



Capacity building for bioinformatics in Latin America



FACULTAD DE
MEDICINA HUMANA

ANÁLISIS Y FILOGENIA DE POBLACIONES HUMANAS

Chiara Barbieri, University of Zurich, Suiza

Max Planck Institute for the Science of Human History, Jena, Alemania



MAX-PLANCK-GESELLSCHAFT



University of
Zurich^{UZH}



Wenner-Gren
Foundation

barbieri.chiara@gmail.com

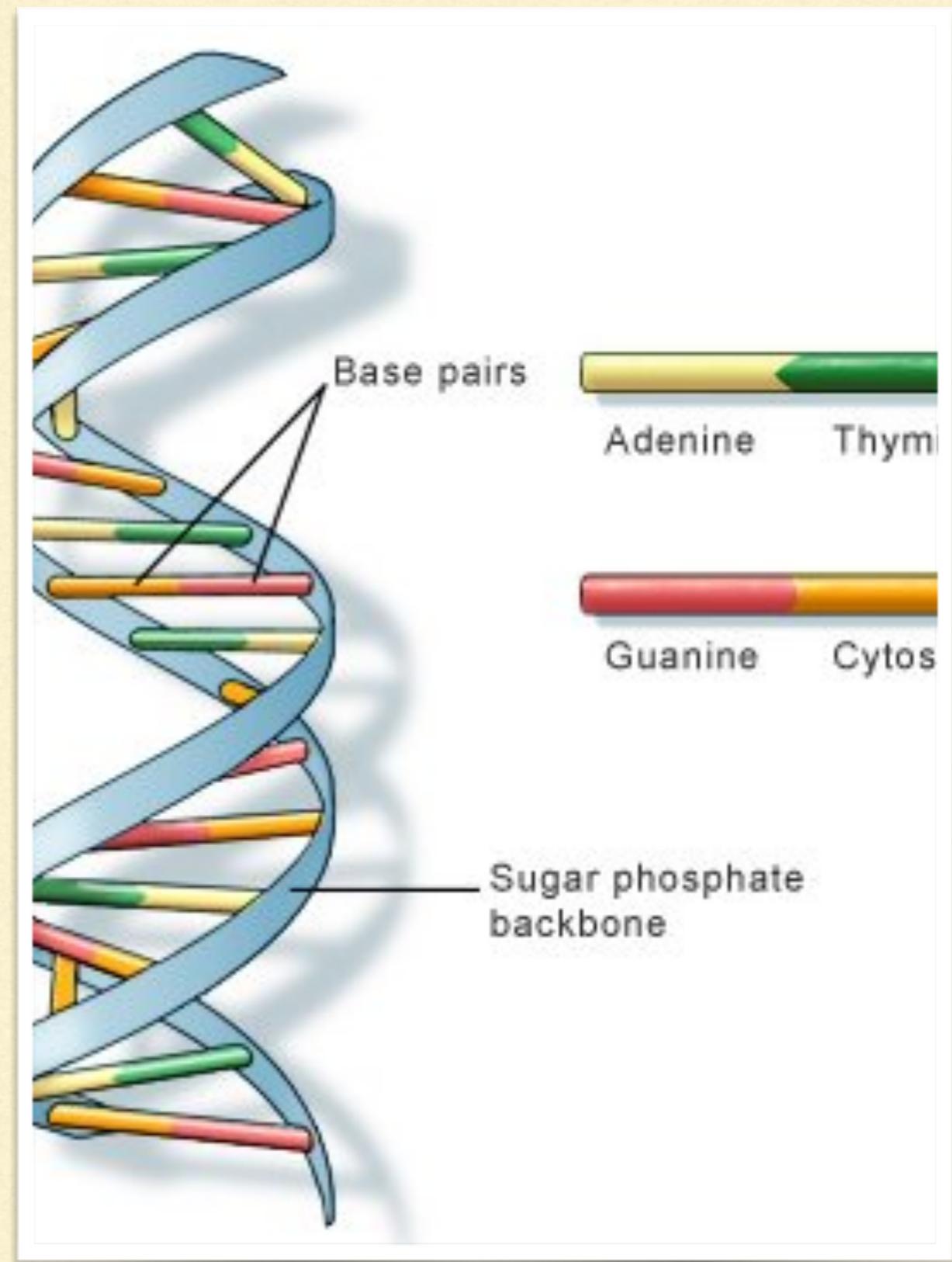
OUTLINE

- Marcadores genéticos en estudios de diversidad humana
- Population structure, diversidad y admixture
- Filogenia
- Case study: historia de poblaciones en Africa

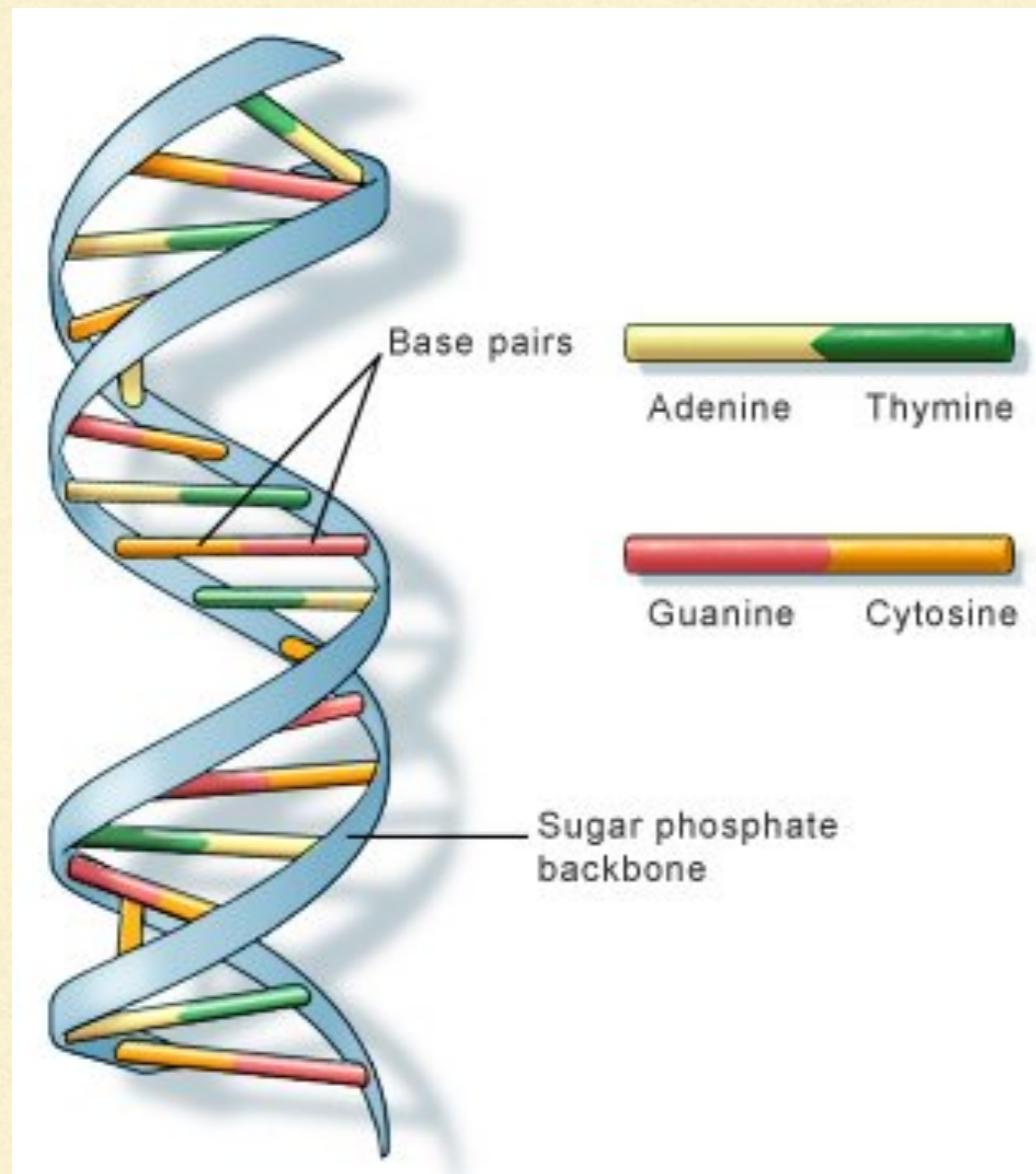
EJERCICIO:

- https://github.com/chiarabarbieri/CABANA_Lima_2019
-

MARCADORES GENÉTICOS EN ESTUDIOS DE DIVERSIDAD HUMANA



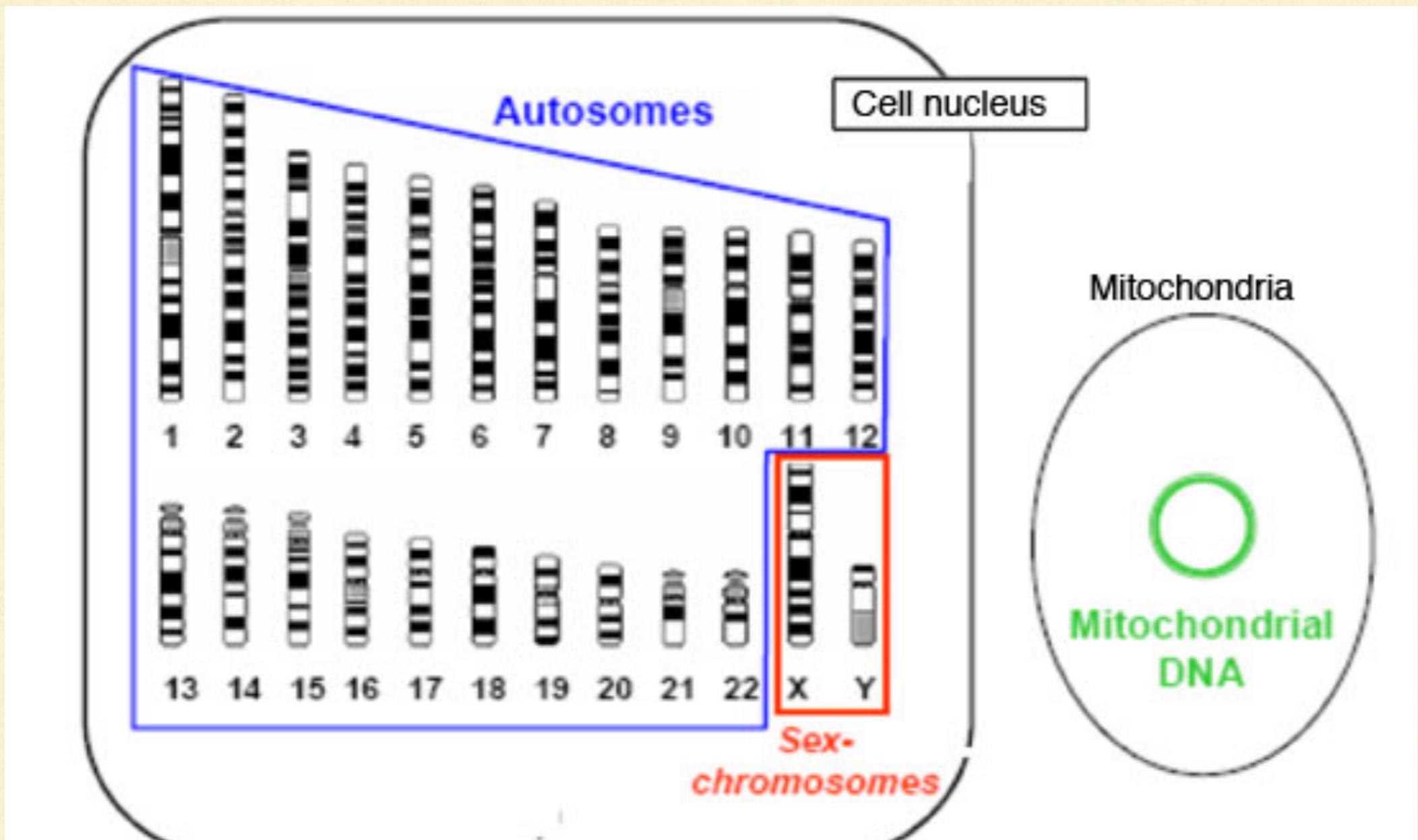
GENETICS AND POPULATION HISTORY



- **Classical markers**
 - little information
- **Molecular markers**
- **Autosomal DNA**
 - ancestors history
 - SNPs
 - Genome sequencing
- **Uniparental DNA (mtDNA, Y chromosome)**
 - mother and father's lineages
 - SNPs
 - Sequences

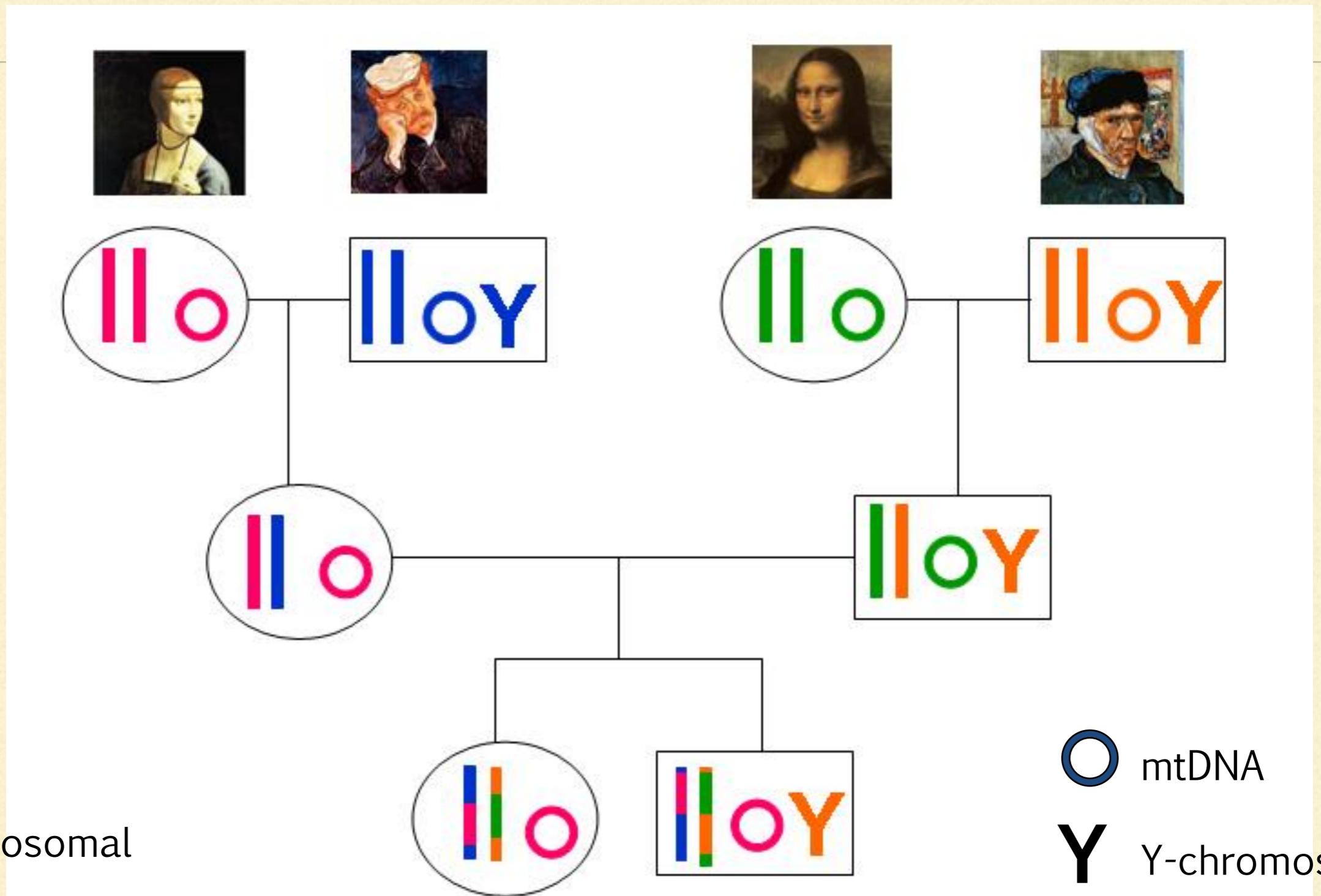
THE HUMAN GENOME

23 pairs of Chromosomes + mtDNA



Butler, J.M. (2006) *Forensic DNA Typing*, 2nd Edition, Figure 2.3, ©Elsevier Science/Academic P

GENETIC MARKERS: TRANSMISSION



MARCADORES GENÉTICOS AUTOSOMICOS

Autosomal markers:

- Large amount of genetic information
- Unbiased view of population prehistory (*virtually all the ancestors are taken into account*)

LIMITS

- Recombination prevents the tracing back of individual mutations through space and time
- Can be more cost-intensive

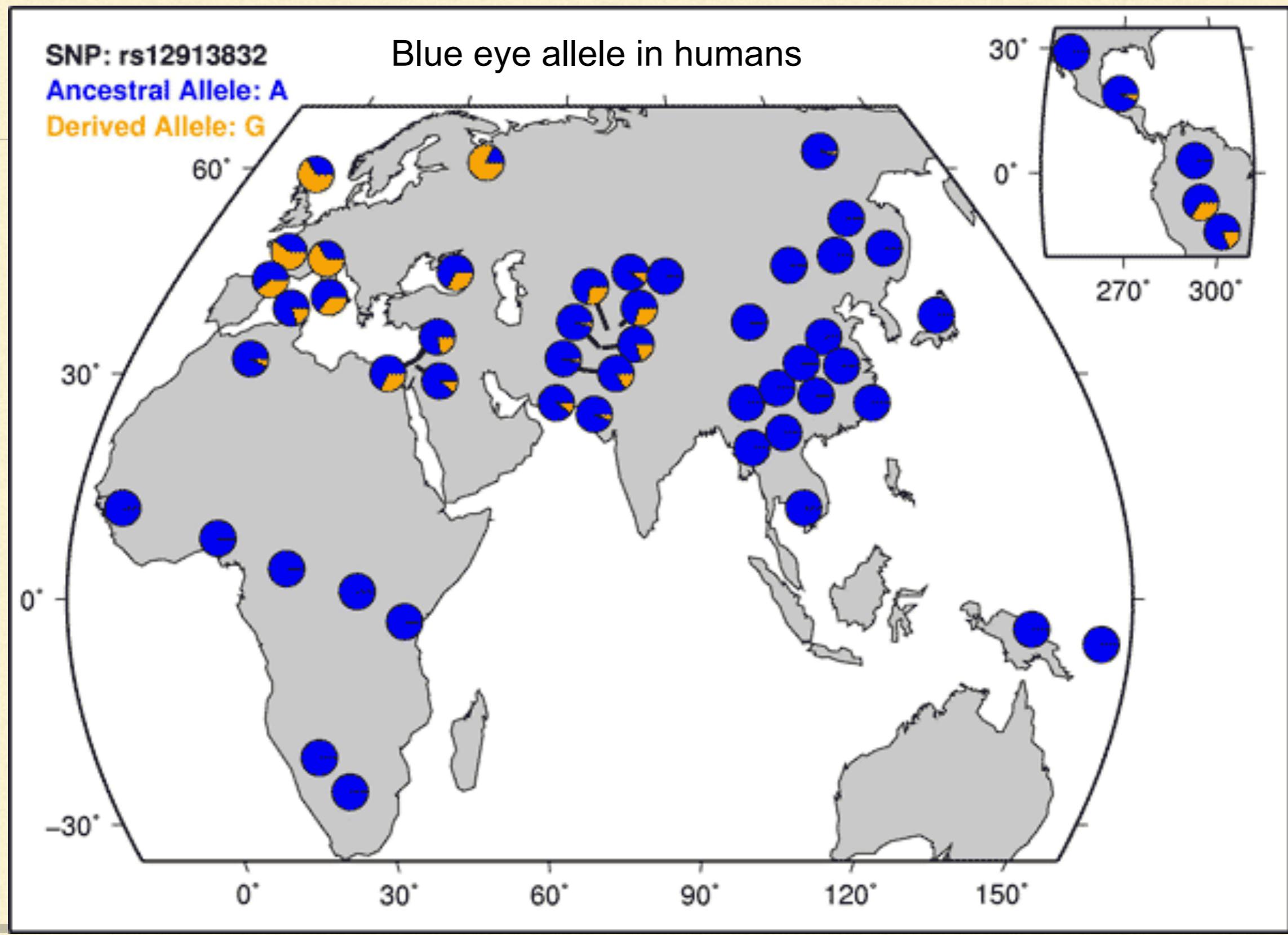
AUTOSOMAL DATA

- Focus on one gene or variant
 - Population variation
 - Reconstruct the phylogeny
- Capture a set of independent positions through the genome: SNP chip
 - Ascertained to be variable in human population
-

AUTOSOMAL DATA

- Focus on one gene or variant
 - Population variation
 - Reconstruct the phylogeny
- Capture a set of independent positions through the genome: SNP chip
 - Ascertained to be variable in human population
-

PRESENT DAY VARIATION



PHYLOGENY OF THE LACTASE ALLELE IN SOUTHERN AFRICA

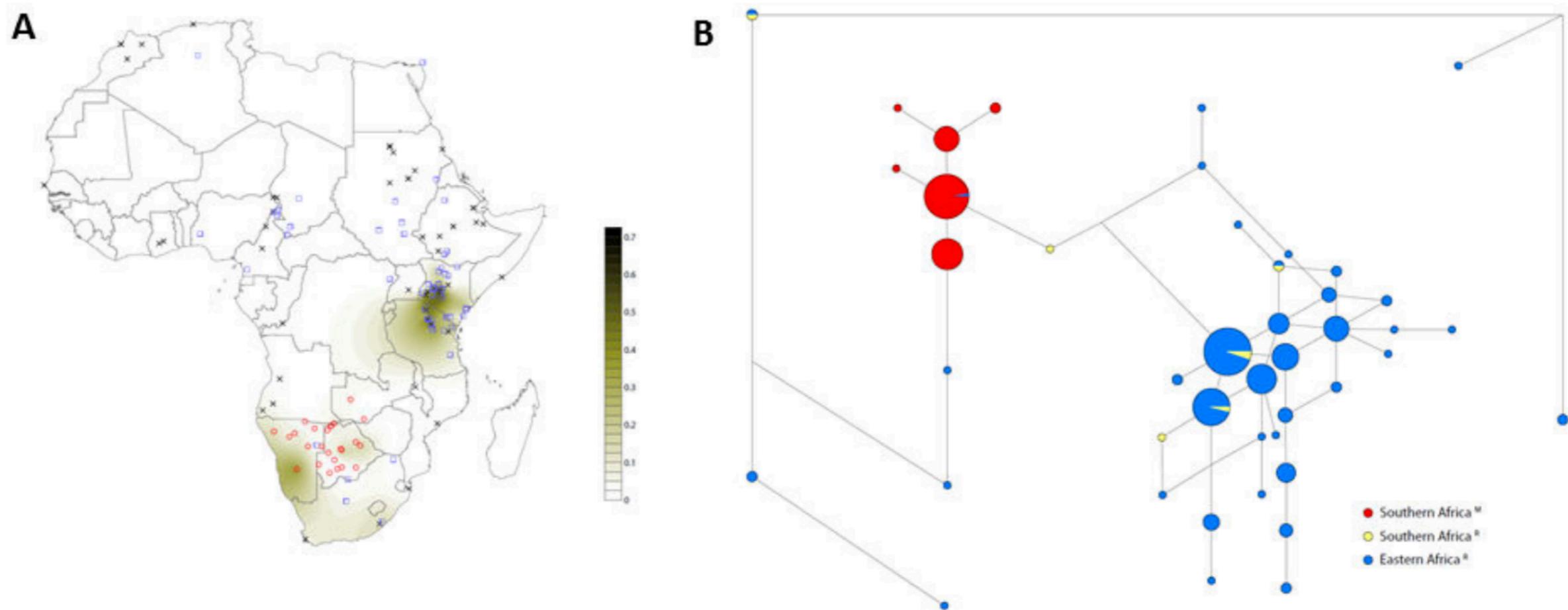


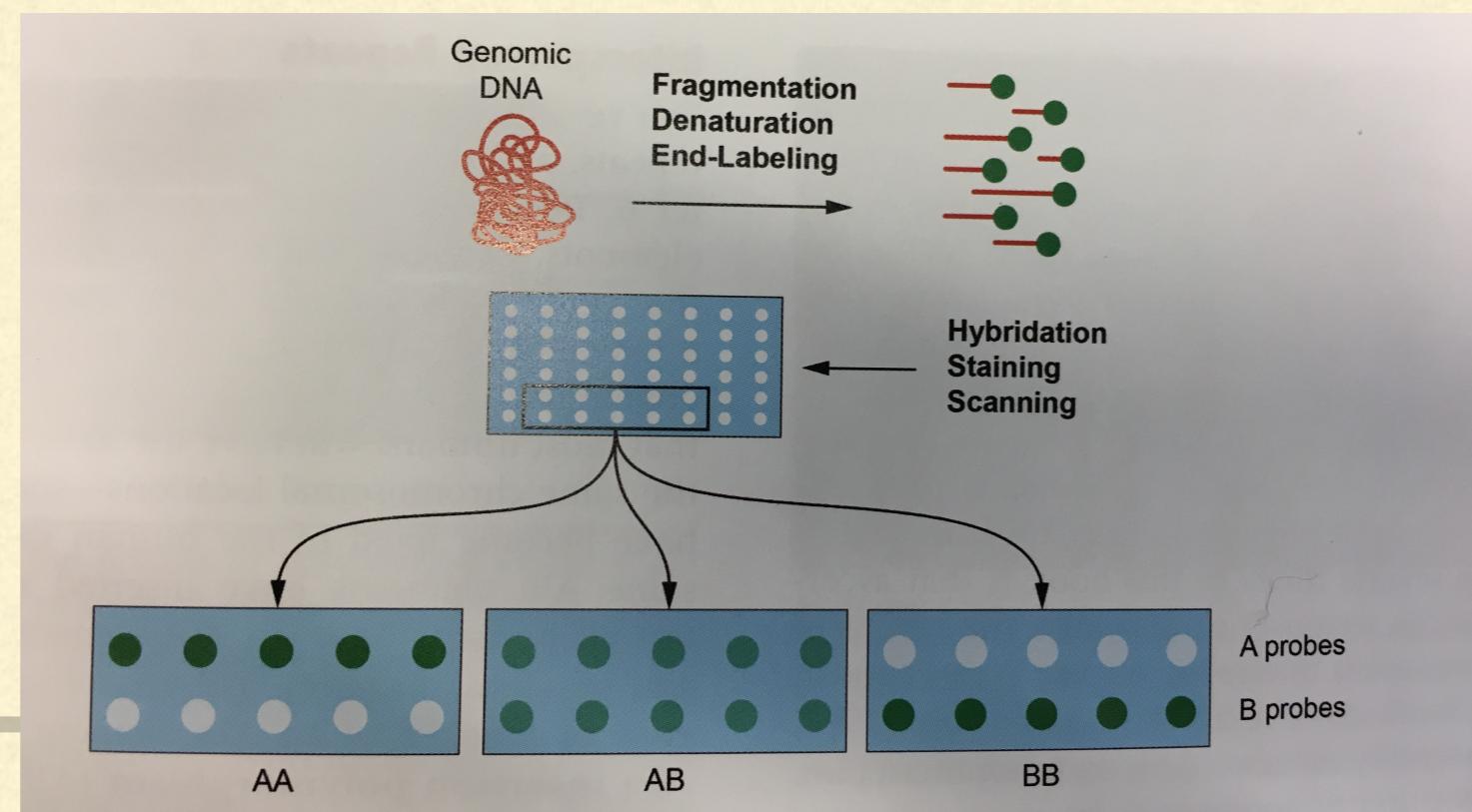
Fig. 1. Analyses of the C-14010 LP variant and associated STR haplotypes in Eastern and Southern African populations. **A:** Surfer map of the C-14010 allele frequency. Red circles denote sampling locations by Macholdt et al.; blue squares denote sampling locations by Ranciaro et al.; black crosses denote data from other published studies (taken from Macholdt et al. and Ranciaro et al.). **B:** Median-joining network of haplotypes associated with the C-14010 variant, based on four STR loci that flank the LP enhancer region. The M superscript denotes data by Macholdt et al., and the R superscript denotes data by Ranciaro et al. [Color]

AUTOSOMAL DATA

- Focus on one gene or variant
 - Population variation
 - Reconstruct the phylogeny
- Capture a set of independent positions through the genome: SNP chip
 - Ascertained to be variable in human population
-

ANALIZAR SNPs

- SNP chip (early 2000) : microarray that contain synthetic DNA probes to the allele for a selection of SNPs
- Fragmented labelled DNA hybridizes to the DNA probes, the rest is washed away
- Only the probes bound to DNA are read and their position in the array is recognised



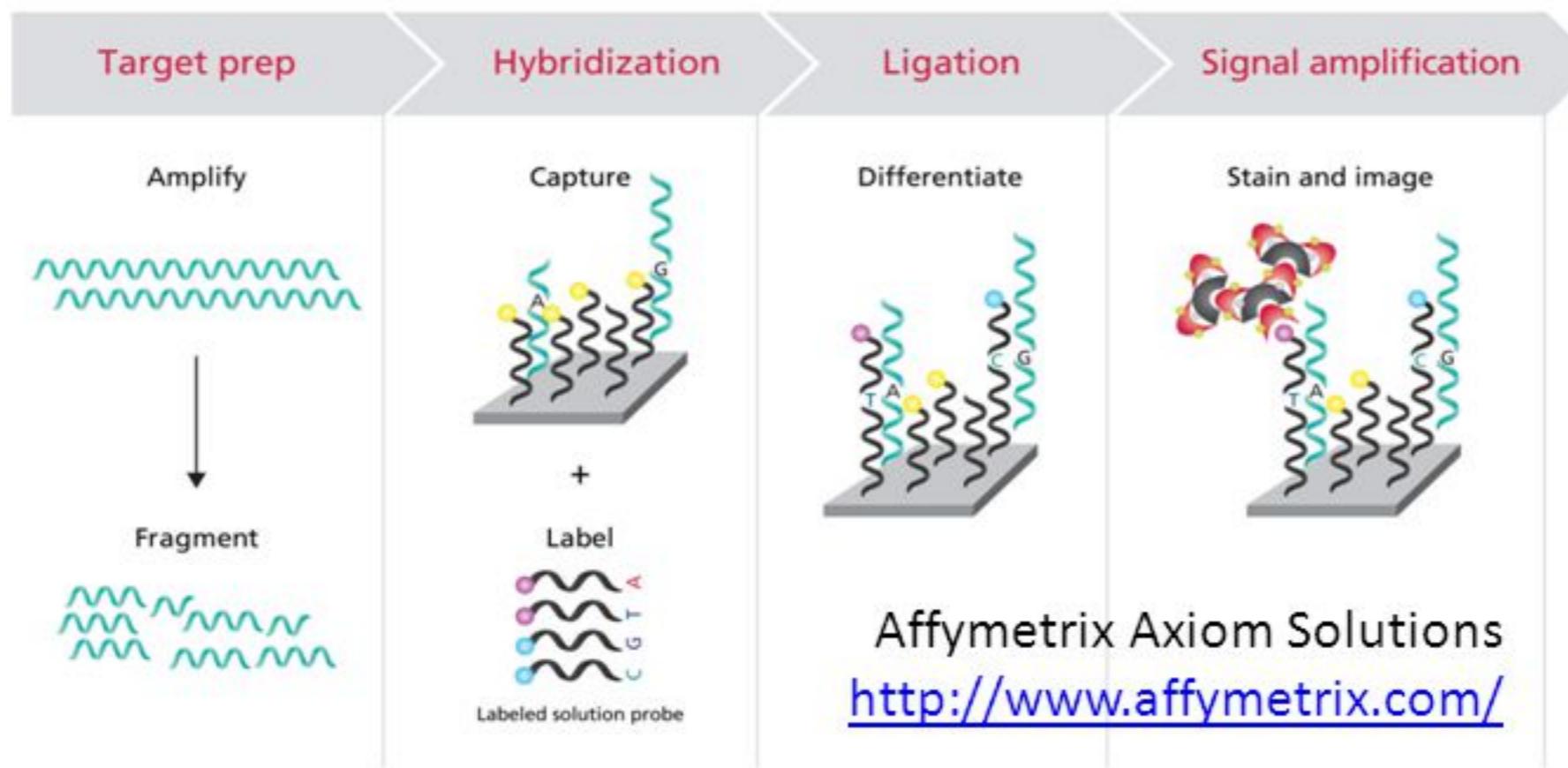
Stoneking 2016

SNP arrays

- Affymetrix

- Illumina

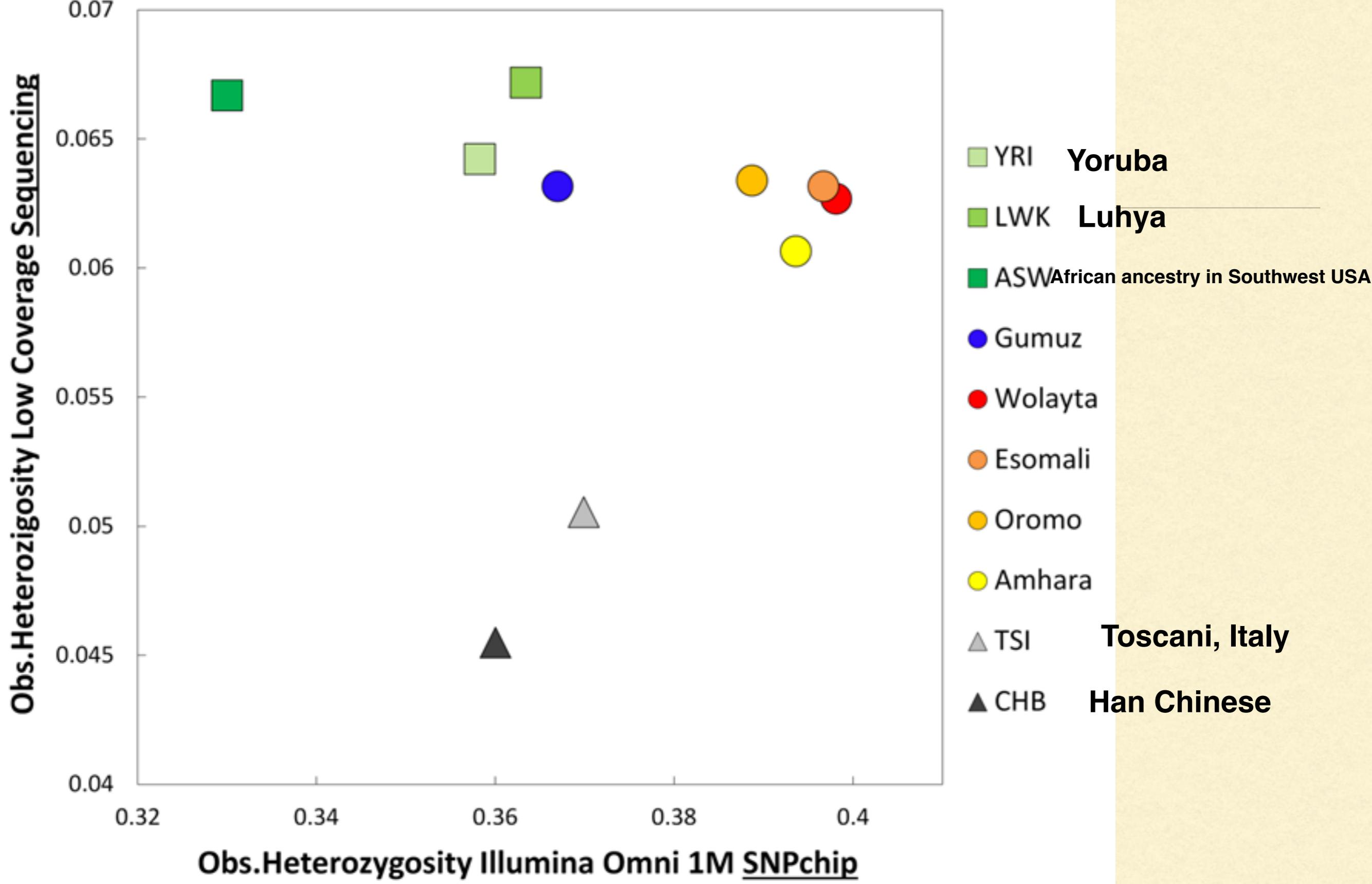
- Probes on microarray technology



Credits: J. Chen

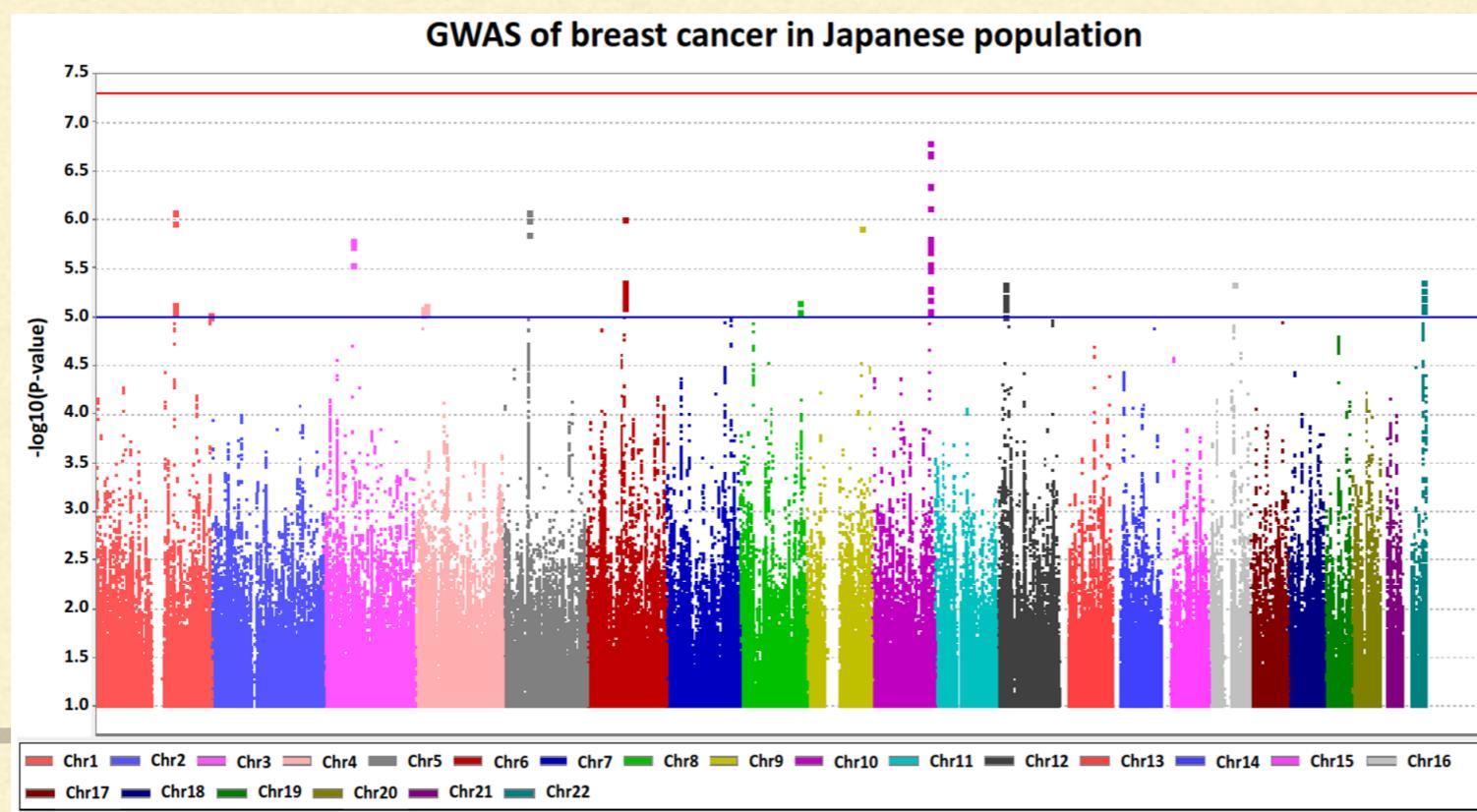
SNP CHIPS AND ASCERTAINMENT BIAS

- The diversity recovered does not correspond to a true measure of diversity in terms of population genetics
- can be caused by sampling a nonrandom set of individuals or by biased SNP discovery protocols.
- SNPs have been identified by sequencing European individuals - only known SNPs can be genotyped!!
- pre-ascertained SNPs are biased toward older SNPs



GENOME WIDE ASSOCIATION STUDIES - GWAS

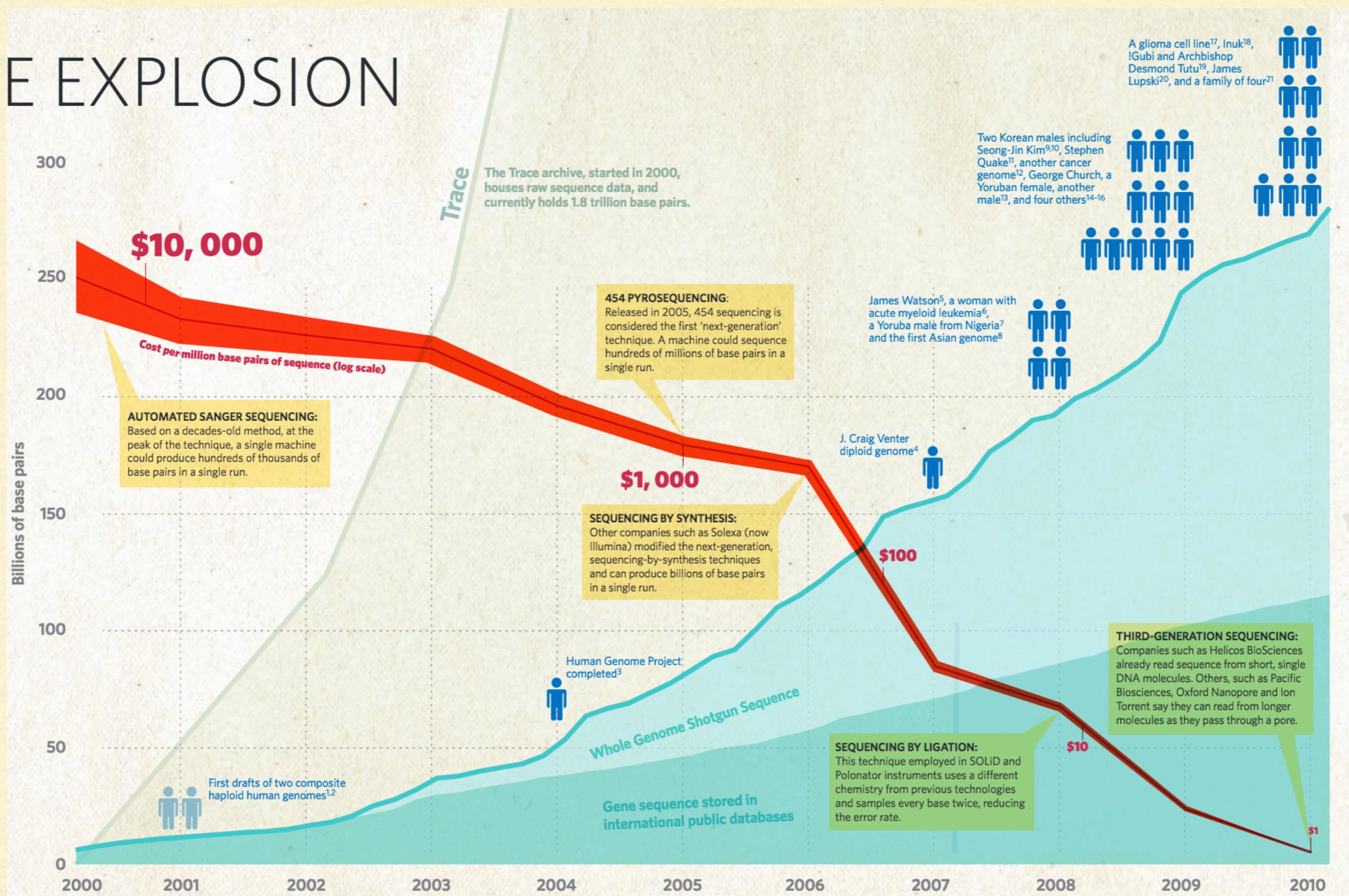
- Based on SNP chip (not really genome wide)
- Hypothesis free testing for association to SNP variants and diseases (or other traits)
- Problems: replicate failures, false positives, missing rare variants, Linkage Disequilibrium, population specific variants



SECUENCIA DE ADN

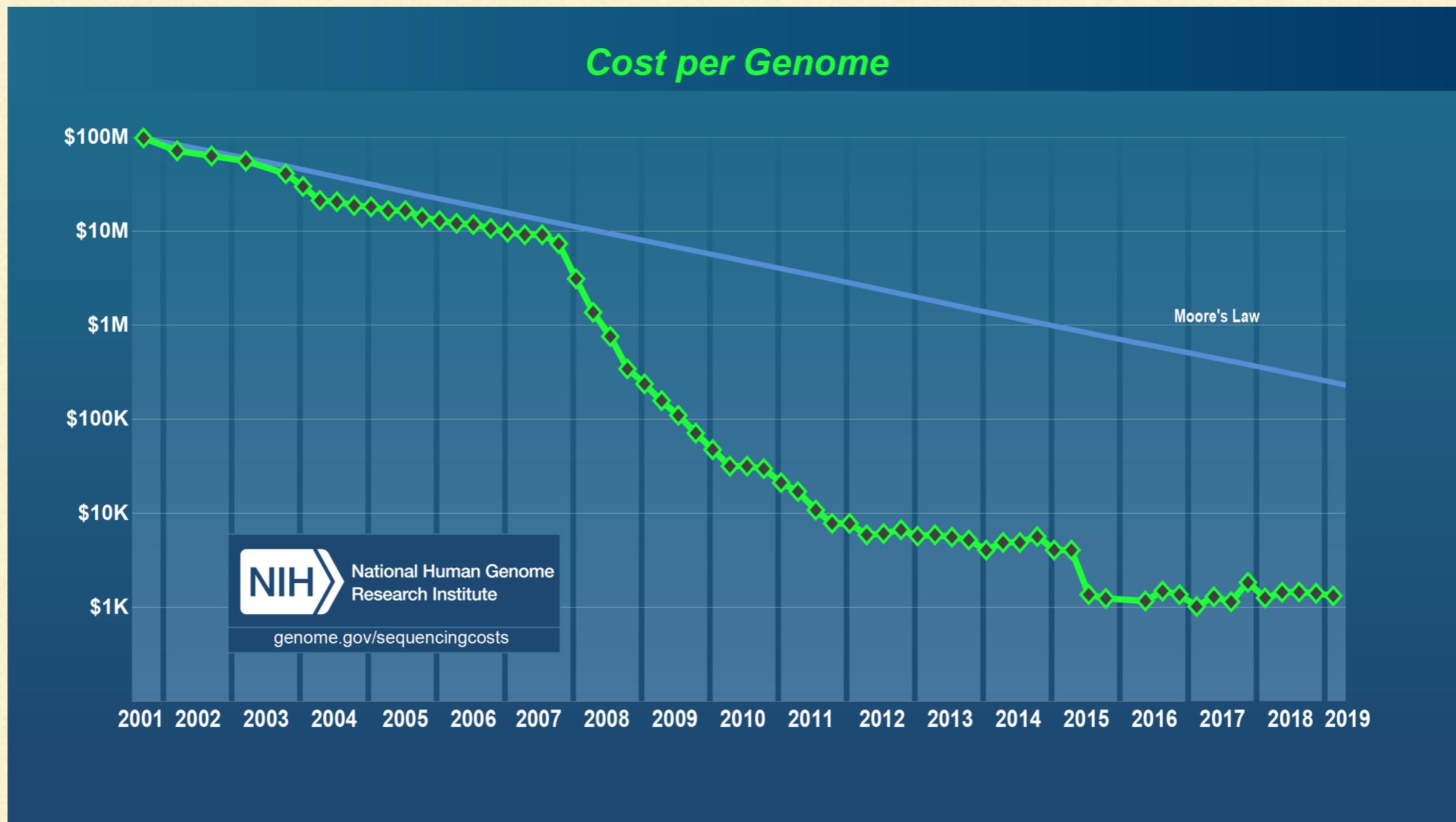
- Next Generation sequencing (first papers out in 2005)
- Amplify DNA on fragments (Illumina: attached to a flow cell)
 - Higher error rate compensated by coverage
 - Shorter read fragments make difficult to cover genomic regions such as high repetitive DNA
 - A complete genome sequence at 30X is only 85% complete

THE EXPLOSION



The sequence explosion - Nature, 2010

CUANTO CUESTA SECUENCIAR UN GENOMA?



“Our ability to sequence human genomes has vastly outpaced our ability to interpret genetic variation”

REVIEW

doi:10.1038/nature24286

DNA sequencing at 40: past, present and future

Jay Shendure^{1,2}, Shankar Balasubramanian^{3,4}, George M. Church⁵, Walter Gilbert⁶, Jane Rogers⁷, Jeffery A. Schloss⁸ & Robert H. Waterston¹

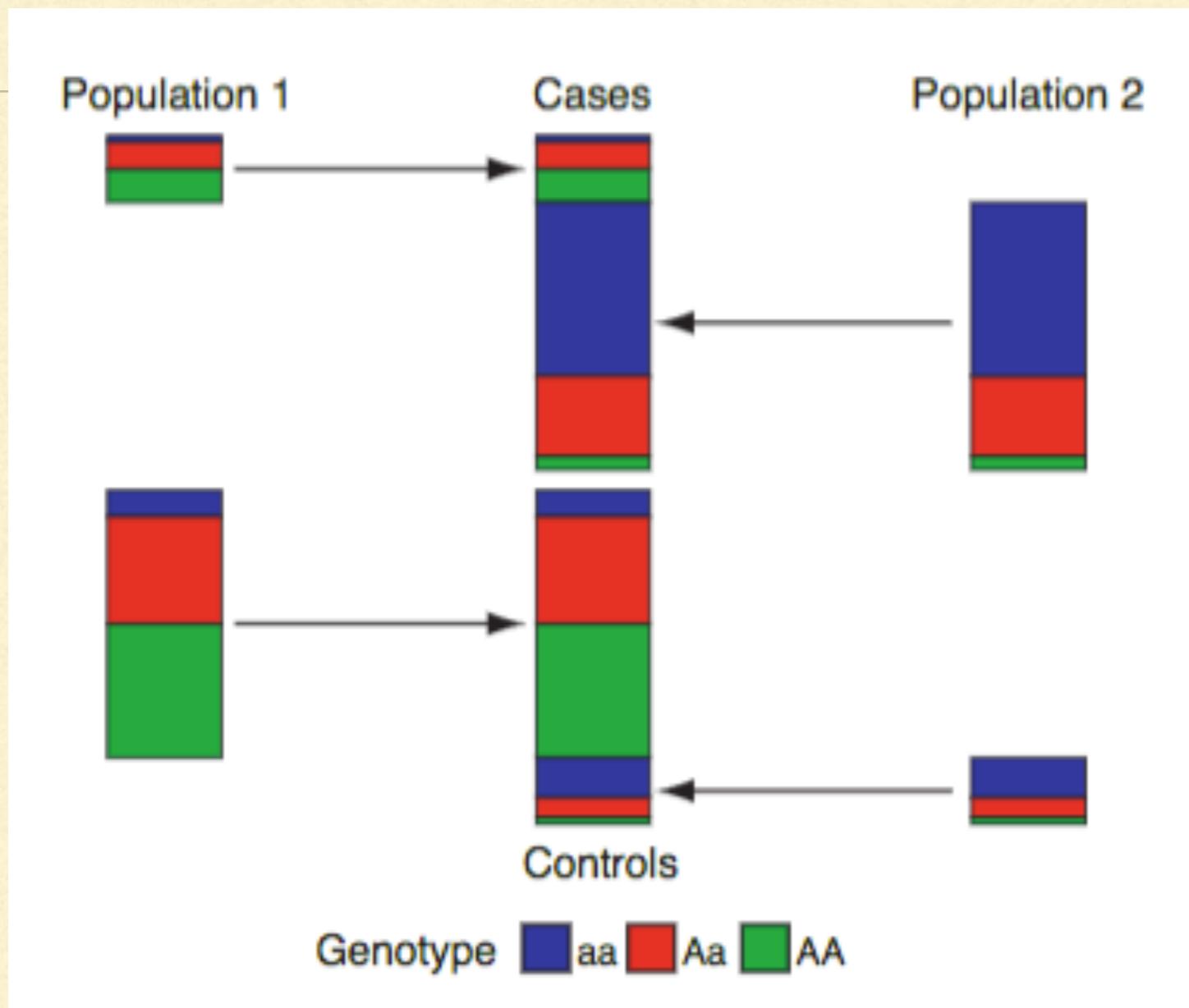
<https://www.nationalgeographic.com/news/2010/3/human-genome-project-tenth-anniversary/>

POPULATION STRUCTURE / STRATIFICATION

- Important for medical association studies
- Target gene(s) associated to an illness, or other phenotypic trait
- Hidden structure (relatedness between separate groups) brings false positive results and failures to detect genuine associations
- Effects increase with sample size

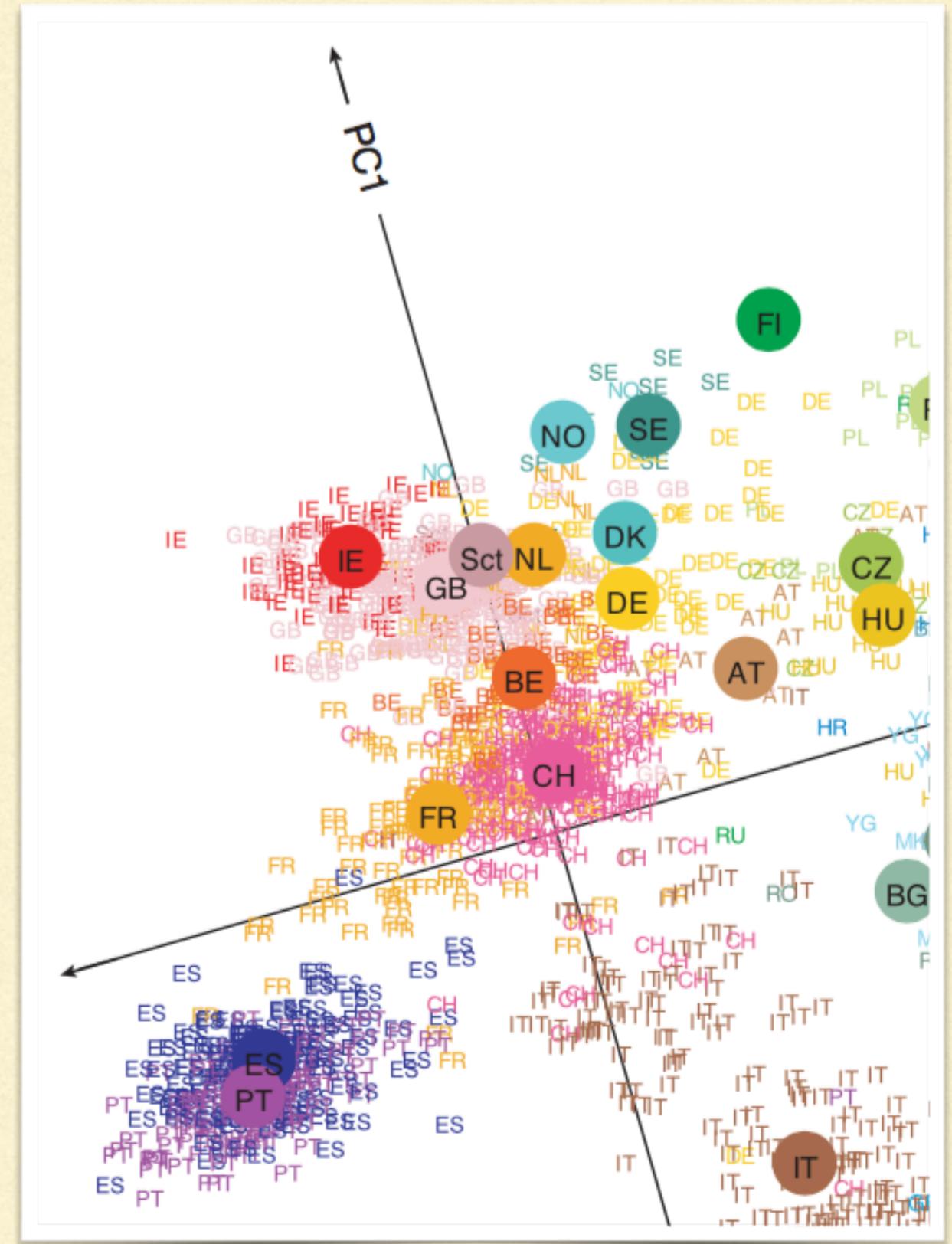


Example of effects of pop structure at a SNP locus. Two populations in which the cases have an excess of individuals from population 2 and population 2 has a lower frequency of allele A than population 1. In this example, the structure mimics the signal of association in that there is a significant difference in allele and genotype frequencies between cases and controls.



Example of effects of pop structure at a SNP locus. Two populations in which the cases have an excess of individuals from population 2 and population 2 has a lower frequency of allele A than population 1. In this example, the structure mimics the signal of association in that there is a significant difference in allele and genotype frequencies between cases and controls.

POPULATION STRUCTURE: PCA AND ADMIXTURE



PCA - PRINCIPAL COMPONENT ANALYSIS

Principal component analysis (PCA) is a statistical procedure that uses an orthogonal transformation to convert a set of observations of possibly correlated variables (entities each of which takes on various numerical values) into a set of values of linearly uncorrelated variables called **principal components**. (Wikipedia)

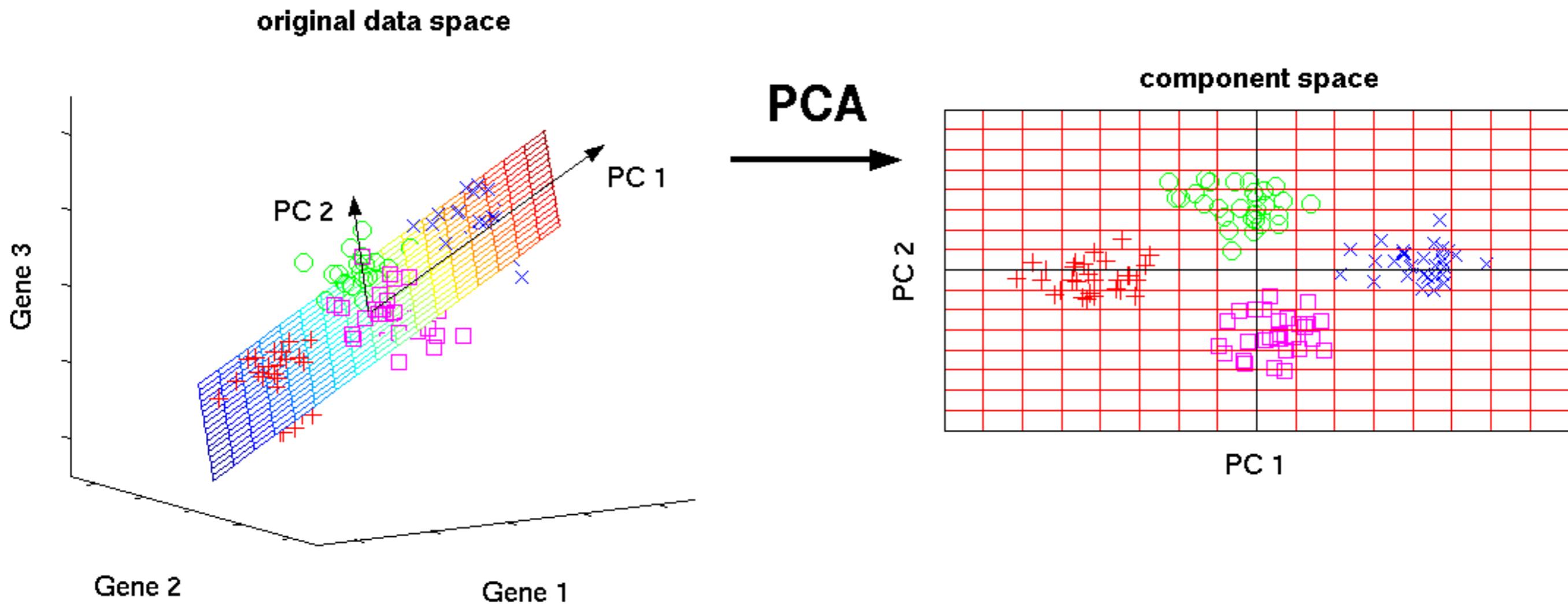
Each PC projection is designed to capture as much variance as possible for a linear transformation. The first dimension is the component which explains the highest variation.

Eigenvalues decomposition of a covariance matrix from uncorrelated source of variation (genetic markers).

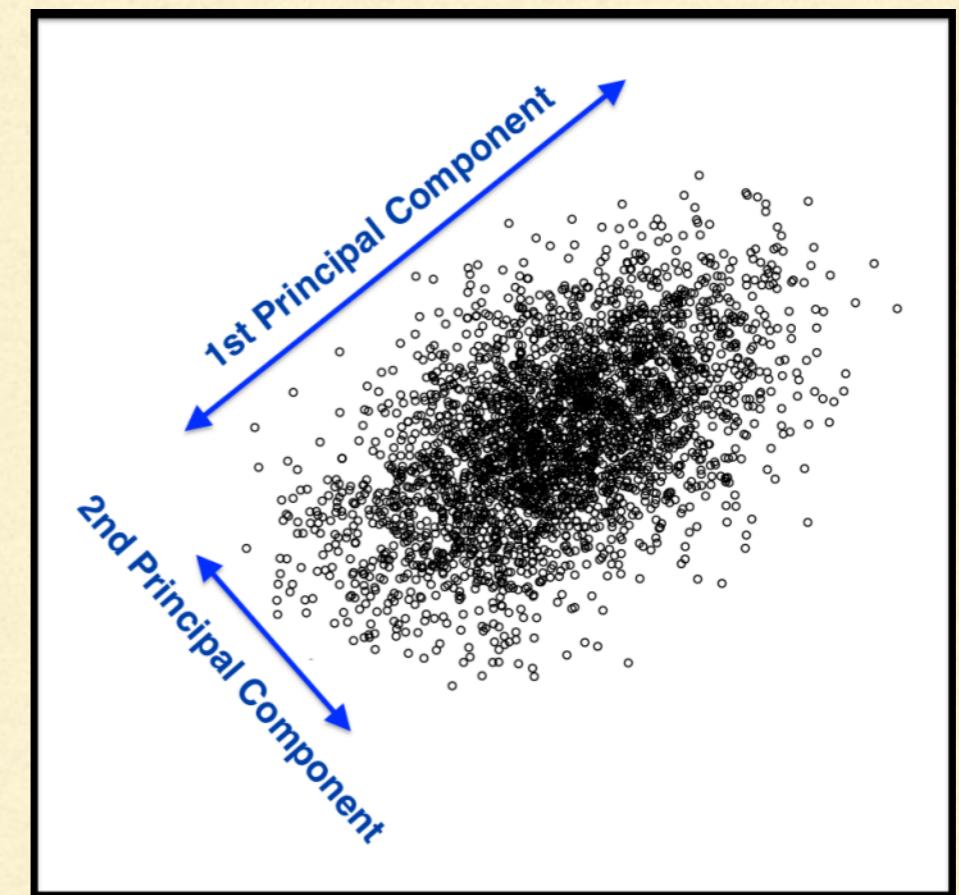
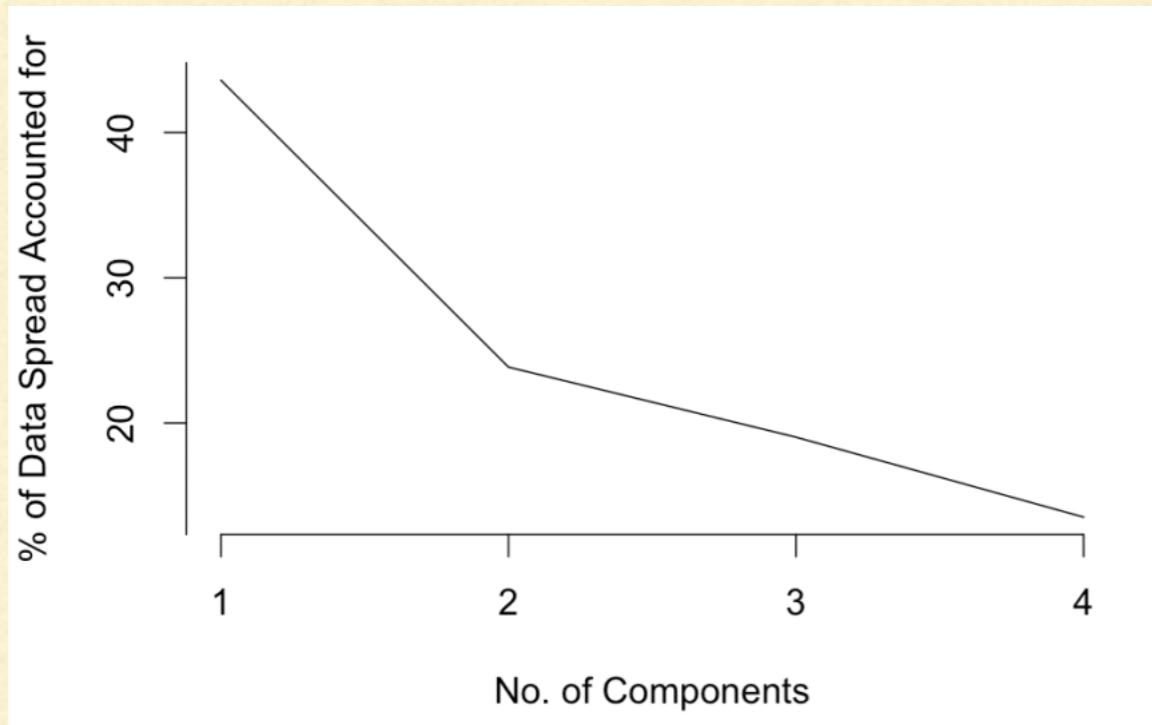
The amount of eigenvectors/values that exist equals the number of dimensions the data set has.

PCA

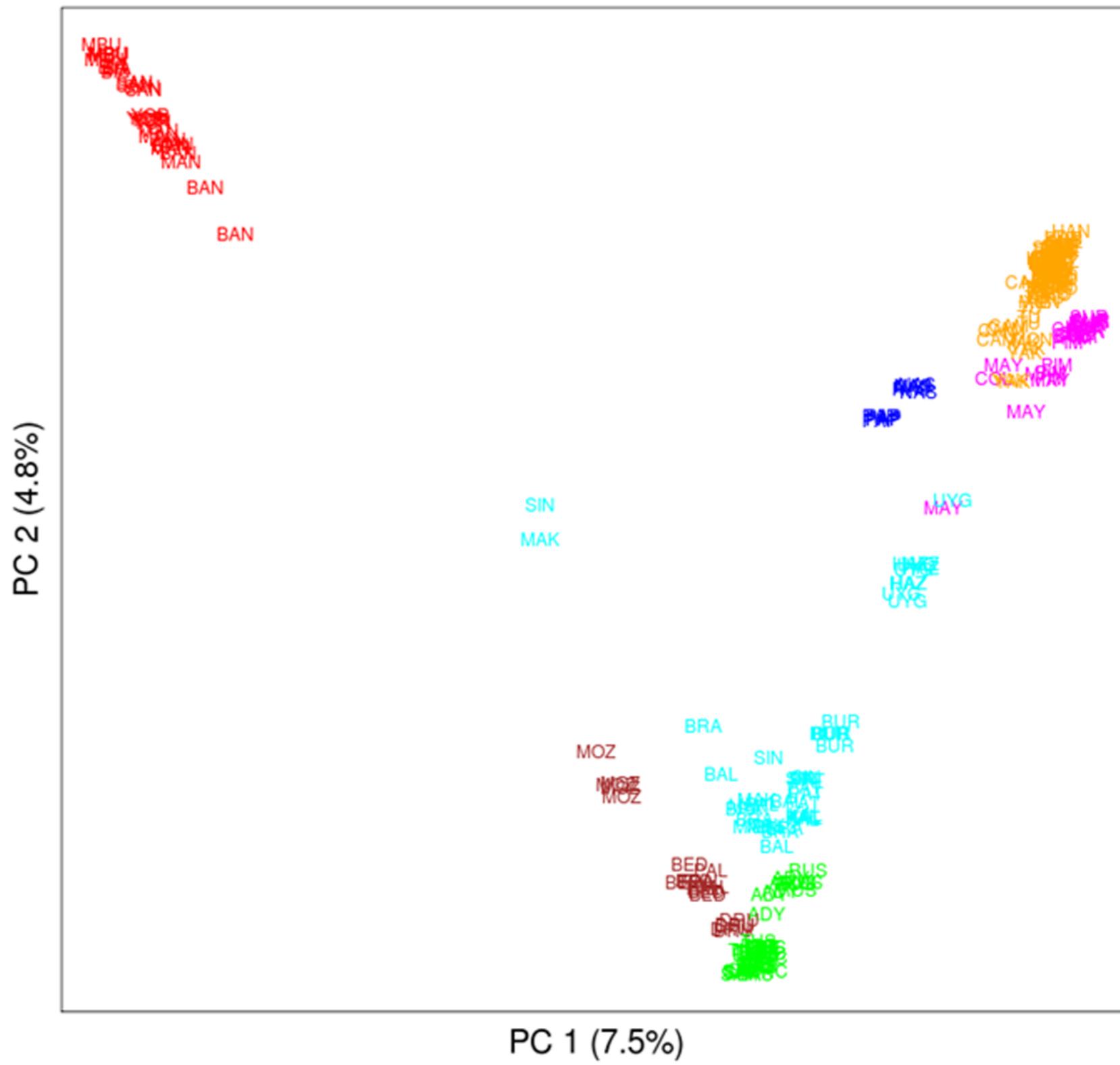
PCA is a standard technique for visualizing high dimensional data and for data pre-processing. PCA reduces the dimensionality (the number of variables) of a data set by maintaining as much variance as possible.



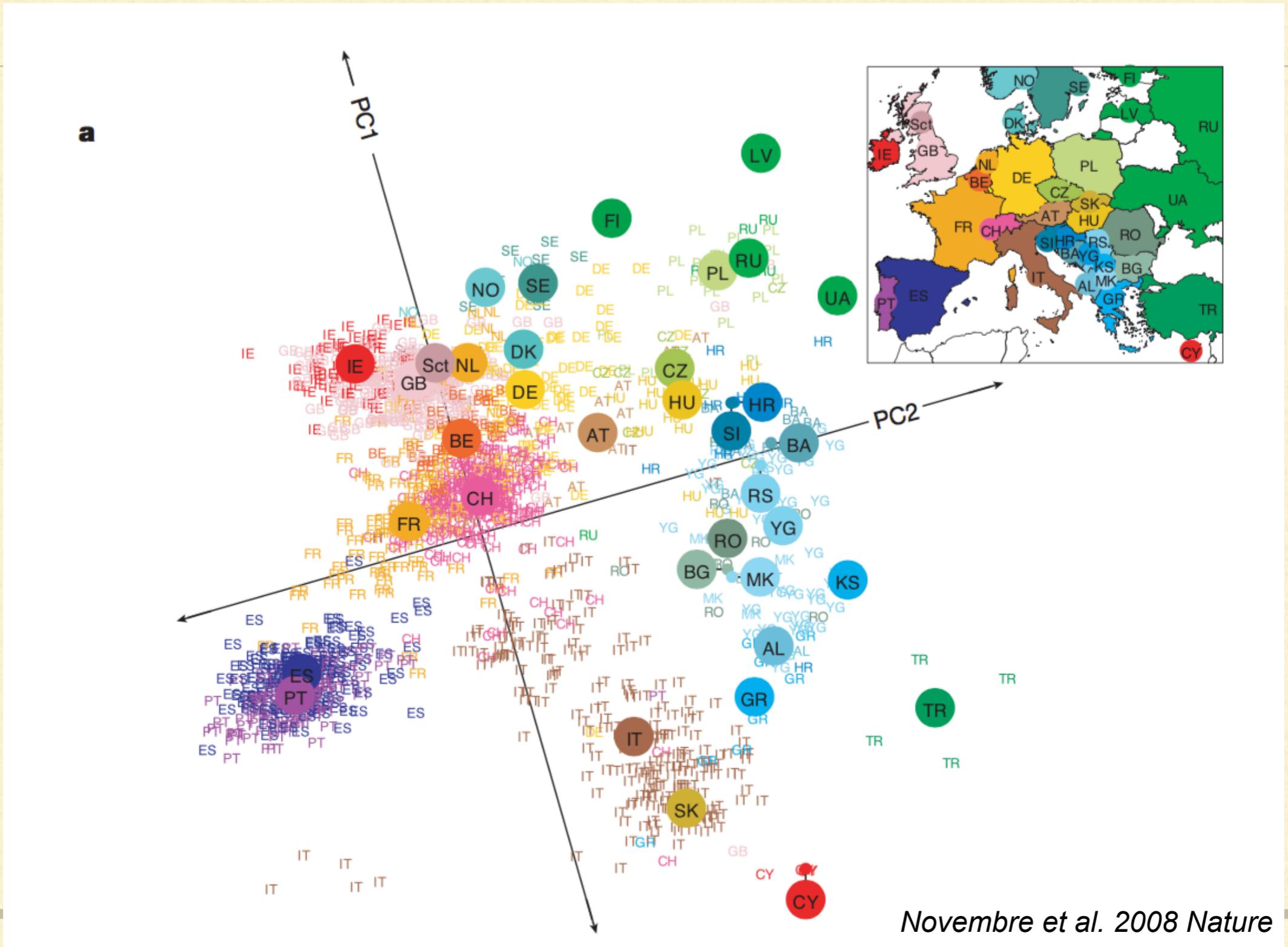
Principal components are also ordered by their effectiveness in differentiating data points, with the first principal component doing so to the largest degree. To keep results simple and generalizable, only the first few principal components are selected for visualization and further analysis. The number of principal components to consider is determined by something called a scree plot:



COL	AM - Colombian (Arawak)
KAR	AM - Karitiana
MAY	AM - Maya
PIM	AM - Pima
SUR	AM - Surui
BAL	CSA - Balochi
BRA	CSA - Brahui
BUR	CSA - Burusho
HAZ	CSA - Hazara
KAL	CSA - Kalash
MAK	CSA - Makrani
PAT	CSA - Pathan
SIN	CSA - Sindhi
UYG	CSA - Uygur
CAM	EA - Cambodia
DAI	EA - Dai
DAU	EA - Daur
HAN	EA - Han
HEZ	EA - Hezhen
JAP	EA - Japanese
LAH	EA - Lahu
MIA	EA - Miaozi
MON	EA - Mongolia
NAX	EA - Naxi
ORO	EA - Oroqen
SHE	EA - She
TU	EA - Tu
TUJ	EA - Tujia
XIB	EA - Xibo
YAK	EA - Yakut
YIZ	EA - Yizu
ADY	EUR - Adygei
BAS	EUR - Basque
BER	EUR - Bergamo
FRE	EUR - French
ORC	EUR - Orcadian
RUS	EUR - Russian
SAR	EUR - Sardinian
TUS	EUR - Tuscan
BED	ME - Bedouin
DRU	ME - Druze
MOZ	ME - Mozabite
PAL	ME - Palestinian
NAS	OC - Nasioi
PAP	OC - Papuan
BAN	SSA - Bantu
BIA	SSA - BiakaPygmy
MAN	SSA - Mandenka
MBU	SSA - MbutiPygmy
SAN	SSA - San
YOR	SSA - Yoruba

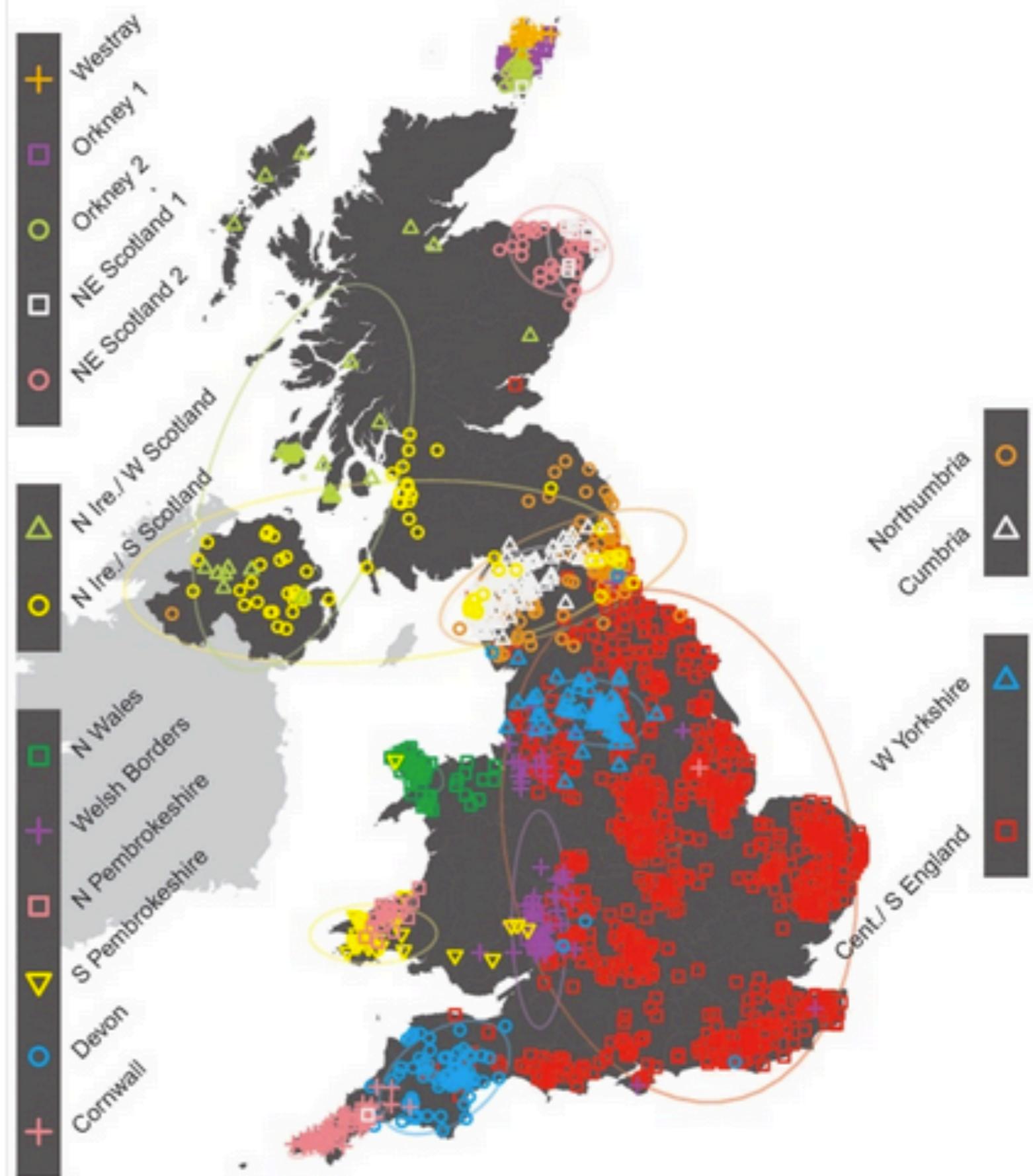


RELATIONSHIP BETWEEN INDIVIDUALS



RELATIONSHIP BET

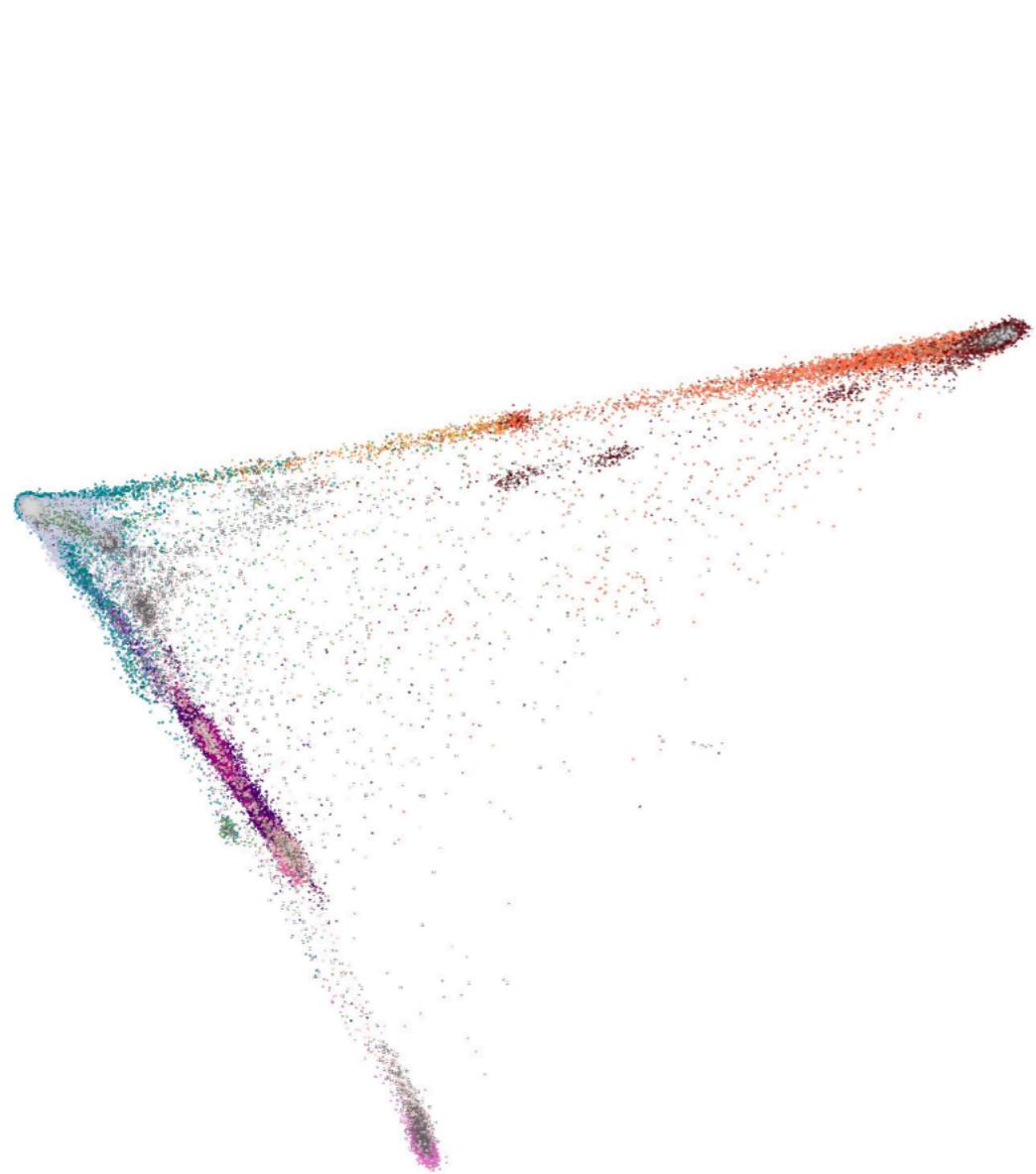
- ...At finer scale



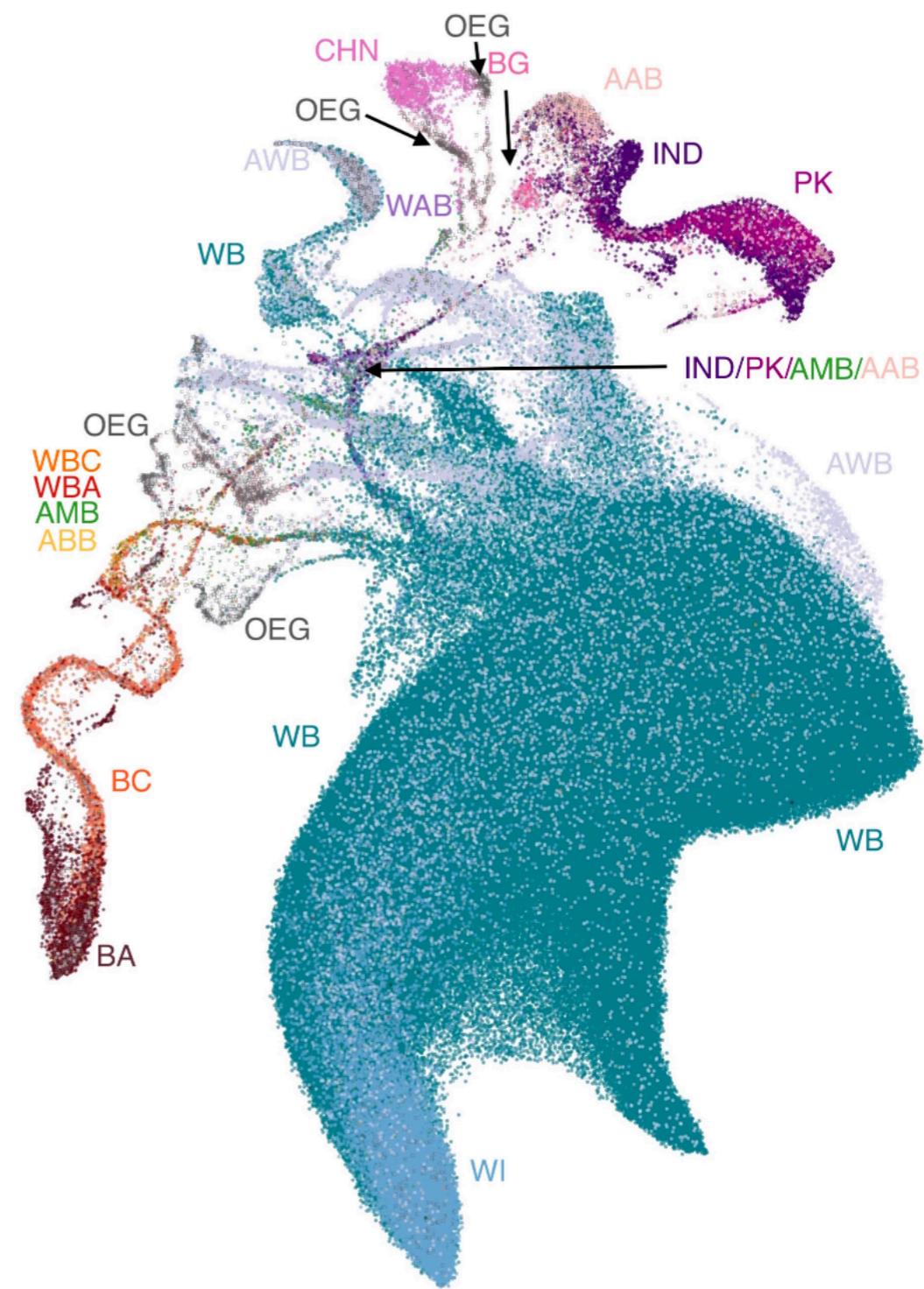
Stephen Leslie

A map of the United Kingdom shows how individuals cluster based on their genetics, with a striking relationship to the geography of the country.

Leslie et al. 2015 *Nature*



(a) Principal components 1 and 2



(b) UMAP on first 10 principal components

UK BioBank dataset plotted with different PCA methods

Diaz-Papkovich et al. 2018 BioRxiv

OTHER NON-LINEAR DIMENSIONAL REDUCTION METHODS

- MDS - multidimensional scaling
- Starts from a distance matrix to force individual relationships onto a fixed sets of dimensions (usually bi or tridimensional)
- The geometrical transformation aims at minimising the Stress (an estimate of how the data is fitting the space)

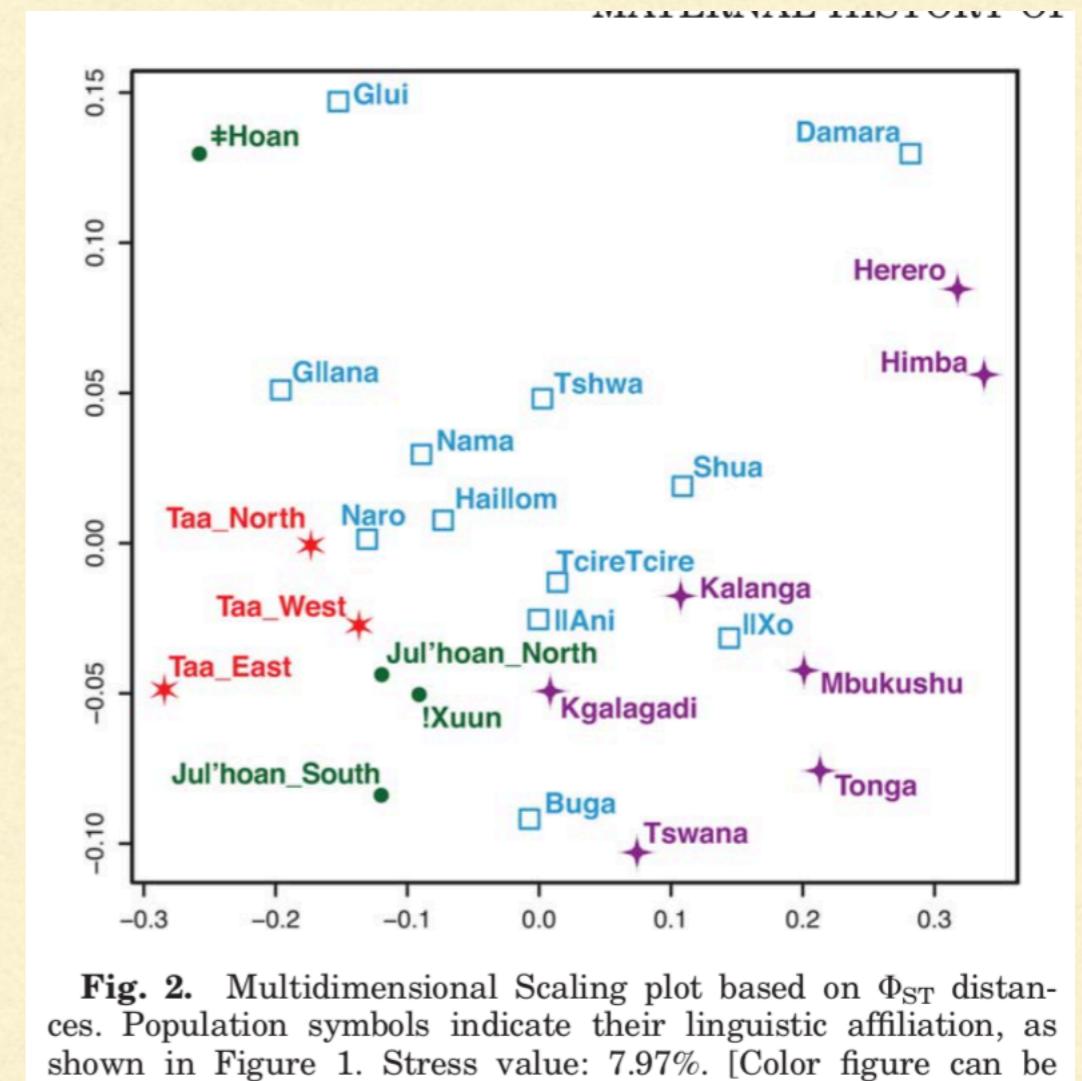
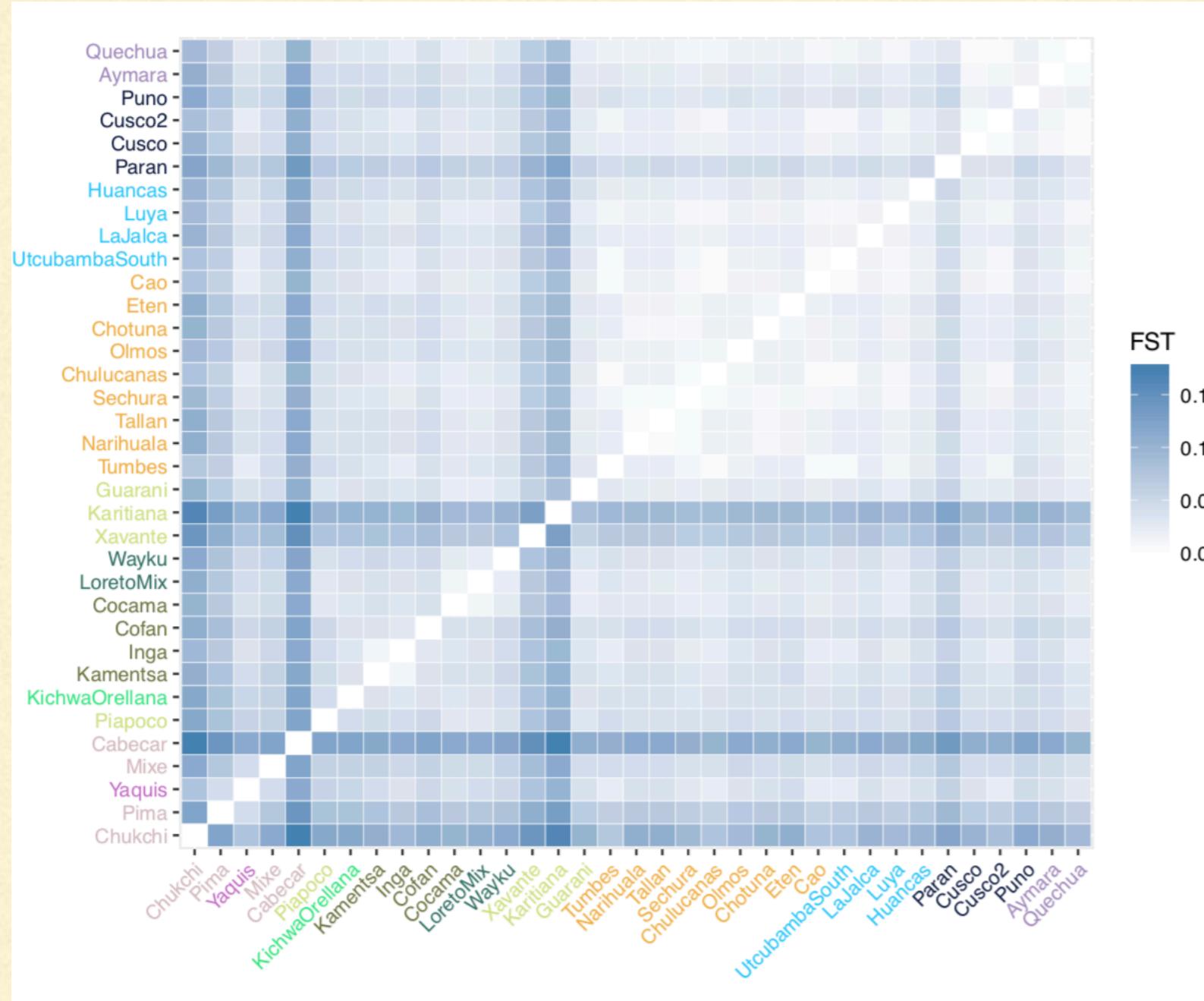


Fig. 2. Multidimensional Scaling plot based on Φ_{ST} distances. Population symbols indicate their linguistic affiliation, as shown in Figure 1. Stress value: 7.97%. [Color figure can be seen in the online version of this article.]

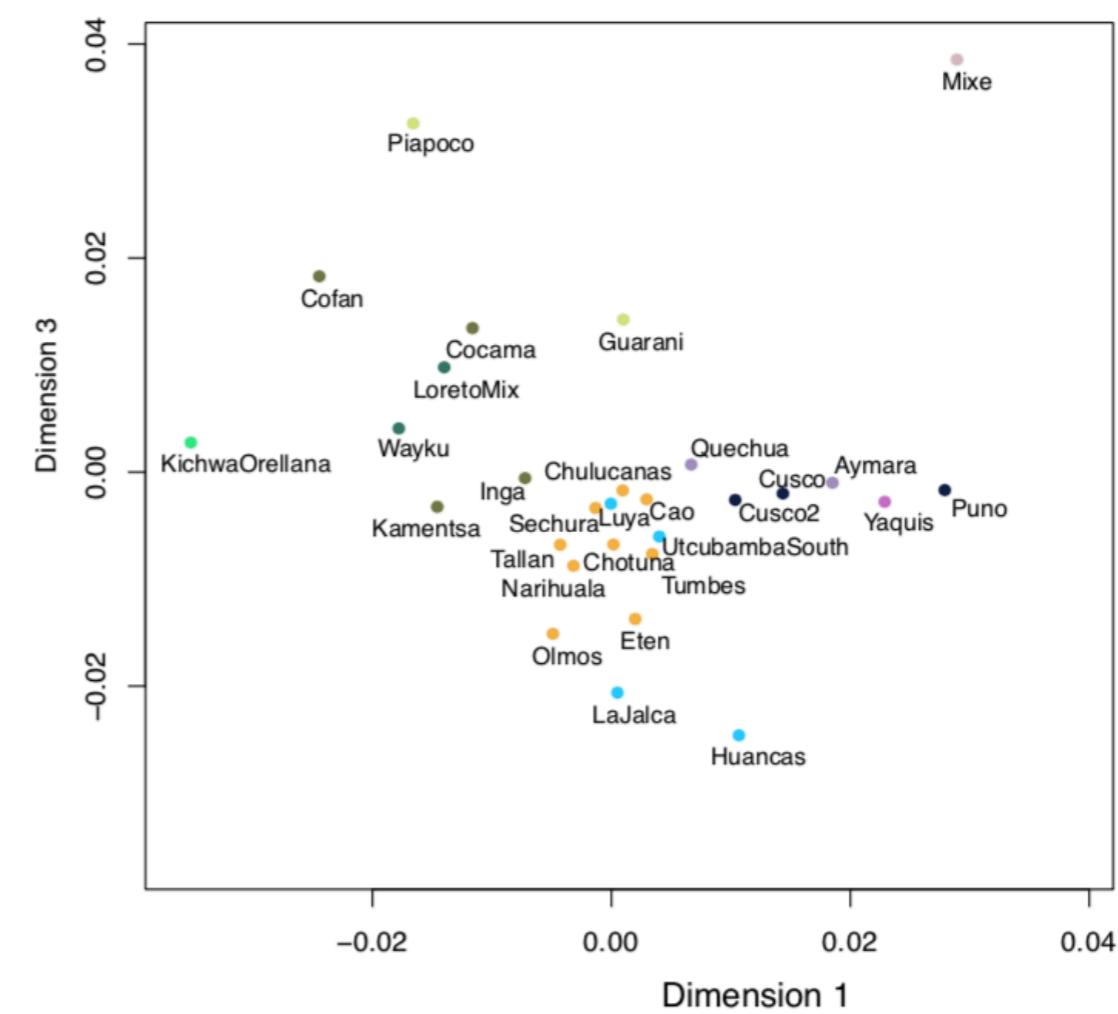
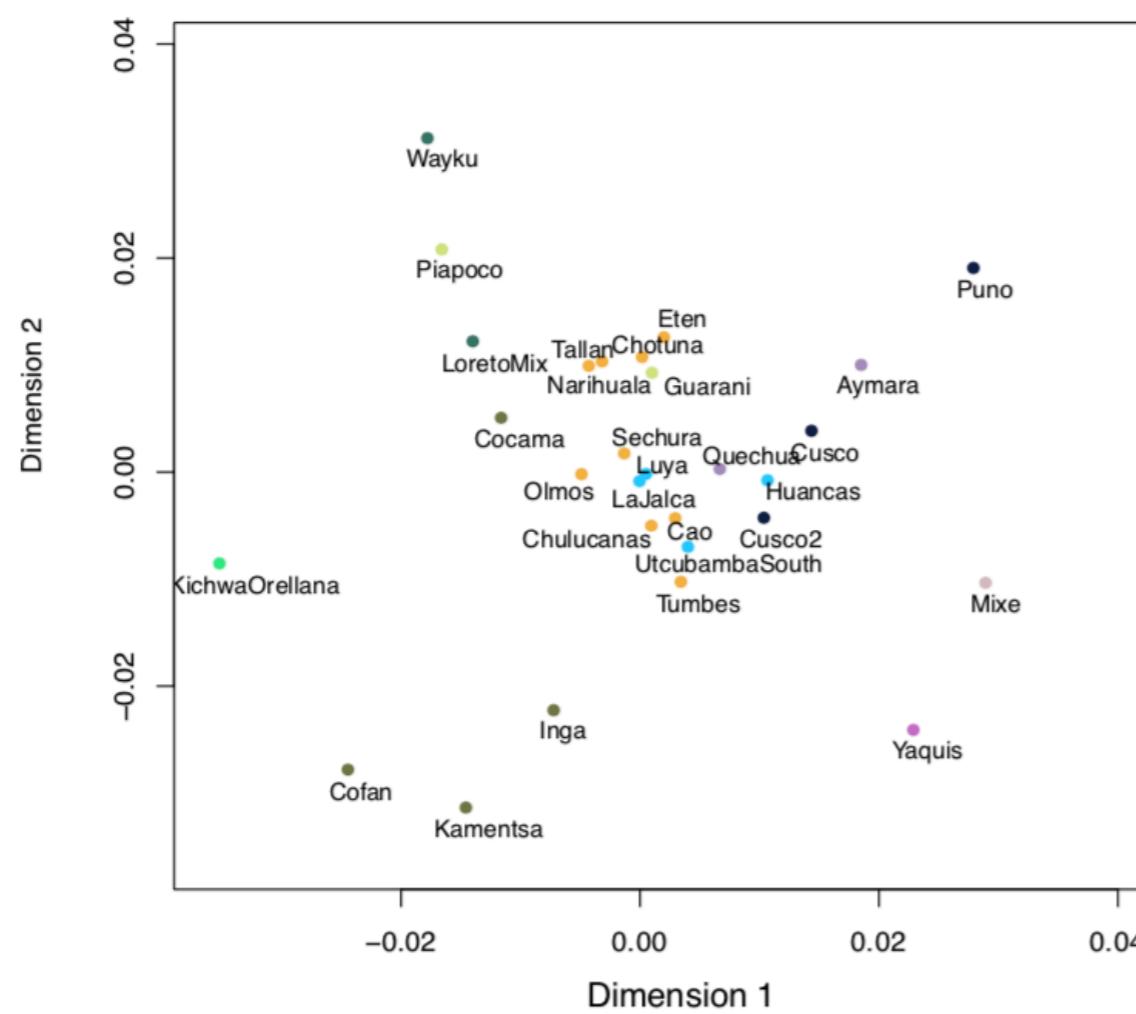
FST DISTANCE MATRIX RELATIONSHIP BETWEEN POPULATIONS



RELATIONSHIP BETWEEN POPULATIONS

■ MultiDimensional Scaling (MDS)

A



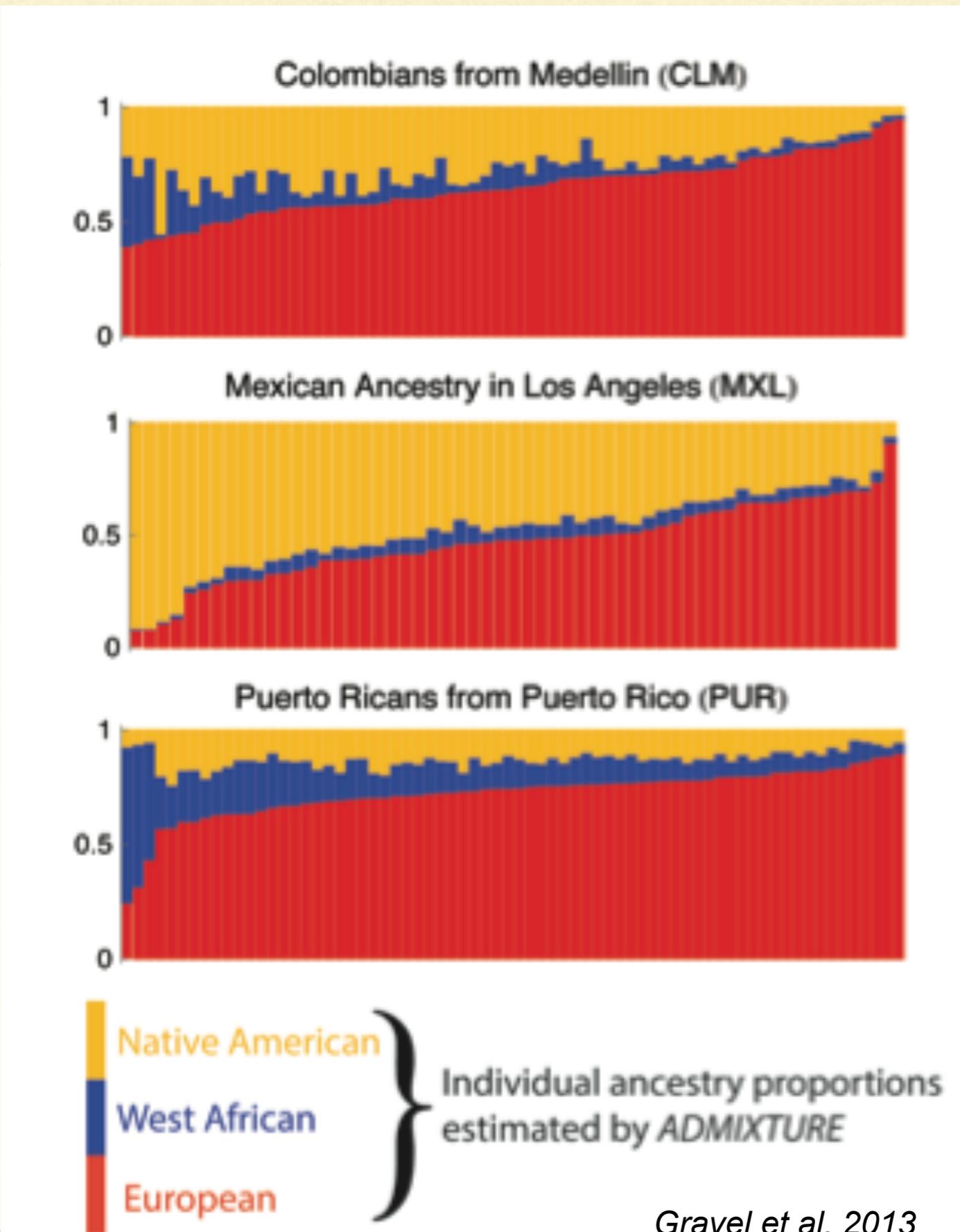
ADMIXTURE

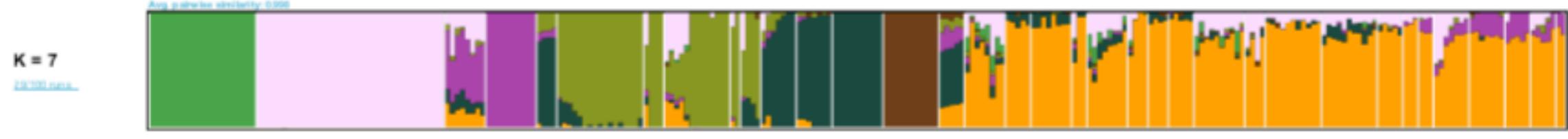
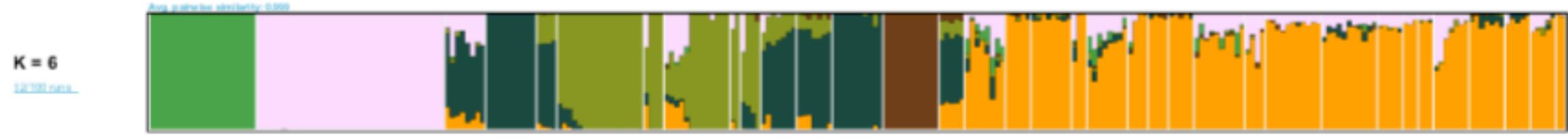
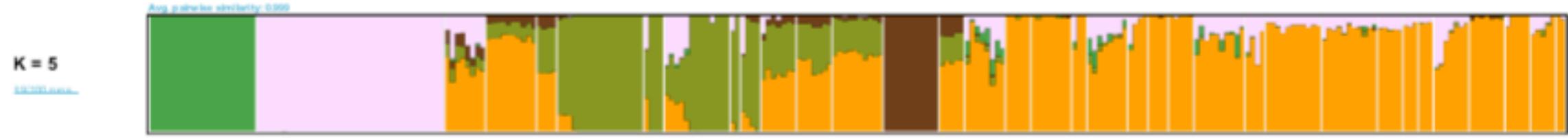
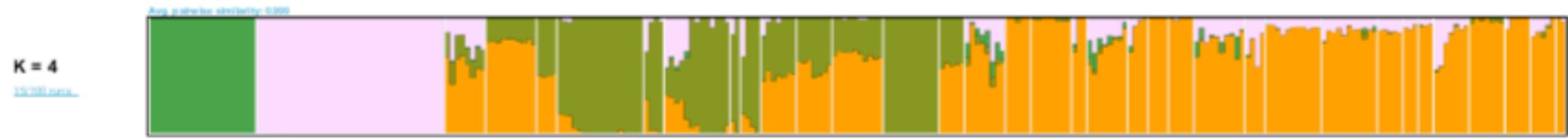
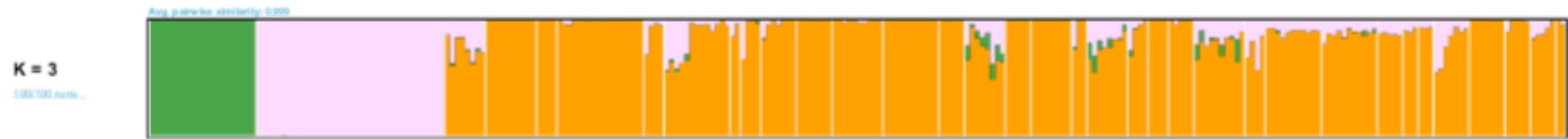
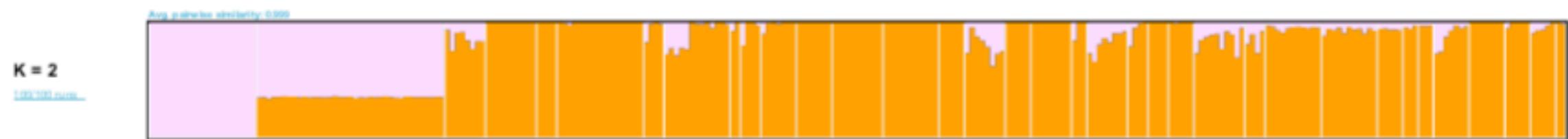
Admixture is a very useful and popular tool to analyse SNP data. It performs an unsupervised clustering of large numbers of samples, and allows each individual to be a mixture of clusters.

STRUCTURE barplot has become a de-facto standard used as a non-parametric description of genetic data alongside a Principle Components Analysis.

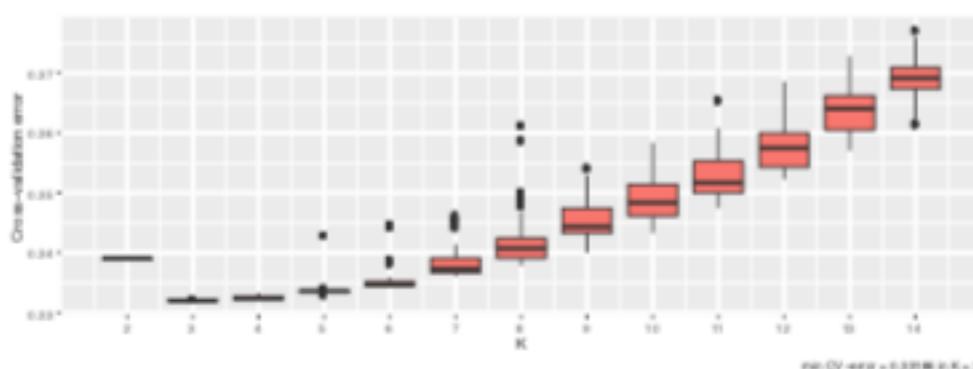
ADMIXTURE

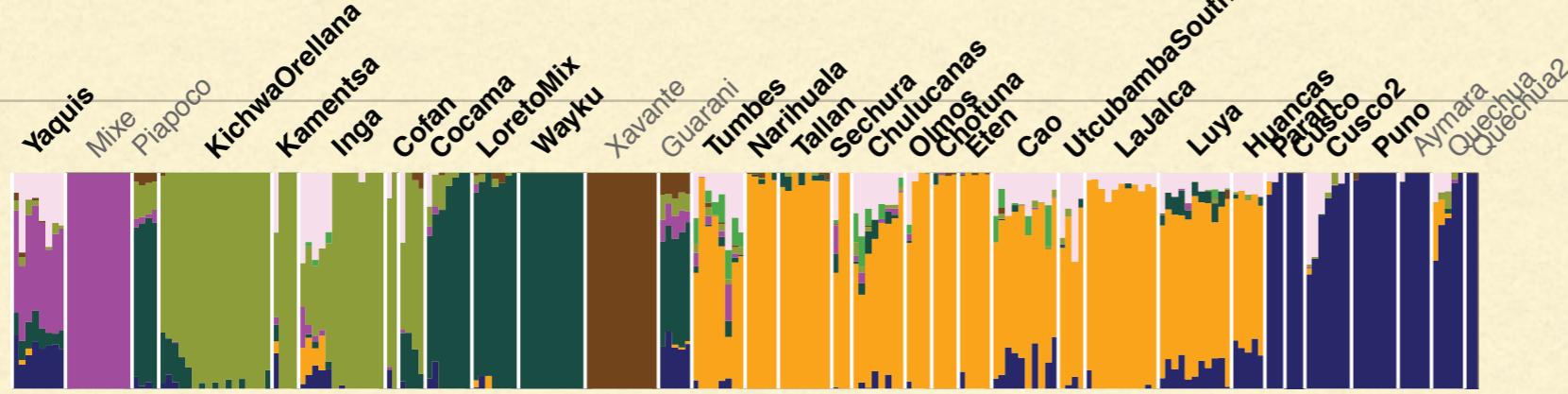
- Explain genetic diversity in terms of ancestry blocks



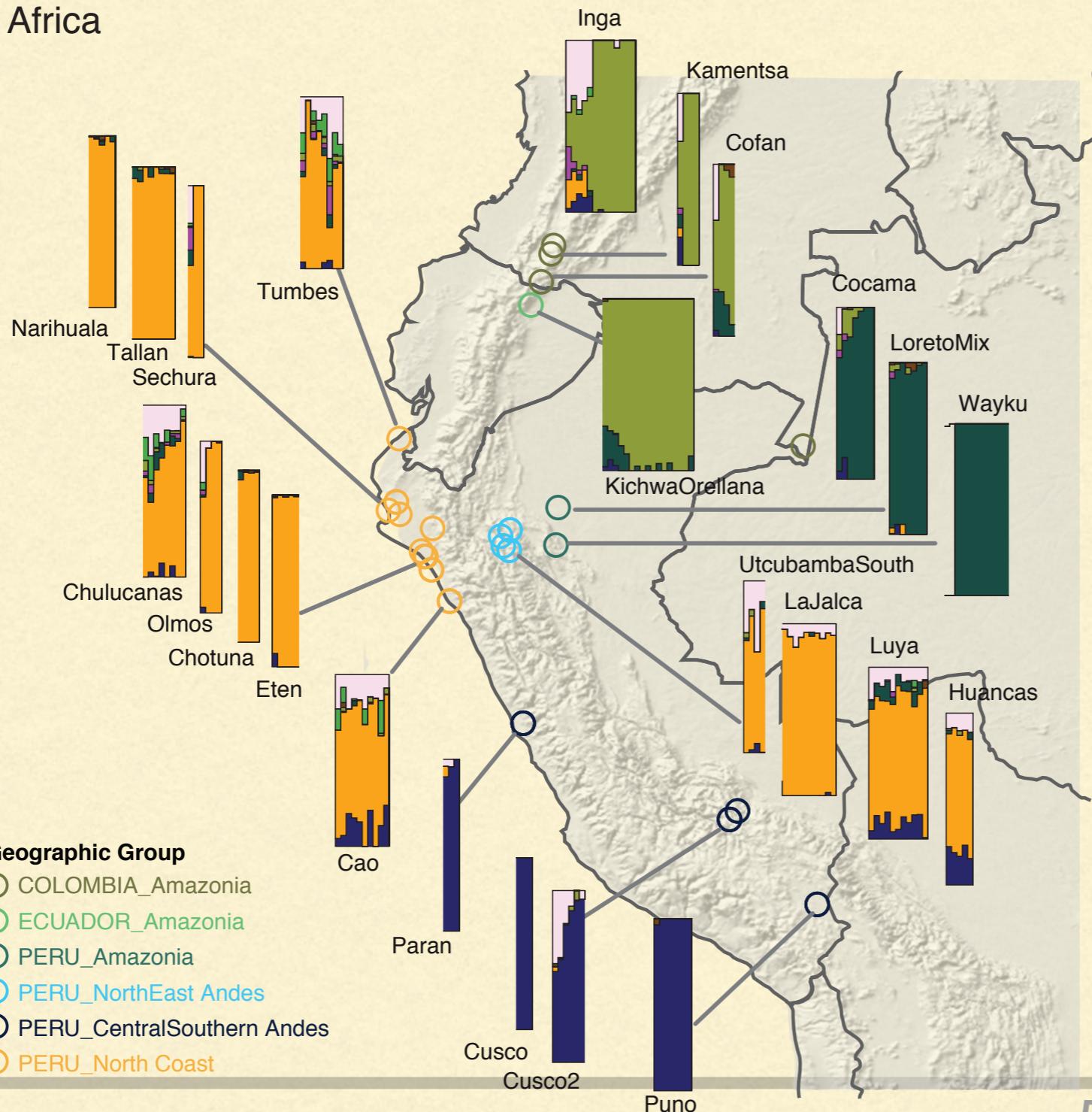


Cross-Validation Error associated to each K





Europe
Africa



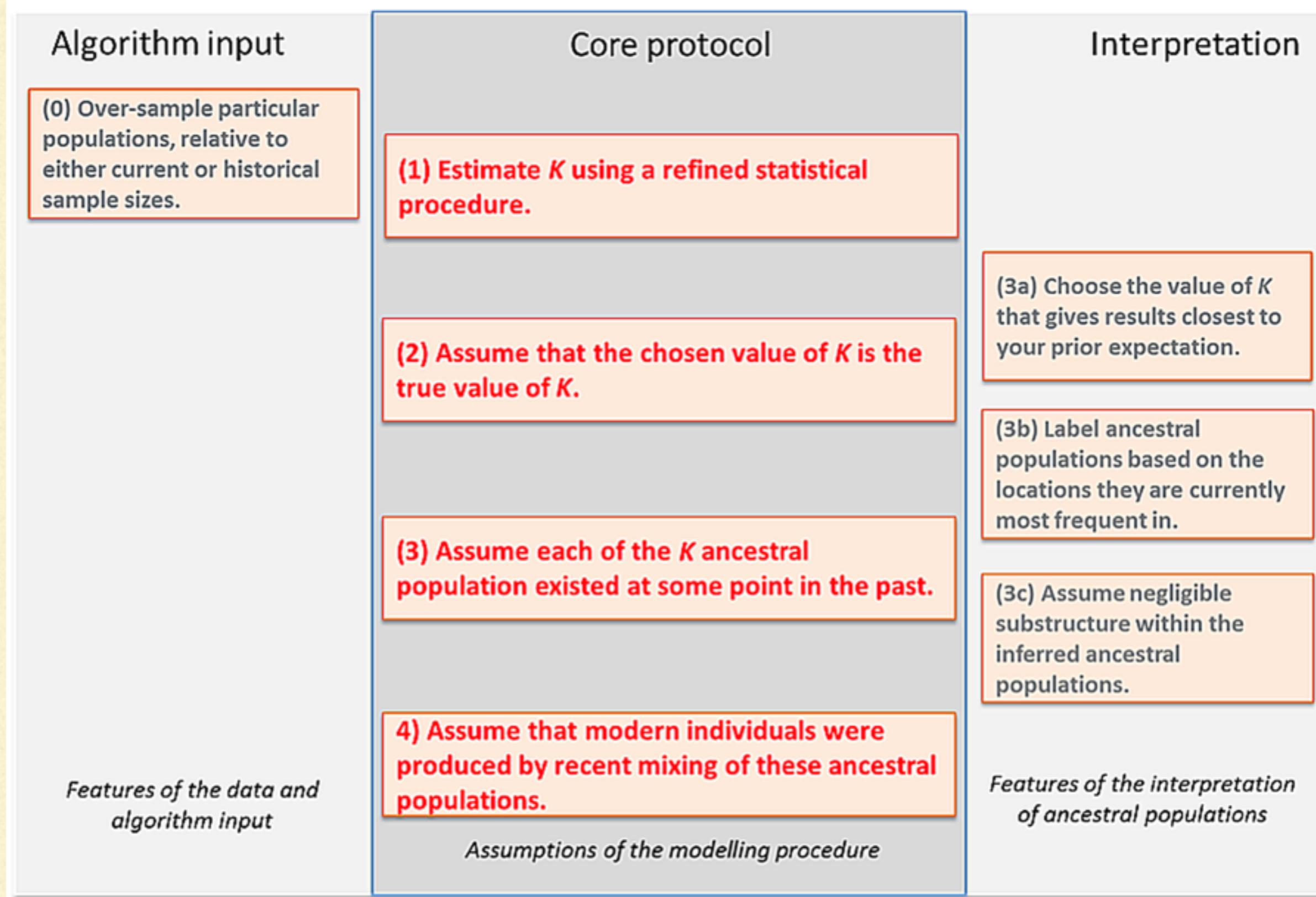
LIMITS OF ADMIXTURE

Histories can also be reconstructed using the same procedure for groups that do not have admixture in their recent history, where recent genetic drift is strong or that deviate in other ways from the underlying inference model. Unfortunately, such histories can be misleading.



The image shows a thumbnail of a research article from Nature Communications. The header features the journal logo 'nature COMMUNICATIONS' with a stylized orange and red wavy line graphic. Below the header, the text 'Article | OPEN | Published: 14 August 2018' is displayed. The main title of the article is 'A tutorial on how not to over-interpret STRUCTURE and ADMIXTURE bar plots'. Below the title, the authors are listed as 'Daniel J. Lawson, Lucy van Dorp & Daniel Falush' with an envelope icon for email. At the bottom, the citation information is provided: 'Nature Communications 9, Article number: 3258 (2018) | Download Citation ↓'

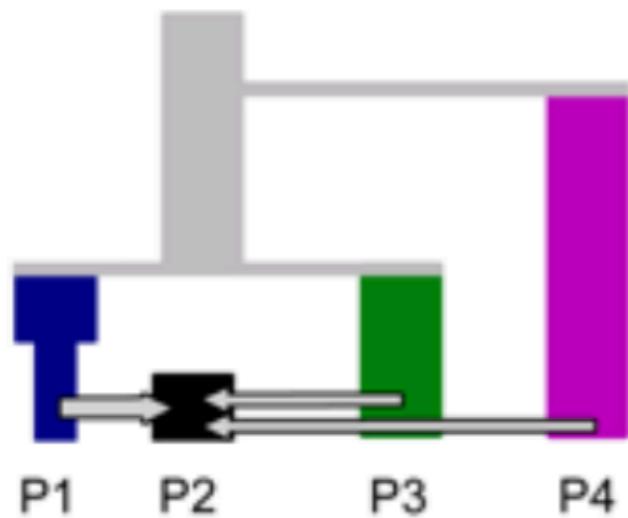
A protocol for (mis)interpreting admixture model results



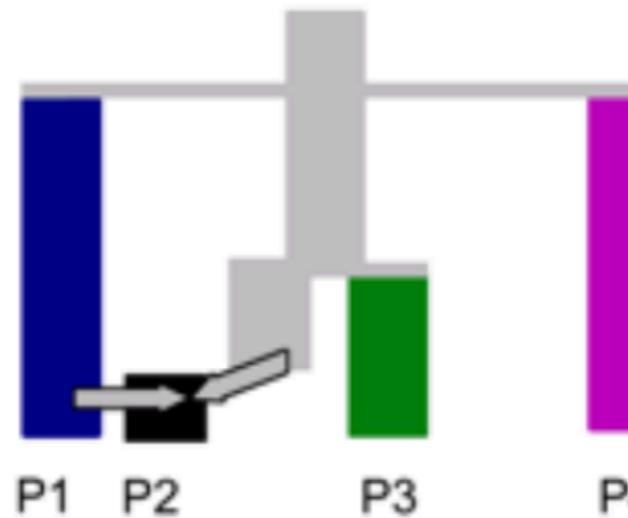
3 different population histories giving the same ADMIXTURE results

a

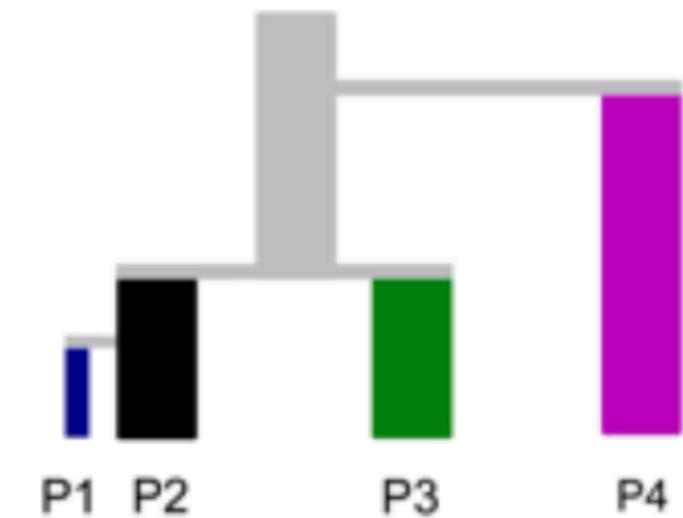
Recent Admixture



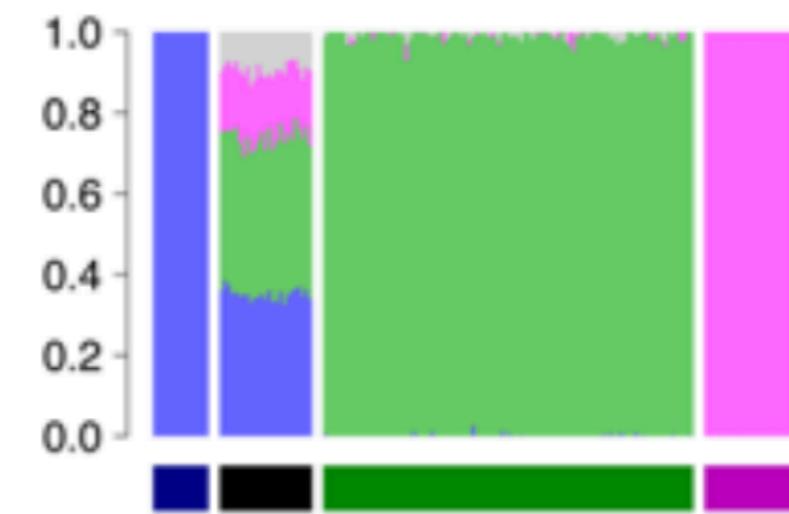
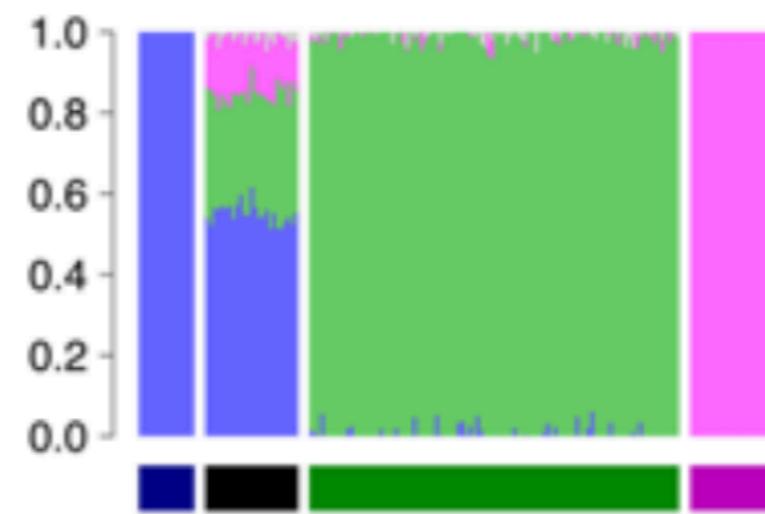
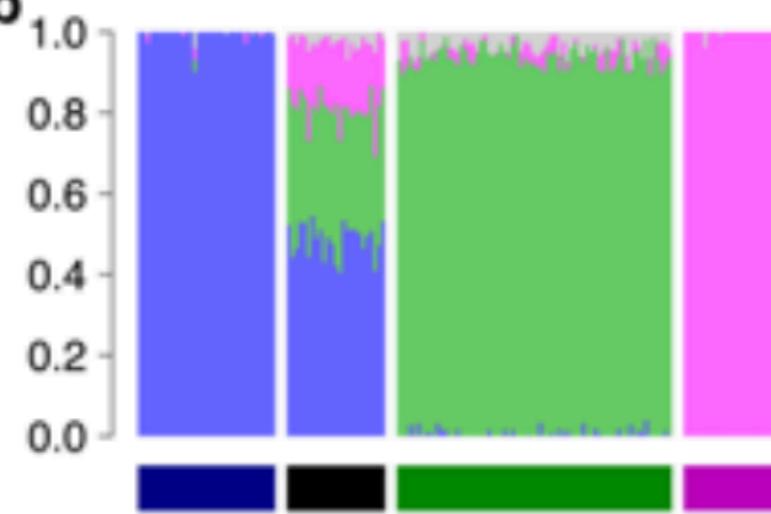
Ghost Admixture



Recent Bottleneck



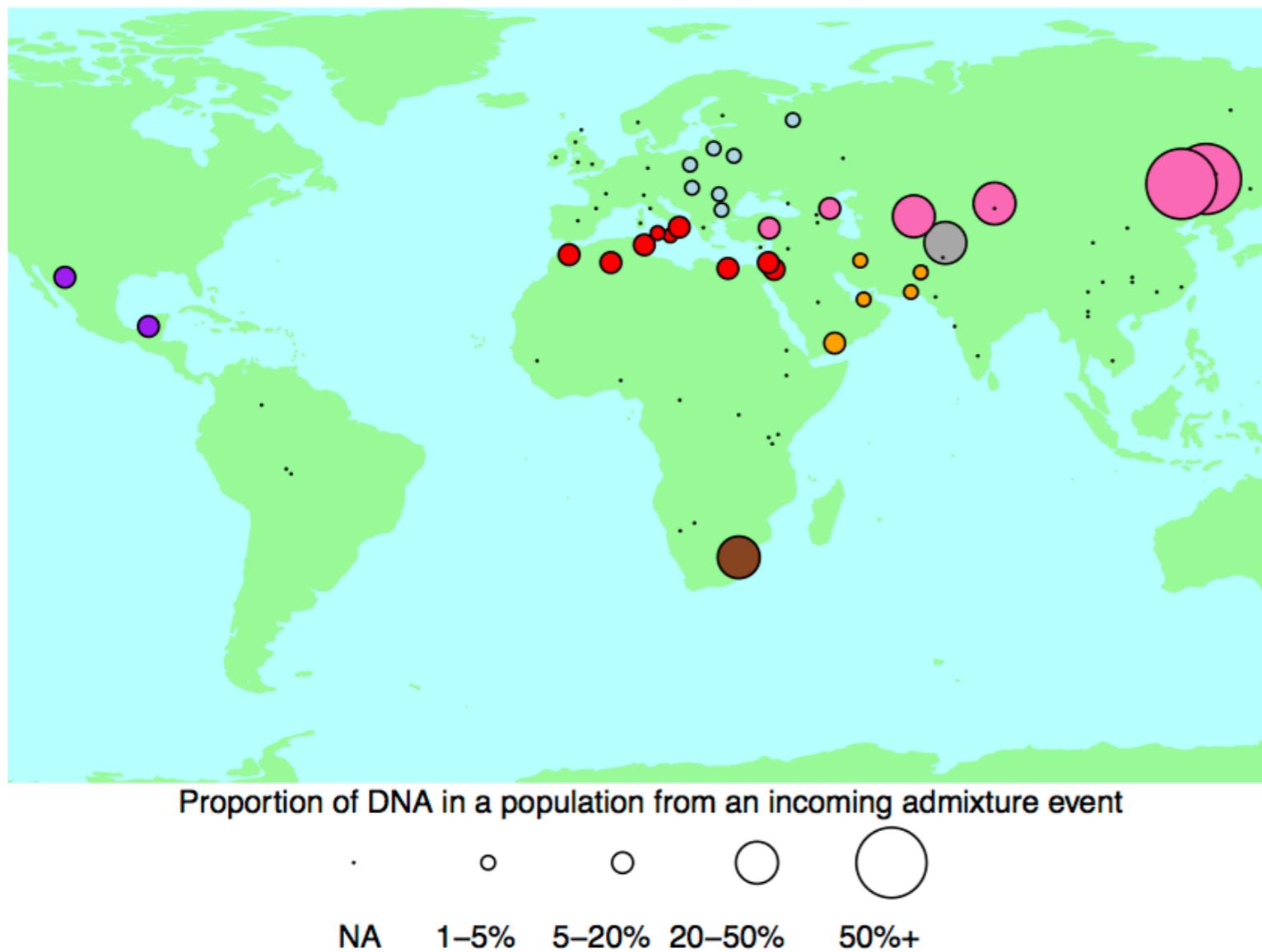
b



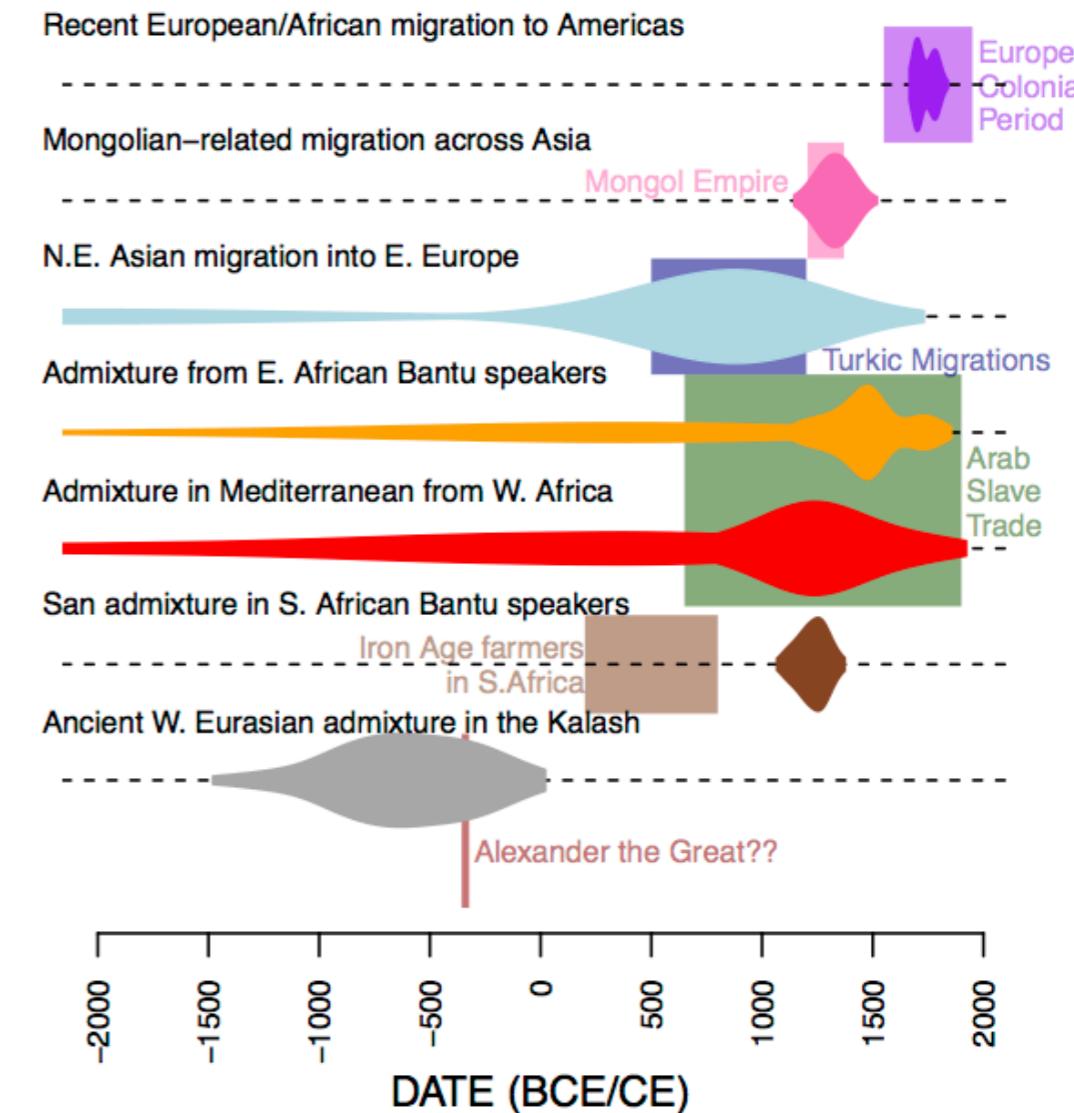
FINDING ADMIXTURE SOURCES AND DATING ADMIXTURE EVENTS

A Genetic Atlas of Human Admixture History

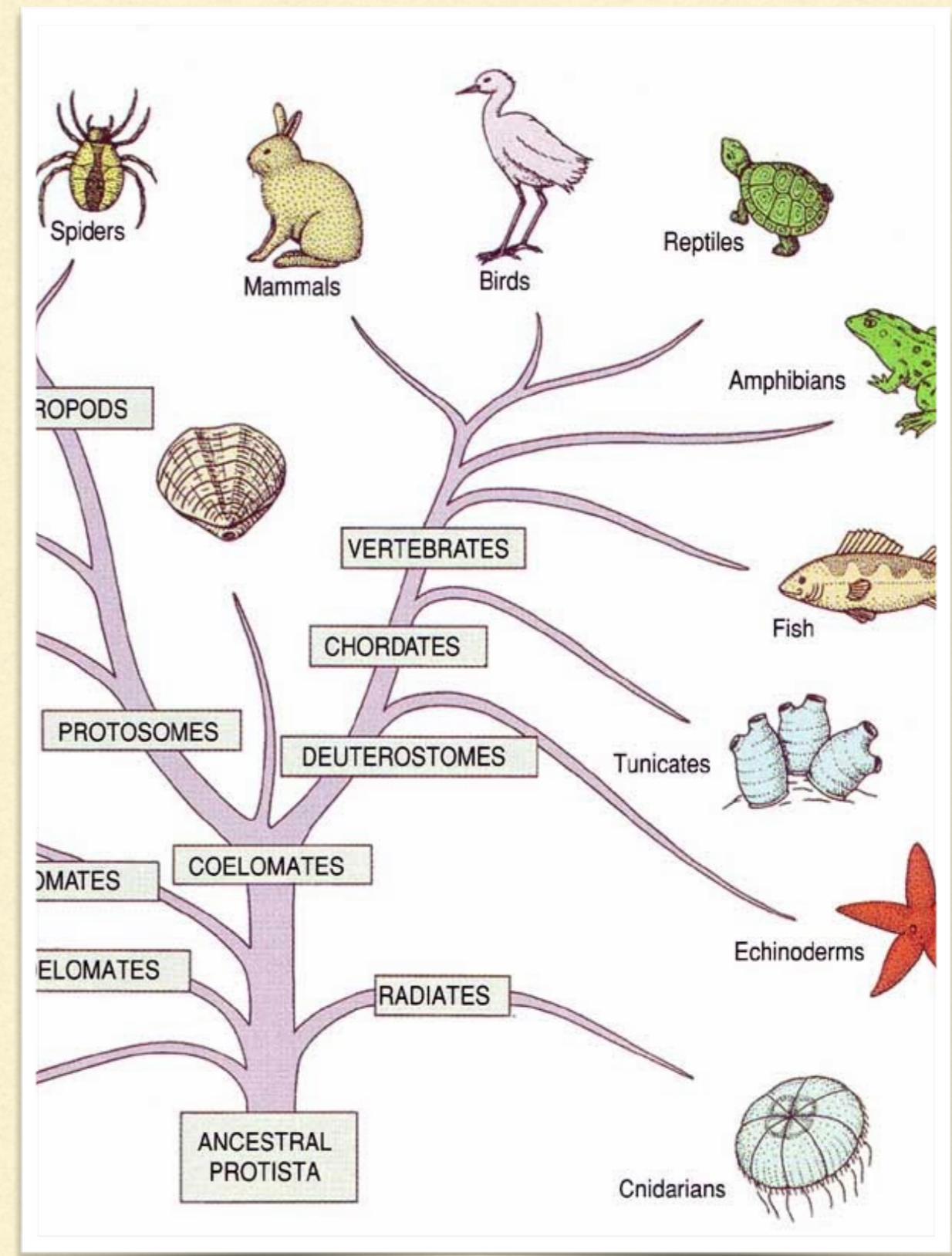
A summary of a selection of human migration events inferred only from genetic data using the new GLOBETROTTER method reported in: G. Hellenthal, G.B.J. Busby, G. Band, J.F. Wilson, C. Capelli, D. Falush & S. Myers Science (14th Feb 2014)



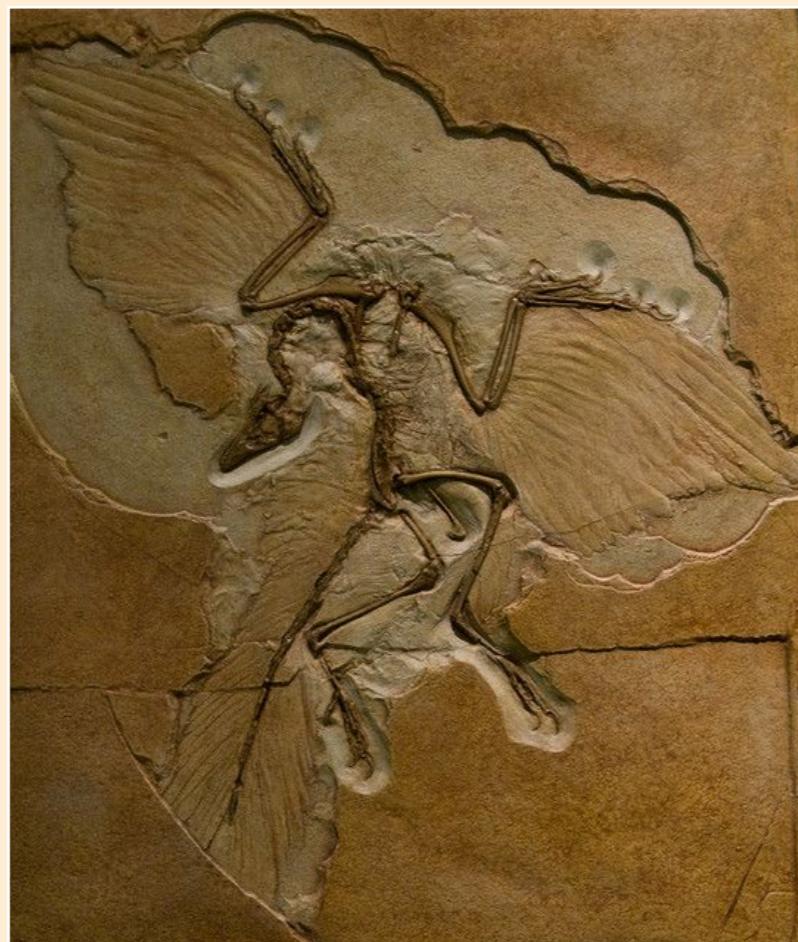
Timeline of events inferred by GLOBETROTTER



PHYLOGENETIC RECONSTRUCTIONS



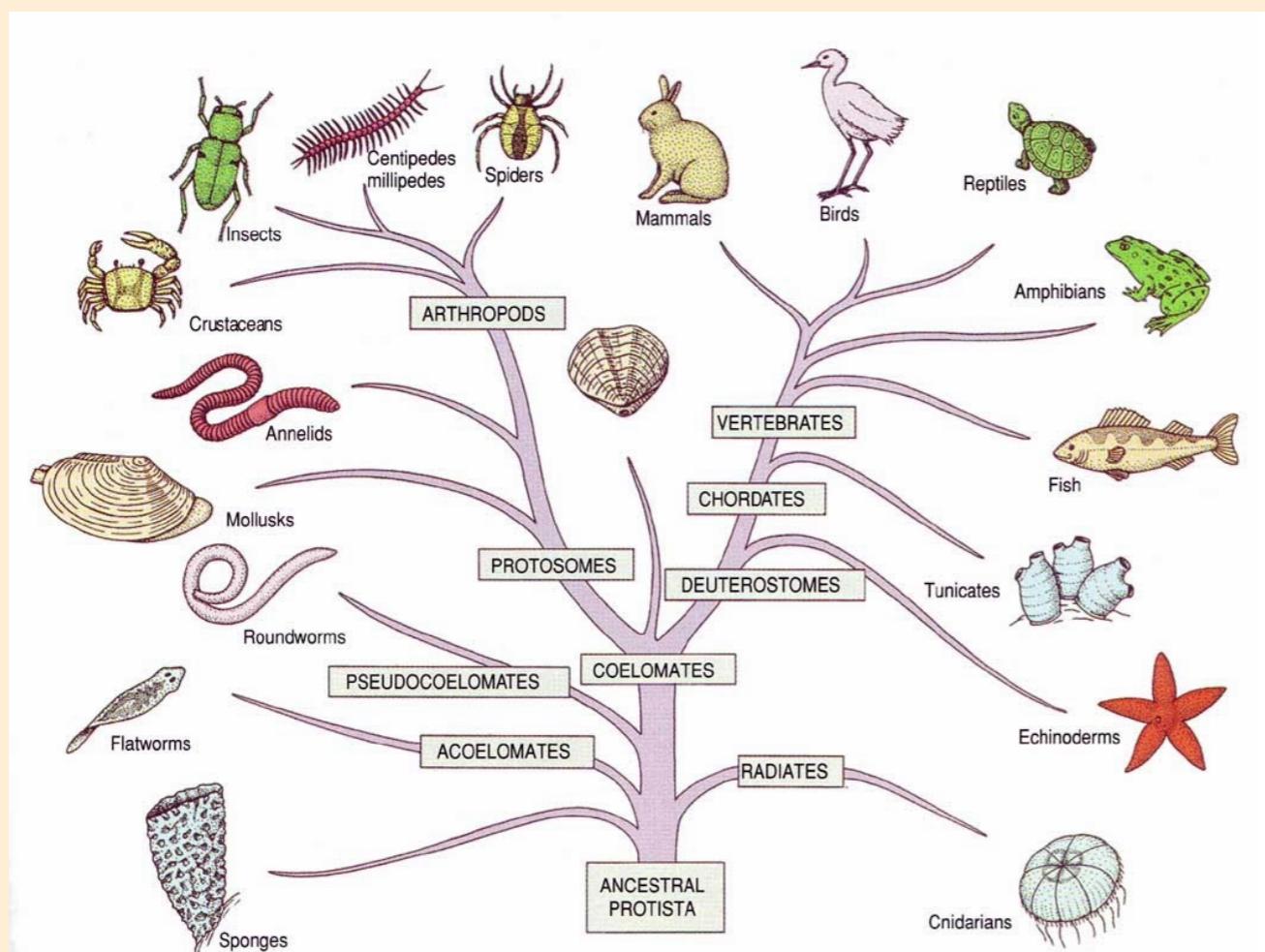
**NOTHING IN BIOLOGY MAKES
SENSE EXCEPT IN THE LIGHT
OF EVOLUTION**



T. Dobzhansky – 1973

EVERYTHING IN BIOLOGY MAKES MORE SENSE WITH A PHYLOGENETIC TREE

(muchos biólogos)

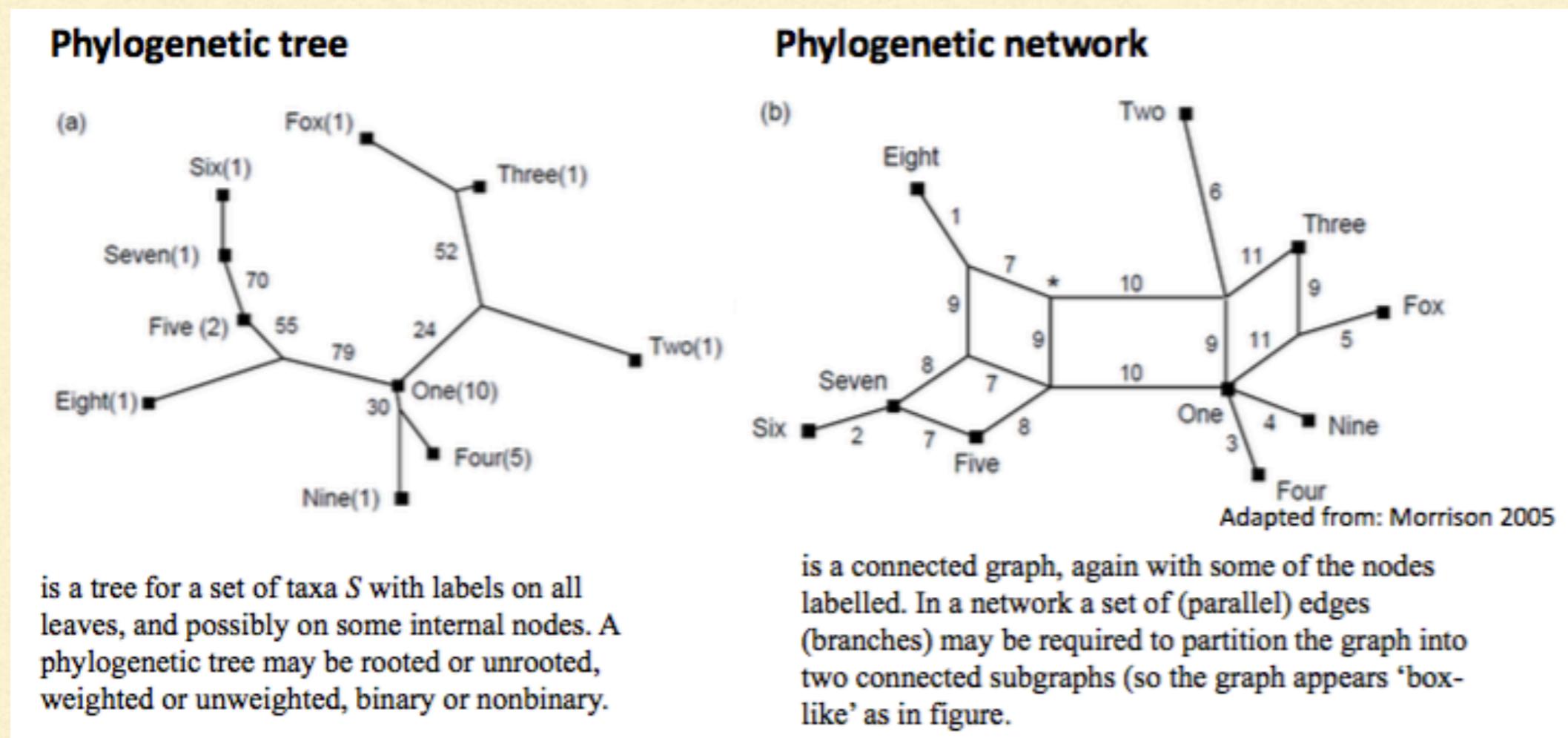


PHYLOGENETIC NETWORKS

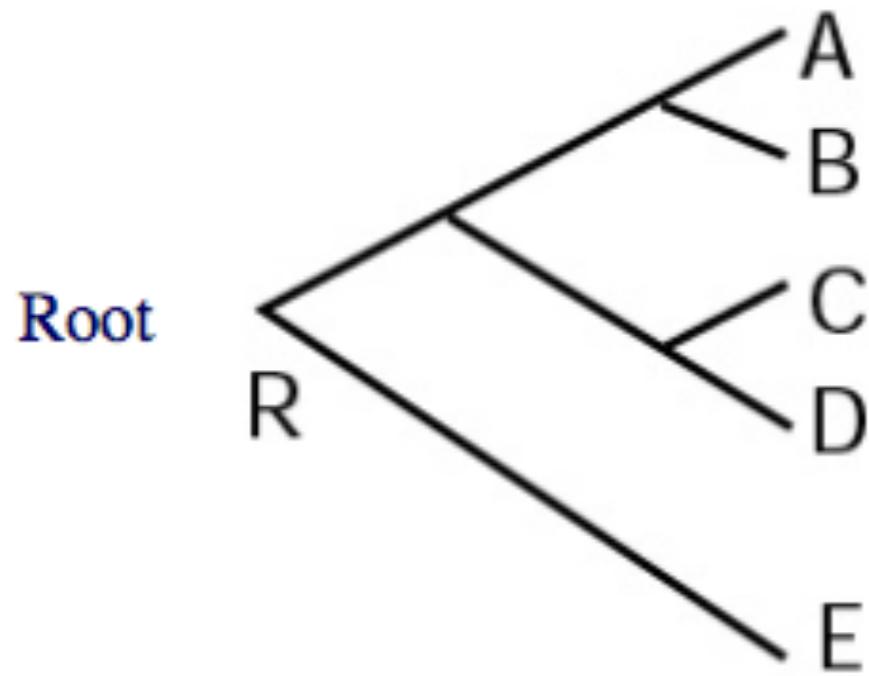
- “any” network in which taxa are represented by nodes and their evolutionary relationships are represented by edges. (For phylogenetic trees, edges are referred to as branches.) *Huson & Bryant 2006, Mol Biol Evol*

PHYLOGENETIC TREES AND NETWORKS

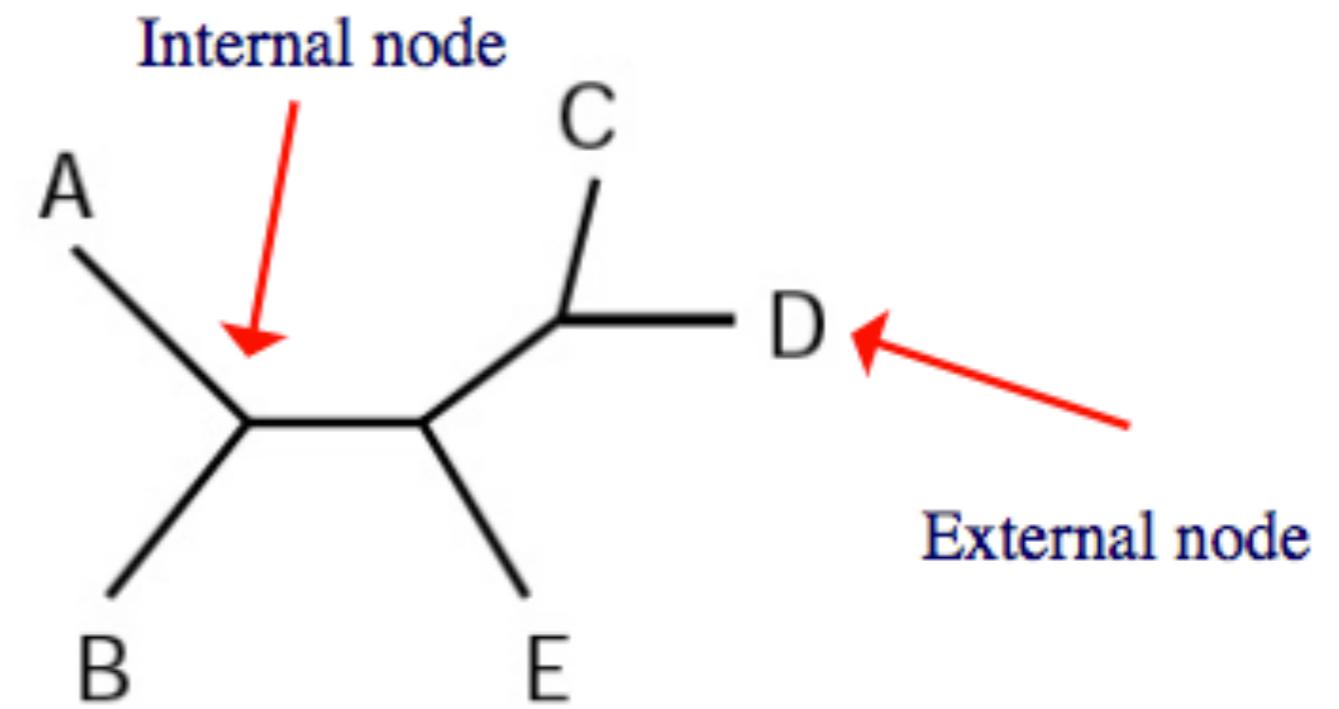
- Trees impose bifurcations
- Networks allow reticulations



ROOTED VS UNROOTED TREES



Rooted tree



Unrooted tree

INFERRING PHYLOGENIES

Inferring a tree is a combination of at least three components:

1. optimality criterion (parsimony, max likelihood, minimum evolution, least-squares fit, etc.).
2. search strategy (cluster methods, branch-and-bound, quartets, heuristic searches, etc.)
3. Assumptions about the mechanisms of evolution (JC, K2P, HKY, etc.)

METHODS FOR BUILDING TREES

- Distance-based
 - UPGMA
 - Neighbour Joining (NJ)
- Character-based
 - Maximum Parsimony (MP)
 - Maximum Likelihood (ML)
 - Bayesian methods (Markov Chain Monte Carlo MCMC)

DISTANCE-BASED TREES

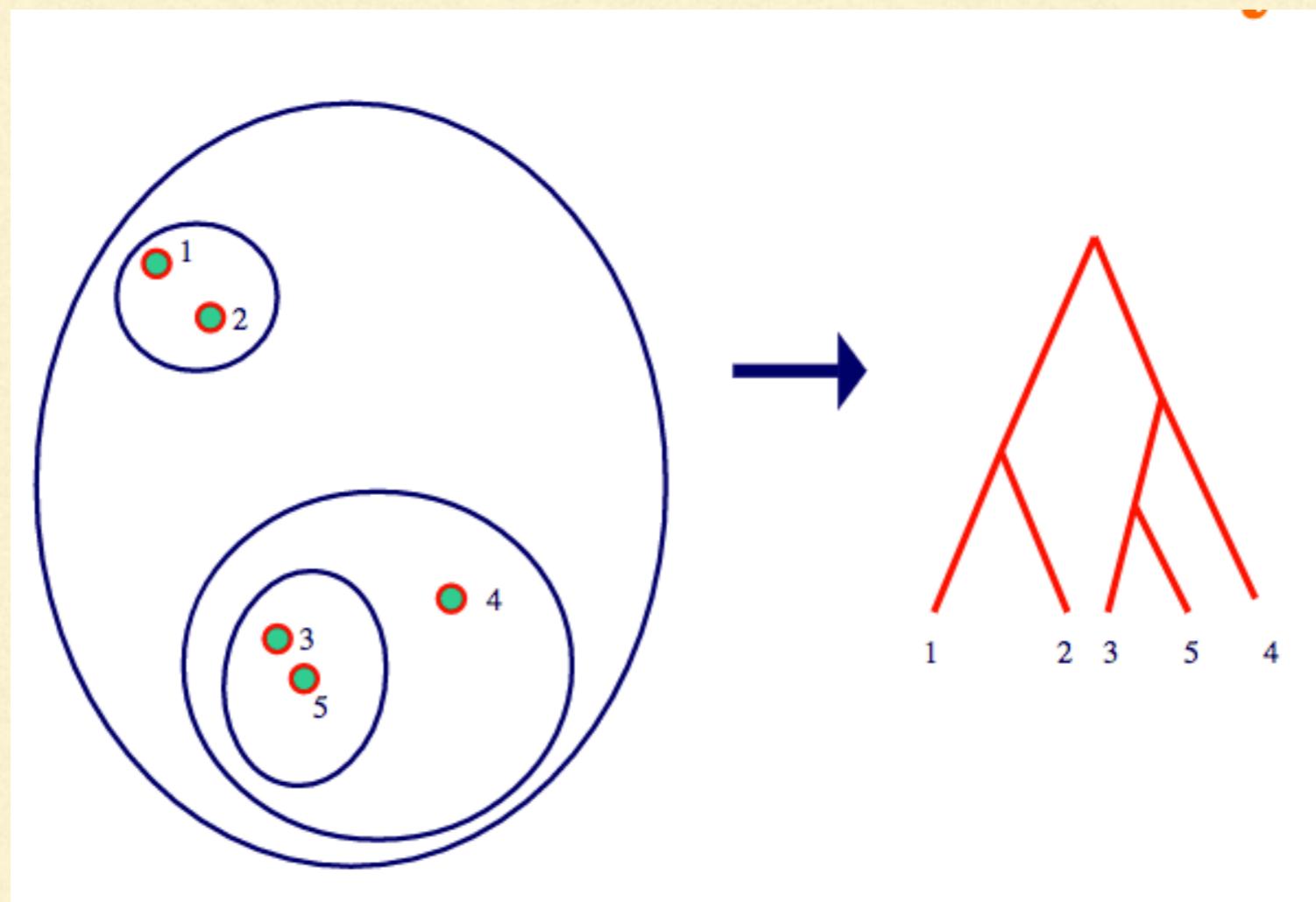
- First calculate distance matrix between pairs of sequences or populations
- Then build a tree

DISTANCE-BASED TREES

- UPGMA (Sokal and Sneath 1963): based on the molecular clock assumption generates ultrametric trees.
 - Rooted tree
 - all the end nodes are equidistant from the root
 - assuming a **molecular clock**.
- agglomerative (bottom-up) hierarchical clustering method. Picks the closest pair of neighbors, and adds the closest, and so on

DISTANCE-BASED TREES

- UPGMA (Sokal and Sneath 1963): based on the molecular clock assumption generates ultrametric trees.

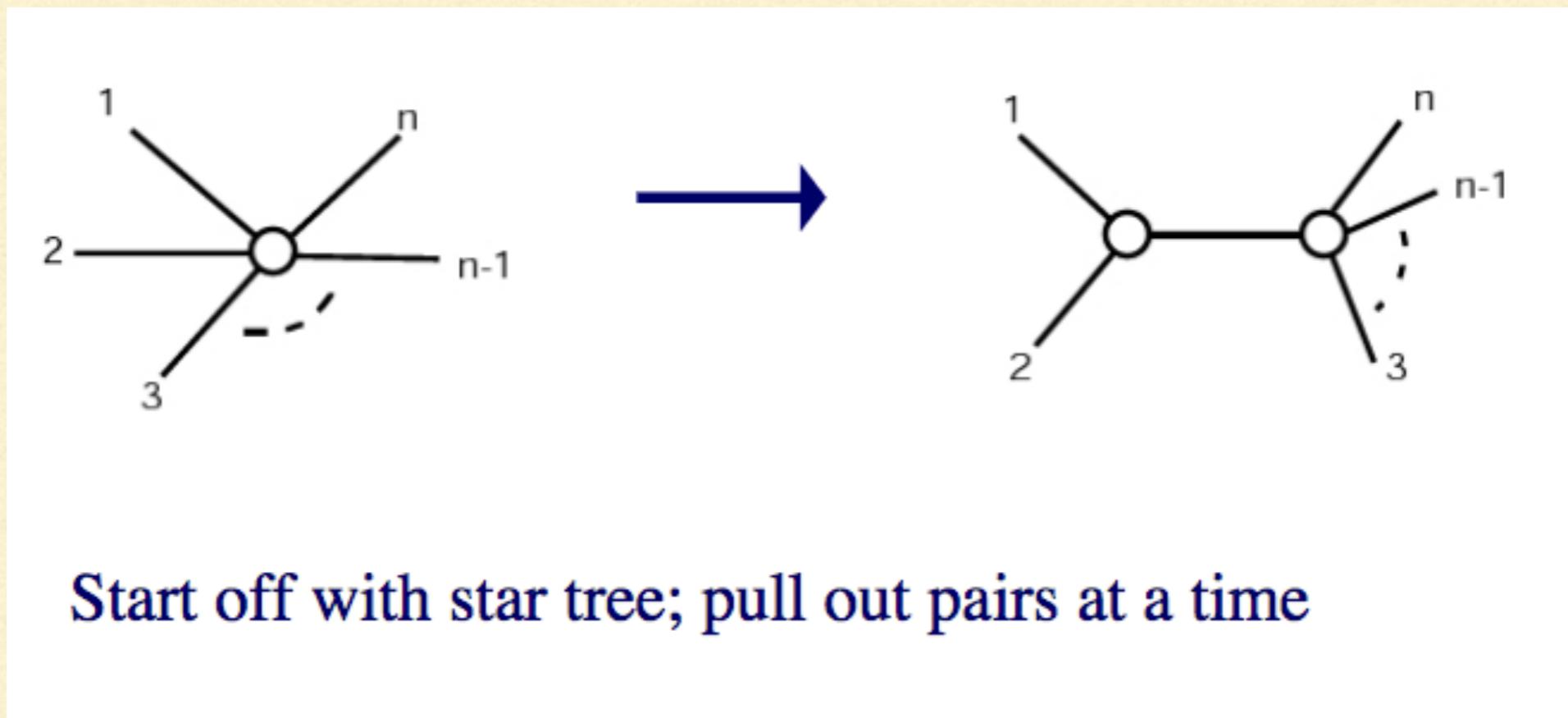


DISTANCE-BASED TREES

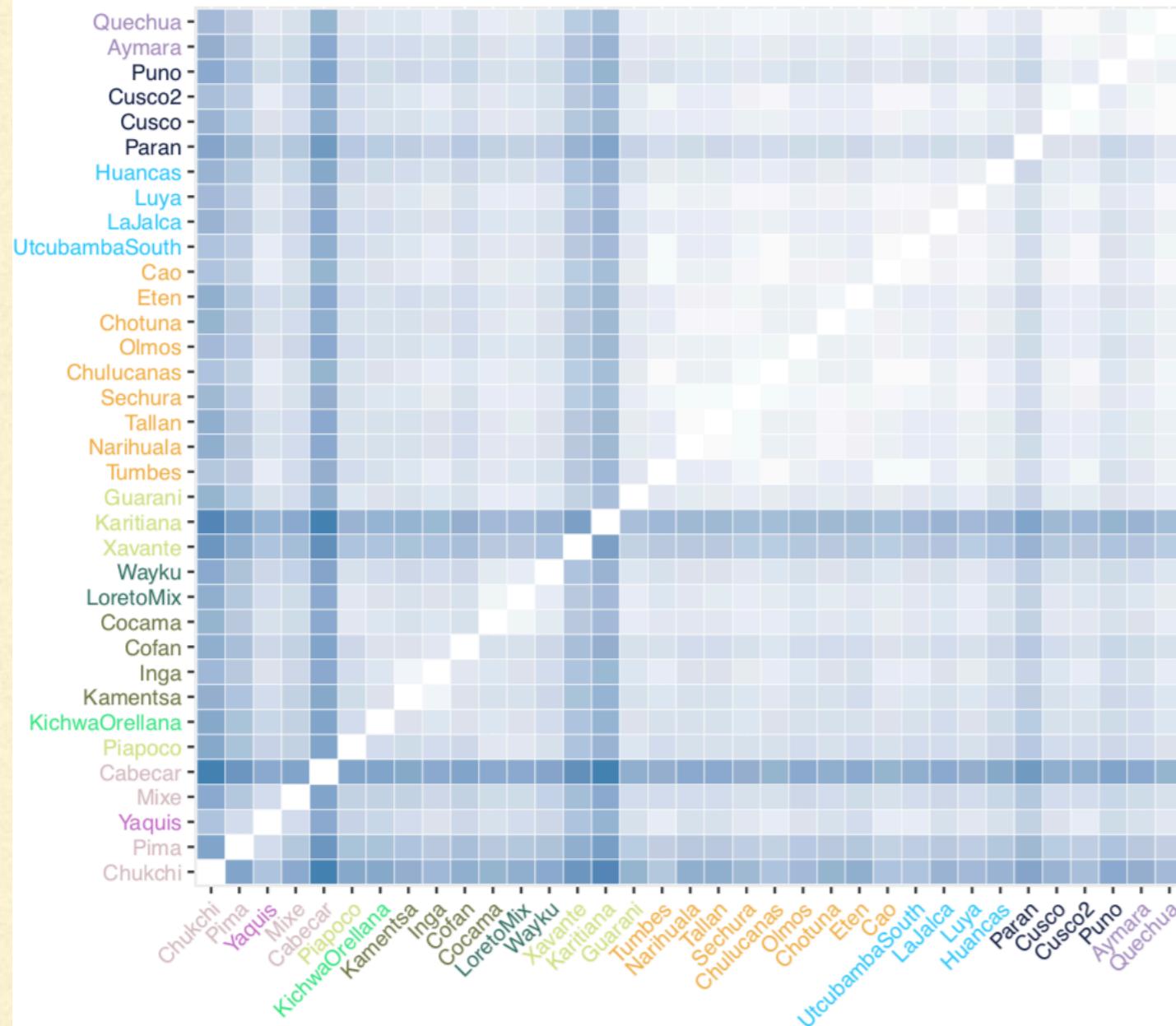
- **Neighbor-Joining NJ** (Saitou and Nei 1987)
 - Unrooted tree
 - Does not assume a **molecular clock**.
- Local search strategy using a Minimum Evolution (ME) optimality criteria
- Starts with an unresolved star-like tree, calculate the sum of branch length. Joins the pair with the closest branch length. And so on

DISTANCE-BASED TREES

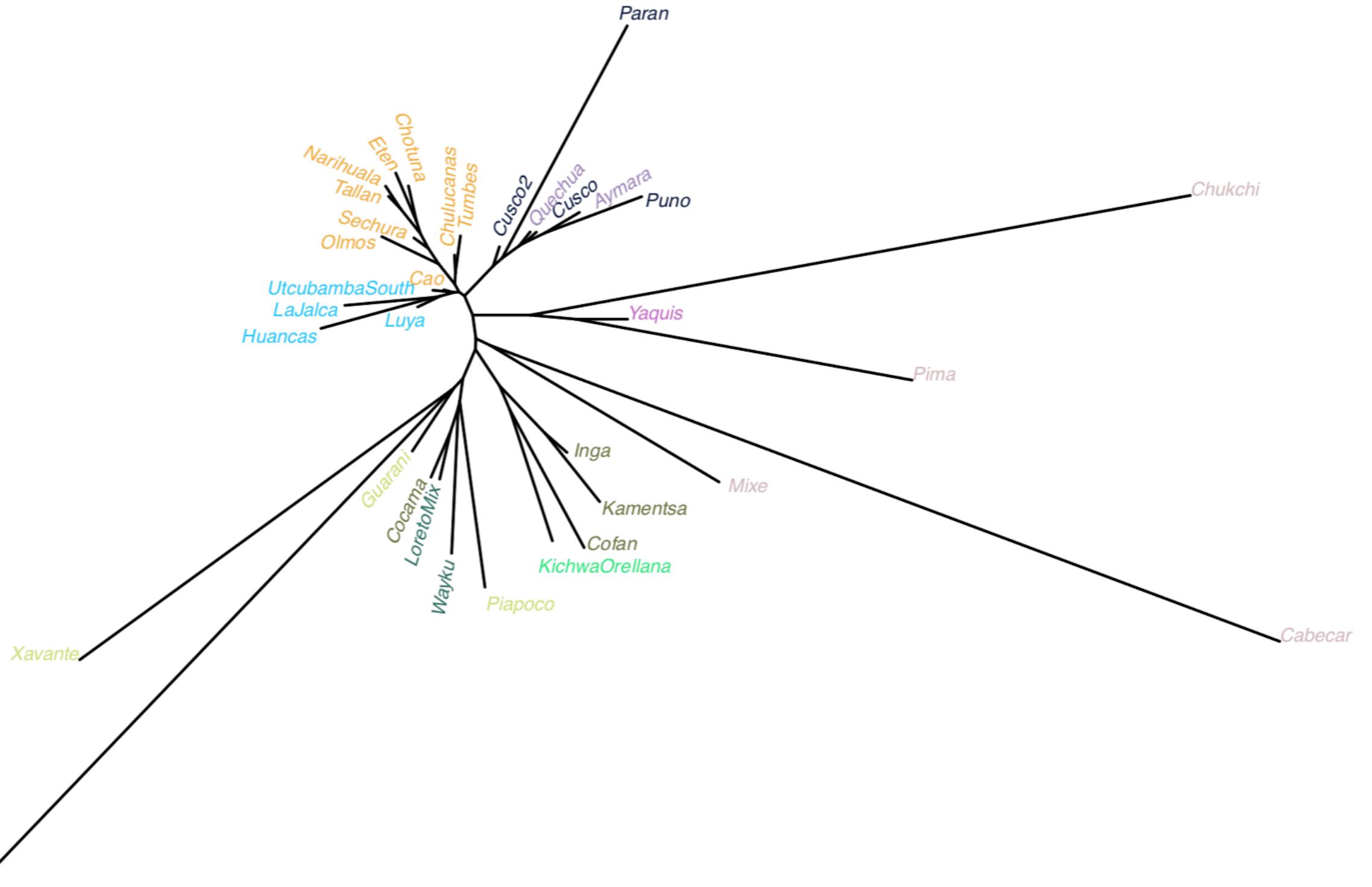
- **Neighbor-Joining NJ** (Saitou and Nei 1987)



EXAMPLES



EXAMPLES



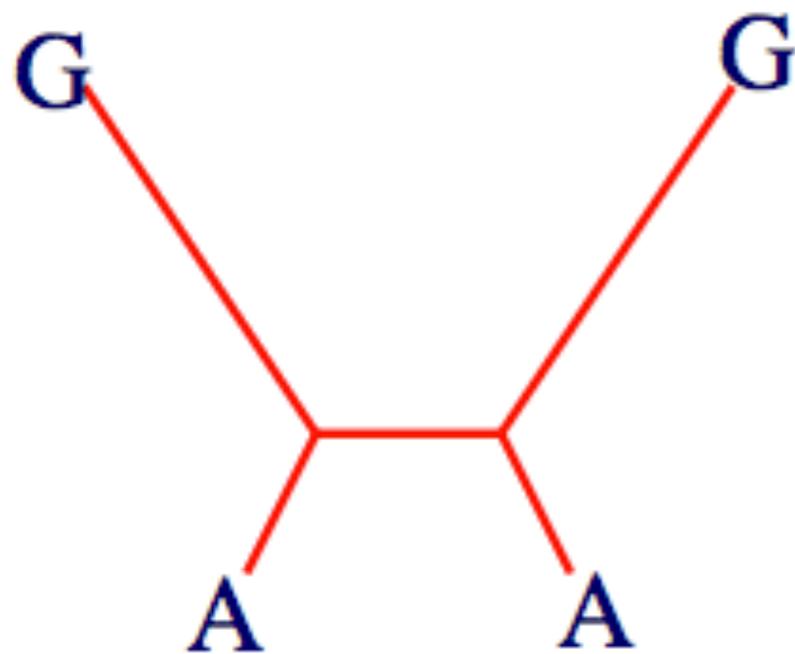
CHARACTER BASED TREES

- Maximum parsimony (MP): choose tree that minimizes number of changes from a common ancestor
 - MP yields more than one tree with the same score
- Maximum likelihood (ML): find the tree which gives the highest likelihood of the observed data

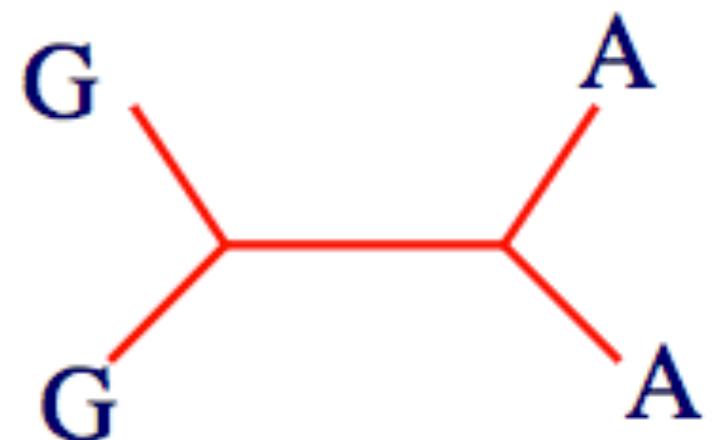
They both imply model of evolution

PARSIMONY WEAKNESS: LONG BRANCH ATTRACTION

- Parsimony analysis implicitly assumes that rate of change along branches are similar



**Real tree: two long branches
where G has turned to A independently**



Inferred tree

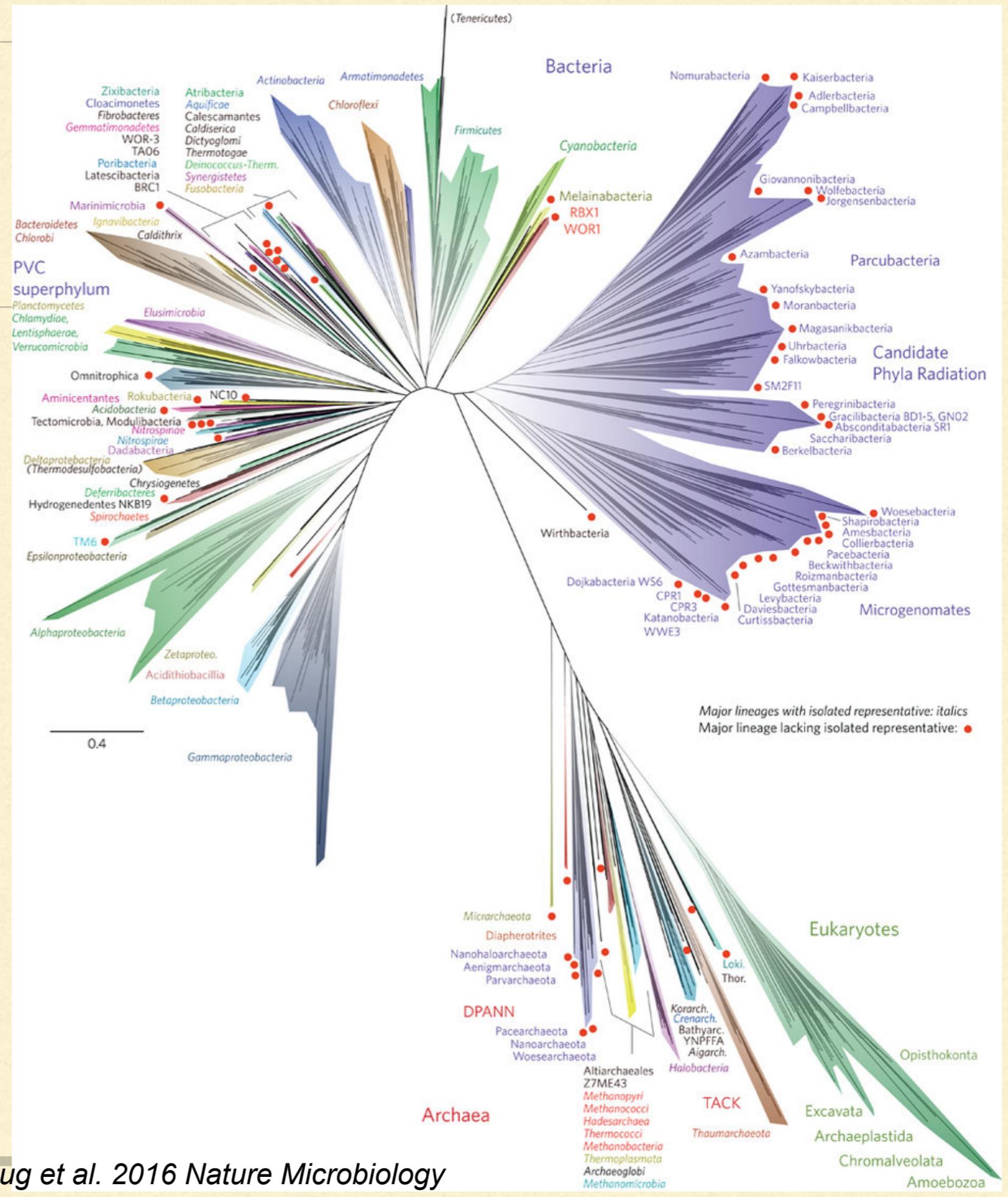
SUMMARY:TREES

- Distance methods are good for large data sets of highly similar sequences
- **Likelihood** and **Bayesian** methods often have more power and are more robust, especially for inferring deep phylogenies

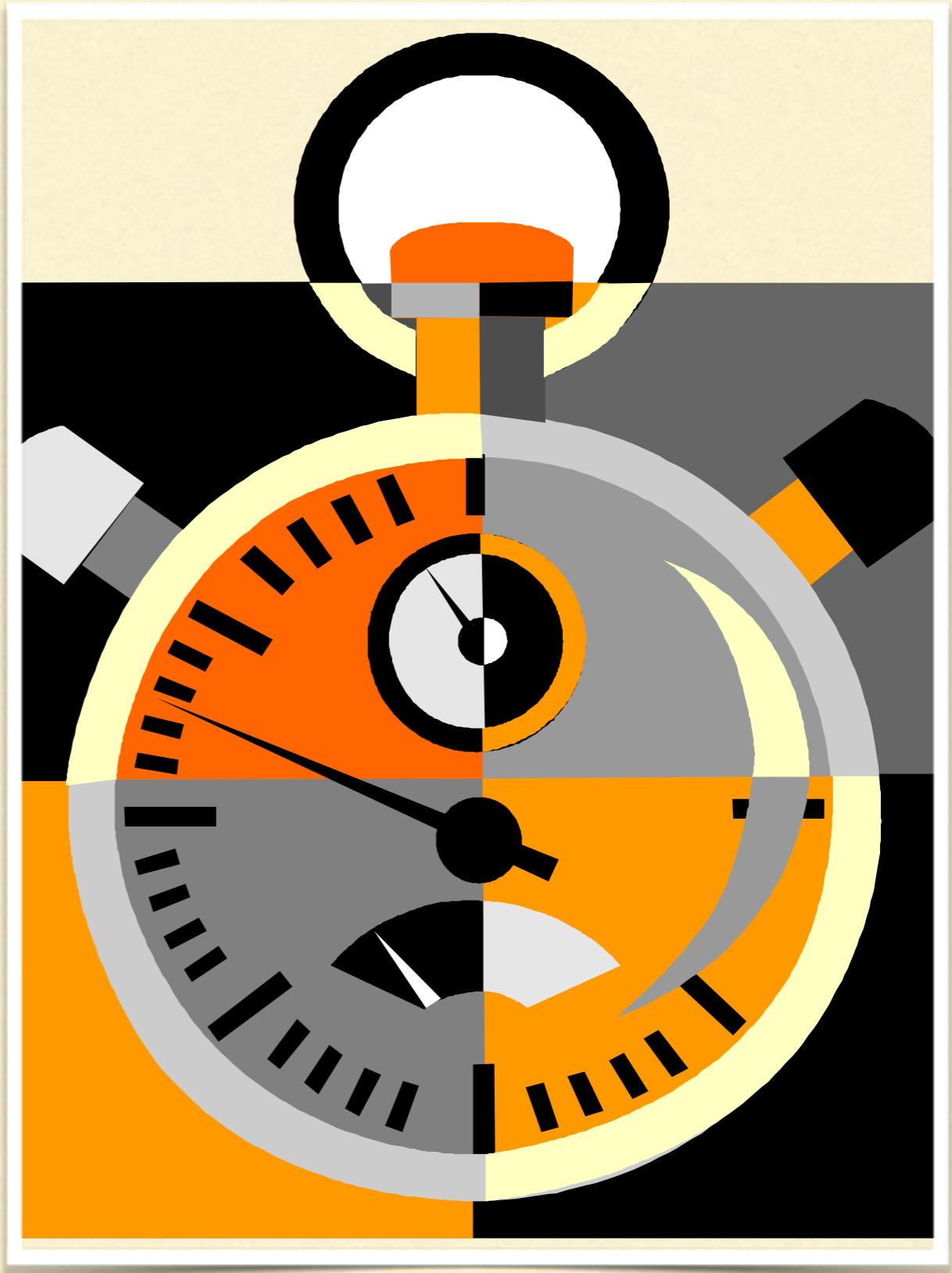
Battle between preferences:

- Many people like **Max Likelihood** based methods:
 - sensitive at large evolutionary distances
- Often a **BEAST** tree is the answer
 - But takes computational time
 - Advantage of including complex models with priori assumptions

ML TREE OF LIFE



MOLECULAR CLOCK



THE MOLECULAR CLOCK HYPOTHESIS

- Amount of genetic difference between sequences is a function of time since separation.
- Rate of molecular change is constant (enough) to predict times of divergence

HOW TO CALIBRATE HUMAN GENETICS

- Deep pedigree data
- Count the mutations between nth grade cousins

I'm going to name you after your father and grandfather so genealogists have a heck of a time trying to research you in the next century.

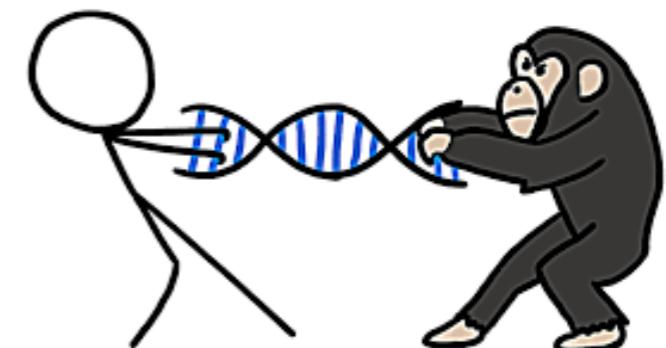
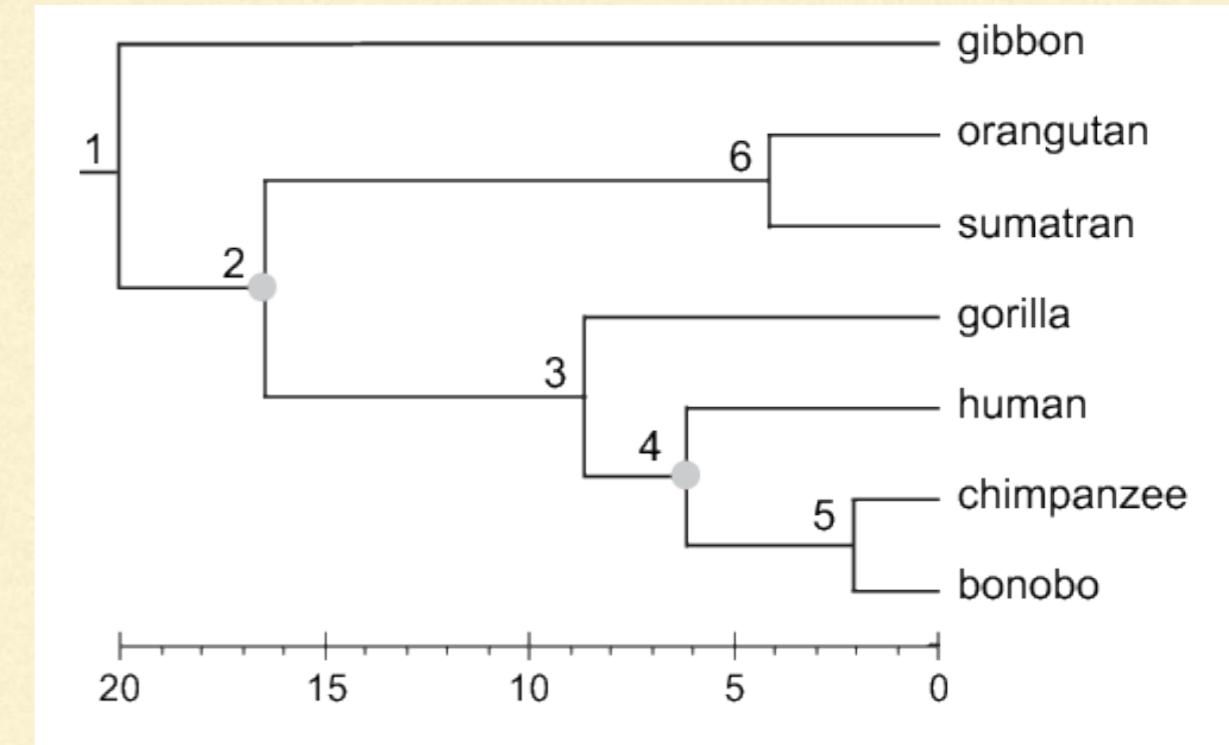
som~~e~~ecards
user card



HOW TO CALIBRATE HUMAN GENETICS

■ ARCHAEOLOGY

- Species divergence - E.g Human-chimp split
- Historical events E.g Colonization of the pacific



HOW TO CALIBRATE HUMAN GENETICS

- Tip fossils
 - aDNA from dated fossils



GENETIC MARKERS

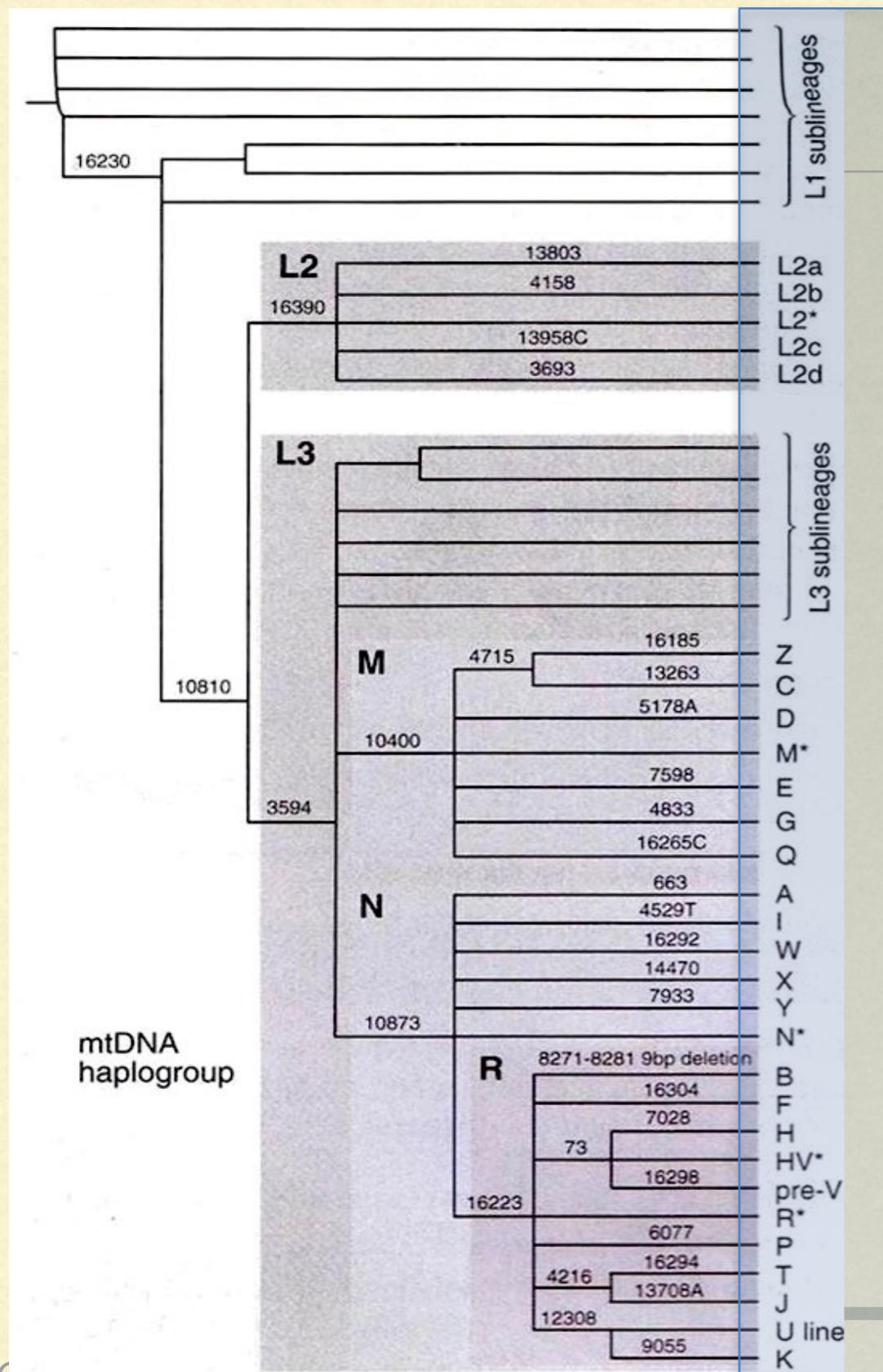
Uniparental markers:

- No recombination (HAPLOID)
- Mutations accumulate with time alone
- Shared mutations indicate shared ancestry (if mutations are rare)

LIMITS

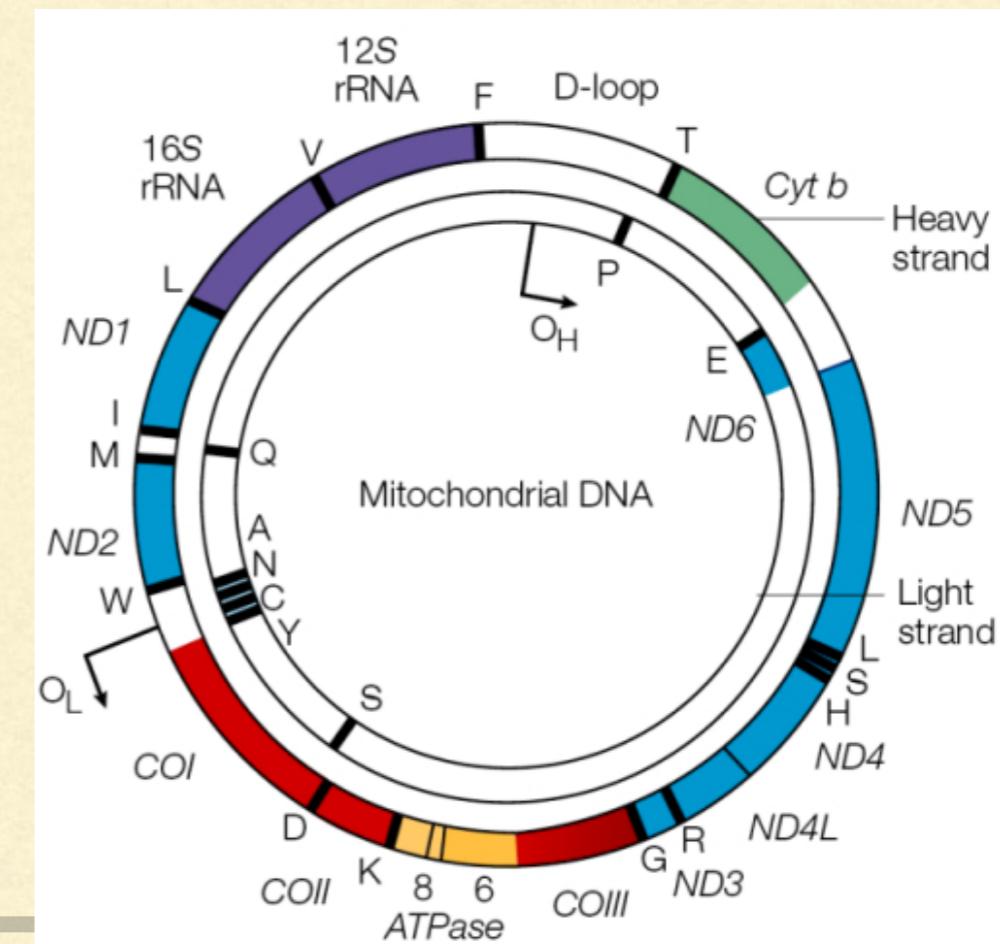
- Small portion of human genome
- Only maternal/paternal view of history

HAPLOGROUP PHYLOGENY

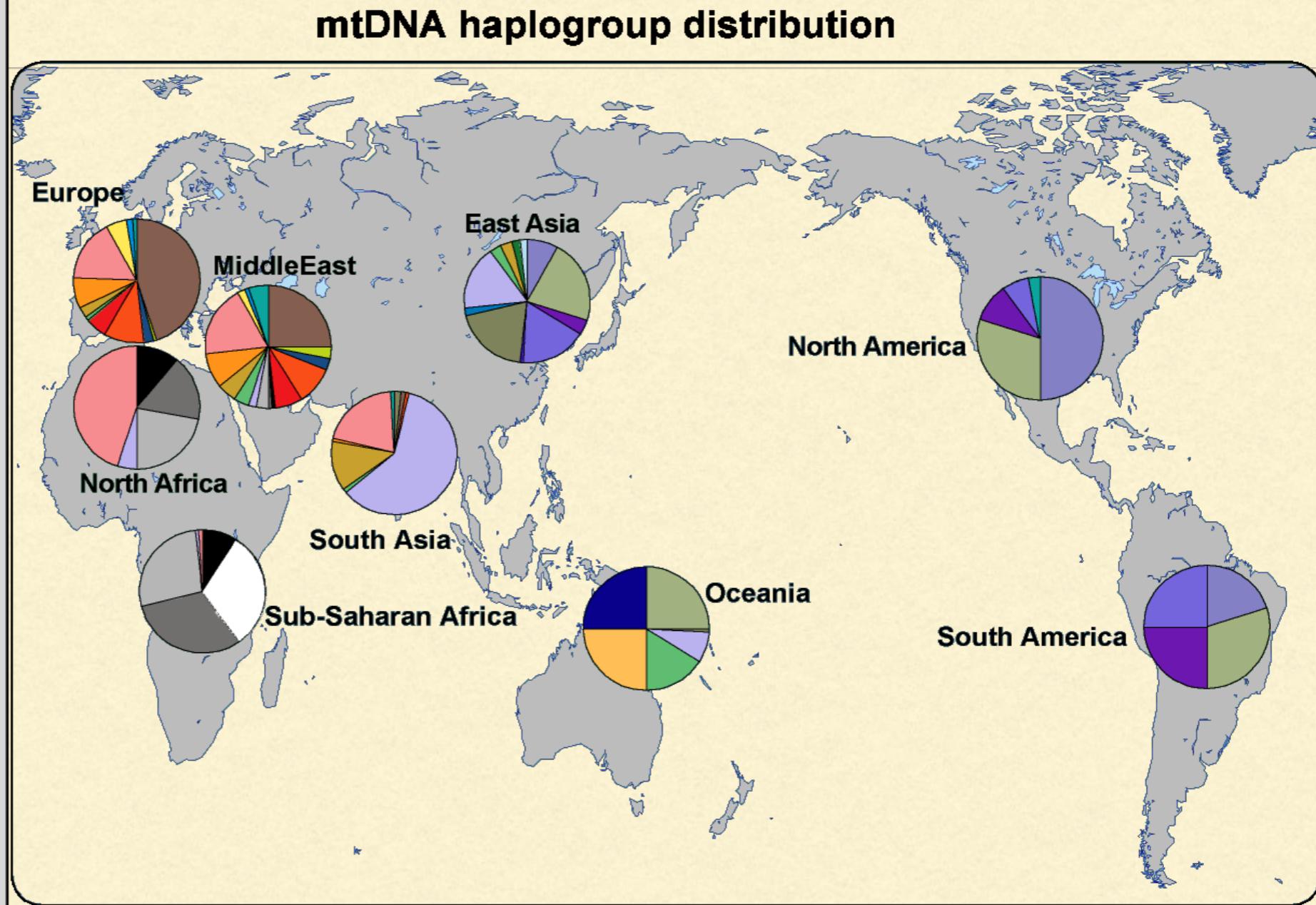
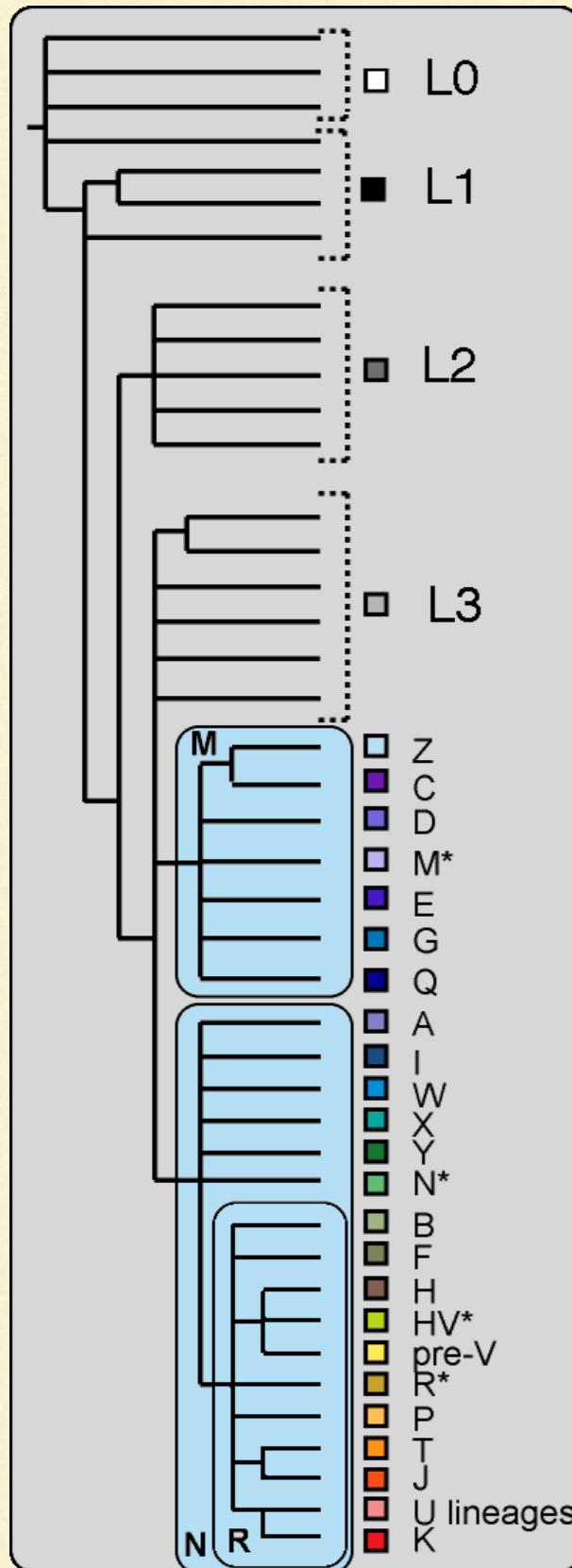


Phylogeny
Tree structure representing evolutionary relationships between clades

Haplogroups indicated with a capital letter

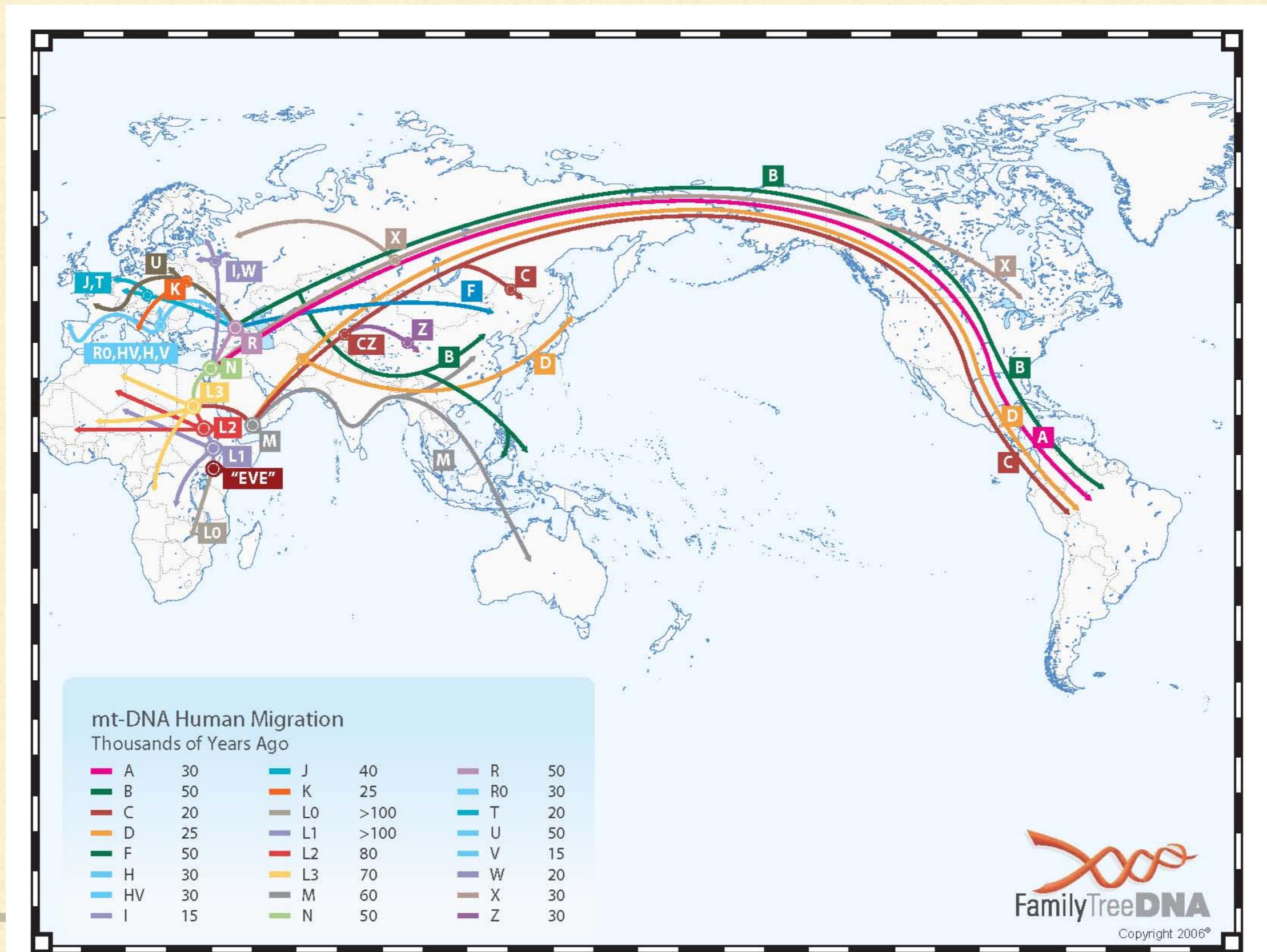


WORLD MTDNA PHYLOGEOGRAPHY



It combines temporal and evolutionary dimension of phylogeny
with geographic distribution of haplogroups

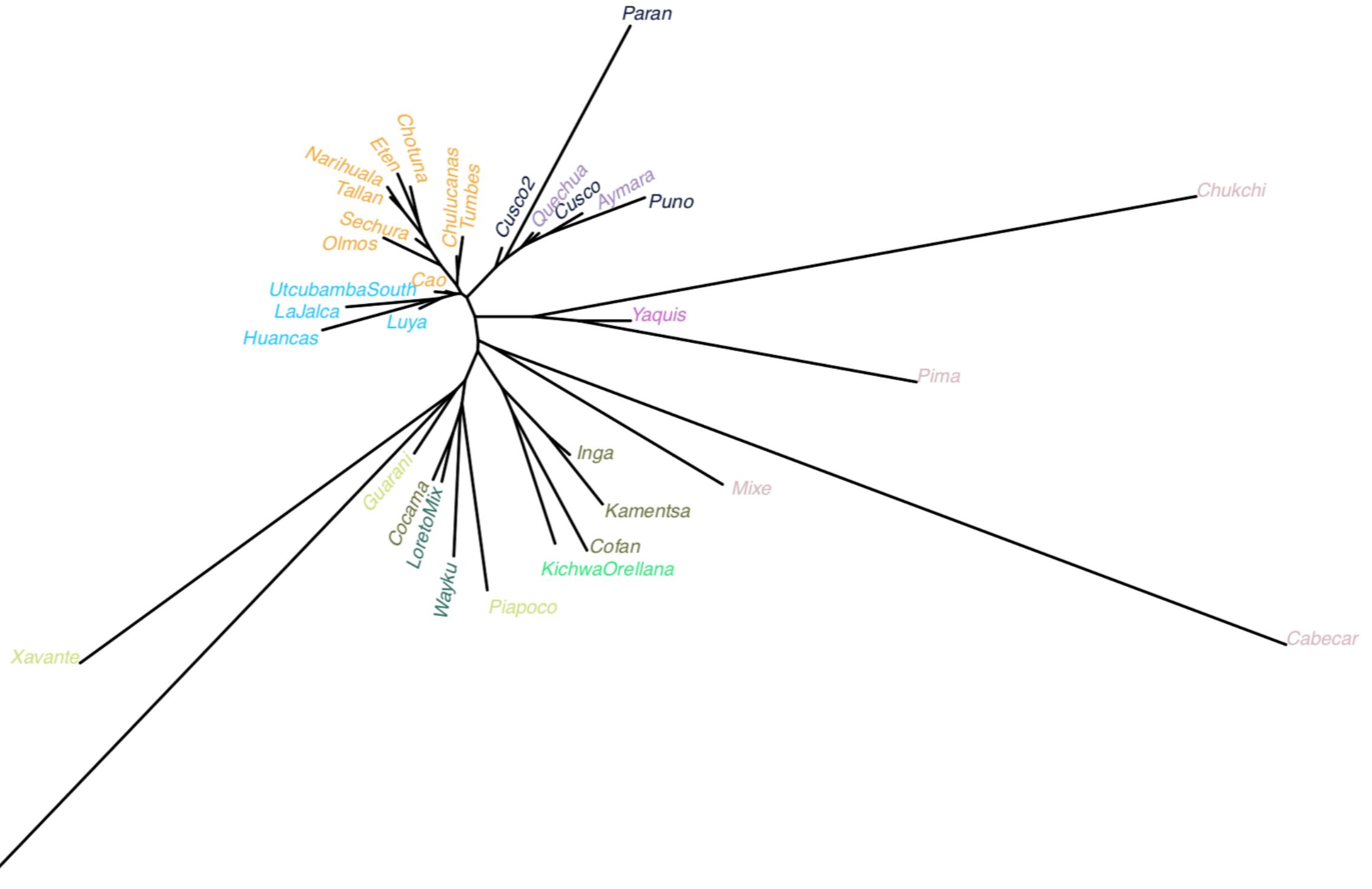
WORLD MTDNA PHYLOGEOGRAPHY



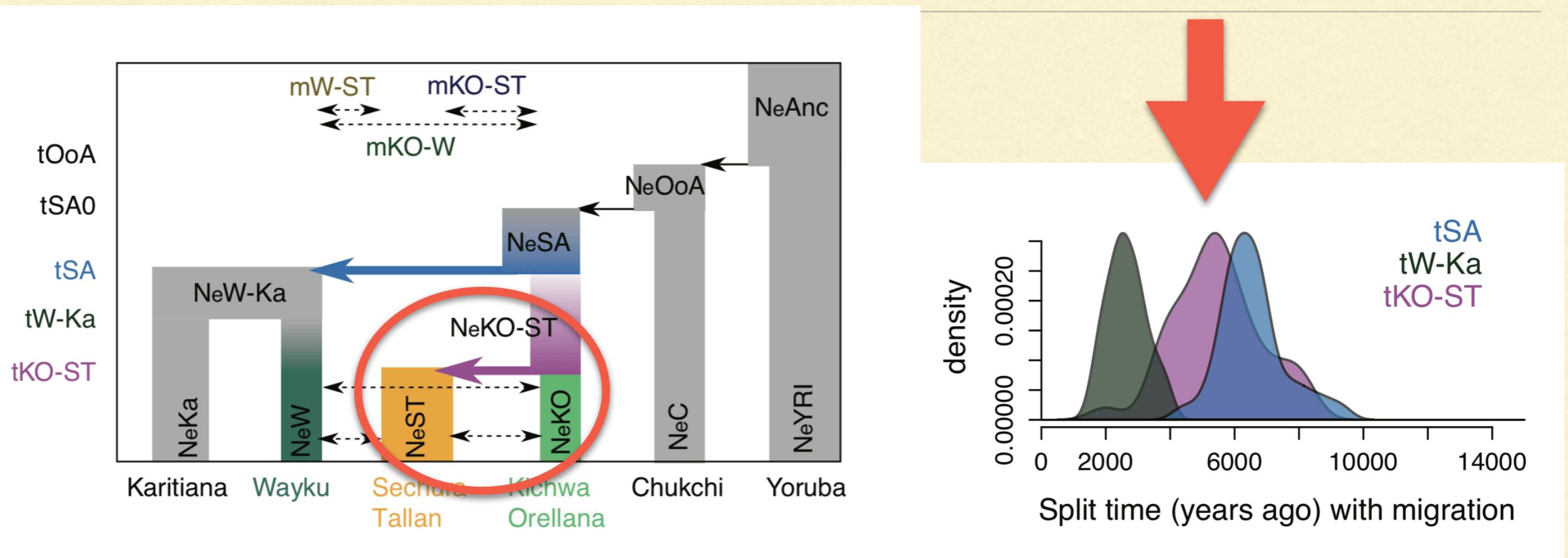
COMO RECONSTRUIR LA HISTORIA FILOGENETICA DE POBLACIONES?

- Distancia genetica entre poblaciones (FST)
- Simulaciones de escenarios demográficos
- Admixture graphs, TreeMix

EXAMPLES

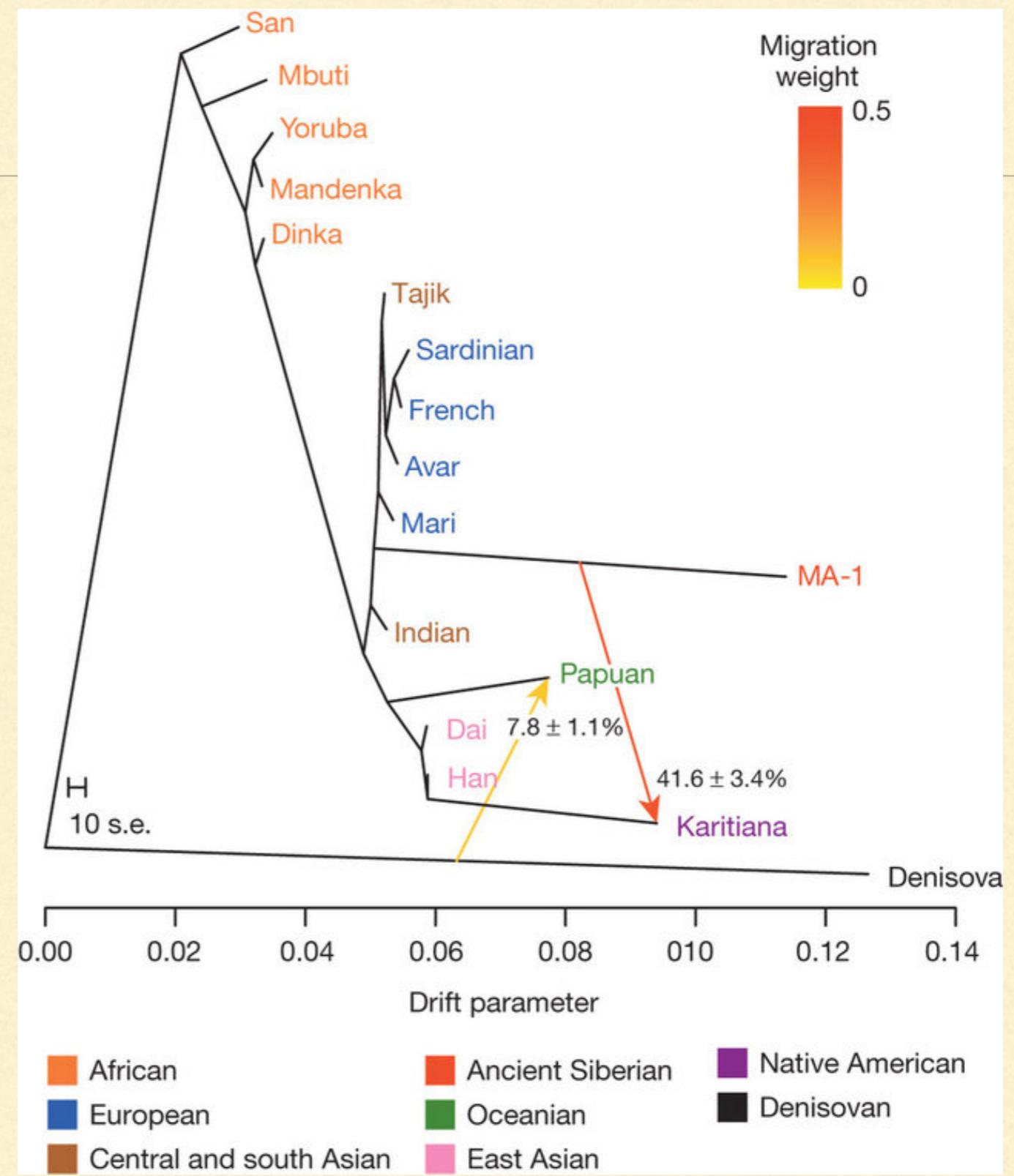


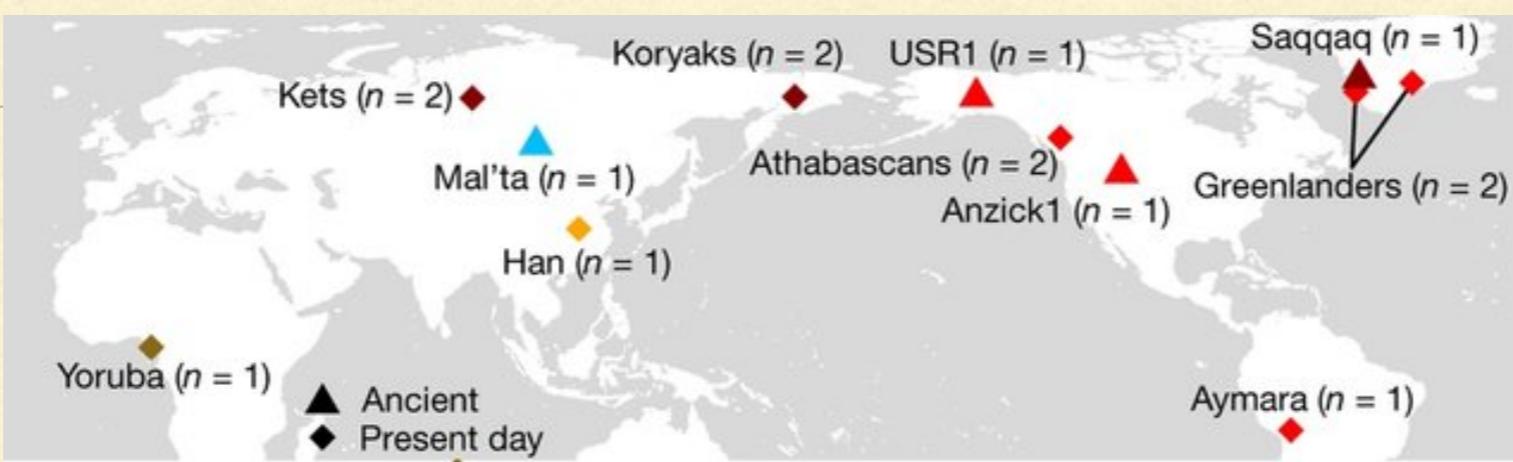
ANCESTRIES IN SOUTH AMERICA: DIVERGENCE TIMES



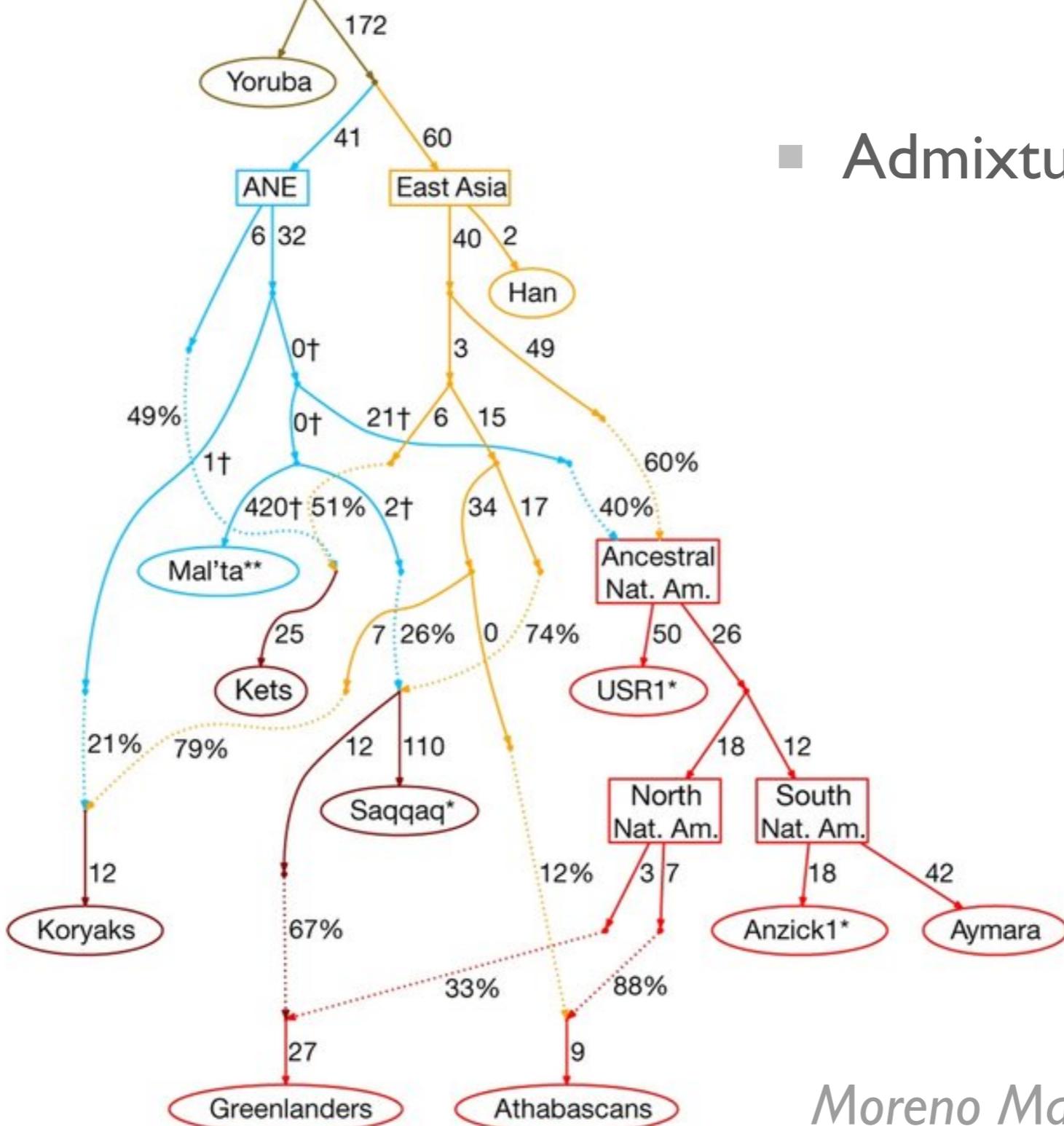
Simulating demographic scenarios, estimating parameters with ABC framework

TreeMix

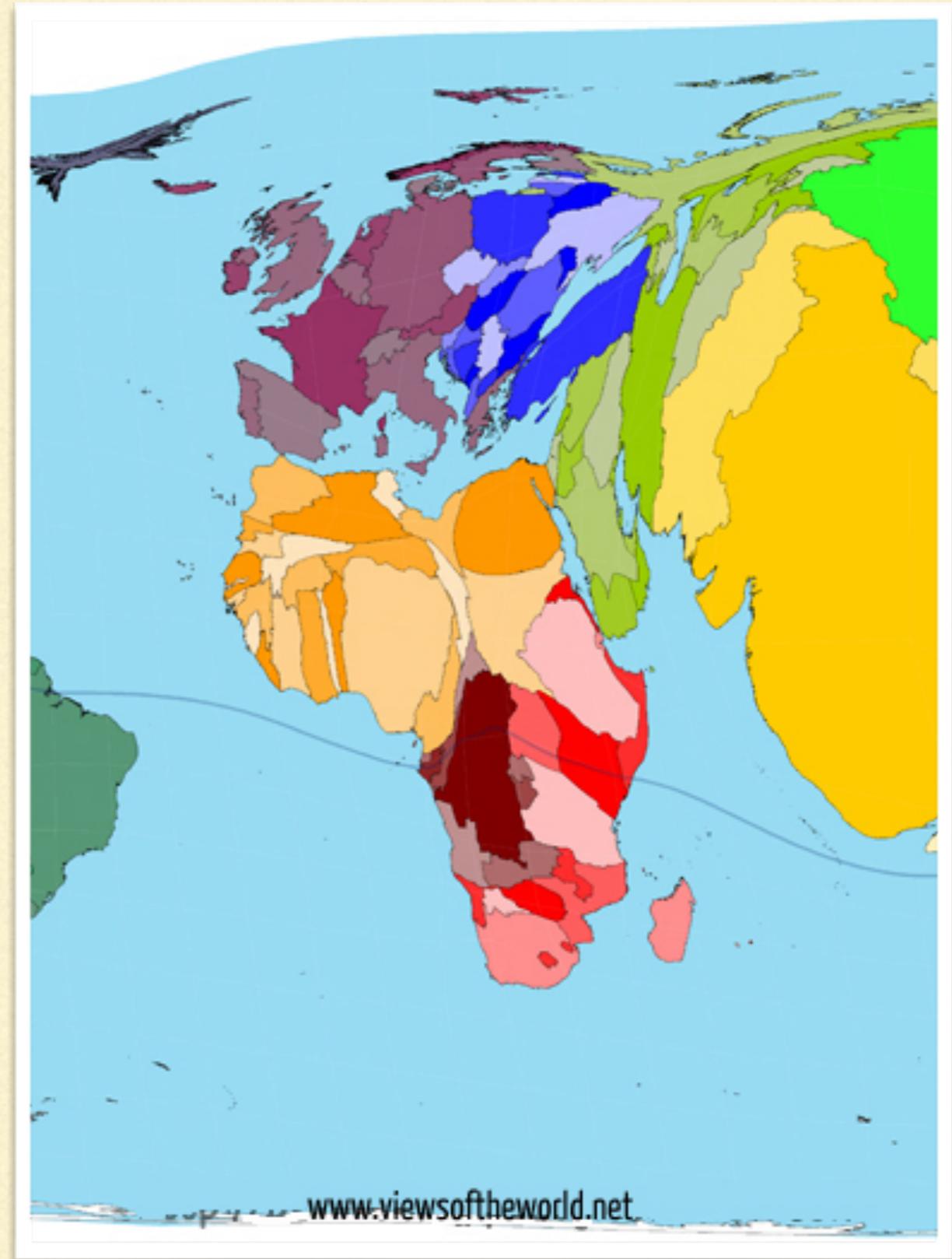




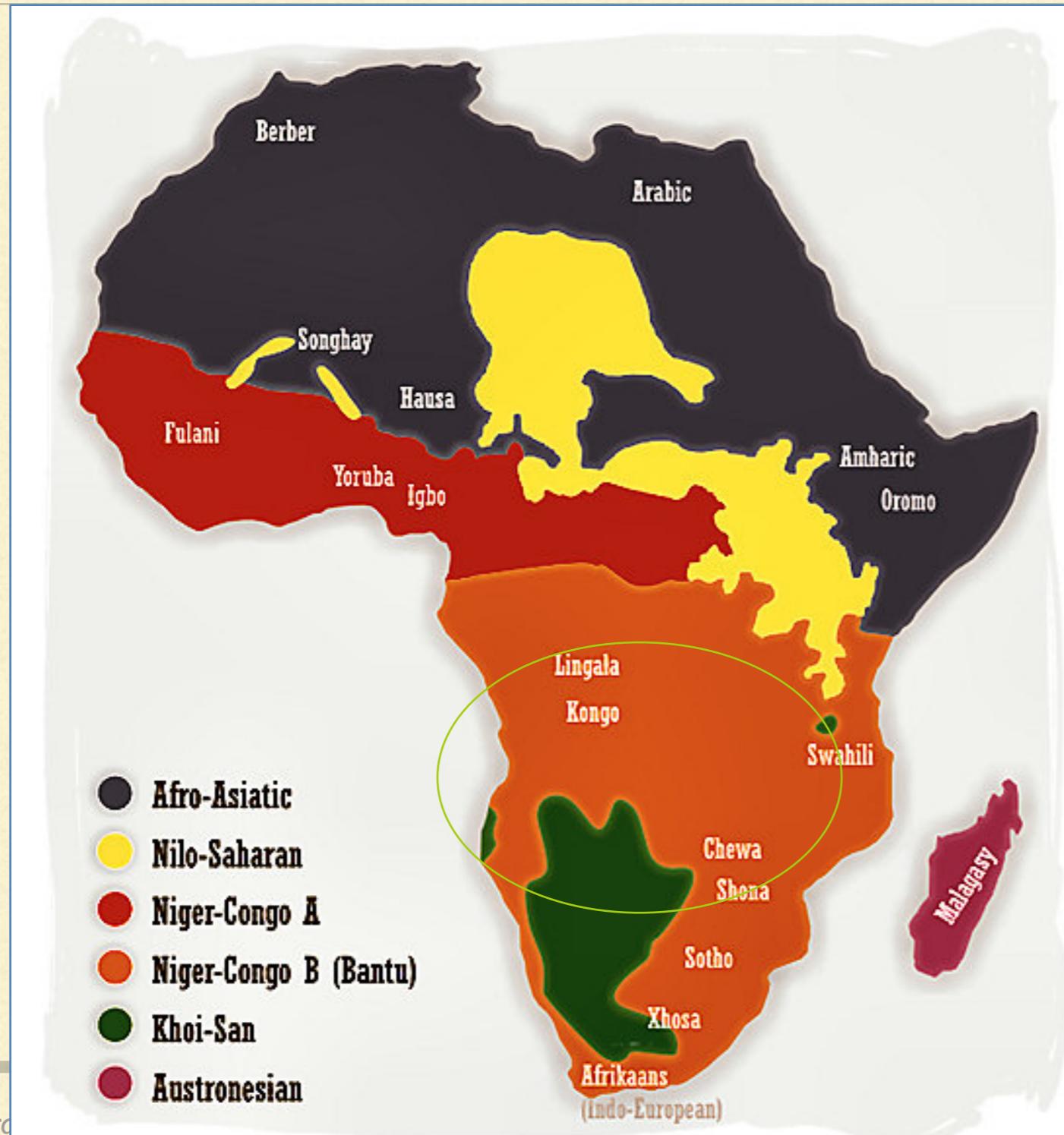
■ Admixture graph



HUMAN POPULATION HISTORY AND RELEVANCE FOR GENOMIC STUDIES



LINGUISTIC DIVERSITY IN AFRICA



- 4 linguistic stocks
- **Bantu** is the most diffused language family

BANTU-SPEAKING PEOPLE

- Language spread with human migration
- Starting ~5.000 kya from Cameroon
- Diffusion of agriculture

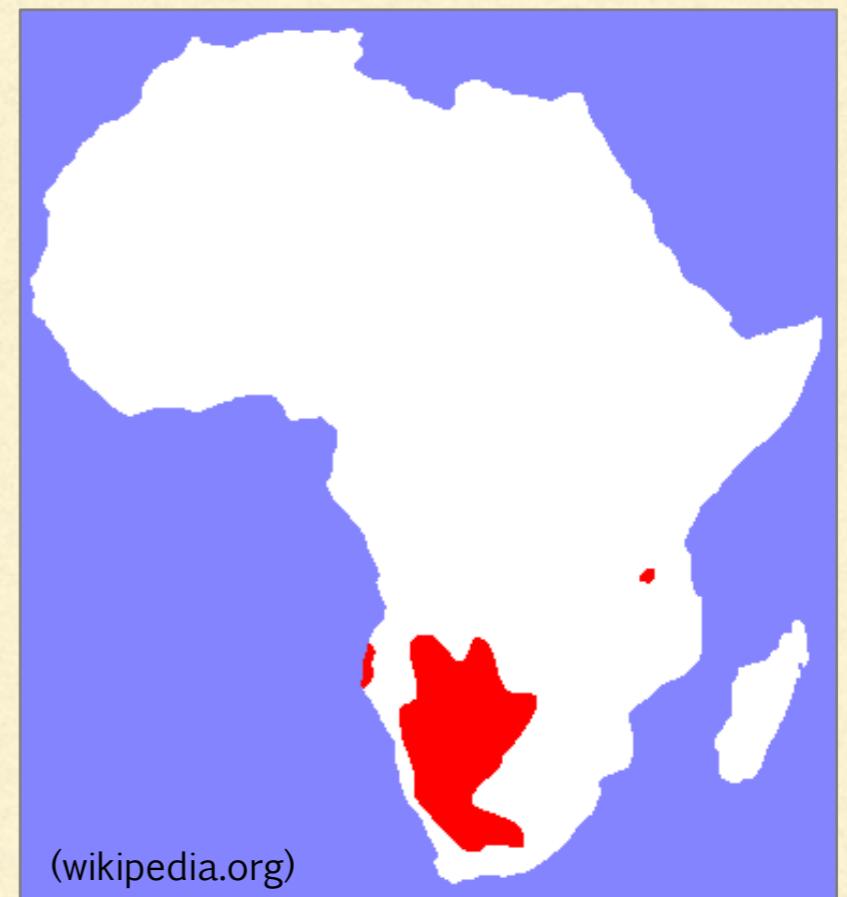


KHOISAN PEOPLE



(Credits: Kure B)

- Non-Bantu groups
- click consonant as prominent linguistic feature
- pastoralist Khoi and hunter-gatherer San



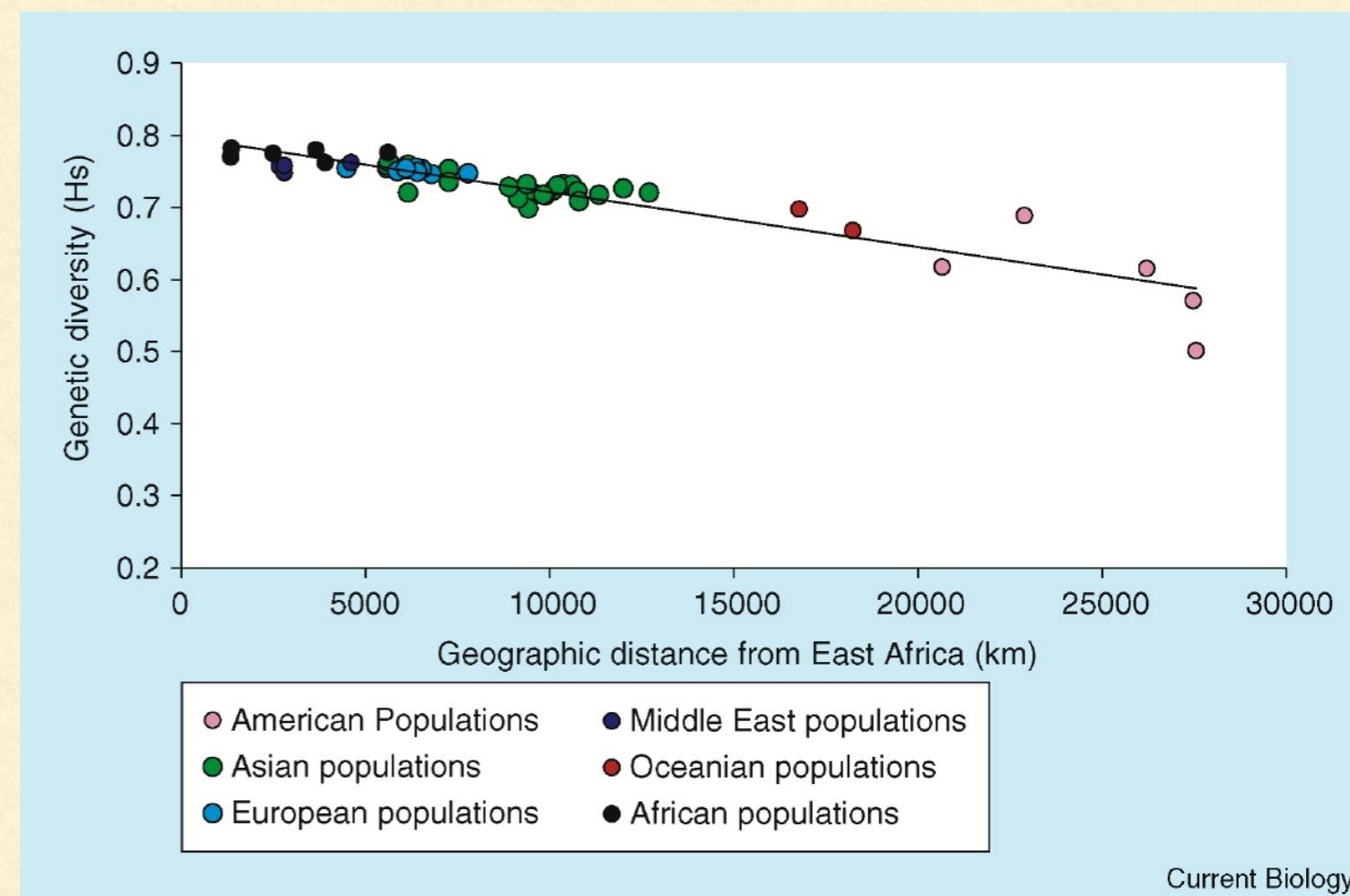
(wikipedia.org)

KHOISAN

- Three major lineages:
 - **Tuu**
 - **Kx'aa**
 - **Khoe-Kwadi**
- Shared linguistic features (clicks) might be result of contact

HUMAN ORIGINS IN AFRICA

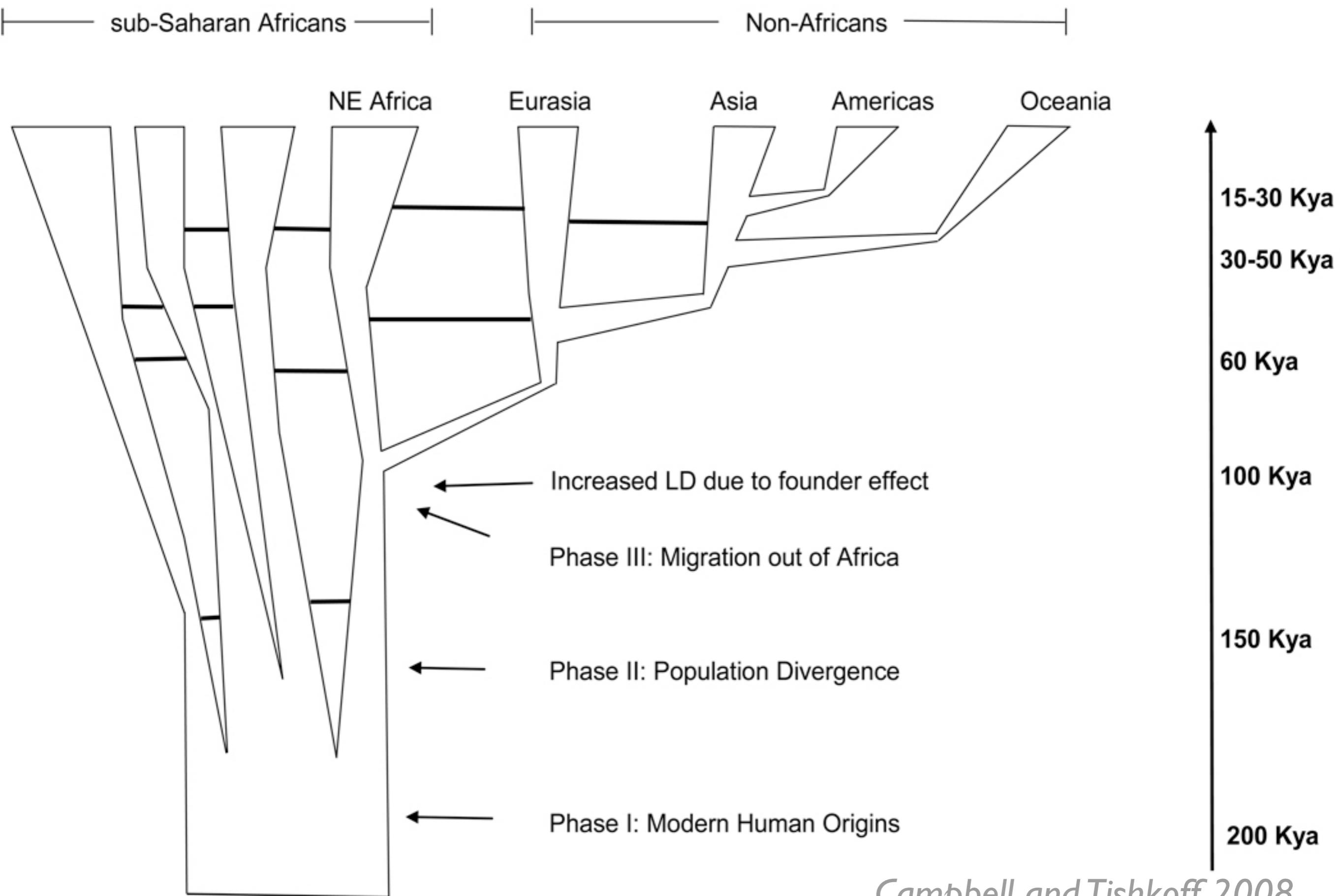
- Africa harbors the highest genetic diversity
- Out of Africa: Major bottleneck in the migration which colonised the rest of the continents ~70 cya



Current Biology

Divergent pattern of LD

Shared Pattern of LD

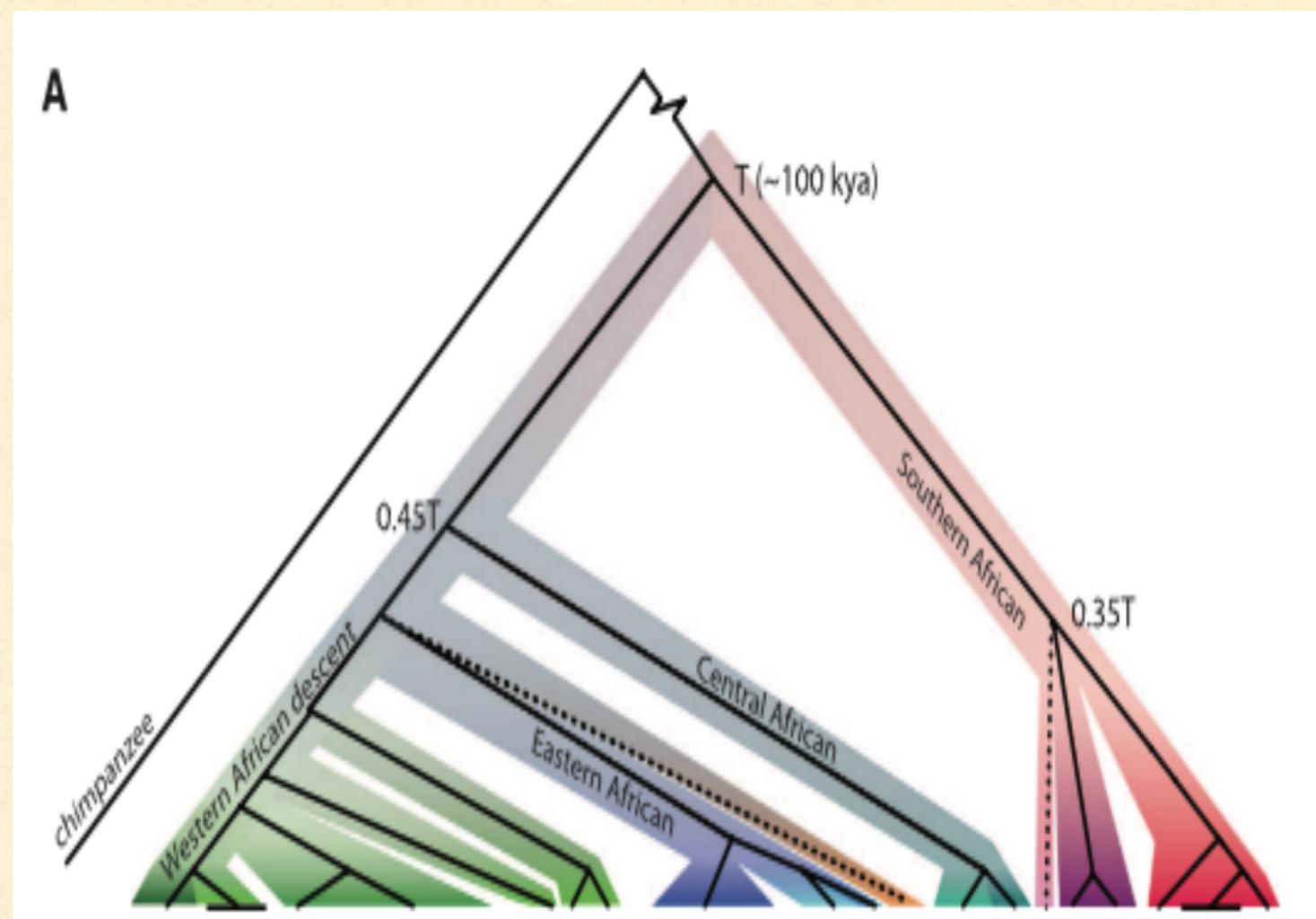


WHERE IN AFRICA?

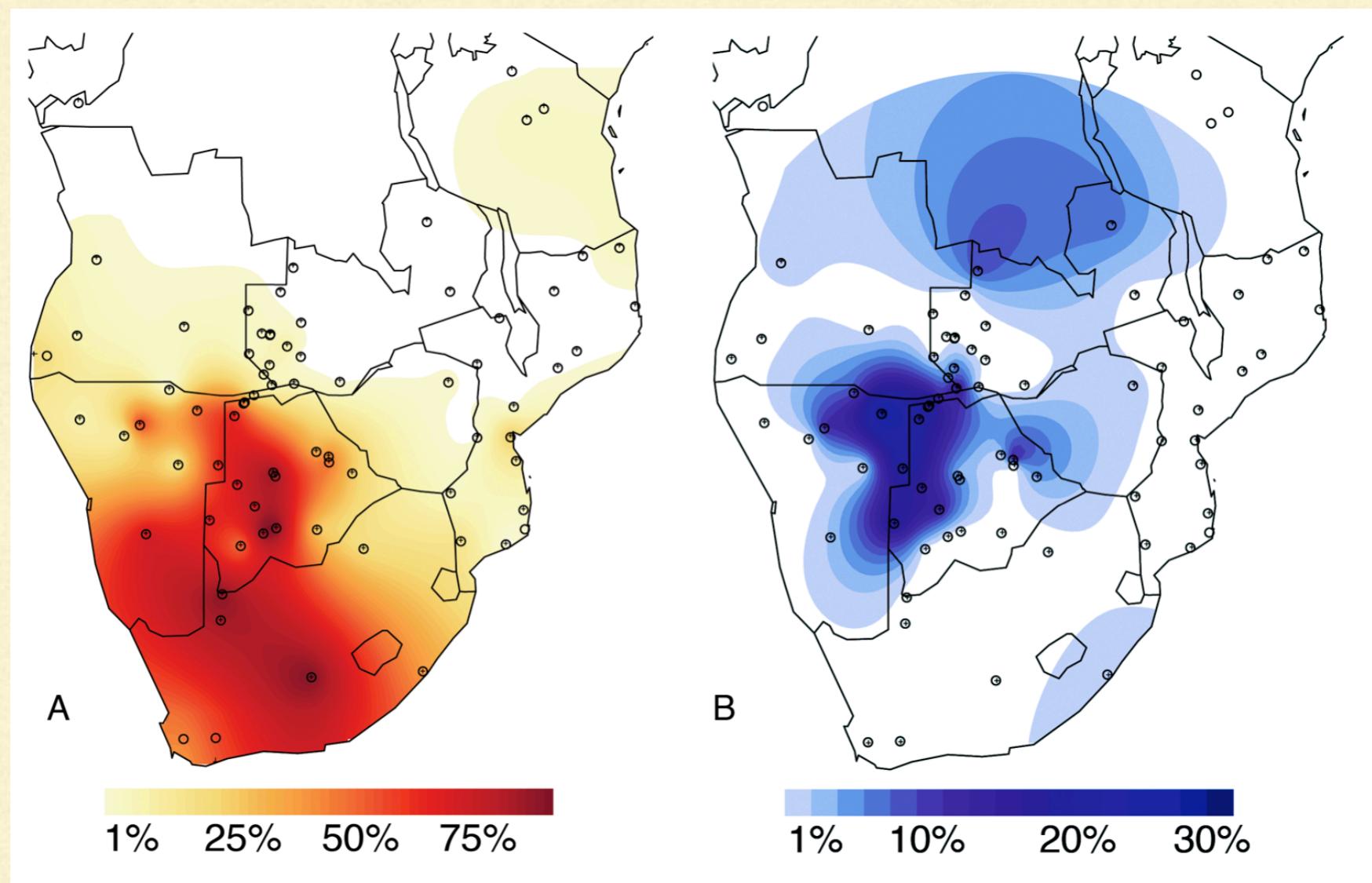
- San hunter gatherer of Southern Africa harbour the most divergent genetic lineages and early structure



Photo: C. Barbieri

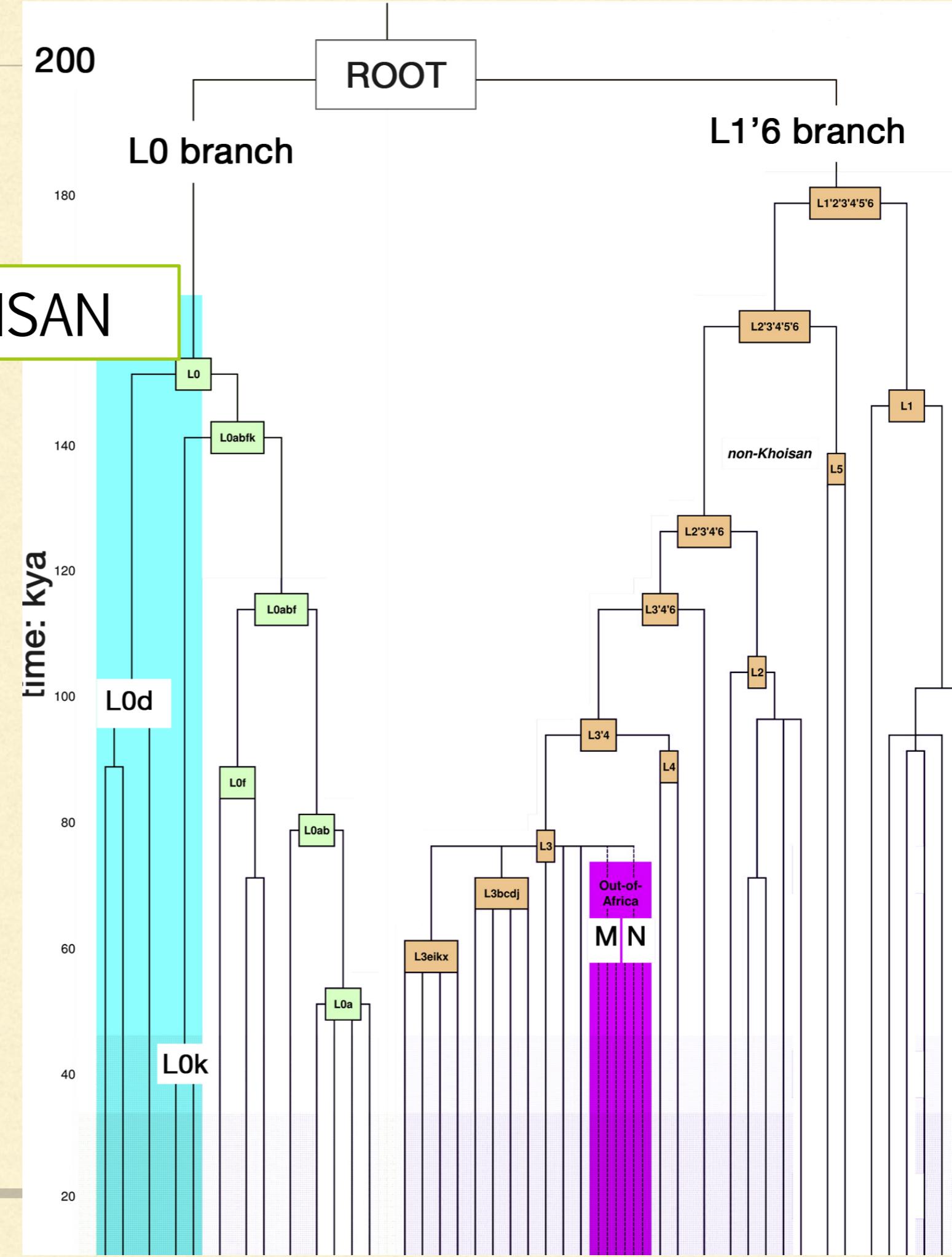


EARLY DIVERGING MTDNA LINEAGES



The first branching mtDNA lineages are found in southern Africa

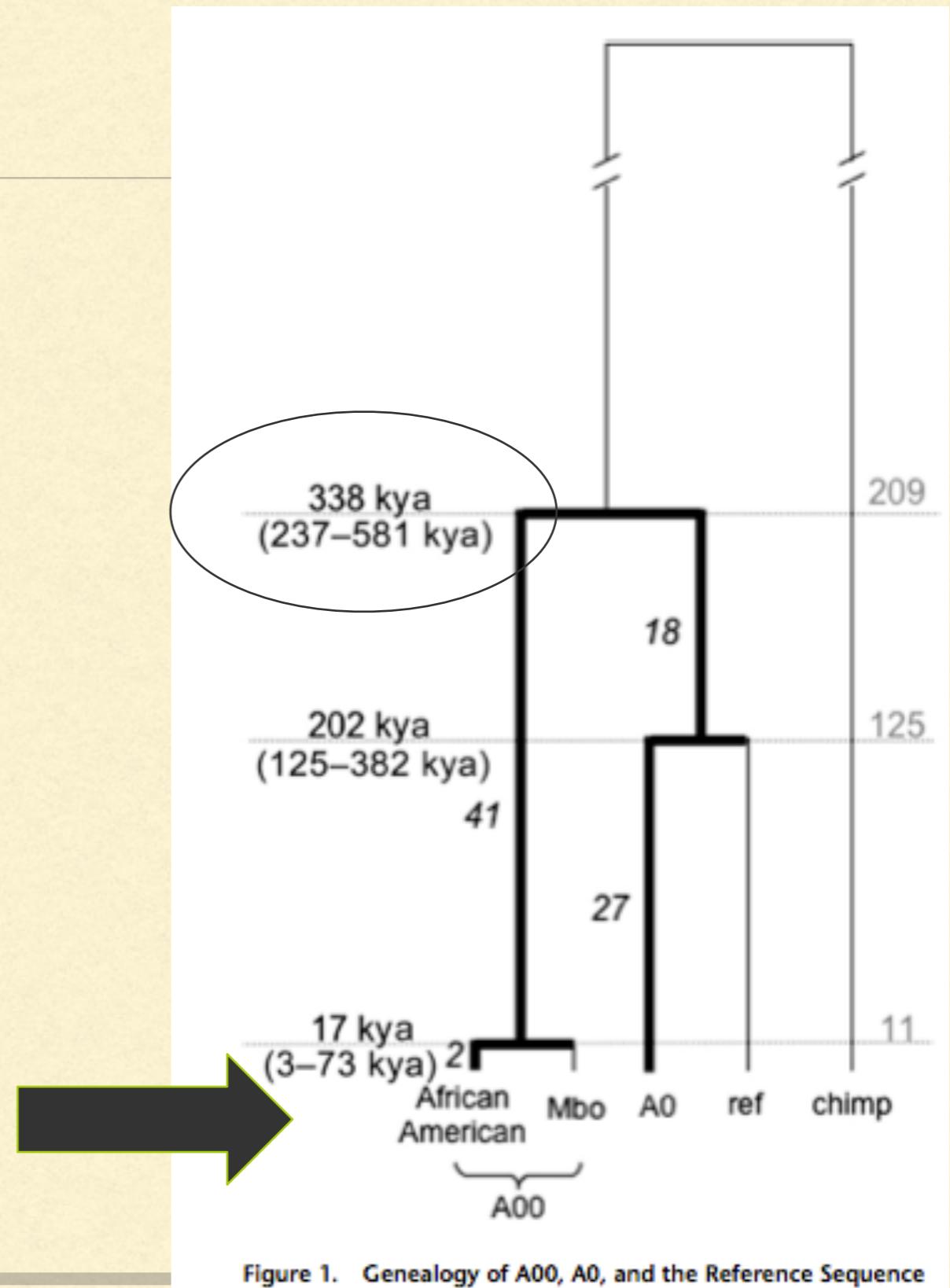
mtDNA tree

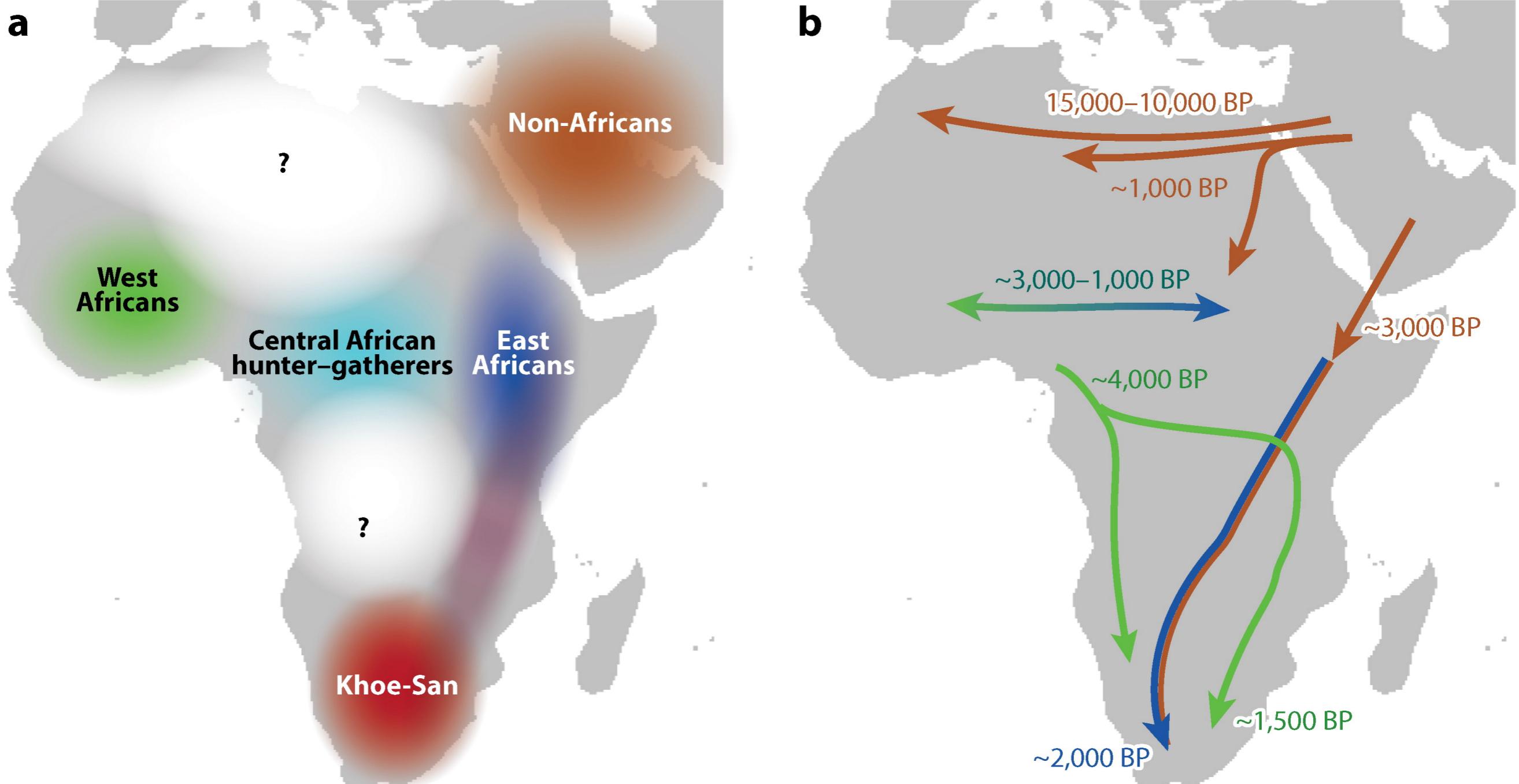


Modified from
Behar et al. 2008, Am J Hum Genet

EARLY DIVERGING Y CHROMOSOME

But the earliest branch of
the Y chromosome
phylogeny is found in
Cameroon





Schlebusch CM, Jakobsson M. 2018.
Annu. Rev. Genom. Hum. Genet. 19:405–28

PHYLOGENY OF THE LACTASE ALLELE IN SOUTHERN AFRICA

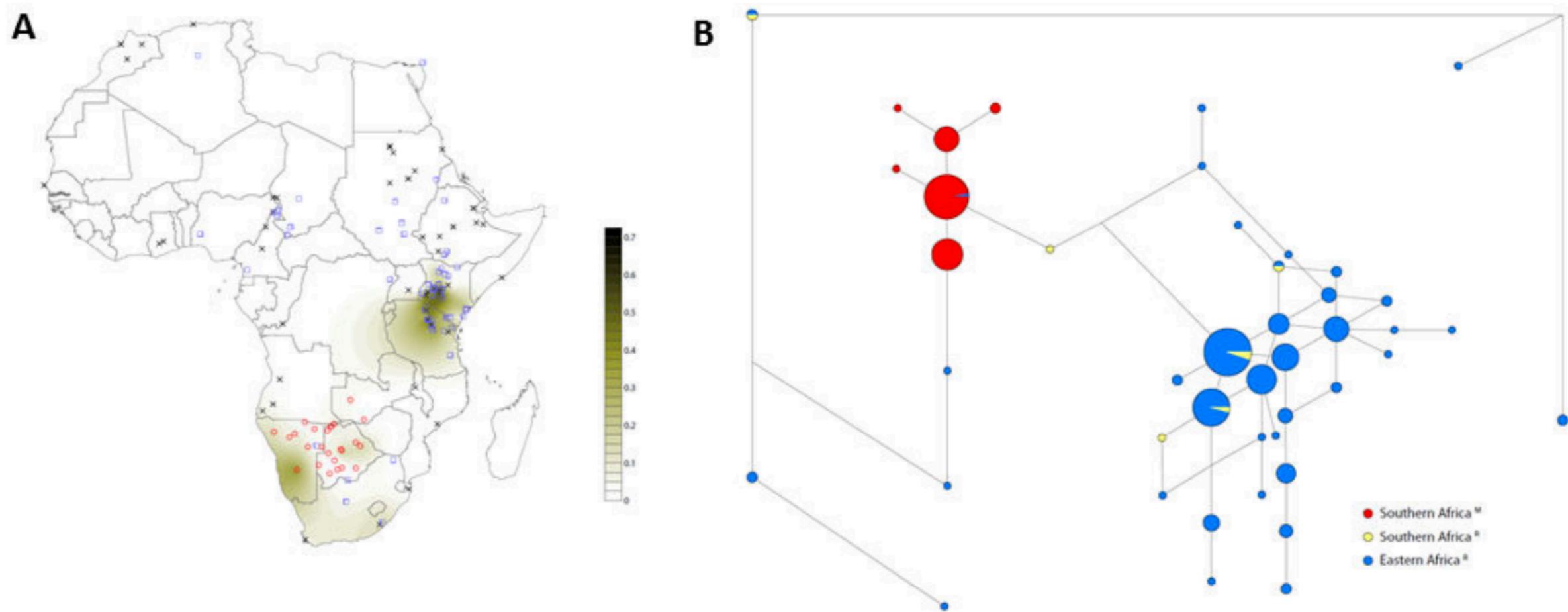
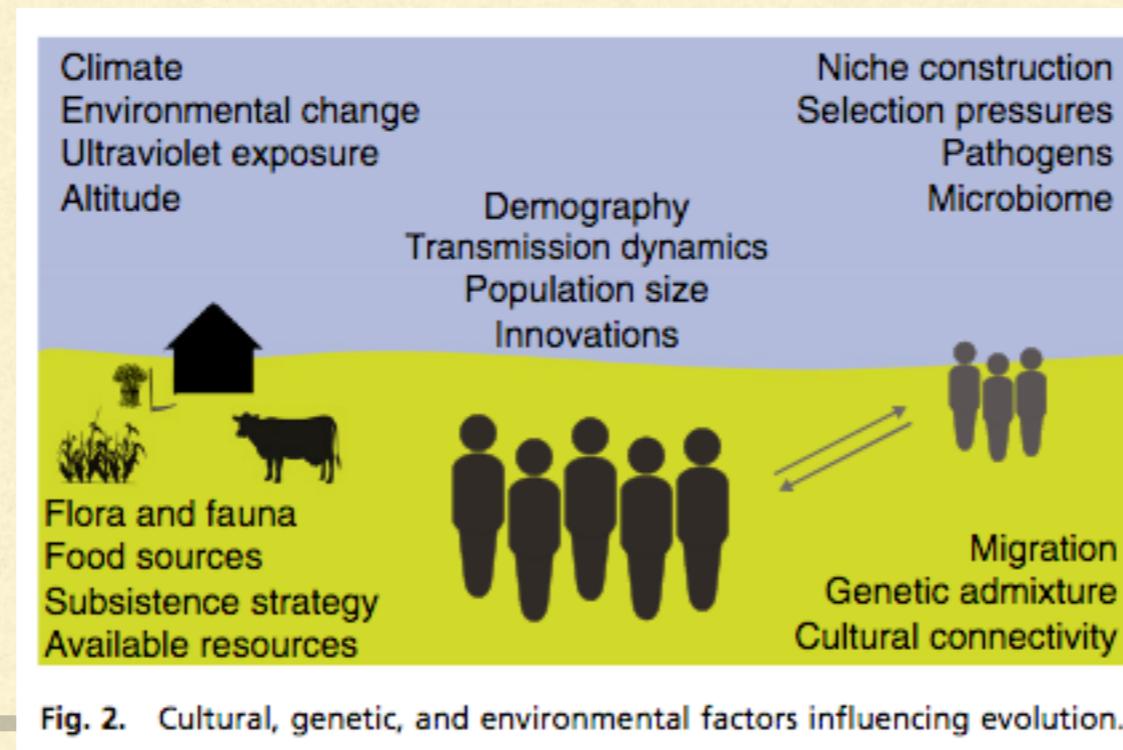


Fig. 1. Analyses of the C-14010 LP variant and associated STR haplotypes in Eastern and Southern African populations. **A:** Surfer map of the C-14010 allele frequency. Red circles denote sampling locations by Macholdt et al.; blue squares denote sampling locations by Ranciaro et al.; black crosses denote data from other published studies (taken from Macholdt et al. and Ranciaro et al.). **B:** Median-joining network of haplotypes associated with the C-14010 variant, based on four STR loci that flank the LP enhancer region. The M superscript denotes data by Macholdt et al., and the R superscript denotes data by Ranciaro et al. [Color]

CULTURE: TRANSMISSION WITH MODIFICATION

- Human cultural traits (behaviors, ideas, and technologies that can be learned from other individuals) can exhibit complex patterns of transmission and evolution
- Cultural and genetic evolution can interact with one another and influence both transmission and selection



BIOLOGY + CULTURE = MULTIDISCIPLINARY STUDIES

- Evolving under similar demographic processes
- Addressing questions about human past and present diversity, ecology and evolution
- The whole study is greater than the sum of its discipline parts



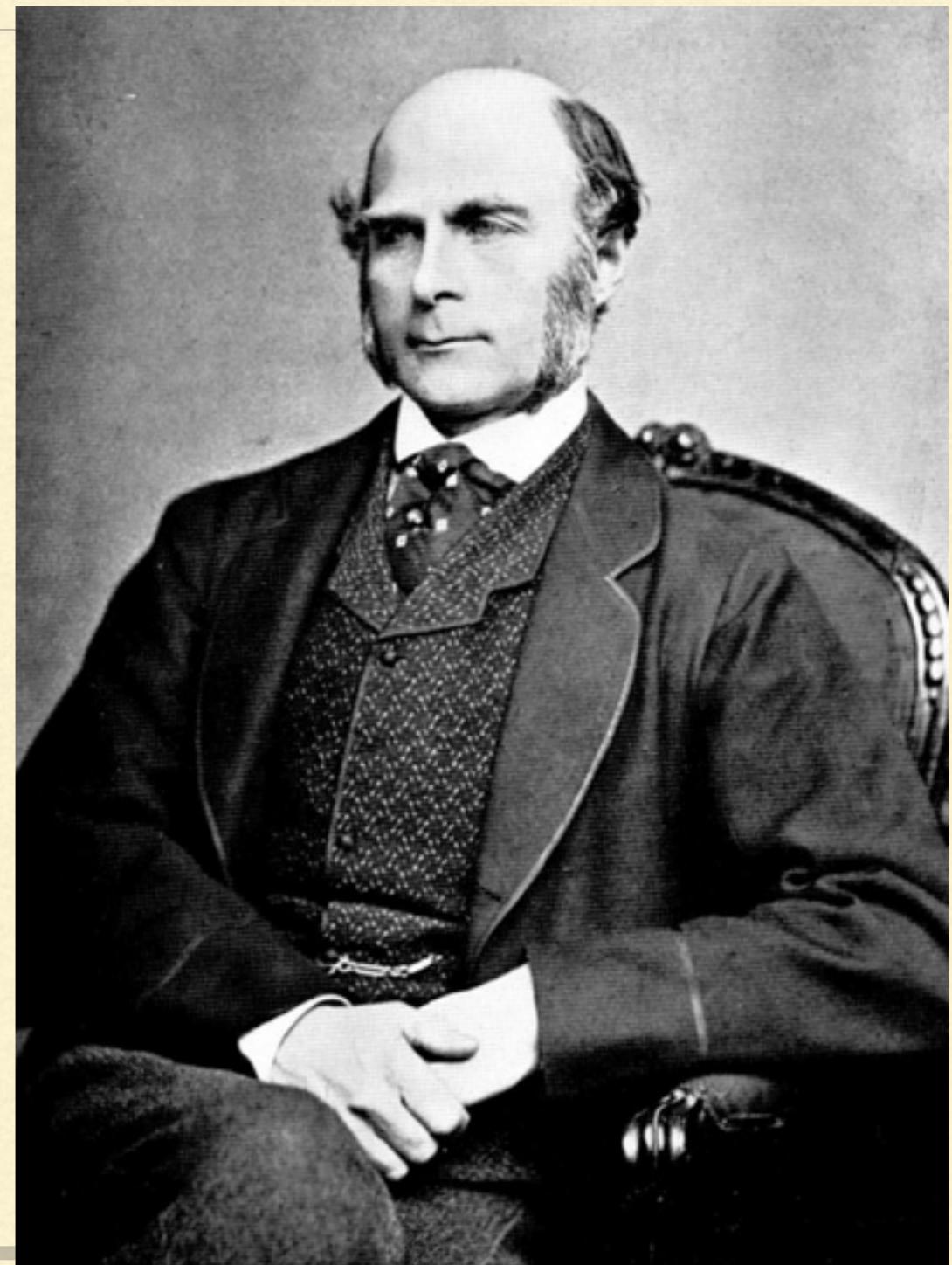


Image: Minna Sundberg

EVOLUTIONARY PARALLELS

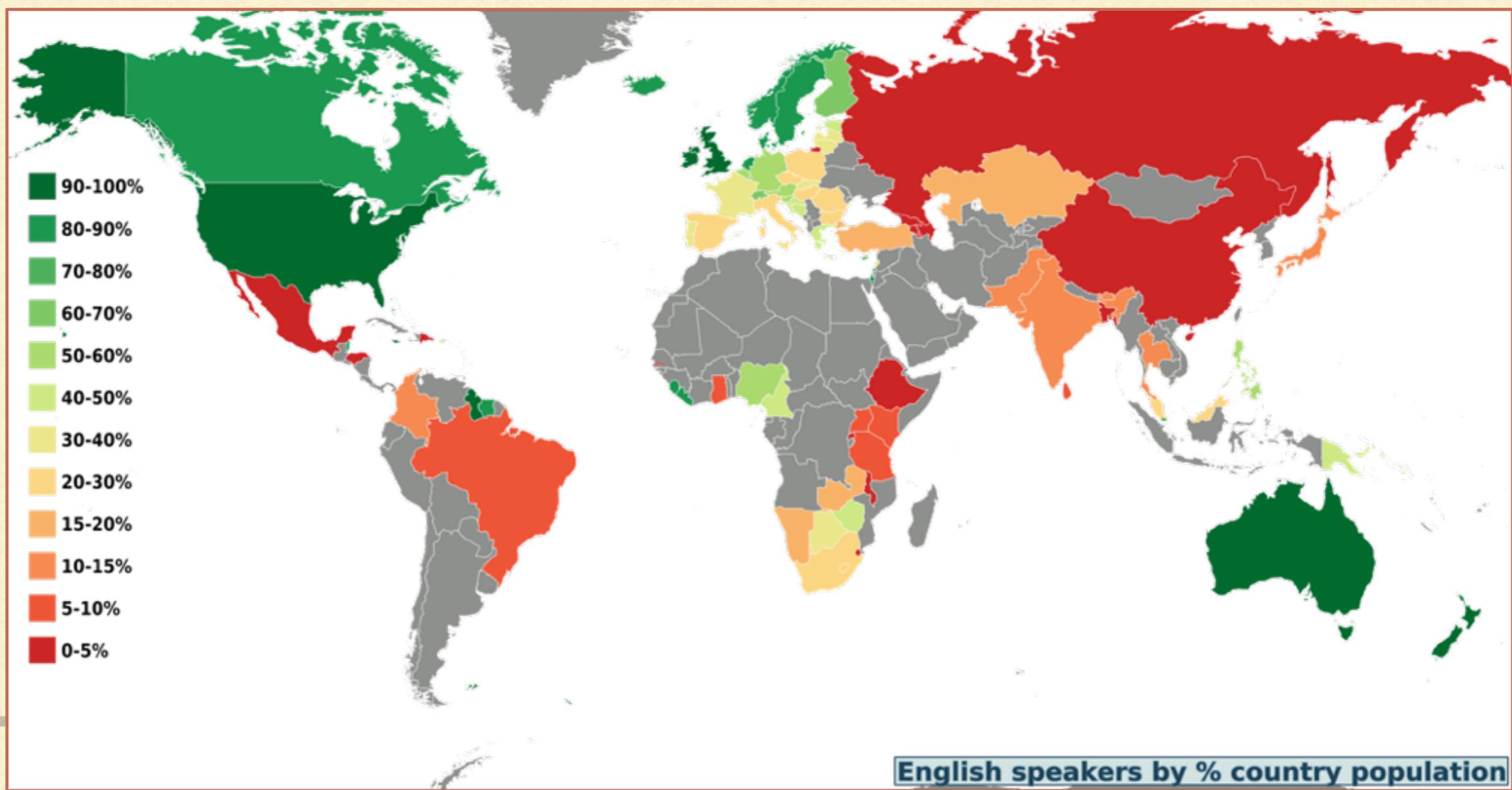
TREES AND GALTON'S PROBLEM

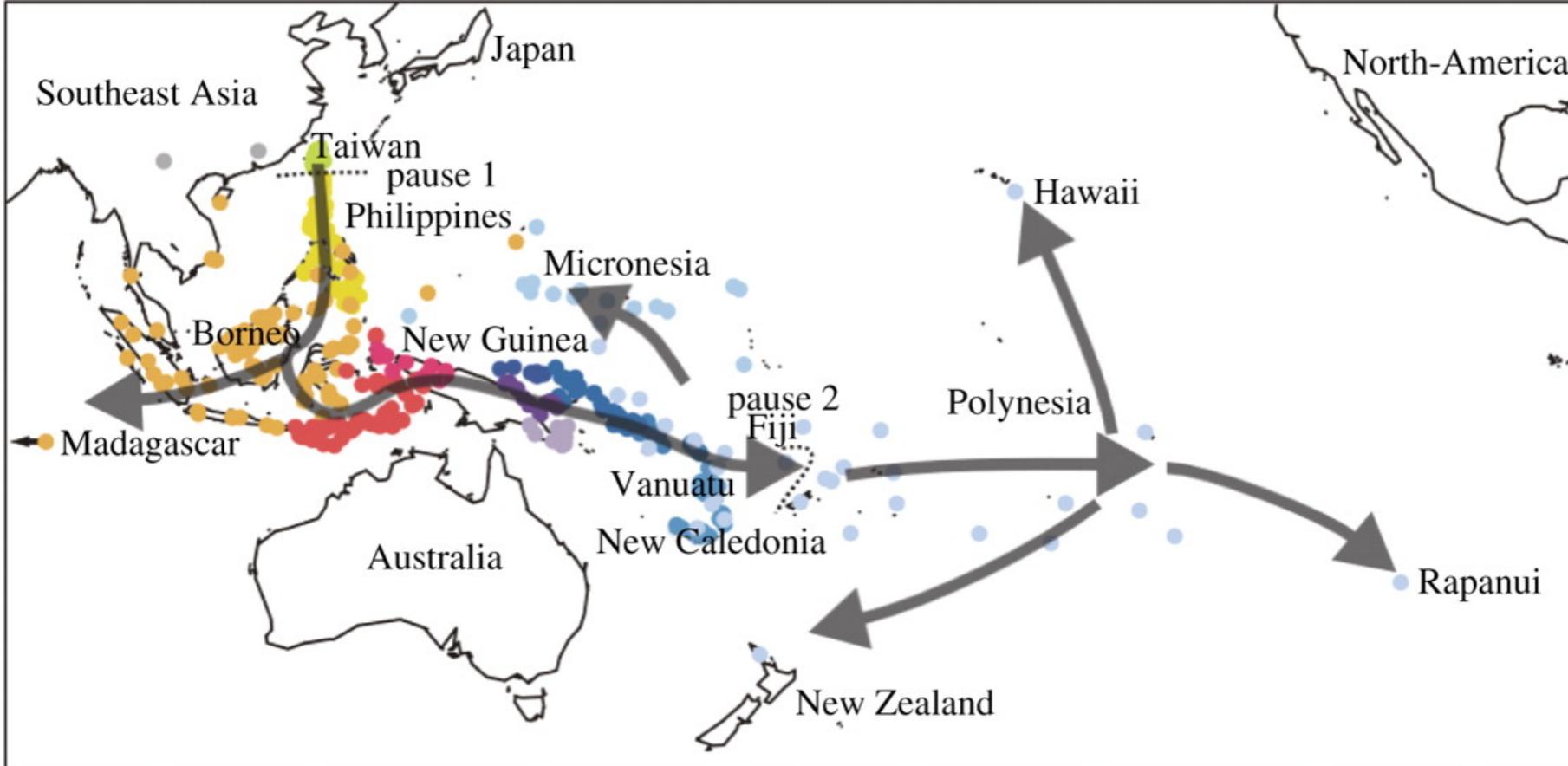
- Spurious associations:
- Relatedness by descent explains correlation between variables in two populations



LIMITATIONS OF THE LINGUISTIC AND CULTURAL COMPARISON

■ Horizontal transfer





An elegant model of language spread: The case of the Austronesian (Gray et al. 2009)

