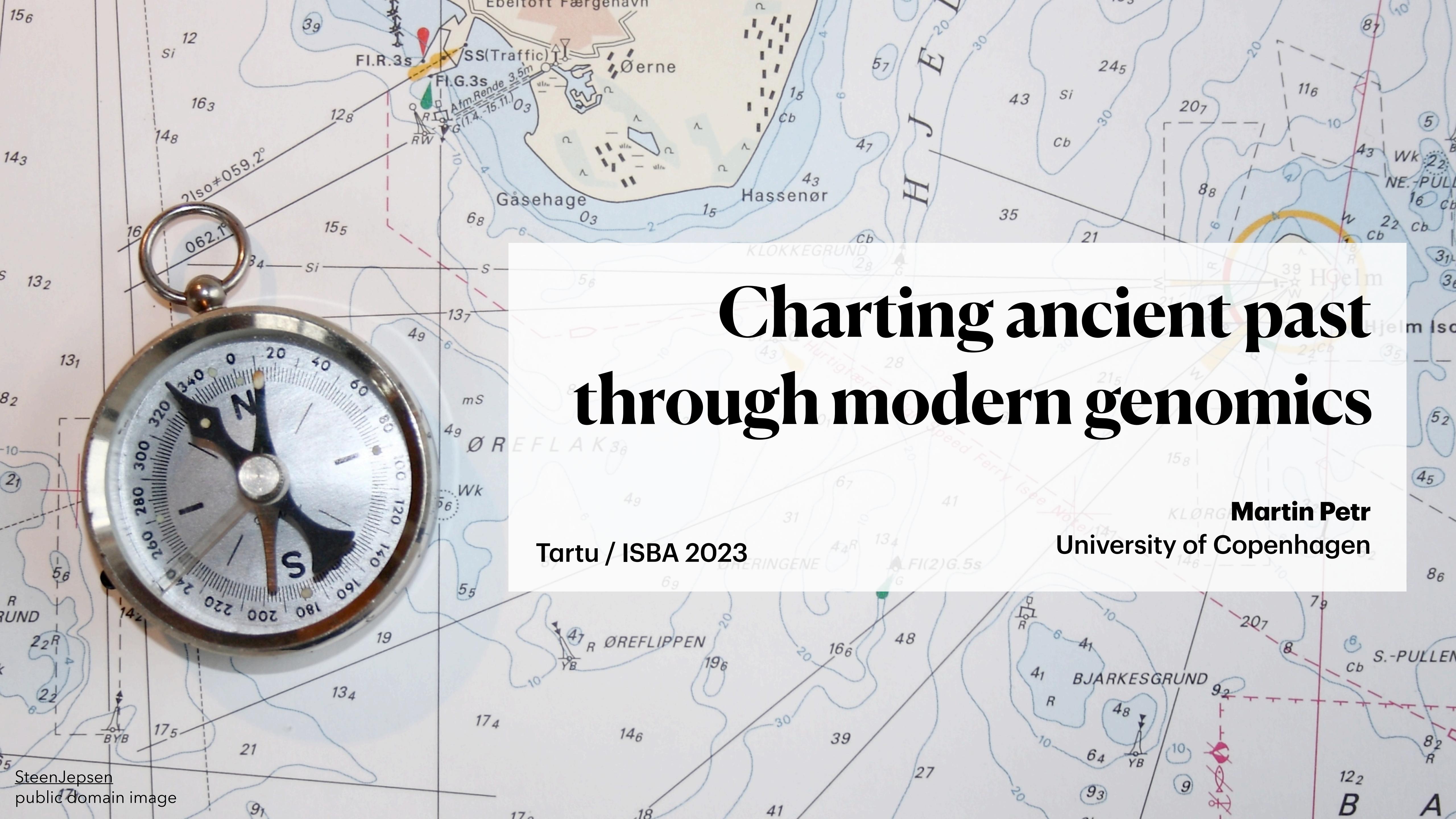


Charting ancient past through modern genomics

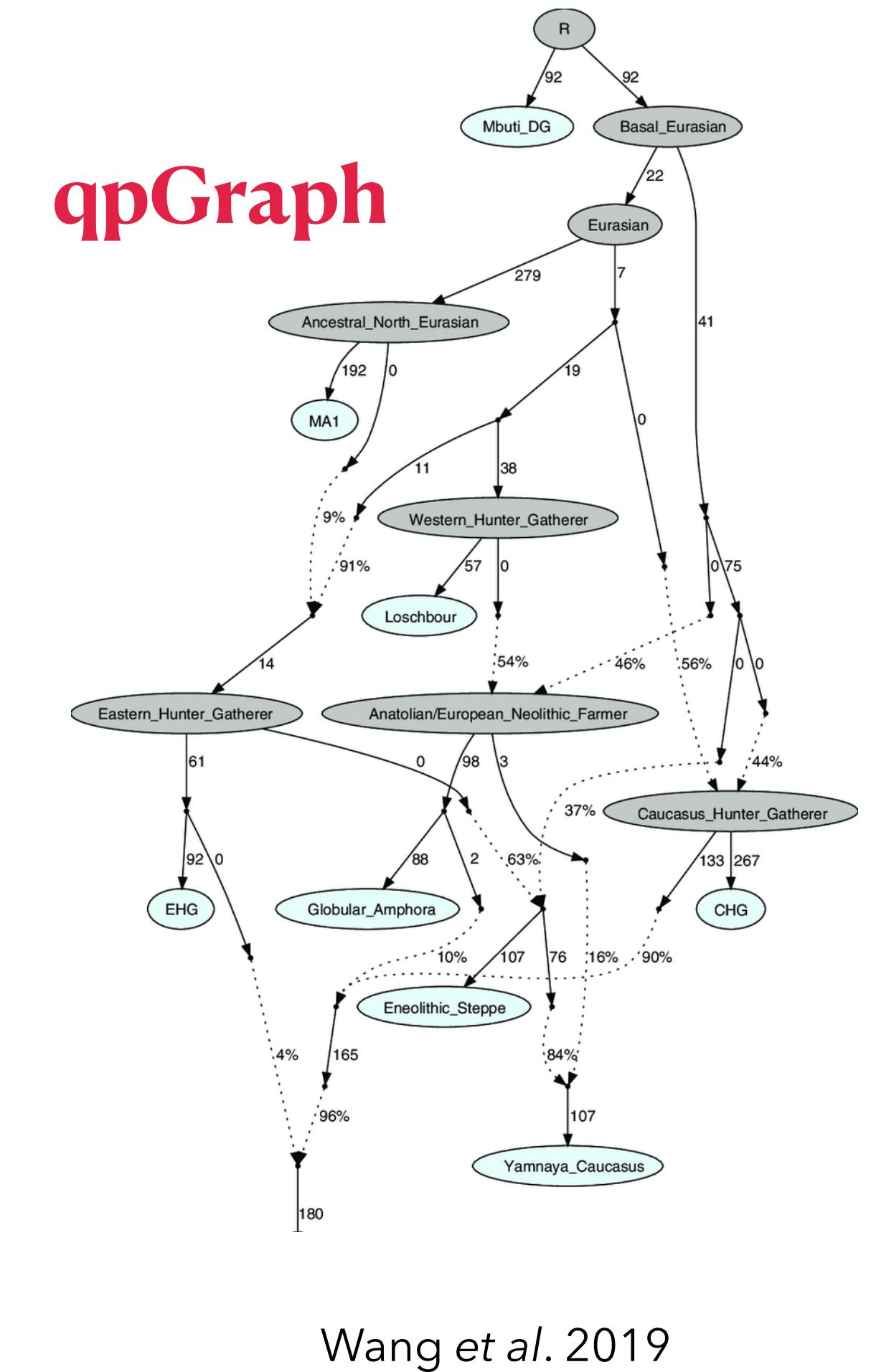
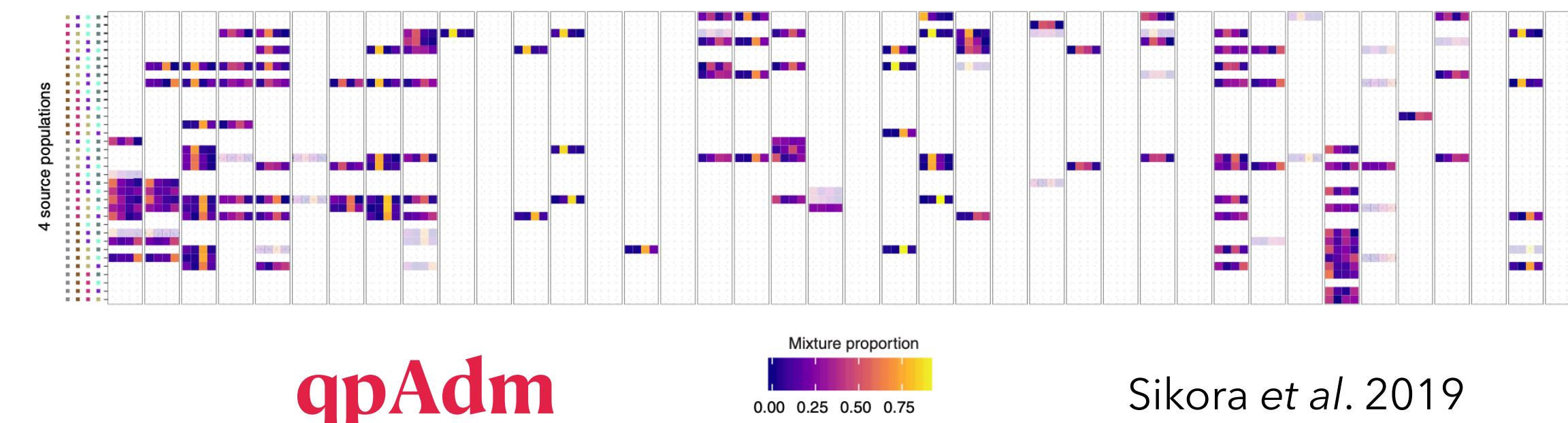
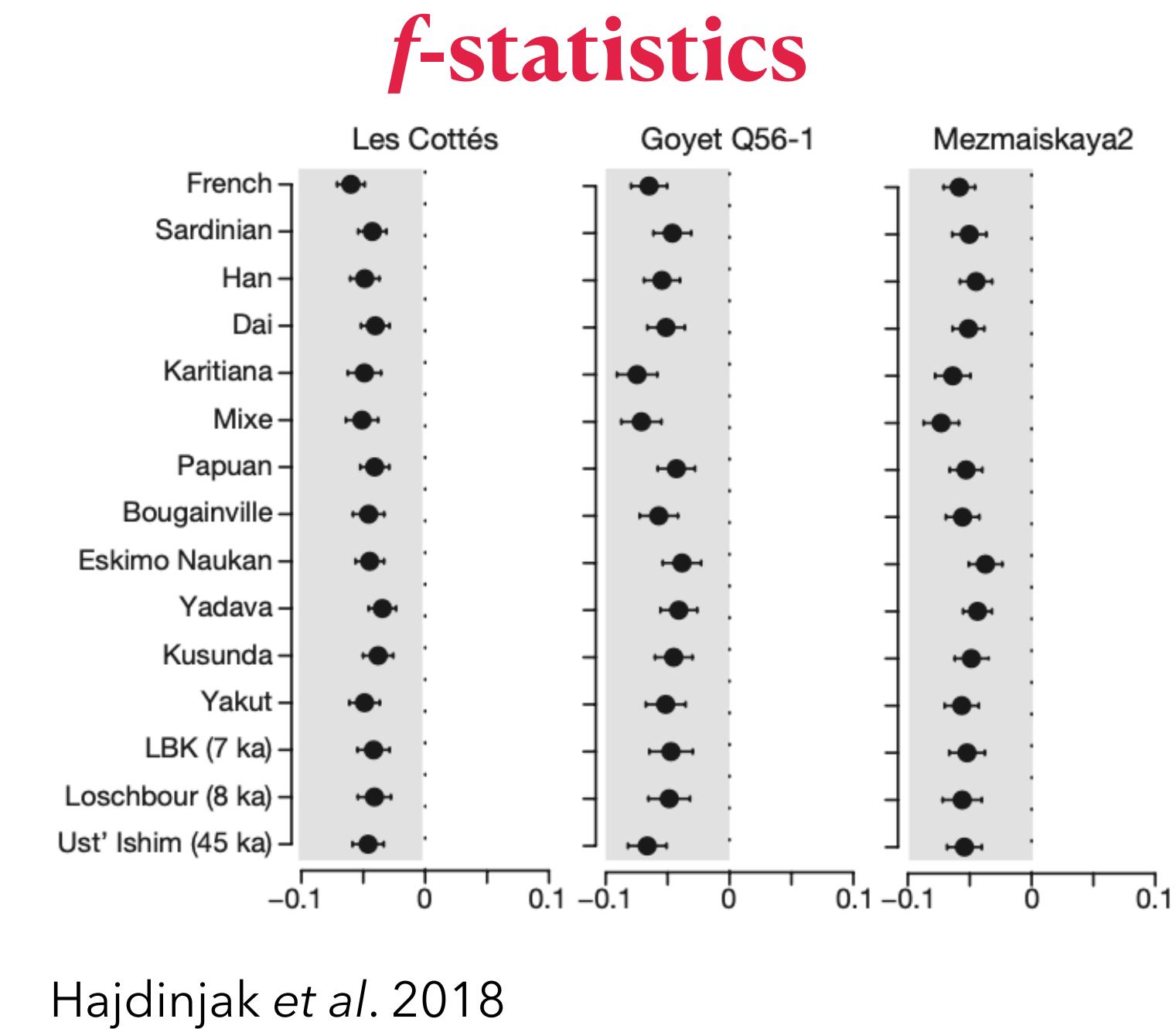
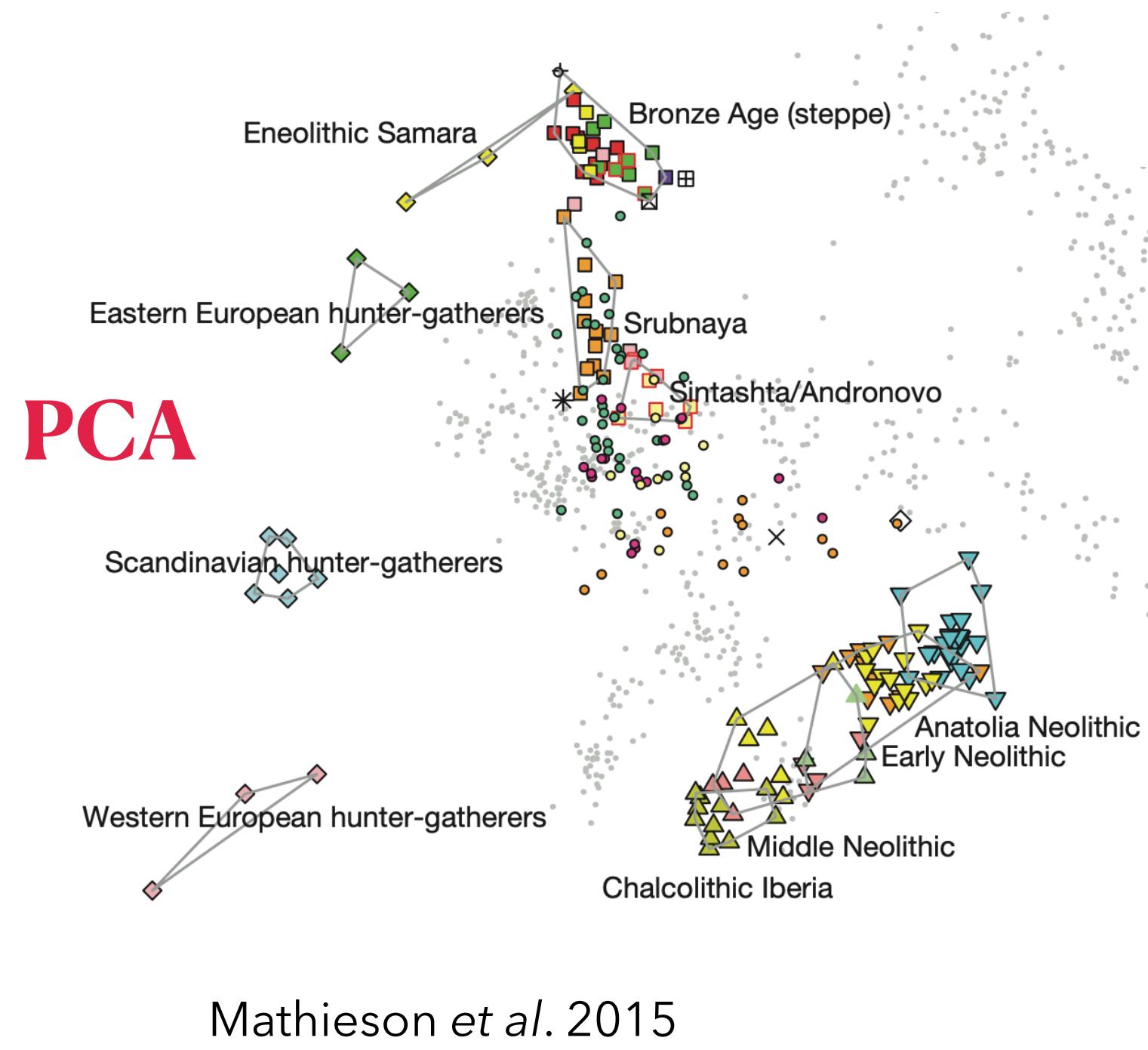
Tartu / ISBA 2023

Martin Petr

University of Copenhagen

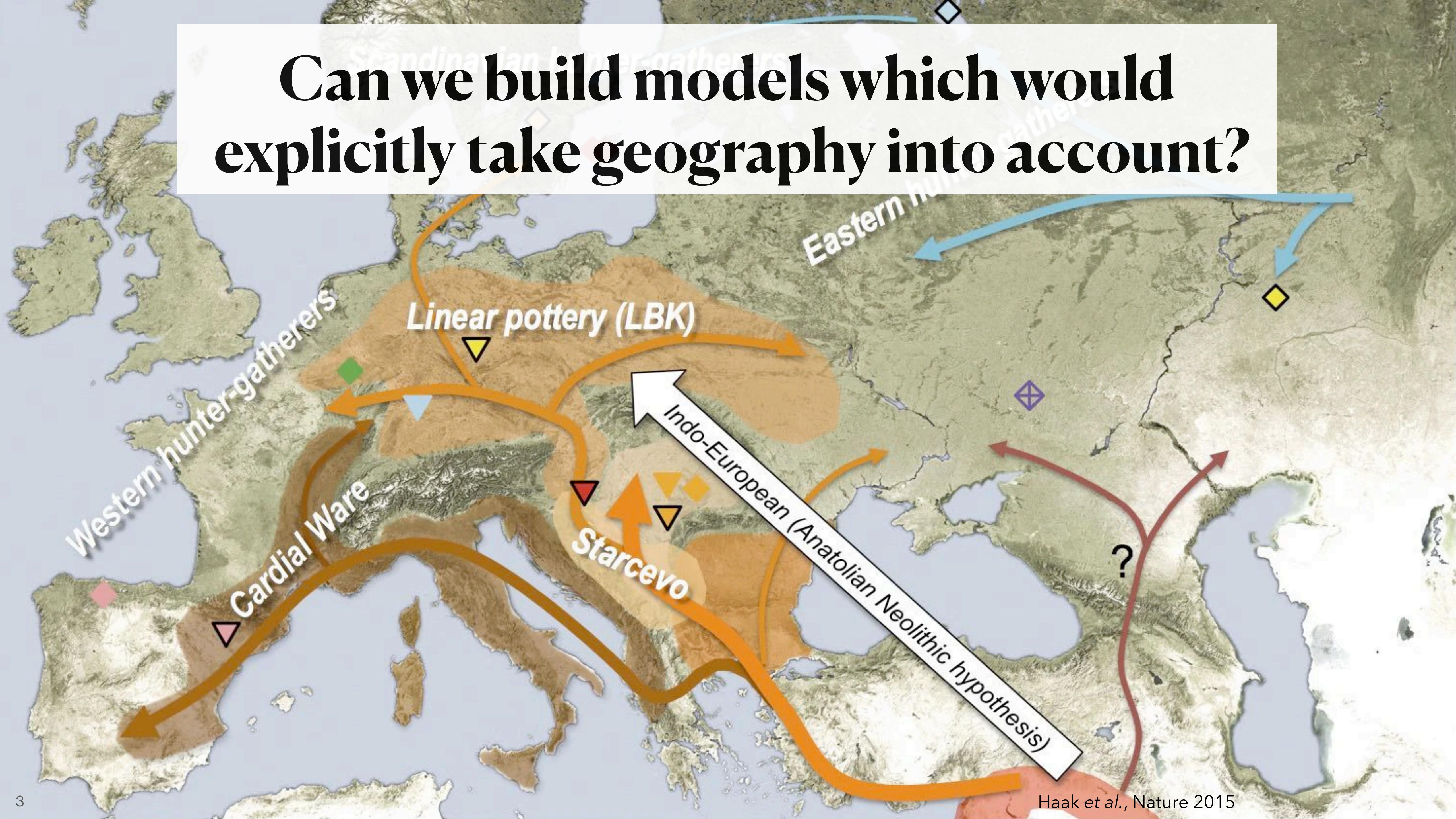


The all-stars of ancient population genetics





Can we build models which would explicitly take geography into account?



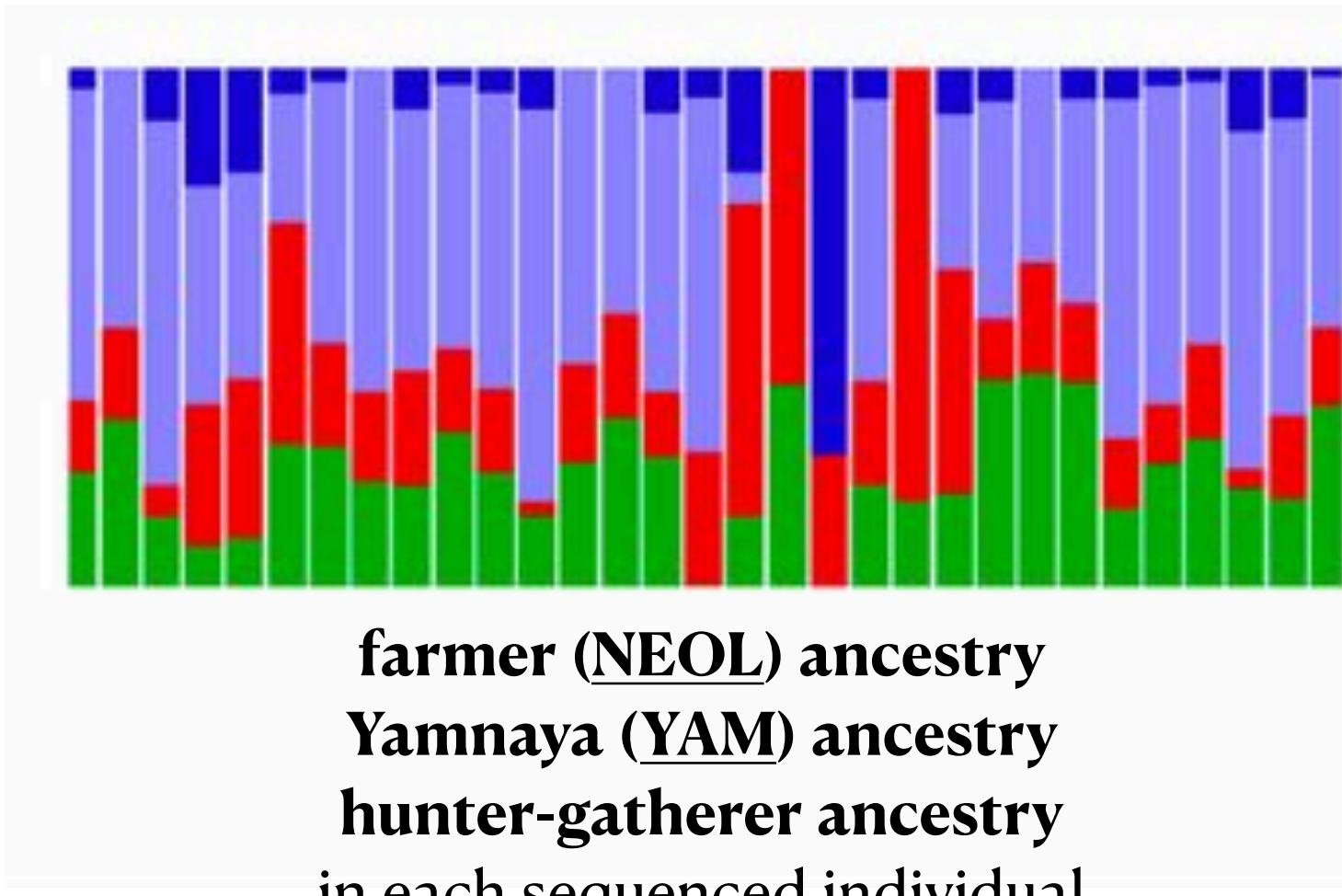
The spatiotemporal spread of human migrations during the European Holocene

Fernando Racimo^{a,1} , Jessie Woodbridge^b , Ralph M. Fyfe^b , Martin Sikora^a, Karl-Göran Sjögren^c , Kristian Kristiansen^c, and Marc Vander Linden^d 

^aLundbeck GeoGenetics Centre, The Globe Institute, University of Copenhagen, 1350 Copenhagen, Denmark; ^bSchool of Geography, Earth, and Environmental Sciences, University of Plymouth, Plymouth PL4 8AA, United Kingdom; ^cDepartment of Historical Studies, University of Gothenburg, 405 30 Gothenburg, Sweden; and ^dDepartment of Archaeology, University of Cambridge, Cambridge CB2 1TN, United Kingdom

Spatio-temporal mapping of genetic ancestry

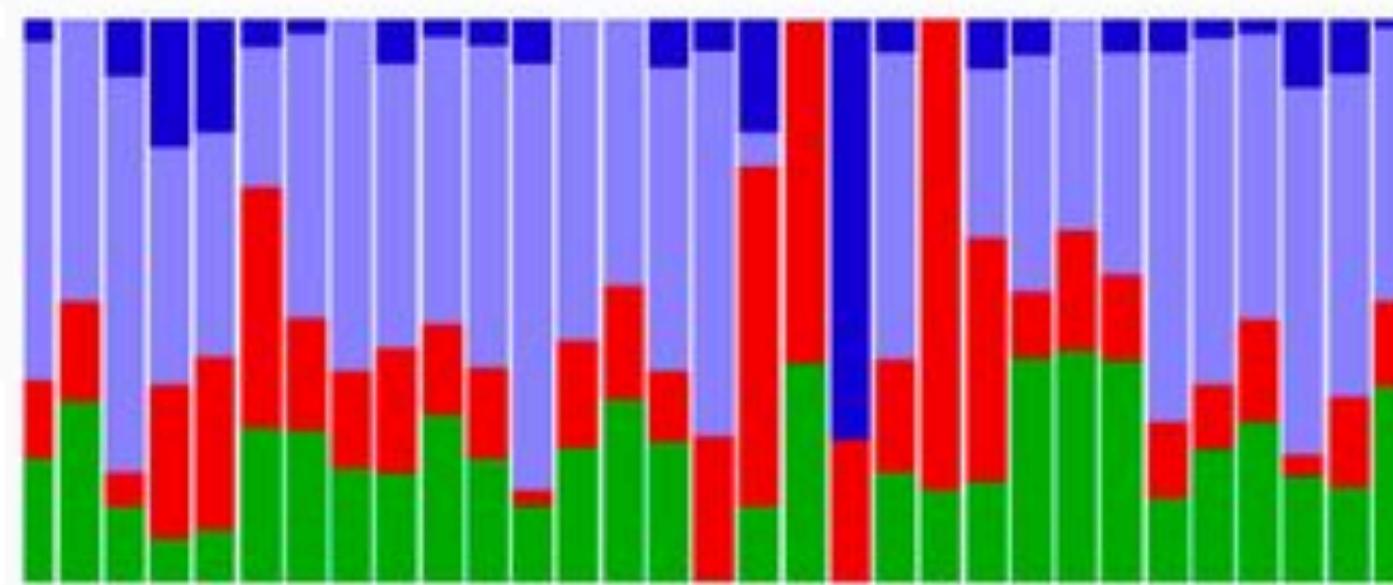
ADMIXTURE analysis



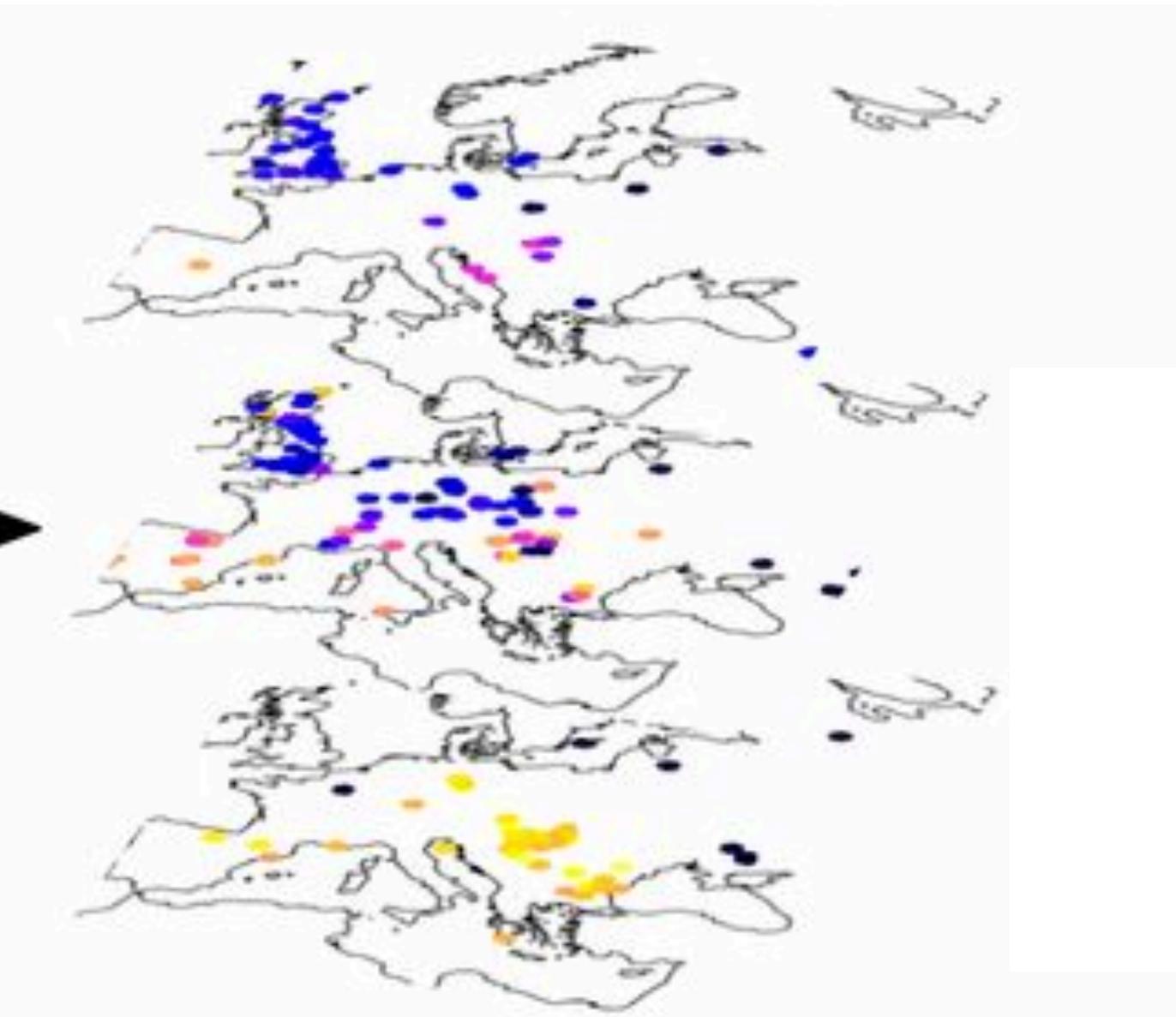
Spatio-temporal mapping of genetic ancestry

position of each individual
(and their ancestry) on a map

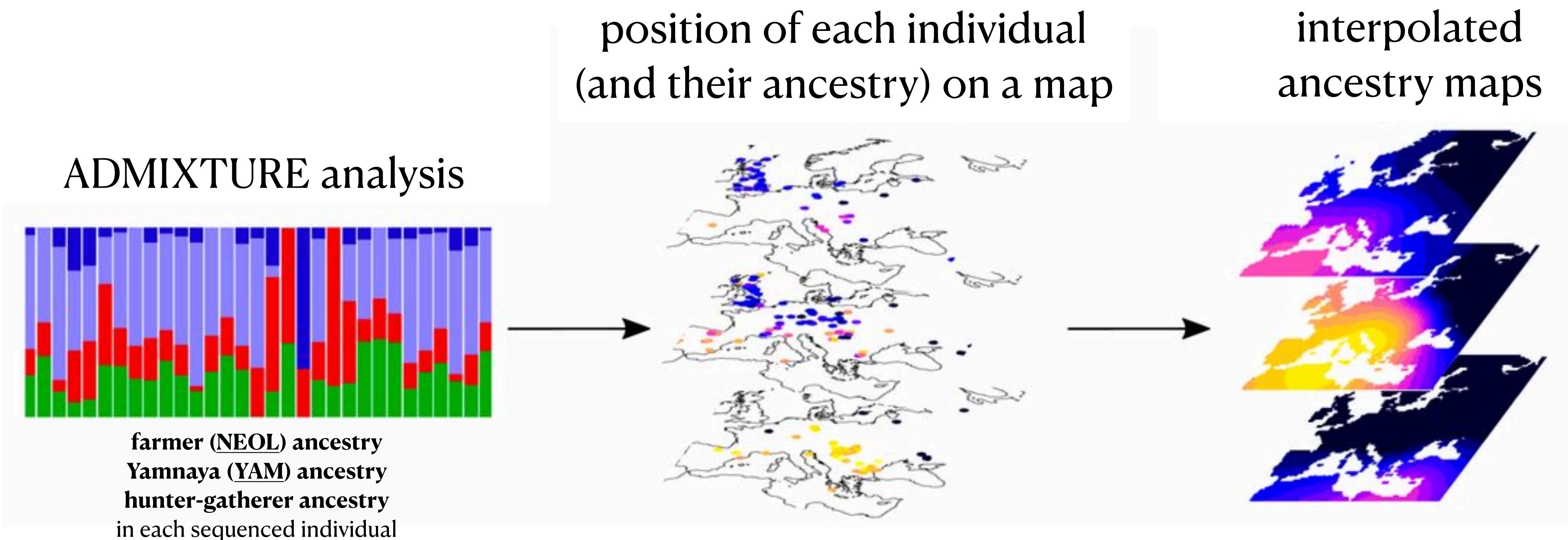
ADMIXTURE analysis



farmer (NEOL) ancestry
Yamnaya (YAM) ancestry
hunter-gatherer ancestry
in each sequenced individual



Spatio-temporal mapping of genetic ancestry



Interpolated ancestry maps

"NEOL" ancestry

-10800



"YAM" ancestry

-10800



Interpolated ancestry maps

"NEOL" ancestry

-10800



"YAM" ancestry

-10800



How fast did an "ancestry wave" move from its origin?

"NEOL" ancestry

-10800 BP



"YAM" ancestry

-10800 BP



How fast did an "ancestry wave" move from its origin?

"NEOL" ancestry

-10800 BP



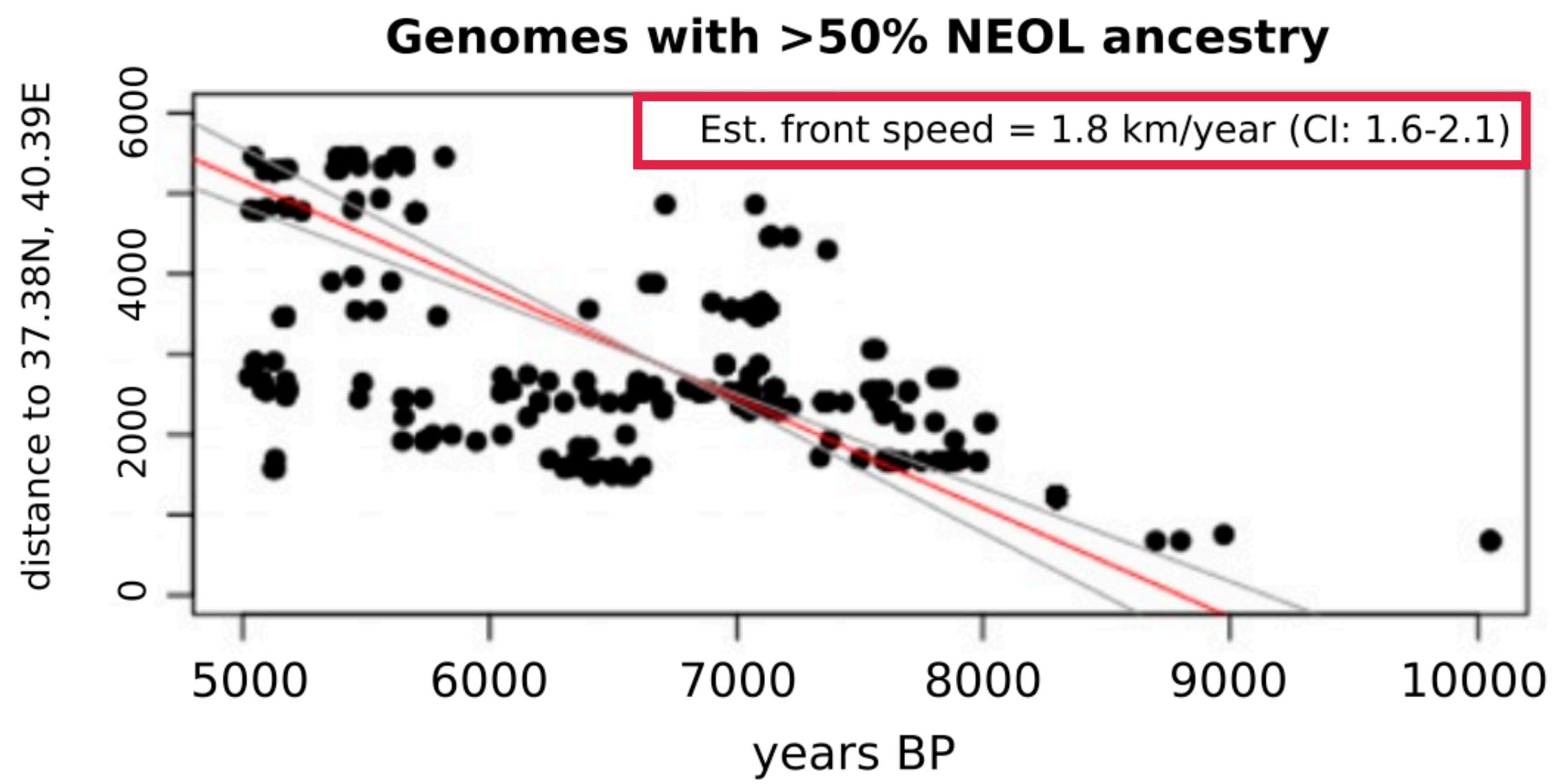
"YAM" ancestry

-10800 BP

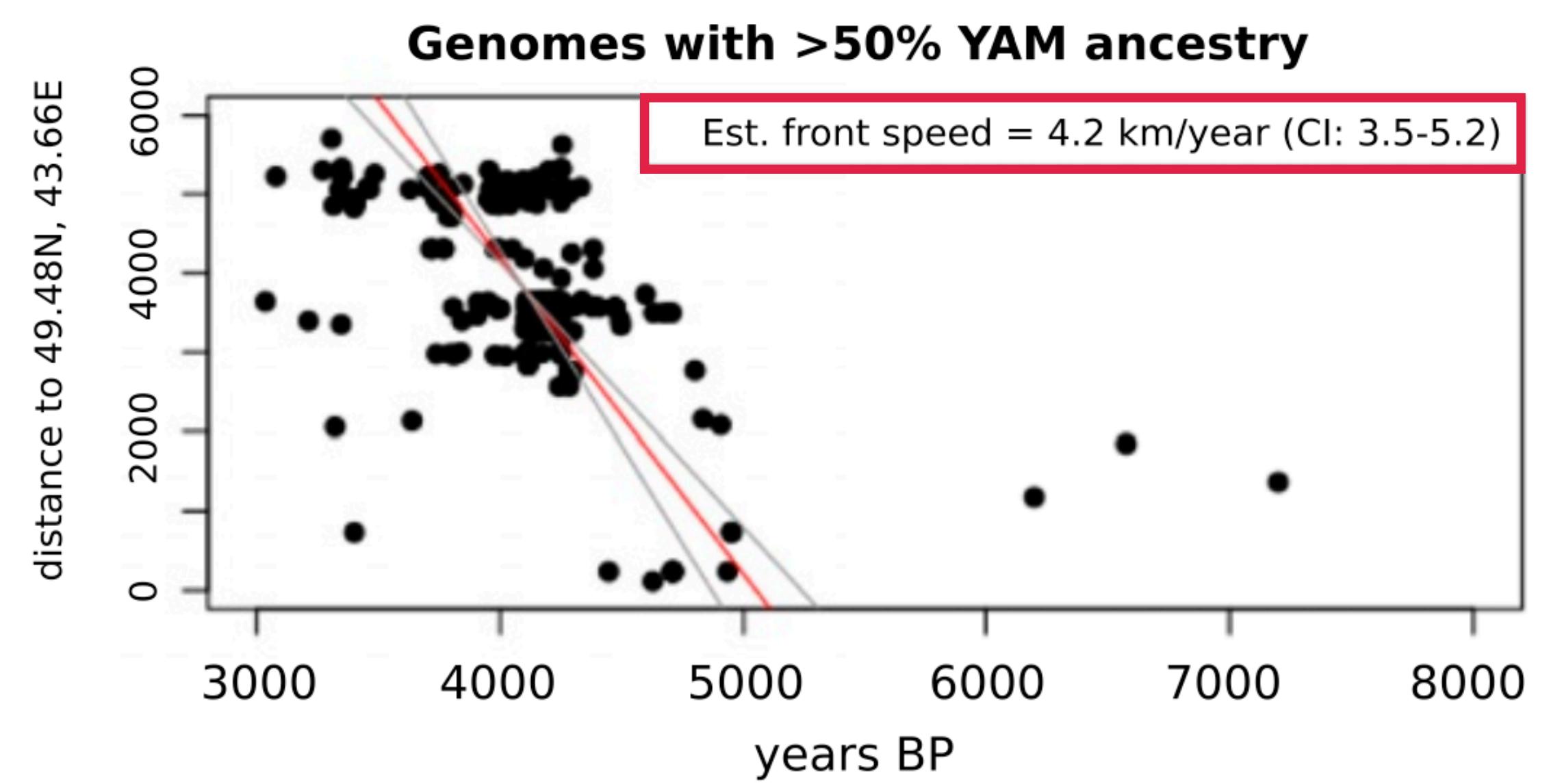


How fast did an "ancestry wave" move from its origin?

"NEOL" ancestry



"YAM" ancestry



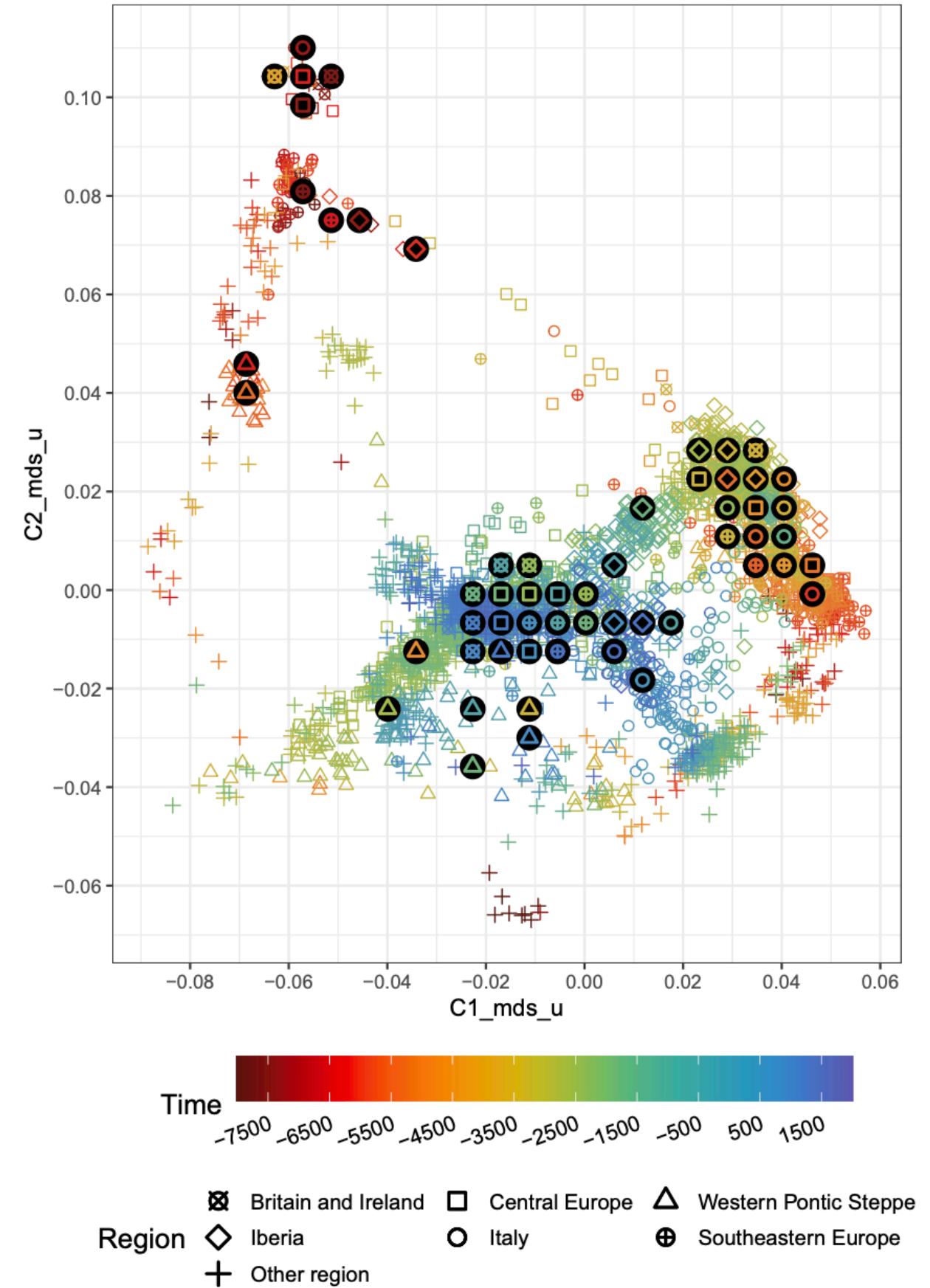
Estimating human mobility in Holocene Western Eurasia with large-scale ancient genomic data

Clemens Schmid^{a,b}  and Stephan Schiffels^{a,1} 

Edited by Liisa Loog, University of Oxford, Cambridge, UK; received November 1, 2022; accepted December 5, 2022, by Editorial Board Member Richard G. Klein

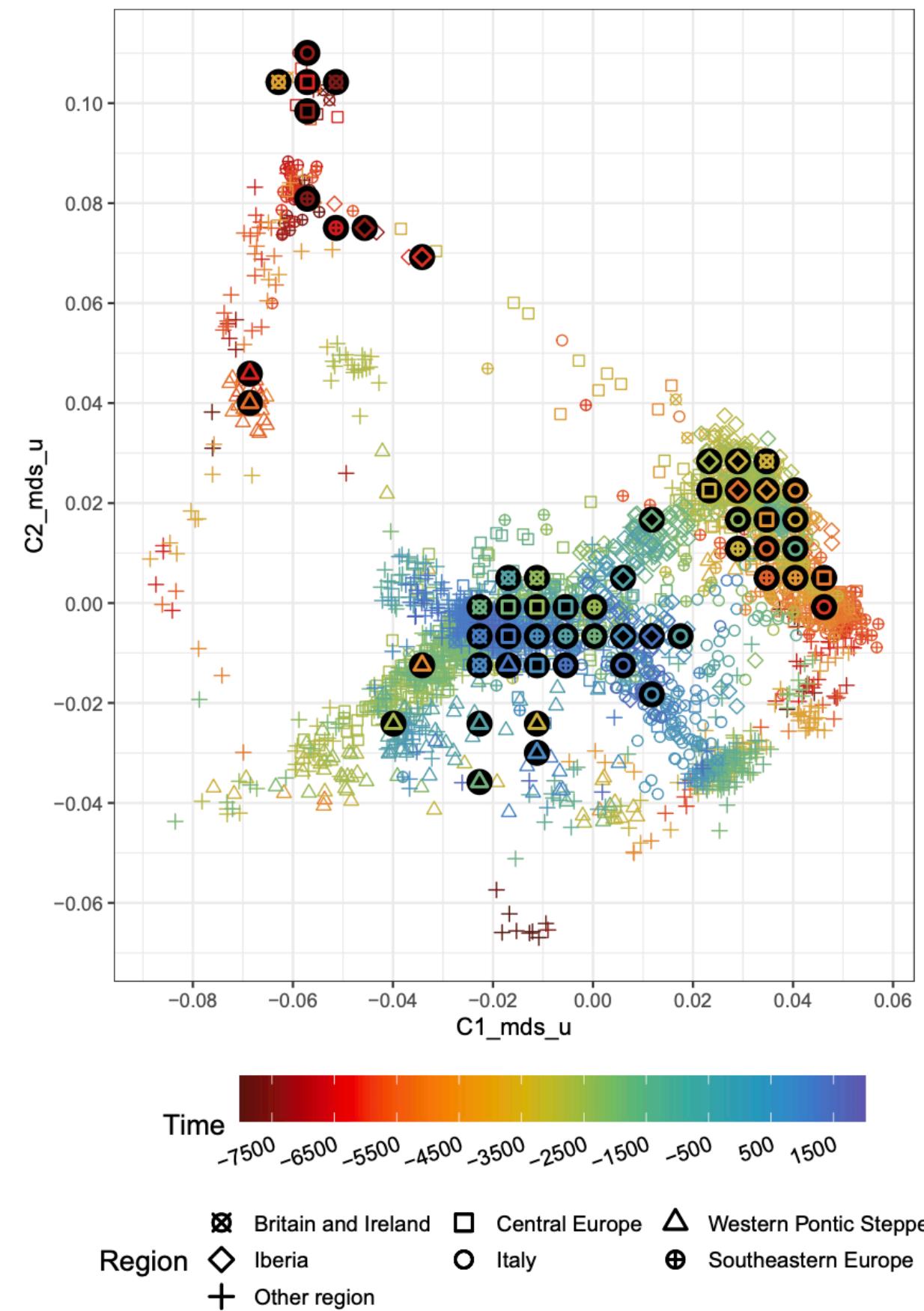
Where does a sample trace its ancestral origin?

MDS / PCA

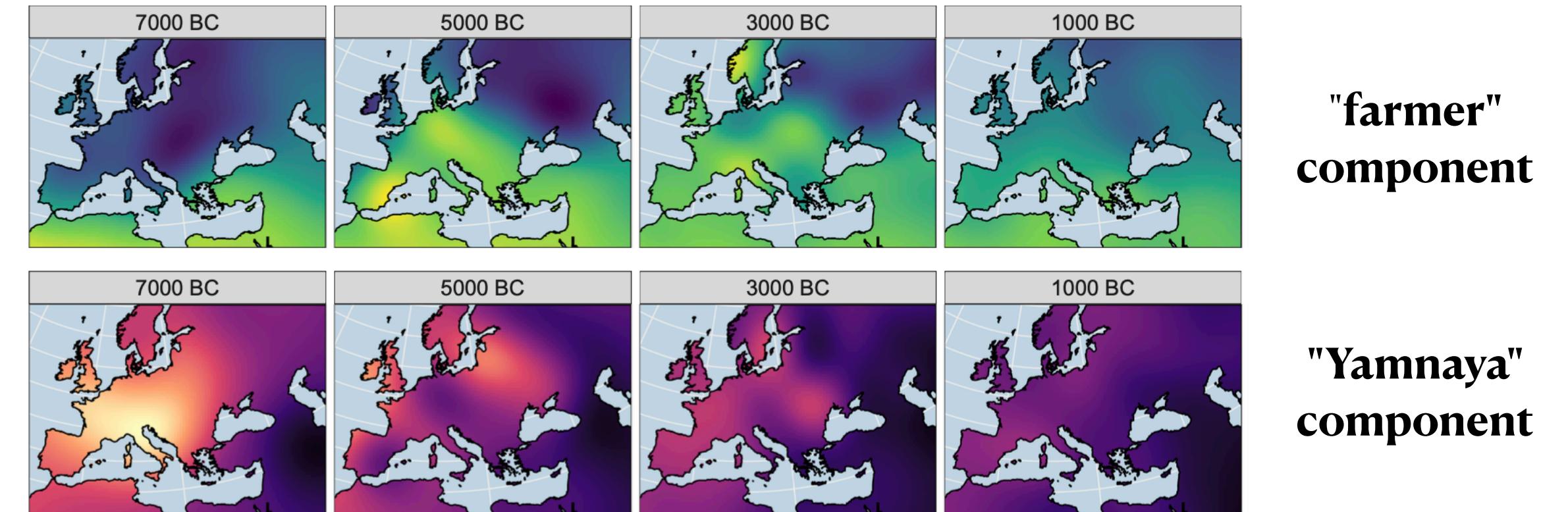


Where does a sample trace its ancestral origin?

MDS / PCA

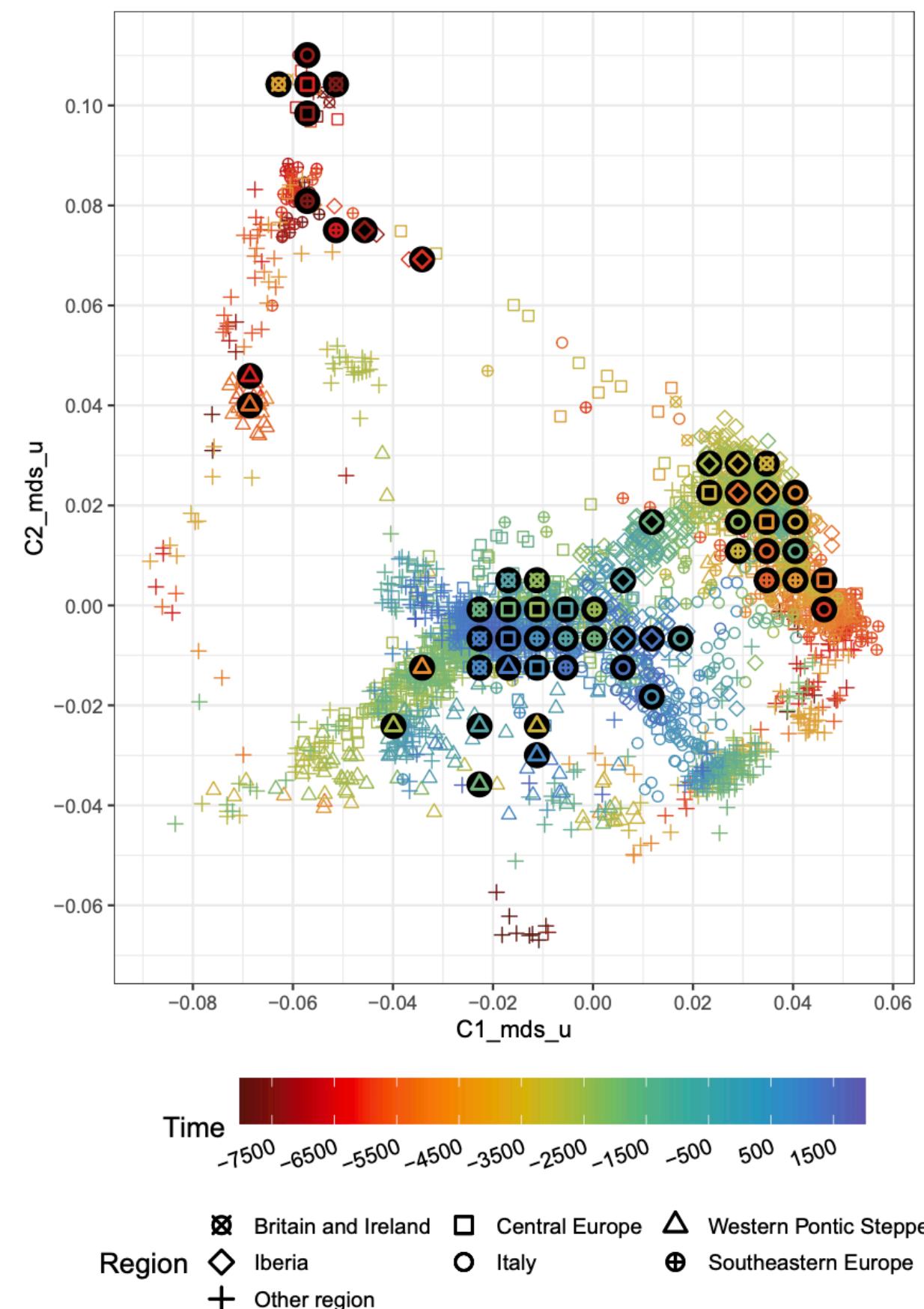


spatio-temporal interpolated "PCA field"

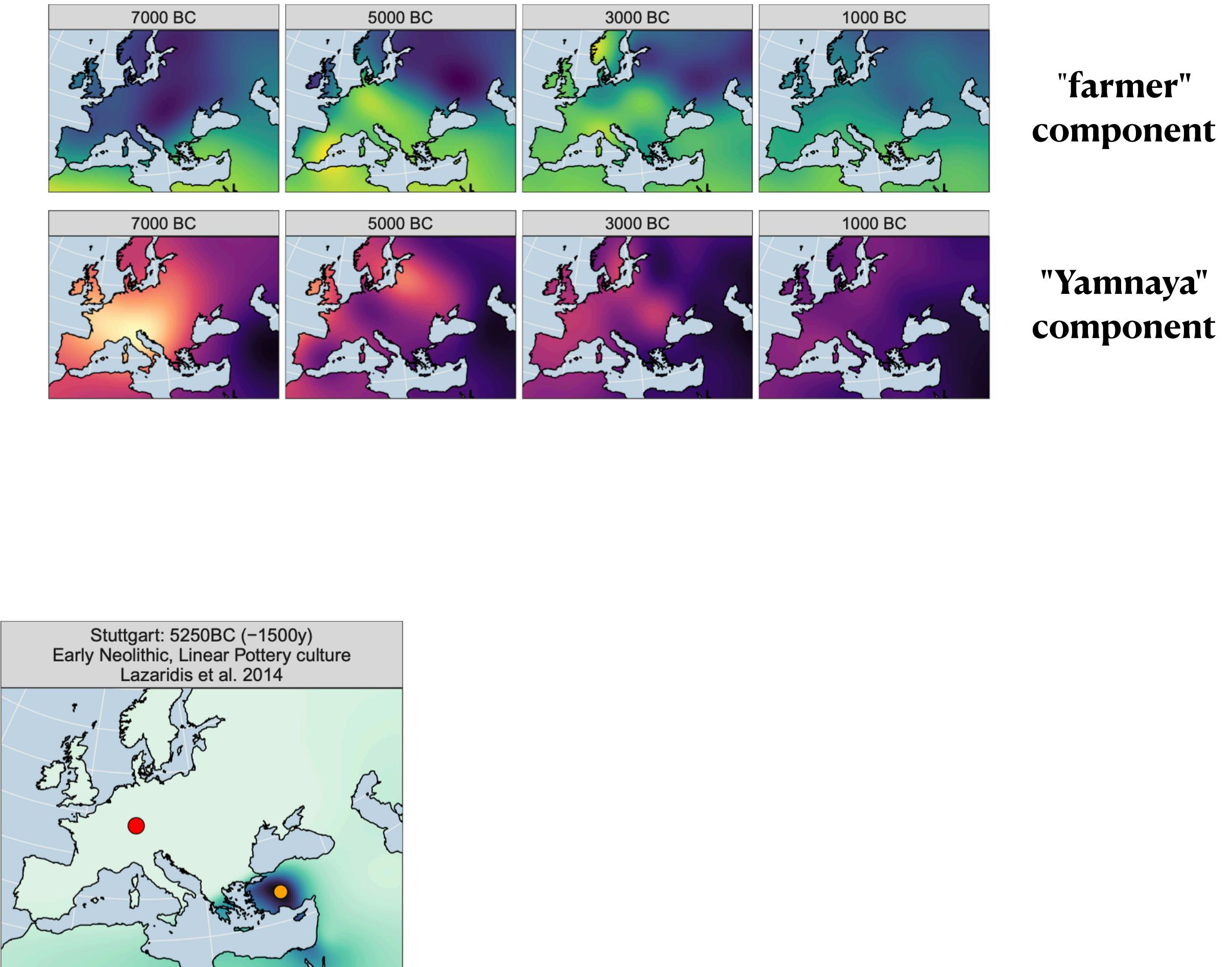


Where does a sample trace its ancestral origin?

MDS / PCA

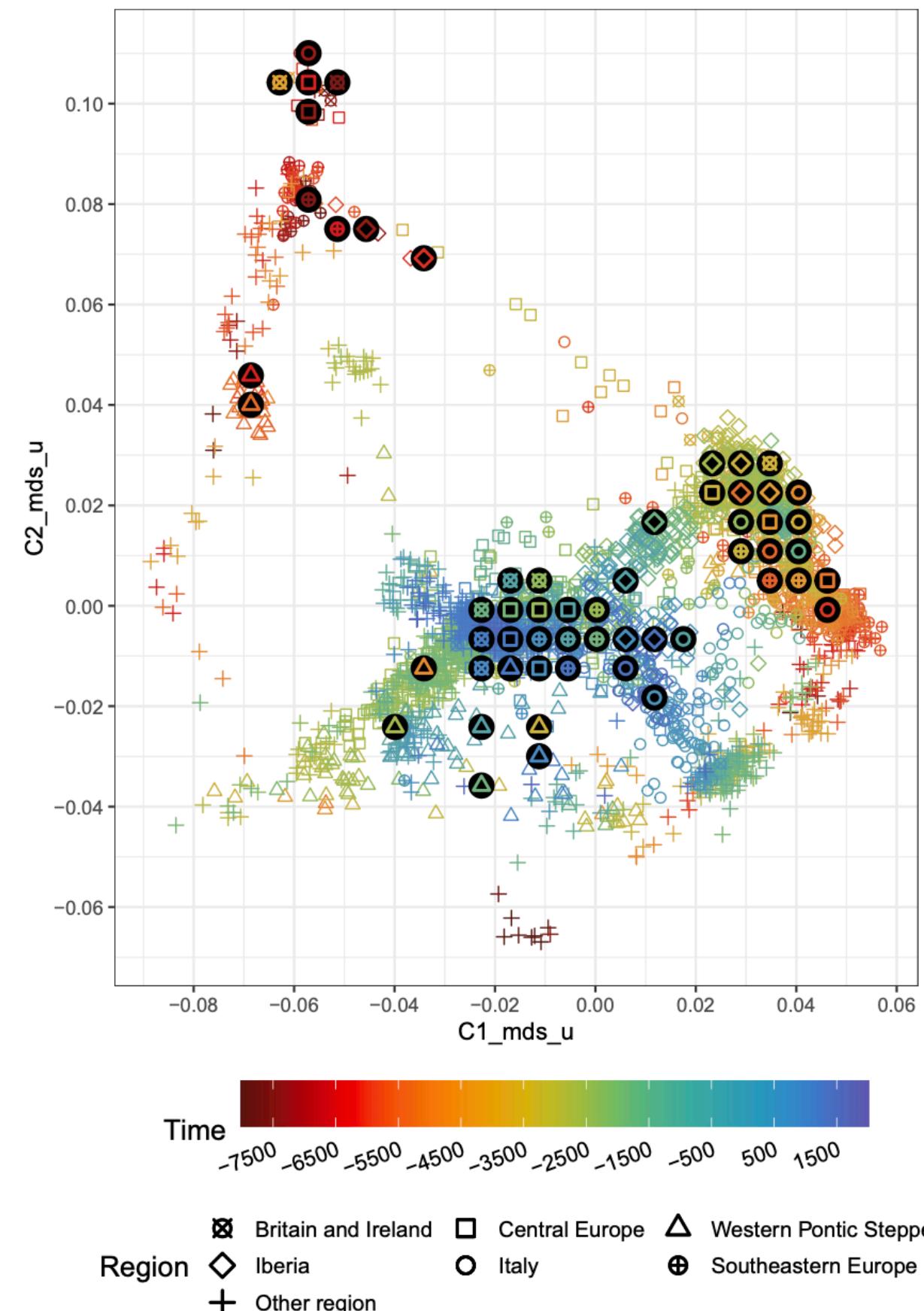


spatio-temporal interpolated "PCA field"

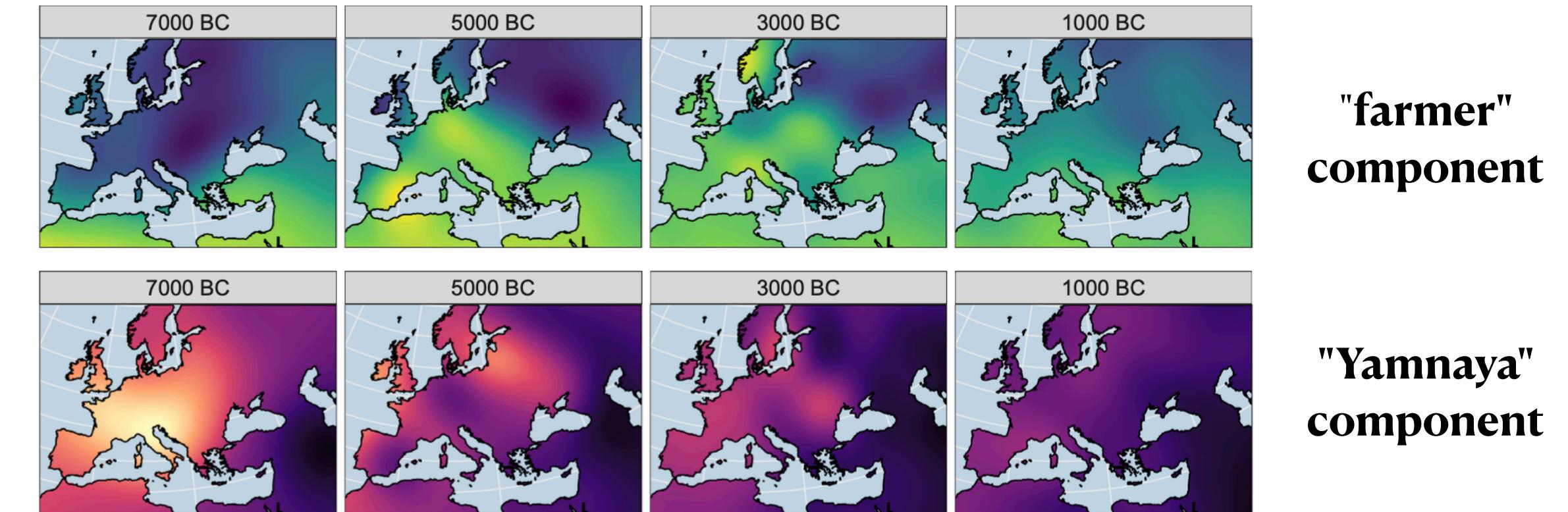


Where does a sample trace its ancestral origin?

MDS / PCA



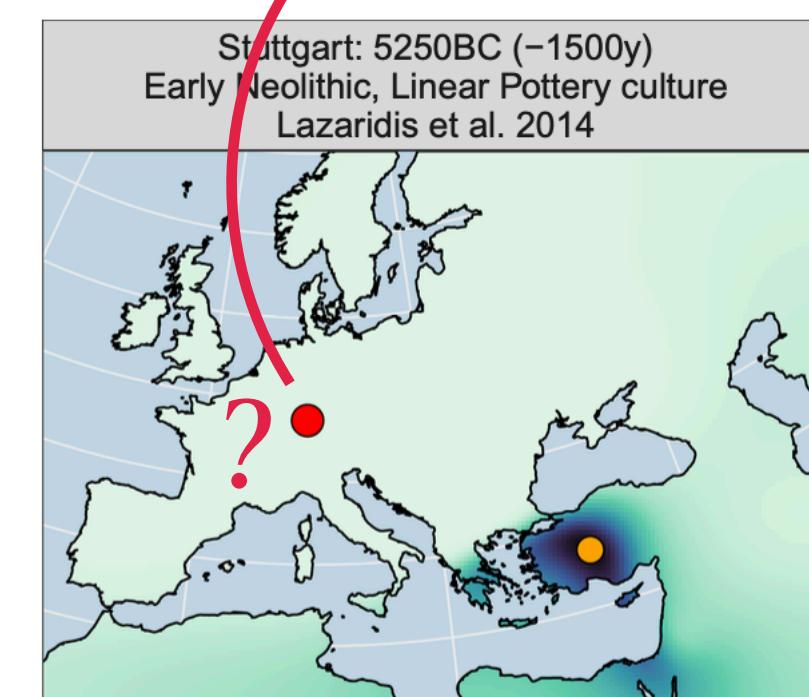
spatio-temporal interpolated "PCA field"



"farmer"
component

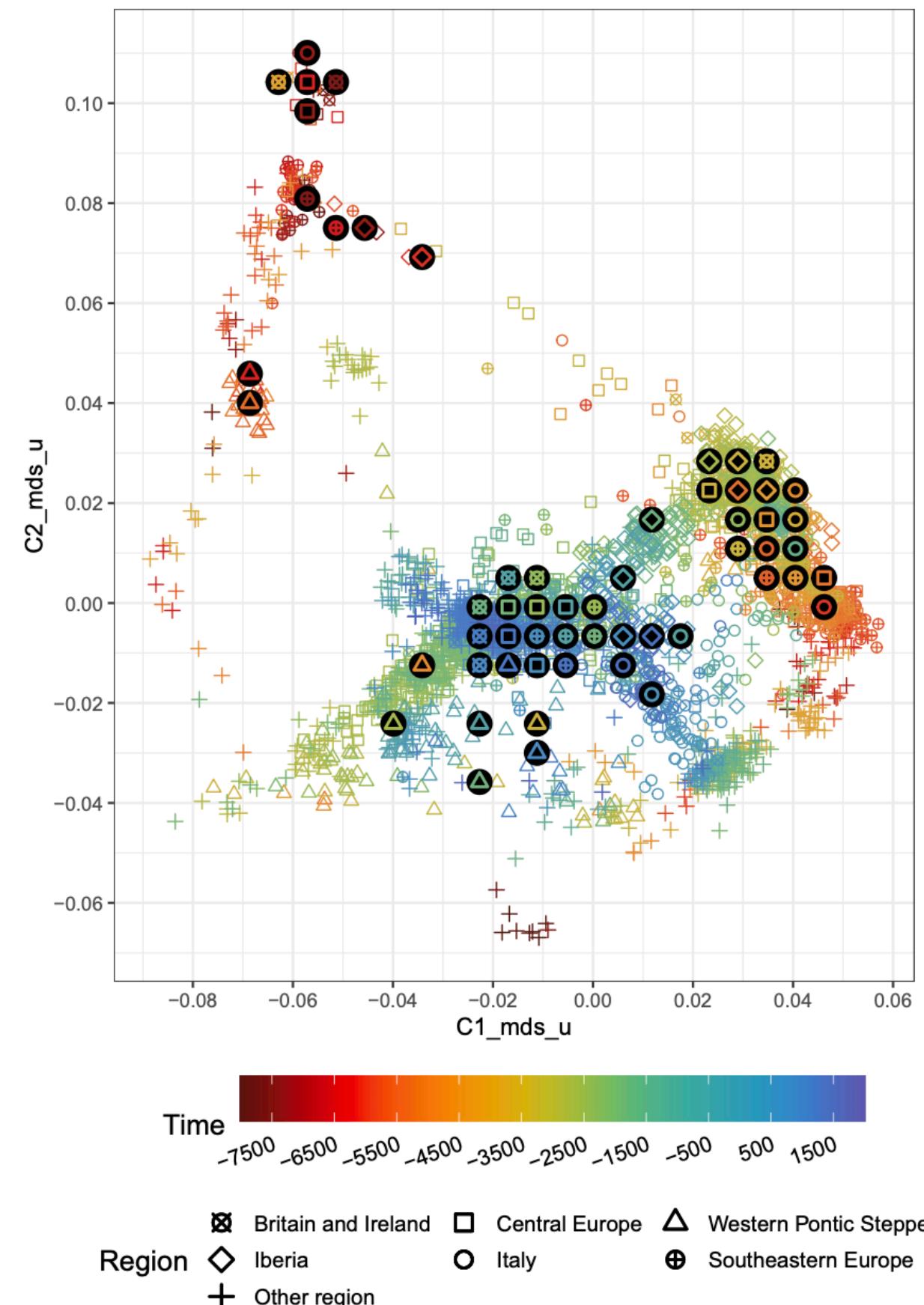
"Yamnaya"
component

where in the PCA
field are the closest
component values?

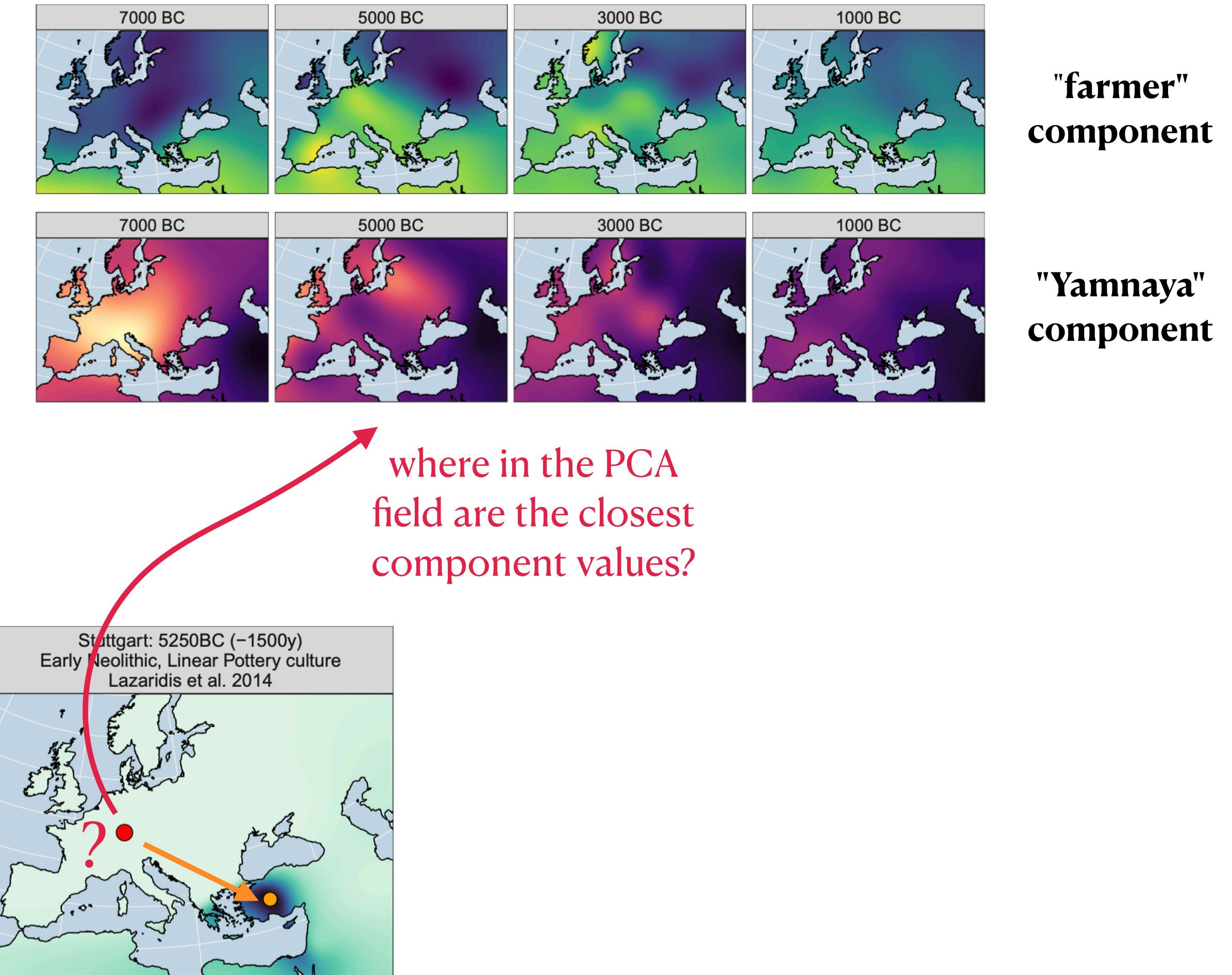


Where does a sample trace its ancestral origin?

MDS / PCA

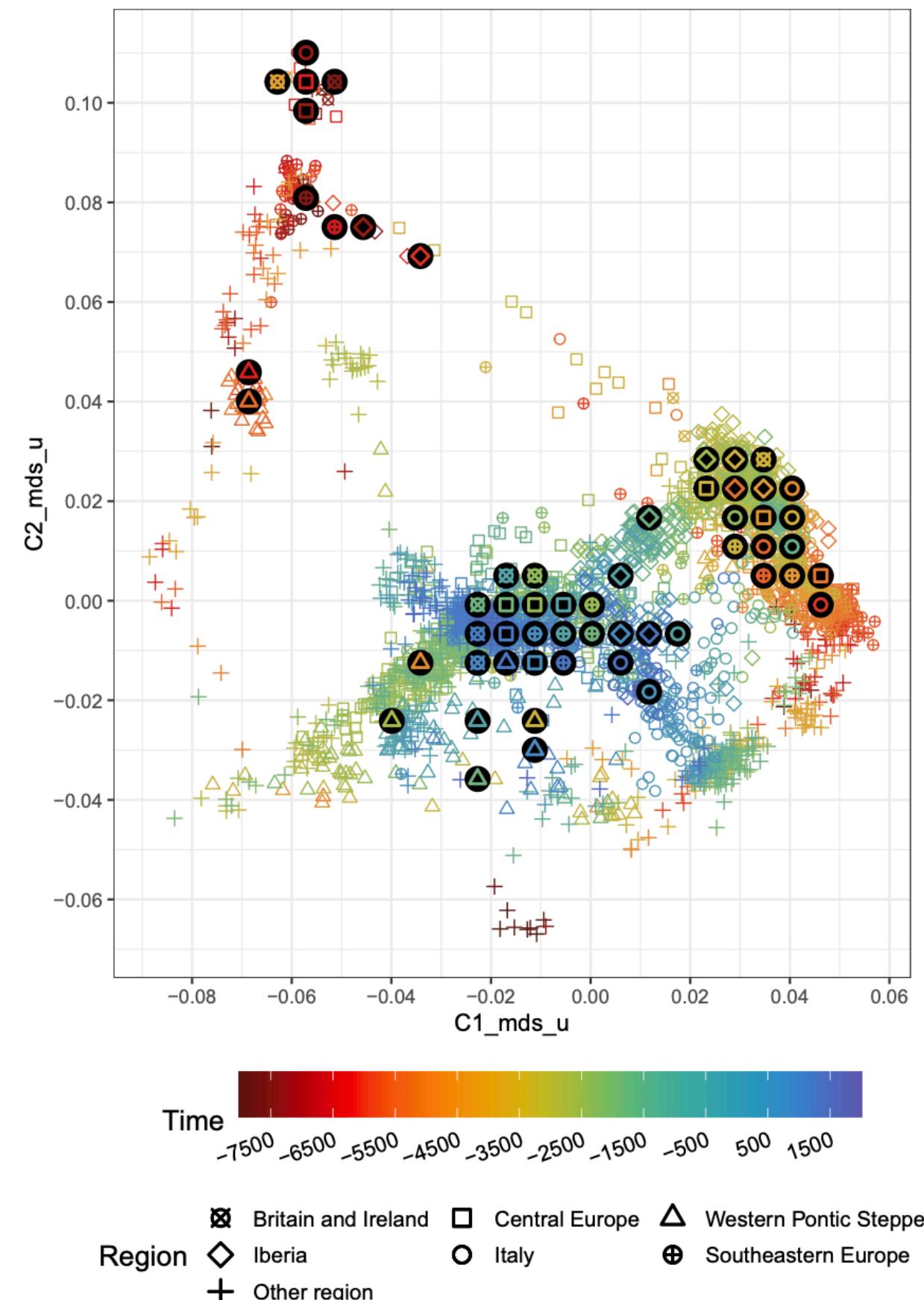


spatio-temporal interpolated "PCA field"

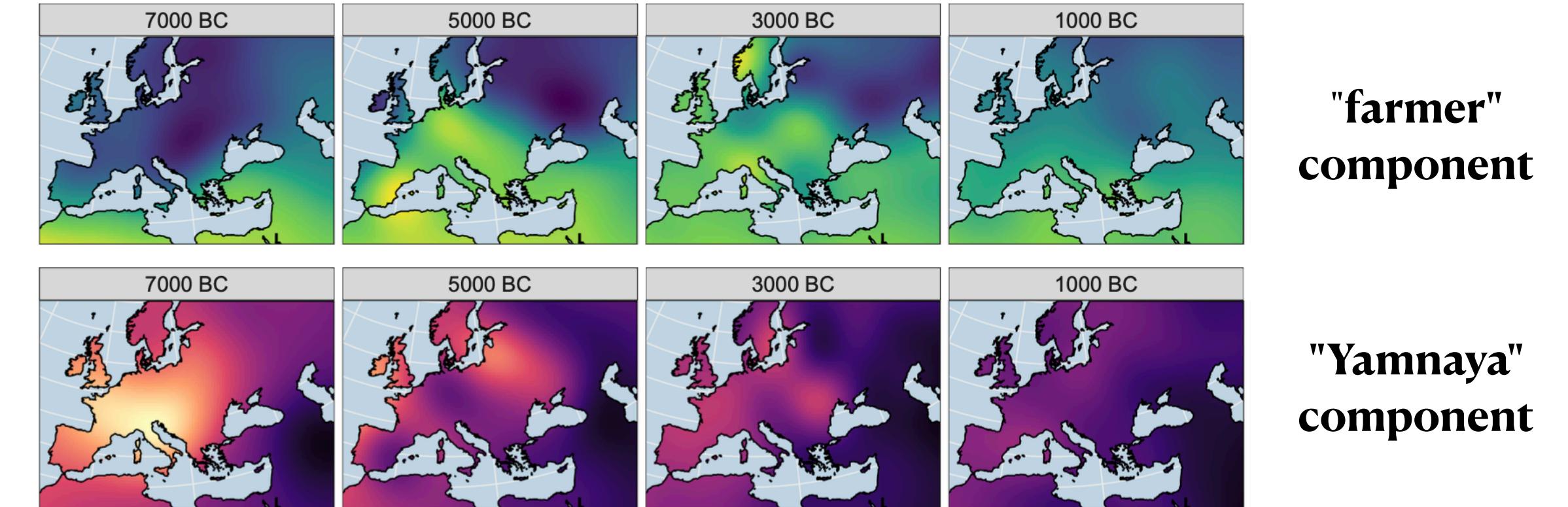


Where does a sample trace its ancestral origin?

MDS / PCA



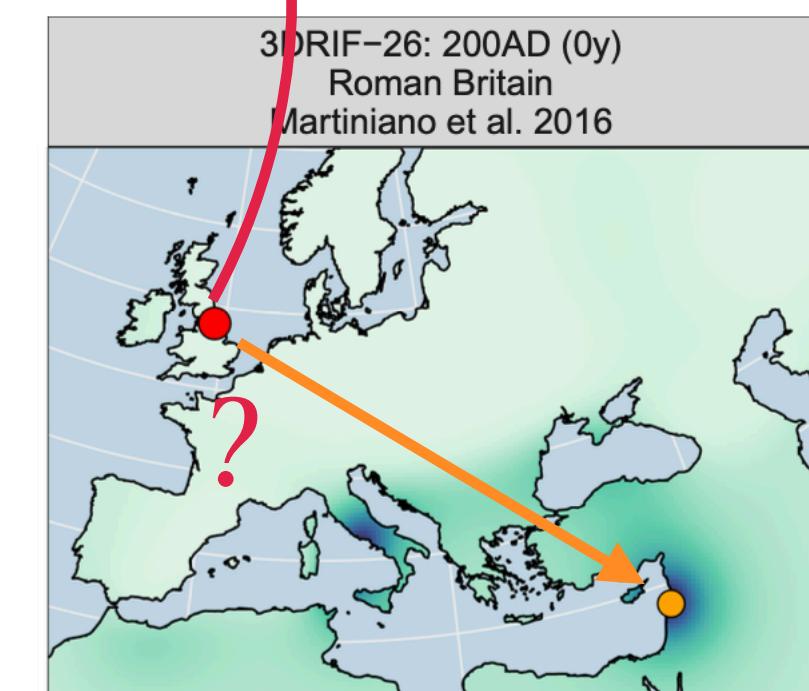
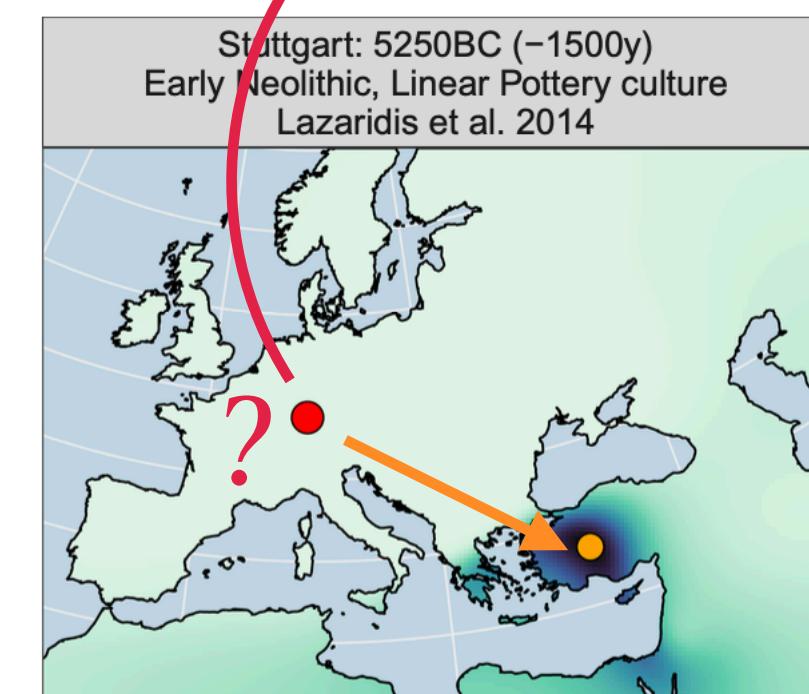
spatio-temporal interpolated "PCA field"



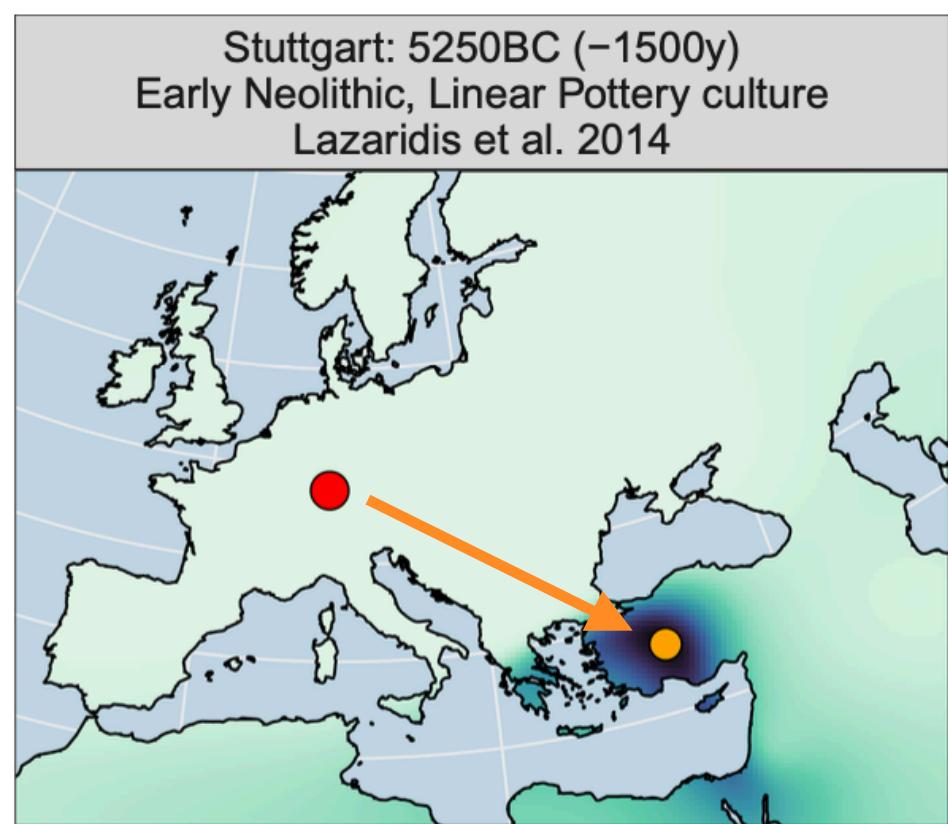
"farmer"
component

"Yamnaya"
component

where in the PCA
field are the closest
component values?



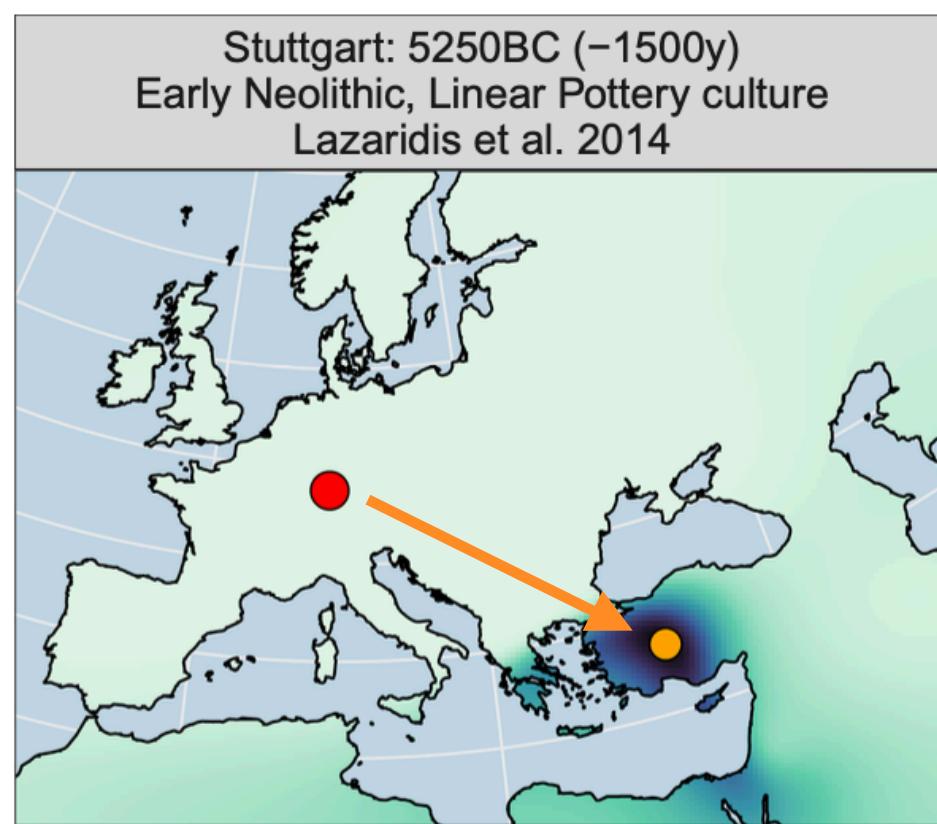
A "global" summary of individual mobilities across time



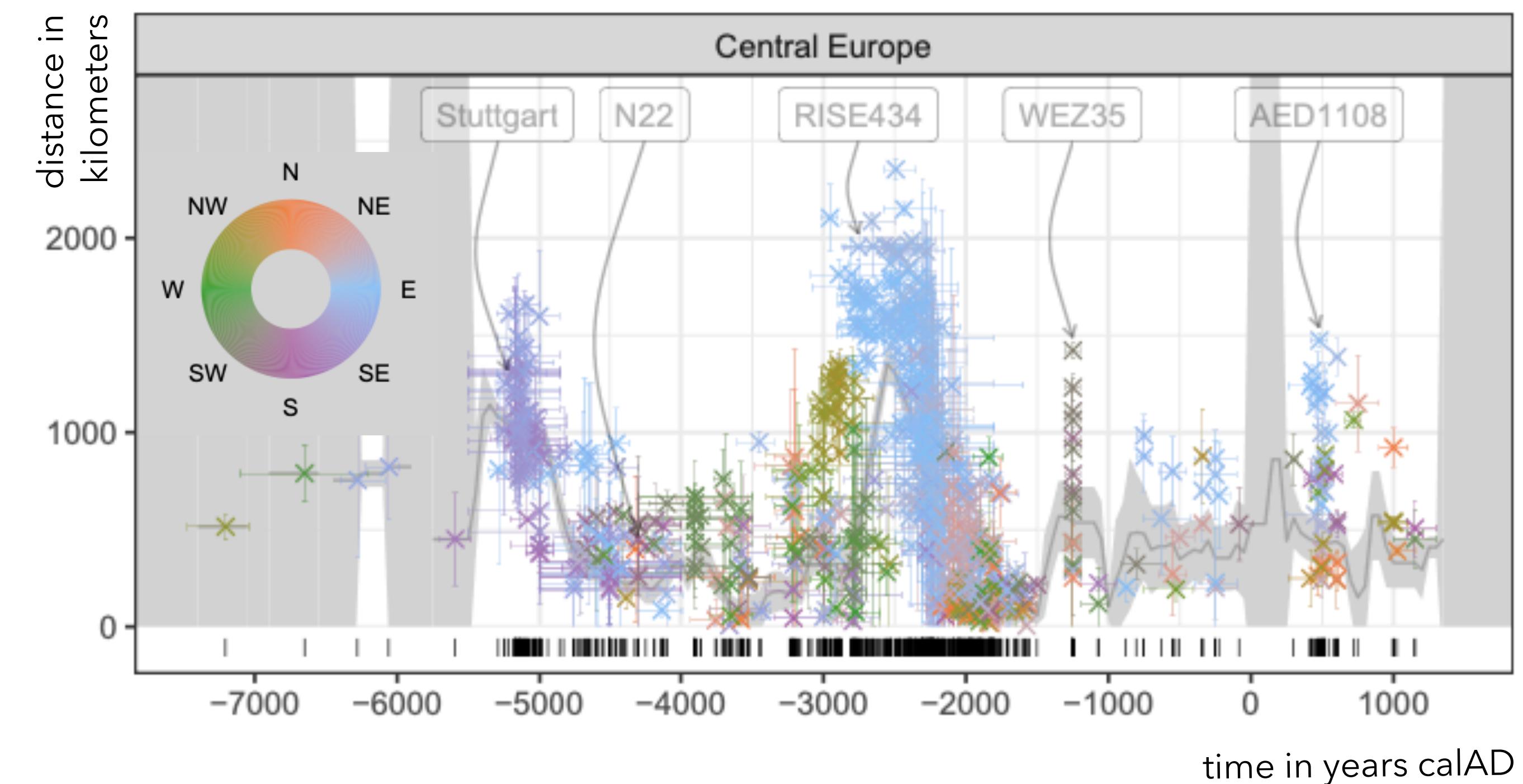
aggregation of **lengths**
and **directions** of many
mobility vectors



A "global" summary of individual mobilities across time



aggregation of **lengths**
and **directions** of many
mobility vectors





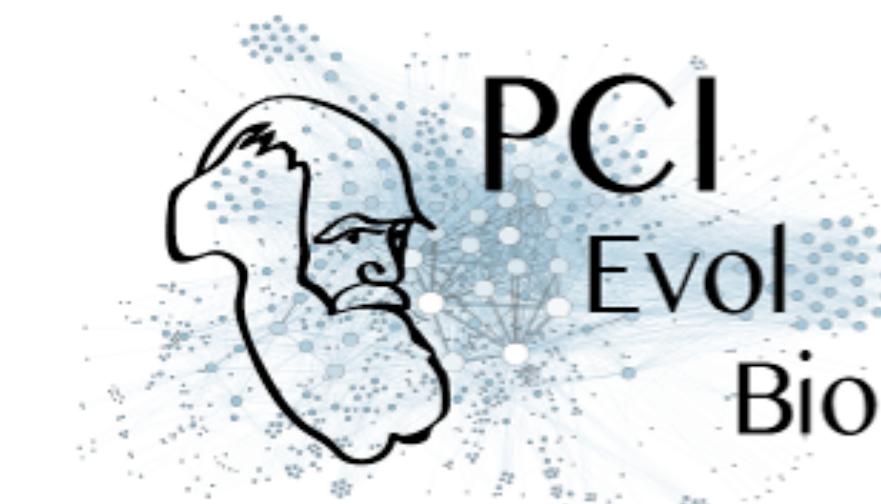
slendr: a framework for spatio-temporal population genomic simulations on geographic landscapes

Martin Petr, Benjamin C. Haller, Peter L. Ralph, Fernando Racimo

(2023), *bioRxiv*, ver.5, peer-reviewed and recommended by PCI Evol Biol

<https://doi.org/10.1101/2022.03.20.485041>

www.slendr.net

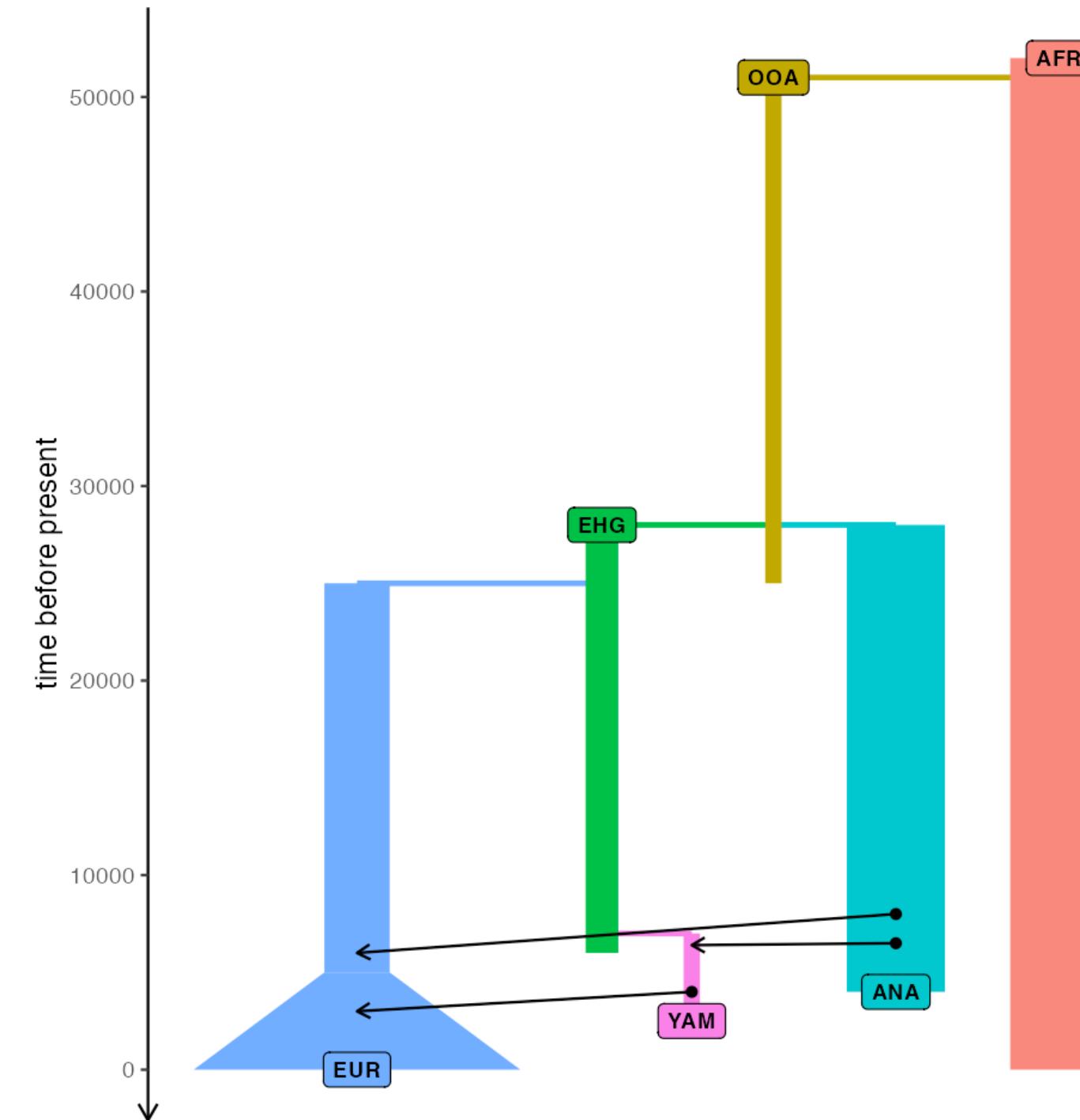


<https://peercommunityin.org/>

slendr can simulate genomic data from...

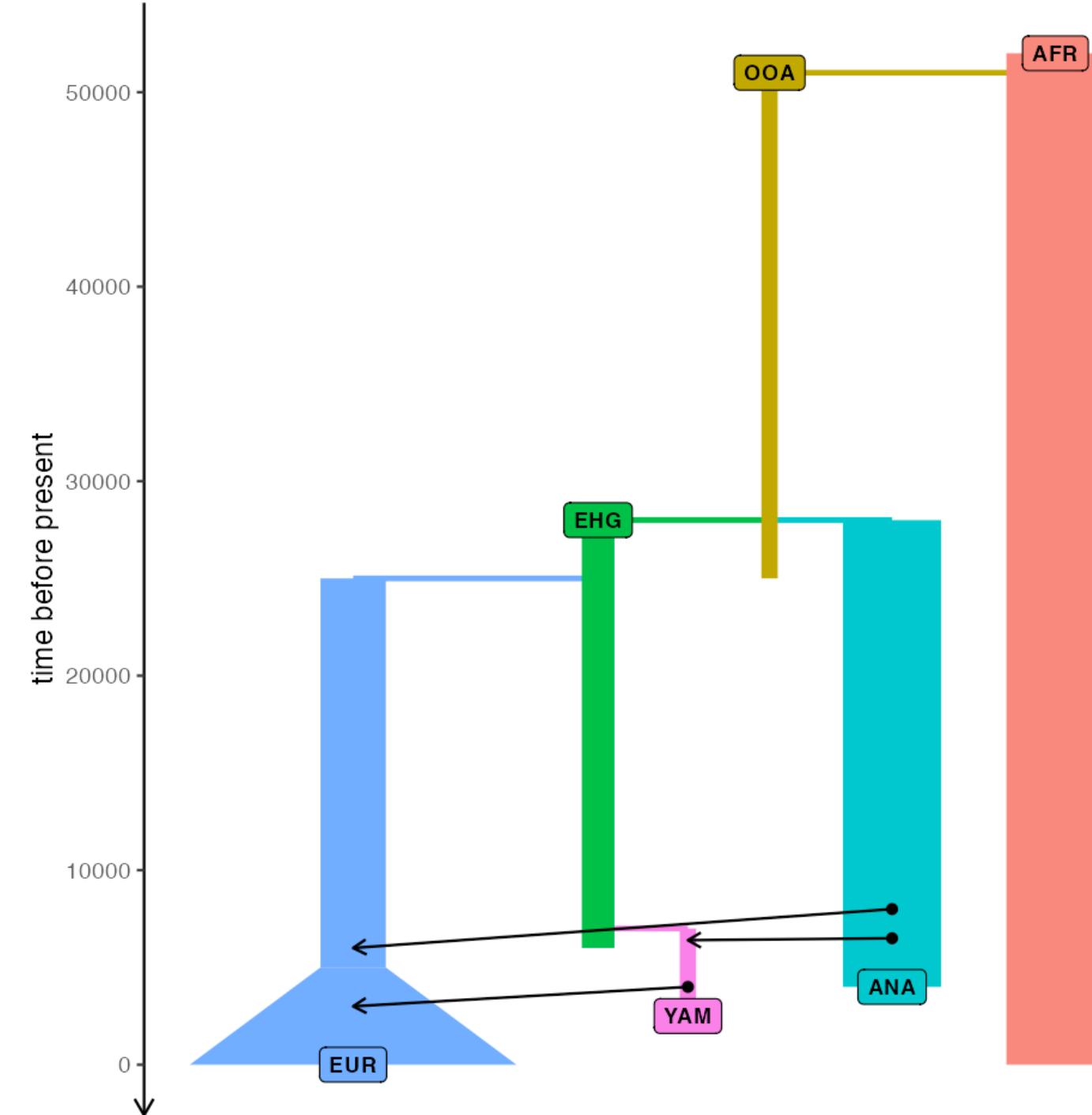
slendr can simulate genomic data from...

... traditional demographic models...

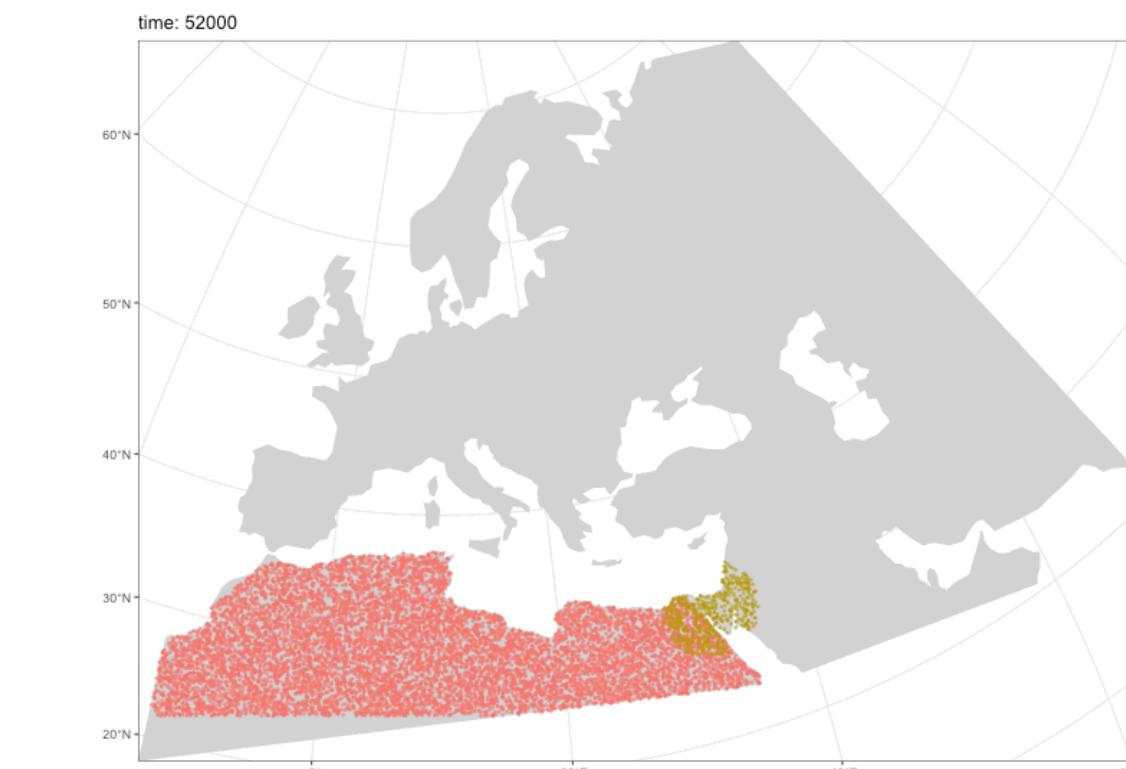
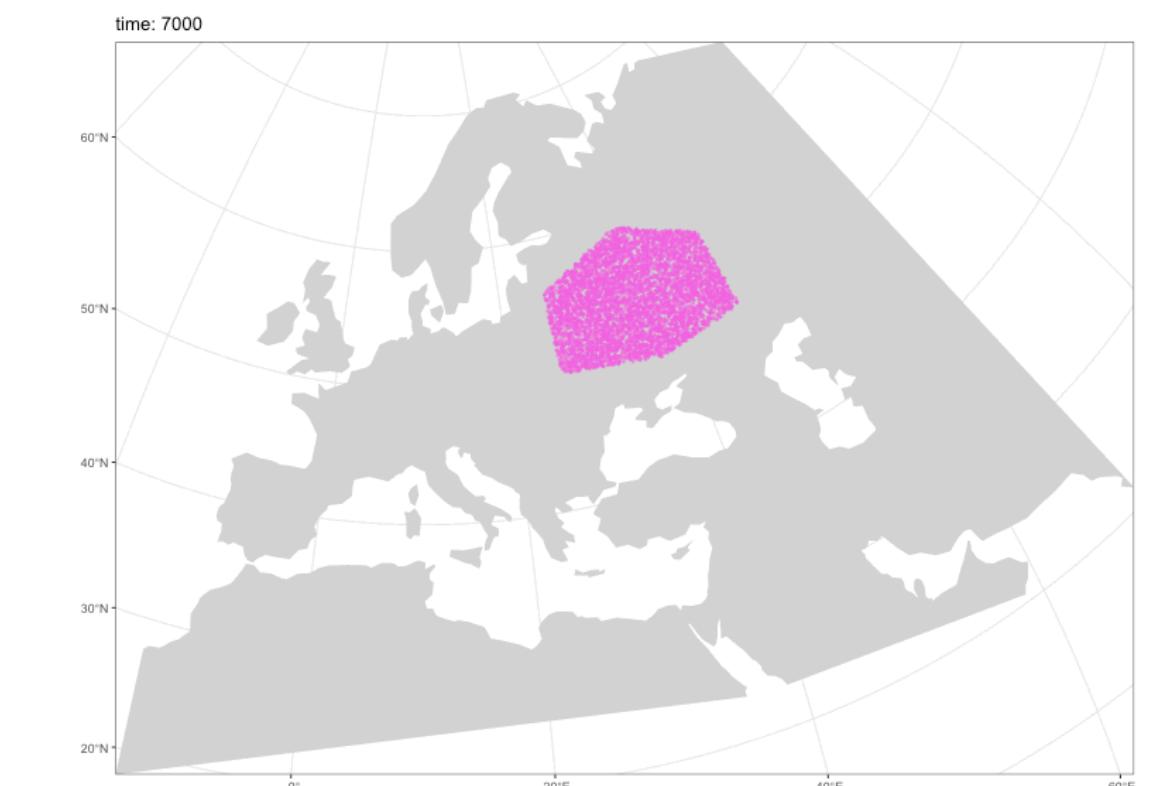


slendr can simulate genomic data from...

... traditional demographic models...



... but also explicitly spatial models!



Tree sequence / Ancestral Recombination Graph (ARG)

tsinfer

Inferring whole-genome histories in large population datasets

Jerome Kelleher^{ID}*, Yan Wong, Anthony W. Wohns^{ID}, Chaimaa Fadil^{ID}, Patrick K. Albers^{ID} and Gil McVean^{ID}



2019

Relate

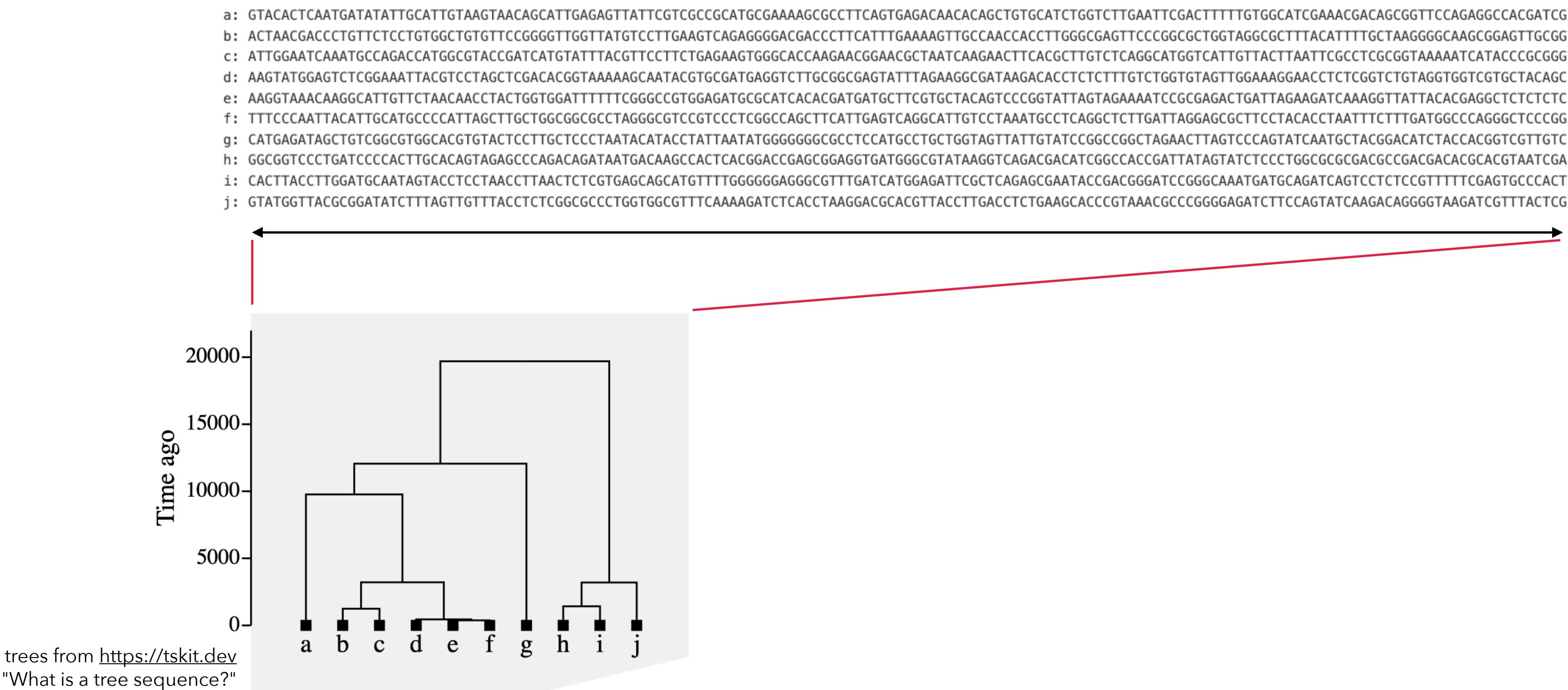
A method for genome-wide genealogy estimation for thousands of samples

Leo Speidel^{ID}¹, Marie Forest², Sinan Shi¹ and Simon R. Myers^{ID}^{1,3*}



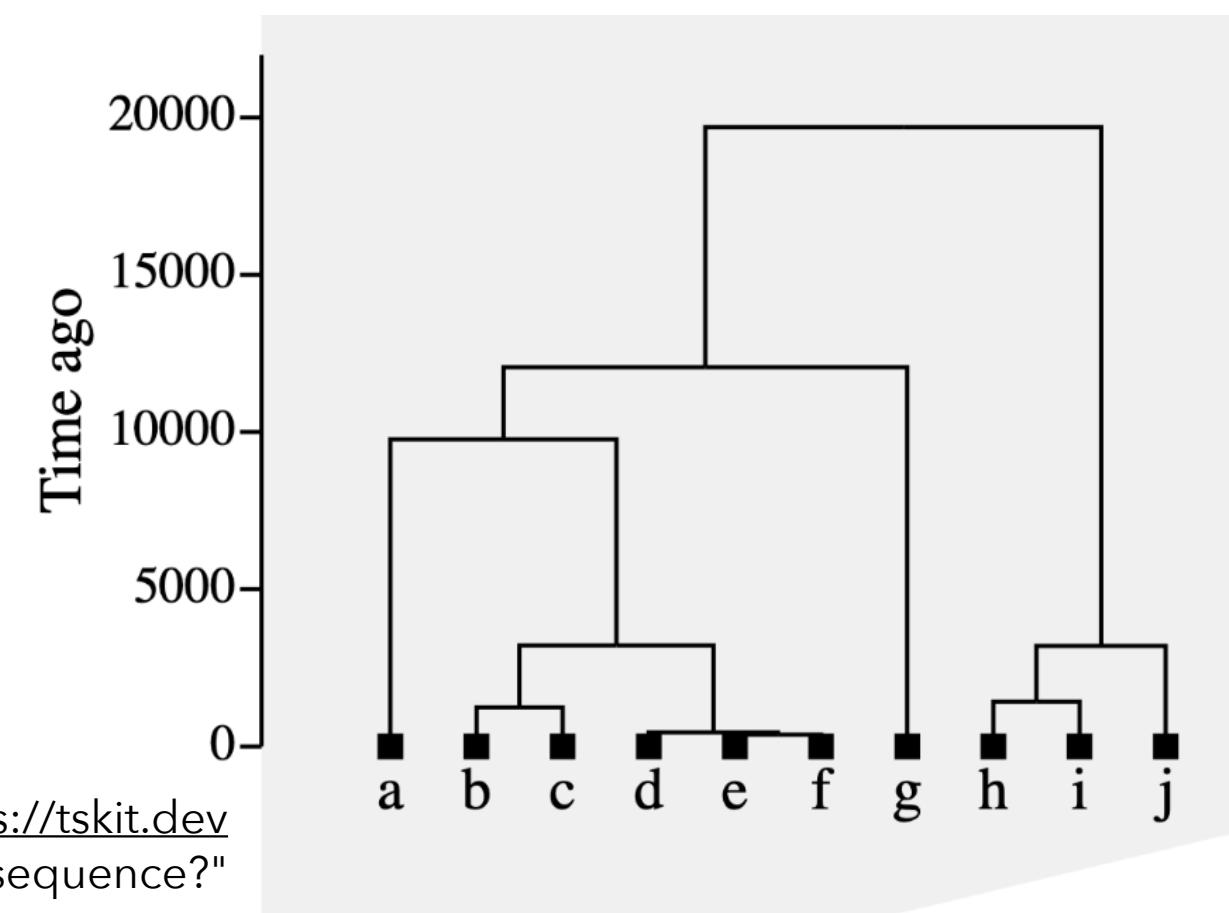
2019

Tree sequence / Ancestral Recombination Graph (ARG)

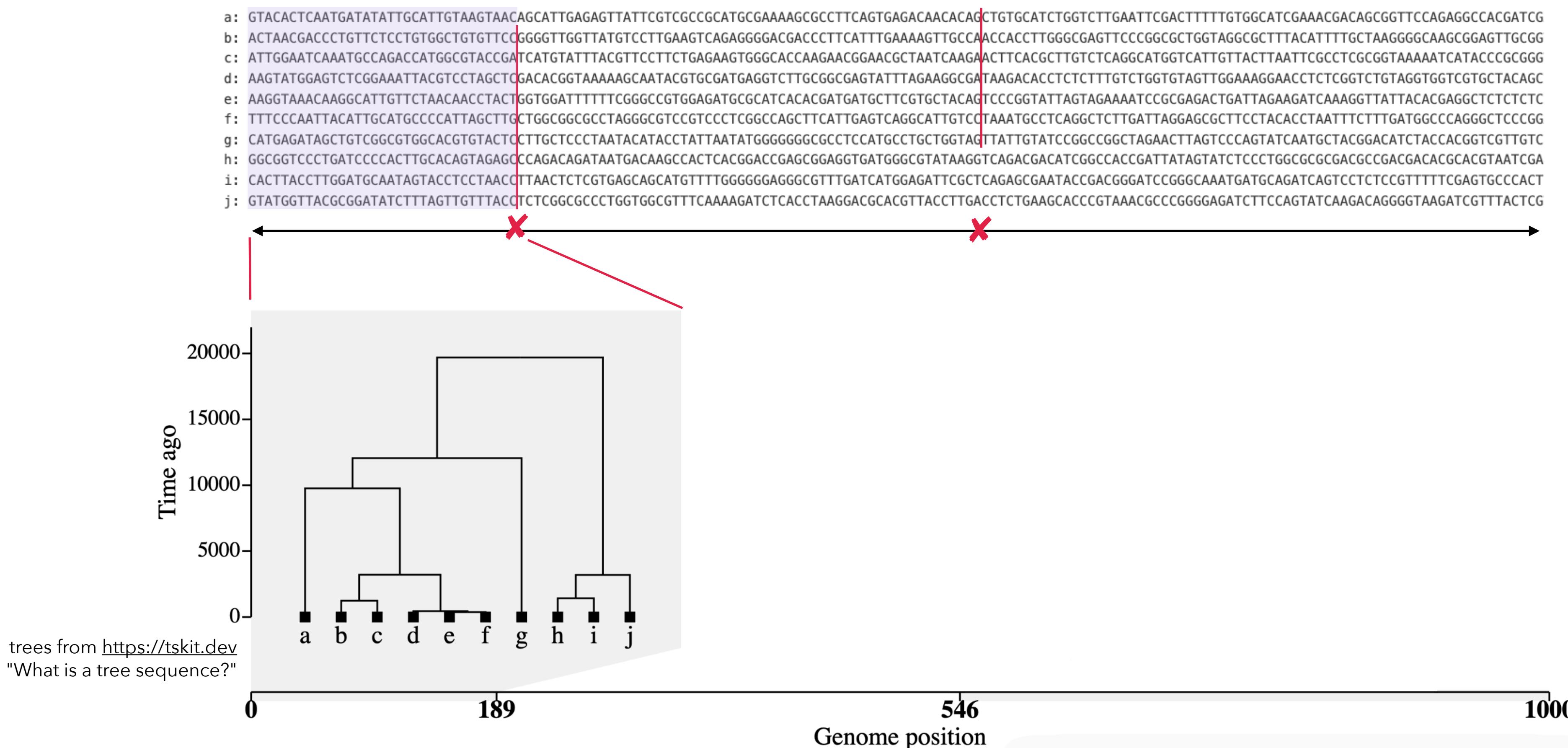


Tree sequence / Ancestral Recombination Graph (ARG)

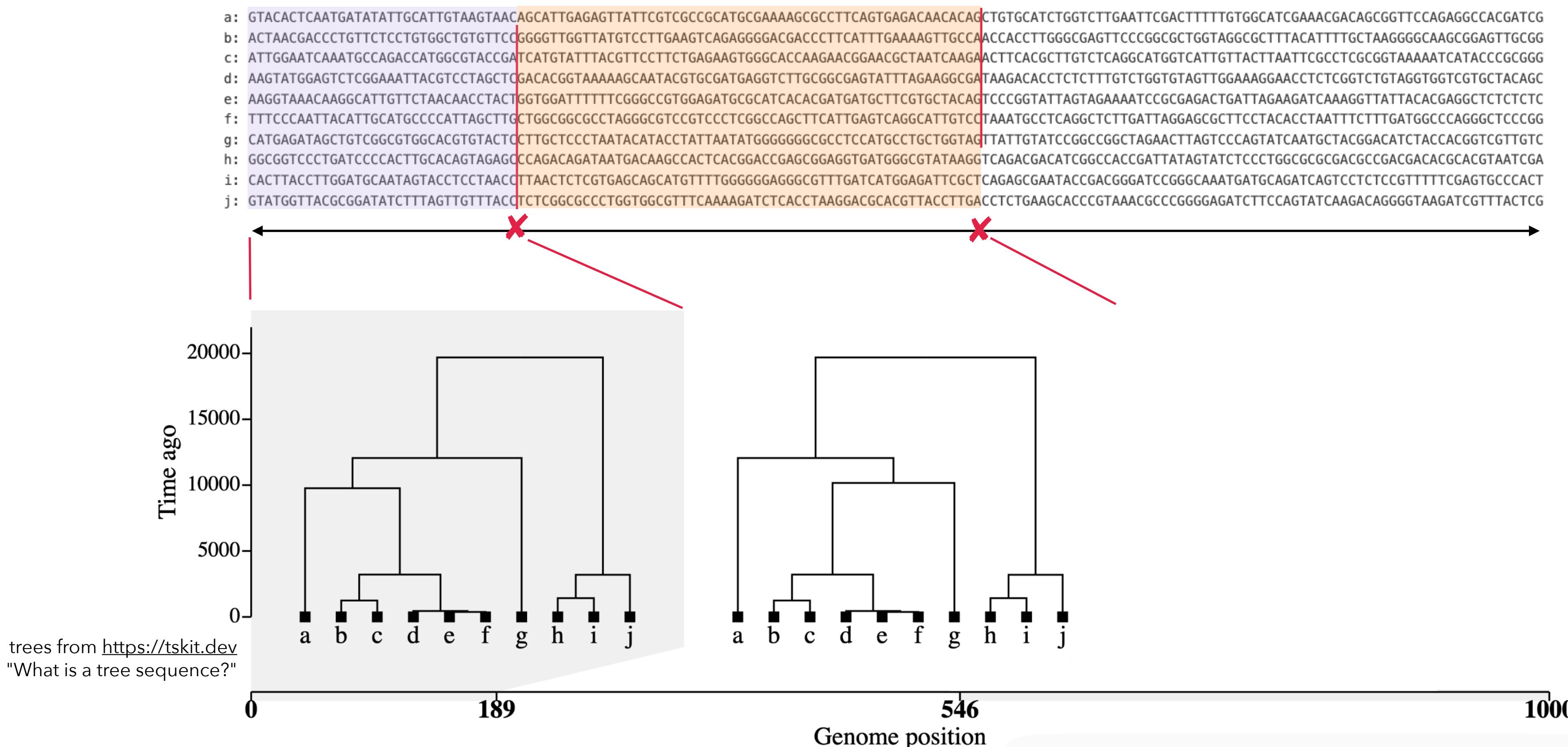
a: GTACACTCAATGATATATTGCATTGTAAGTAACGATTGAGAGTTTCGTGCCCATCGAAAAGCGCTTCAGTGAGACAACACAGCTGTGCATCTGGTCTTGAATTGACTTTGTGGCATCGAACGACAGCGTTCCAGAGGCCACGATCG
b: ACTAACGACCTGTTCTCTGTGGCTGTTCGGGTTGGTTATGCTTGAAGTCAGAGGGACGACCTTCATTGAAAGTTGCAACCACCTTGGCGAGTCCCGCGCTGGTAGGCCTTACATTGCTAAGGGCAAGCGAGTTGCGG
c: ATTGGAATCAAATGCCAGACCATGGTACCGATCATGTATTACGTTCTCTGAGAAAGTGGCACCAGAACGAAACGCTAATCAAAGAACCTCACGCTTGTCTCAGGCATGGTCAATTGTTACTTAATTGCTCGCGTAAAGATCACCGCGG
d: AAGTATGGAGTCTCGAAATTACGCTCTAGCTGACACGGTAAAAAGCAATACGTGCATGAGGTCTTGGCGAGTATTAGAAGGCATAAGACACCTCTTGTCTGGTAGTTGGAAAGGAACCTCTCGGTCTGTAGGTGGTGTACAGC
e: AAGGAAACAAAGGCAATTGTTCTAACAACTACTGGTGGATTTCGGCCGTGGAGATGCGCATCACAGATGATGCTCGTACAGTCCCCTGATTAGAAGGTTAAAGGTTACACGAGGCTCTCTC
f: TTTCCAATTACATTGCATGCCCATAGCTTGTGGCGCCCTAGGGCGTCCGTCCCTGGCCAGCTTCAATTGAGTCAGGCATTGCTAAATGCTCAGGCCTTGATTAGGAGCGCTTACACCTAATTGATGGCCAGGGCTCCGG
g: CATGAGATAGCTGCGCGTGGCAGCTGTACTCTTGTCTCCCTAACACCTATTAAATATGGGGGGCGCCTCATGCCTGCTGGTAGTTATTGTATCCGGCCGGCTAGAAACTTAGTCCAGTATCAATGCTACGGACATCTACCACGGTGTGTC
h: GGCAGTCCCTGATCCCCACTTGCACTGAGAGCCAGACAGATAATGACAAGCCACTCACGGACCGAGCGGAGGTGATGGCGTATAAGGTCAAGACGACATGGCCACCGATTATAGTATCTCCCTGGCGCGACGCCACGCAGTAATCGA
i: CACTTACCTGGATGCAATAGTACCTCTAACCTTAACCTCGTGAACGAGCATGTTGGGGGAGGGCGTTGATCATGGAGATTGCTCAGAGCGAATACCGACGGGATCAGGCAAATGATGCAAGTCAGTCCCTCCGTTTCGAGTCCCCACT
j: GTATGGTTACGCGGATATCTTAGTTACCTCTCGGCCCTGGTGGCTTCAAAAGATCTCACCTAACGACGTTACCTTGACCTCTGAAGCACCCGAAACGCCGGGAGATCTTCAGTATCAAGACAGGGTAAGATCGTTACTCG



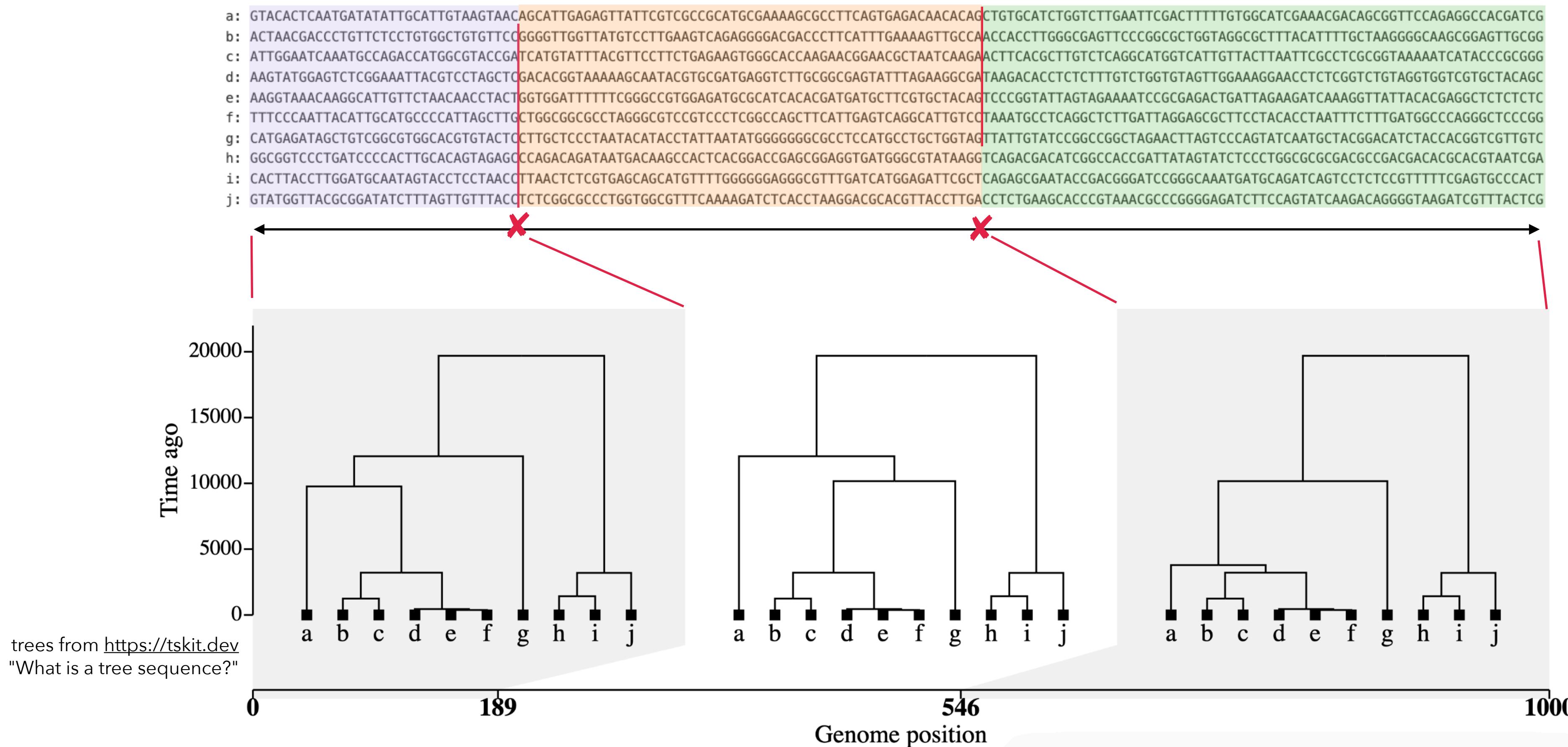
Tree sequence / Ancestral Recombination Graph (ARG)



Tree sequence / Ancestral Recombination Graph (ARG)

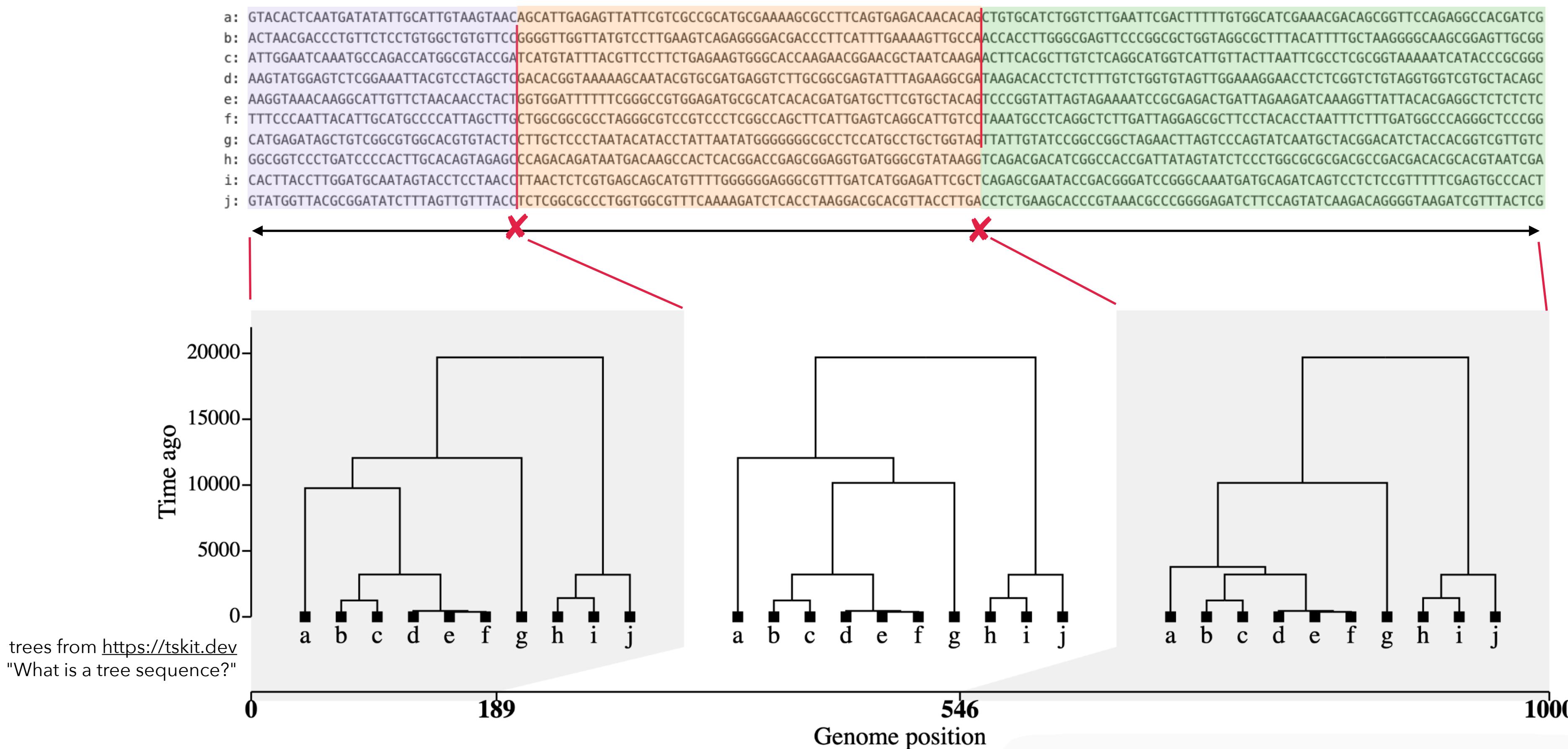


Tree sequence / Ancestral Recombination Graph (ARG)



Recombination produces a sequence of (correlated) trees along the genome

Tree sequence / Ancestral Recombination Graph (ARG)



A tree sequence represents
everything knowable about an evolutionary process.

What we (usually) have

#CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO	FORMAT	HG00096	HG00099	HG00101
12	60076	.	A	C	100	PASS	.	GT	1 0	0 0	0 0
12	60252	.	A	G	100	PASS	.	GT	0 0	0 0	0 1
12	60317	.	C	T	100	PASS	.	GT	0 0	0 0	0 0
12	60344	.	C	A	100	PASS	.	GT	0 0	0 0	0 0
12	60383	.	G	A	100	PASS	.	GT	0 0	0 0	0 0
12	60405	.	T	C	100	PASS	.	GT	0 0	0 0	0 0
12	60474	.	G	A	100	PASS	.	GT	0 0	0 1	0 1
12	60614	.	C	A	100	PASS	.	GT	0 0	0 0	0 0
12	60628	.	T	C	100	PASS	.	GT	0 0	0 0	0 0
12	60654	.	G	A	100	PASS	.	GT	0 0	0 0	0 0
12	61021	.	C	T	100	PASS	.	GT	0 0	0 0	0 0
12	61107	.	G	T	100	PASS	.	GT	0 0	0 0	0 0
12	61172	.	G	A	100	PASS	.	GT	0 0	0 0	0 0
12	61220	.	G	A	100	PASS	.	GT	0 0	0 0	0 1
12	61258	.	C	T	100	PASS	.	GT	0 0	0 0	0 1
12	61272	.	T	C	100	PASS	.	GT	0 0	0 0	0 0
12	61329	.	C	T	100	PASS	.	GT	0 0	0 0	0 0
12	61341	.	G	A	100	PASS	.	GT	0 0	0 1	0 1
12	61368	.	C	T	100	PASS	.	GT	0 0	0 1	0 1
12	61392	.	T	A	100	PASS	.	GT	0 0	0 0	0 0
12	61405	.	G	C	100	PASS	.	GT	0 0	0 0	0 0
12	61411	.	C	A	100	PASS	.	GT	0 0	0 0	0 0
12	61416	.	G	A	100	PASS	.	GT	0 0	0 0	0 0
12	61422	.	C	T	100	PASS	.	GT	0 0	0 0	0 0
12	61476	.	C	G	100	PASS	.	GT	0 0	0 0	0 0
12	61510	.	G	A	100	PASS	.	GT	0 0	0 0	0 0
12	61516	.	C	T	100	PASS	.	GT	0 0	0 0	0 0
12	61552	.	C	T	100	PASS	.	GT	0 0	0 0	0 0
12	61604	.	T	G	100	PASS	.	GT	0 0	0 0	0 0
12	61687	.	G	A	100	PASS	.	GT	1 0	0 1	0 1
12	61700	.	C	T	100	PASS	.	GT	0 0	0 0	0 0

What we (usually) have

positions along a chromosome

#CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO	FORMAT	HG00096	HG00099	HG00101
12	60076	.	A	C	100	PASS	.	GT	1 0	0 0	0 0
12	60252	.	A	G	100	PASS	.	GT	0 0	0 0	0 1
12	60317	.	C	T	100	PASS	.	GT	0 0	0 0	0 0
12	60344	.	C	A	100	PASS	.	GT	0 0	0 0	0 0
12	60383	.	G	A	100	PASS	.	GT	0 0	0 0	0 0
12	60405	.	T	C	100	PASS	.	GT	0 0	0 0	0 0
12	60474	.	G	A	100	PASS	.	GT	0 0	0 1	0 1
12	60614	.	C	A	100	PASS	.	GT	0 0	0 0	0 0
12	60628	.	T	C	100	PASS	.	GT	0 0	0 0	0 0
12	60654	.	G	A	100	PASS	.	GT	0 0	0 0	0 0
12	61021	.	C	T	100	PASS	.	GT	0 0	0 0	0 0
12	61107	.	G	T	100	PASS	.	GT	0 0	0 0	0 0
12	61172	.	G	A	100	PASS	.	GT	0 0	0 0	0 0
12	61220	.	G	A	100	PASS	.	GT	0 0	0 0	0 1
12	61258	.	C	T	100	PASS	.	GT	0 0	0 0	0 1
12	61272	.	T	C	100	PASS	.	GT	0 0	0 0	0 0
12	61329	.	C	T	100	PASS	.	GT	0 0	0 0	0 0
12	61341	.	G	A	100	PASS	.	GT	0 0	0 1	0 1
12	61368	.	C	T	100	PASS	.	GT	0 0	0 1	0 1
12	61392	.	T	A	100	PASS	.	GT	0 0	0 0	0 0
12	61405	.	G	C	100	PASS	.	GT	0 0	0 0	0 0
12	61411	.	C	A	100	PASS	.	GT	0 0	0 0	0 0
12	61416	.	G	A	100	PASS	.	GT	0 0	0 0	0 0
12	61422	.	C	T	100	PASS	.	GT	0 0	0 0	0 0
12	61476	.	C	G	100	PASS	.	GT	0 0	0 0	0 0
12	61510	.	G	A	100	PASS	.	GT	0 0	0 0	0 0
12	61516	.	C	T	100	PASS	.	GT	0 0	0 0	0 0
12	61552	.	C	T	100	PASS	.	GT	0 0	0 0	0 0
12	61604	.	T	G	100	PASS	.	GT	0 0	0 0	0 0
12	61687	.	G	A	100	PASS	.	GT	1 0	0 1	0 1
12	61700	.	C	T	100	PASS	.	GT	0 0	0 0	0 0

What we (usually) have

#CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO	FORMAT	ind. #1	ind. #2	ind. #3	...
12	60076	.	A	C	100	PASS	.	GT	1 0	0 0	0 0	
12	60252	.	A	G	100	PASS	.	GT	0 0	0 0	0 1	
12	60317	.	C	T	100	PASS	.	GT	0 0	0 0	0 0	
12	60344	.	C	A	100	PASS	.	GT	0 0	0 0	0 0	
12	60383	.	G	A	100	PASS	.	GT	0 0	0 0	0 0	
12	60405	.	T	C	100	PASS	.	GT	0 0	0 0	0 0	
12	60474	.	G	A	100	PASS	.	GT	0 0	0 1	0 1	
12	60614	.	C	A	100	PASS	.	GT	0 0	0 0	0 0	
12	60628	.	T	C	100	PASS	.	GT	0 0	0 0	0 0	
12	60654	.	G	A	100	PASS	.	GT	0 0	0 0	0 0	
12	61021	.	C	T	100	PASS	.	GT	0 0	0 0	0 0	
12	61107	.	G	T	100	PASS	.	GT	0 0	0 0	0 0	
12	61172	.	G	A	100	PASS	.	GT	0 0	0 0	0 0	
12	61220	.	G	A	100	PASS	.	GT	0 0	0 0	0 1	
12	61258	.	C	T	100	PASS	.	GT	0 0	0 0	0 1	
12	61272	.	T	C	100	PASS	.	GT	0 0	0 0	0 0	
12	61329	.	C	T	100	PASS	.	GT	0 0	0 0	0 0	
12	61341	.	G	A	100	PASS	.	GT	0 0	0 1	0 1	
12	61368	.	C	T	100	PASS	.	GT	0 0	0 1	0 1	
12	61392	.	T	A	100	PASS	.	GT	0 0	0 0	0 0	
12	61405	.	G	C	100	PASS	.	GT	0 0	0 0	0 0	
12	61411	.	C	A	100	PASS	.	GT	0 0	0 0	0 0	
12	61416	.	G	A	100	PASS	.	GT	0 0	0 0	0 0	
12	61422	.	C	T	100	PASS	.	GT	0 0	0 0	0 0	
12	61476	.	C	G	100	PASS	.	GT	0 0	0 0	0 0	
12	61510	.	G	A	100	PASS	.	GT	0 0	0 0	0 0	
12	61516	.	C	T	100	PASS	.	GT	0 0	0 0	0 0	
12	61552	.	C	T	100	PASS	.	GT	0 0	0 0	0 0	
12	61604	.	T	G	100	PASS	.	GT	0 0	0 0	0 0	
12	61687	.	G	A	100	PASS	.	GT	1 0	0 1	0 1	
12	61700	.	C	T	100	PASS	.	GT	0 0	0 0	0 0	

positions along a chromosome



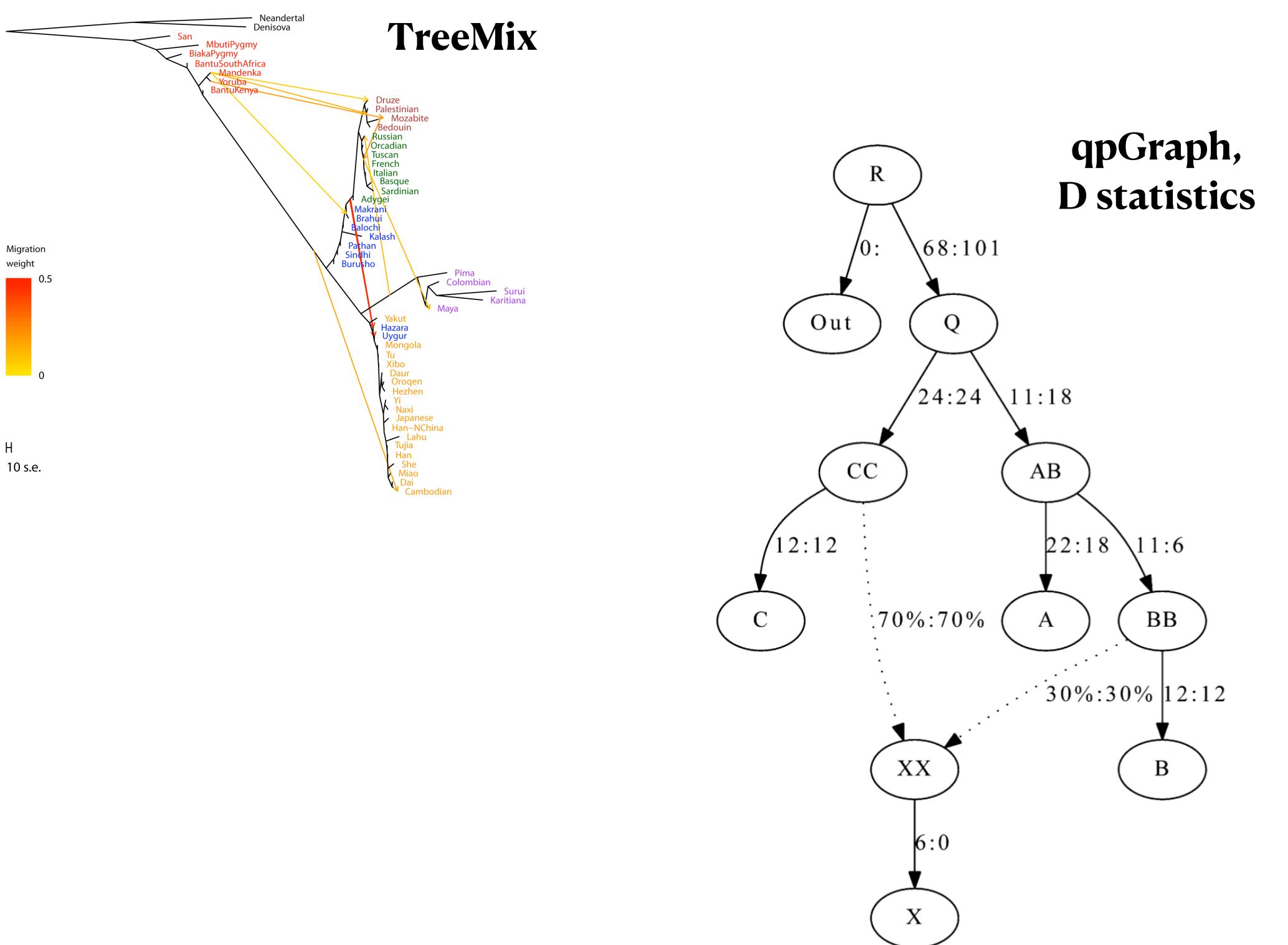
What we (usually) want

- estimate split times between populations



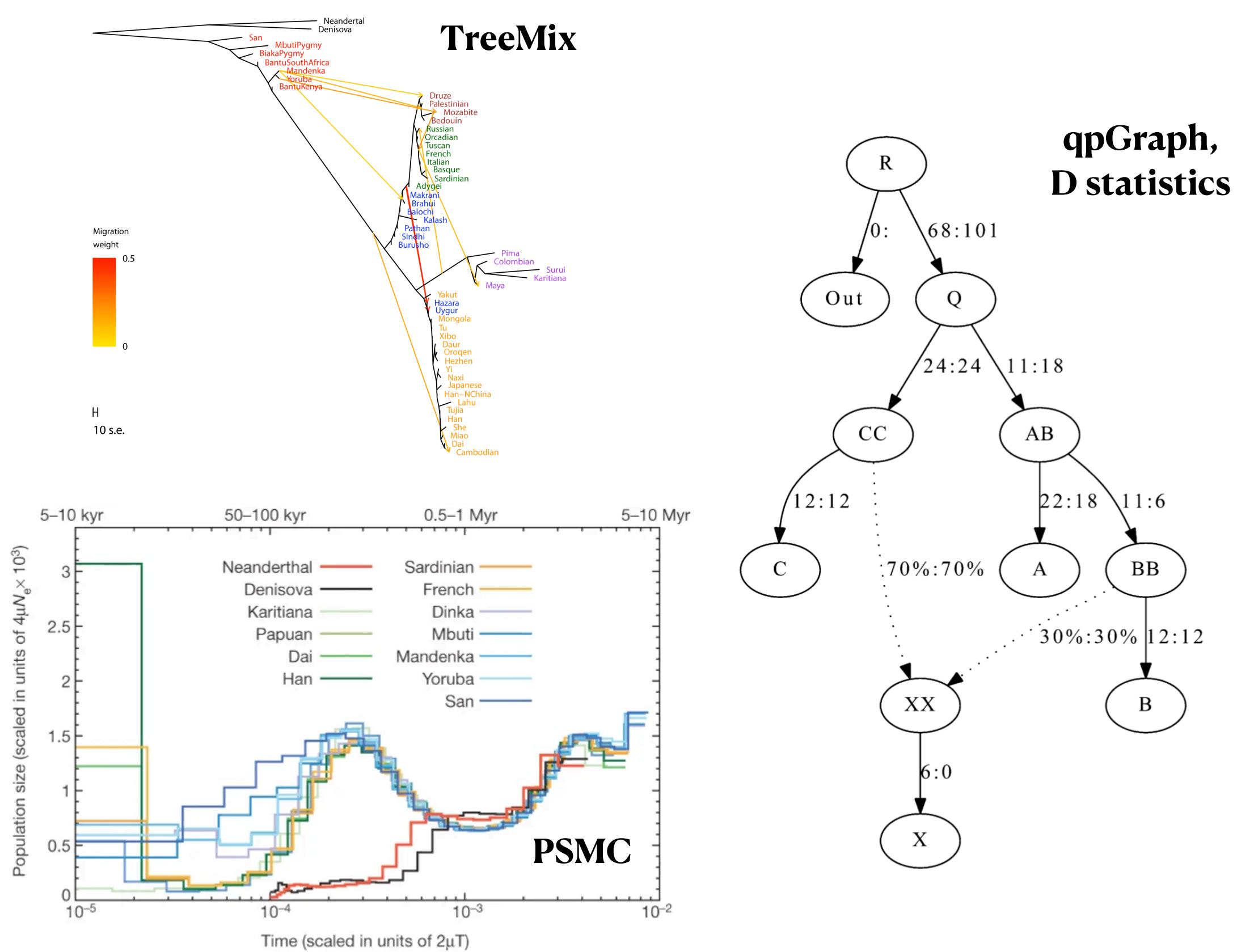
What we (usually) want

- estimate split times between populations
- investigate admixture hypotheses



What we (usually) want

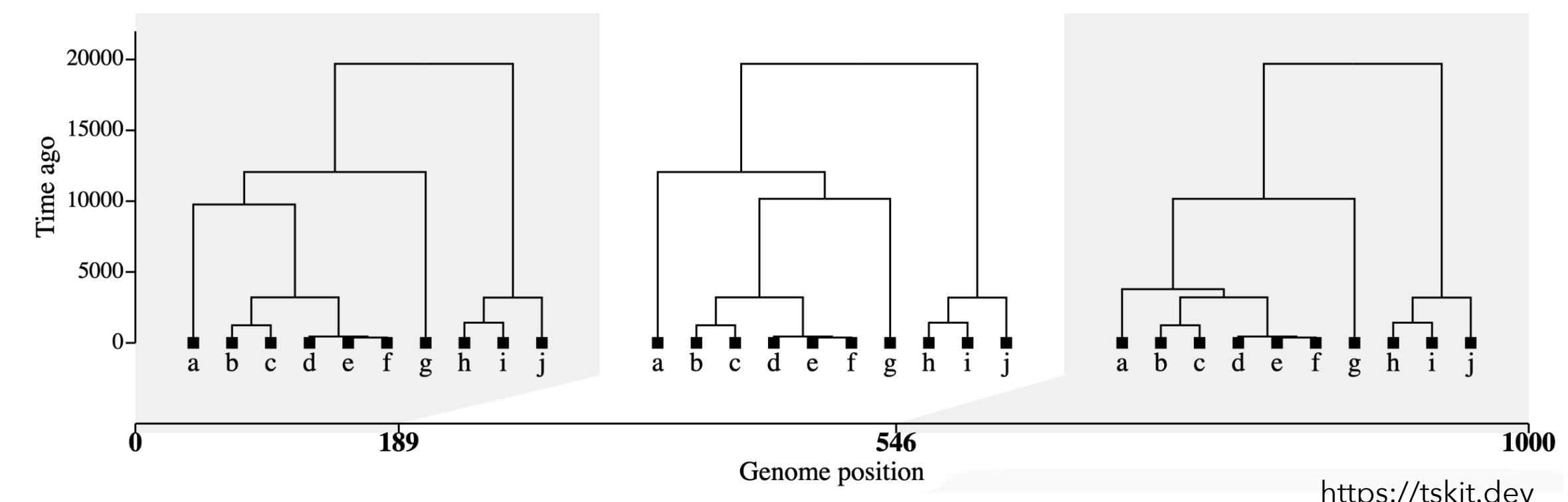
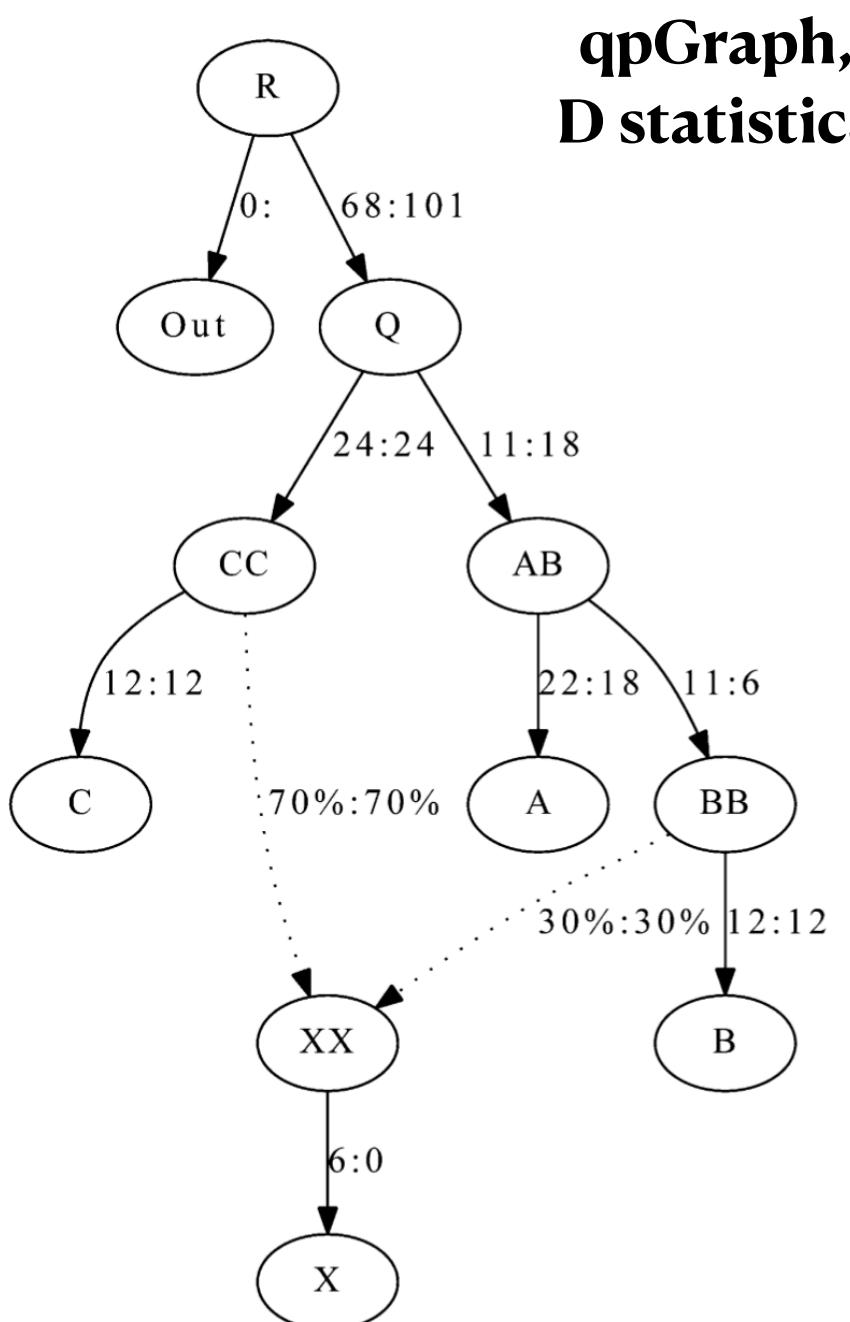
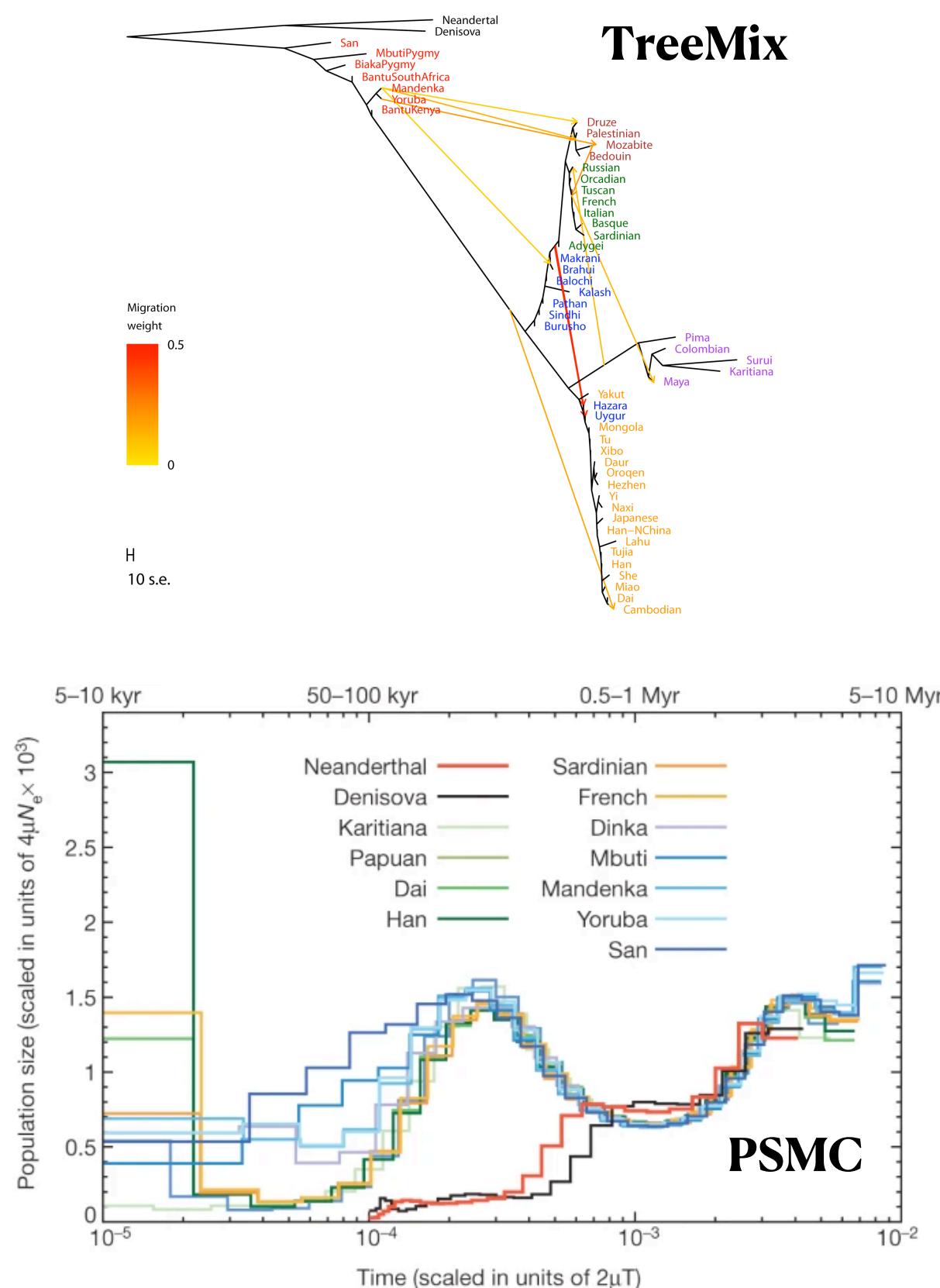
- estimate split times between populations
- investigate admixture hypotheses
- infer population sizes changes over time



What we (usually) want

- estimate split times between populations
 - investigate admixture hypotheses
 - infer population sizes changes over time

These are all features
of a real (but hidden)
genealogical process.

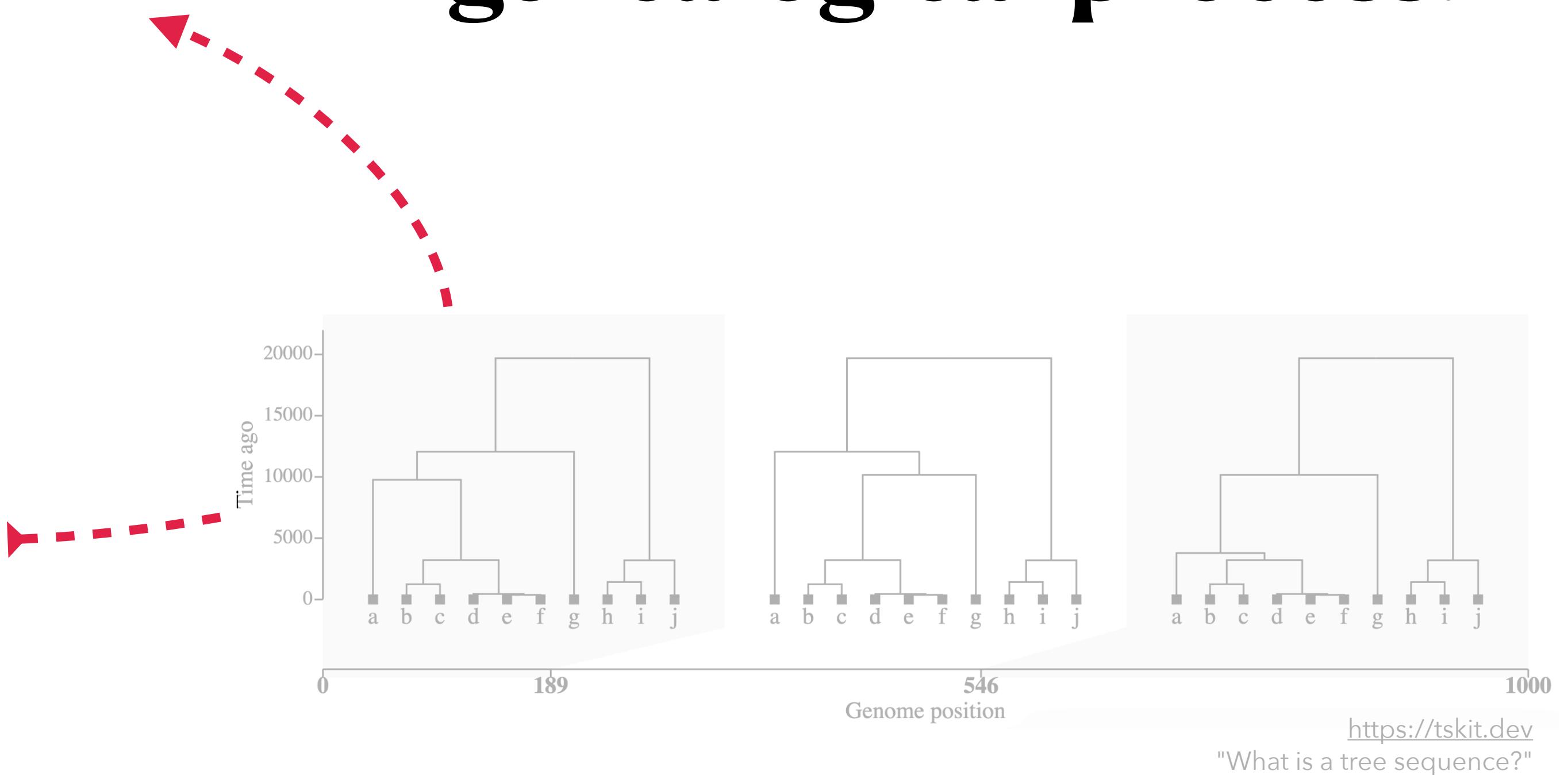


What we (usually) want

- estimate split times between populations
- investigate admixture hypotheses
- infer population sizes changes over time

These are all features
of a real (but hidden)
genealogical process.

#CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO	FORMAT	HG00099	HG00099	HG00101
12	60076	.	A	C	100	PASS	.	GT	1 0	0 0	0 0
12	60252	.	A	G	100	PASS	.	GT	0 0	0 0	0 1
12	60317	.	C	T	100	PASS	.	GT	0 0	0 0	0 0
12	60344	.	C	A	100	PASS	.	GT	0 0	0 0	0 0
12	60383	.	G	A	100	PASS	.	GT	0 0	0 0	0 0
12	60405	.	T	C	100	PASS	.	GT	0 0	0 0	0 0
12	60474	.	G	A	100	PASS	.	GT	0 0	0 1	0 1
12	60614	.	C	A	100	PASS	.	GT	0 0	0 0	0 0
12	60628	.	T	C	100	PASS	.	GT	0 0	0 0	0 0
12	60654	.	G	A	100	PASS	.	GT	0 0	0 0	0 0
12	61021	.	C	T	100	PASS	.	GT	0 0	0 0	0 0
12	61187	.	G	T	100	PASS	.	GT	0 0	0 0	0 0
12	61172	.	G	A	100	PASS	.	GT	0 0	0 0	0 0
12	61220	.	G	A	100	PASS	.	GT	0 0	0 1	0 1
12	61258	.	C	T	100	PASS	.	GT	0 0	0 0	0 1
12	61272	.	T	C	100	PASS	.	GT	0 0	0 0	0 0
12	61329	.	C	T	100	PASS	.	GT	0 0	0 0	0 0
12	61341	.	G	A	100	PASS	.	GT	0 0	0 1	0 1
12	61368	.	C	T	100	PASS	.	GT	0 0	0 1	0 1
12	61392	.	T	A	100	PASS	.	GT	0 0	0 0	0 0
12	61405	.	G	C	100	PASS	.	GT	0 0	0 0	0 0
12	61411	.	C	A	100	PASS	.	GT	0 0	0 0	0 0
12	61416	.	G	A	100	PASS	.	GT	0 0	0 0	0 0
12	61422	.	C	T	100	PASS	.	GT	0 0	0 0	0 0
12	61476	.	C	G	100	PASS	.	GT	0 0	0 0	0 0
12	61510	.	G	A	100	PASS	.	GT	0 0	0 0	0 0
12	61516	.	C	T	100	PASS	.	GT	0 0	0 0	0 0
12	61552	.	C	T	100	PASS	.	GT	0 0	0 0	0 0
12	61604	.	T	G	100	PASS	.	GT	0 0	0 0	0 0
12	61687	.	G	A	100	PASS	.	GT	1 0	0 1	0 1
12	61700	.	C	T	100	PASS	.	GT	0 0	0 0	0 0

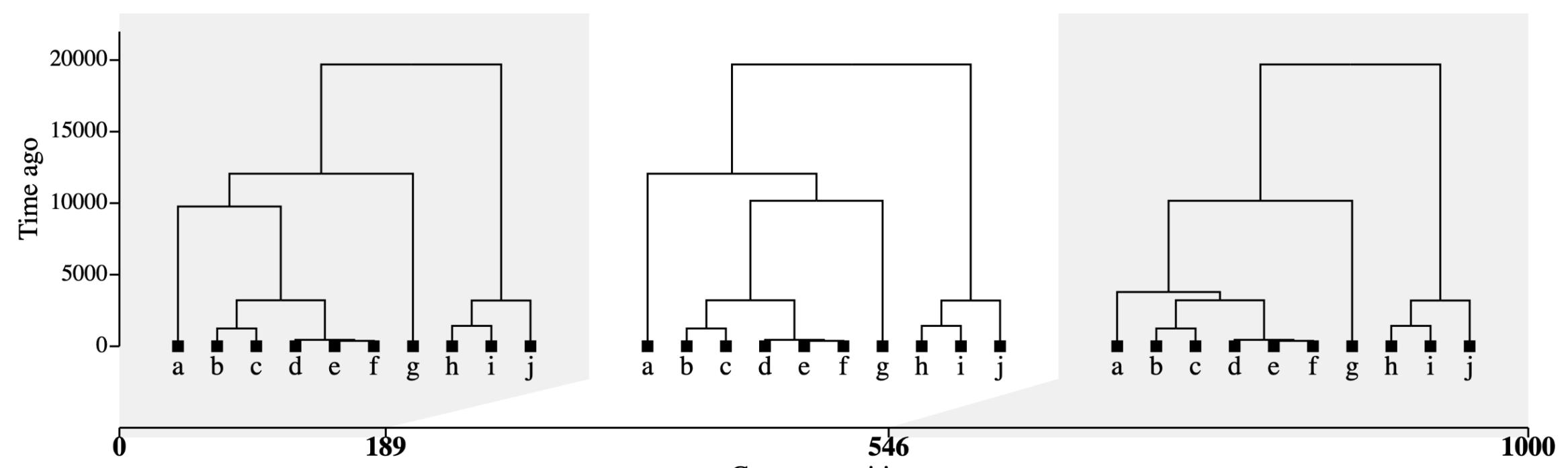
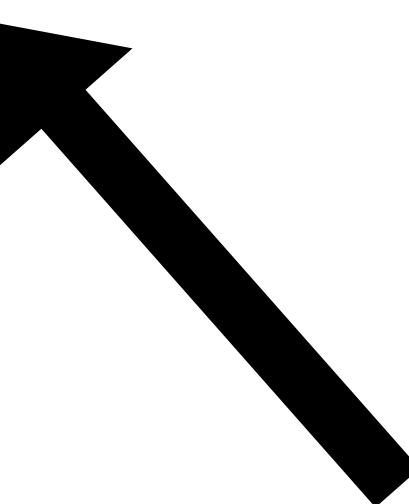


What we (usually) want

- estimate split times between populations
- investigate admixture hypotheses
- infer population sizes changes over time

If we have a tree sequence,
we get direct answers
almost "for free".

These are all features
of a real (but hidden)
genealogical process.



<https://tskit.dev>

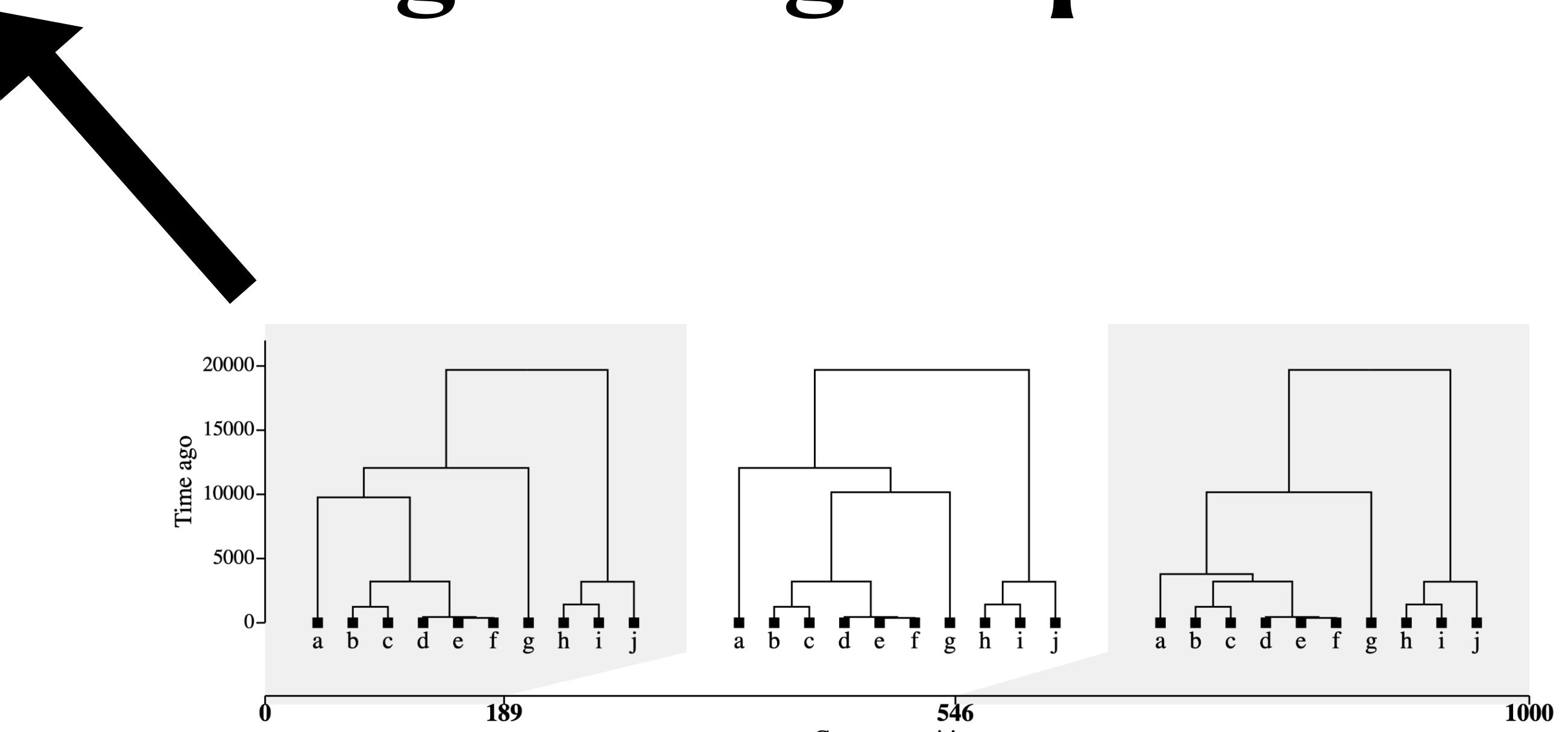
"What is a tree sequence?"

What we (usually) want

- estimate split times between populations
- investigate admixture hypotheses
- infer population sizes changes over time

If we have a tree sequence,
we get direct answers
almost "for free".

These are all features
of a real (but hidden)
genealogical process.



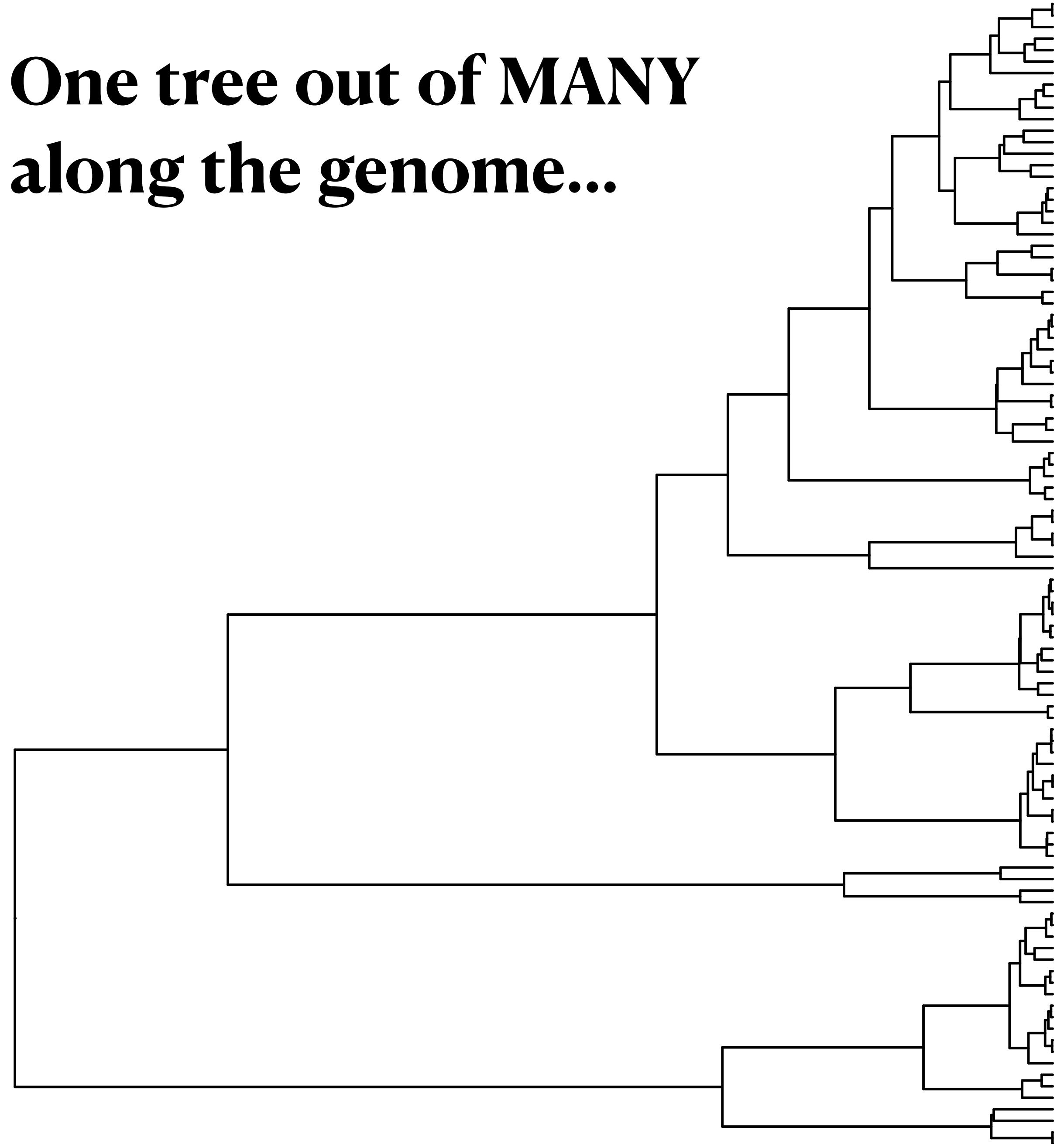
<https://tskit.dev>
"What is a tree sequence?"

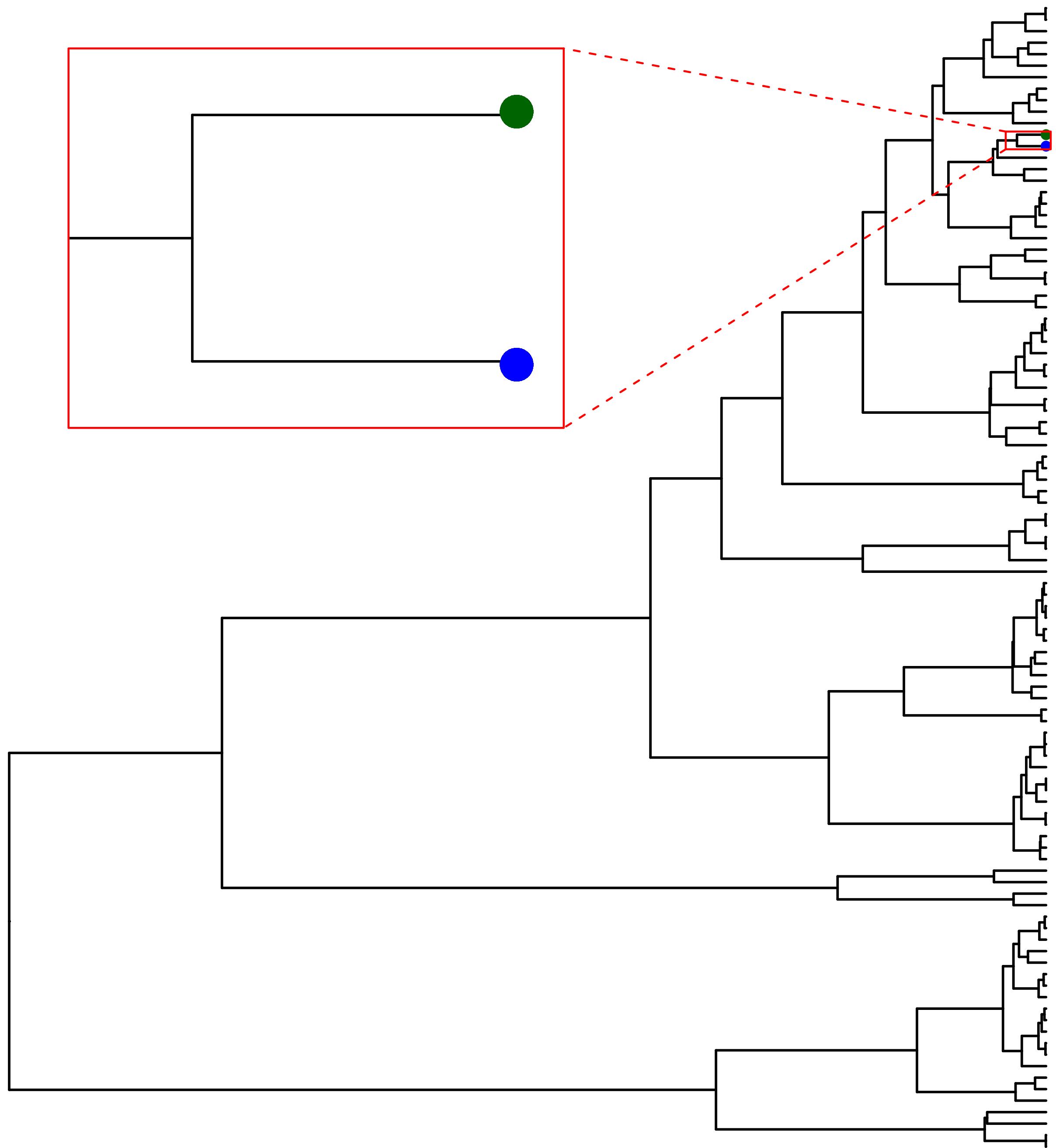
[... although the tree sequence must first be inferred from real data.]

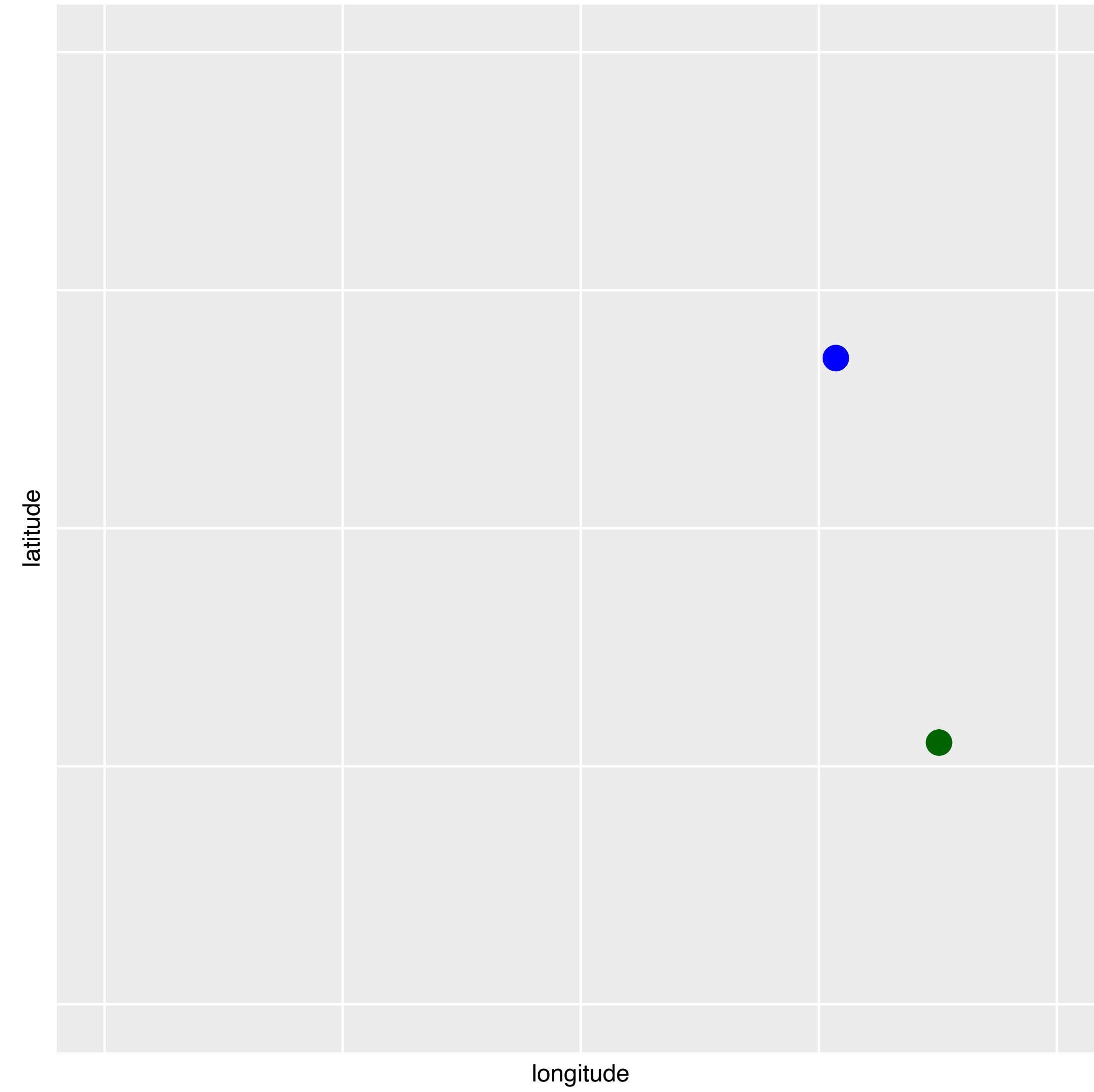
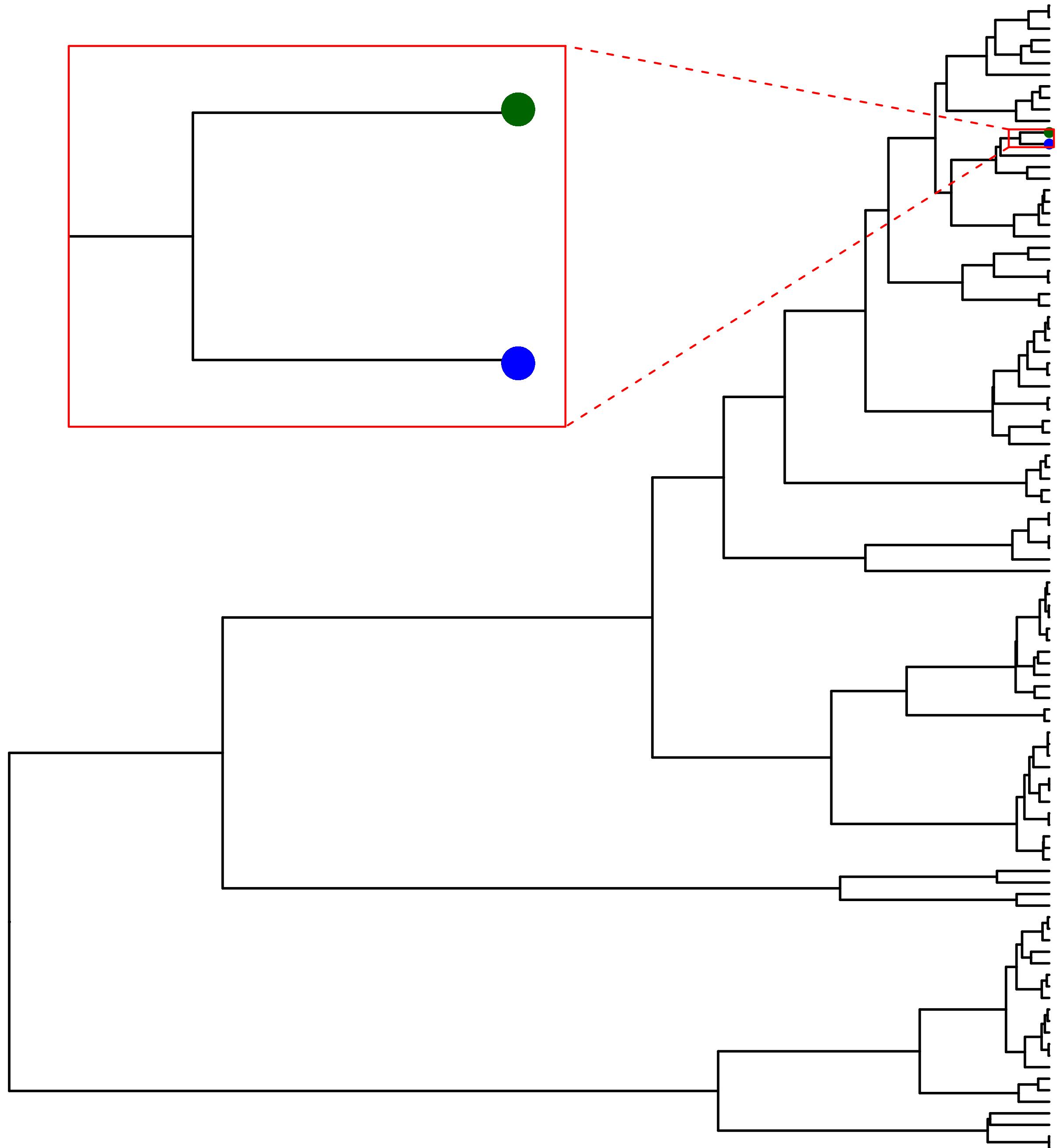
**Tree sequences will
push spatial inference
to the next level.**

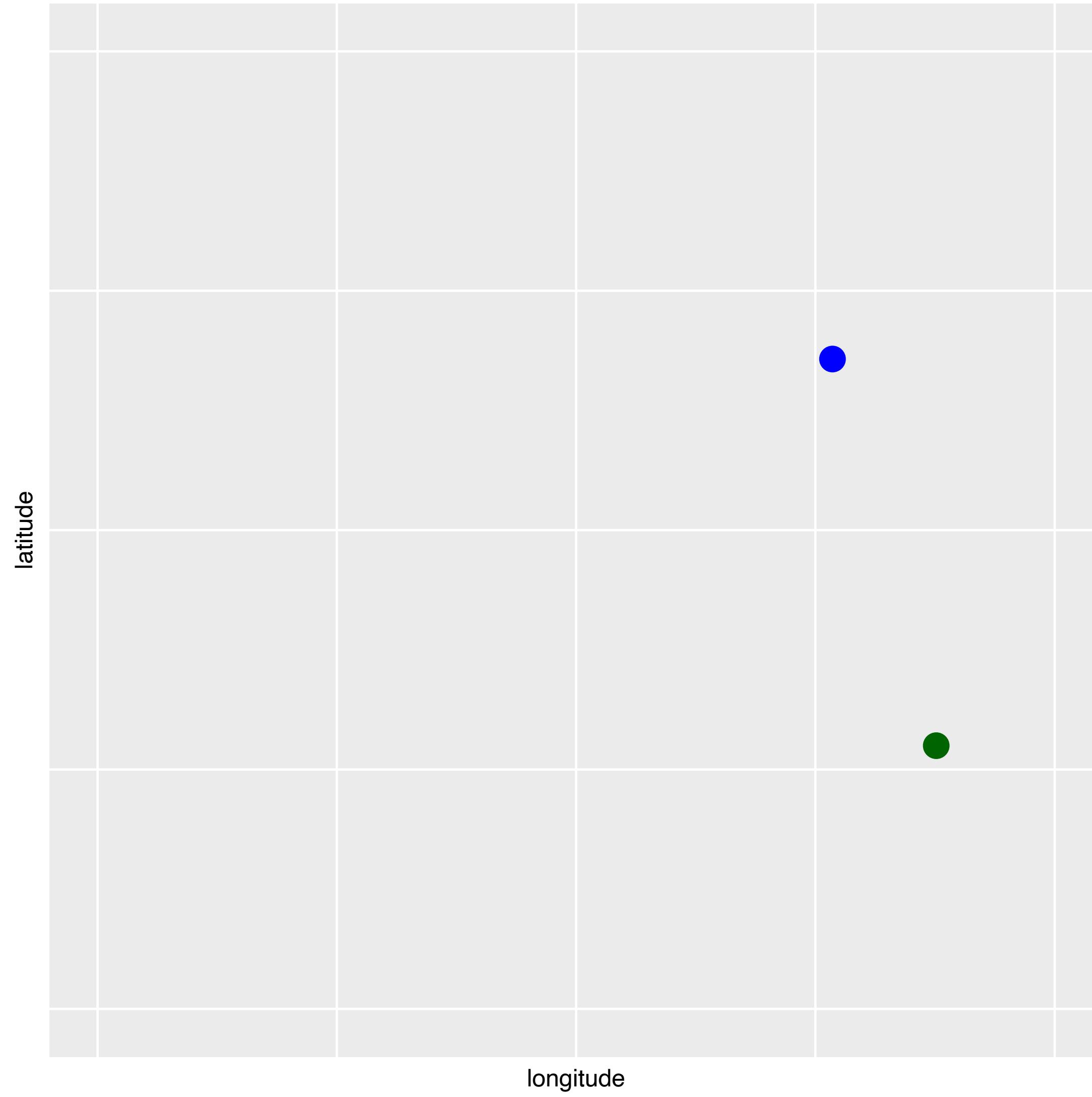
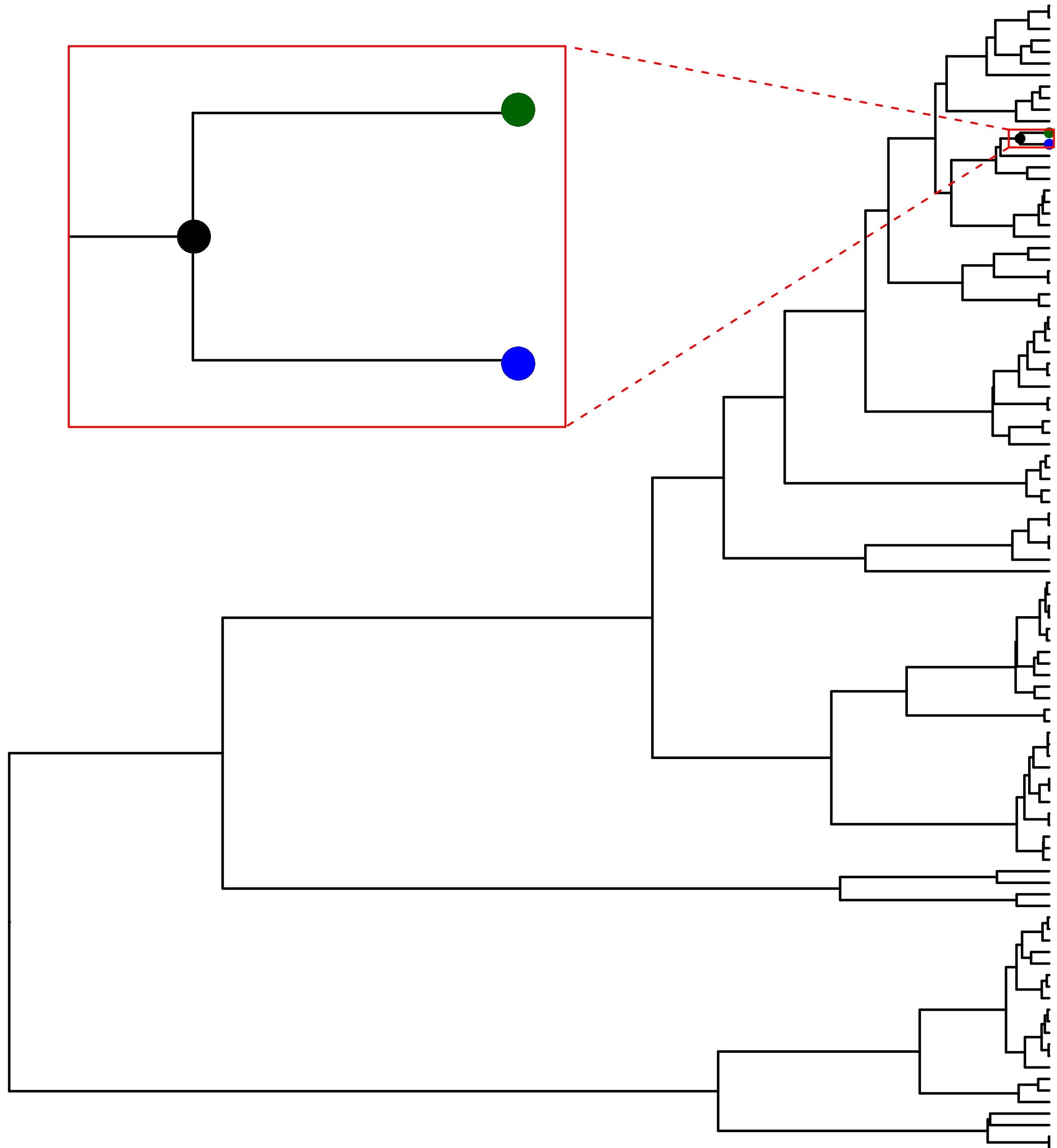


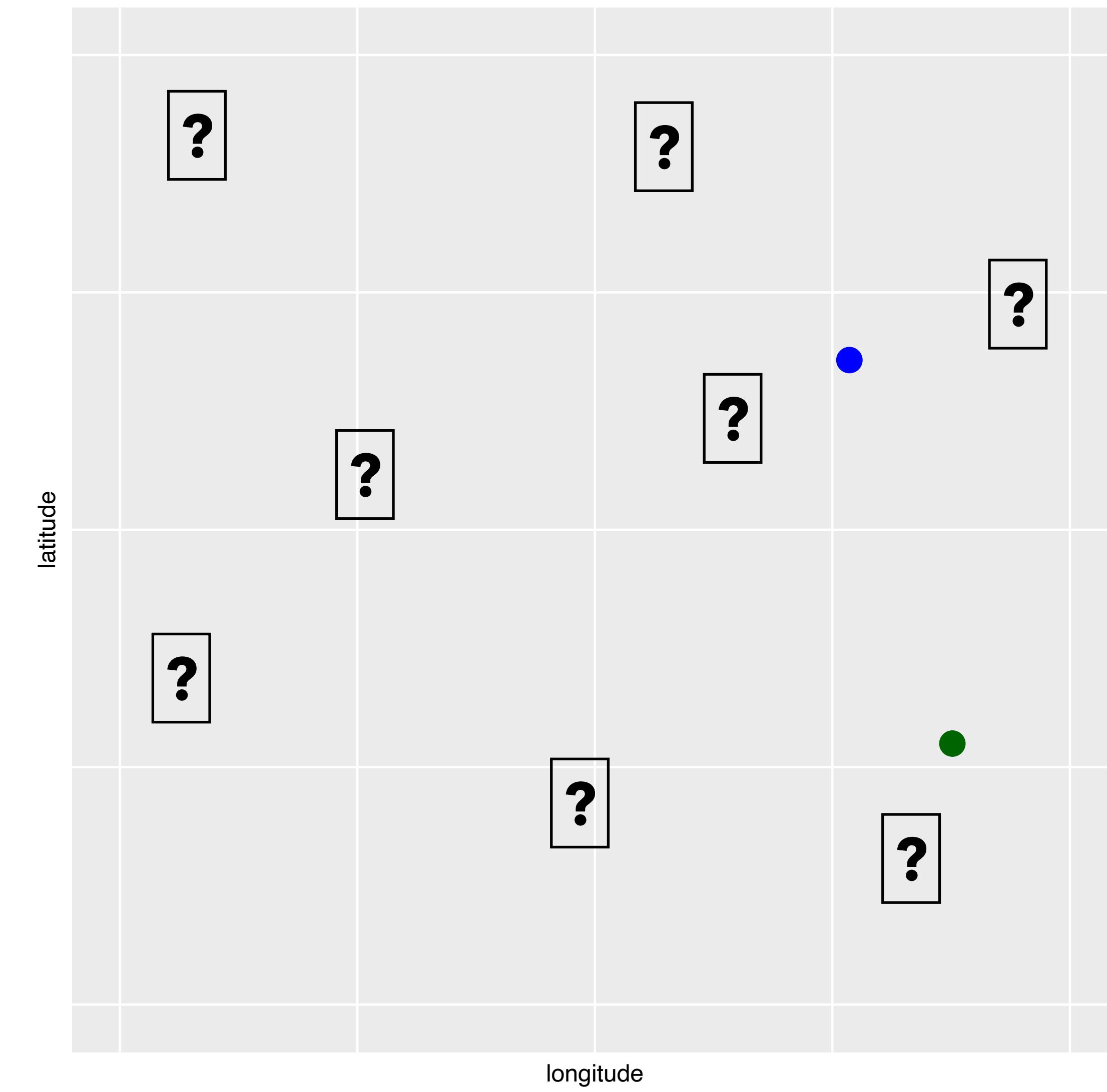
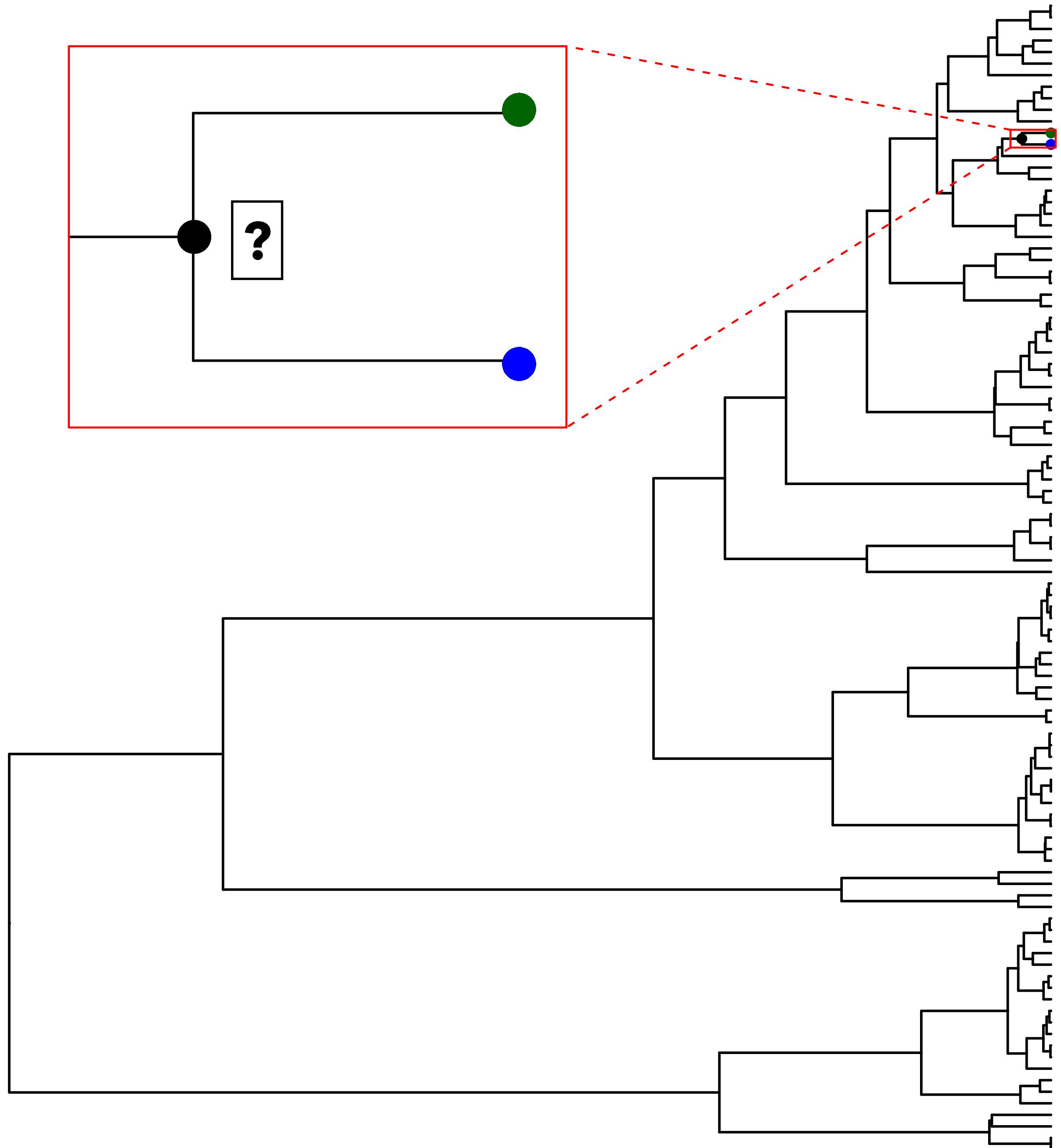
**One tree out of MANY
along the genome...**

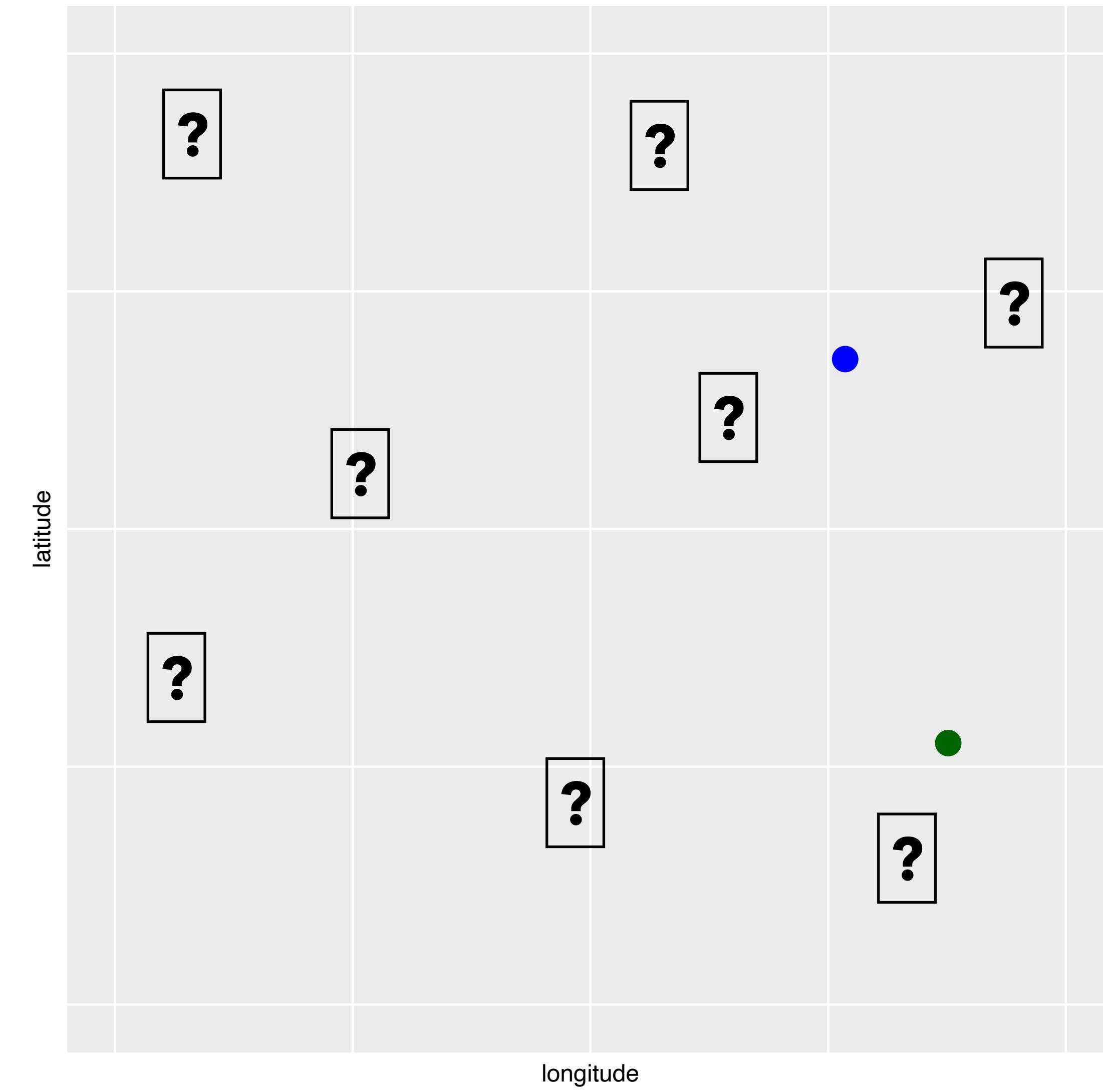
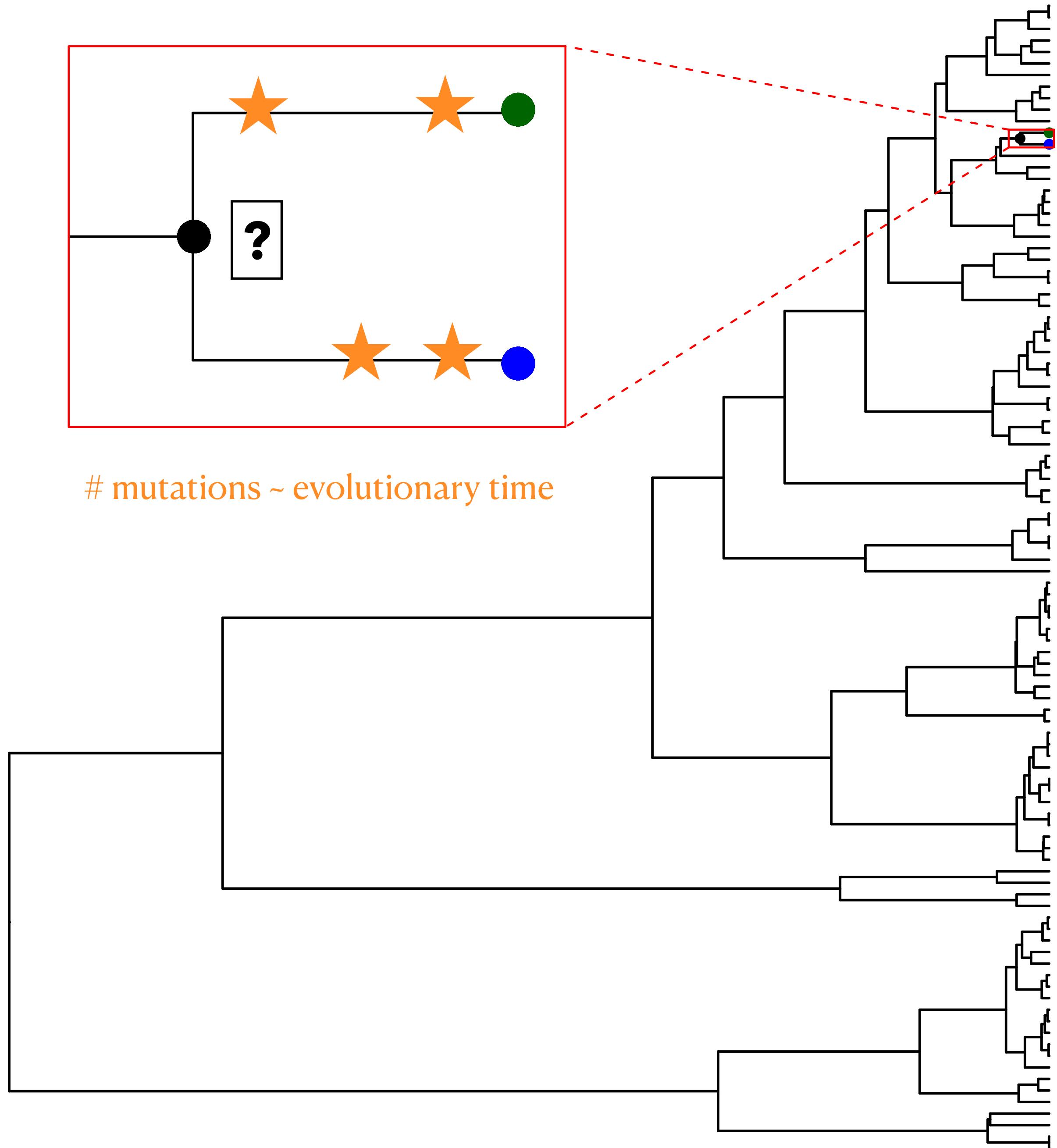


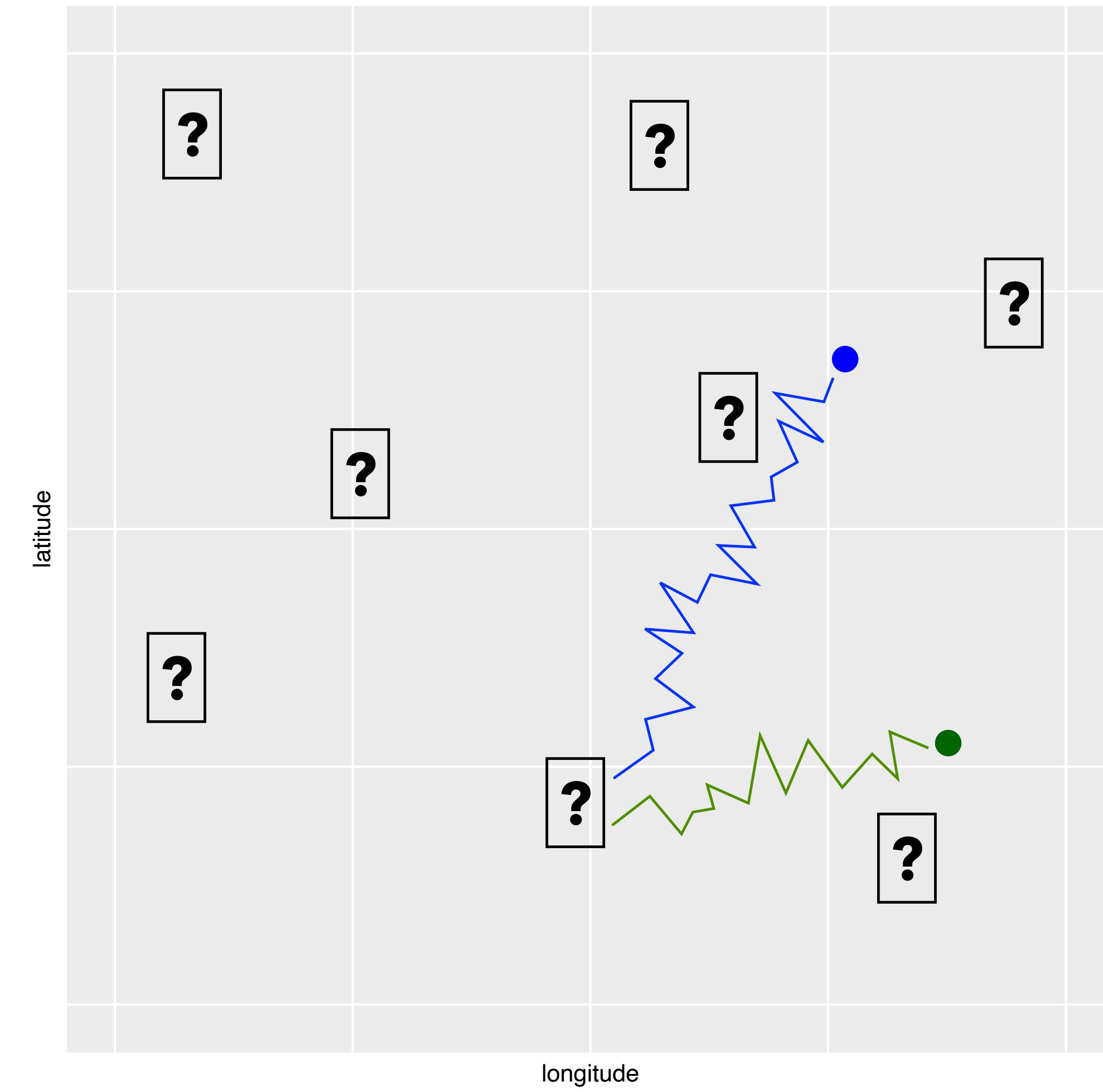
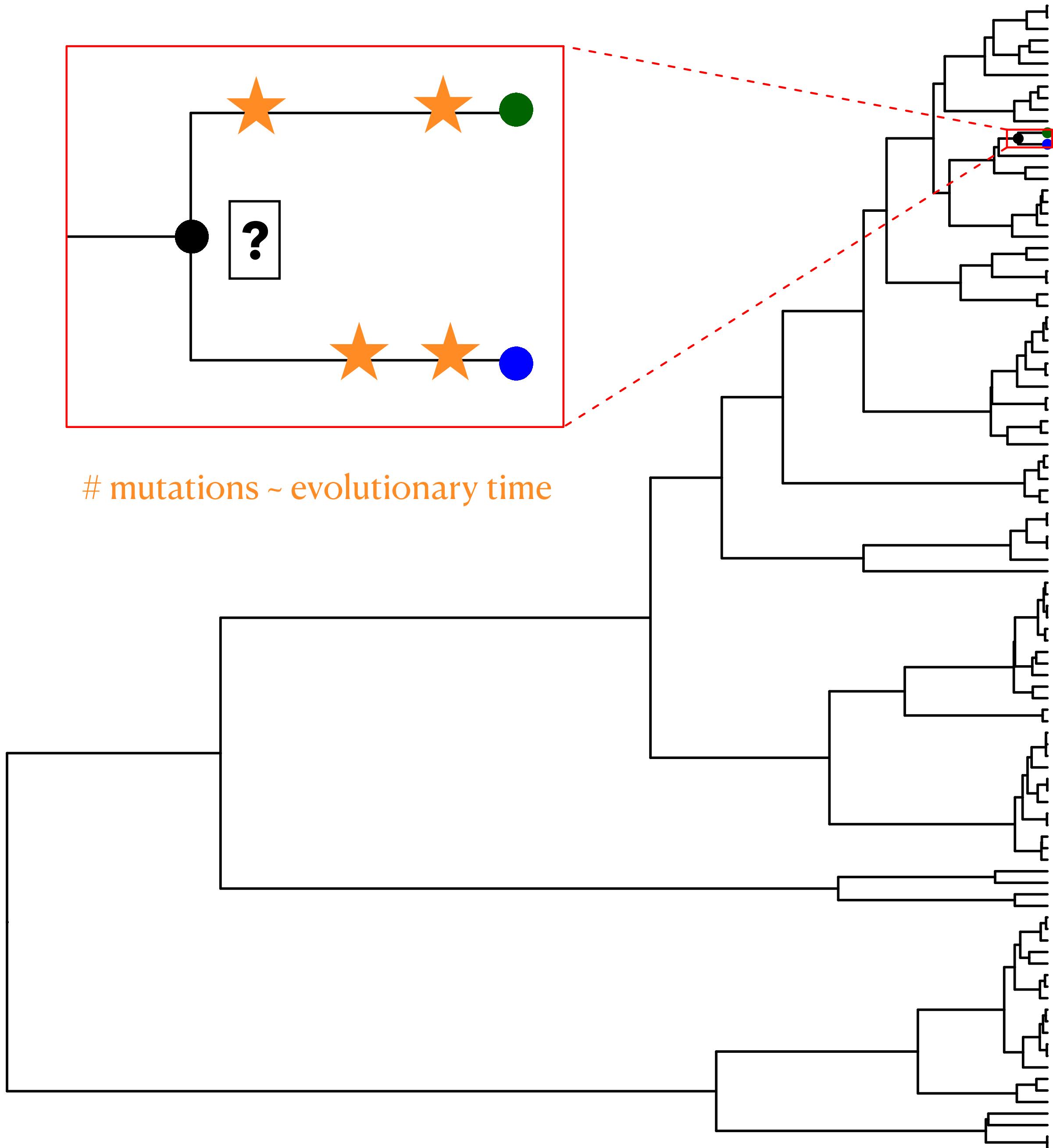


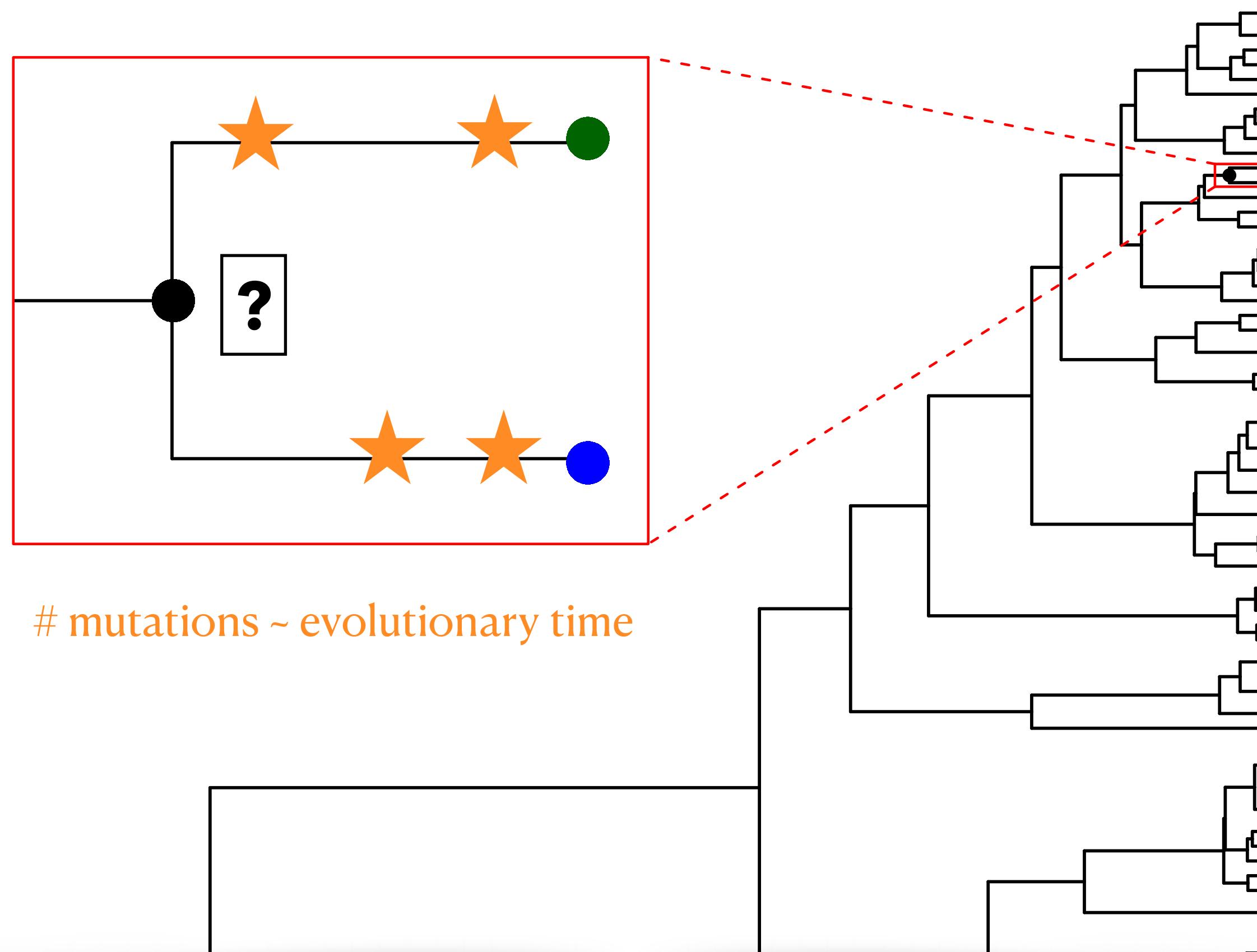










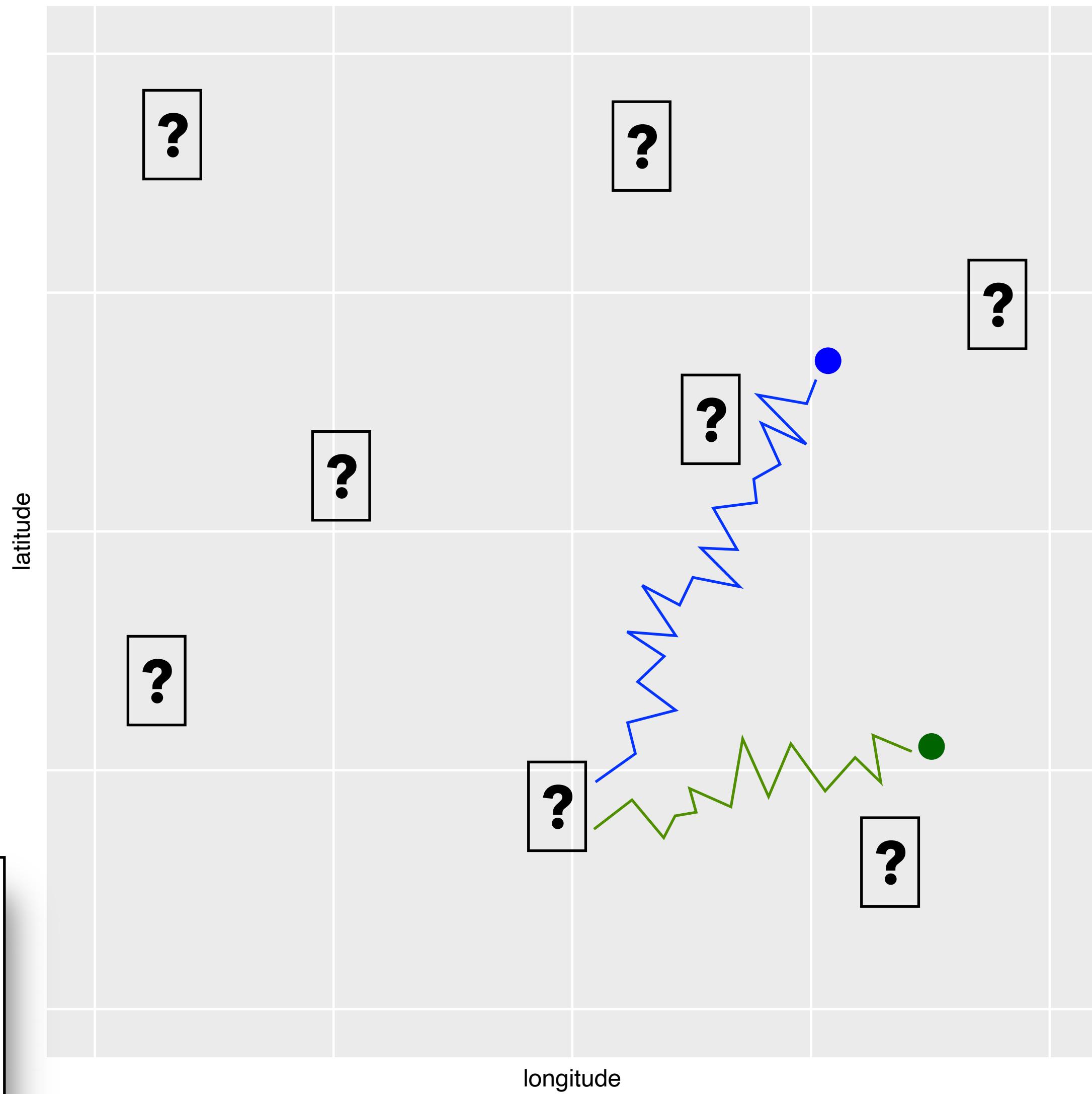


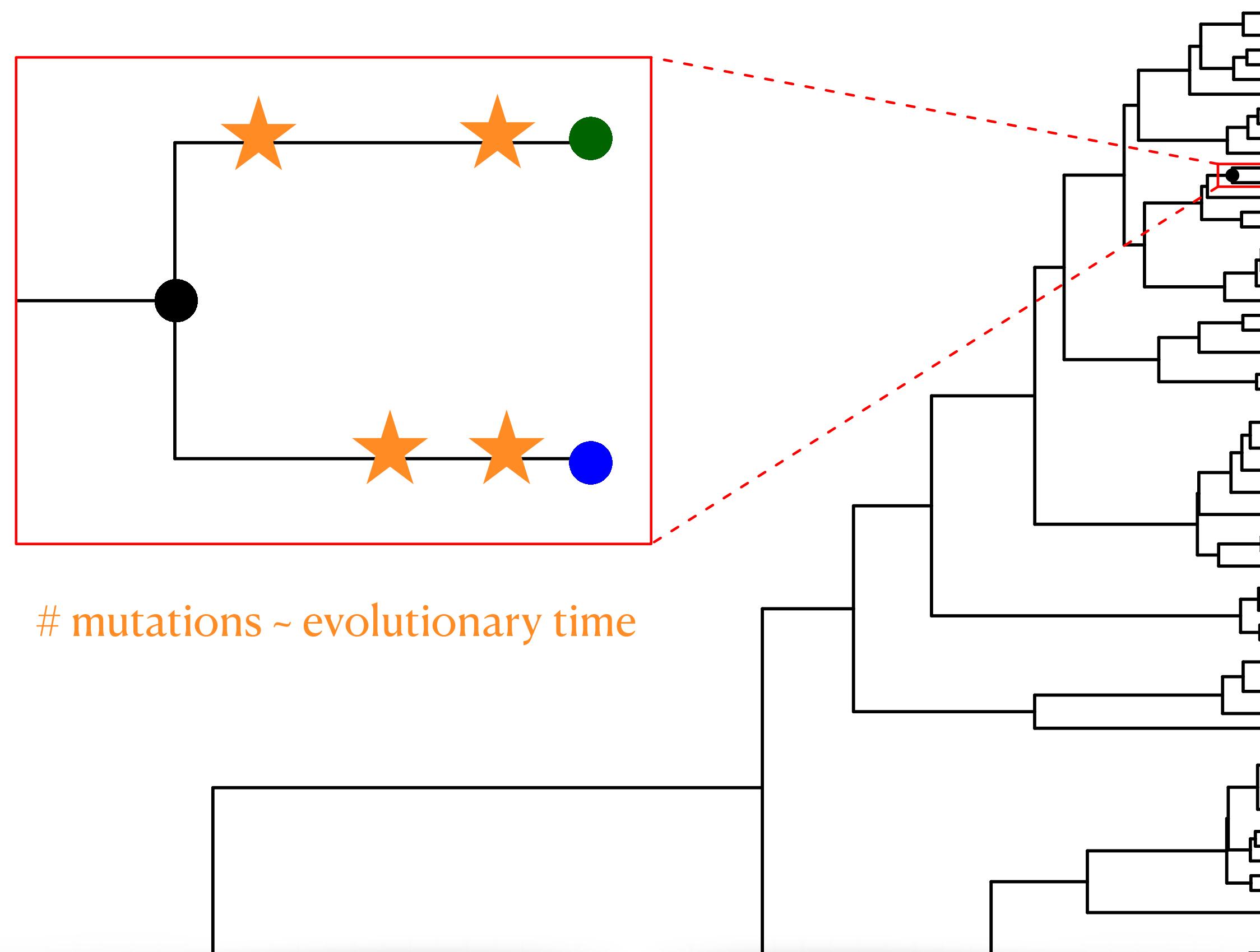
Bayesian model:

"What is the joint posterior distribution of ancestral positions & ages

given known locations of samples and mutation counts on branches?"

MCMC using NumPyro 



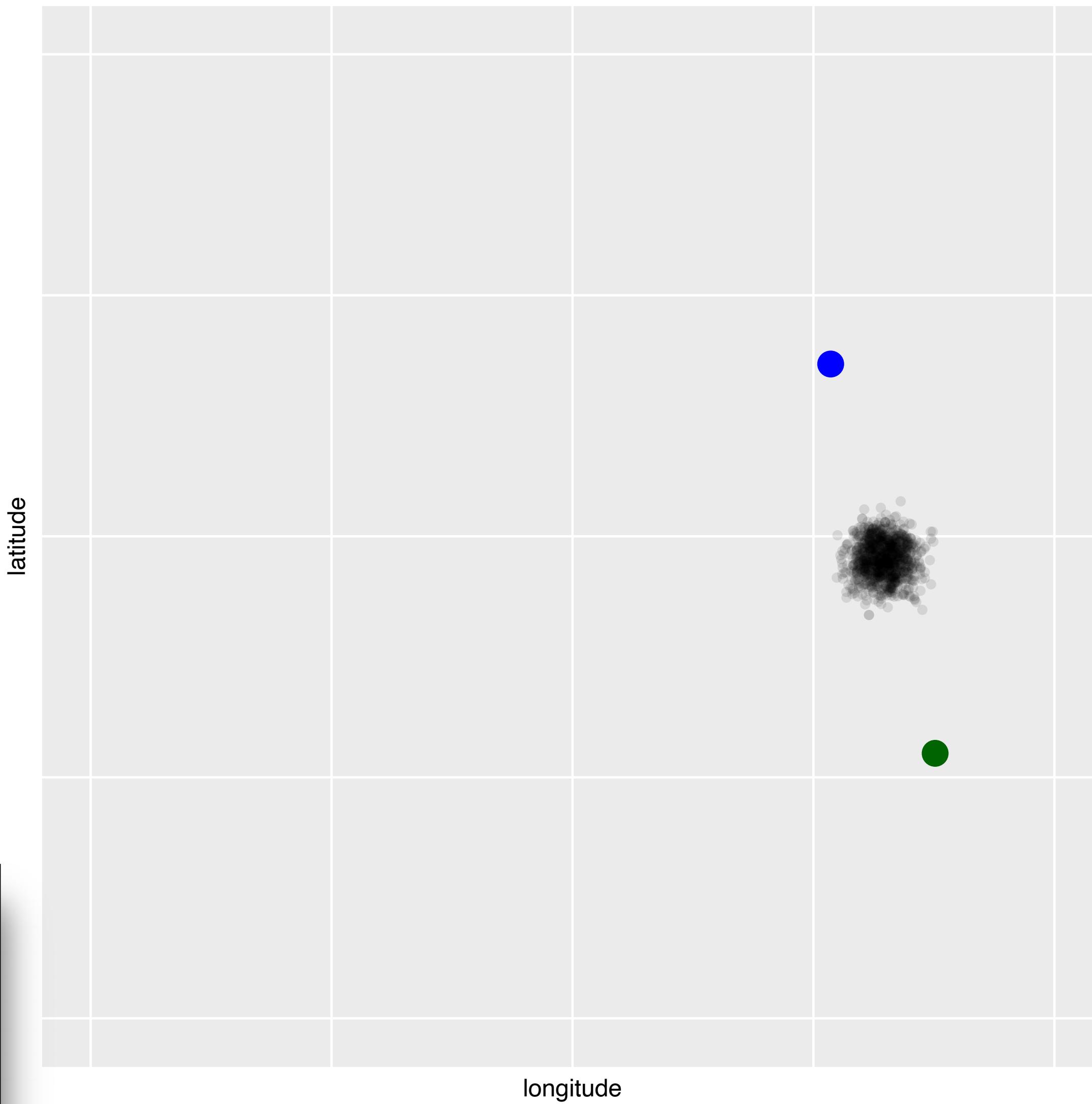


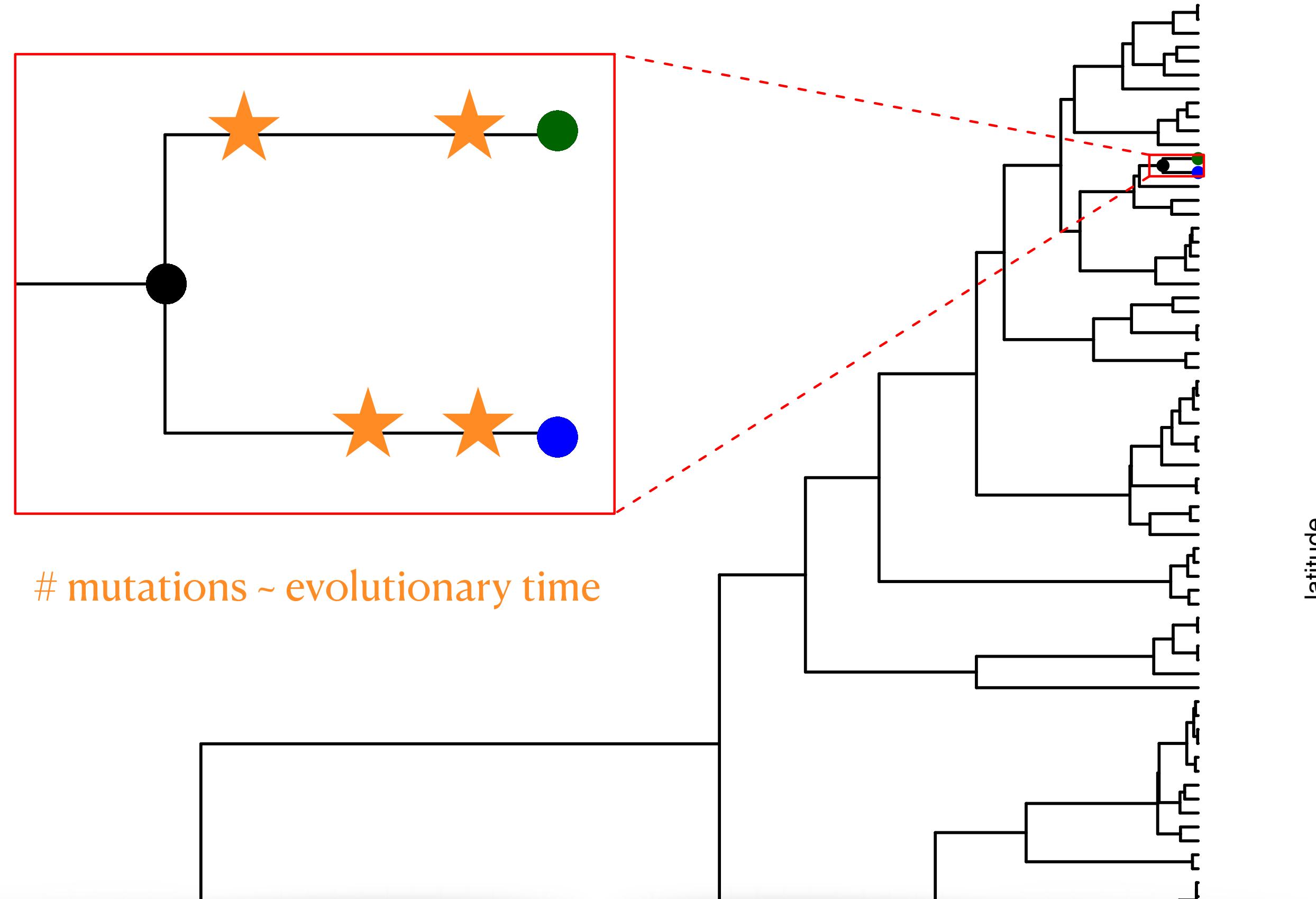
Bayesian model:

"What is the joint posterior distribution of ancestral positions & ages

given known locations of samples and mutation counts on branches?"

MCMC using NumPyro 



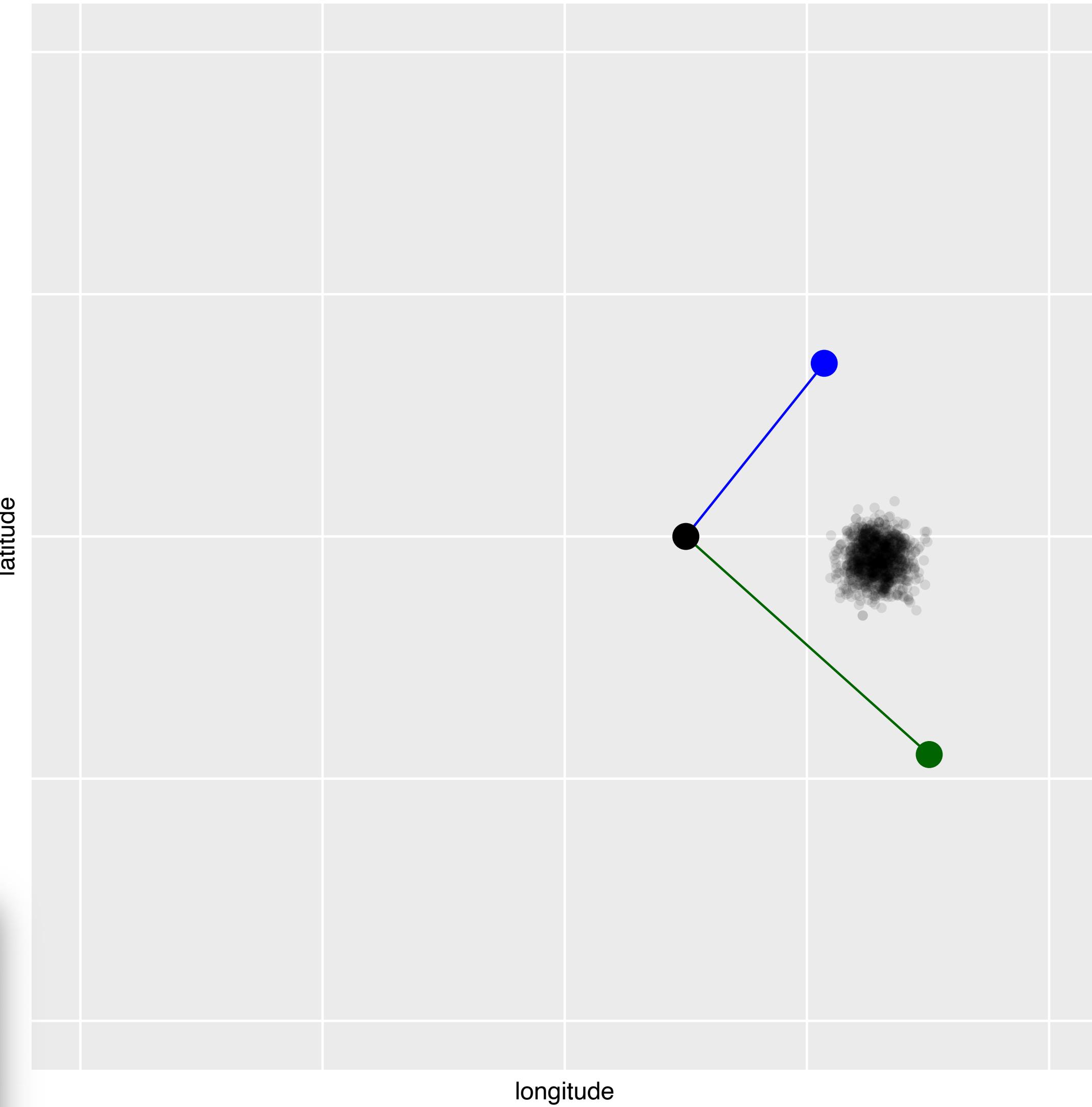


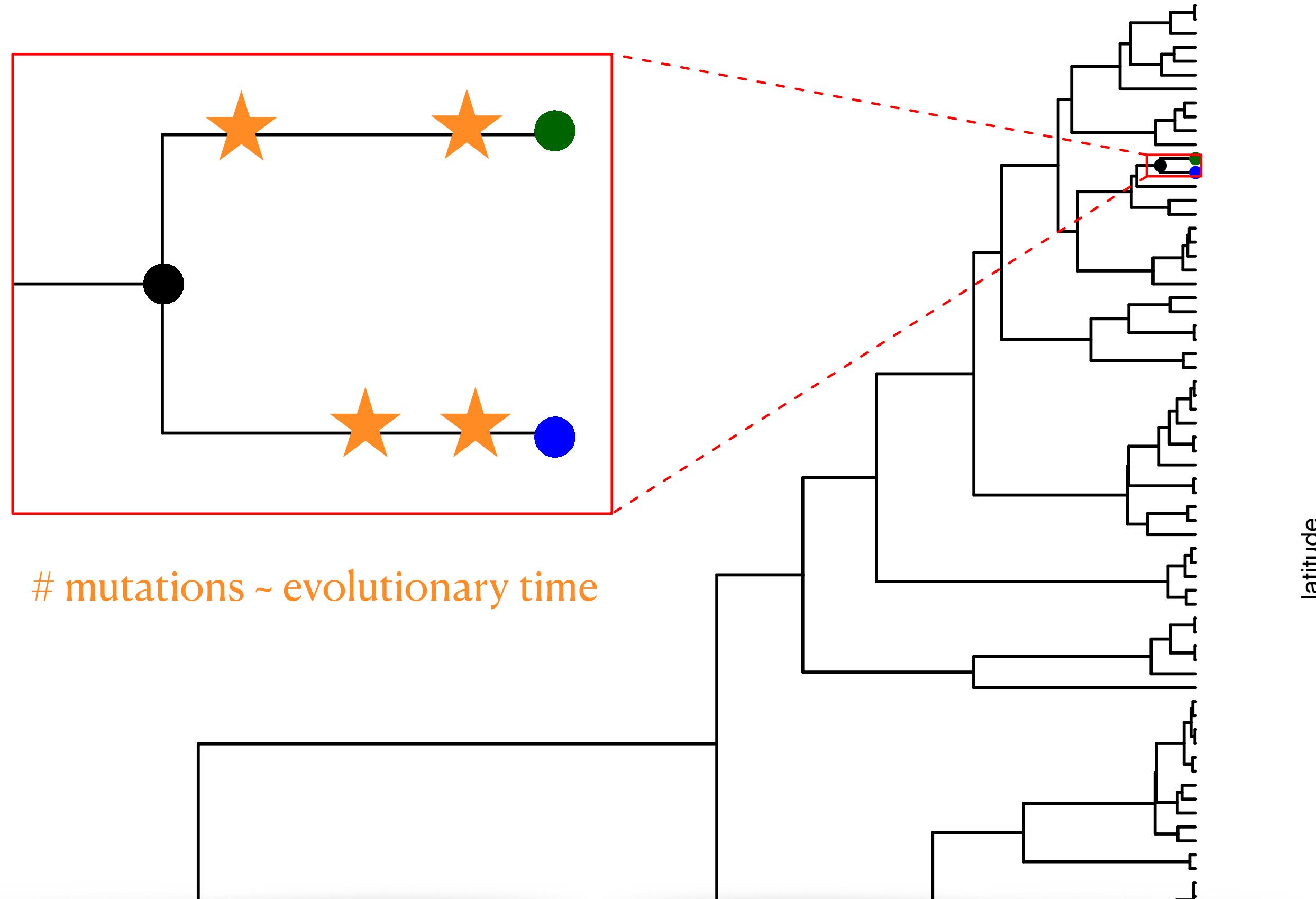
Bayesian model:

"What is the joint posterior distribution of ancestral positions & ages

given known locations of samples and mutation counts on branches?"

MCMC using NumPyro 



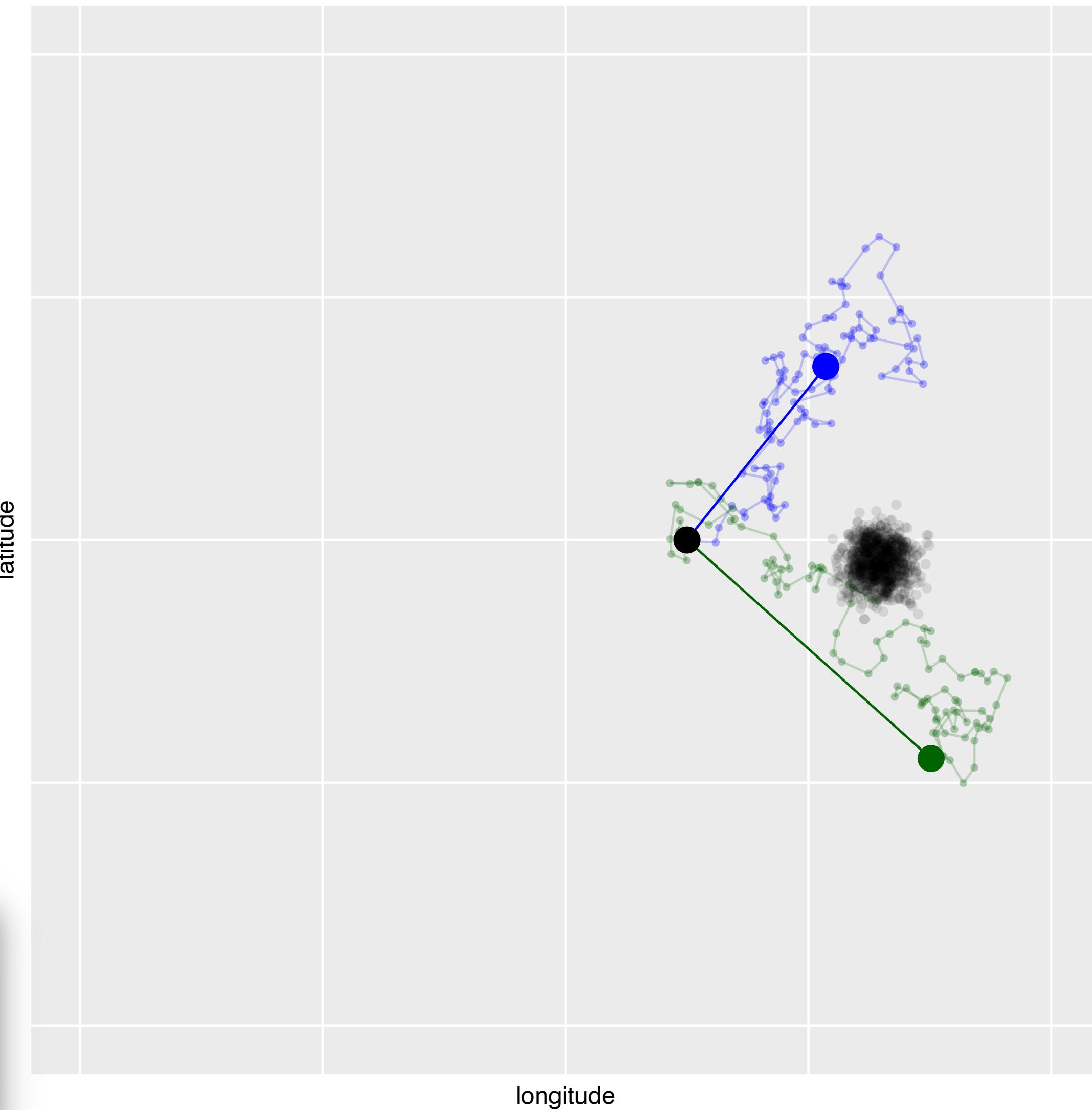


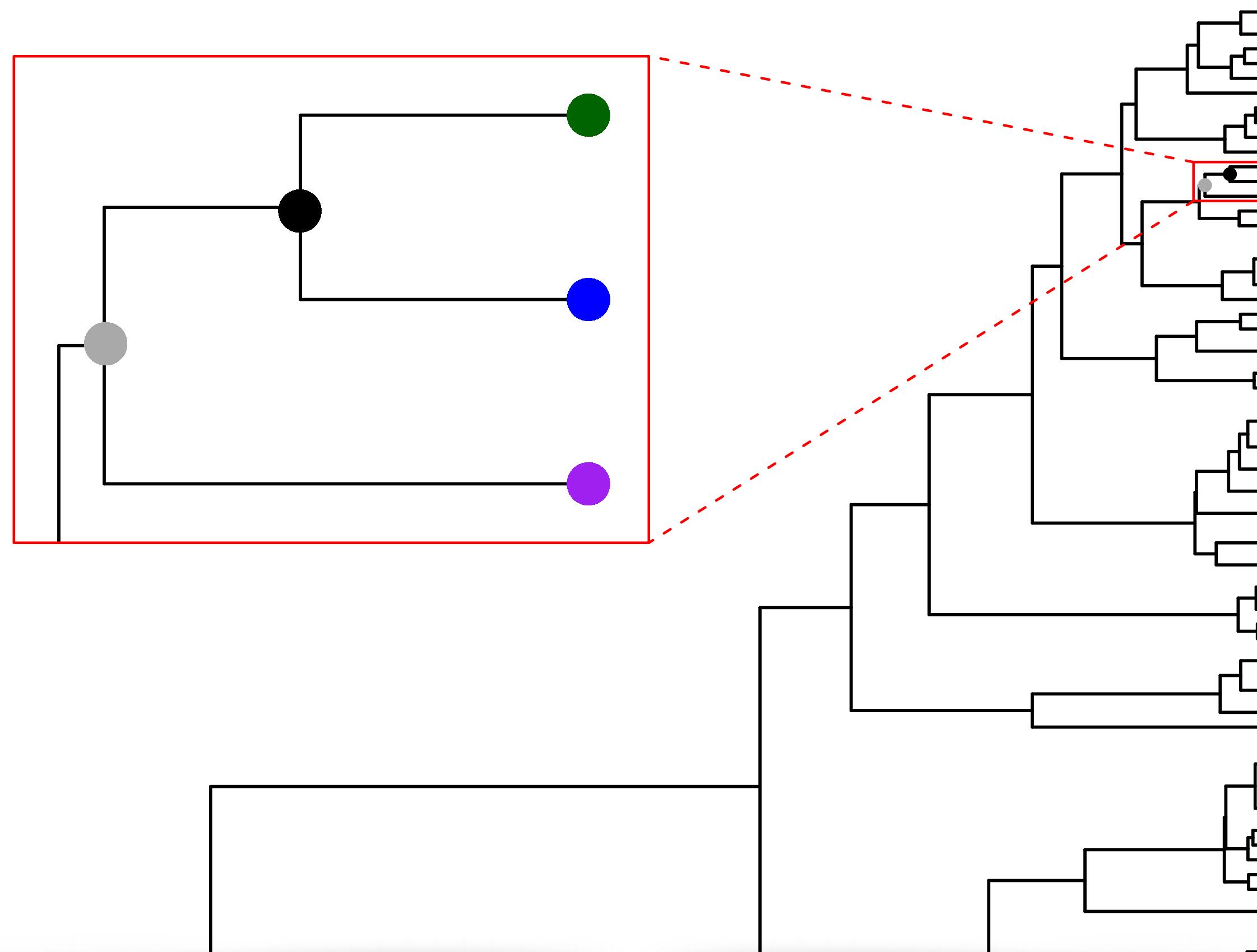
Bayesian model:

"What is the joint posterior distribution of ancestral positions & ages

given known locations of samples and mutation counts on branches?"

MCMC using NumPyro 



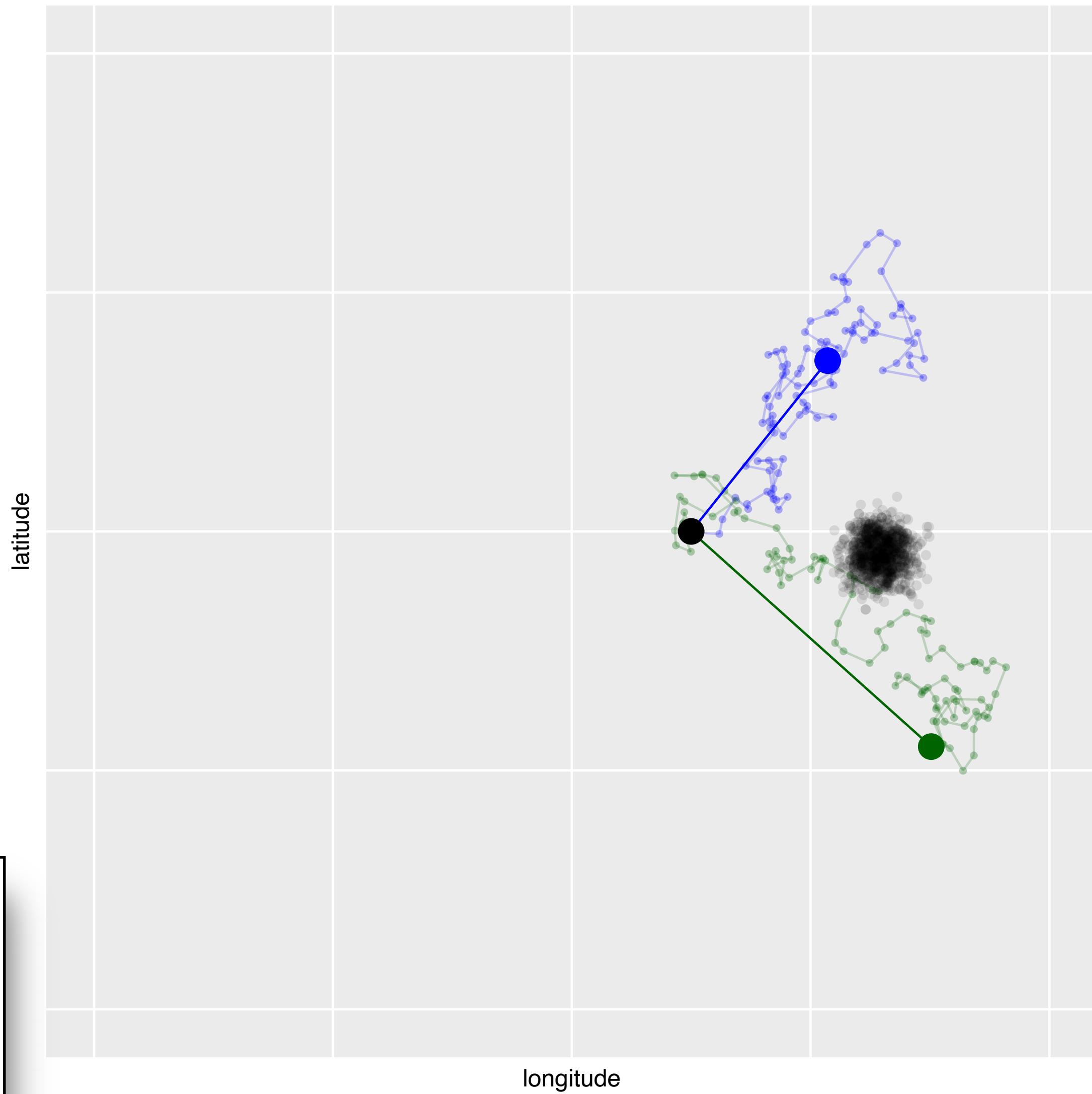


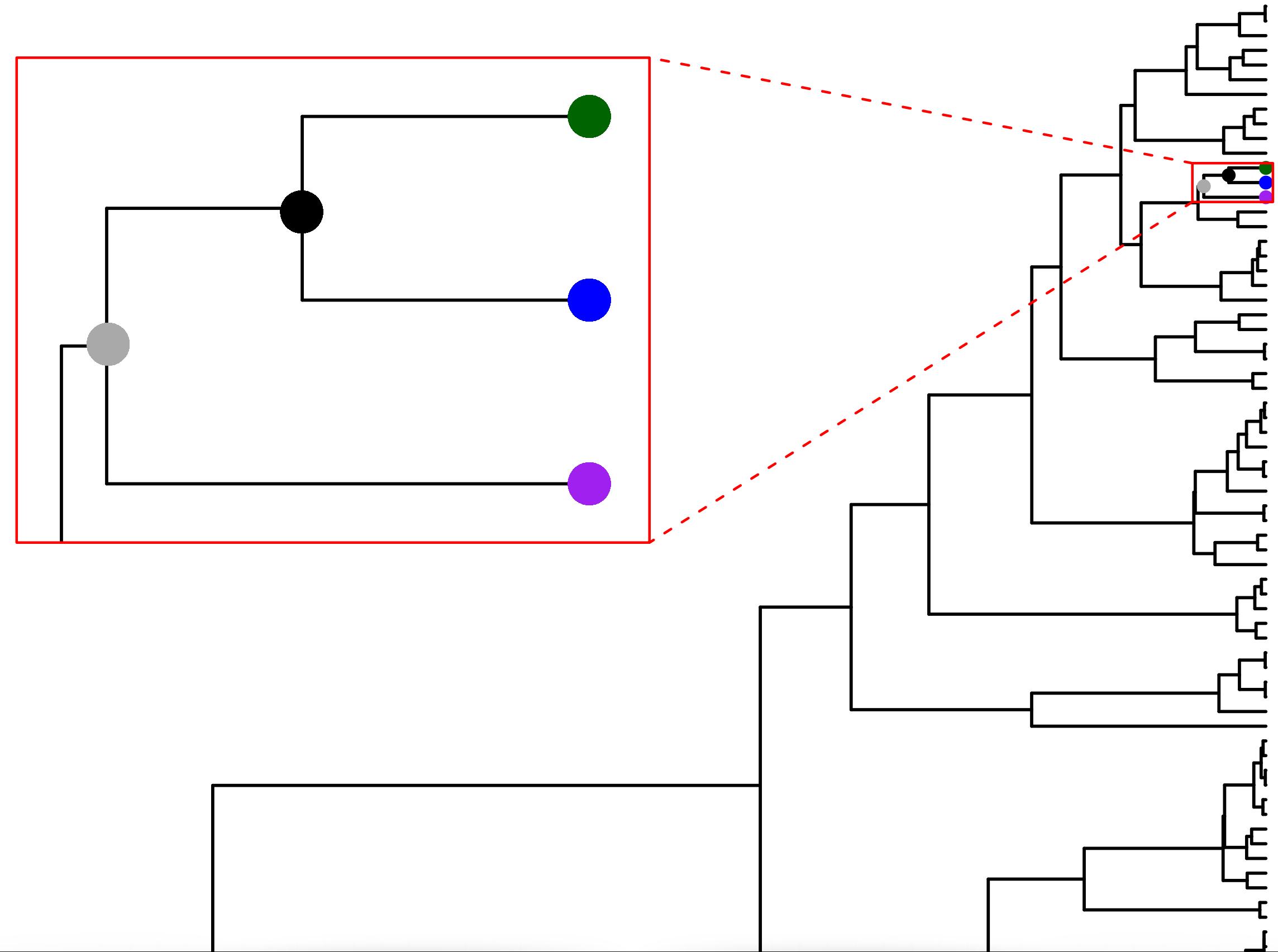
Bayesian model:

"What is the joint posterior distribution of ancestral positions & ages

given known locations of samples and mutation counts on branches?"

MCMC using NumPyro 



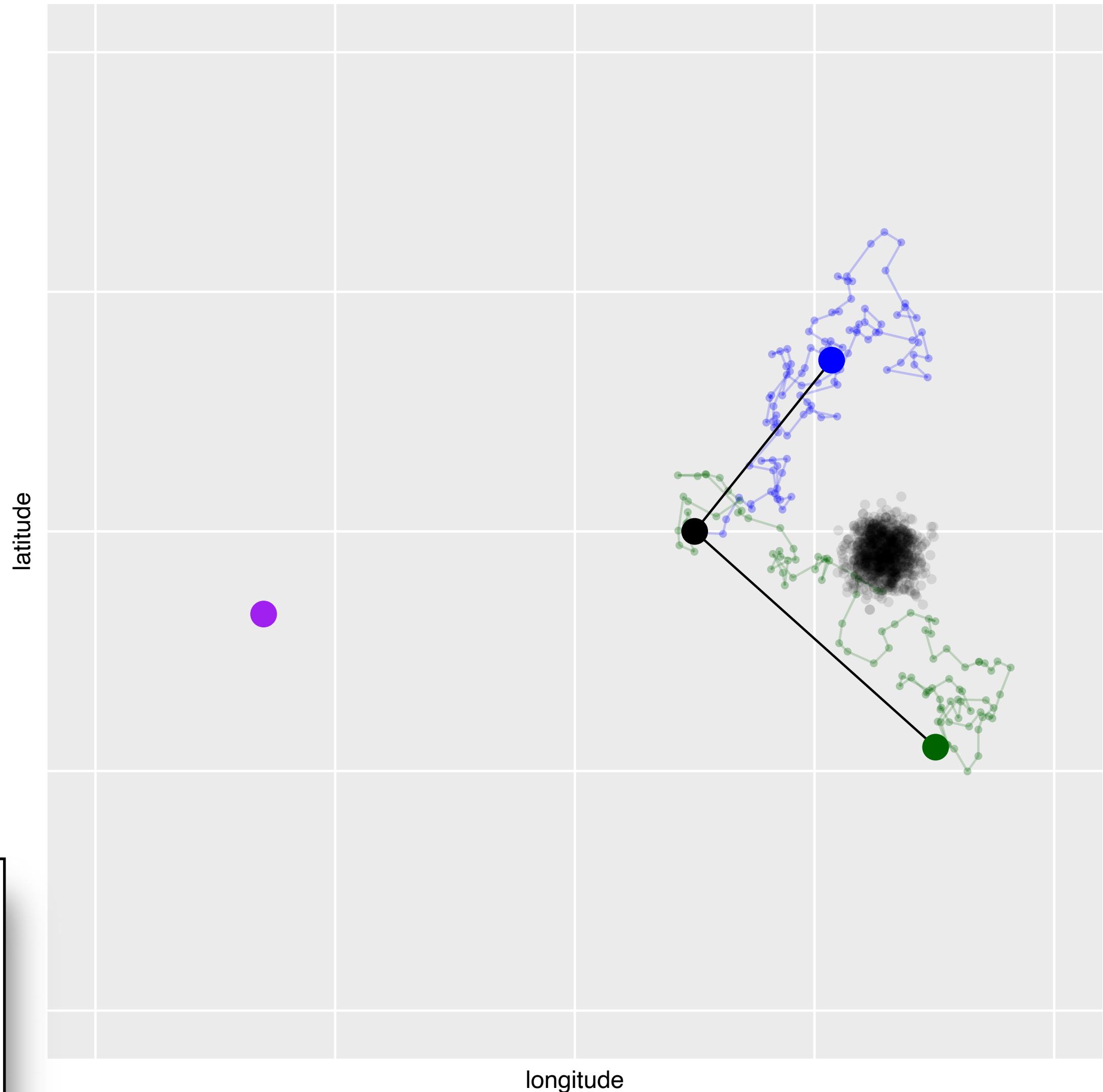


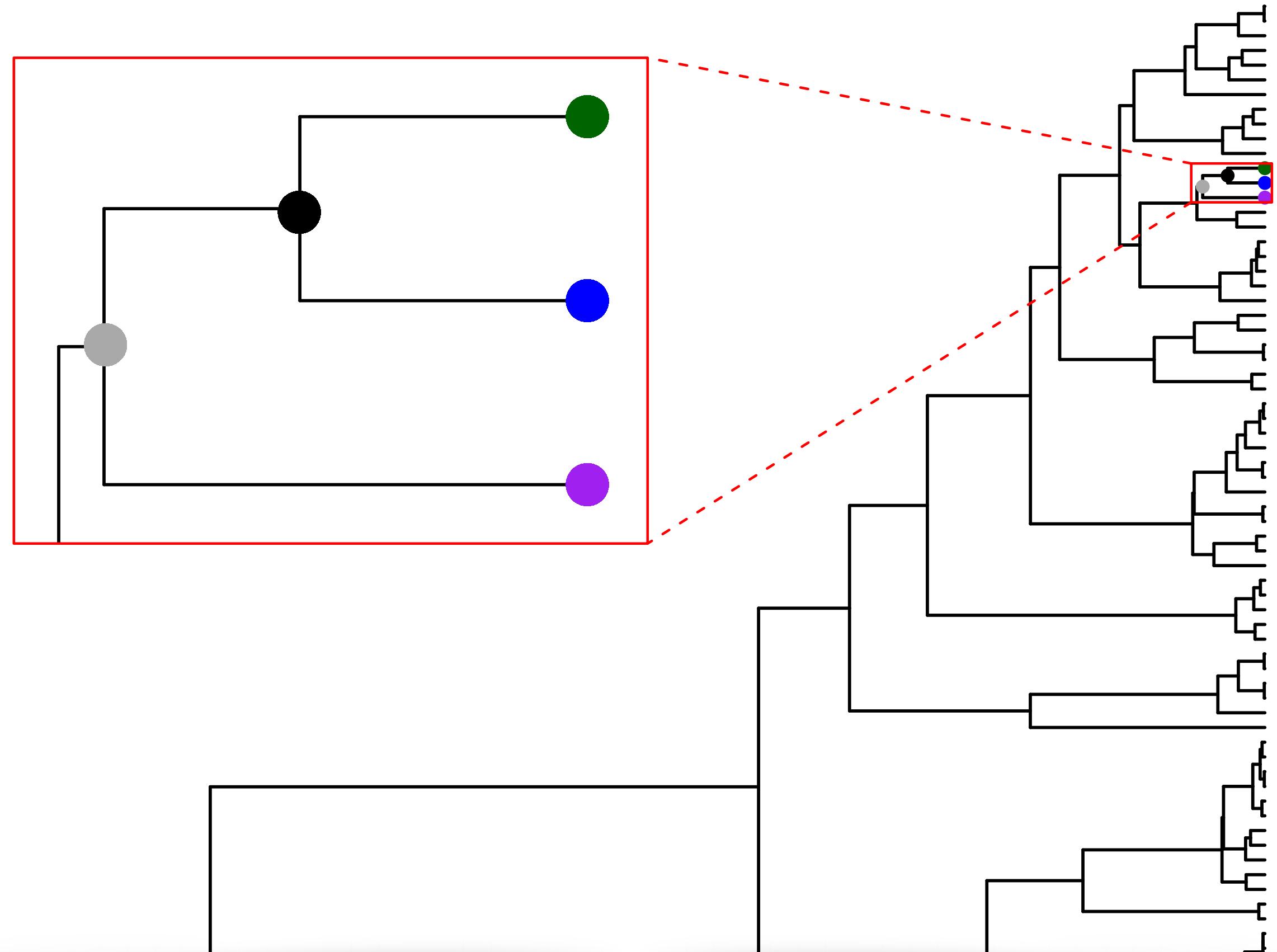
Bayesian model:

"What is the joint posterior distribution of ancestral positions & ages

given known locations of samples and mutation counts on branches?"

MCMC using NumPyro 



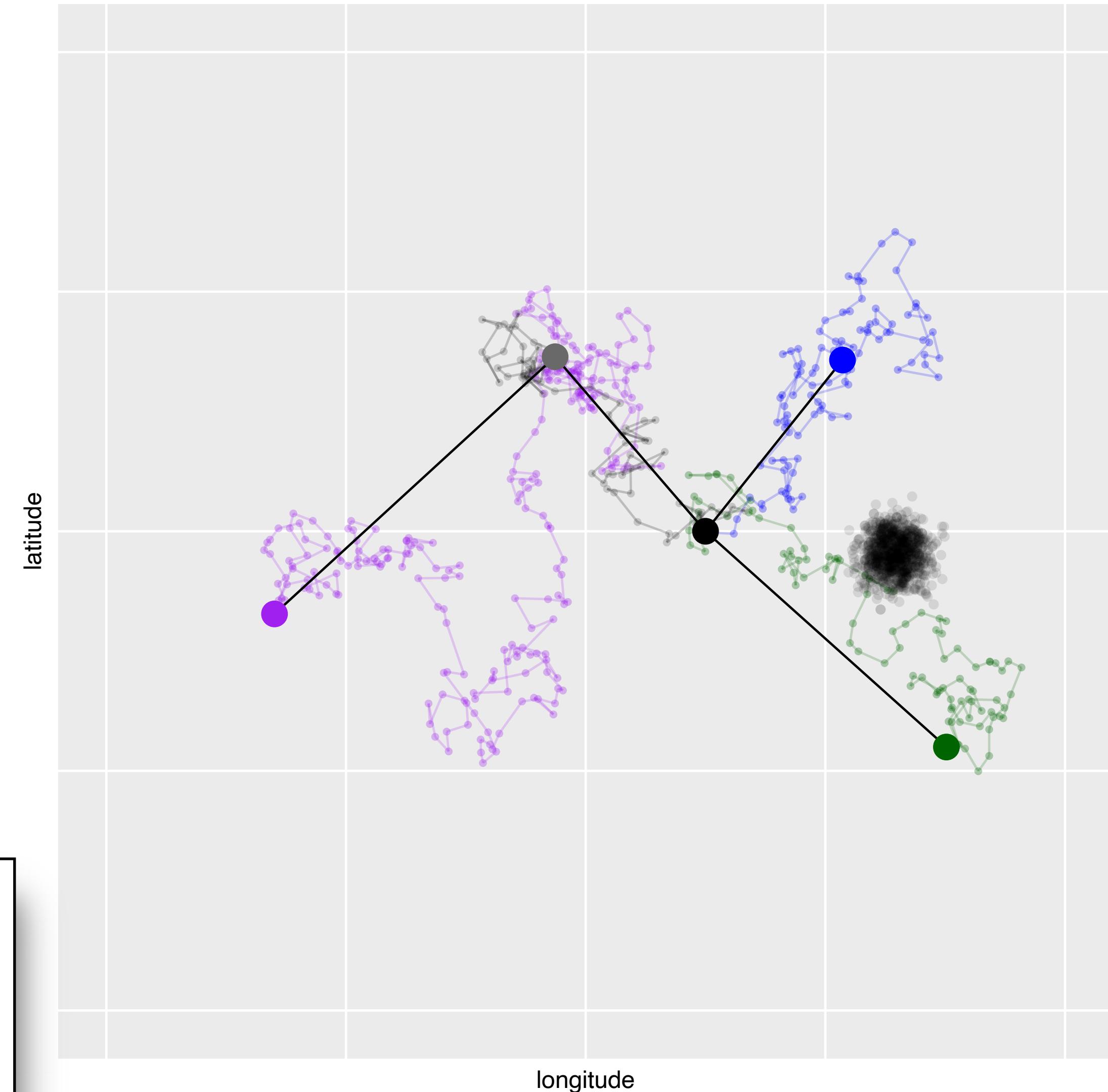


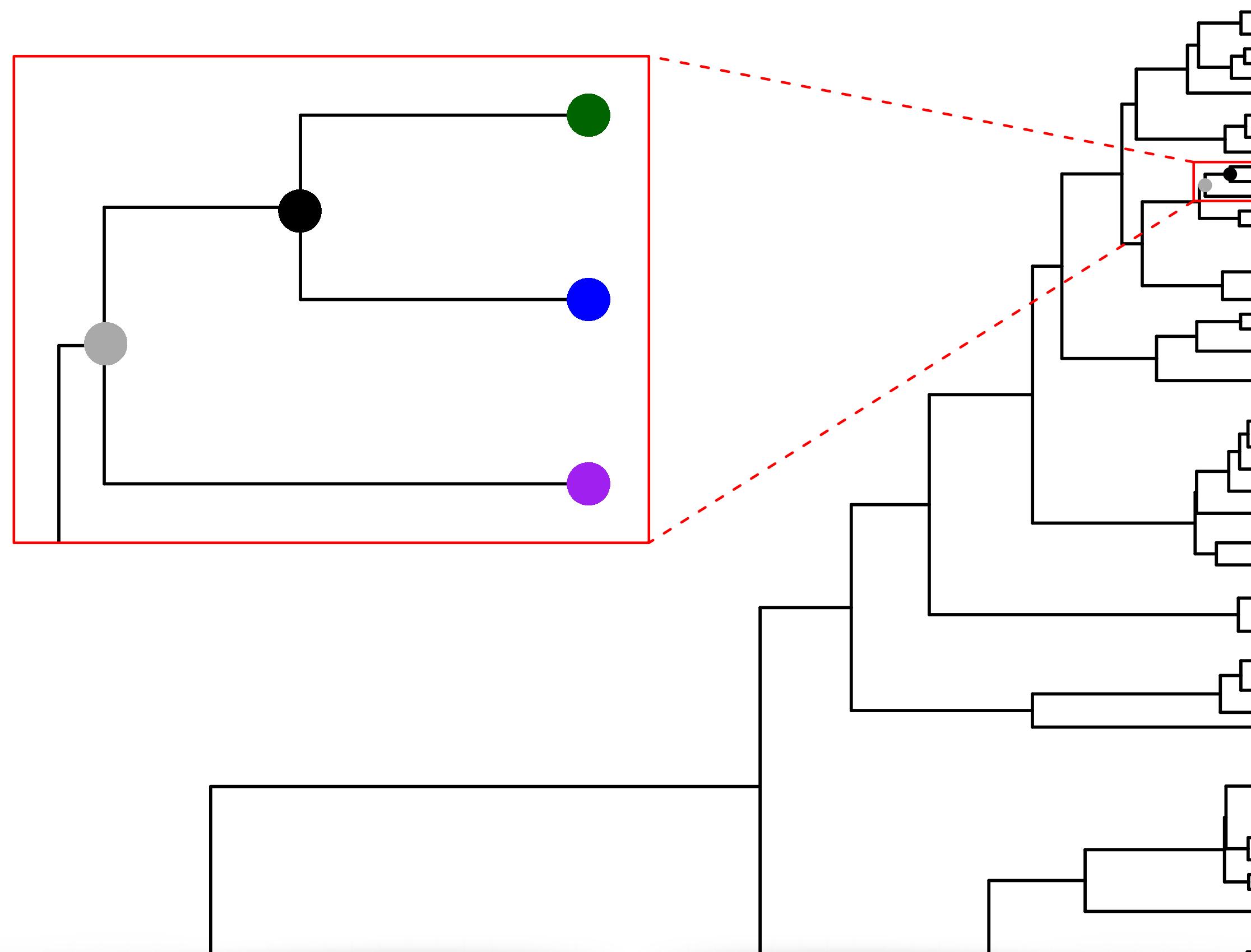
Bayesian model:

"What is the joint posterior distribution of ancestral positions & ages

given known locations of samples and mutation counts on branches?"

MCMC using NumPyro 



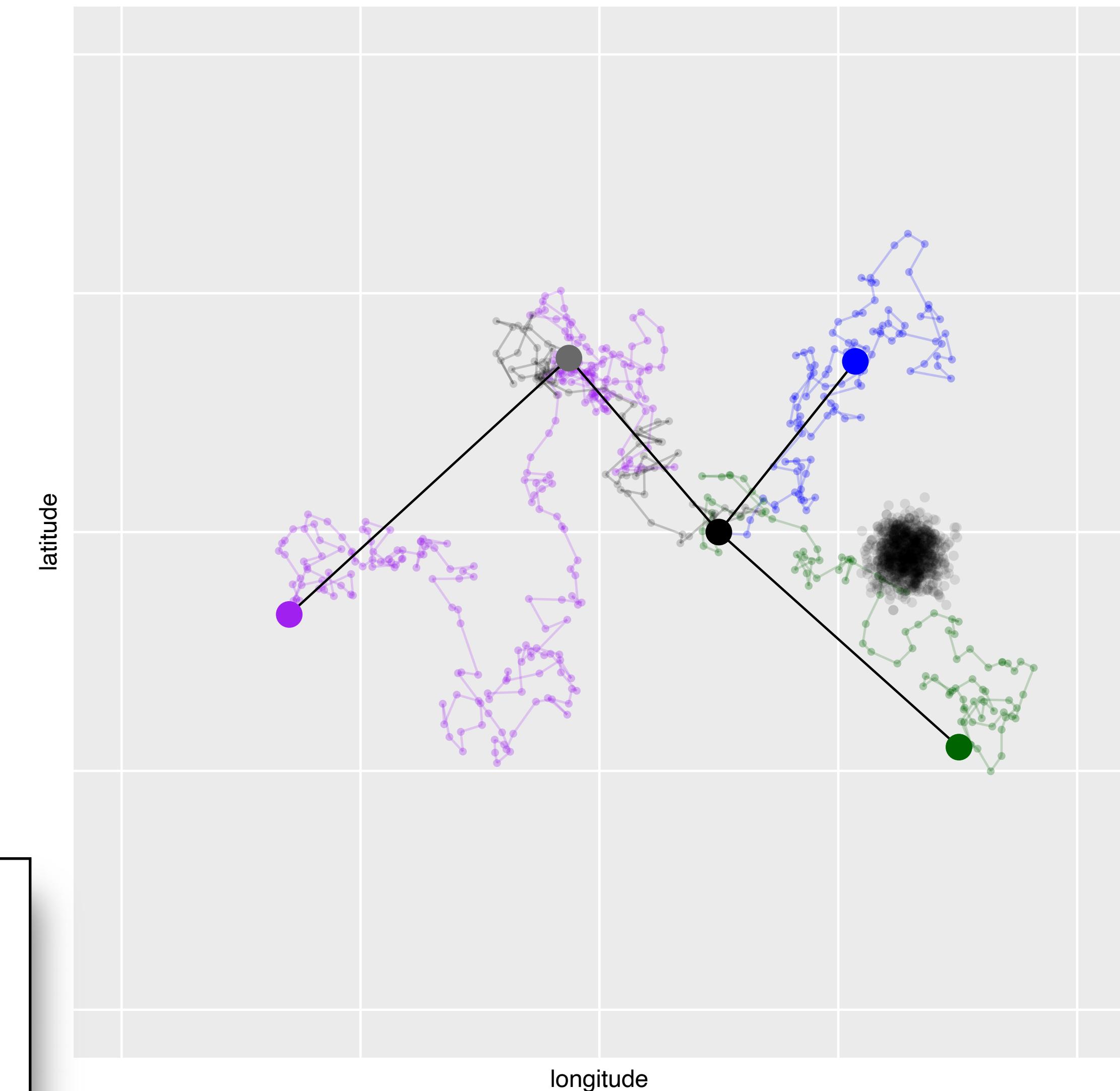


Bayesian model:

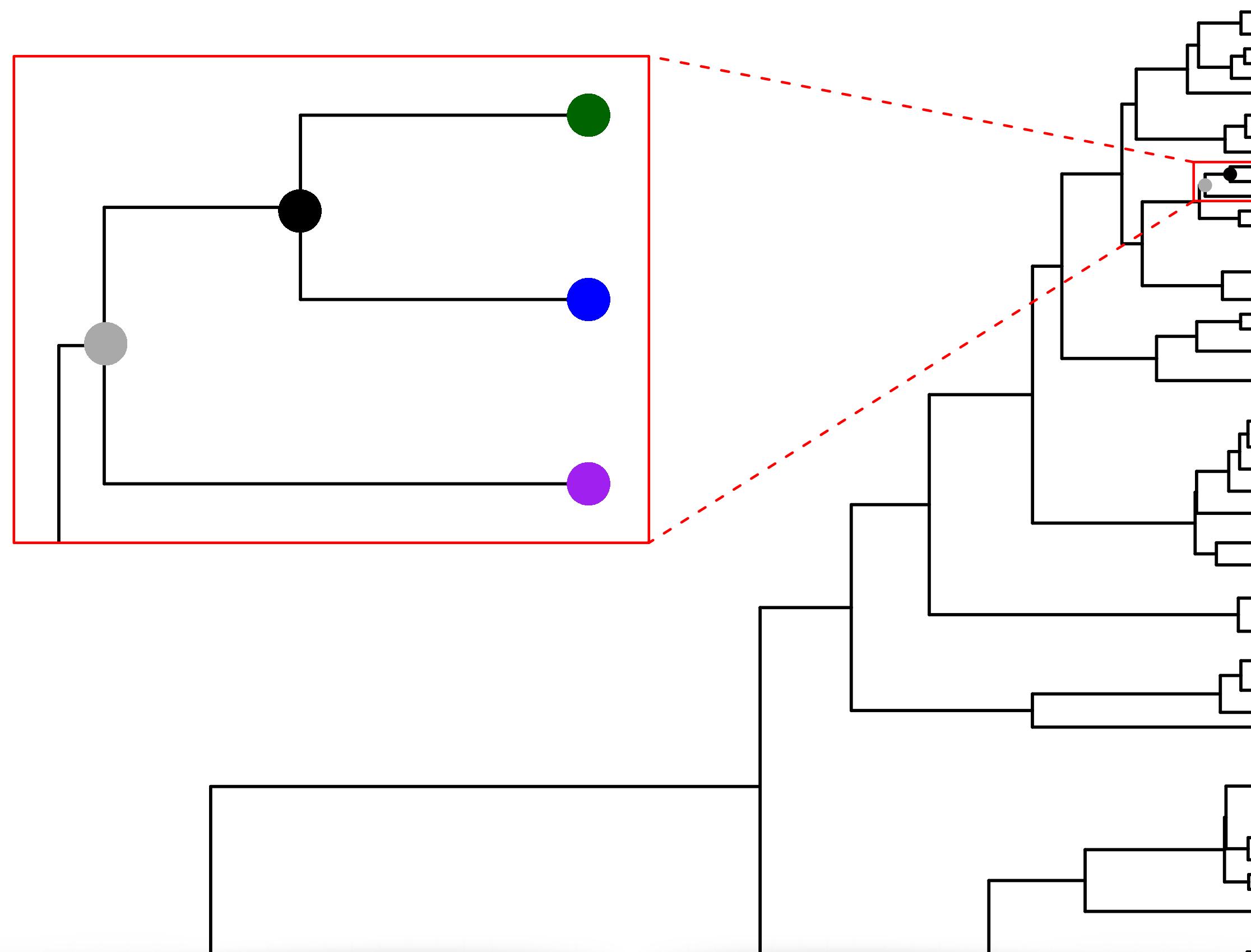
"What is the joint posterior distribution of ancestral positions & ages

given known locations of samples and mutation counts on branches?"

MCMC using NumPyro 



**... all locations and times serve
as constraints on one another.**

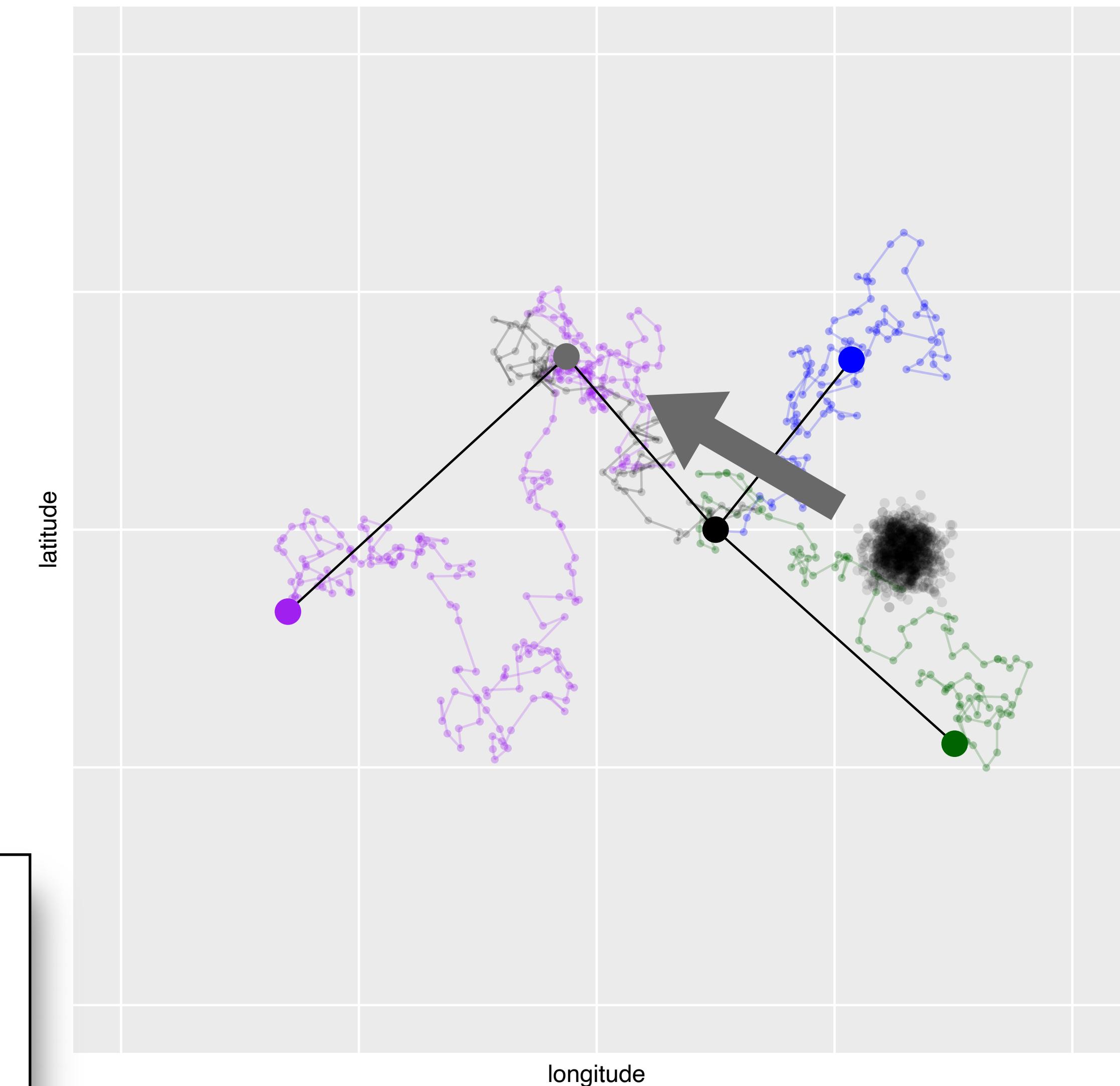


Bayesian model:

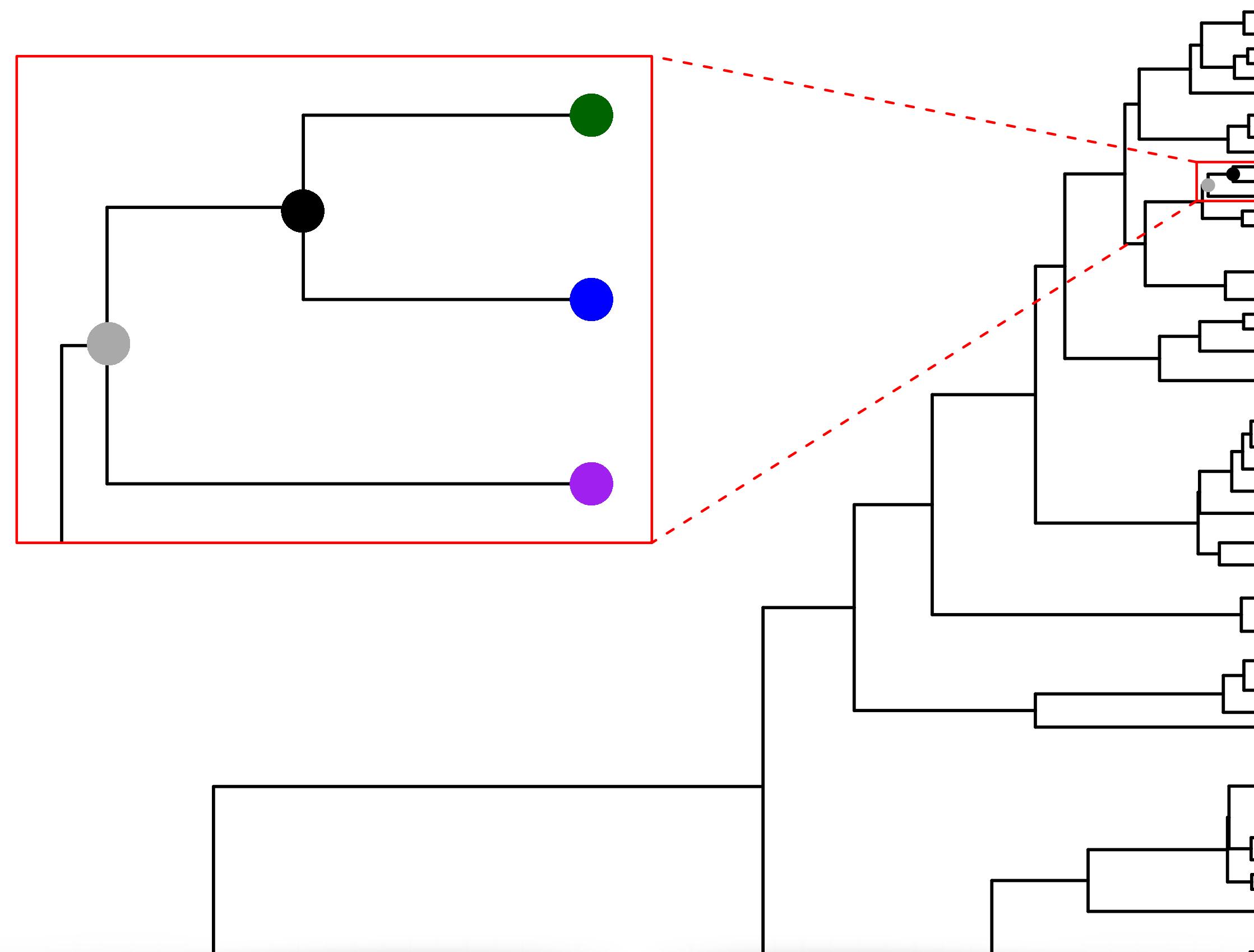
"What is the joint posterior distribution of ancestral positions & ages

given known locations of samples and mutation counts on branches?"

MCMC using NumPyro 



**... all locations and times serve
as constraints on one another.**

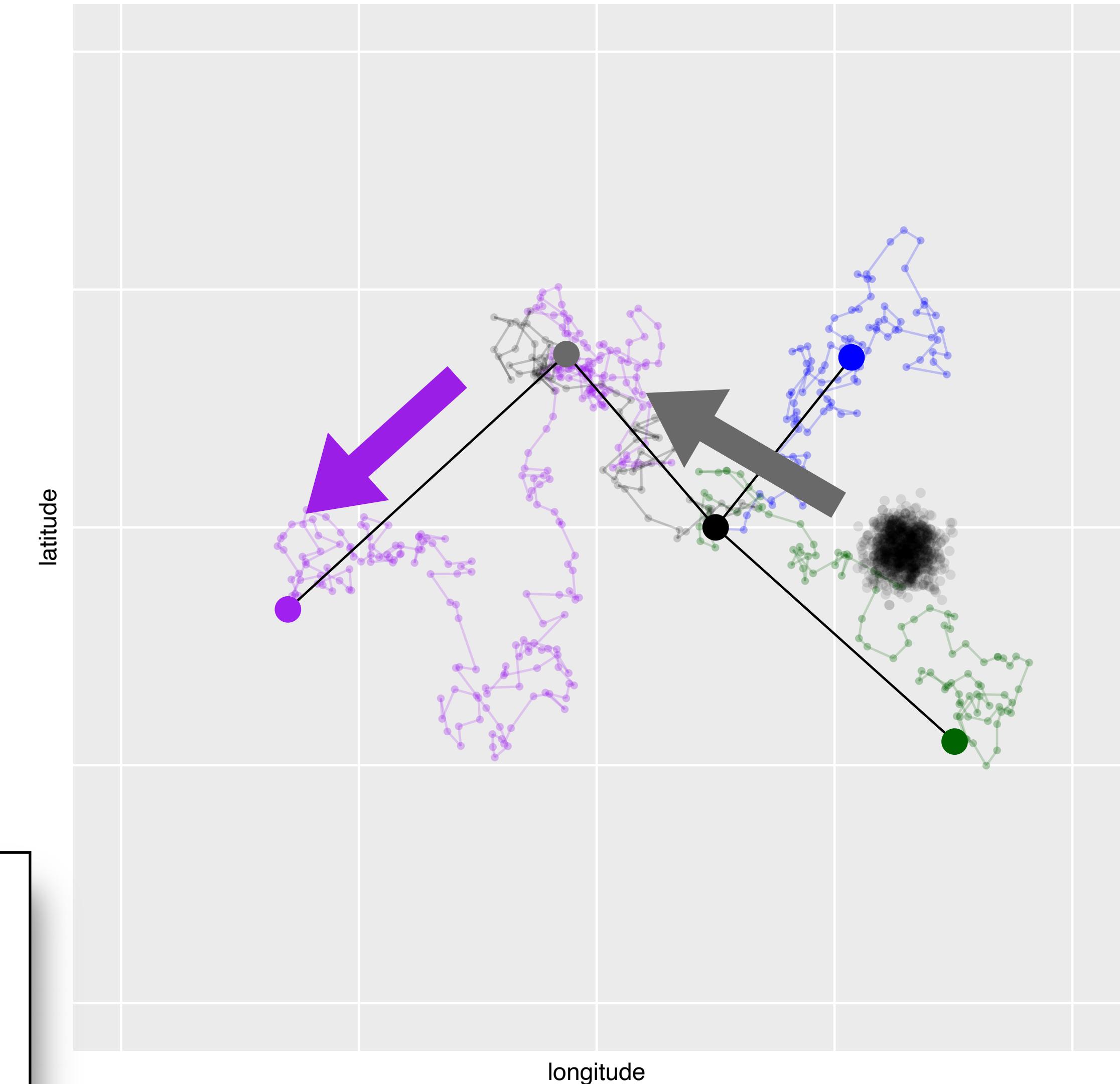


Bayesian model:

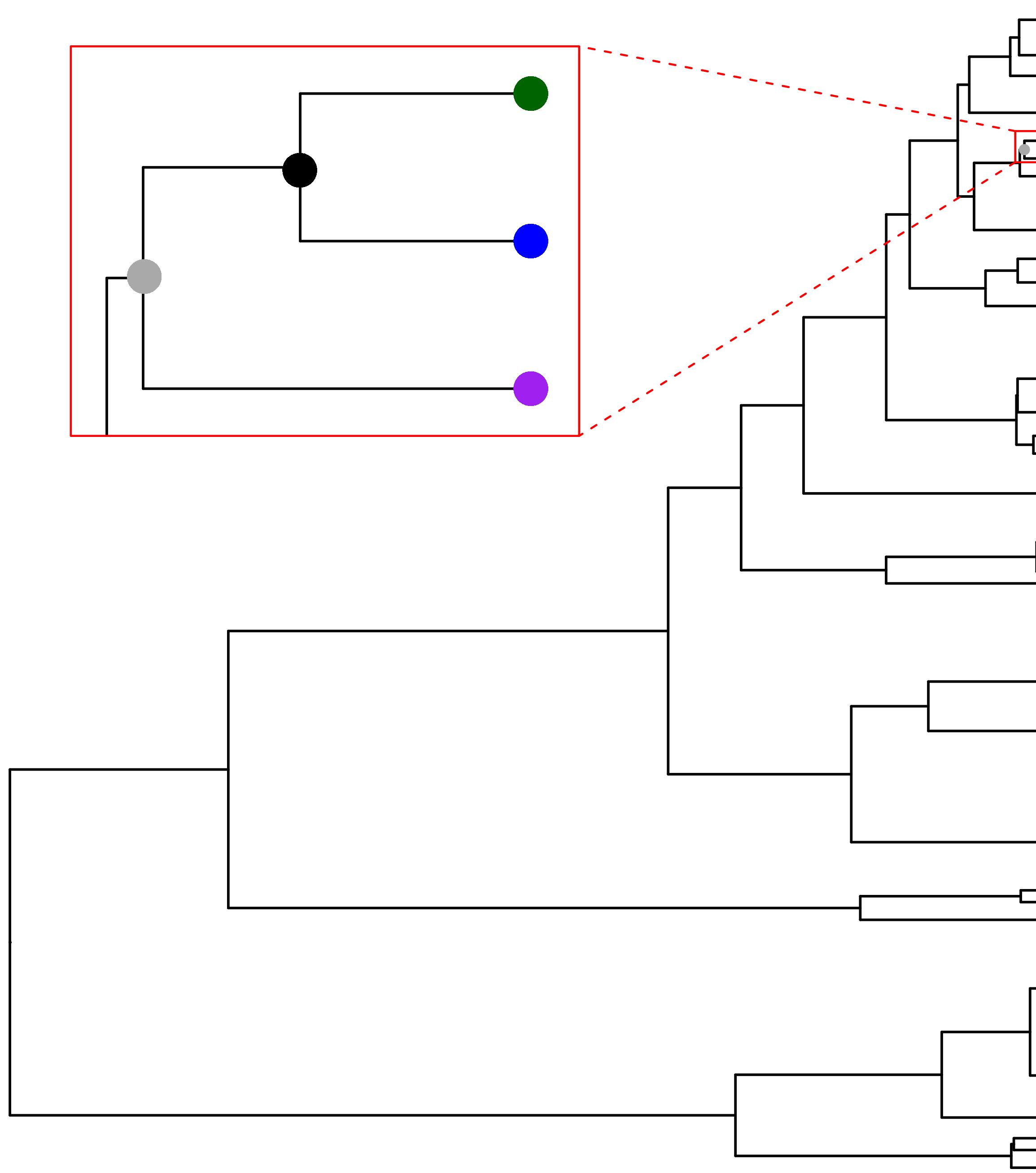
"What is the joint posterior distribution of ancestral positions & ages

given known locations of samples and mutation counts on branches?"

MCMC using NumPyro 



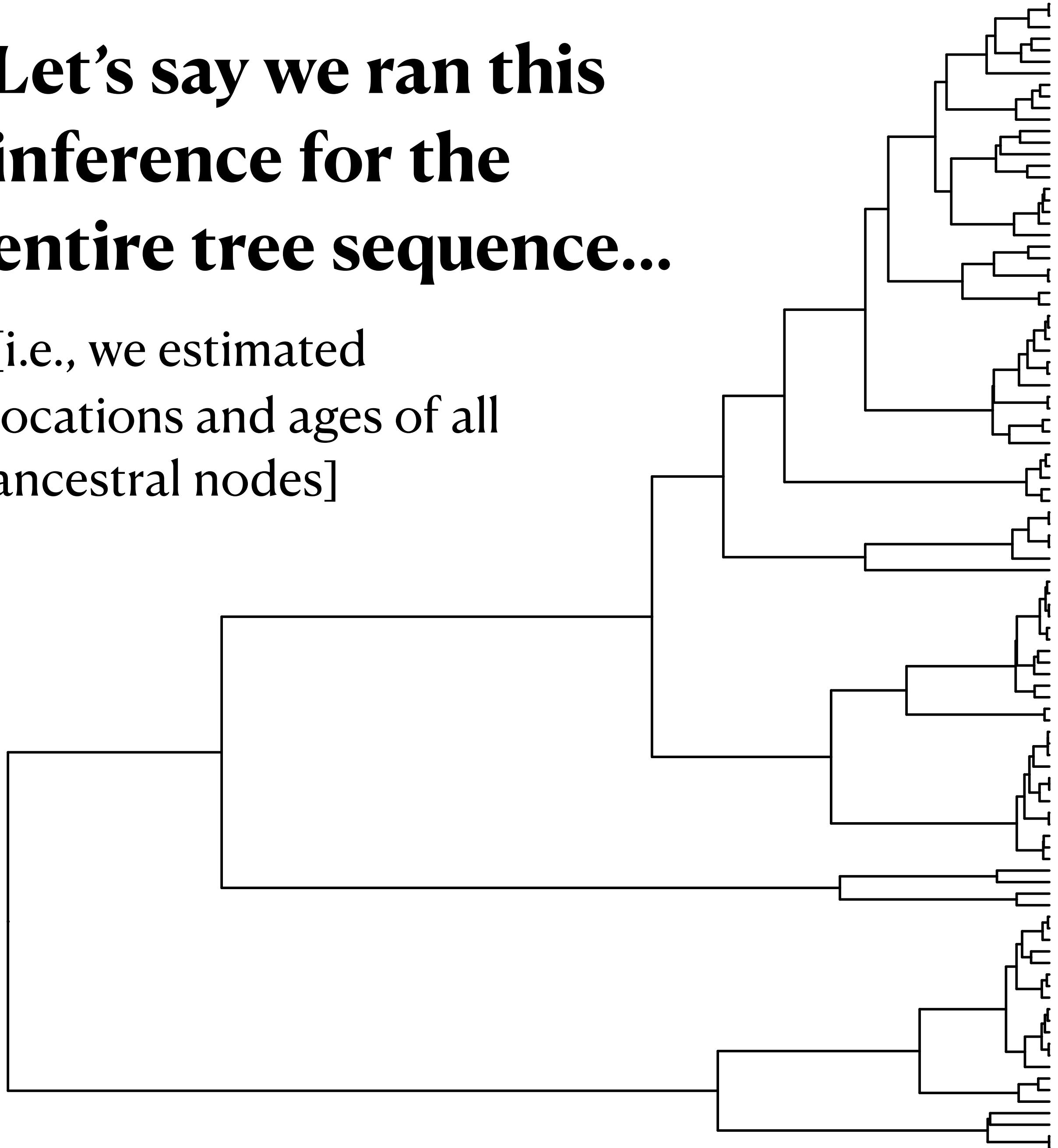
**... all locations and times serve
as constraints on one another.**



**... all locations and times serve
as constraints on one another.**

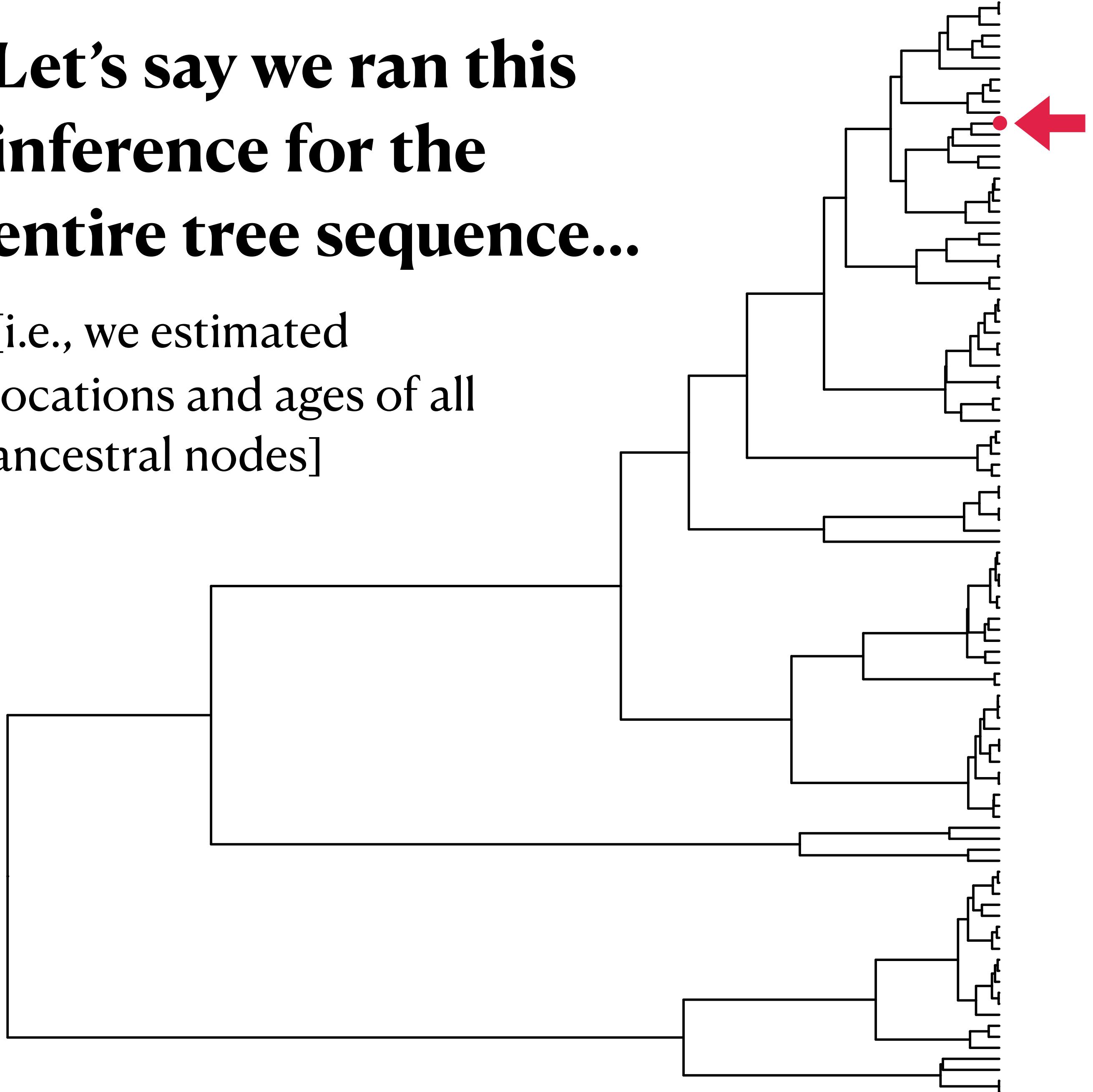
**Let's say we ran this
inference for the
entire tree sequence...**

[i.e., we estimated
locations and ages of all
ancestral nodes]



**Let's say we ran this
inference for the
entire tree sequence...**

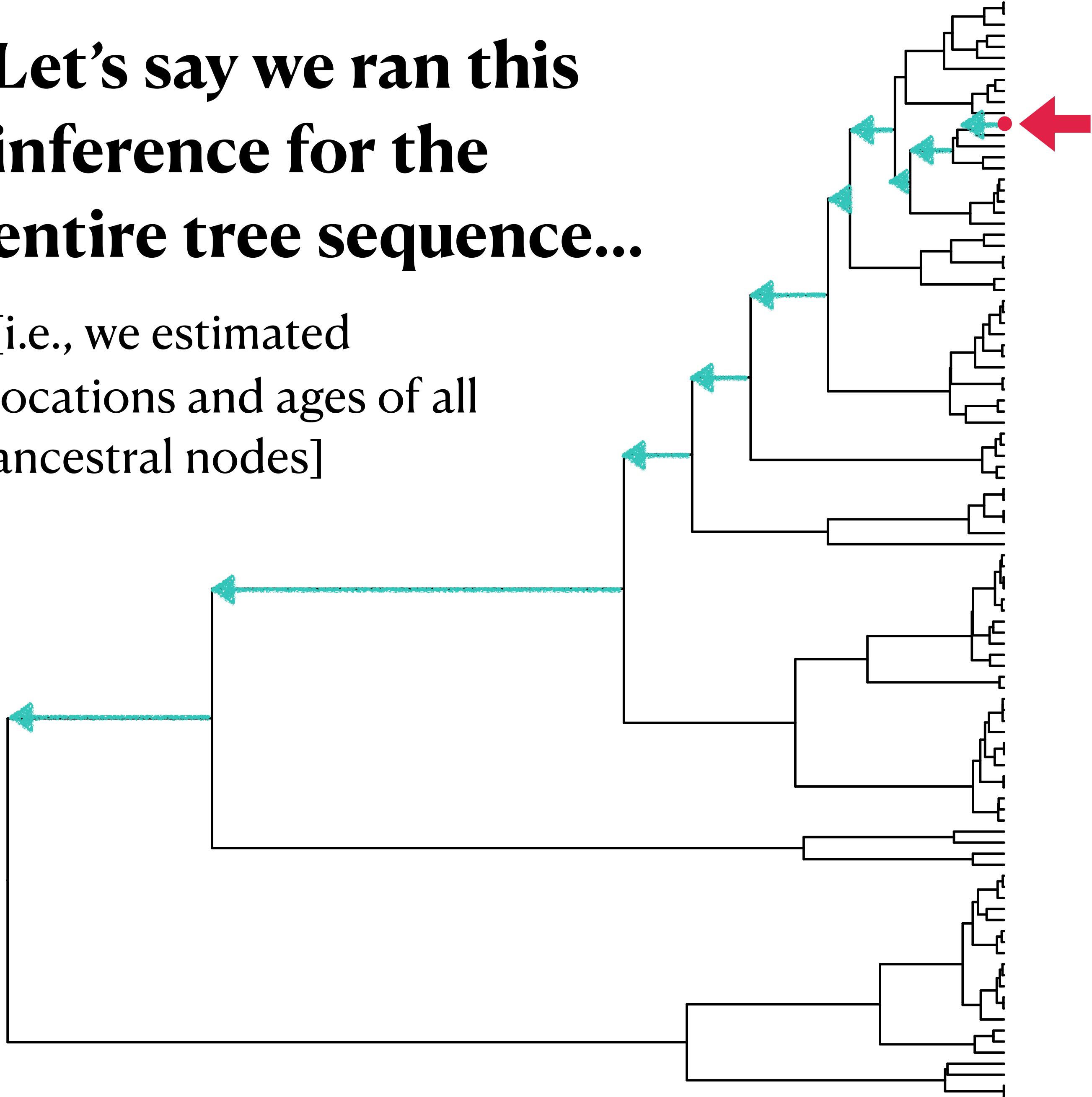
[i.e., we estimated
locations and ages of all
ancestral nodes]



For a **given individual**

Let's say we ran this inference for the entire tree sequence...

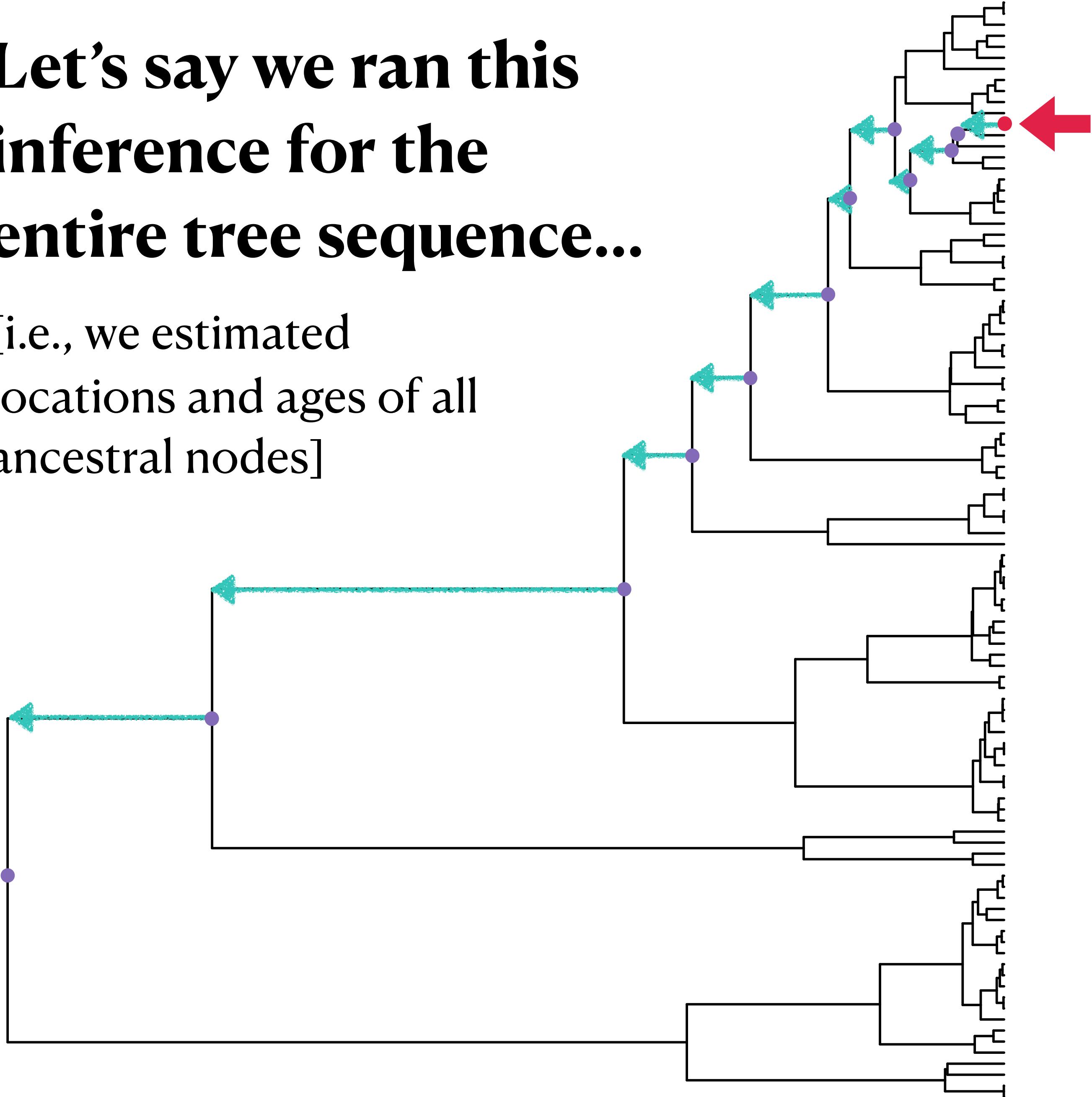
[i.e., we estimated locations and ages of all ancestral nodes]



For a **given individual**
we can "**climb up**" each tree

Let's say we ran this inference for the entire tree sequence...

[i.e., we estimated locations and ages of all ancestral nodes]



For a **given individual**
we can "**climb up**" each tree
and collect "paths" through
inferred **ancestral locations**.

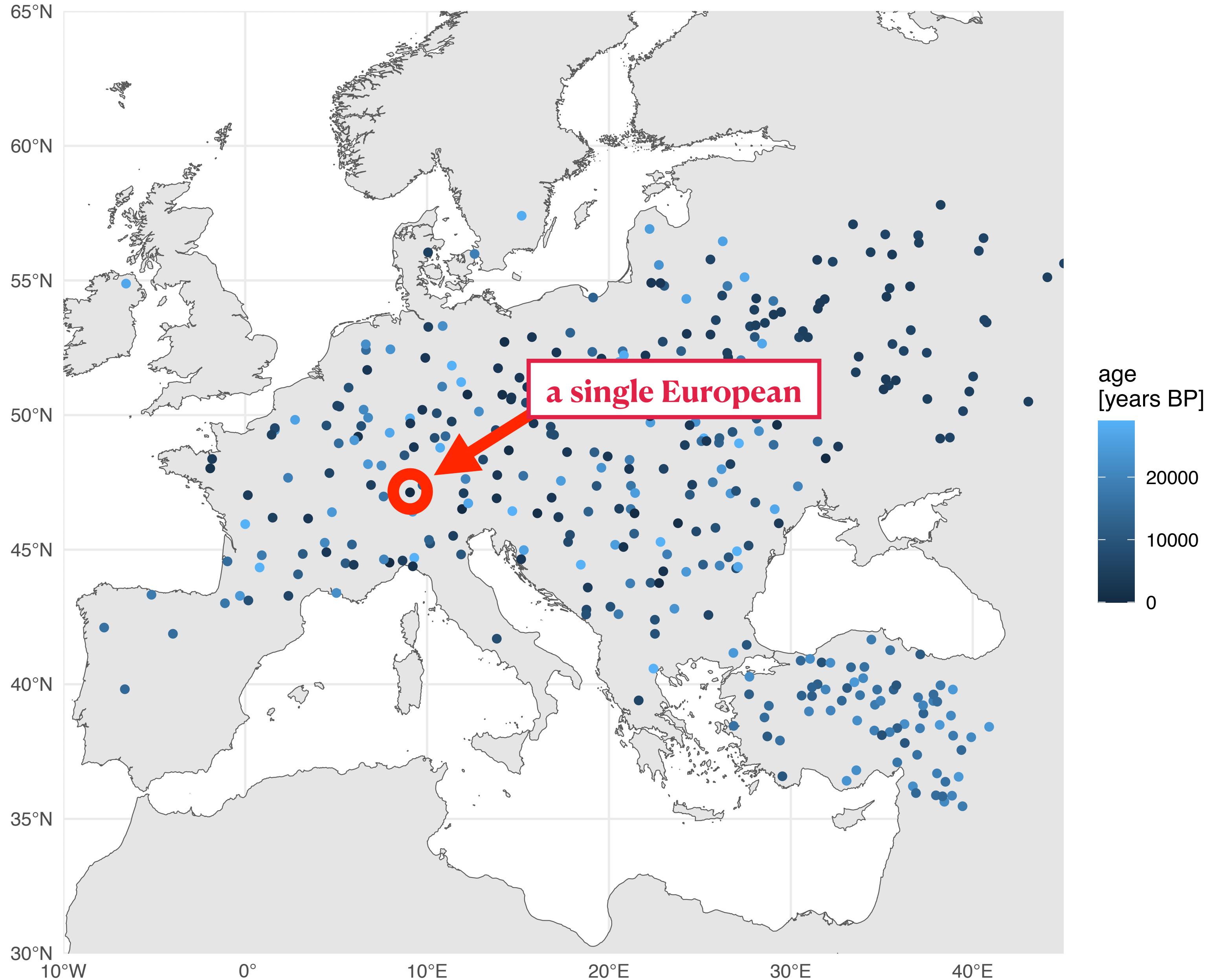
The Holy Grail?

Locations of individuals sampled between 30 kya and present-day



The Holy Grail?

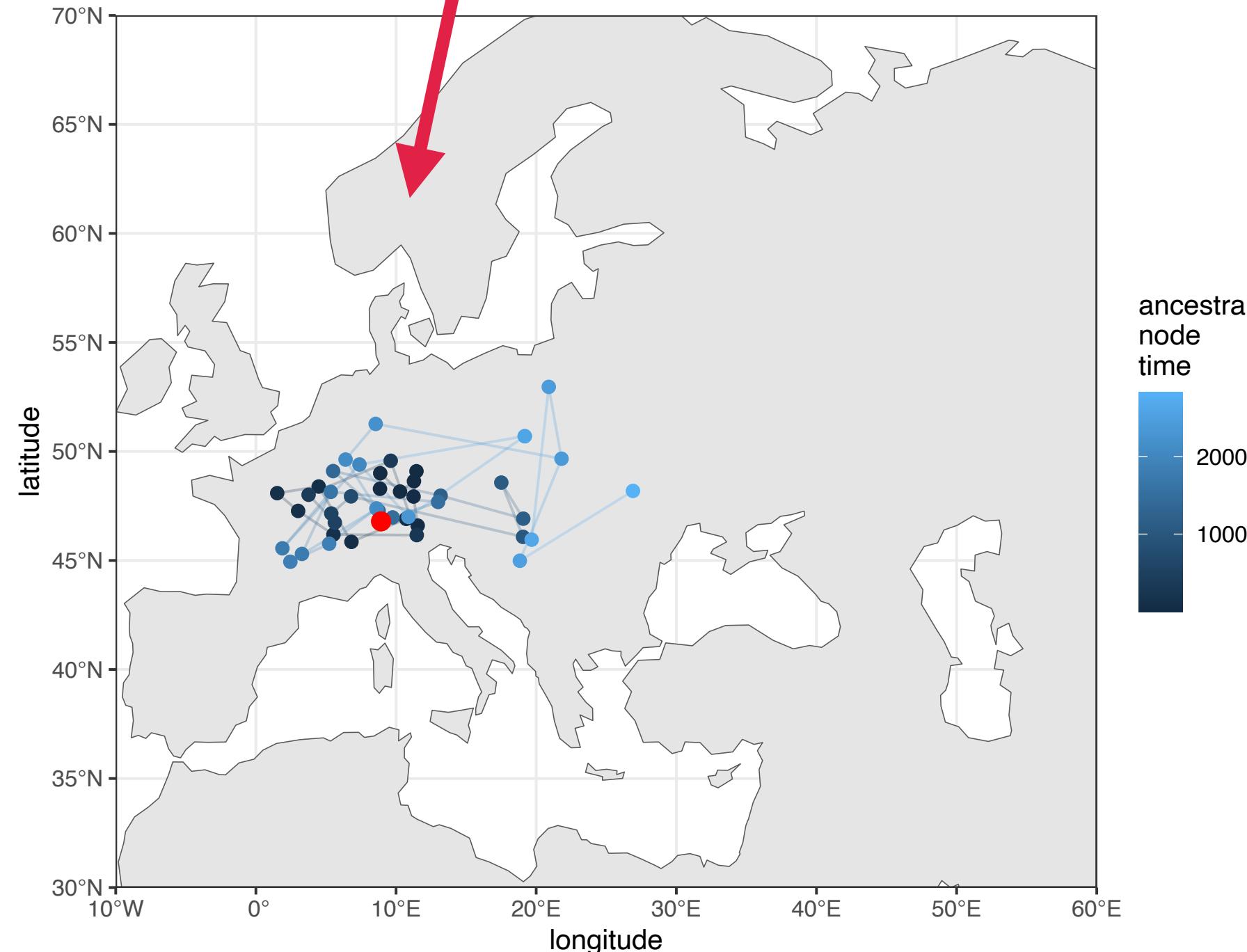
Locations of individuals sampled between 30 kya and present-day



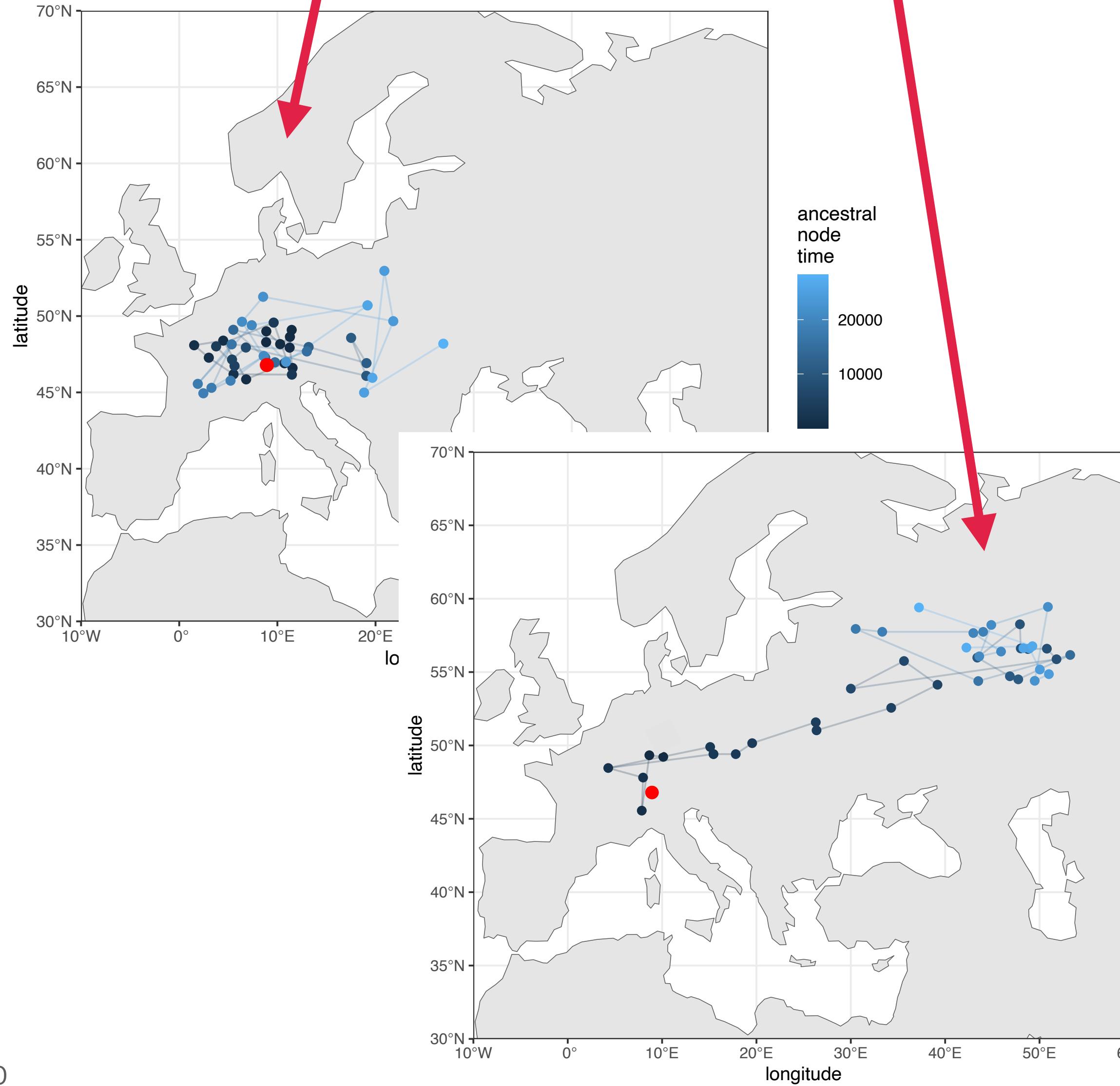
DNA sequence of that individual's genome (just 1 Mb)



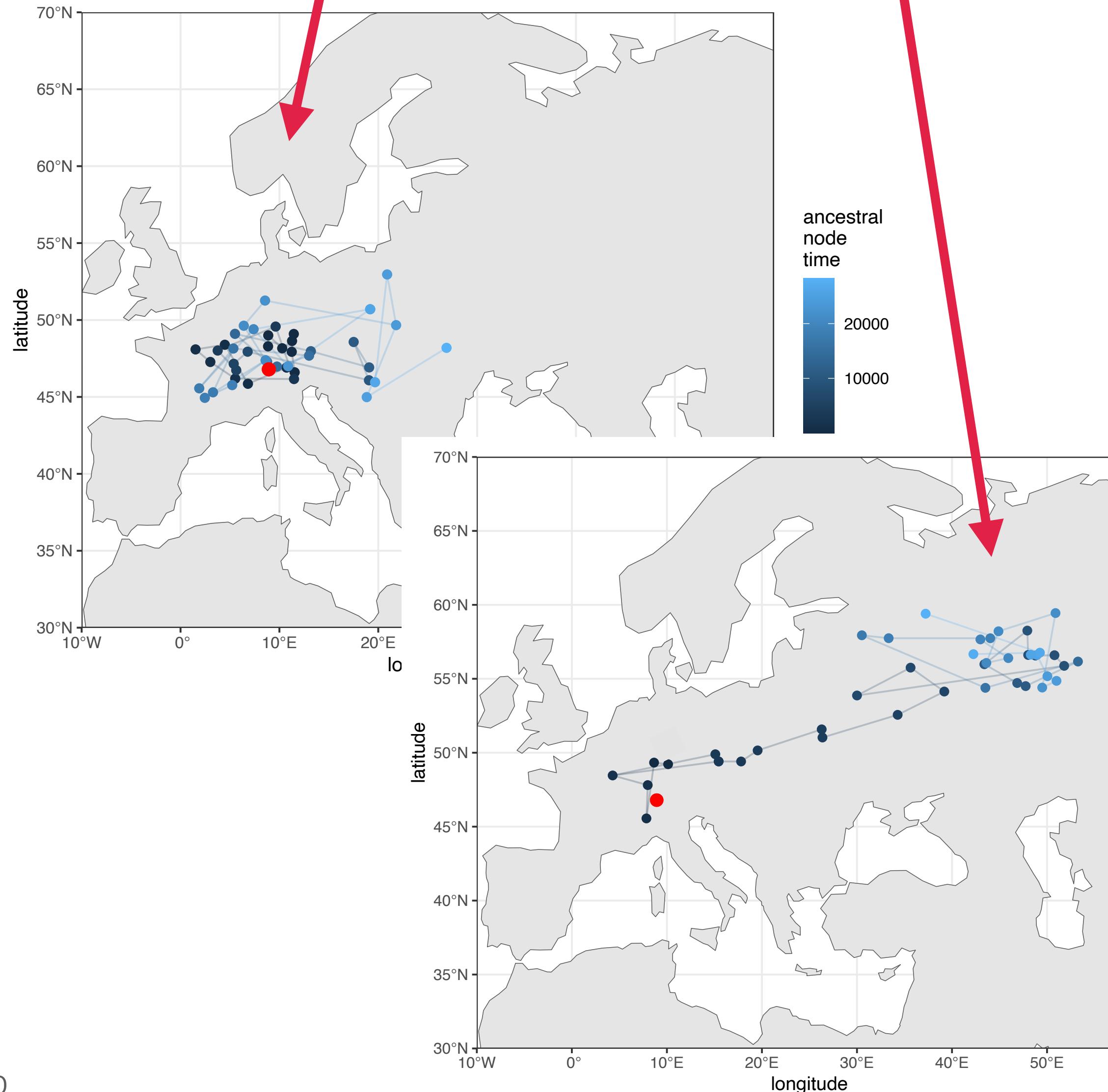
DNA sequence of that individual's genome (just 1 Mb)



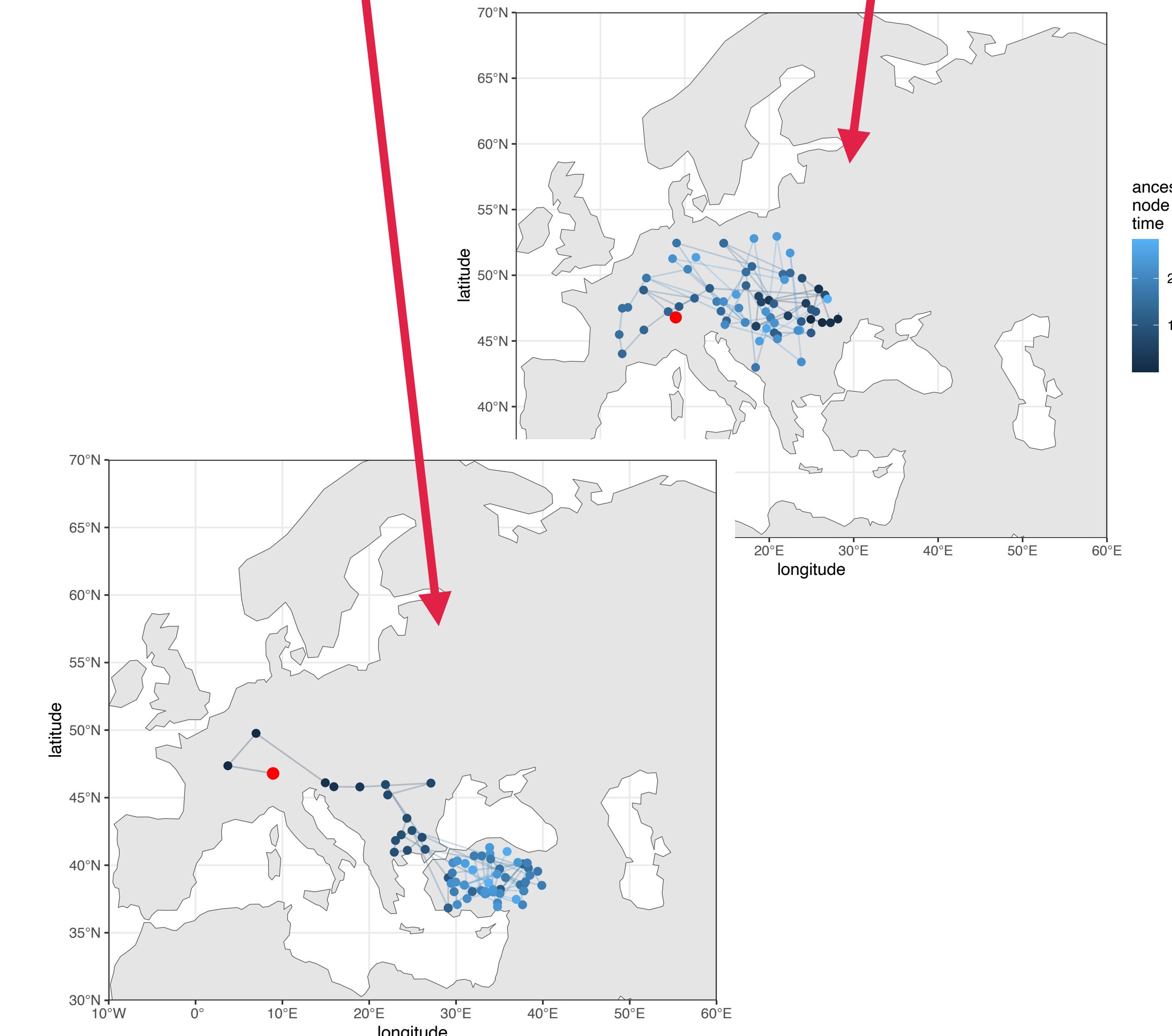
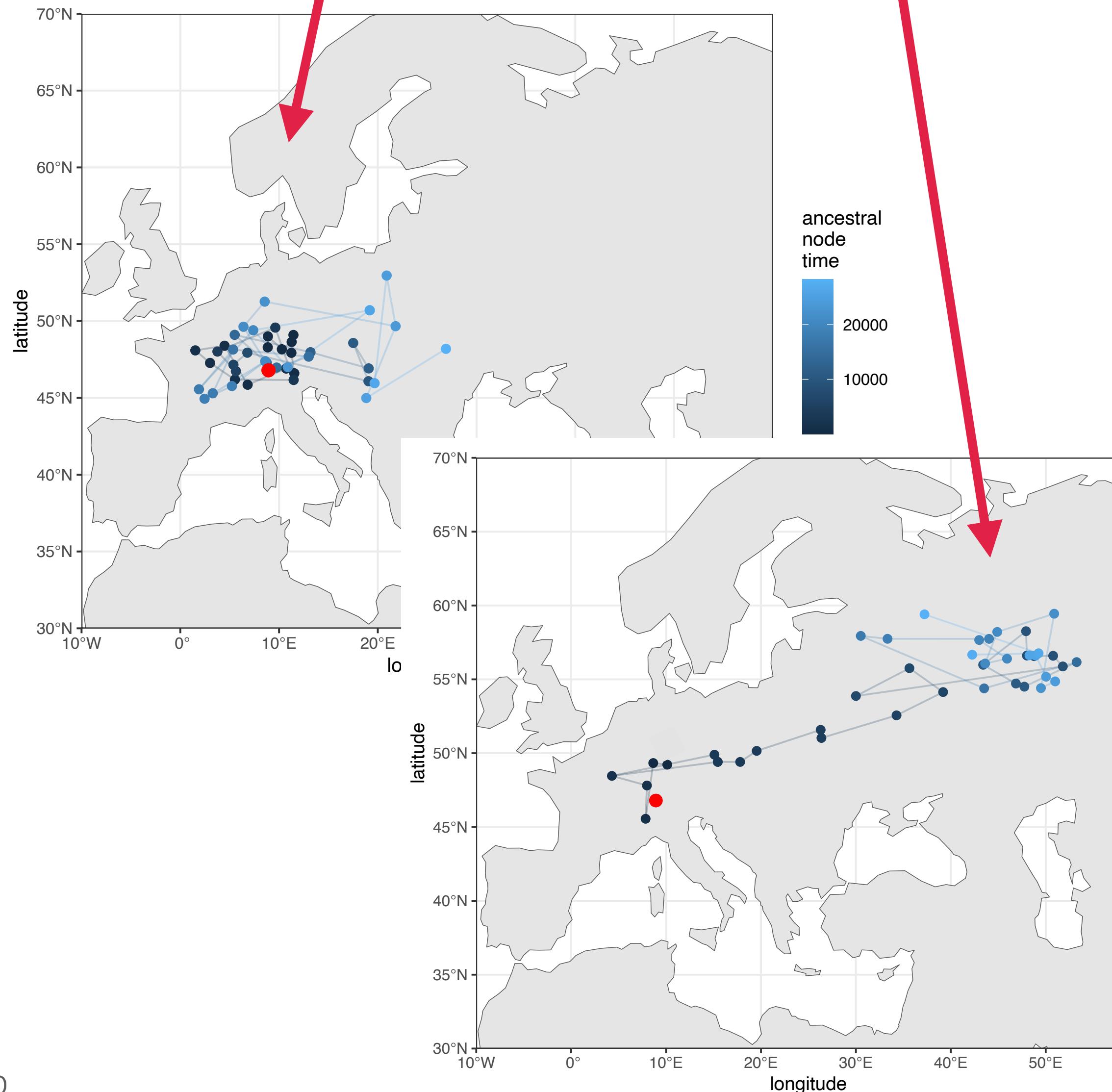
DNA sequence of that individual's genome (just 1 Mb)



DNA sequence of that individual's genome (just 1 Mb)



DNA sequence of that individual's genome (just 1 Mb)



Summary

Summary

- 1 ● We have reached sufficient density of aDNA sampling to facilitate **spatiotemporal analyses** of our evolutionary history.

Summary

- 1 ● We have reached sufficient density of aDNA sampling to facilitate **spatiotemporal analyses** of our evolutionary history.
- 2 ● **Tree sequence** is a true game-changer in population genetic simulation and facilitates novel approaches for analysis of large-scale data sets.

Summary

- 1 • We have reached sufficient density of aDNA sampling to facilitate **spatiotemporal analyses** of our evolutionary history.
- 2 • **Tree sequence** is a true game-changer in population genetic simulation and facilitates novel approaches for analysis of large-scale data sets.
- 3 • We are close to entirely **new kinds of geographically-explicit studies** of the history of genetic ancestry across space and time.

Thanks!

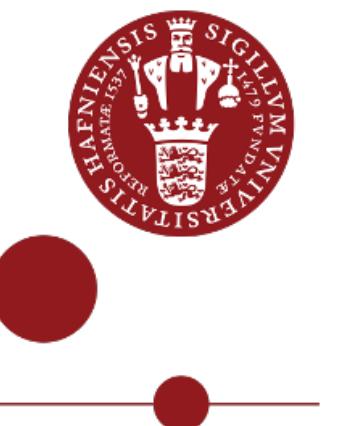
Fernando Racimo

& the whole Racimo group

Ben Haller, Peter L. Ralph



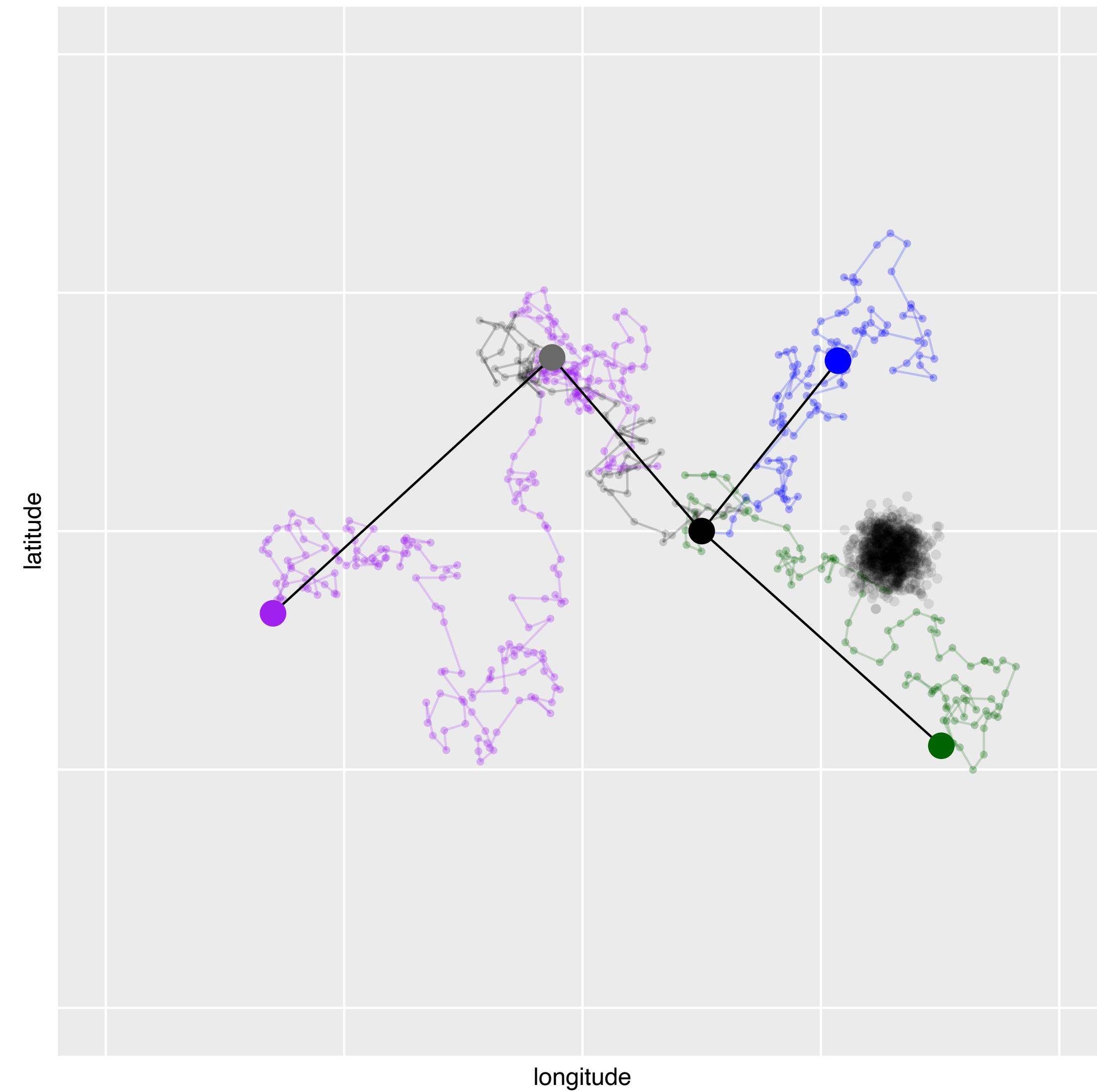
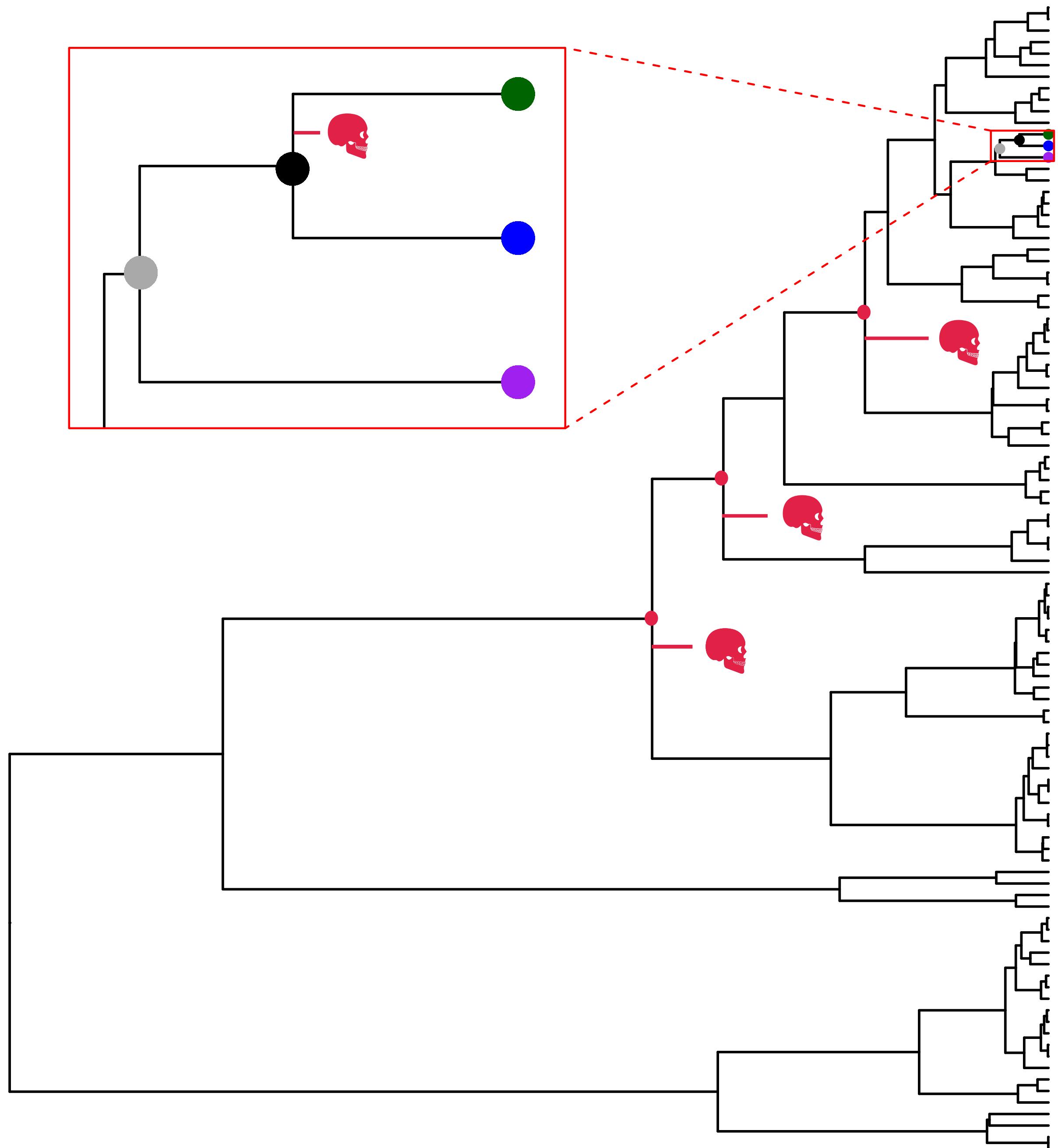
UNIVERSITY OF
COPENHAGEN



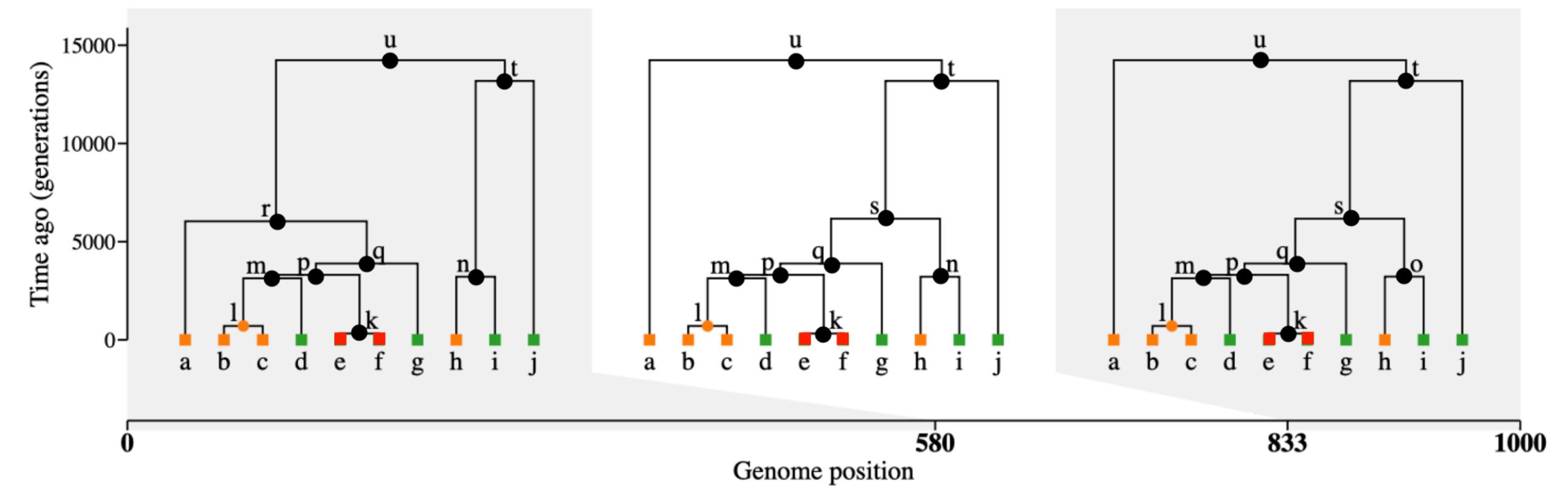
LUNDBECKFONDEN

novo
nordisk
fonden

email: mp@bodkan.net

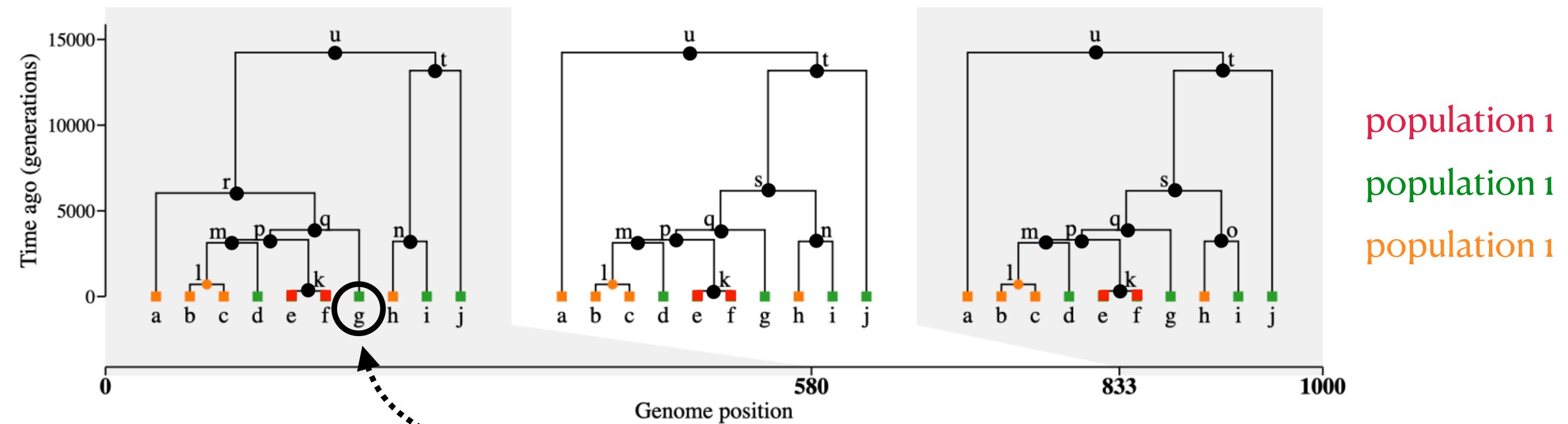


Genealogical nearest neighbours (GNN)



population 1
population 1
population 1

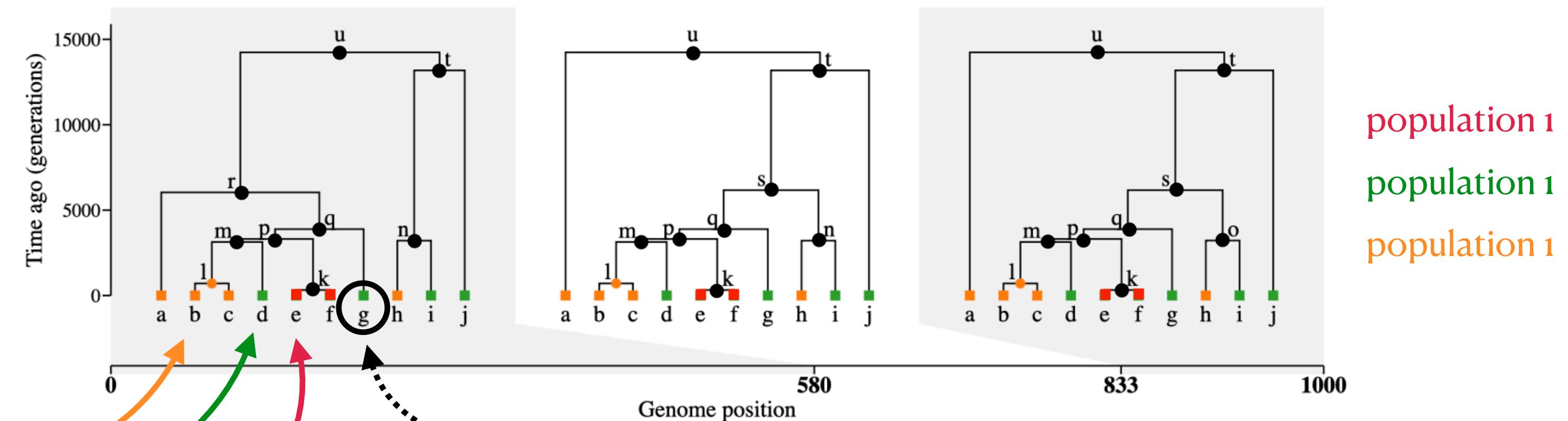
Genealogical nearest neighbours (GNN)



"What proportion of a given chromosome from one population has chromosomes **from all sampled populations** as immediate siblings?

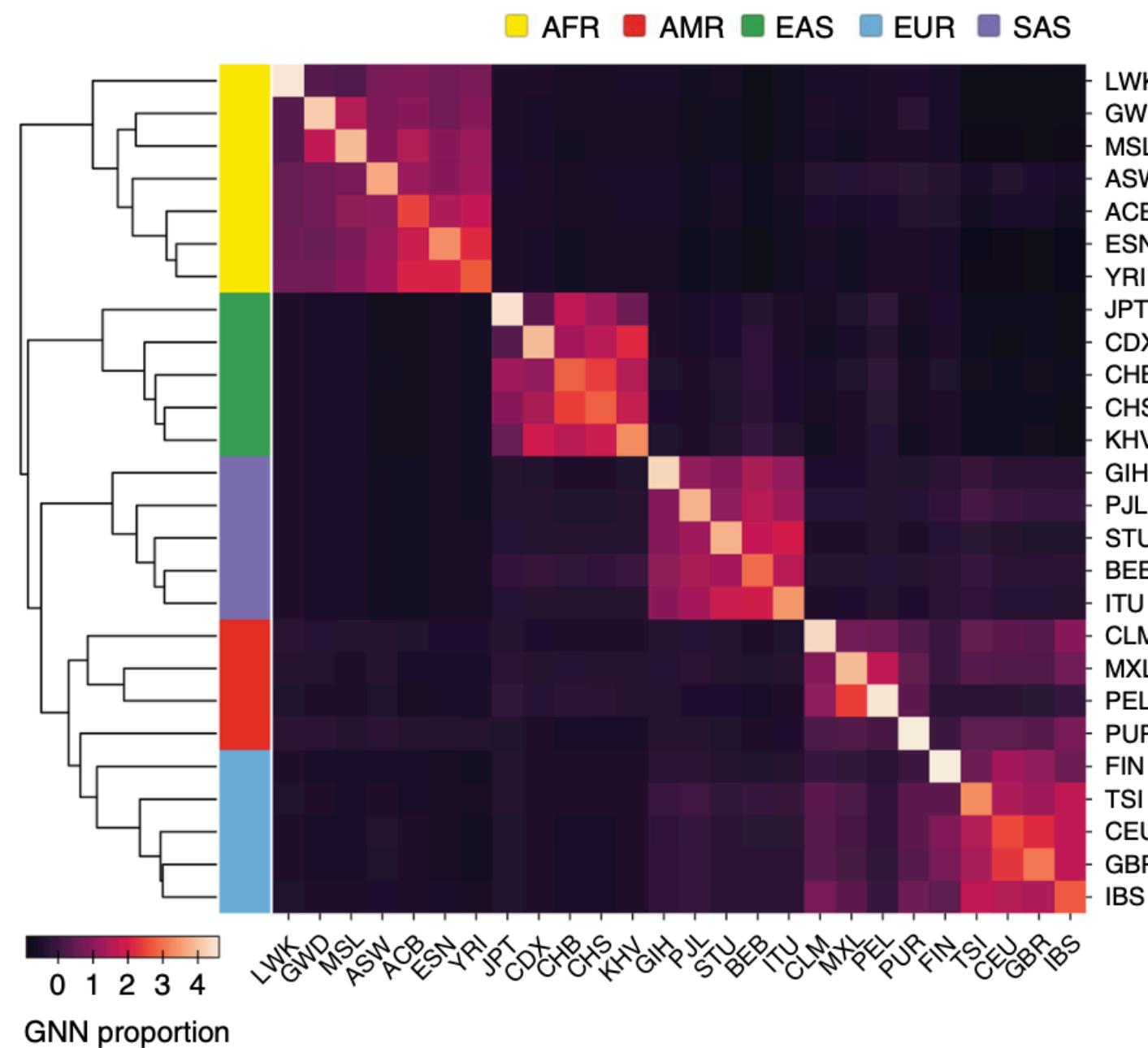
Genealogical nearest neighbours (GNN)

"What proportion of a given chromosome from one population has chromosomes from all sampled populations as immediate siblings?

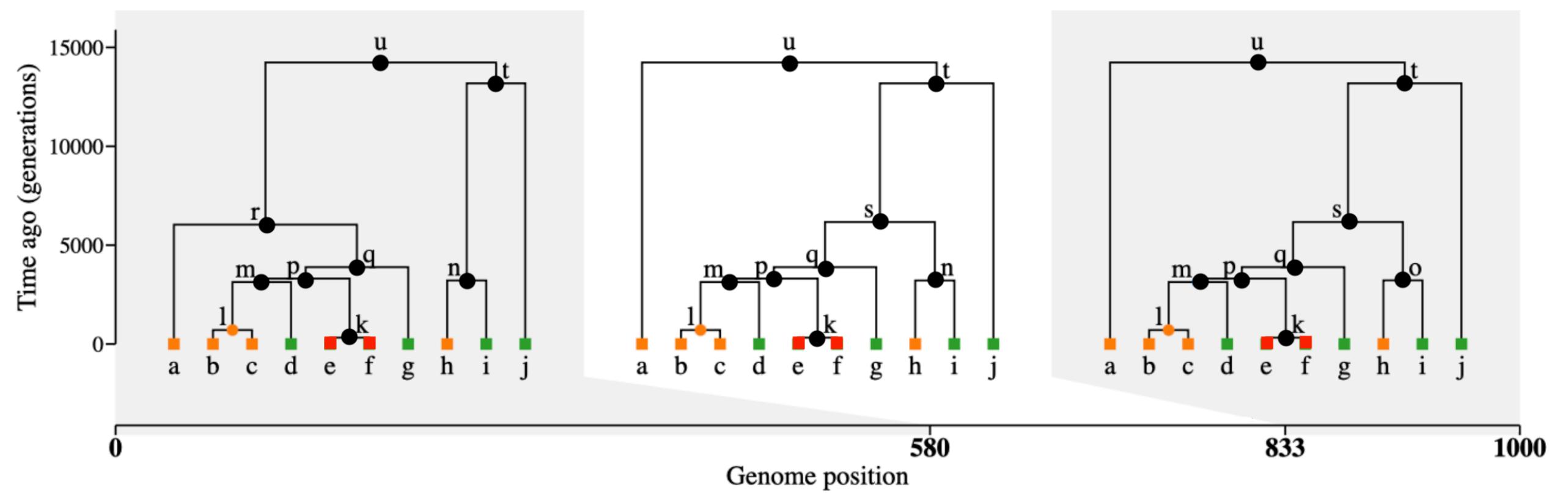


Genealogical nearest neighbours (GNN)

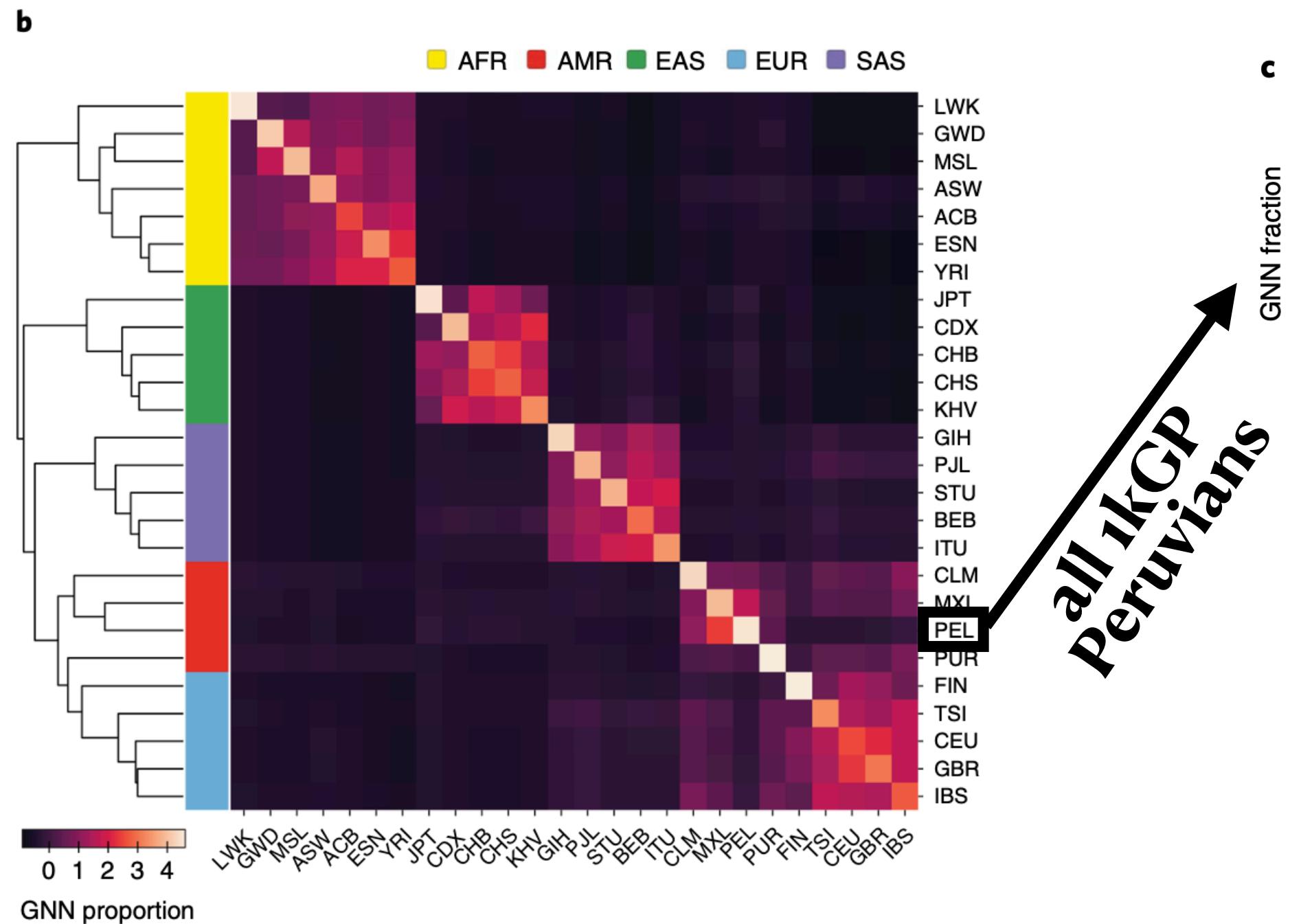
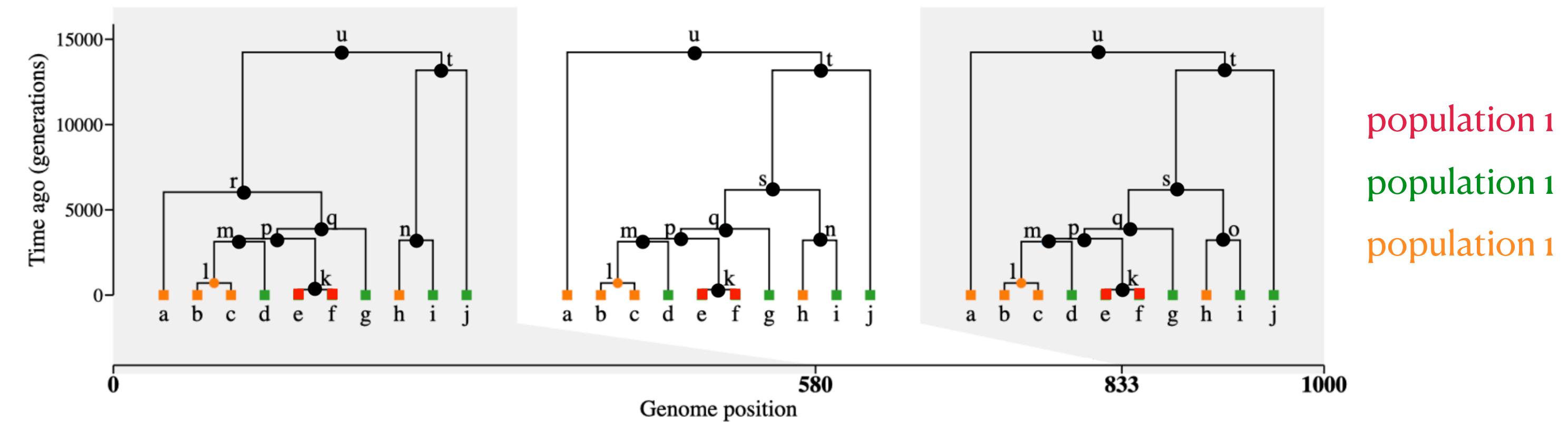
b



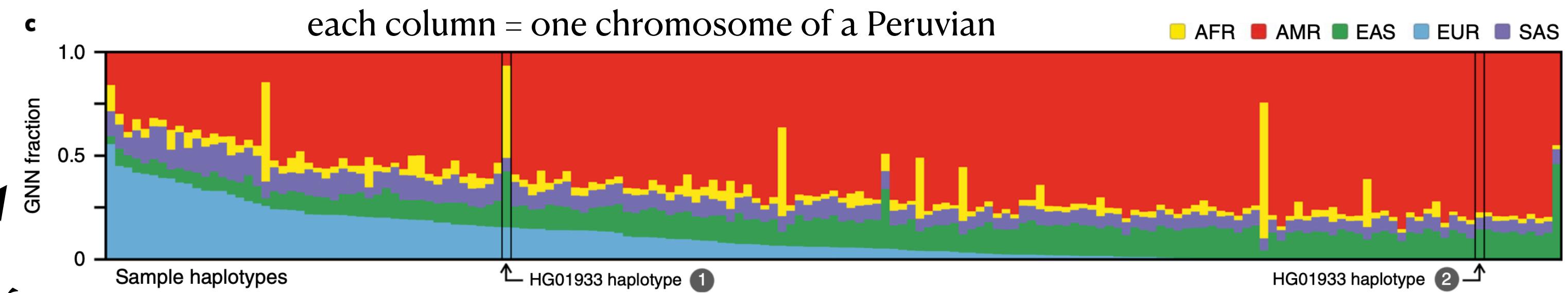
GNN clustering for all 2504
1000 Genomes Project individuals



Genealogical nearest neighbours (GNN)

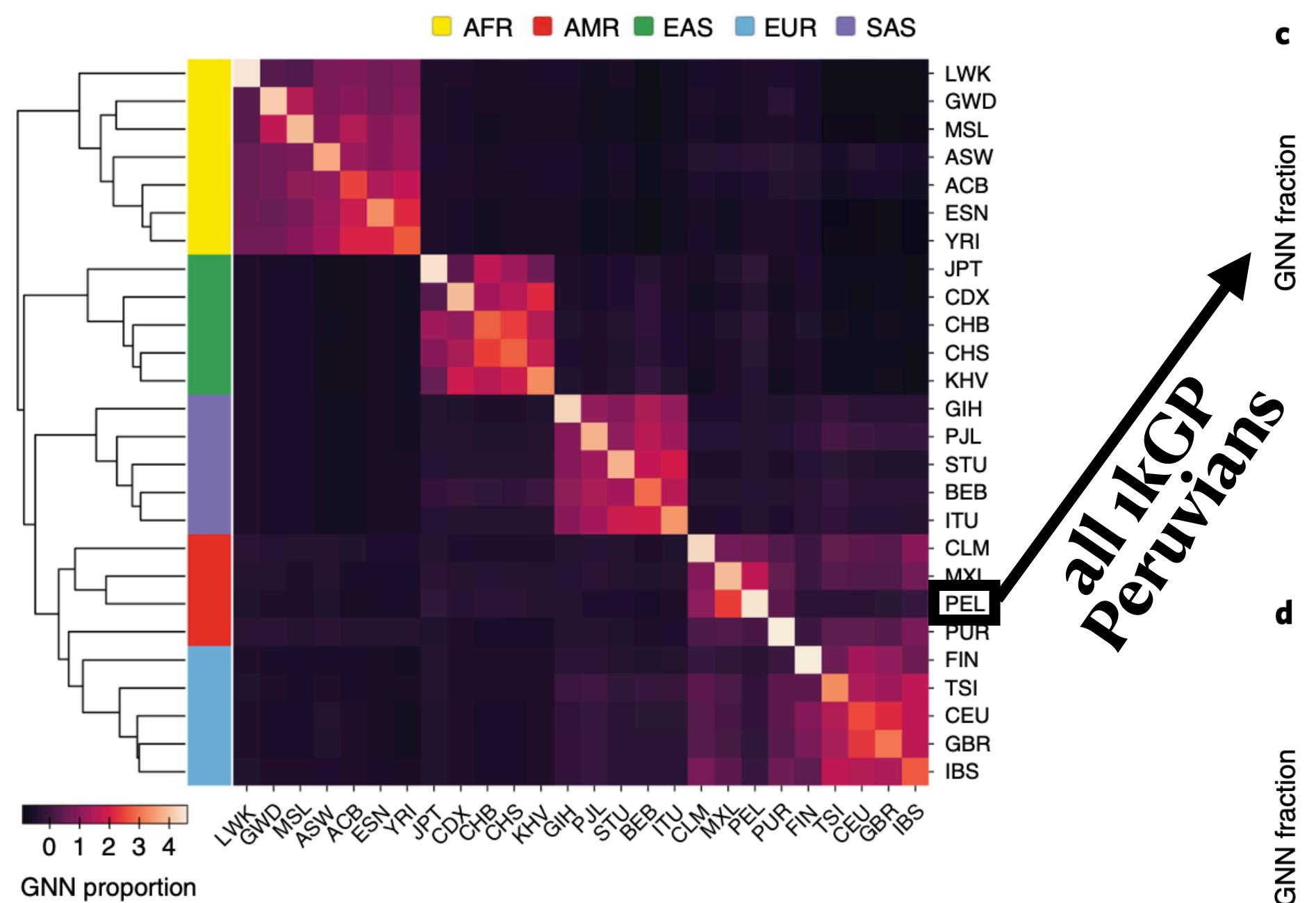


GNN clustering for all 2504
1000 Genomes Project individuals

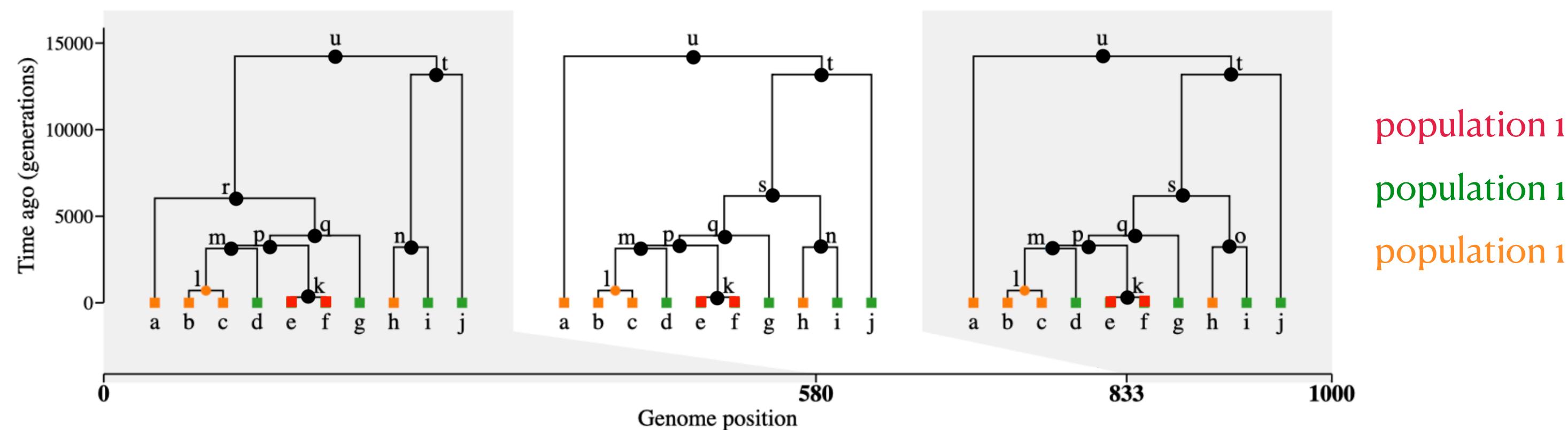


Genealogical nearest neighbours (GNN)

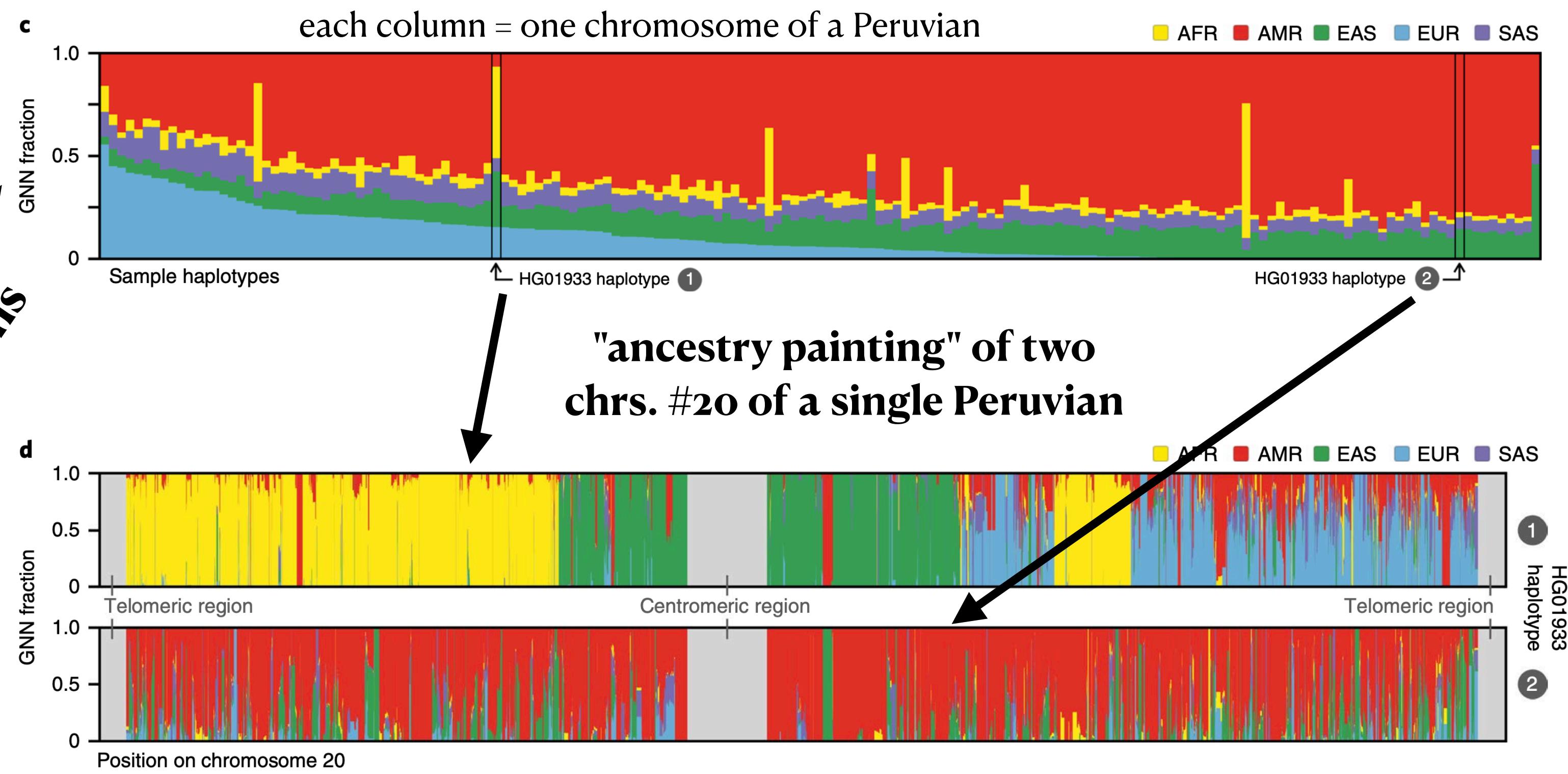
b



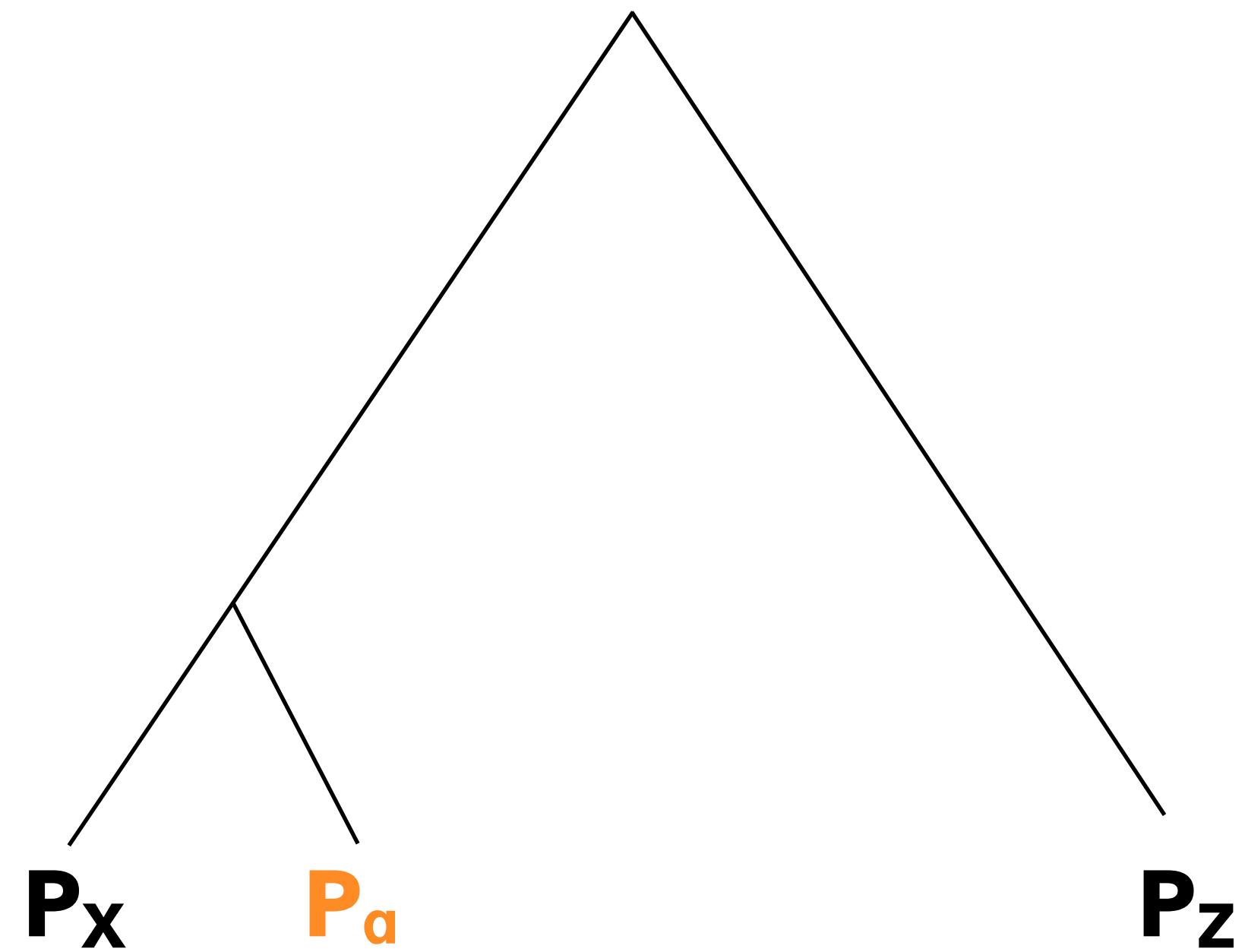
GNN clustering for all 2504
1000 Genomes Project individuals



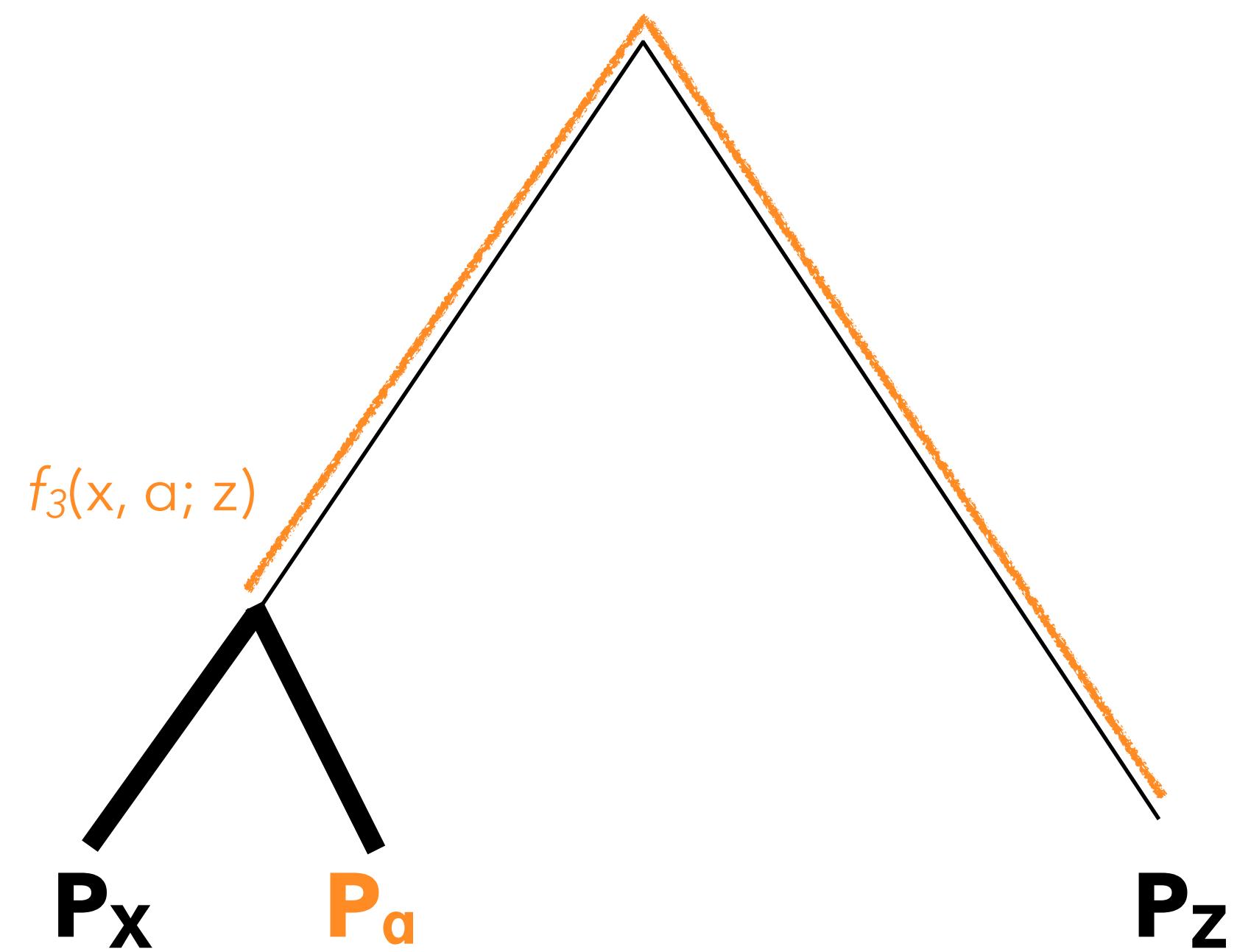
c



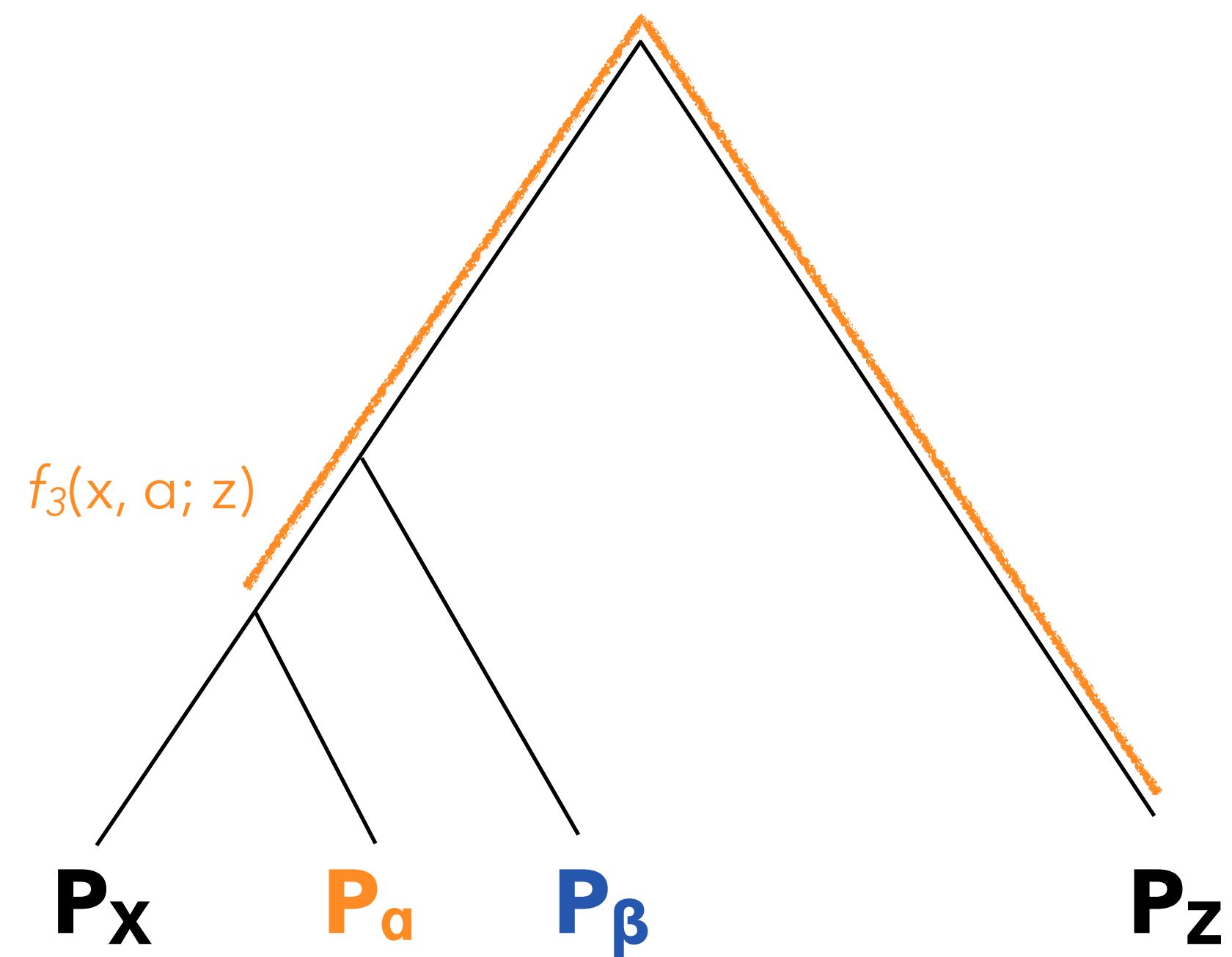
How do traditional popgen statistics work? – *f*-statistics



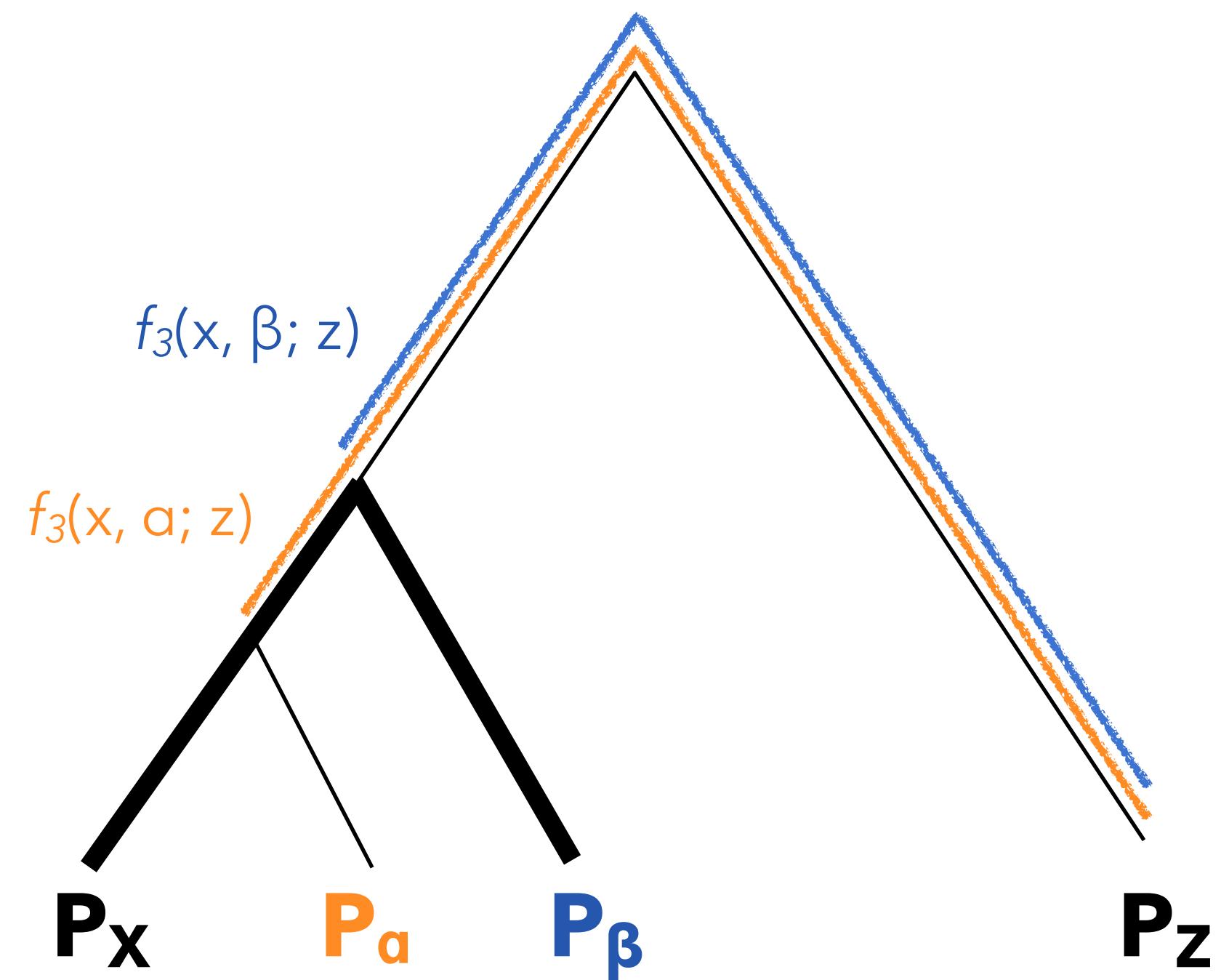
How do traditional popgen statistics work? – *f*-statistics



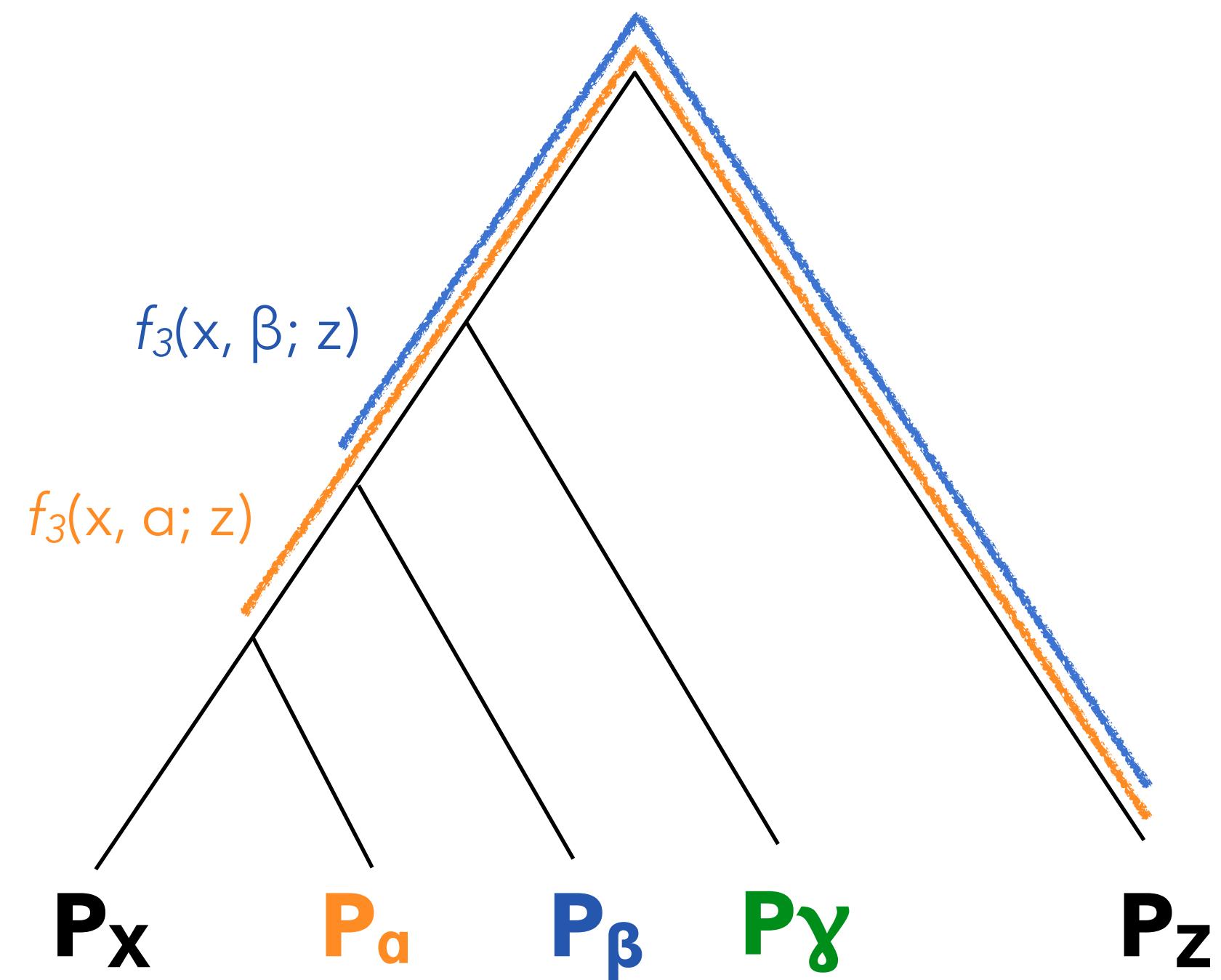
How do traditional popgen statistics work? – *f*-statistics



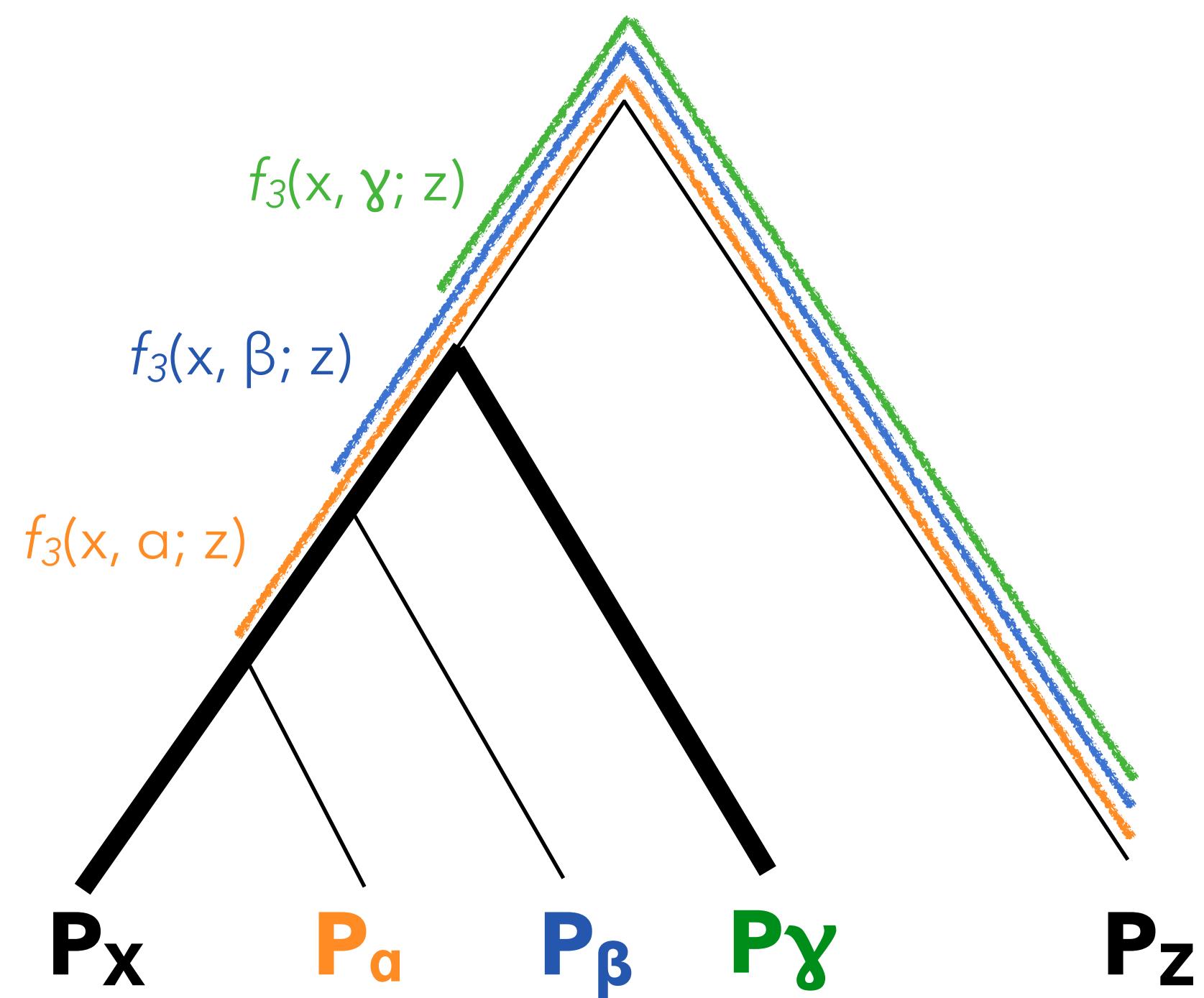
How do traditional popgen statistics work? – f -statistics



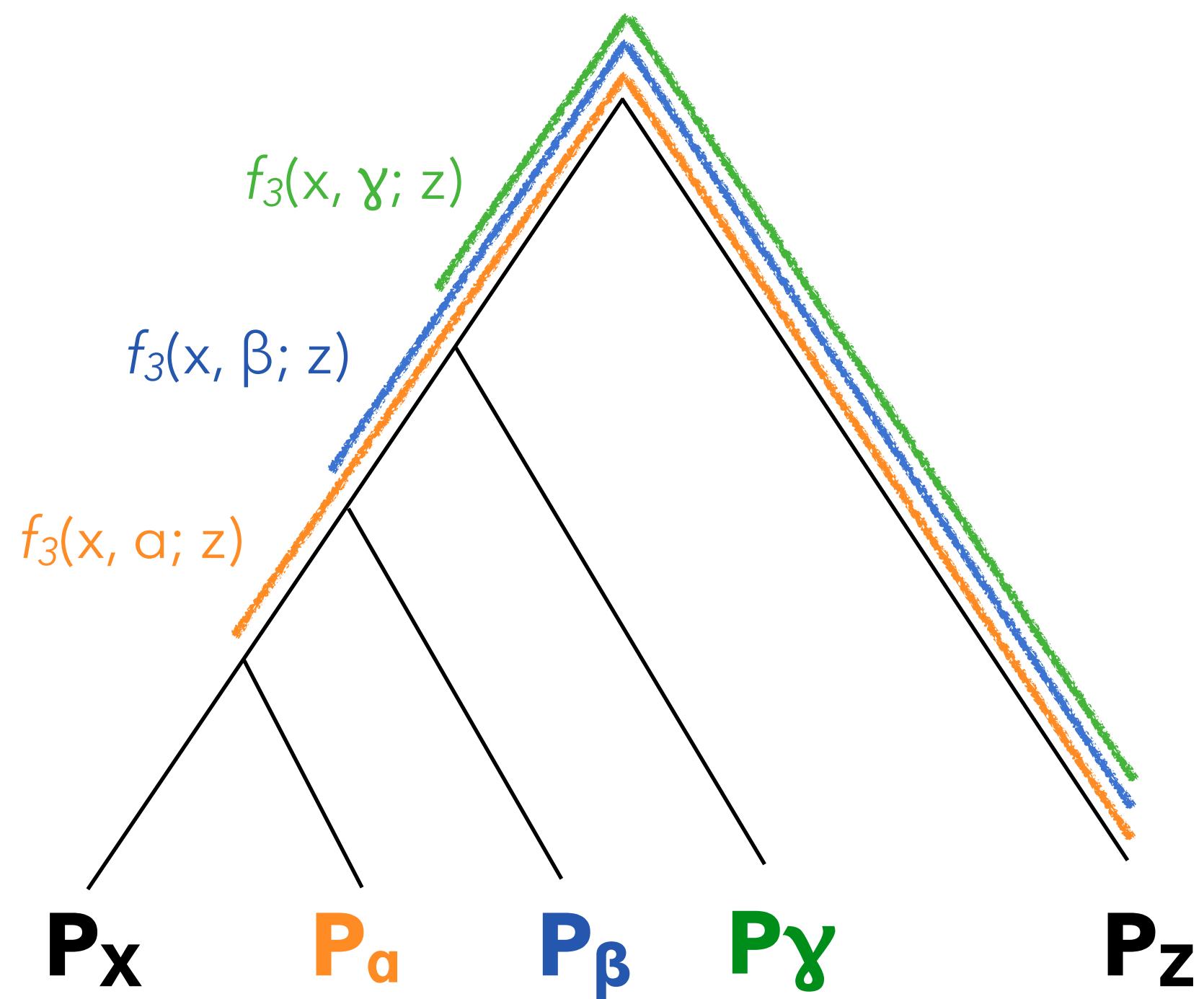
How do traditional popgen statistics work? – *f*-statistics



How do traditional popgen statistics work? – f -statistics

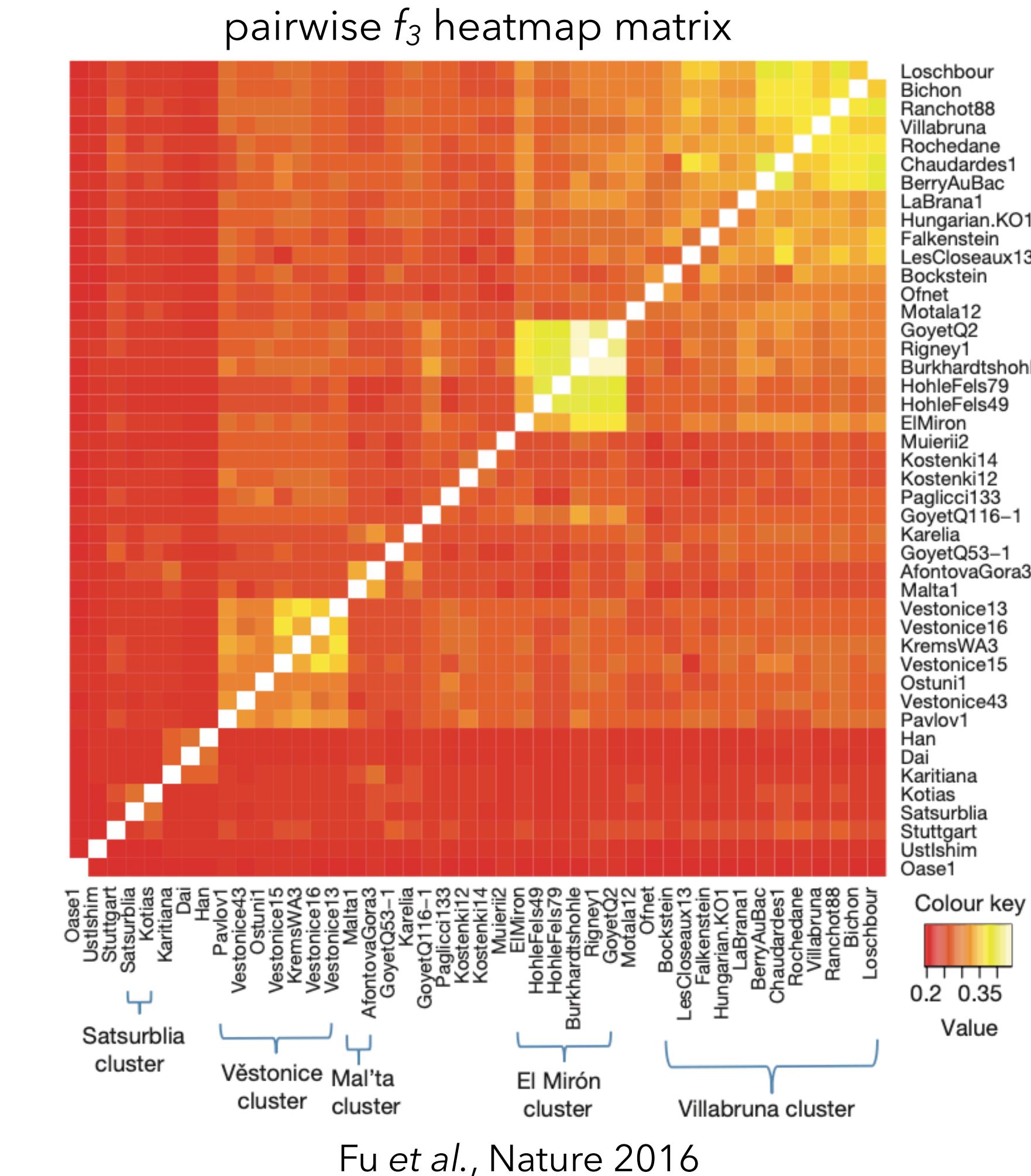
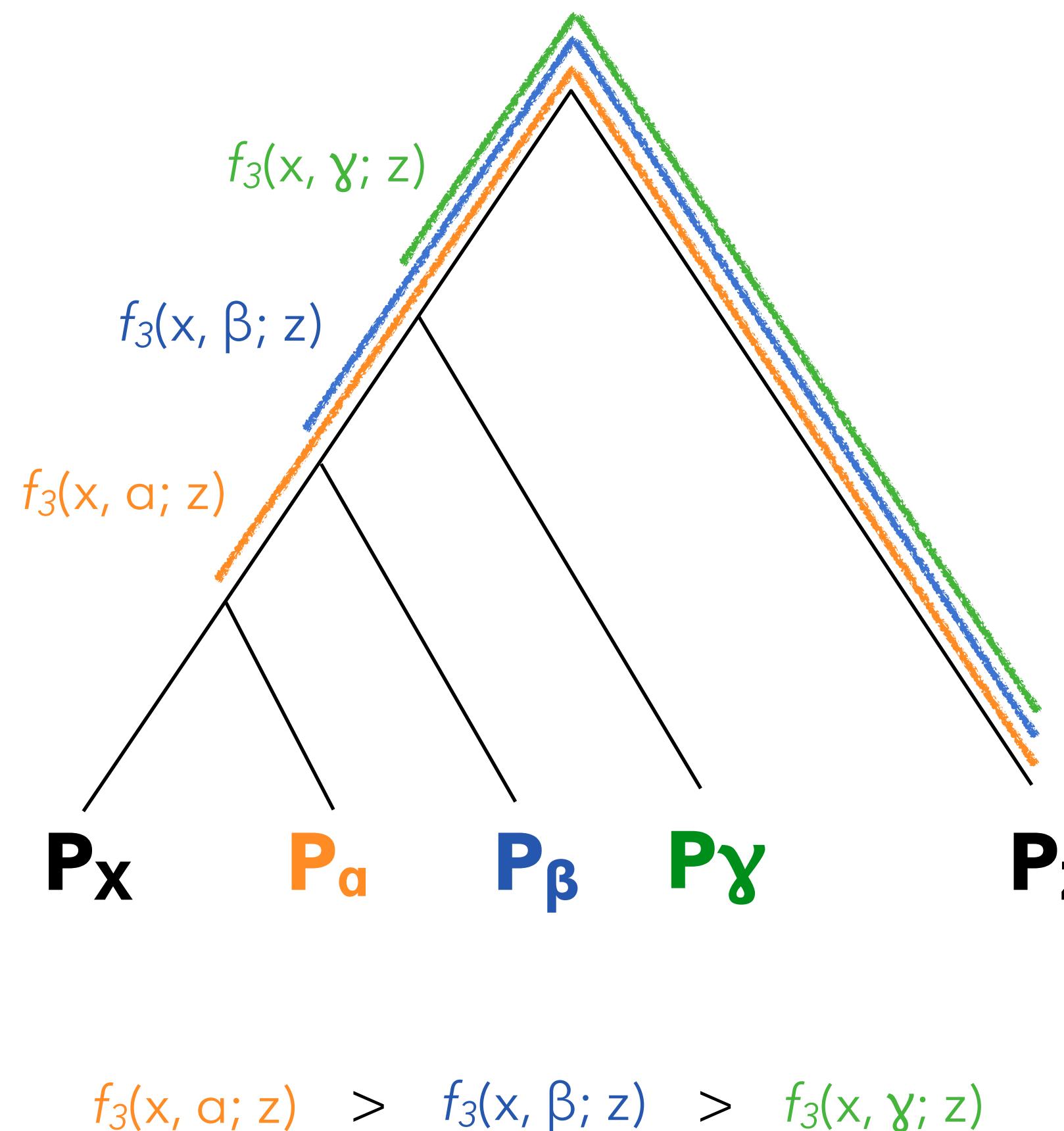


How do traditional popgen statistics work? – f -statistics

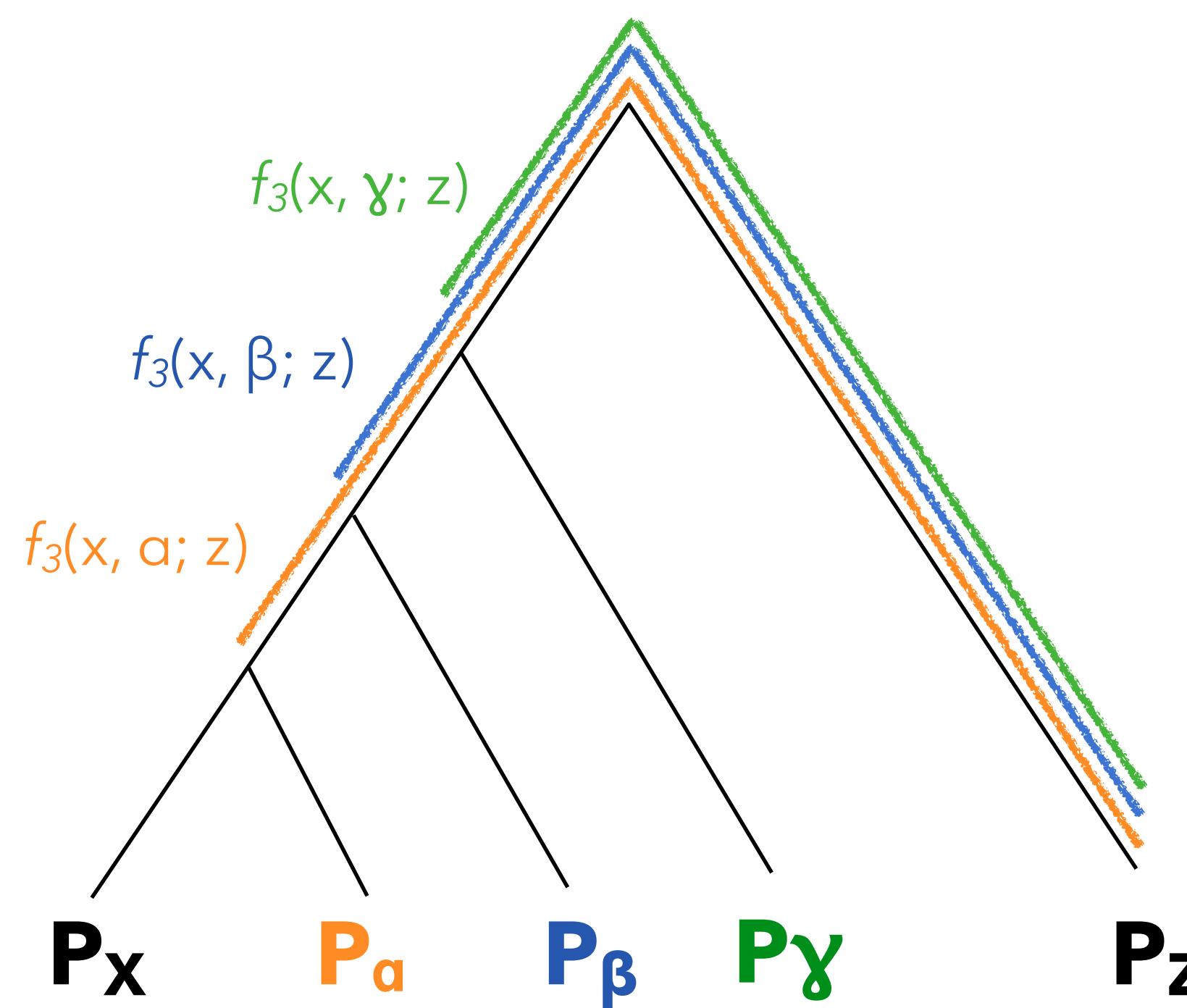


$$f_3(x, \alpha; z) > f_3(x, \beta; z) > f_3(x, \gamma; z)$$

How do traditional popgen statistics work? — f -statistics



How do traditional popgen statistics work? – f -statistics

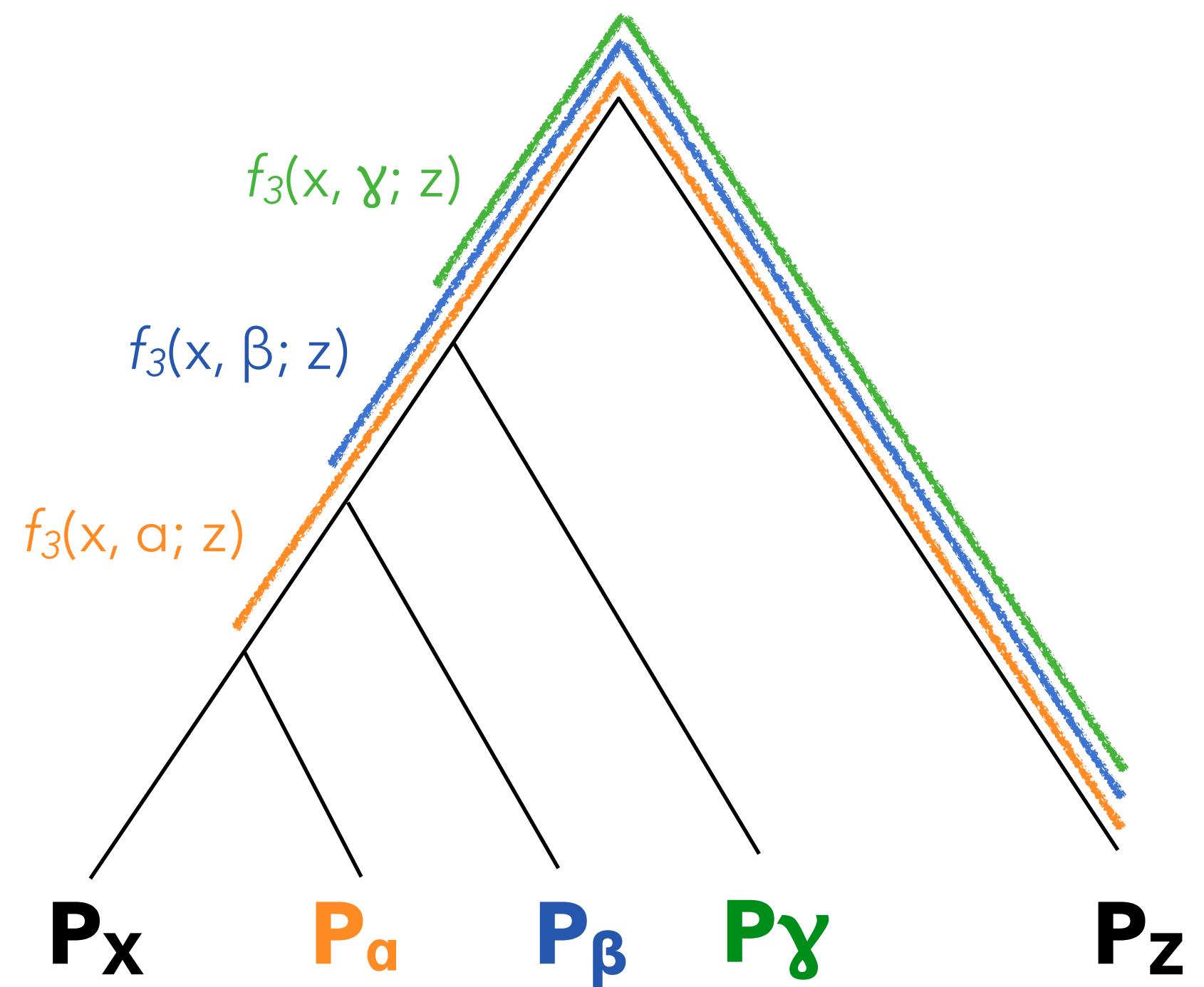


$$f_3(x, \alpha; z) > f_3(x, \beta; z) > f_3(x, \gamma; z)$$

#CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO	FORMAT	HG00096	HG00099	HG00101
12	60076	.	A	C	100	PASS	.	GT	1 0	0 0	0 0
12	60252	.	A	G	100	PASS	.	GT	0 0	0 0	0 1
12	60317	.	C	T	100	PASS	.	GT	0 0	0 0	0 0
12	60344	.	C	A	100	PASS	.	GT	0 0	0 0	0 0
12	60383	.	G	A	100	PASS	.	GT	0 0	0 0	0 0
12	60405	.	T	C	100	PASS	.	GT	0 0	0 0	0 0
12	60474	.	G	A	100	PASS	.	GT	0 0	0 1	0 1
12	60614	.	C	A	100	PASS	.	GT	0 0	0 0	0 0
12	60628	.	T	C	100	PASS	.	GT	0 0	0 0	0 0
12	60654	.	G	A	100	PASS	.	GT	0 0	0 0	0 0
12	61021	.	C	T	100	PASS	.	GT	0 0	0 0	0 0
12	61107	.	G	T	100	PASS	.	GT	0 0	0 0	0 0
12	61172	.	G	A	100	PASS	.	GT	0 0	0 0	0 0
12	61220	.	G	A	100	PASS	.	GT	0 0	0 0	0 1
12	61258	.	C	T	100	PASS	.	GT	0 0	0 0	0 1
12	61272	.	T	C	100	PASS	.	GT	0 0	0 0	0 0
12	61329	.	C	T	100	PASS	.	GT	0 0	0 0	0 0
12	61341	.	G	A	100	PASS	.	GT	0 0	0 1	0 1
12	61368	.	C	T	100	PASS	.	GT	0 0	0 1	0 1
12	61392	.	T	A	100	PASS	.	GT	0 0	0 0	0 0
12	61405	.	G	C	100	PASS	.	GT	0 0	0 0	0 0
12	61411	.	C	A	100	PASS	.	GT	0 0	0 0	0 0
12	61416	.	G	A	100	PASS	.	GT	0 0	0 0	0 0
12	61422	.	C	T	100	PASS	.	GT	0 0	0 0	0 0
12	61476	.	C	G	100	PASS	.	GT	0 0	0 0	0 0
12	61510	.	G	A	100	PASS	.	GT	0 0	0 0	0 0
12	61516	.	C	T	100	PASS	.	GT	0 0	0 0	0 0
12	61552	.	C	T	100	PASS	.	GT	0 0	0 0	0 0
12	61604	.	T	G	100	PASS	.	GT	0 0	0 0	0 0
12	61687	.	G	A	100	PASS	.	GT	1 0	0 1	0 1
12	61700	.	C	T	100	PASS	.	GT	0 0	0 0	0 0

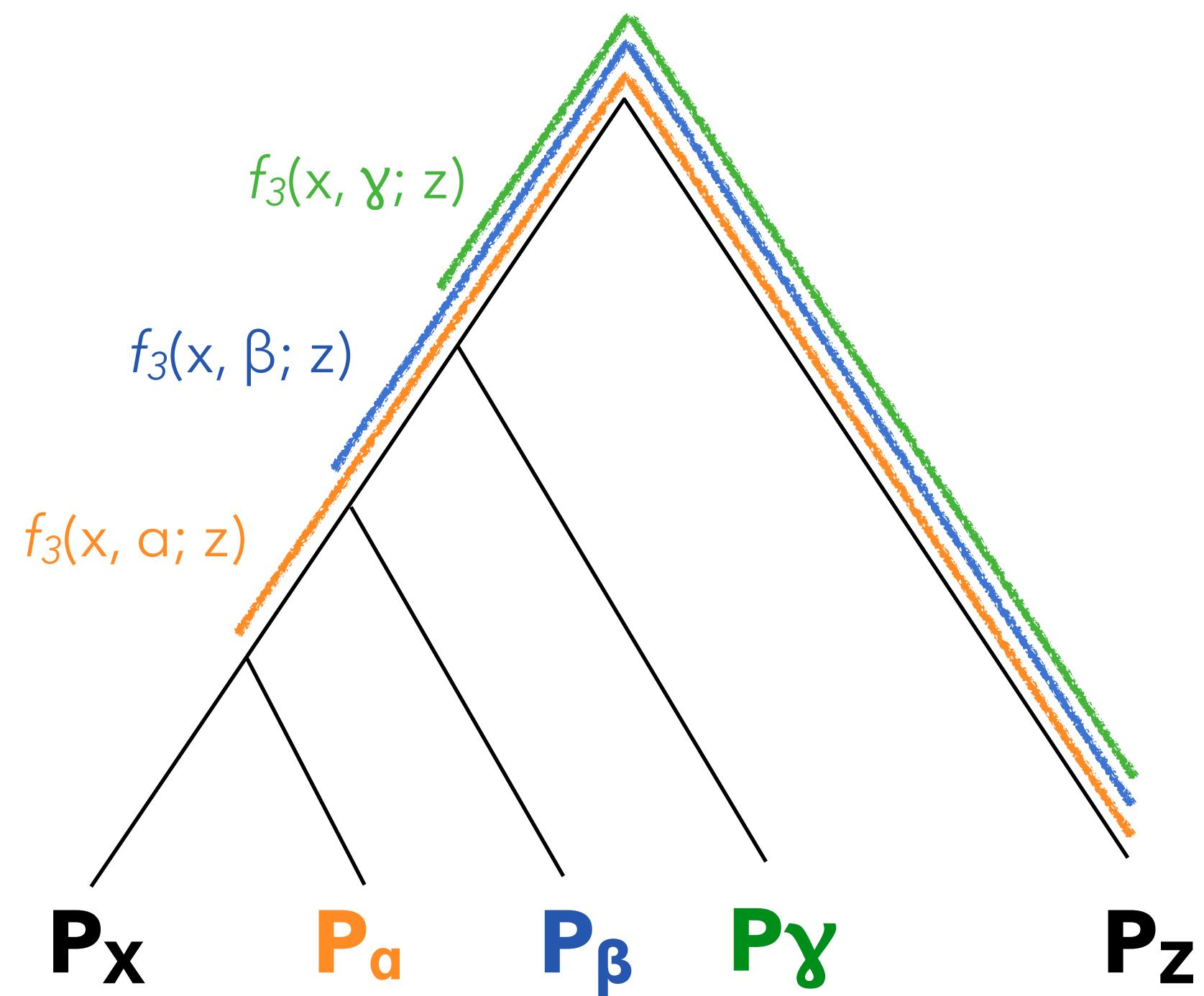
This evaluation of "tree-ness" requires lots of time, disk space, and memory.

With a tree sequence, computing statistics is trivial...

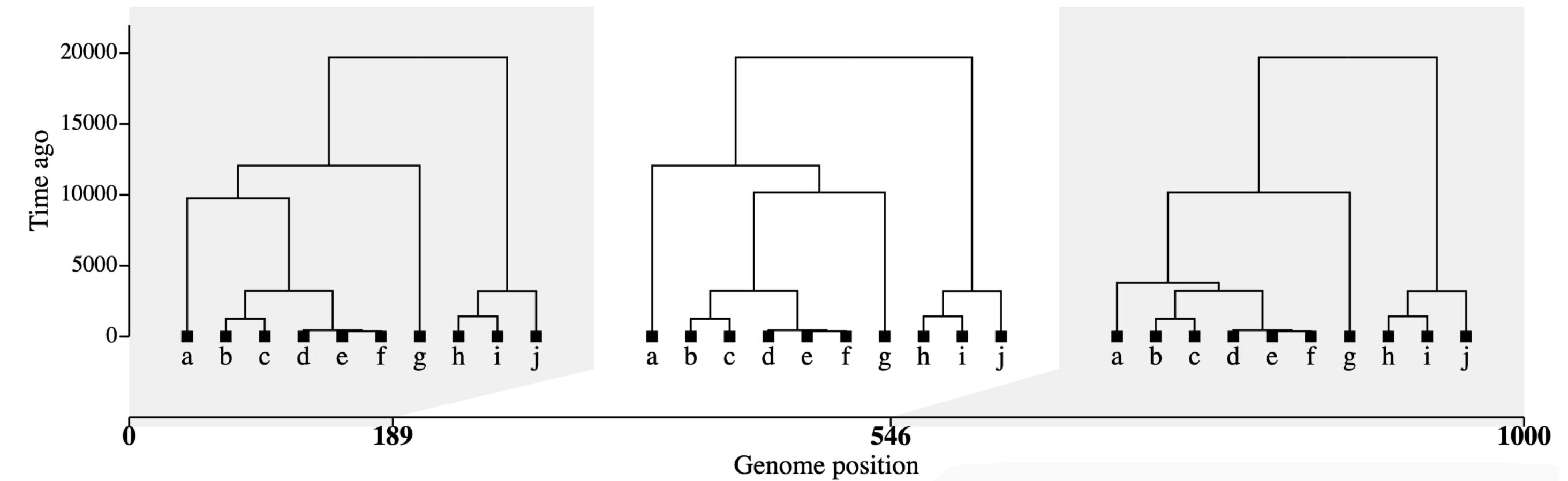


$$f_3(x, \alpha; z) > f_3(x, \beta; z) > f_3(x, \gamma; z)$$

With a tree sequence, computing statistics is trivial... ... because we have all trees to begin with!

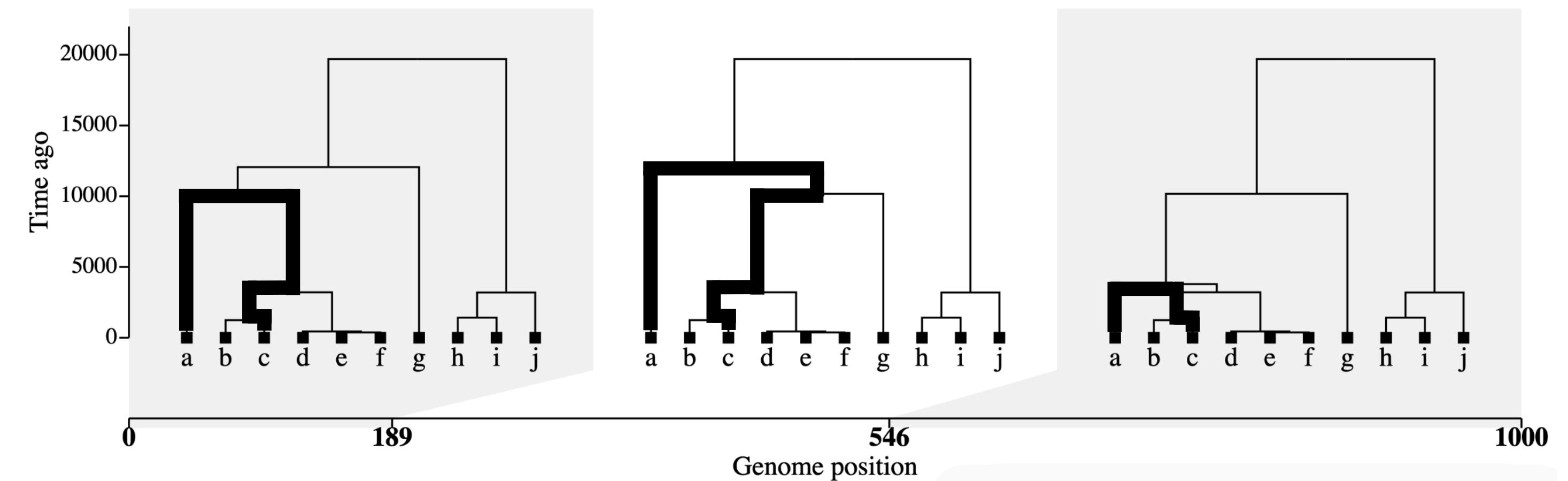
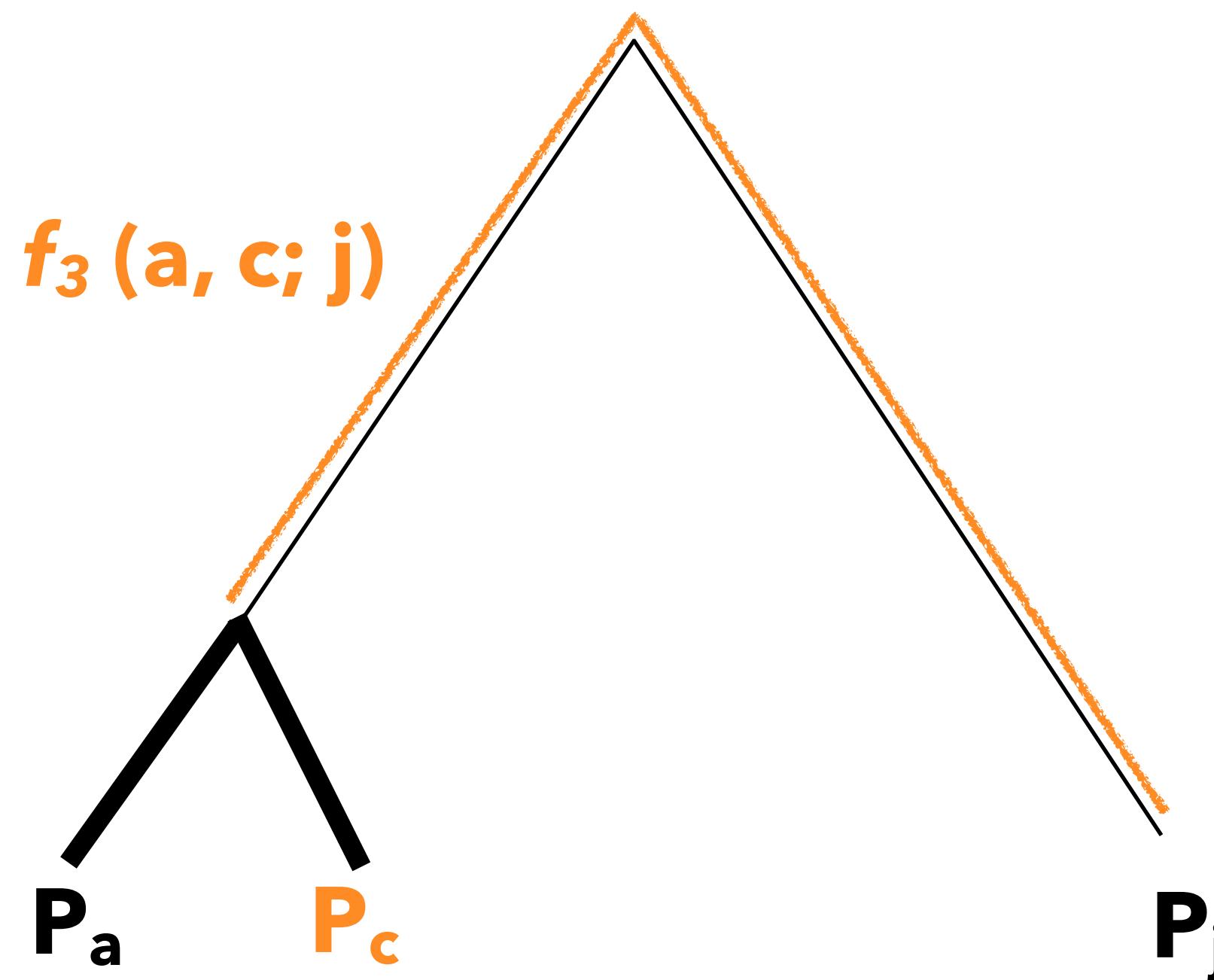


$$f_3(x, \alpha; z) > f_3(x, \beta; z) > f_3(x, \gamma; z)$$



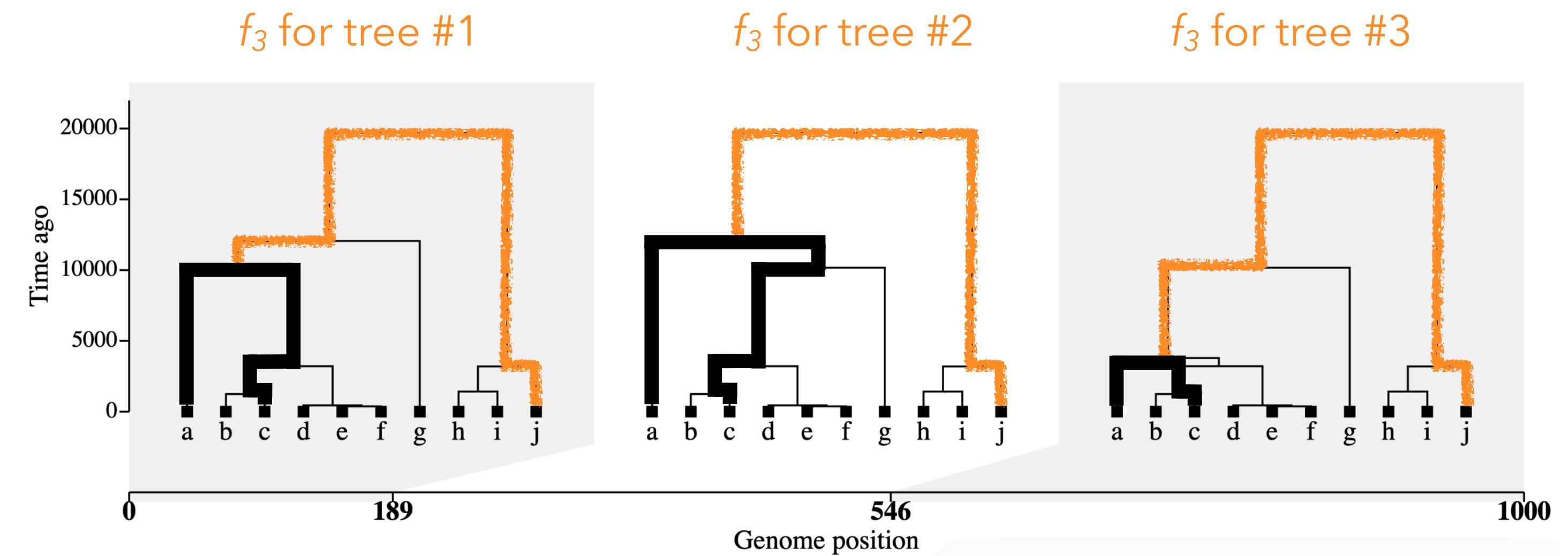
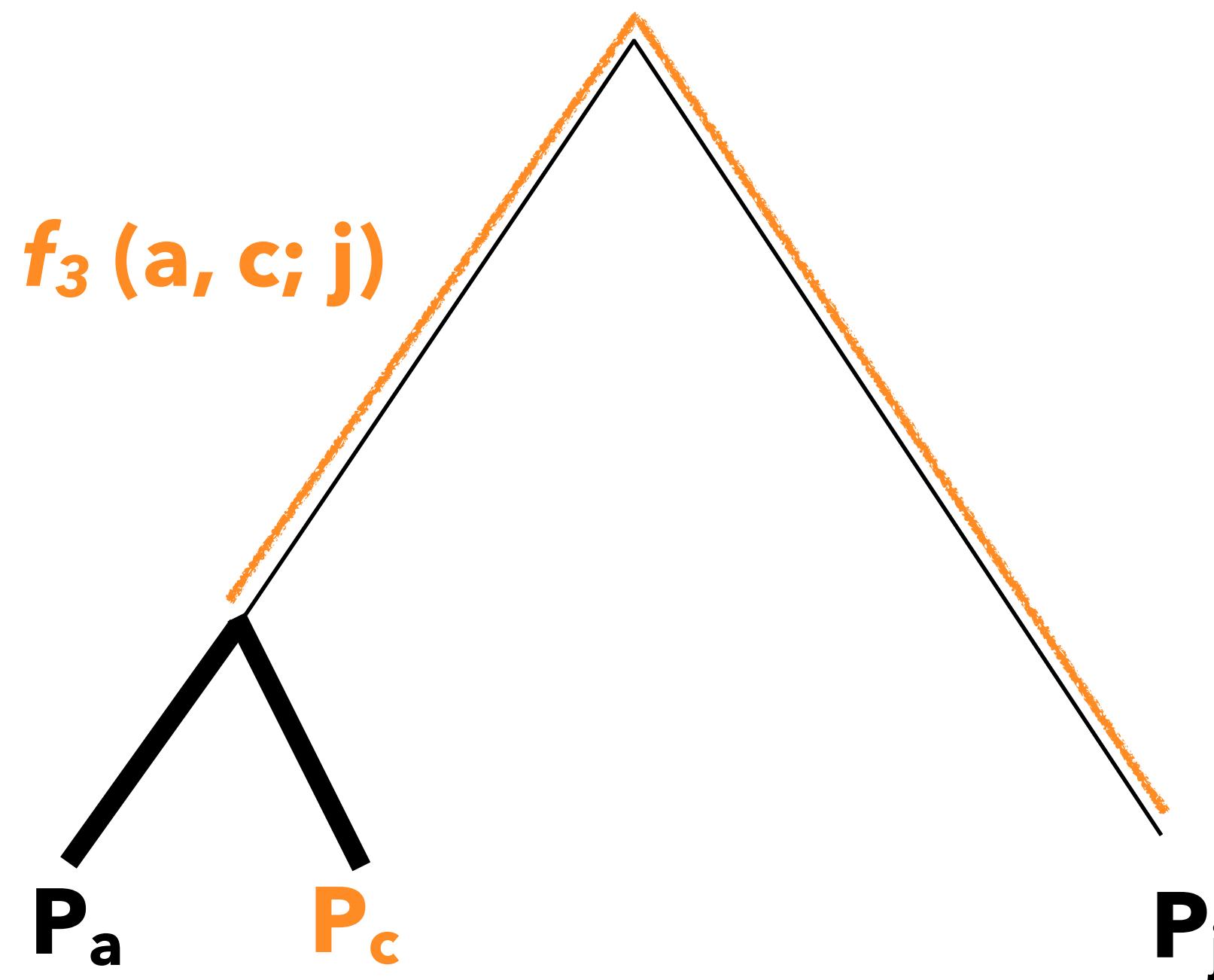
With a tree sequence, computing statistics is trivial...

... because we have all trees to begin with!



With a tree sequence, computing statistics is trivial...

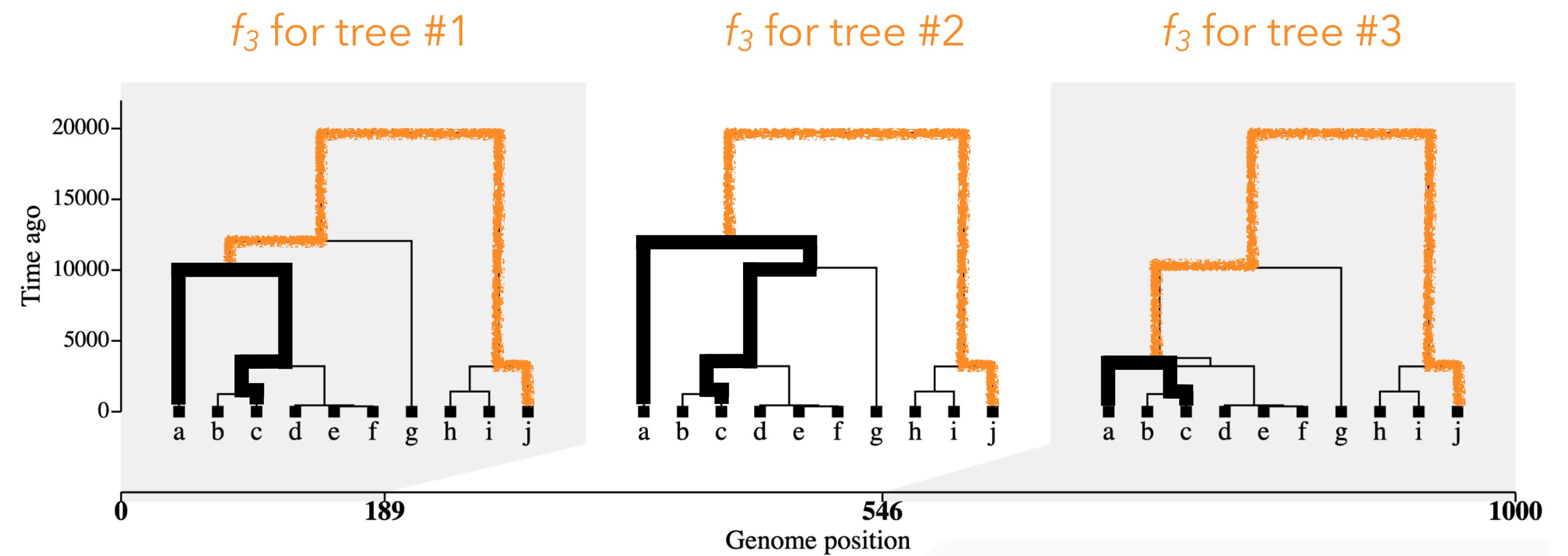
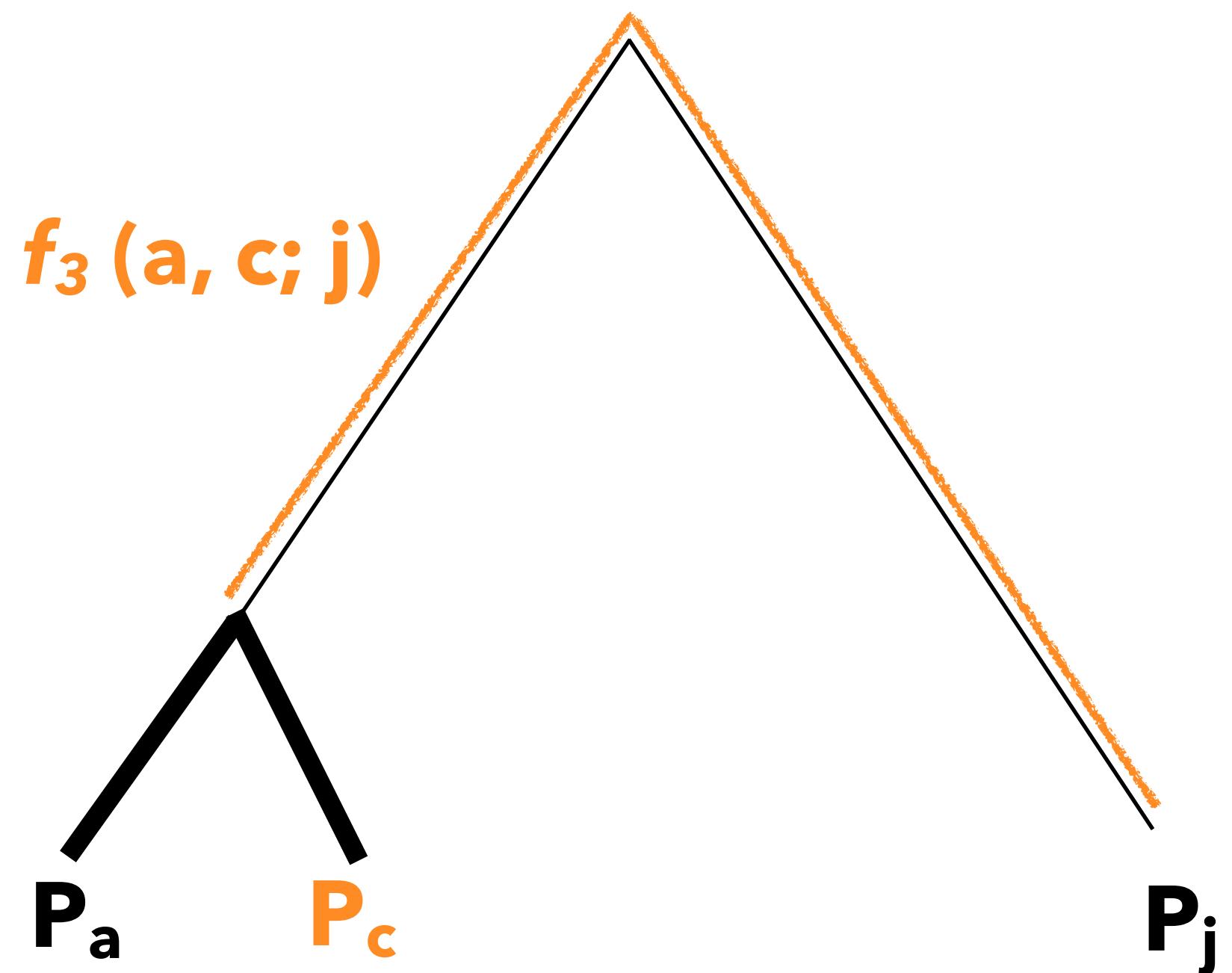
... because we have all trees to begin with!



Then aggregate f_3 values across all trees.

With a tree sequence, computing statistics is trivial...

... because we have all trees to begin with!



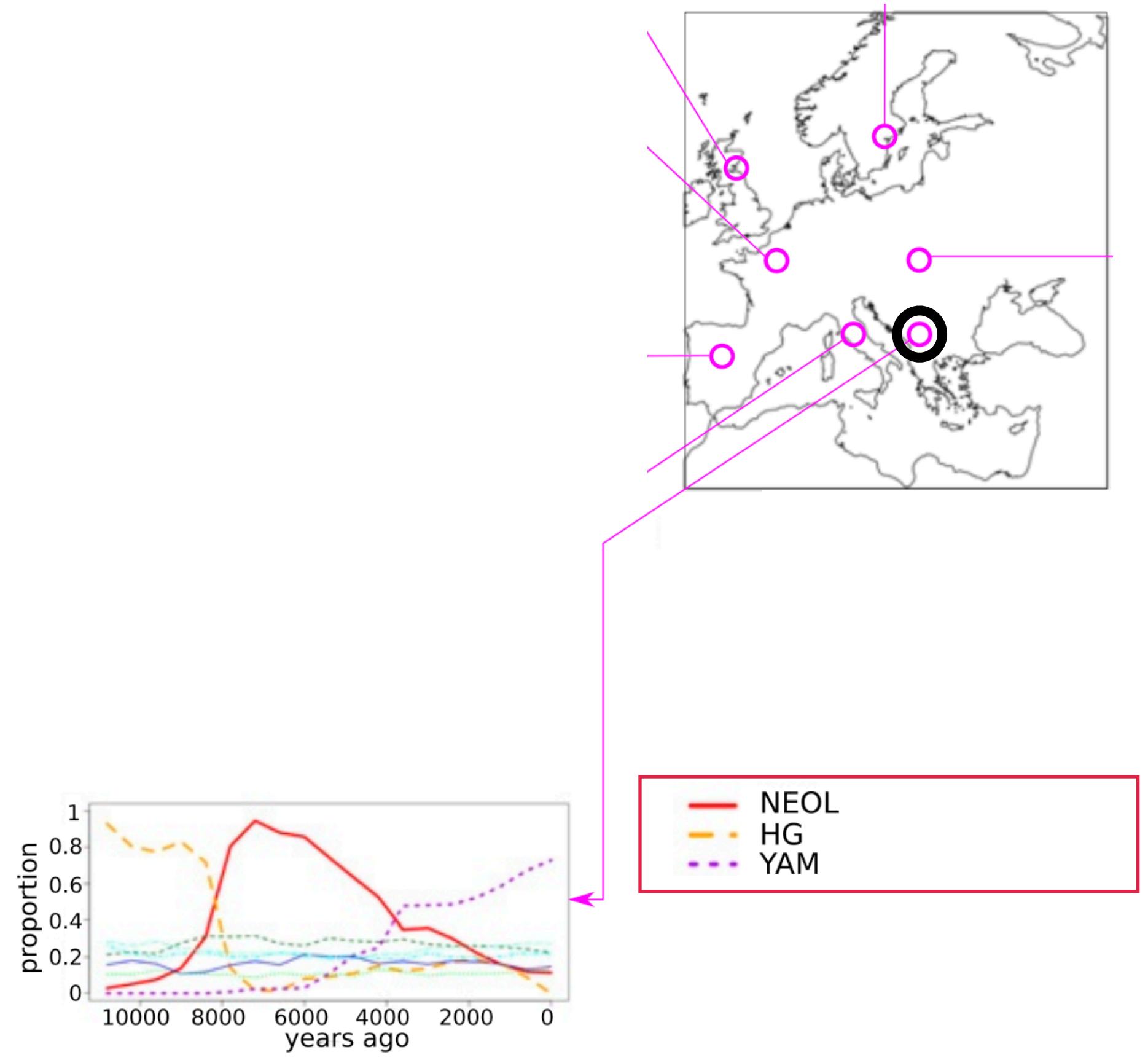
Then aggregate f_3 values across all trees.

This is implemented in **tskit** (tskit.dev) & the R package **slendr** (slendr.net).

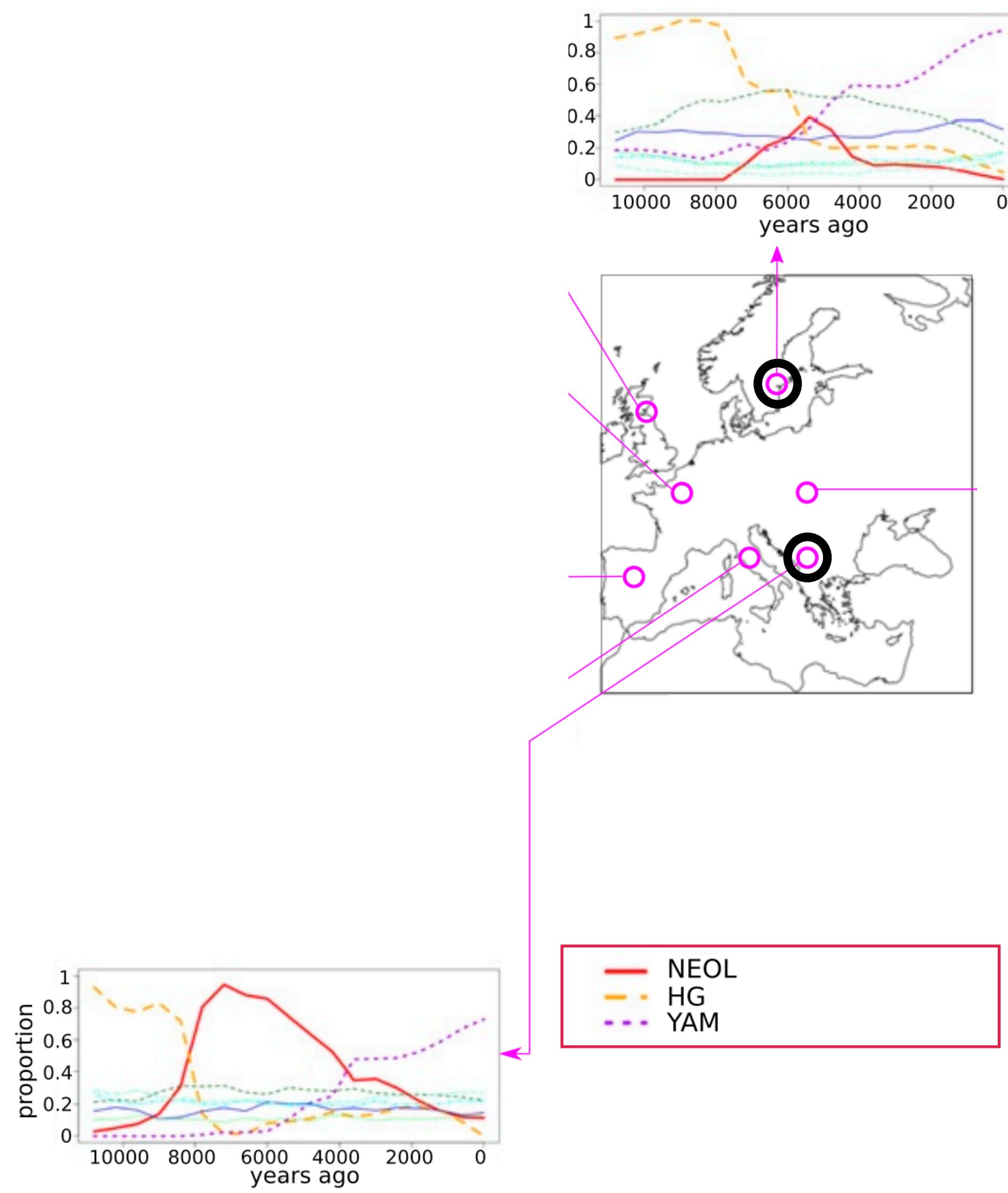
Were human migrations associated with changes in landscape?



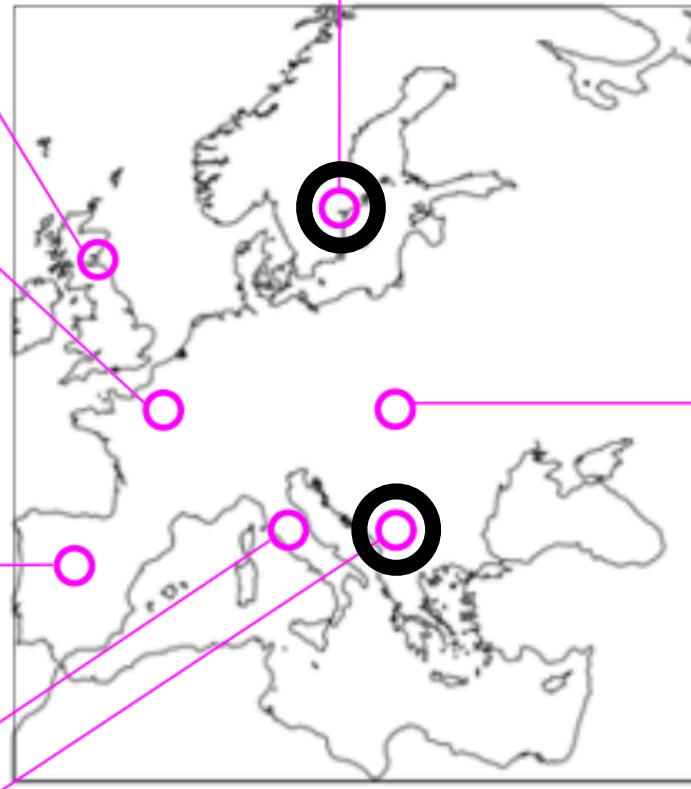
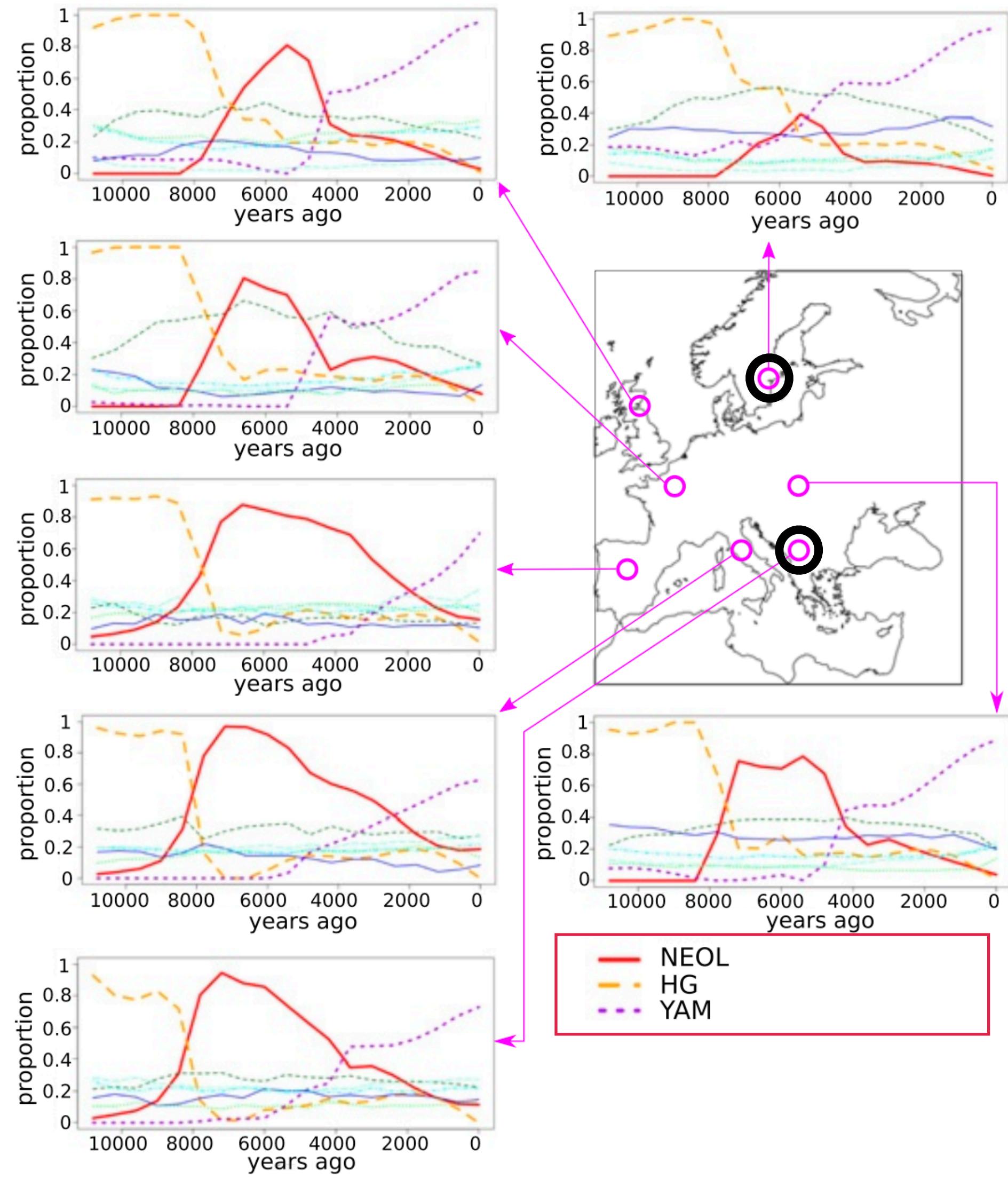
Were human migrations associated with changes in landscape?



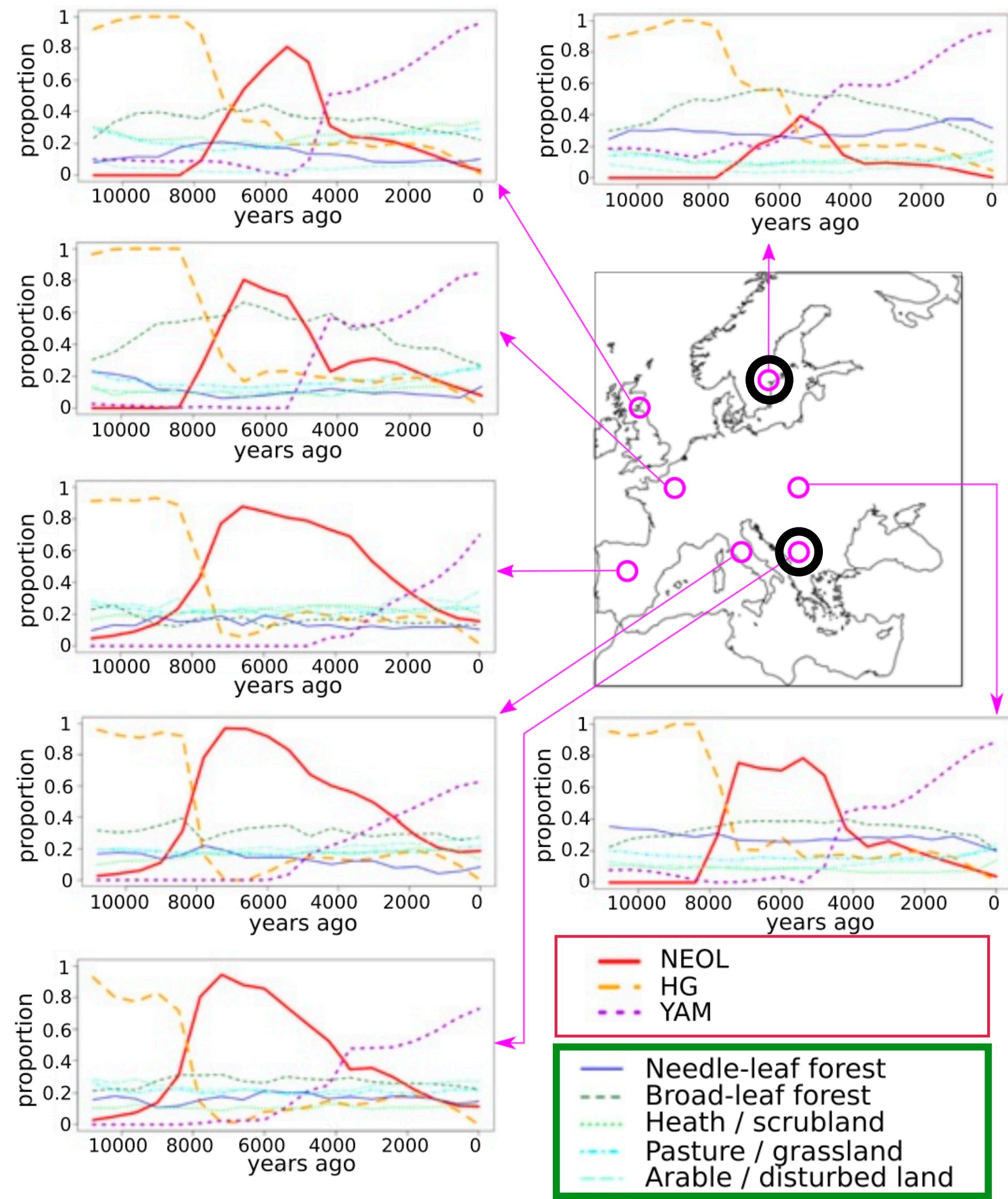
Were human migrations associated with changes in landscape?



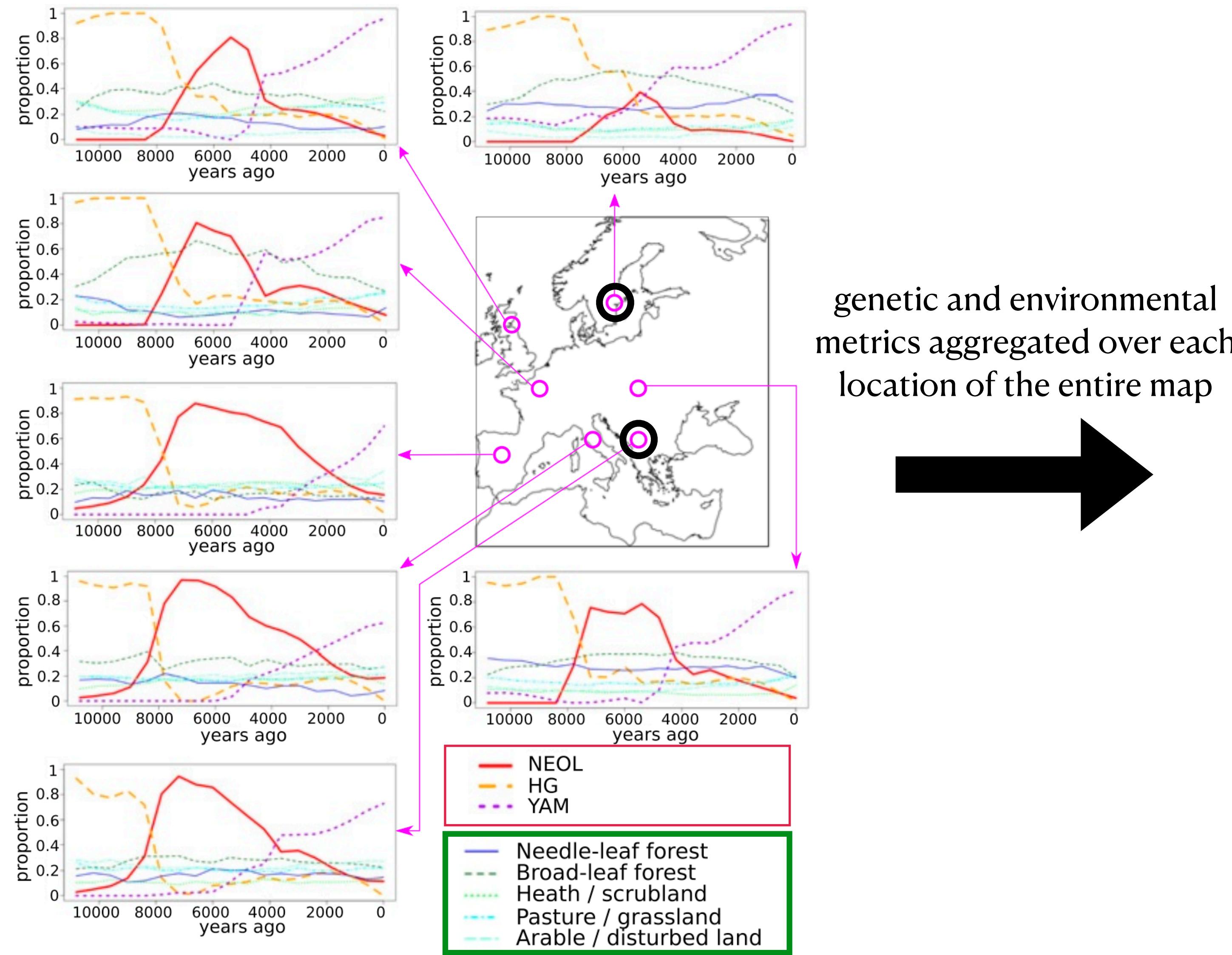
Were human migrations associated with changes in landscape?



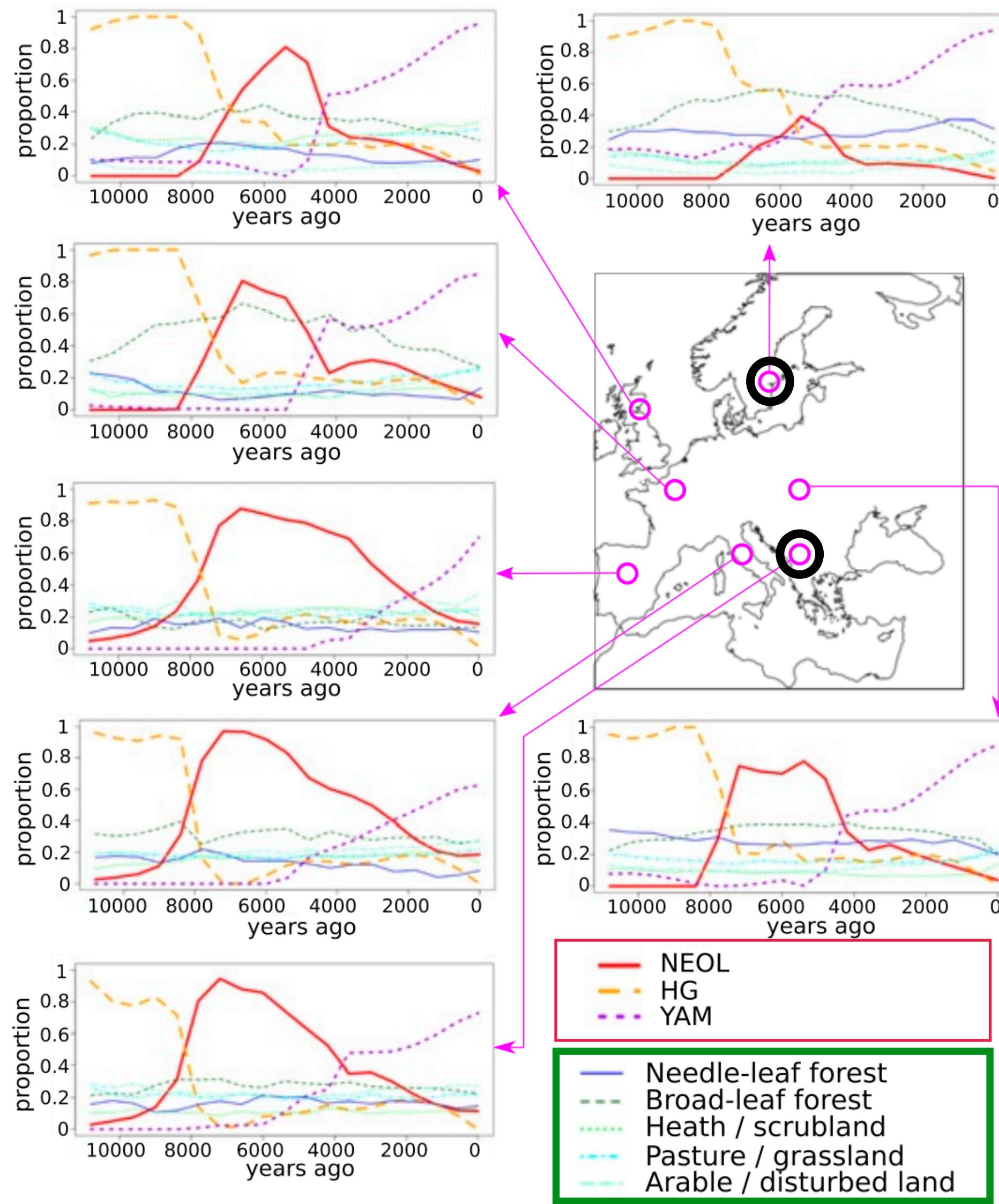
Were human migrations associated with changes in landscape?



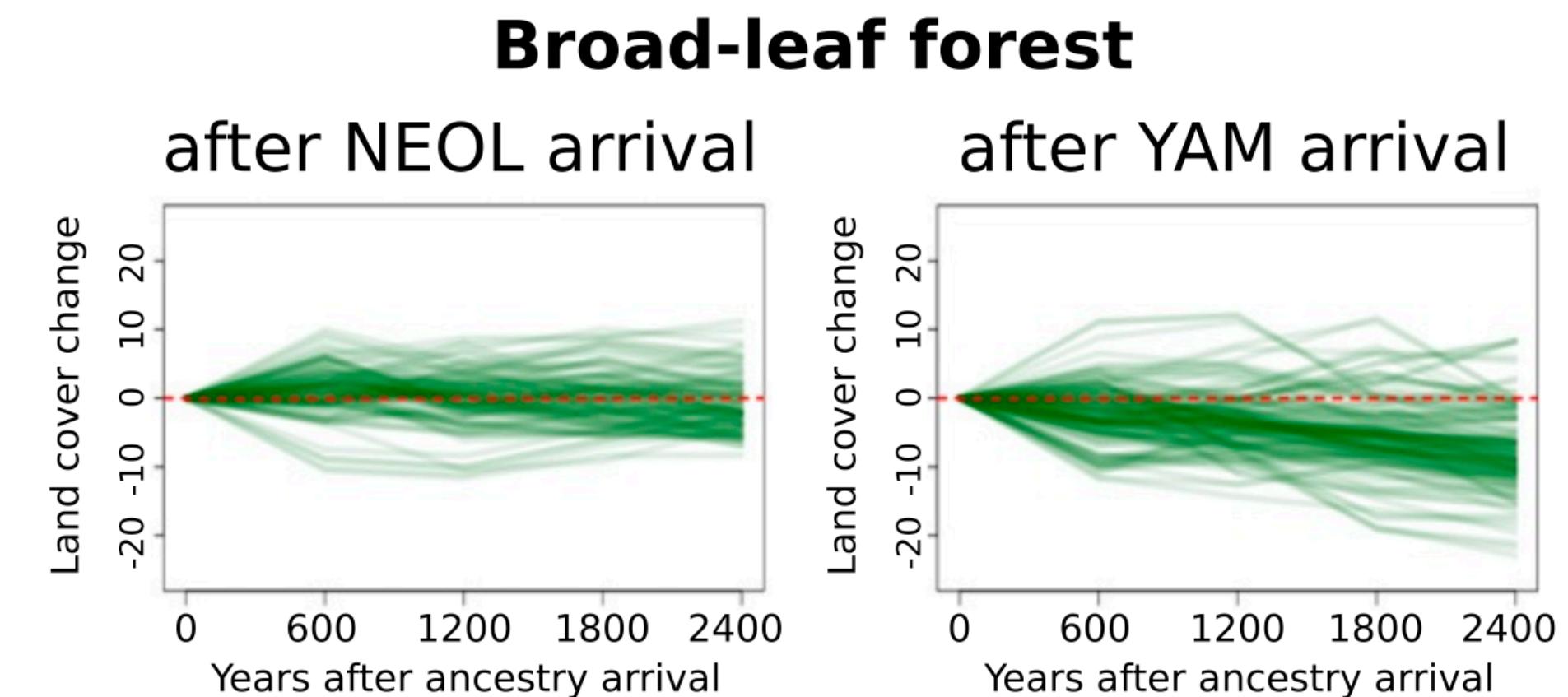
Were human migrations associated with changes in landscape?



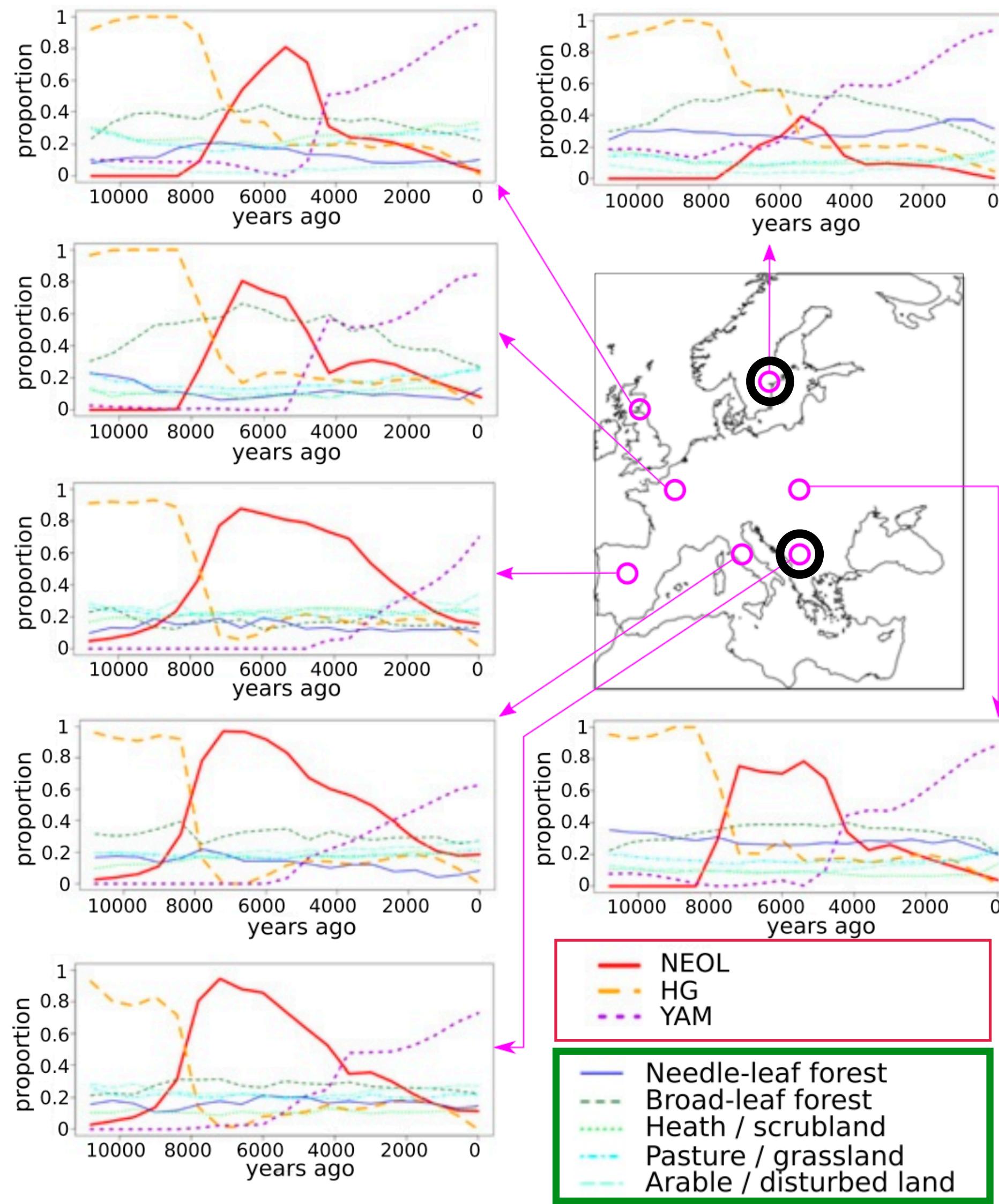
Yamnaya coincided with strong vegetation changes



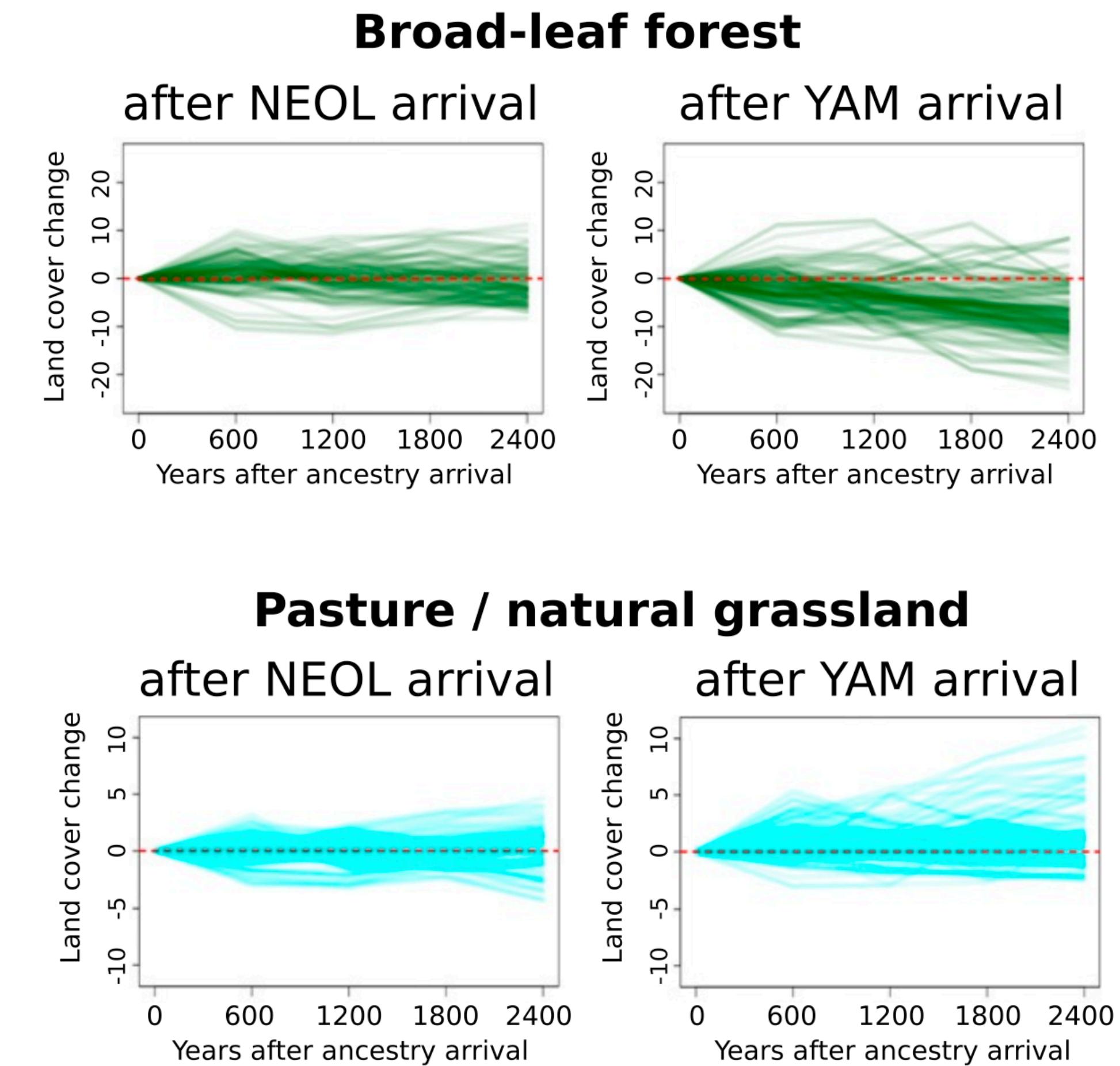
genetic and environmental
metrics aggregated over each
location of the entire map



Yamnaya coincided with strong vegetation changes



genetic and environmental
metrics aggregated over each
location of the entire map



Every other statistic

— nucleotide diversity, IBD, split times, N_e , D statistic, ... —

works in the same way.

Every other statistic

— nucleotide diversity, IBD, split times, N_e , D statistic, ... —

works in the same way.

Inferring whole-genome histories in large population datasets

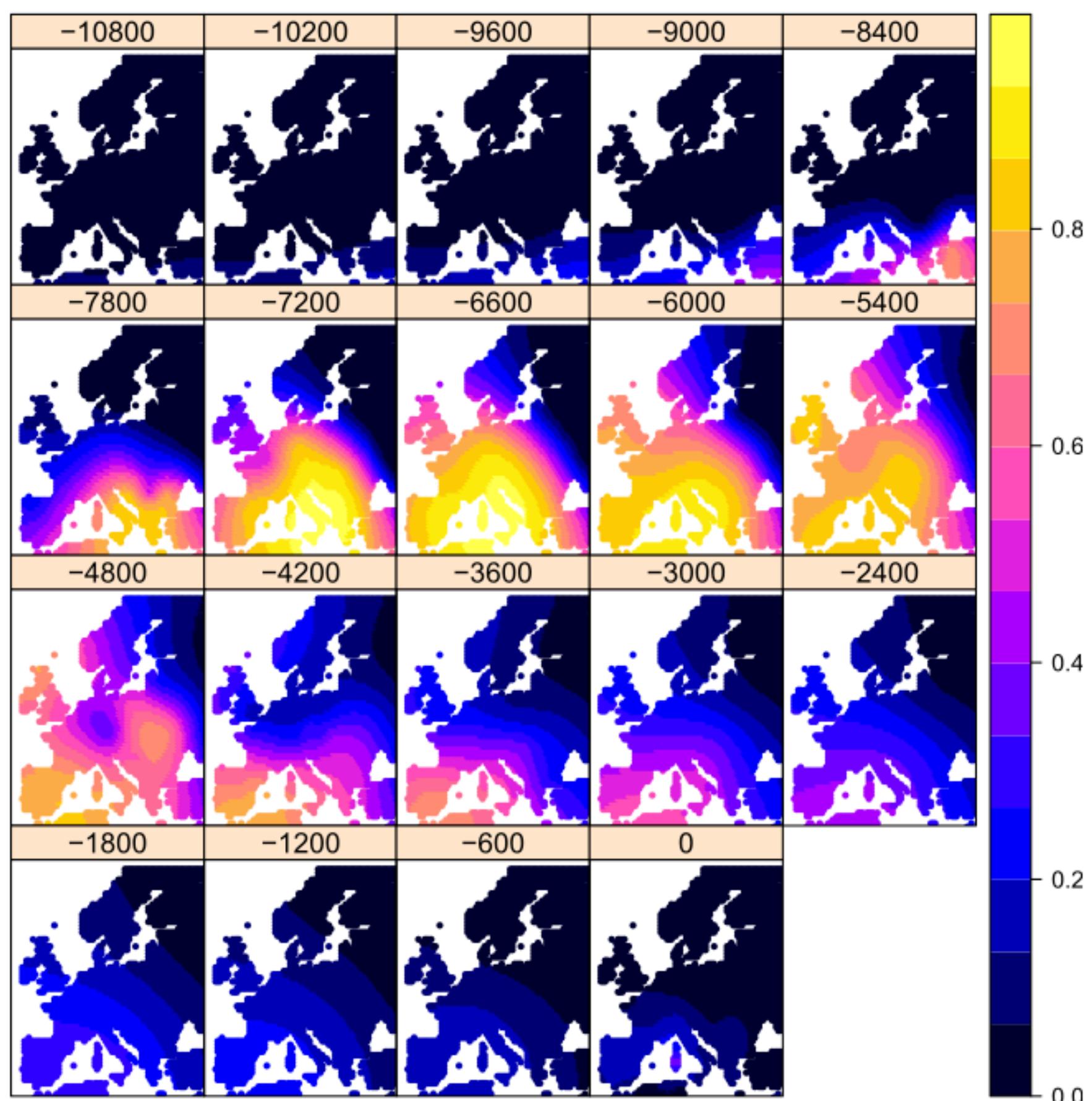
Jerome Kelleher^{ID}*, Yan Wong, Anthony W. Wohns^{ID}, Chaimaa Fadil^{ID}, Patrick K. Albers^{ID}
and Gil McVean^{ID}

A method for genome-wide genealogy estimation for thousands of samples

Leo Speidel^{ID}¹, Marie Forest², Sinan Shi¹ and Simon R. Myers^{ID}^{1,3*}

Interpolated ancestry maps

"NEOL" ancestry



"YAM" ancestry

