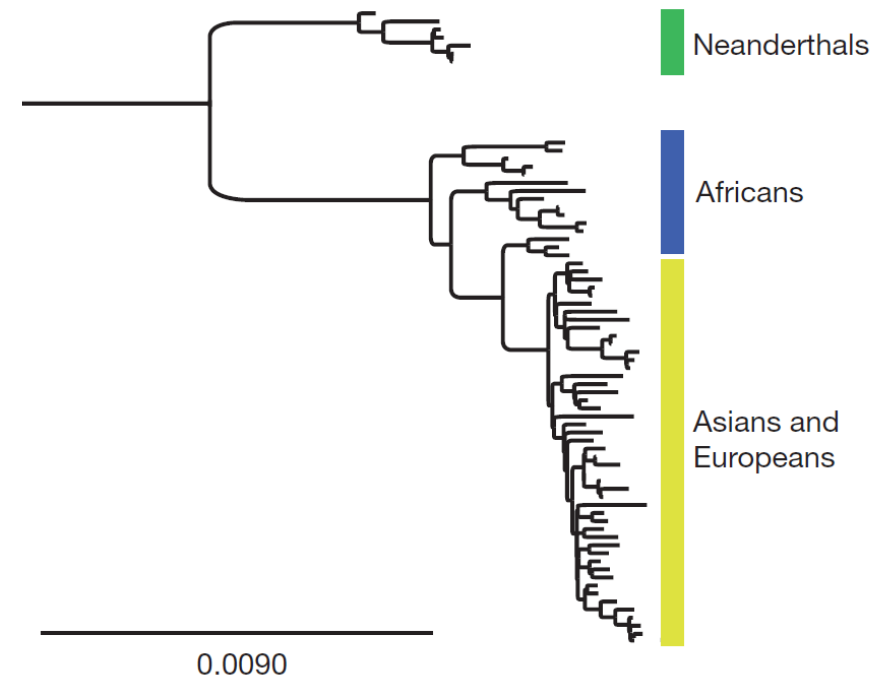
How we can discriminate between the two model:

- Compare the results with different analysis

Y chromosome phylogeny



mtDNA phylogeny



D statistic – Human/Neanderthal admixture

How we can discriminate between the two model:

- Compare the results with different analysis

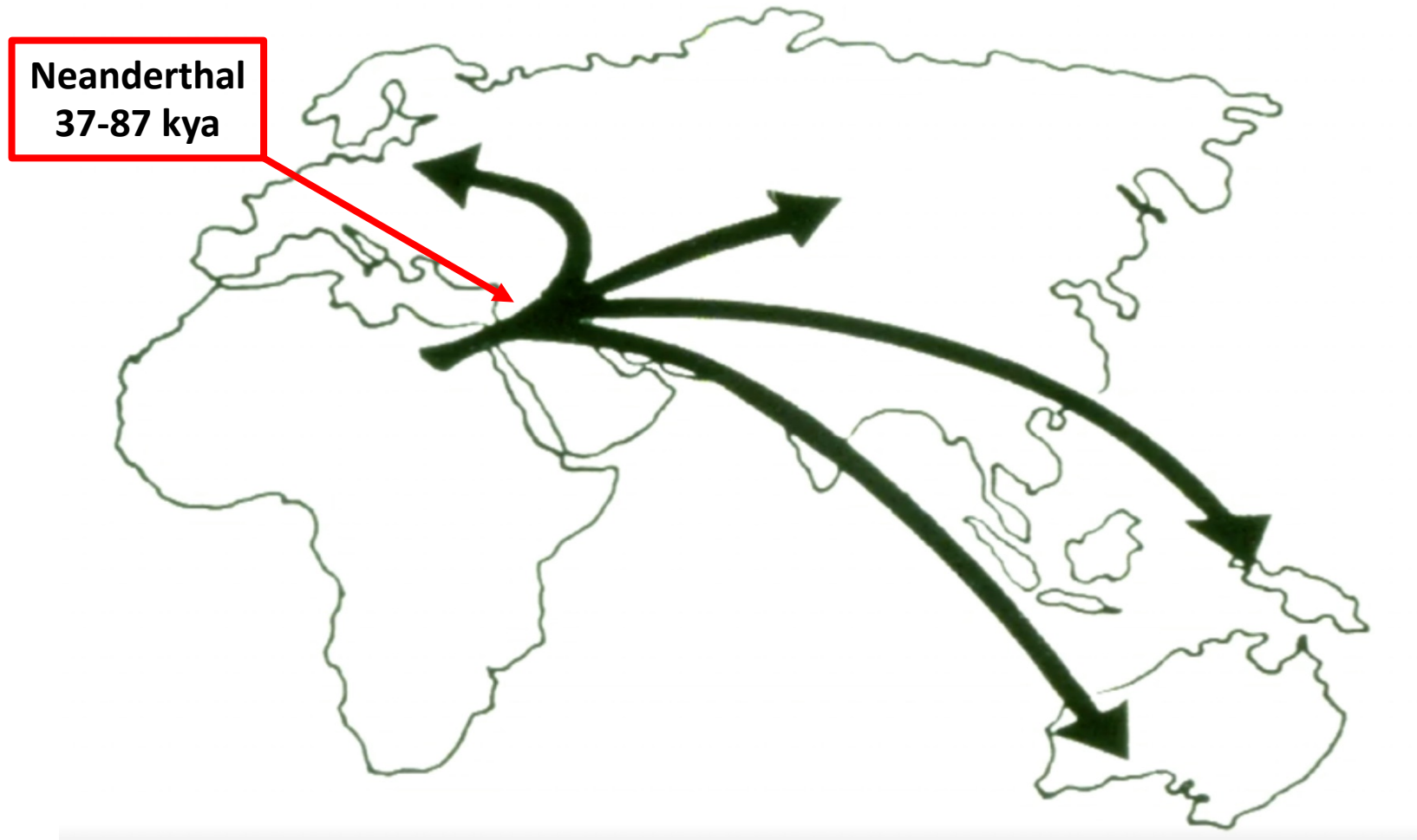
Neanderthal ancestors out of Africa \approx 500 kya



Modern humans out of Africa \approx 100 kya



D statistic – Human/Neanderthal admixture



aDNA

- Hard to analyse (degradation, contamination...)
- Incredibly powerful tool for evolutionary and historical reconstructions
- Insights into onset and evolution of diseases

Report

A 5,000-year-old hunter-gatherer already plagued by *Yersinia pestis*

[Julian Susat](#)^{1,11}, [Harald Lübke](#)^{2,11}, [Alexander Immel](#)¹, [Ute Brinker](#)², [Aija Macāne](#)³, [John Meadows](#)^{2,4}, [Britta Steer](#)⁵, [Andreas Tholey](#)⁵, [Ilga Zagorska](#)⁶, [Guntis Gerhards](#)⁶, [Ulrich Schmölcke](#)², [Mārcis Kalniņš](#)⁶, [Andre Franke](#)¹, [Elīna Pētersone-Gordina](#)⁶, [Barbara Teßman](#)⁷, [Mari Tõrv](#)⁸, [Stefan Schreiber](#)^{1,9}, [Christian Andree](#)¹⁰, [Valdis Bērziņš](#)⁶, [Almut Nebel](#)¹...[Ben Krause-Kyora](#)^{1,12}  

Genotype of a historic strain of *Mycobacterium tuberculosis*

[Abigail S. Bouwman](#)^{a,1,2}, [Sandra L. Kennedy](#)^{a,2}, [Romy Müller](#)^{a,2}, [Richard H. Stephens](#)^a, [Malin Holst](#)^b, [Anwen C. Caffell](#)^c, [Charlotte A. Roberts](#)^c, and [Terence A. Brown](#)^{a,3}

Article

The major genetic risk factor for severe COVID-19 is inherited from Neanderthals

<https://doi.org/10.1038/s41586-020-2818-3> [Hugo Zeberg](#)^{1,2,3} & [Svante Pääbo](#)^{1,3,4}
Received: 3 July 2020

nature communications



Article

<https://doi.org/10.1038/s41467-024-45438-1>

Cases of trisomy 21 and trisomy 18 among historic and prehistoric individuals discovered from ancient DNA

