



CS 68: BIOINFORMATICS

Prof. Sara Mathieson
Swarthmore College
Spring 2018



Outline: Apr 11

Faculty candidate talk right after this!

- Recap working in log-space and posterior mean
- Finish human evolution overview
- Begin: PCA
 - Office hours TODAY 1-2:30pm
 - Interviewing faculty candidate 2:30-3:30pm, available after that
 - Email me if you would like to meet for your project
 - Optional PCA lab released on Thursday, PCA reading posted
 - Graded labs and midterm 2 study guide coming out soon!

Main takeaways from feedback forms

- Needs more work: DP, population genetics, coalescent, Tajima's D, HMMs, tree algorithms
- Lot of diversity in terms of in-class style (some prefer board, some slides, some group work) and lab style (some like problem sets, some don't, etc)
- Overall: recent labs long, too much I/O formatting, lots of material too quickly, lots of math, need more biological motivation + big picture, still like worksheets
- Modifications:
 - *More time for Lab 8, more code for reading files*
 - *Lab 9 optional + human evolution motivation for PCA*
 - *Midterm review session (Monday in-class)*
 - *Project scope small (data X, algorithm Y, results Z)*

Topics to cover in the last three weeks

- Machine learning, deep learning, and neural networks
- Ethics (HeLa, prenatal testing)
- Evolution (natural selection, disprove evolution?)
- More statistics, mathematical models, Bayesian computation and inference
- What is going on in bioinformatics industry
- Conservation genetics, other species (plants)
- Structures (i.e. RNA folding, protein structure) vs. sequences
- Algorithms with broad applications, how do our algorithms fit together
- Epigenetics
- More population genetics
- More runtime analysis

Order roughly corresponds to how many people mentioned a topic

Recap: working in log-space and
posterior mean

Why do we need to work in log-space?

- Multiplying (or adding) many small probabilities will result in underflow
- Underflow: at some point the computer will represent very small numbers as 0
- Overflow: at some point very large numbers can no longer be represented
- Multiplication and division are okay:

$$\log(ab) = \log(a) + \log(b)$$

$$\log(a/b) = \log(a) - \log(b)$$

- The difficulty arises with adding:

$$\log(a + b) = ???$$

In Python3 on my system:

$$e^{709} = 8 \times 10^{307}$$

$$e^{710} = \text{overflow}$$

$$e^{-745} = 5 \times 10^{324}$$

$$e^{-746} = 0.0 \text{ (underflow)}$$

working in log-space

$$\tilde{p} = \log p$$

forward recursion

$$\log(f_k(i) = e_k(x_i) \cdot \sum_l f_l(i-1) a_{lk})$$

$$\tilde{f}_k(i) = \log(e_k(x_i)) + \log \left(\underbrace{f_0(i-1) a_{0,k}}_p + \underbrace{f_1(i-1) a_{1,k}}_q + \underbrace{f_2(i-1) a_{2,k}}_r + \underbrace{f_3(i-1) a_{3,k}}_s \right)$$

$$\rightarrow \log \left(\underbrace{p+q+r+s}_t \right)$$
$$\underbrace{\hspace{1.5cm}}_u$$

$$\begin{array}{c} \log \swarrow \quad \log \swarrow \\ \approx -1600 \quad -1700 \end{array} \rightarrow \tilde{q} - \tilde{p} = -100$$

$$\begin{aligned}
 & \log(p+q) \\
 &= \log(e^{\tilde{p}} + e^{\tilde{q}}) \\
 &= \log\left(e^{\tilde{p}} \cdot \left(1 + \frac{e^{\tilde{q}}}{e^{\tilde{p}}}\right)\right) \\
 &= \log\left(e^{\tilde{p}} \cdot (1 + e^{\tilde{q}-\tilde{p}})\right) \\
 &= \tilde{p} + \log(1 + e^{\tilde{q}-\tilde{p}})
 \end{aligned}$$

def log-add(\tilde{p}, \tilde{q}):

log-add-all(array) $\log(p+q)$

$\tilde{t} = \text{log-add}(\tilde{p}, \tilde{q})$

$\tilde{u} = \text{log-add}(\tilde{t}, \tilde{r})$ } loop!

result = log-add(\tilde{u}, \tilde{s})

variance

x

high variance

x

$(\bar{p})^2$

posterior mean

times = [0.32, 1.75, 4.54, 9.40]

$$\bar{z}_i = \sum_k P(z_i = k | \vec{x}) \cdot g(k)$$

$$\bar{z}_i = \left(\frac{2}{10}\right)(0.32) + \left(\frac{1}{10}\right)(1.75) + \left(\frac{4}{10}\right)(4.54) + \left(\frac{3}{10}\right)(9.40)$$

weighted
avg.

0	$\frac{2}{10}$
1	$\frac{1}{10}$
2	$\frac{4}{10}$
3	$\frac{3}{10}$

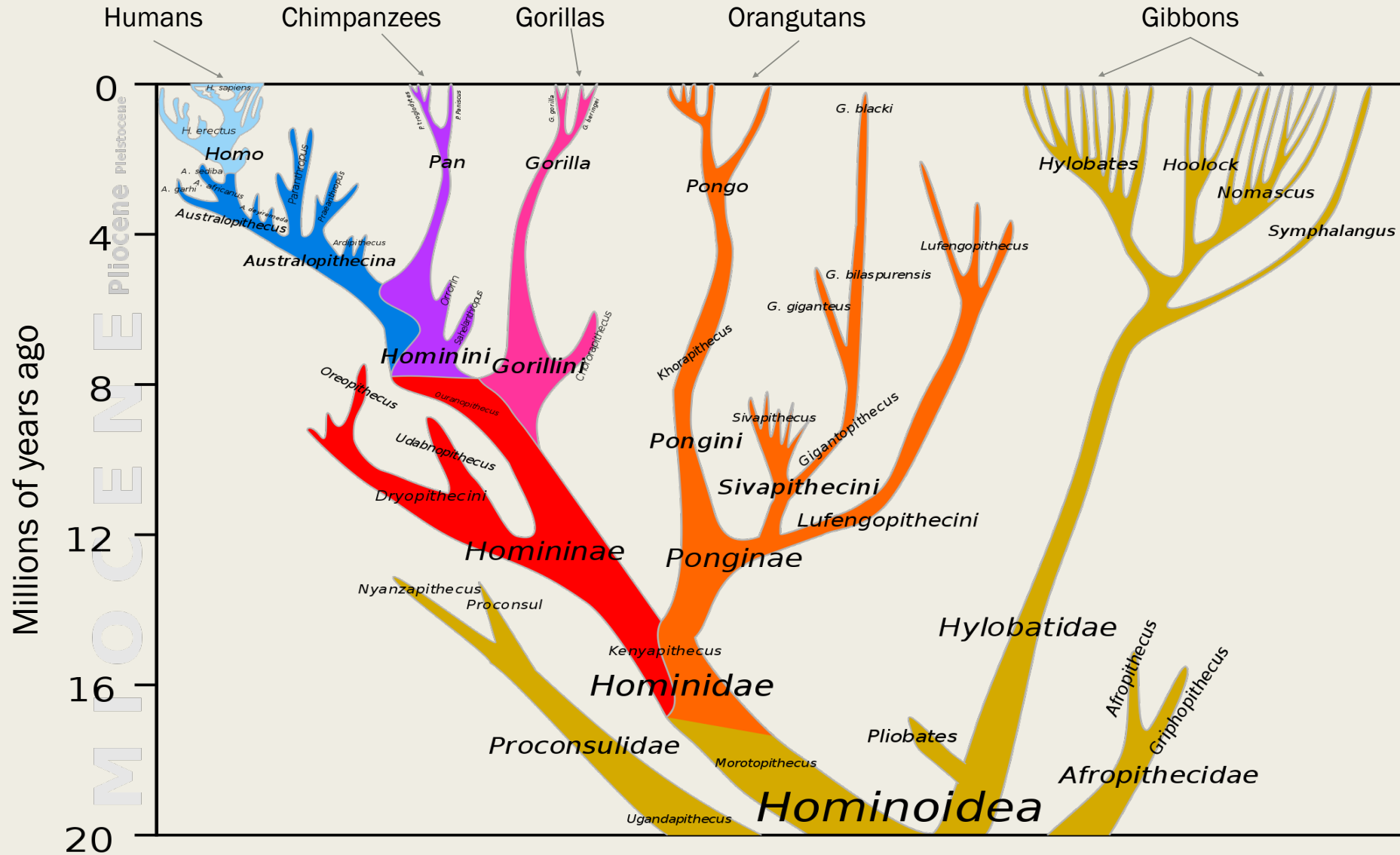
$$\frac{f_k(i) b_k(i)}{P(\vec{x})}$$

$$\sum_k P(z_i = k | \vec{x}) = 1 \approx \pm 10^{-6}$$

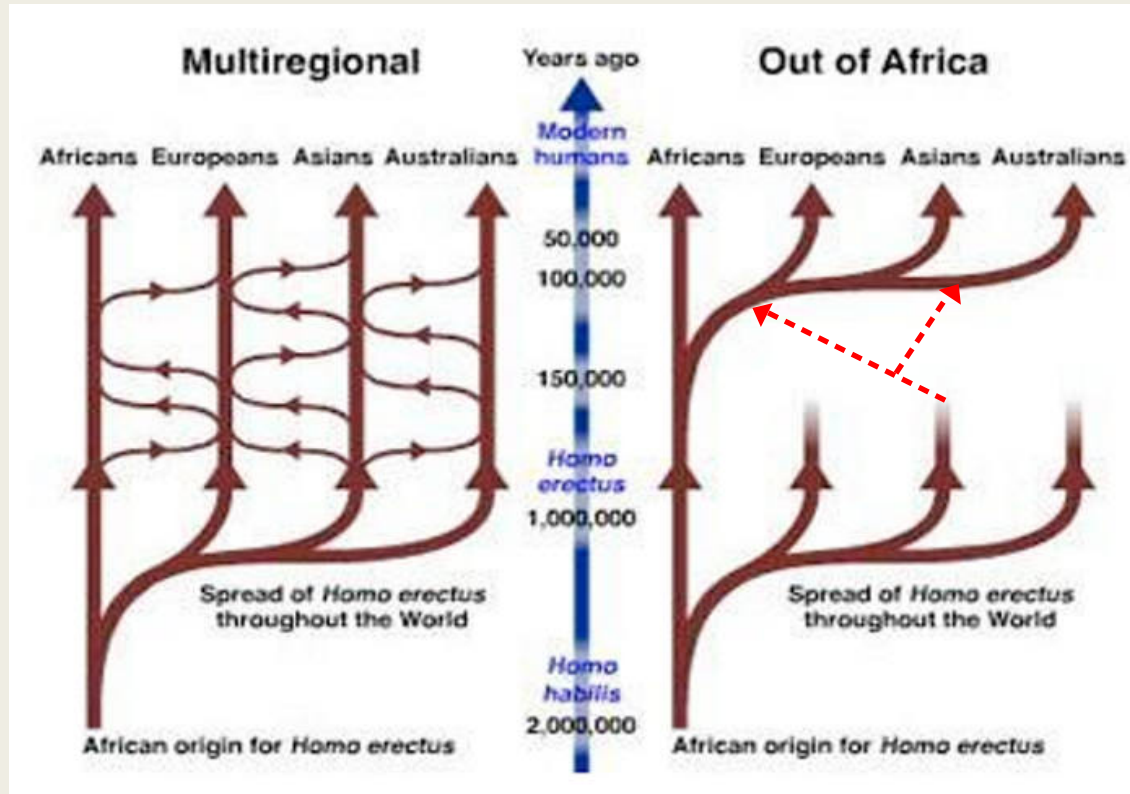
transform
back into non-log
space

Continue: human evolution overview

20 million years of Apes



Multi-regionalism vs Recent African Origin



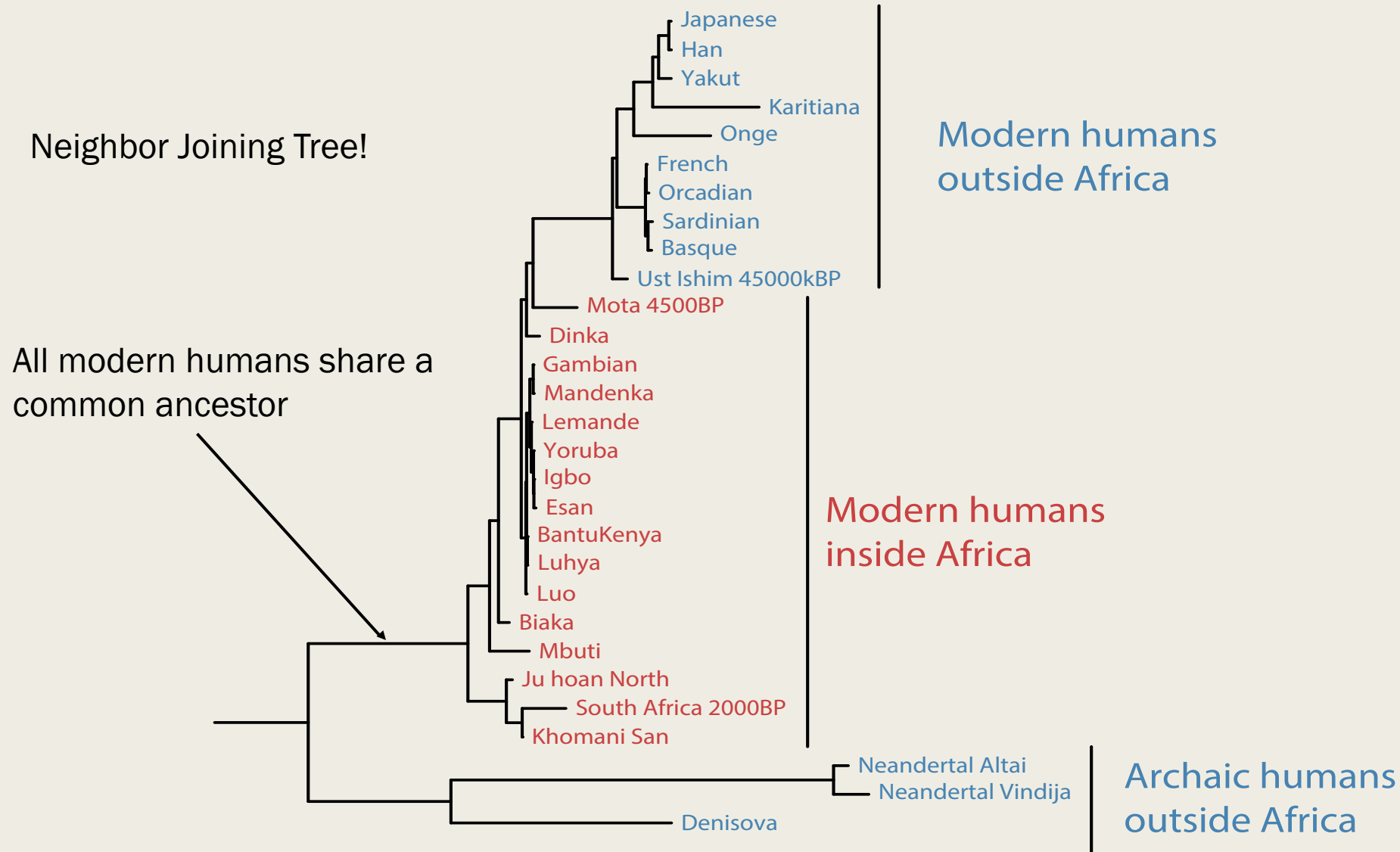
An old debate about human origins

Genetic data largely resolved this question in favor of RAO, starting in the late 1980's

Recent data, particularly ancient DNA slightly modifies this model.

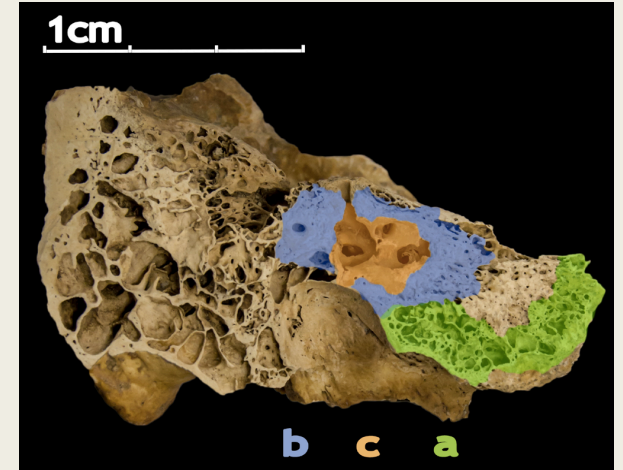
“Leaky replacement” – Svante Pääbo

DNA data confirms RAO (Recent African Origin)

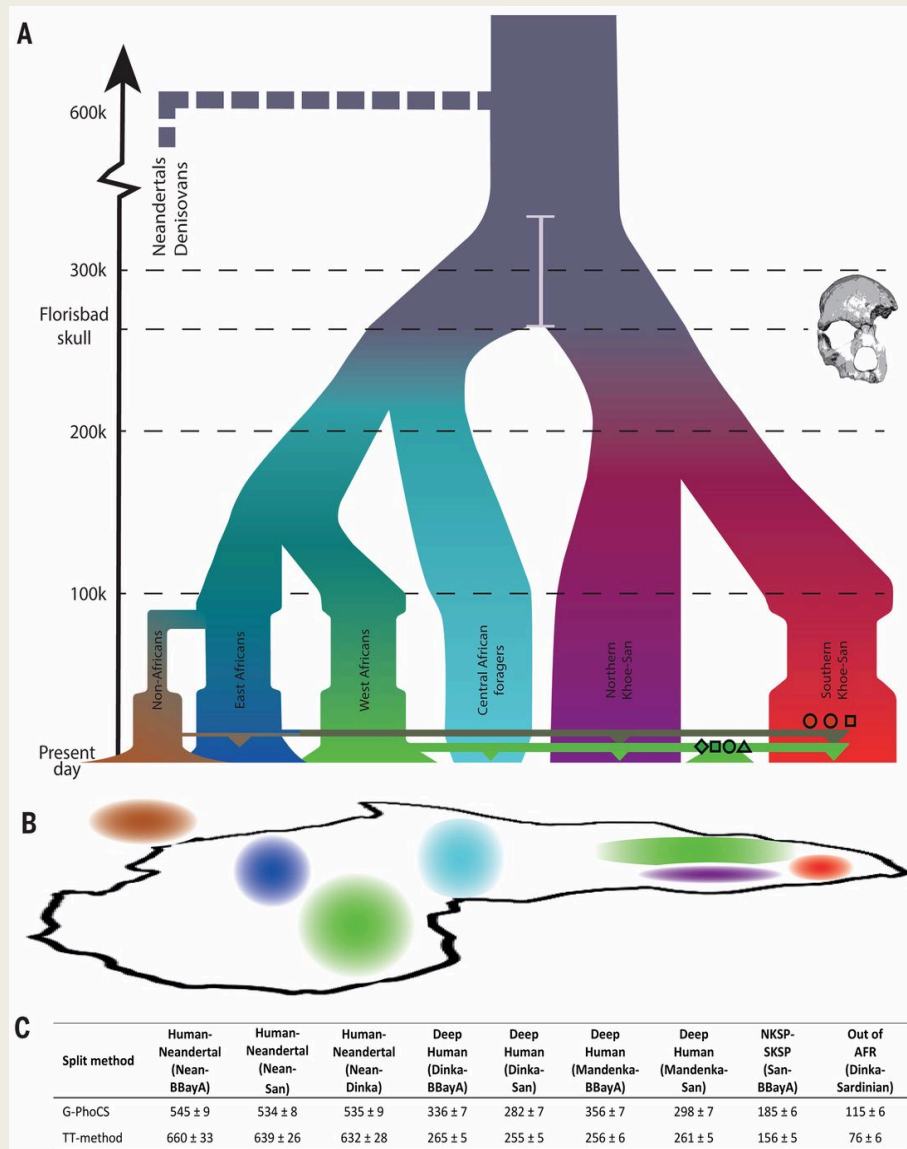


Ancient DNA

- Best preserved in cold and dry conditions, and dense parts of the skeleton (teeth, inner ear)
- Most DNA recovered is bacterial contamination
- When recovered:
 - fragmented to <200 bp molecules
 - cytosine deamination: sequence errors
 - Be paranoid about modern human contamination
- Upper limit?
 - ~700,000 year-old horse from permafrost in Alaska (Orlando et al. 2010)
 - Oldest human DNA – Sima Los Huesos (Spain). 430,000 year-old proto-Neanderthal (Meyer et al. 2016)
- As of today: 2 high coverage (5 low coverage) Neanderthal genomes, 1 high-coverage Denisovan genome, ~1500 ancient modern humans from 45,000 BP to present.



Ancient divergence in Africa



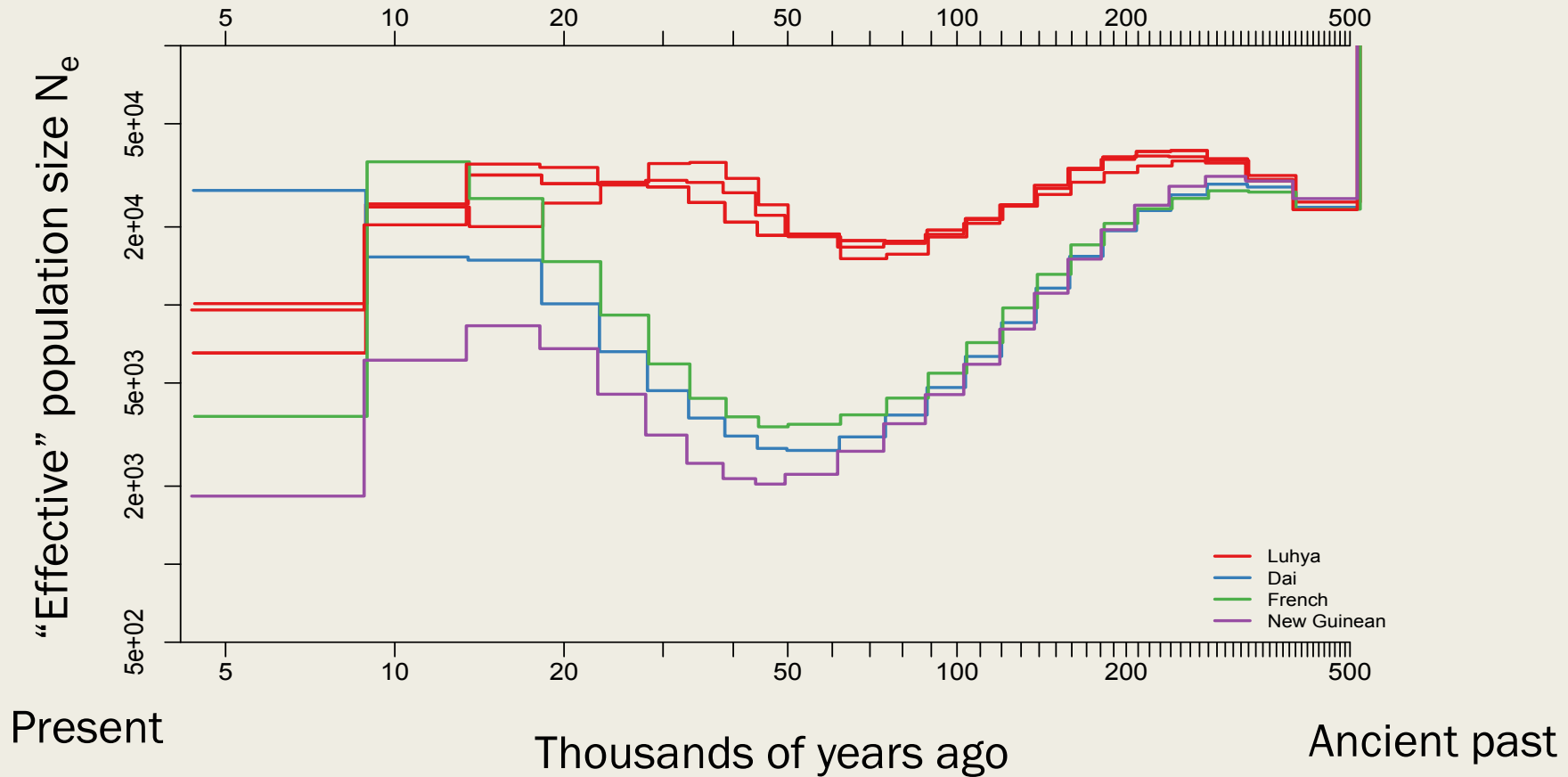
Oldest divergences among human populations found between present-day African populations. >250 KYA separation between Khoe-san and others

Other lineages are also deeply diverged, including those separating East-central and West-central African populations.

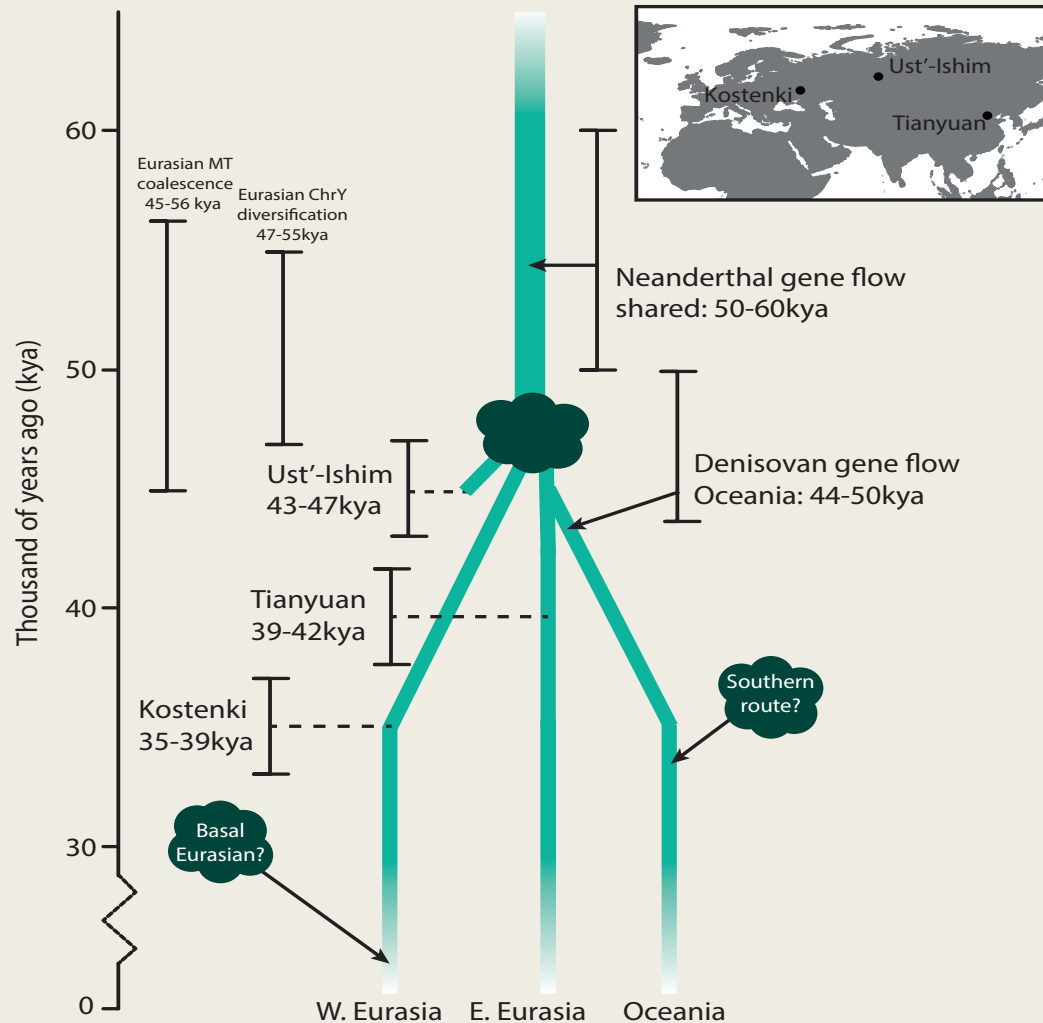
Non-Africans diverge ~70-120 KYA, from a population most closely related to present-day East Africans.

Recent (past 10KYA) population movements have affected this structure

Out-of-Africa bottleneck ~ 50 KYA



Relationship between present-day non-Africans

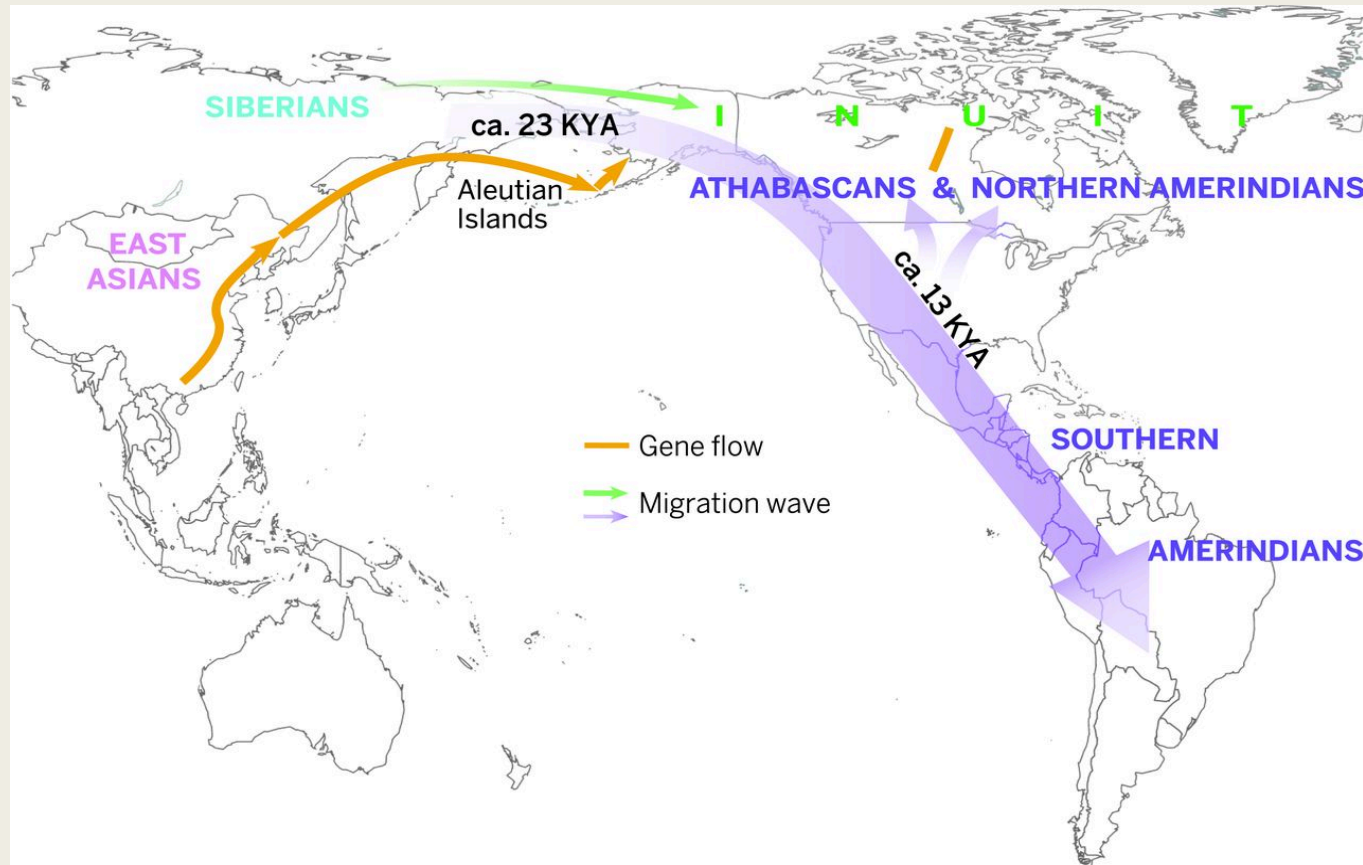


All present-day non-Africans share a common ancestor 45-60 KYA

Oceanian lineage diverged by 45 KYA

East and West Eurasian Lineages diverged by 40 KYA

The peopling of America



Raghavan et al. 2015

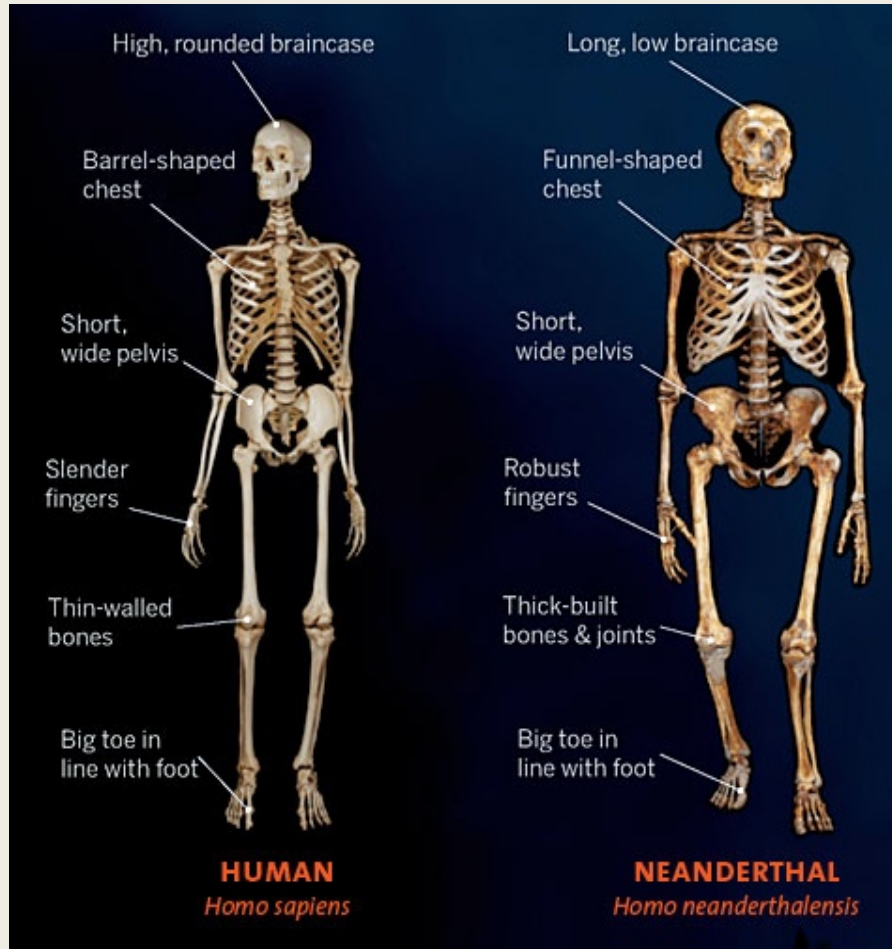
Last major landmass inhabited by humans.
At least 14,500 BP

Current model is that Native Americans derive from an admixture between populations related to present-day East Asians and Siberians, likely 10-25 KYA.

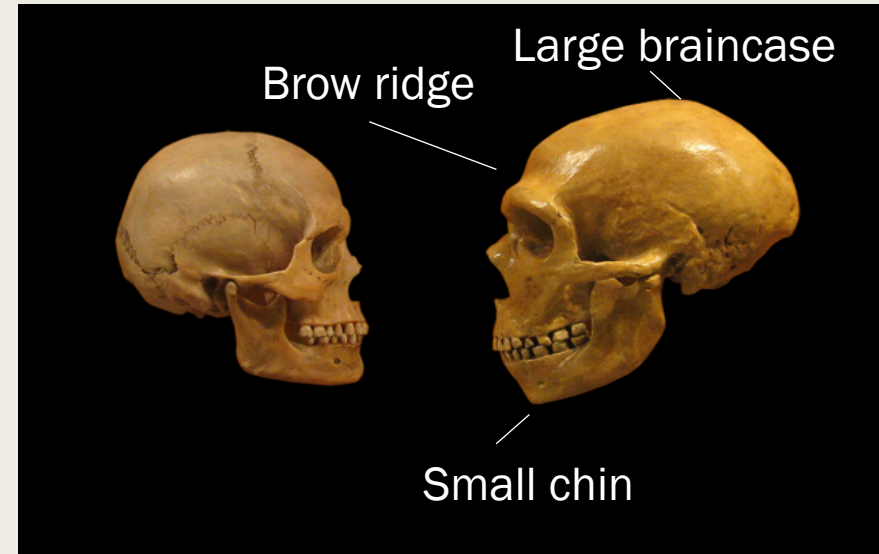
Later gene flow in the Arctic and farther south (Paleo-Eskimo, present-day Inuit, Athabascan-speakers).

Human-Neanderthal divergence

Neanderthals



Smithsonian

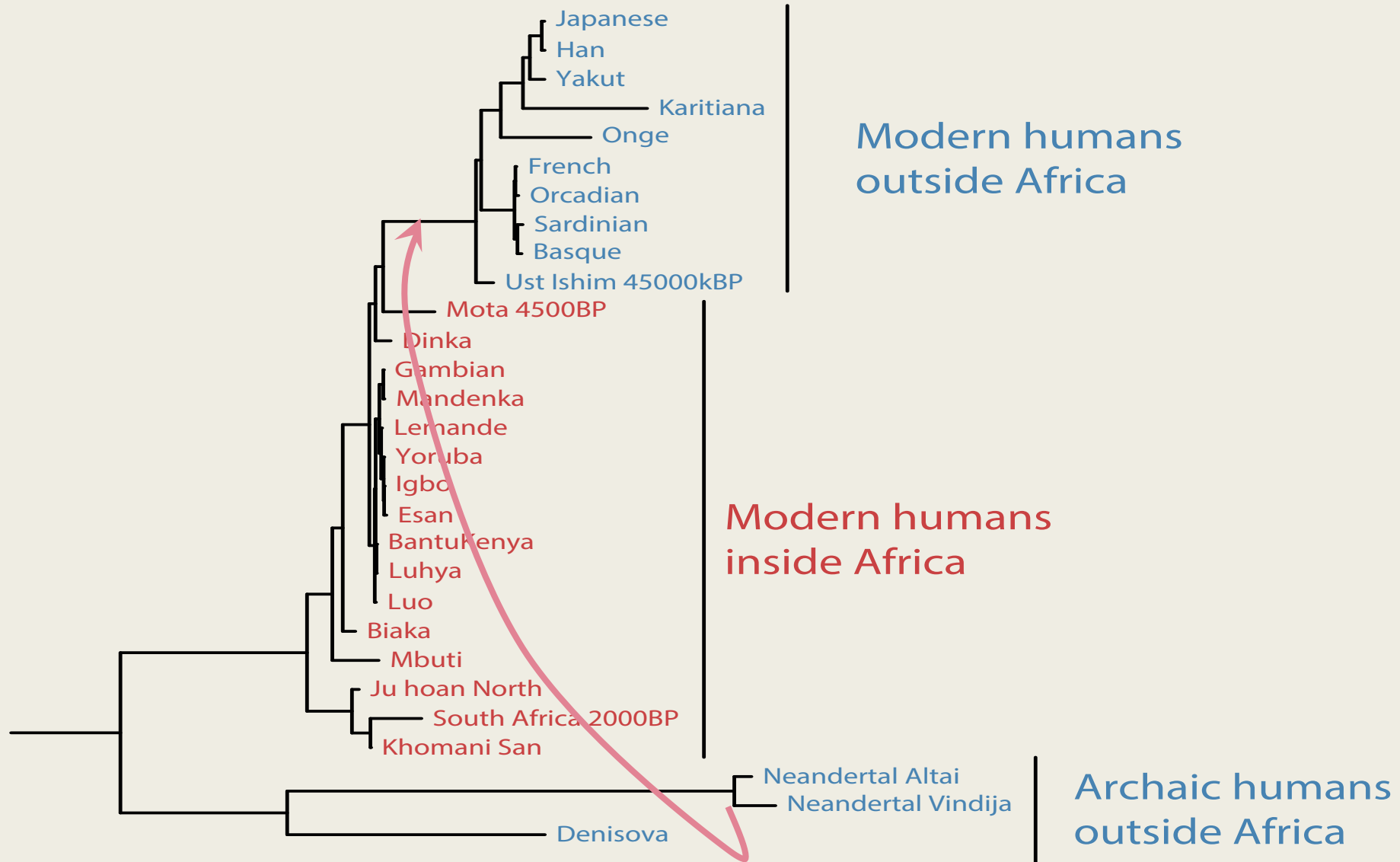


Remains of ~400 Neanderthals identified

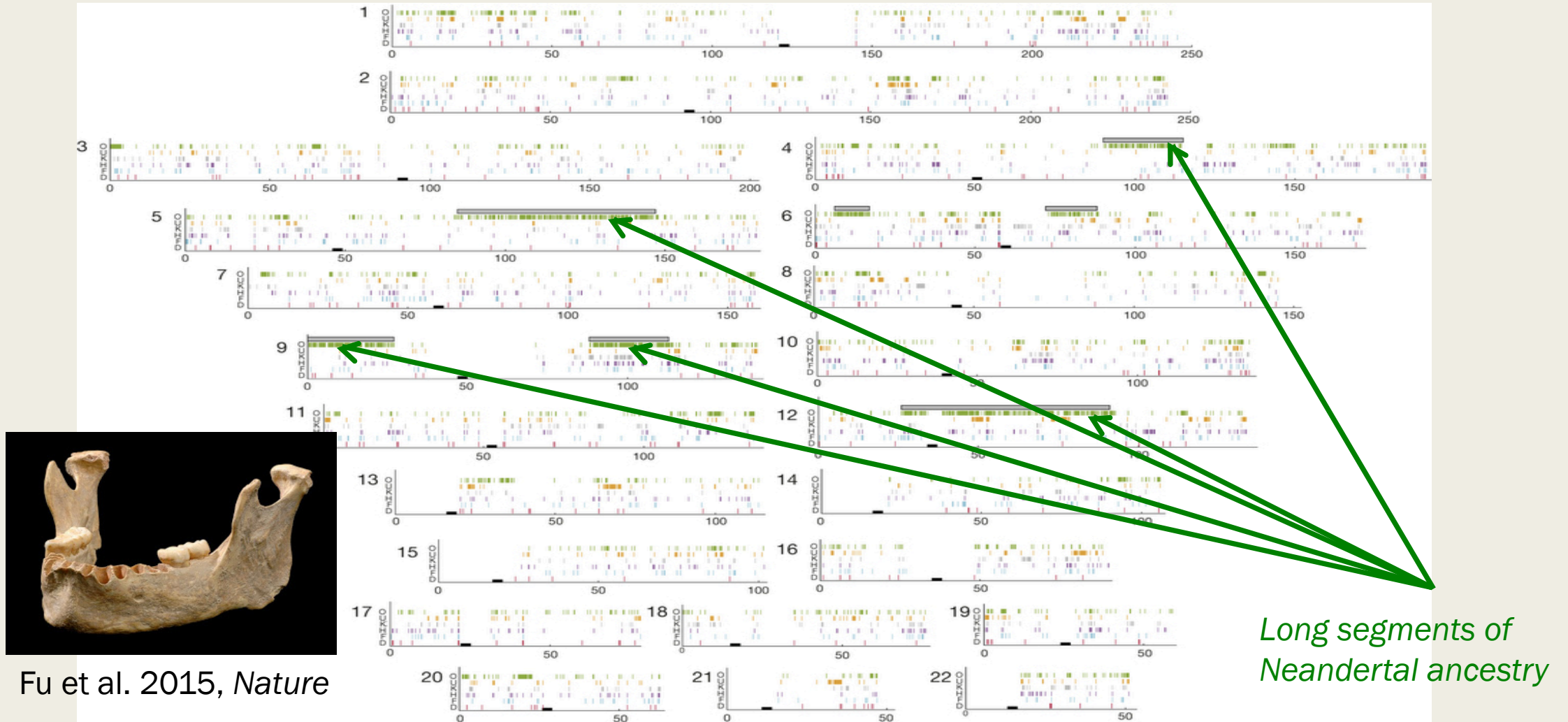
Dated from ~250 KYA to ~40 KYA

Estimated divergence from modern humans 400-800 KYA

Neanderthal variation falls outside modern humans

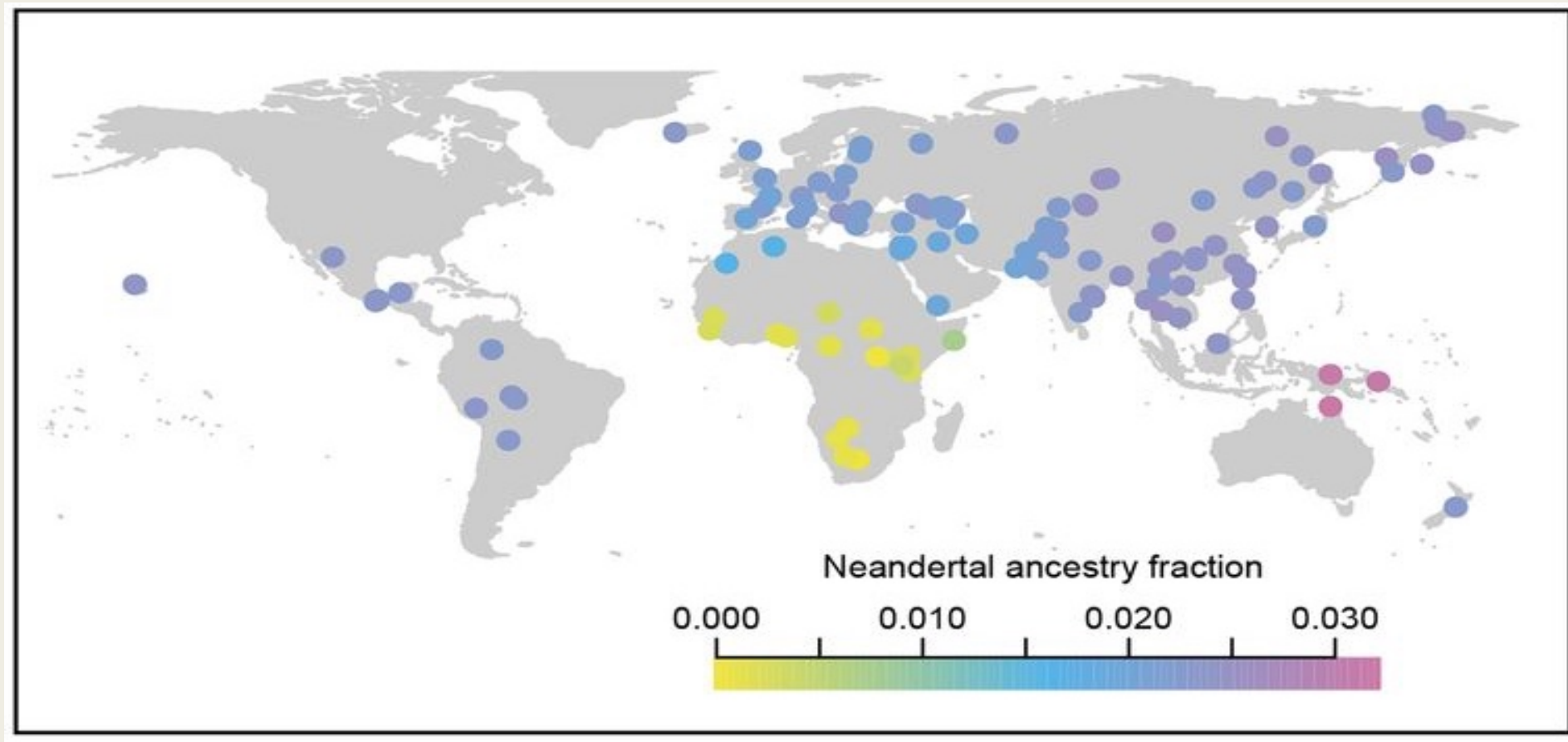


Oase 1: 40,000 year old European has a recent Neanderthal ancestor



Fu et al. 2015, *Nature*

Neanderthal ancestry varies across populations



Western Eurasia: 1.8-2.4%

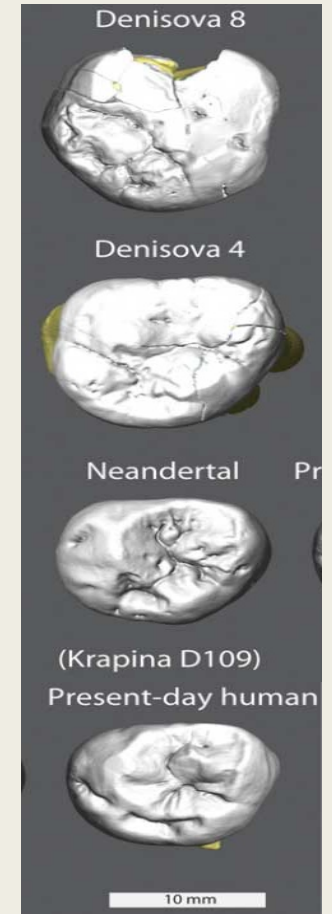
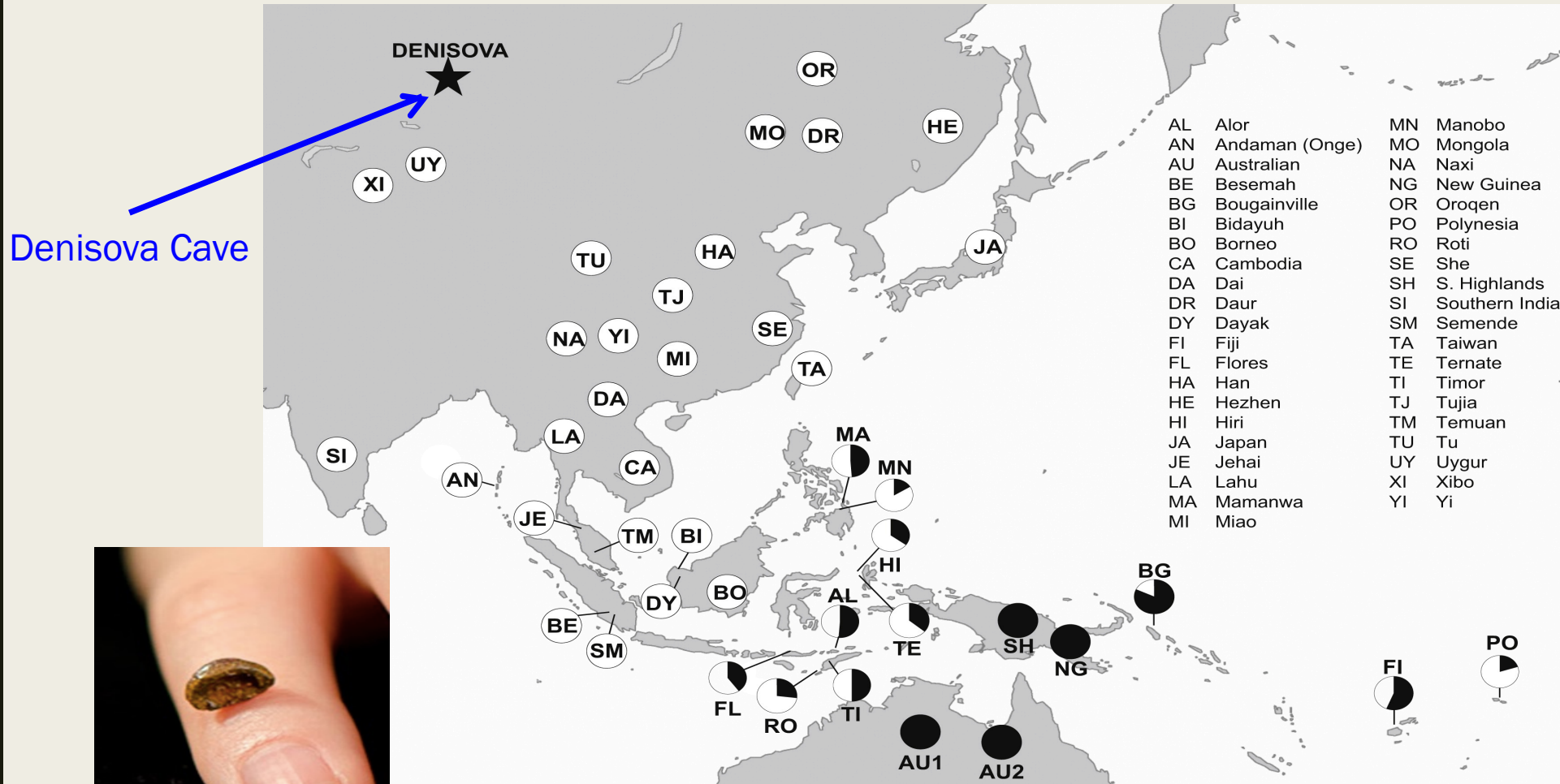
East Asia: 2.3-2.6%

Oceania - > 3% But, possibly confounded with Denisovan ancestry

Denisova Cave

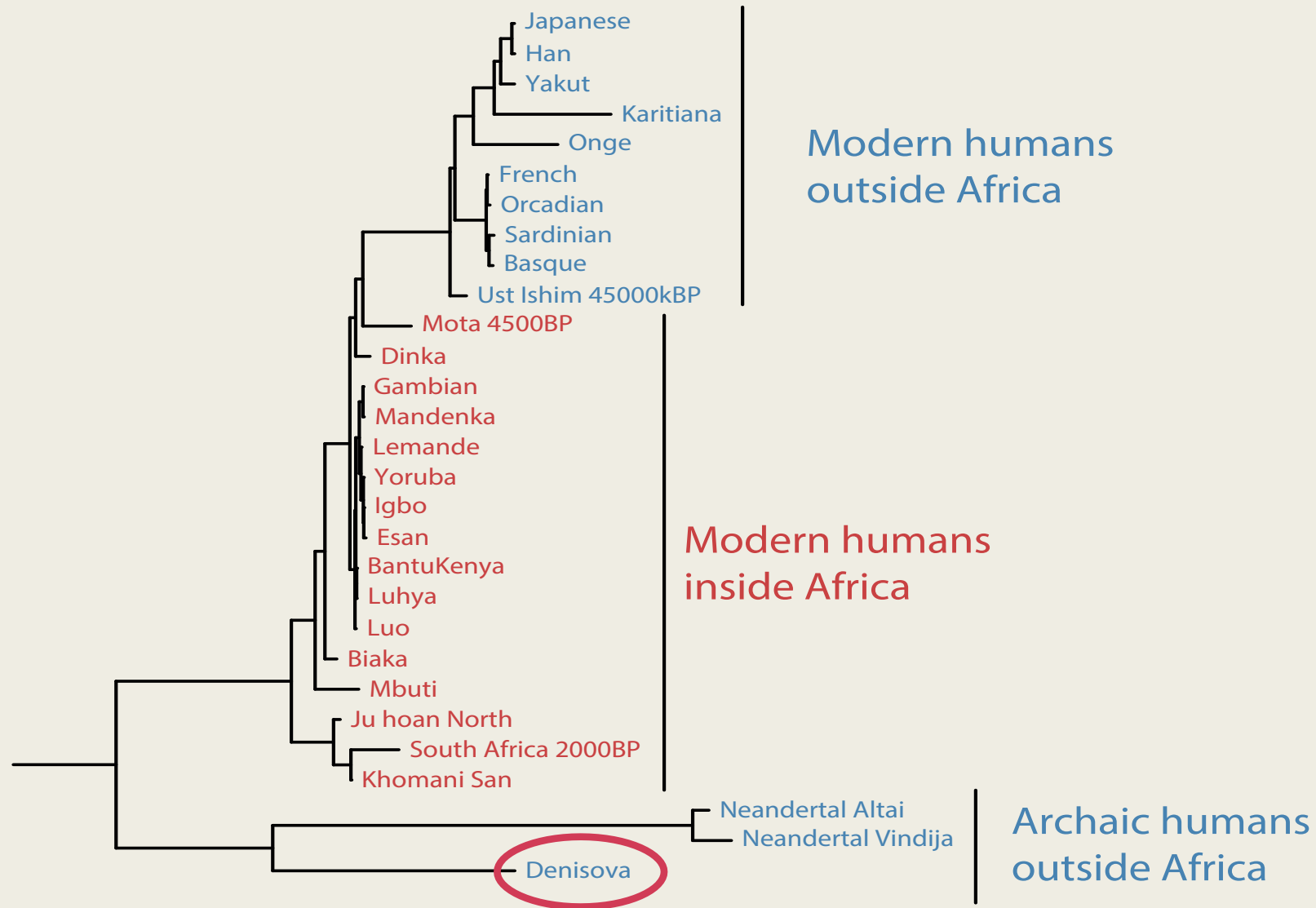


A new human group – Denisovans

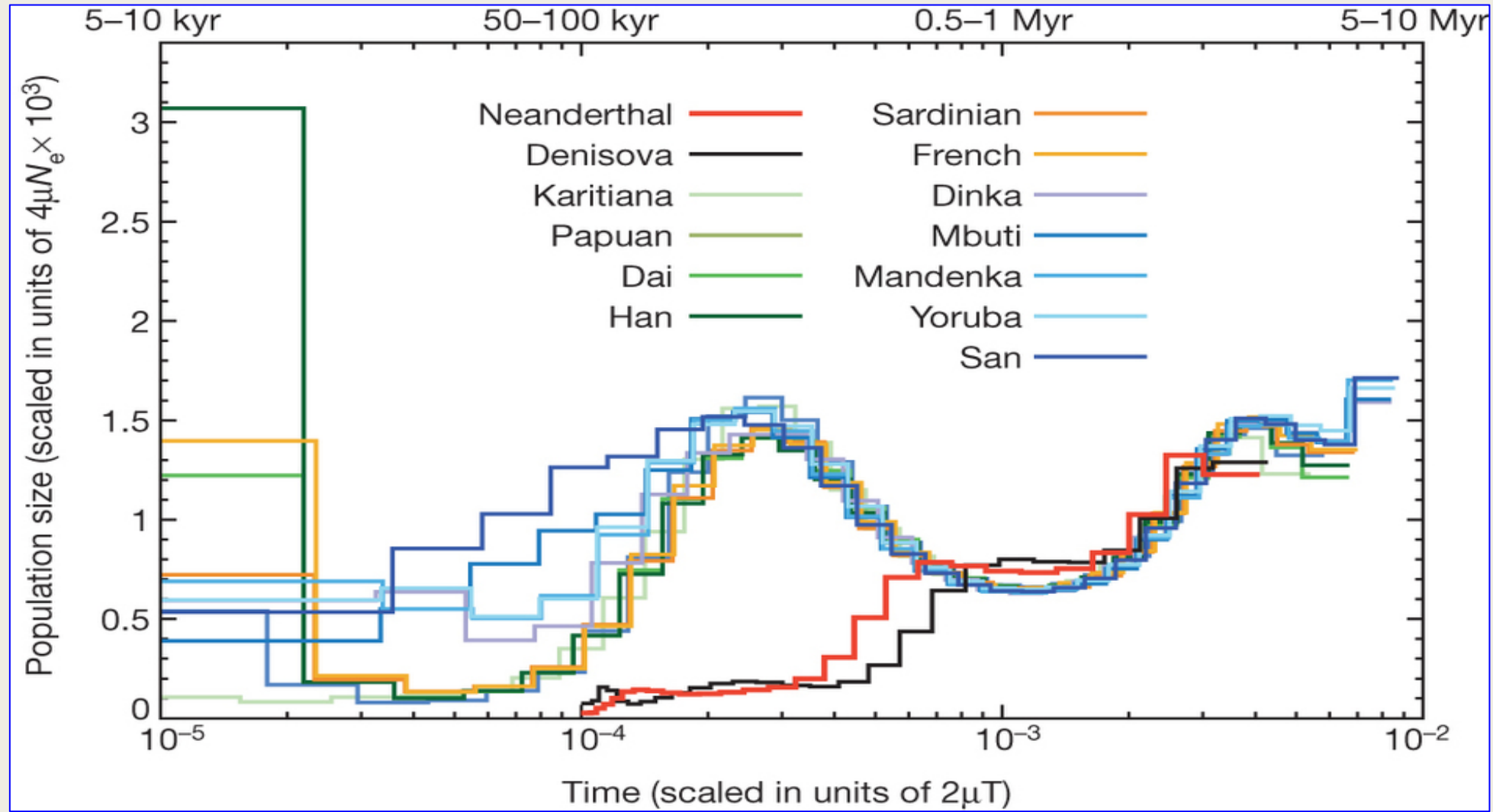


Reich et al. 2010, *Nature*; Reich et al. 2011 *AJHG*, Meyer et al. 2012 *Nature*. Sawyer et al. *PNAS* 2015

Denisovans are a sister group to Neanderthals

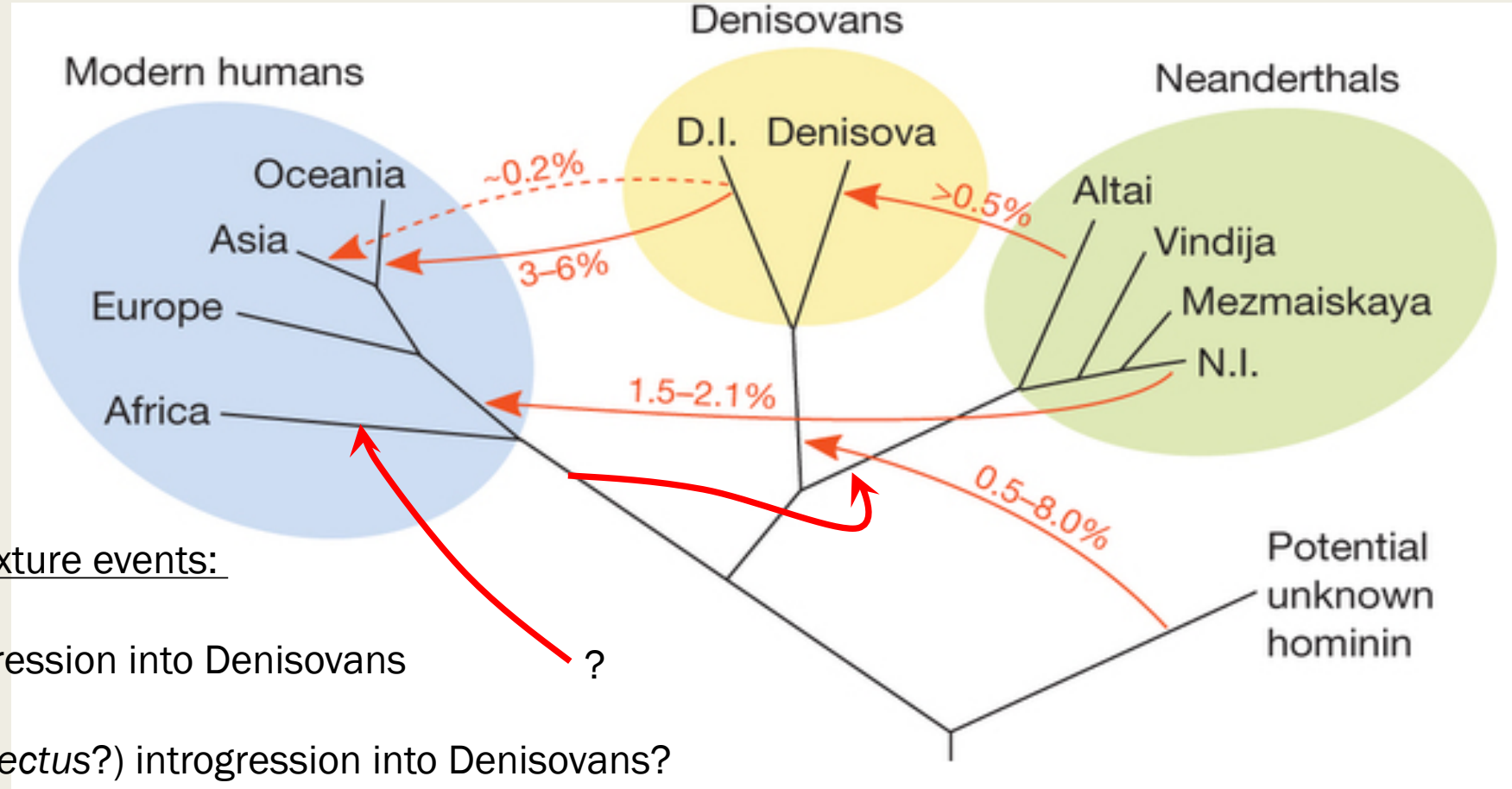


Neanderthals and Denisovans had very small population sizes



- Prüfer et al 2014
- PSMC: Li & Durbin 2011

Ghost populations and additional admixture



More complex admixture events:

- Neanderthal introgression into Denisovans
- Unknown (*Homo erectus*?) introgression into Denisovans?
- Unknown archaic introgression into west Africans?
- Proto-modern human introgression into Neanderthals??

Further reading

Archaic humans:

Prüfer et al. “The complete genome sequence of a Neanderthal from the Altai Mountains”
Nature, 2014

Meyer et al. “A High-Coverage Genome Sequence from an Archaic Denisovan Individual”
Science, 2012

Human history:

Nielsen et al. “Tracing the peopling of the world through genomics”
Nature Reviews Genetics, 2017

Skoglund & Mathieson. “Ancient genomics: a new view into human prehistory and evolution”
<https://t.co/YrZS45Q0xt>, 2017

Natural selection:

Fu & Akey. “Selection and Adaptation in the Human Genome”
Annual Review of Genomics and Human Genetics, 2013

Begin: PCA

P(A)

input $n \times p$ matrix, X_{orig}
 $n = \# \text{ samples } (2 \times \# \text{ individuals})$
 $p = \# \text{ features } (\text{SNPs, sites})$

$$X_{orig} = \begin{bmatrix} 0 & 1 \\ 0 & 1 \\ 0 & 1 \\ \vdots & \vdots \\ \vdots & \vdots \end{bmatrix}$$

$n=6$,
 $p=2$

Step 1: get data

output:

- feature combinations
that best explain
observed variation
in the data

Step 2

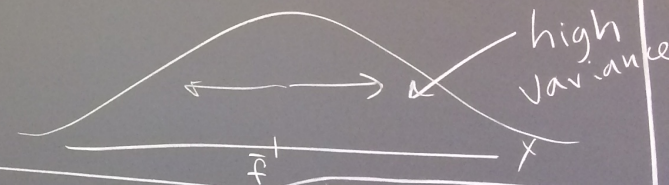
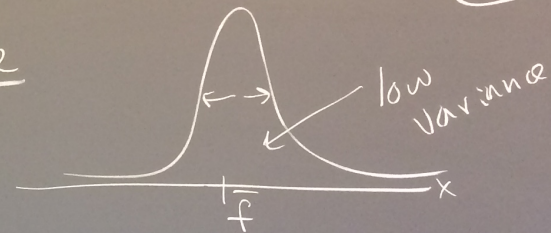
Subtract off column-wise mean

$$\bar{f} = \frac{1}{n} \sum_{i=1}^n f_i, \text{ ex: } \bar{f} = \frac{1}{2}, \bar{g} = \frac{1}{2}$$

$$X = \begin{bmatrix} -\frac{1}{2} & \frac{1}{2} \\ -\frac{1}{2} & \frac{1}{2} \\ -\frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & -\frac{1}{2} \\ \frac{1}{2} & -\frac{1}{2} \\ \frac{1}{2} & -\frac{1}{2} \end{bmatrix}$$

Step 3: Compute
covariance

Variance



$$\text{Var}(f) = \frac{1}{n-1} \sum_{i=1}^n (f_i - \bar{f})^2$$