

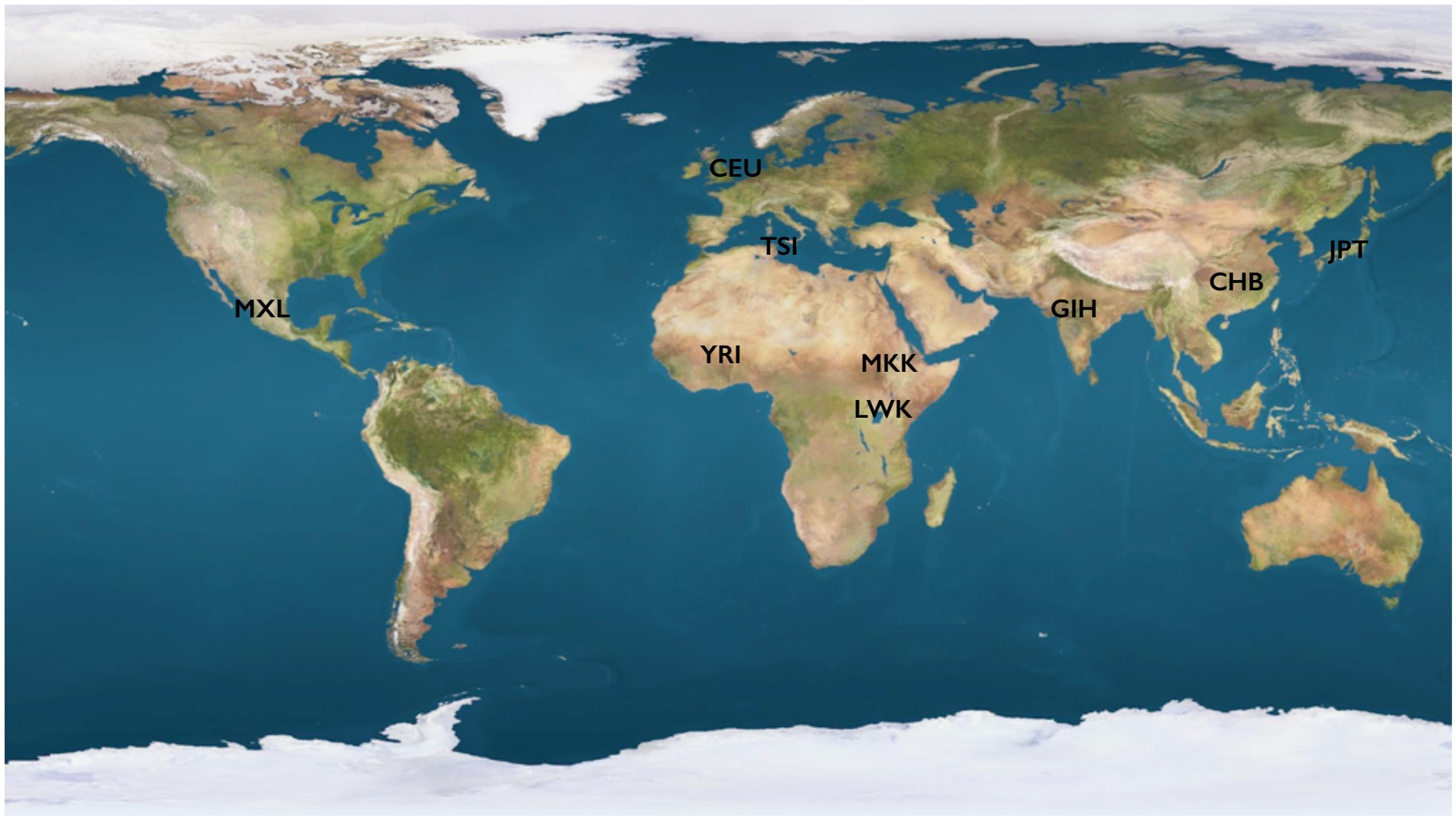
Inferring human population history from multiple genome sequences

Stephan Schiffels

Postdoctoral Fellow with Richard Durbin,
Wellcome Trust Sanger Institute, Cambridge, UK

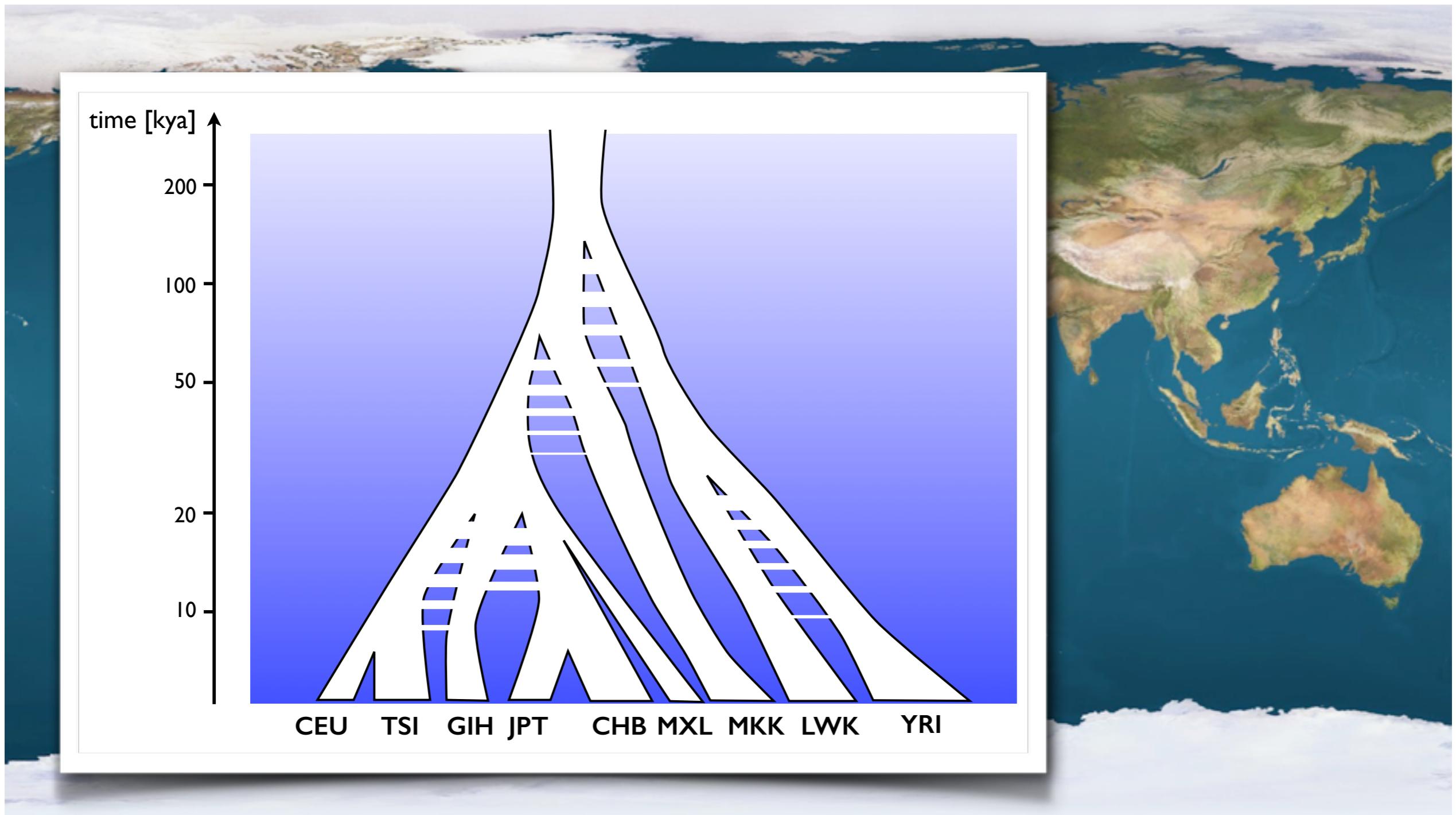


From genome sequences to human history



[Sequence data from Complete Genomics]

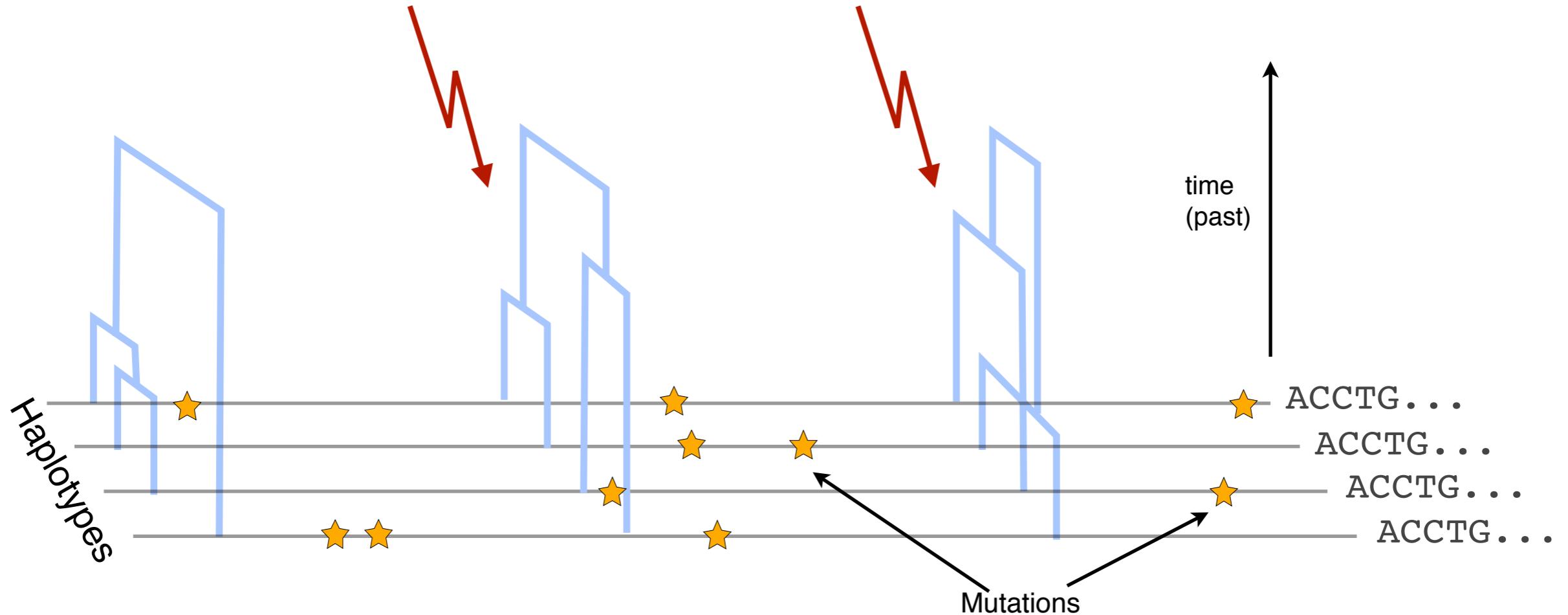
From genome sequences to human history



[Sequence data from Complete Genomics]

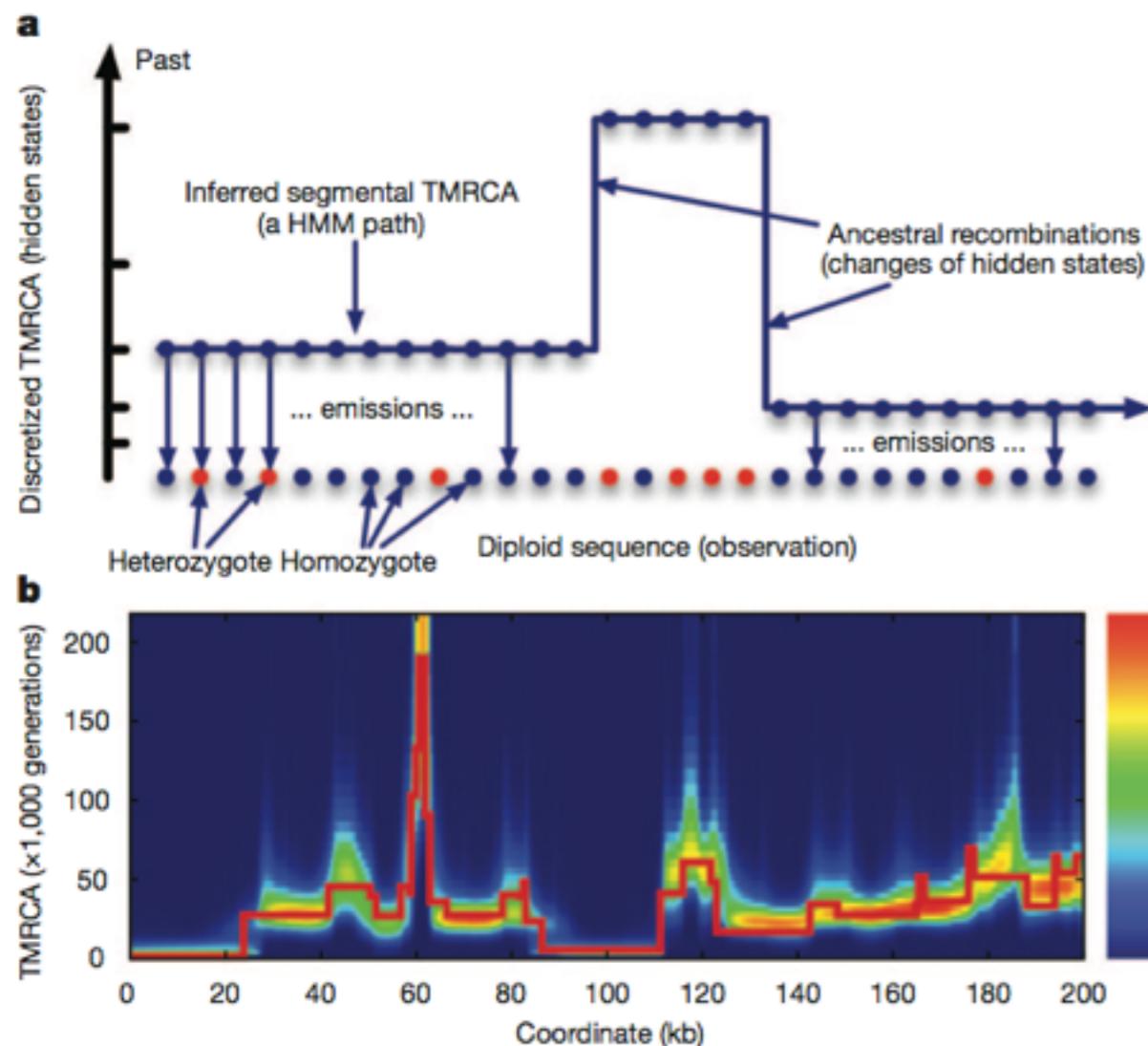
How are sequences related?

Ancestral recombinations change trees along the sequences

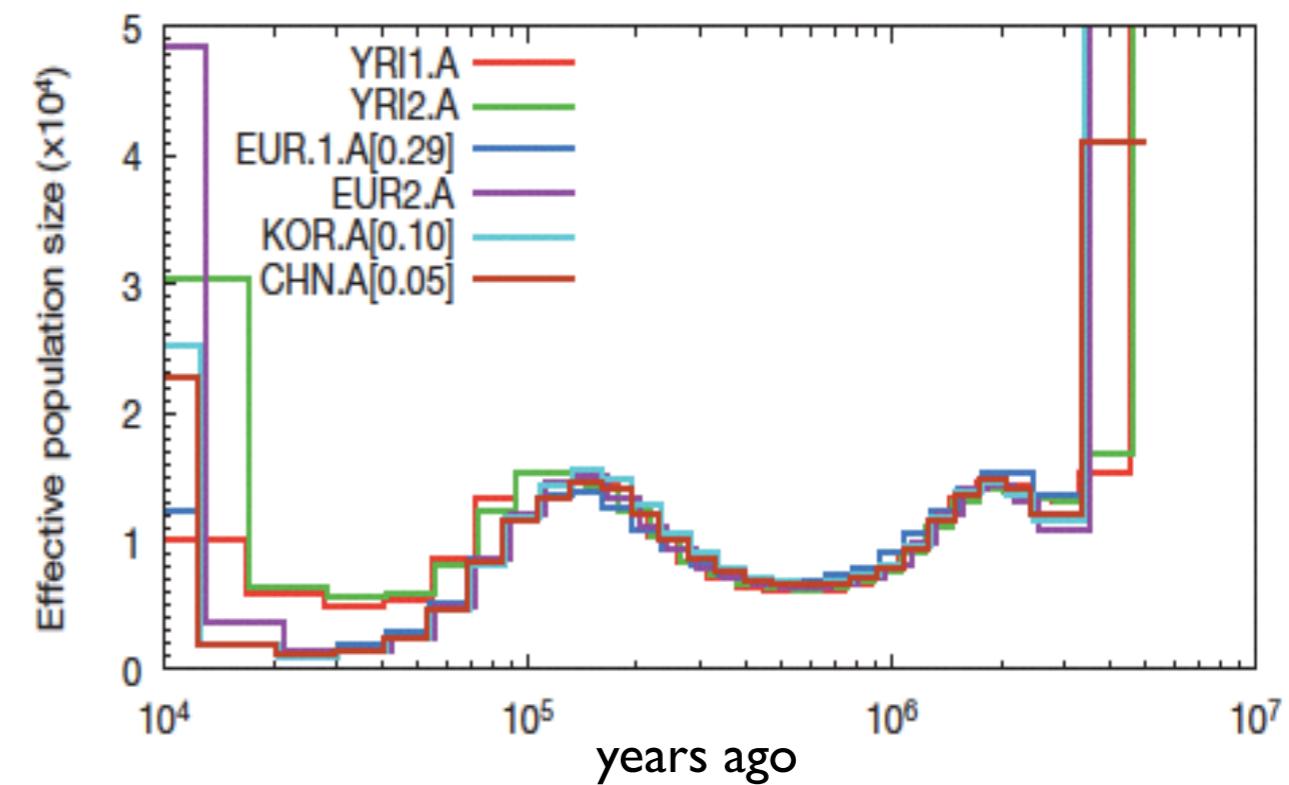
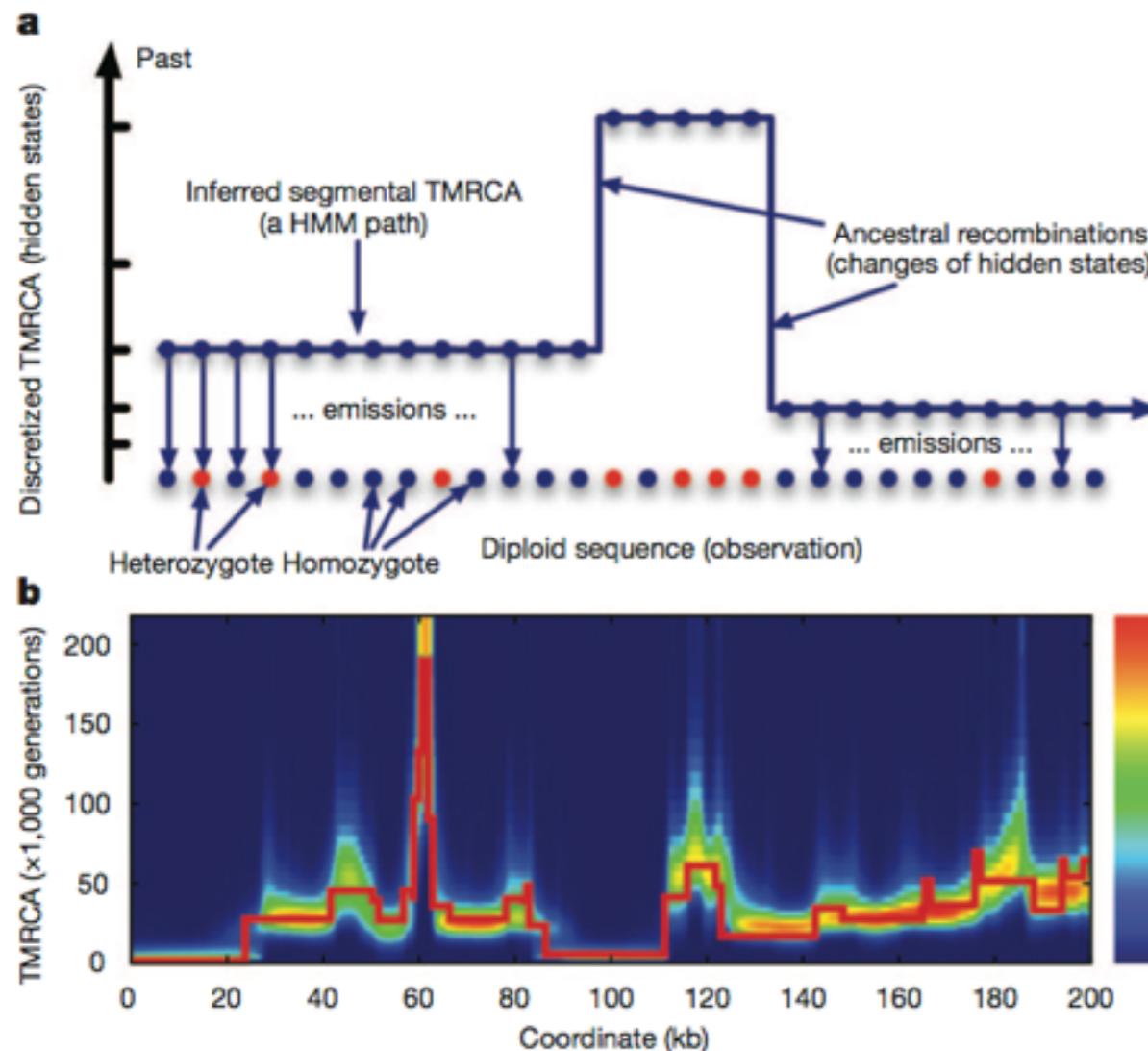


Problem: Estimate trees only from observed mutations

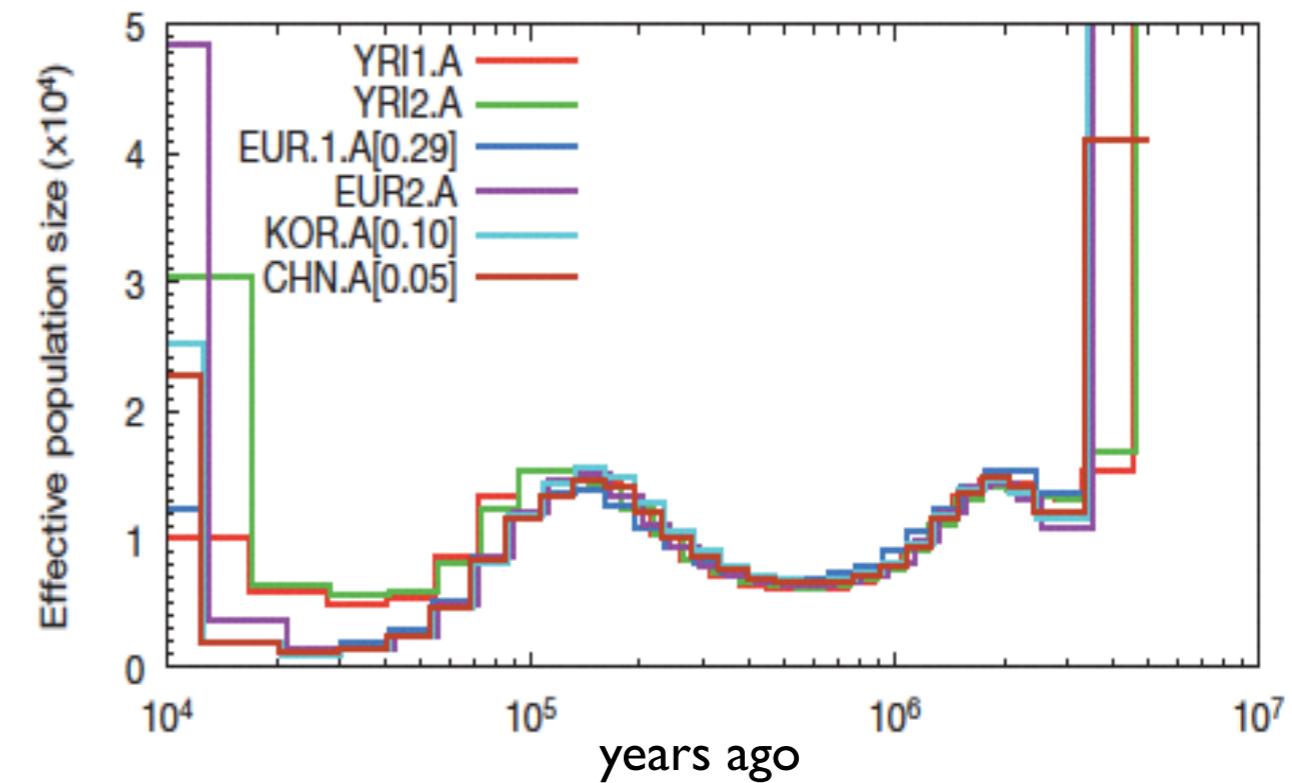
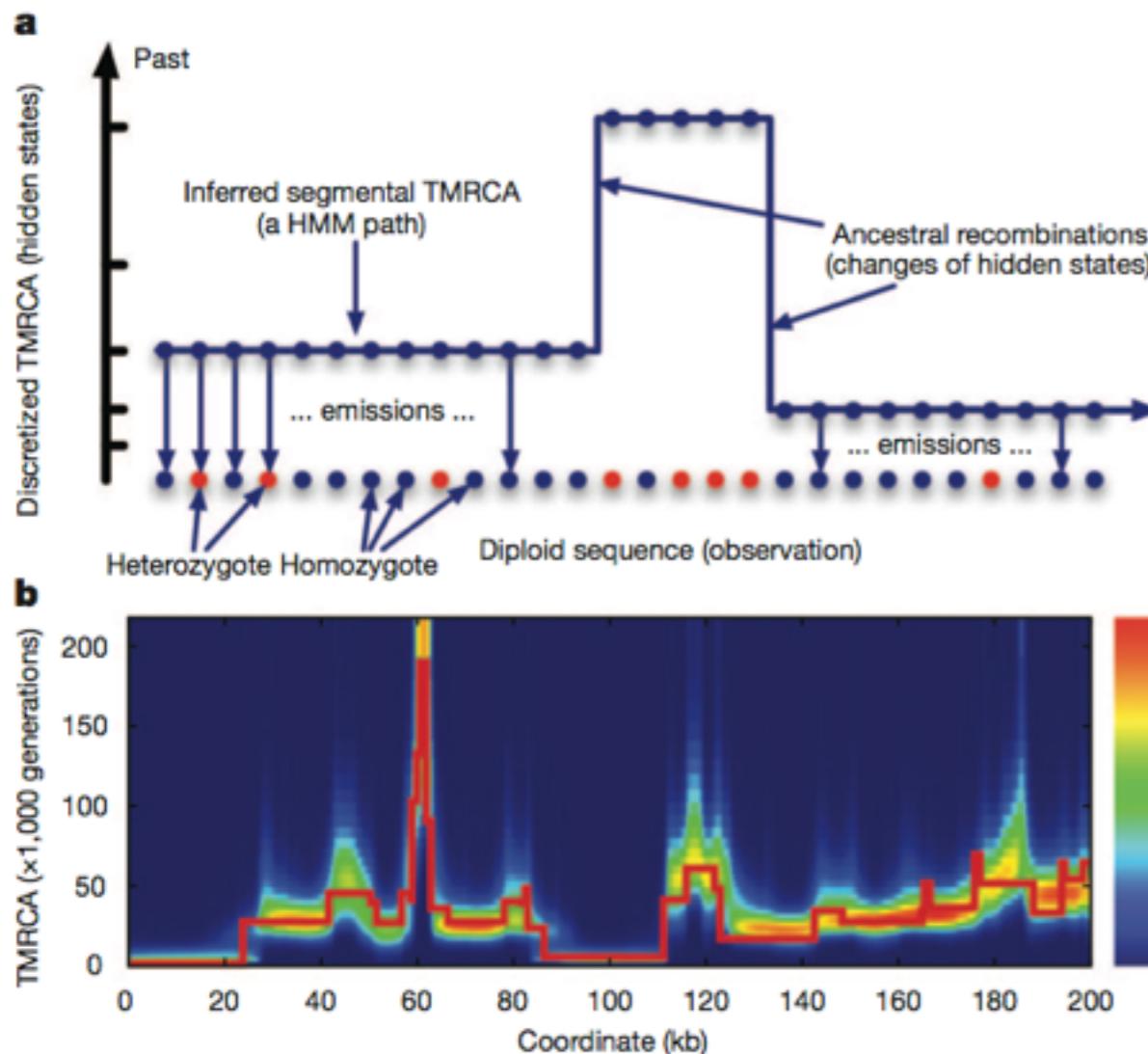
Previous work: PSMC



Previous work: PSMC



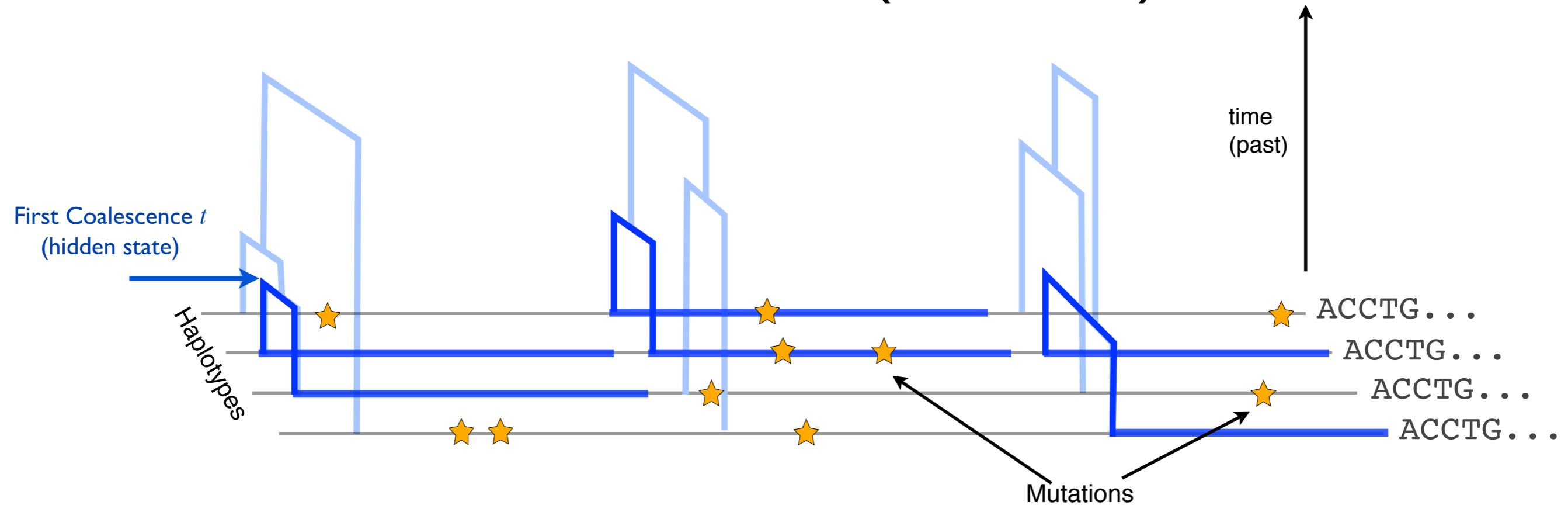
Previous work: PSMC



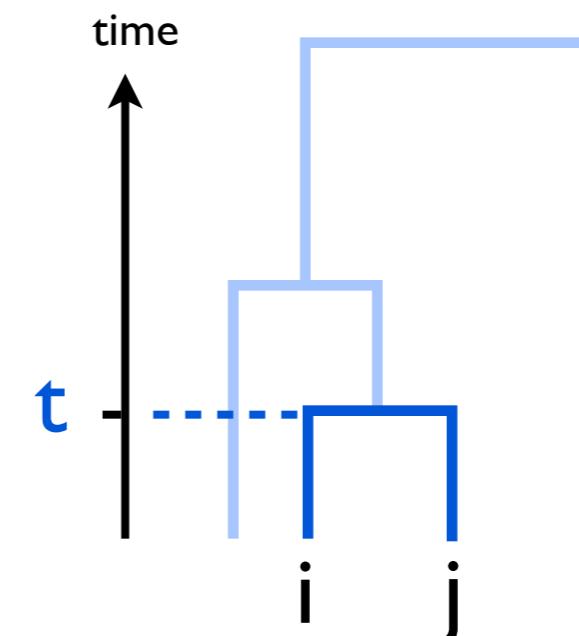
But: Inference from only two sequences limited to times beyond 20kya. Also: population splits difficult to model

[Li and Durbin, 2011]

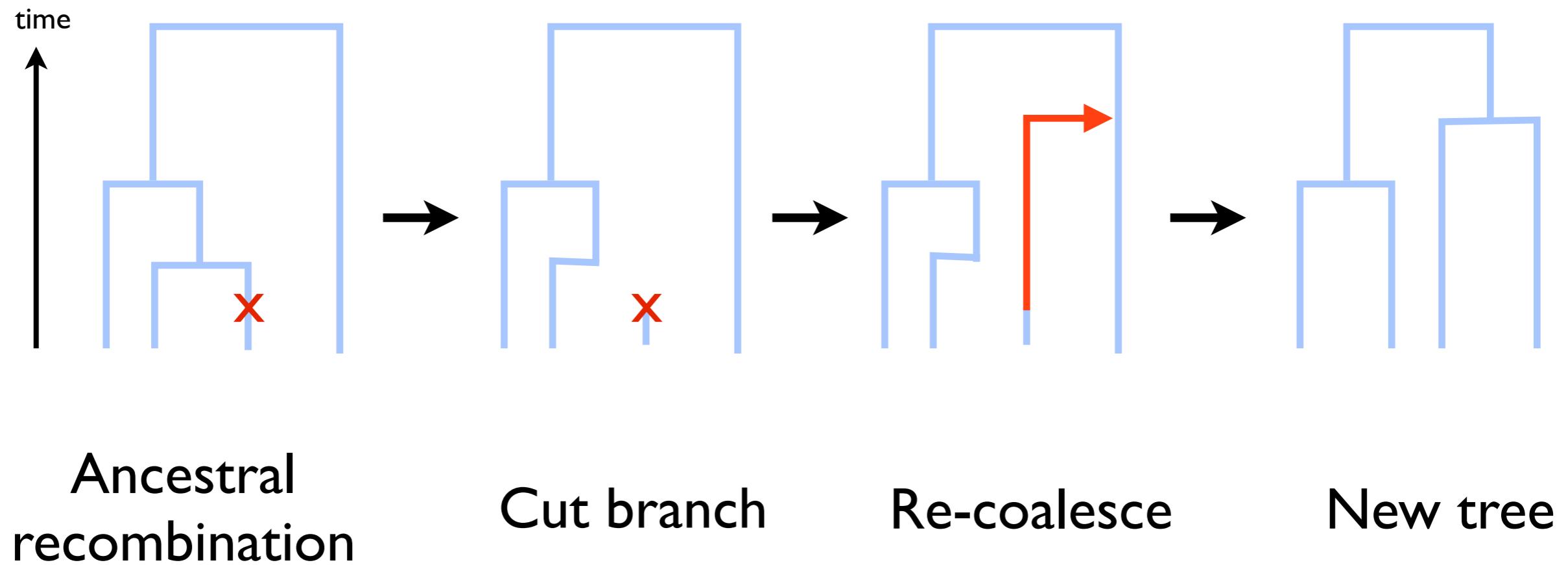
Multiple sequentially Markovian Coalescent (MSMC)



MSMC hidden state:
triple (t, i, j)



Effect of recombination on Genealogies



Ancestral
recombination

Cut branch

Re-coalesce

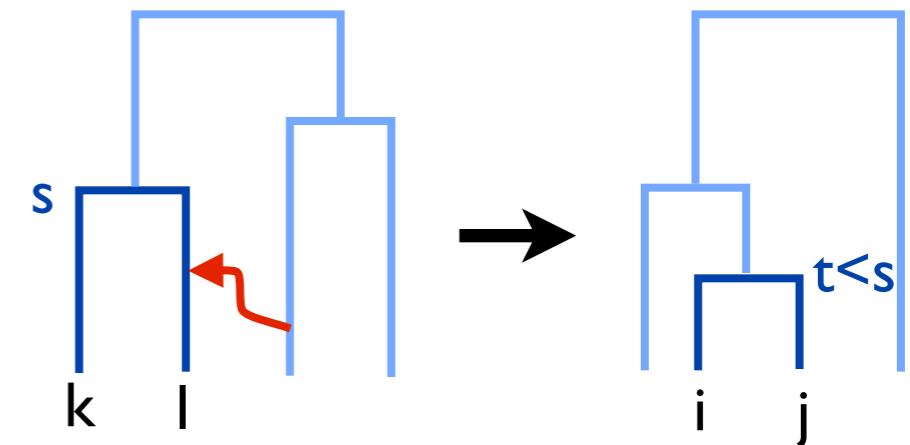
New tree

[*Sequentially Markovian Coalescent*, McVean and Cardin, 2005]

MSMC: state transitions

$(s, k, l) \rightarrow (t, i, j)$

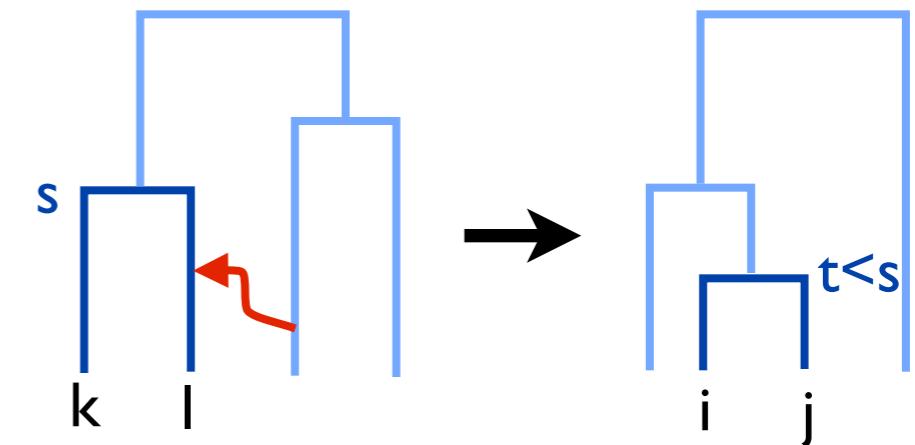
where $t < s$



MSMC: state transitions

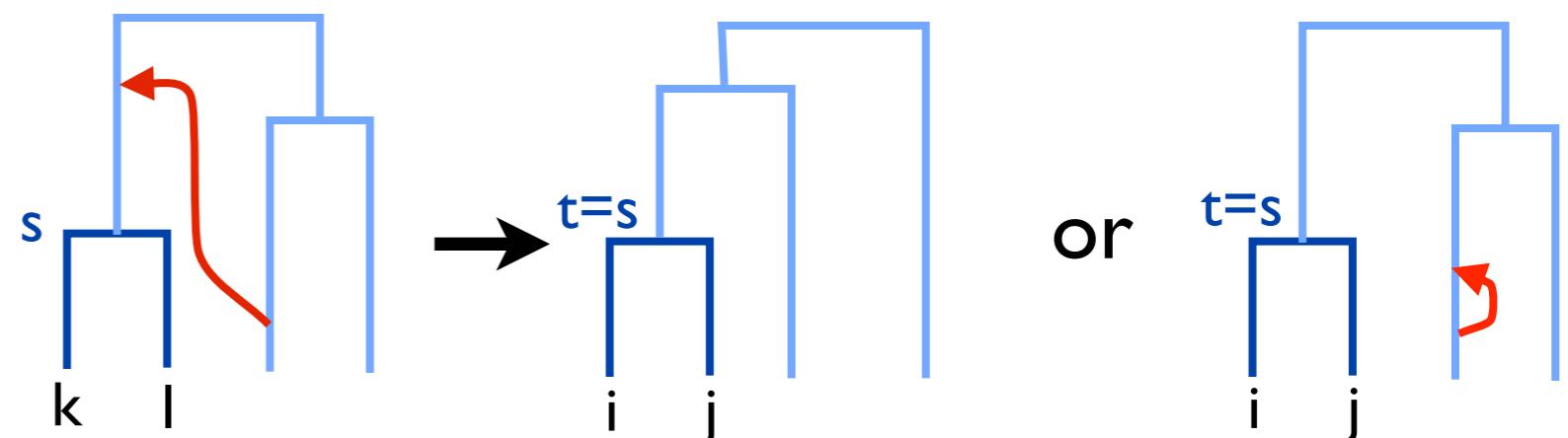
$(s, k, l) \rightarrow (t, i, j)$

where $t < s$



$(s, k, l) \rightarrow (t, i, j)$

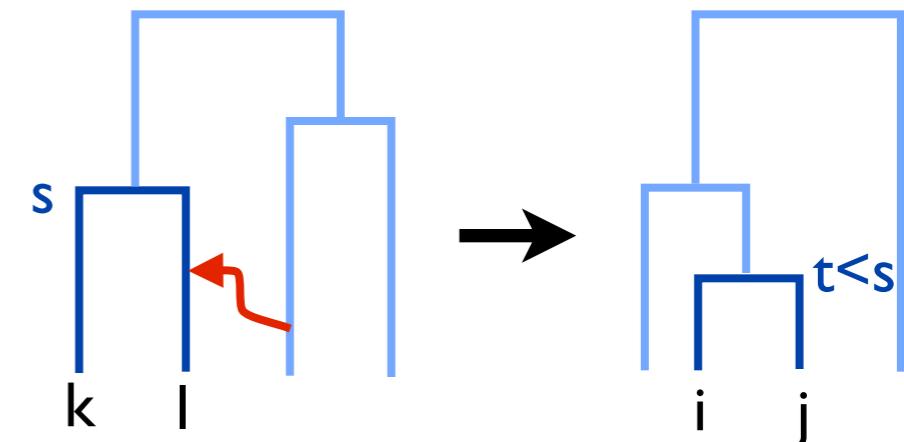
where $t = s$



MSMC: state transitions

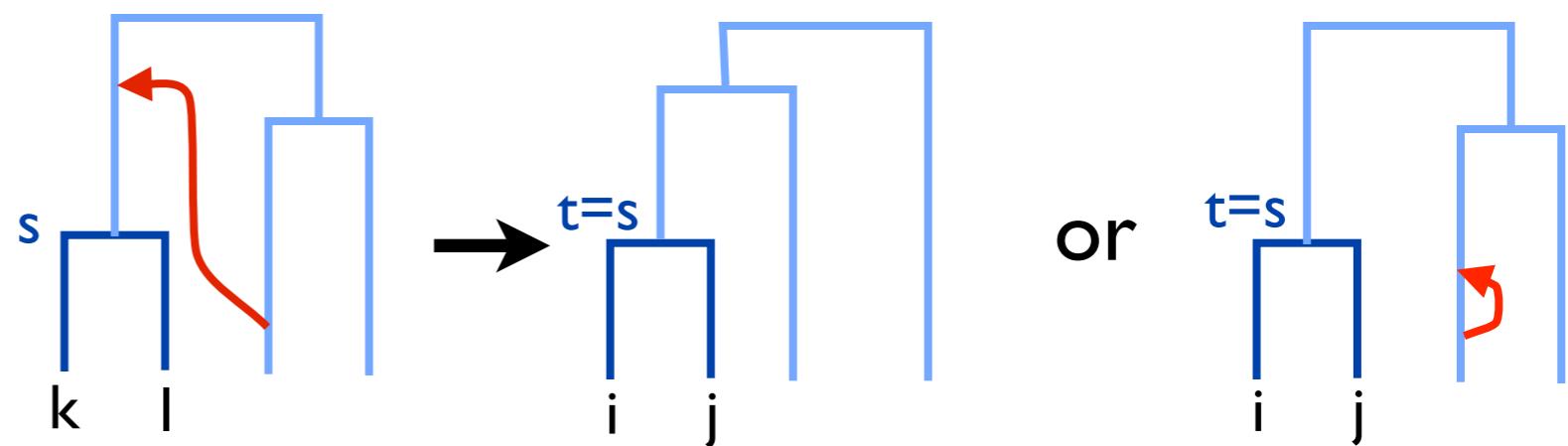
$(s, k, l) \rightarrow (t, i, j)$

where $t < s$



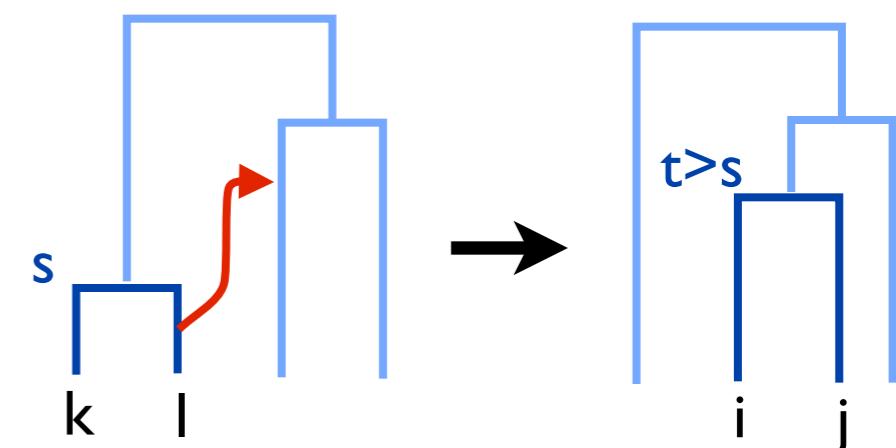
$(s, k, l) \rightarrow (t, i, j)$

where $t = s$

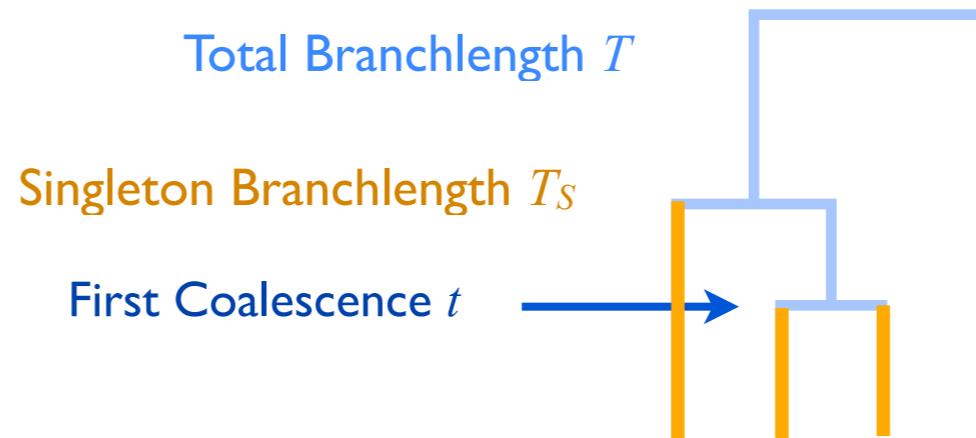


$(s, k, l) \rightarrow (t, i, j)$

where $t > s$



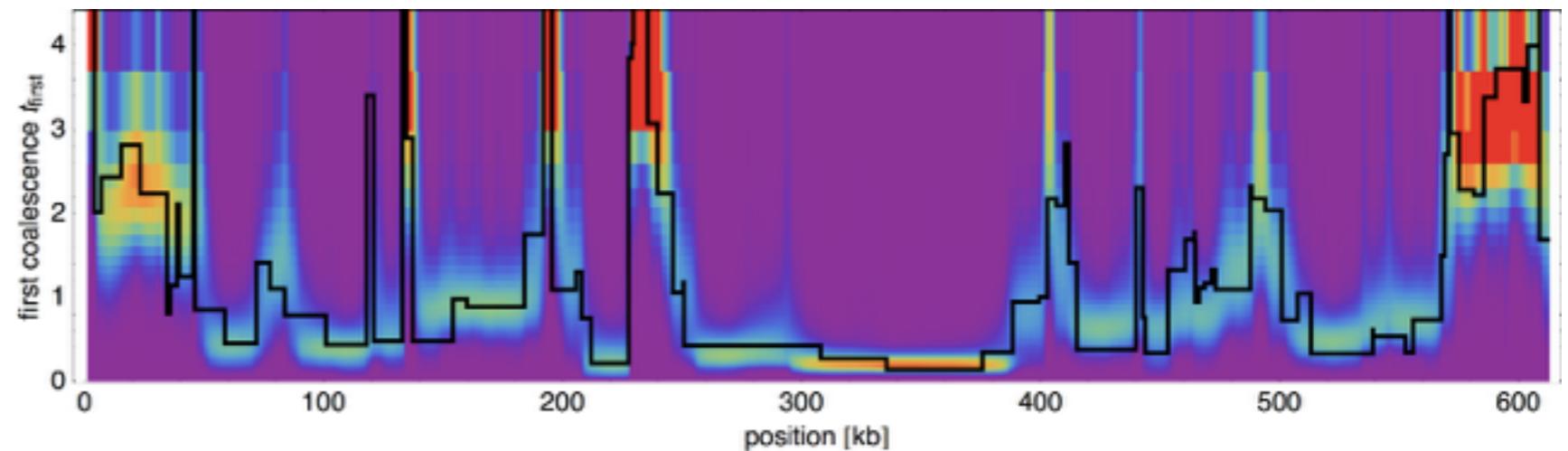
MSMC: mutation probability



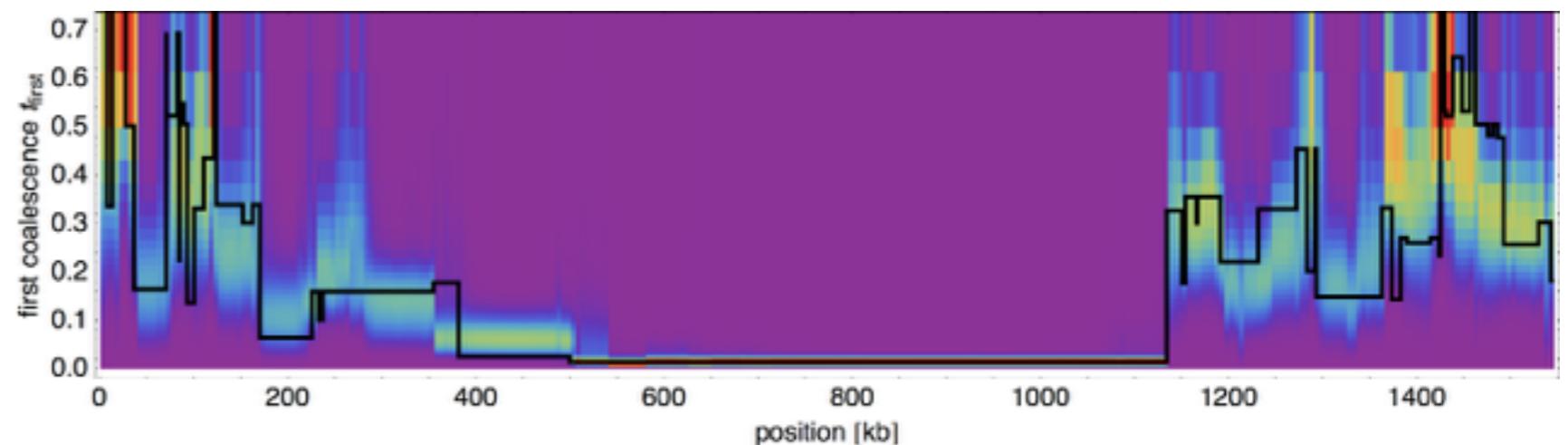
No mutation	$1 - \mu T$	A	A	A	A
Singleton within pair	μt	A	C	A	A
Singleton outside pair	μ	A	A	A	T
Double mutation	0	A	T	A	T
Higher Frequency	μ	A	T	T	A
Missing Data	1	A	A	-	-

Local Inference of first coalescence time

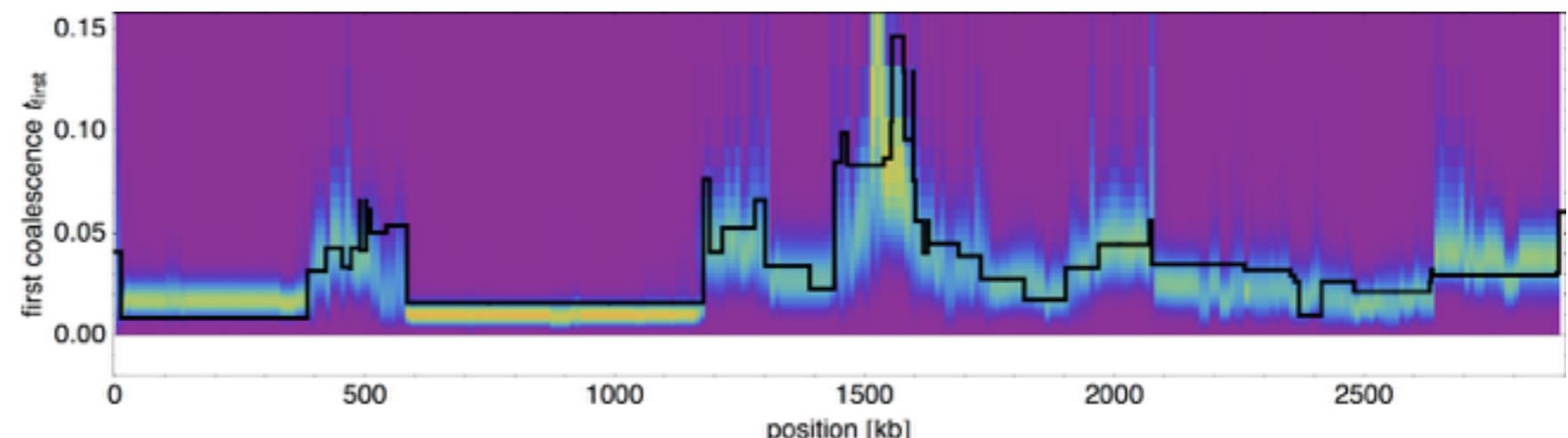
2 haplotypes,
similar to PSMC
[Li and Durbin, 2011]



4 haplotypes

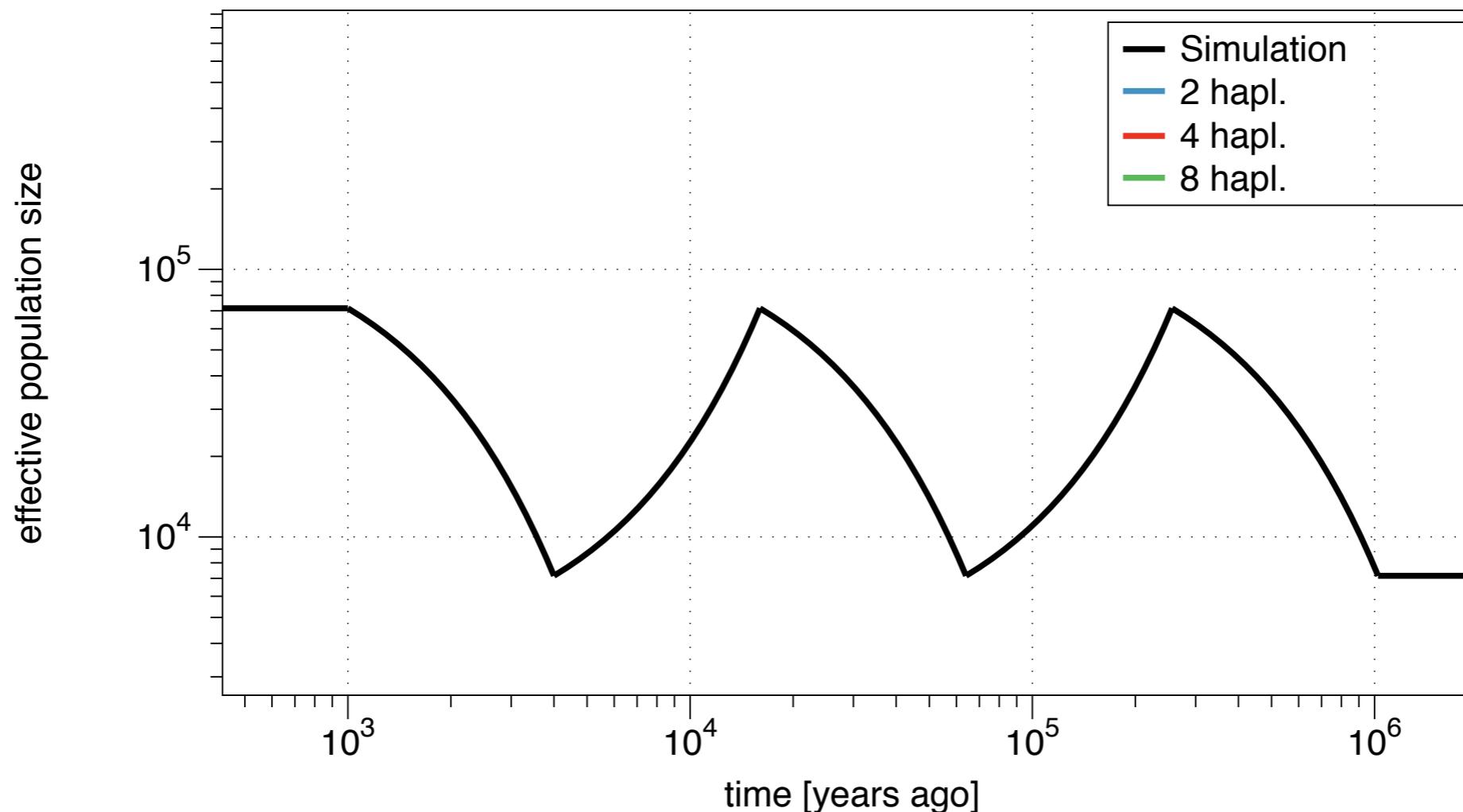


8 haplotypes



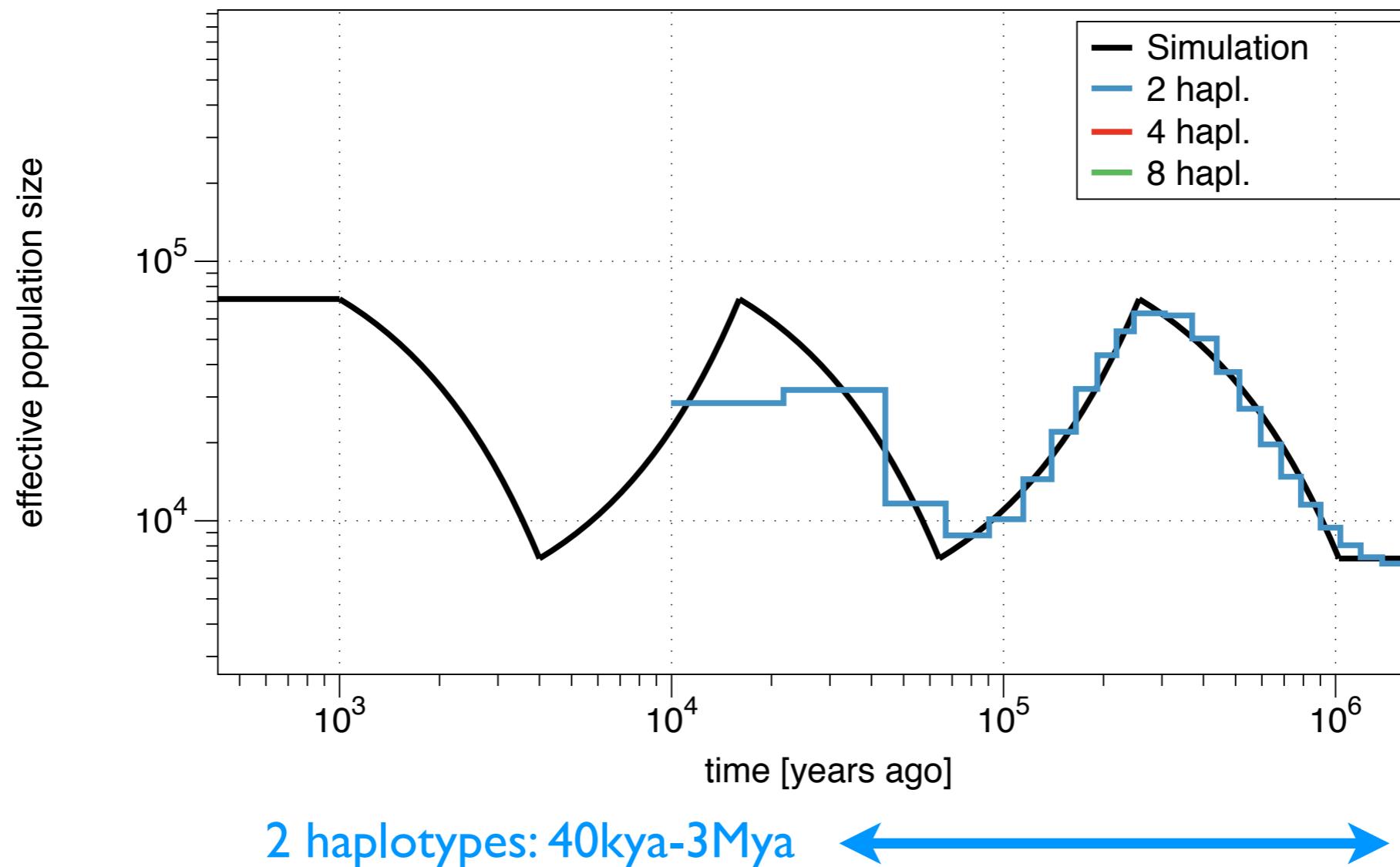
Test with simulations

Example: Exponentially growing and shrinking population size



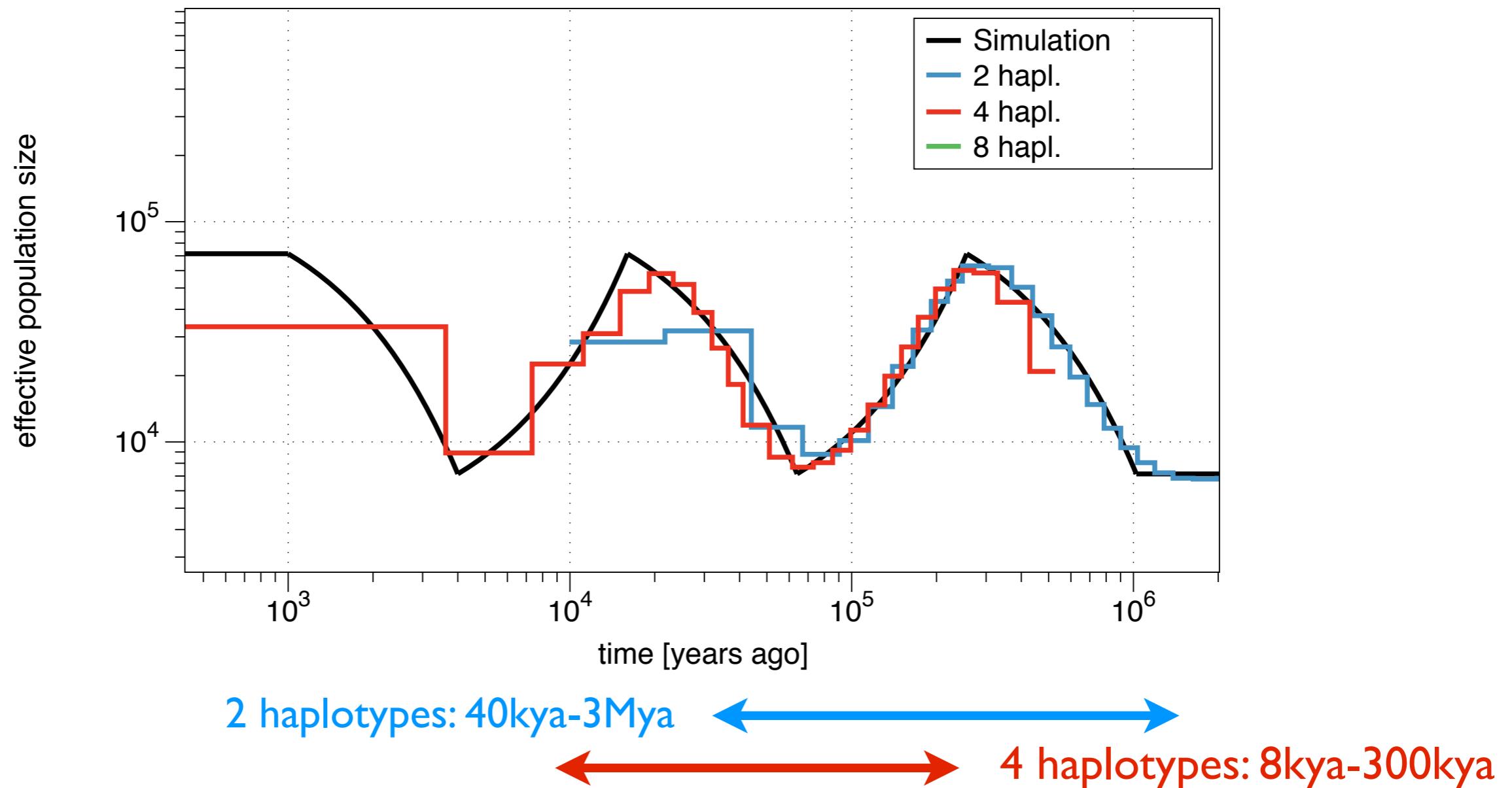
Test with simulations

Example: Exponentially growing and shrinking population size



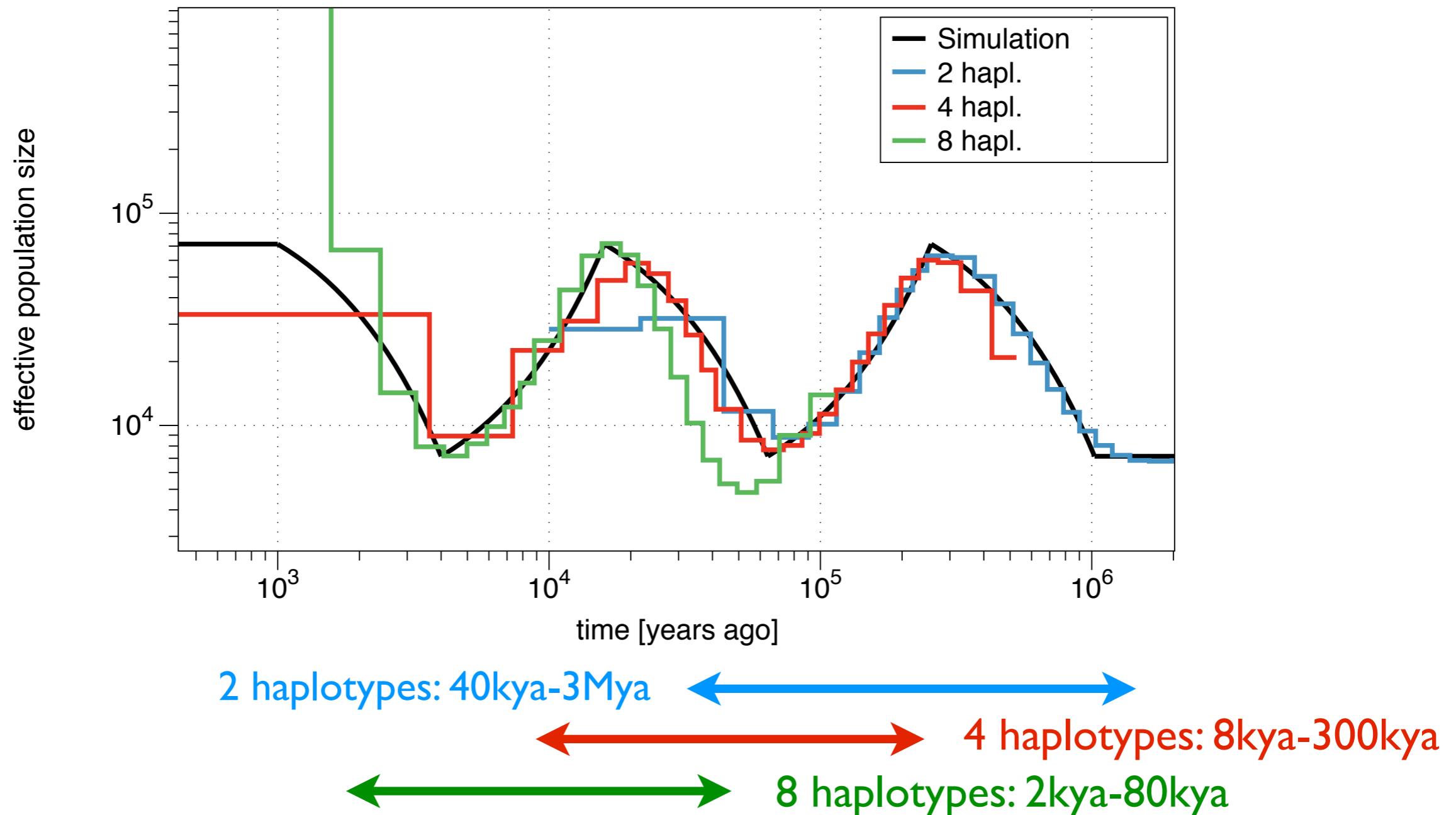
Test with simulations

Example: Exponentially growing and shrinking population size

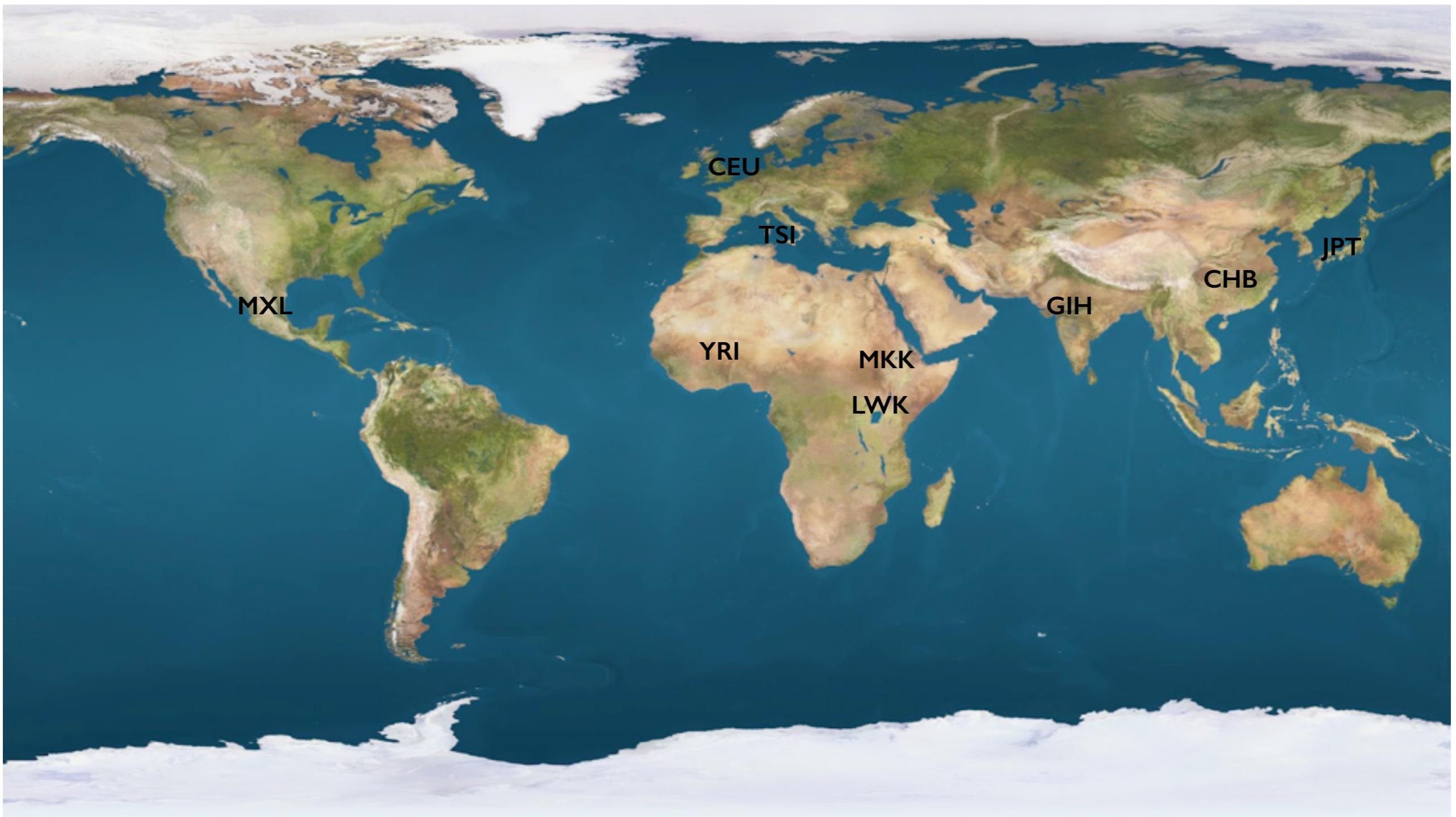


Test with simulations

Example: Exponentially growing and shrinking population size



From genome sequences to human history

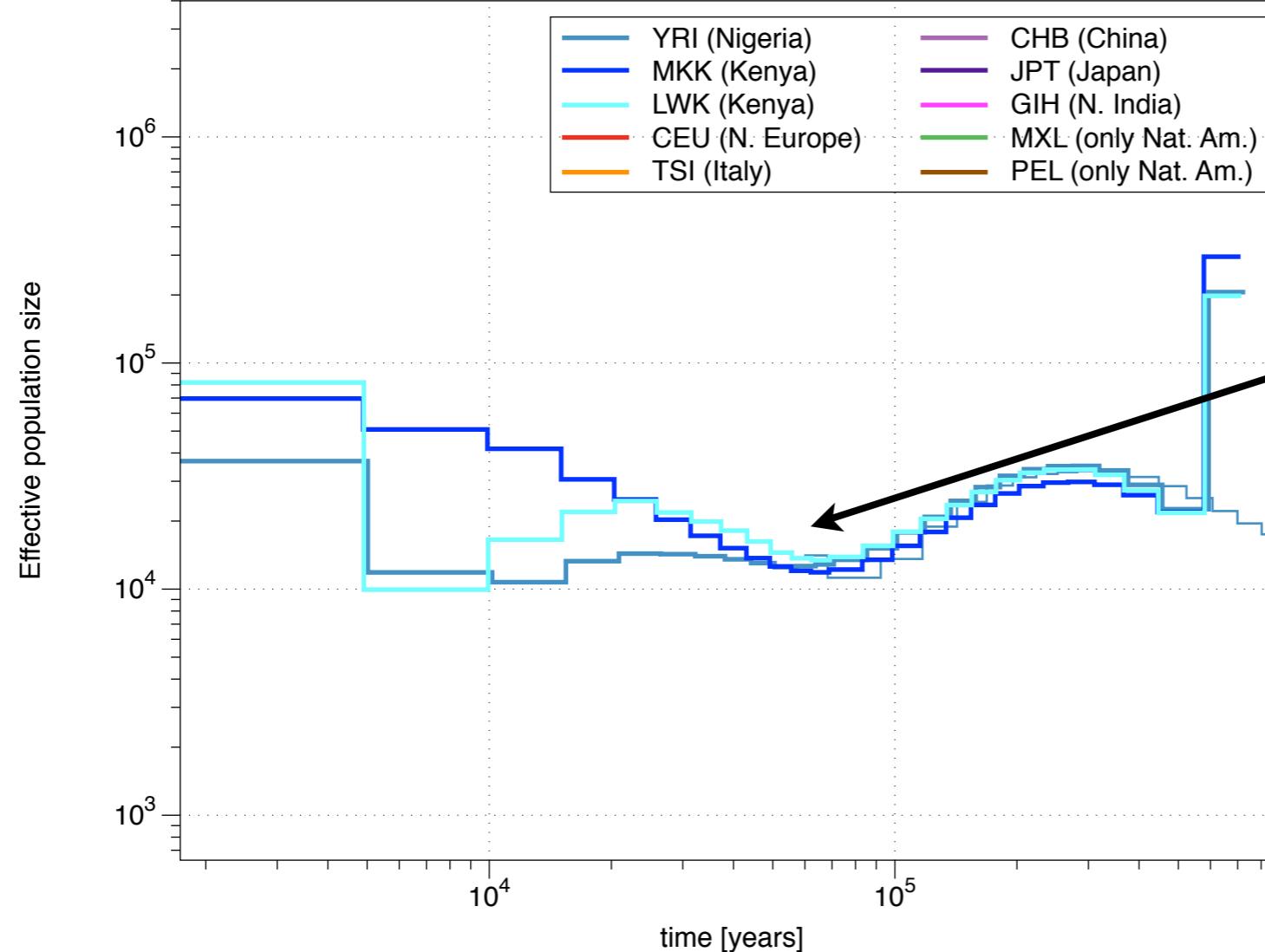


[Sequence data from Complete Genomics]

Inferring historical Population Sizes

real time scaling using
mutation rate per generation
 $\mu=1.25\times10^{-8}$ and a generation
time of 30 years

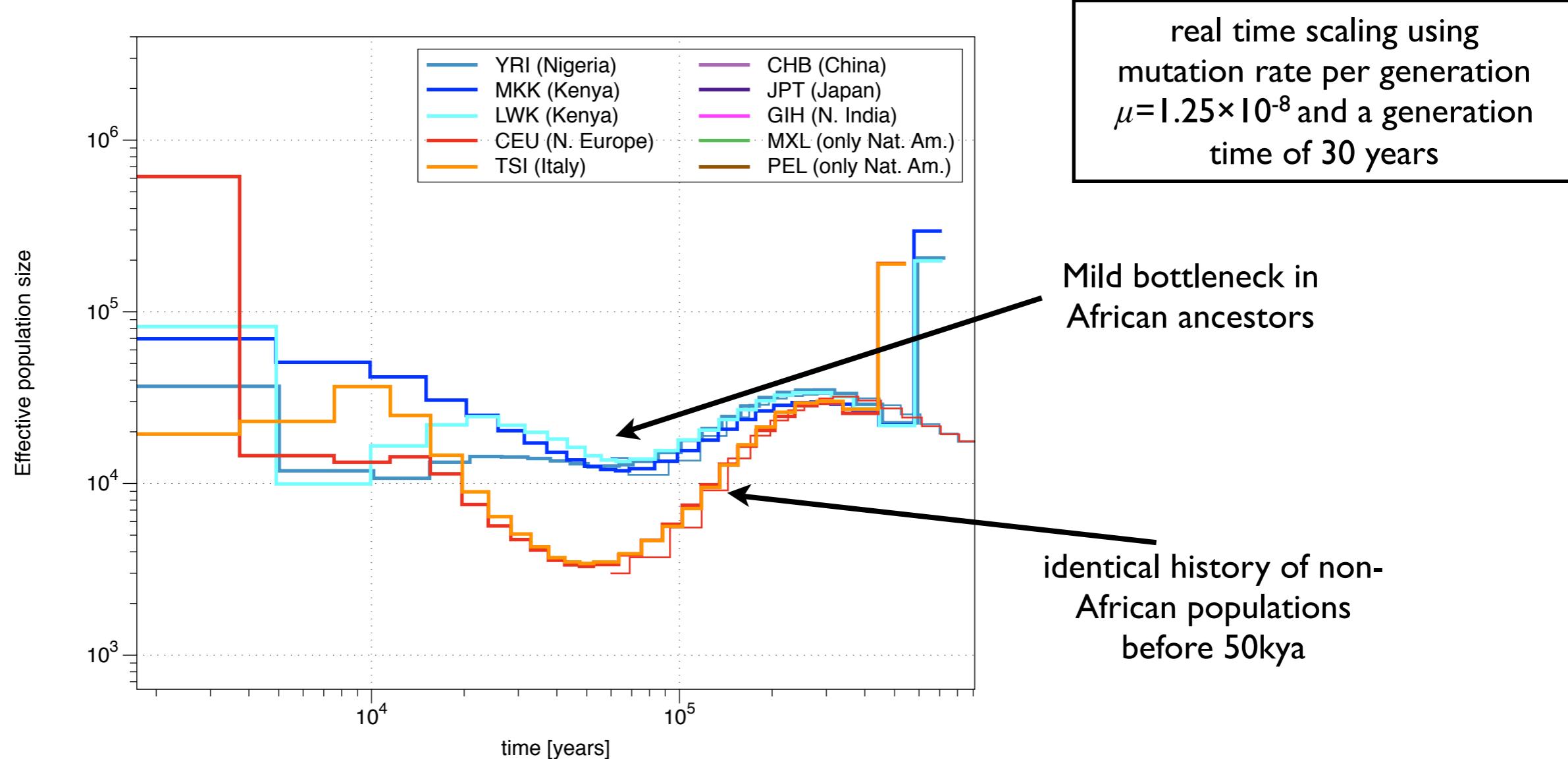
Inferring historical Population Sizes



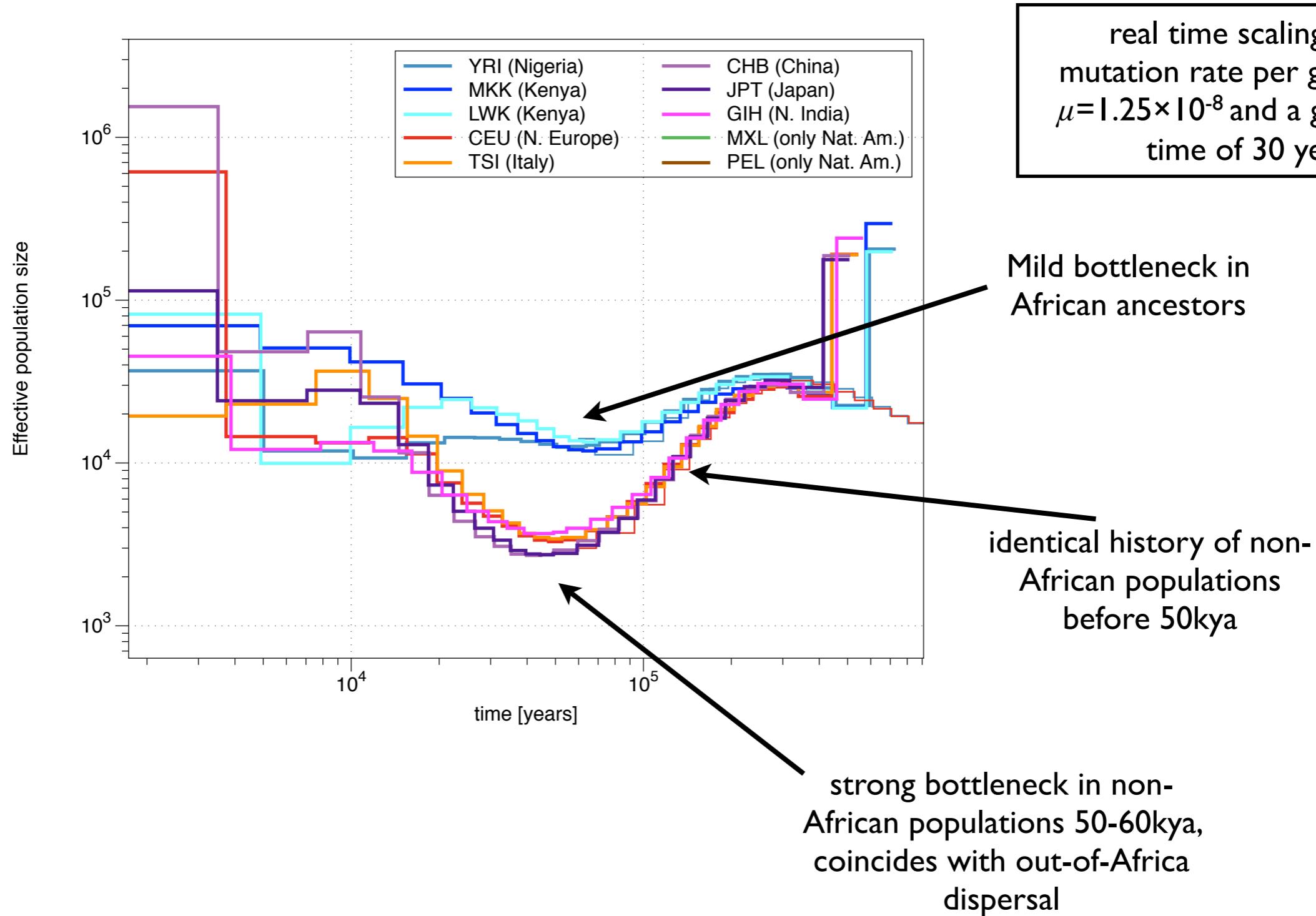
real time scaling using
mutation rate per generation
 $\mu=1.25 \times 10^{-8}$ and a generation
time of 30 years

Mild bottleneck in
African ancestors

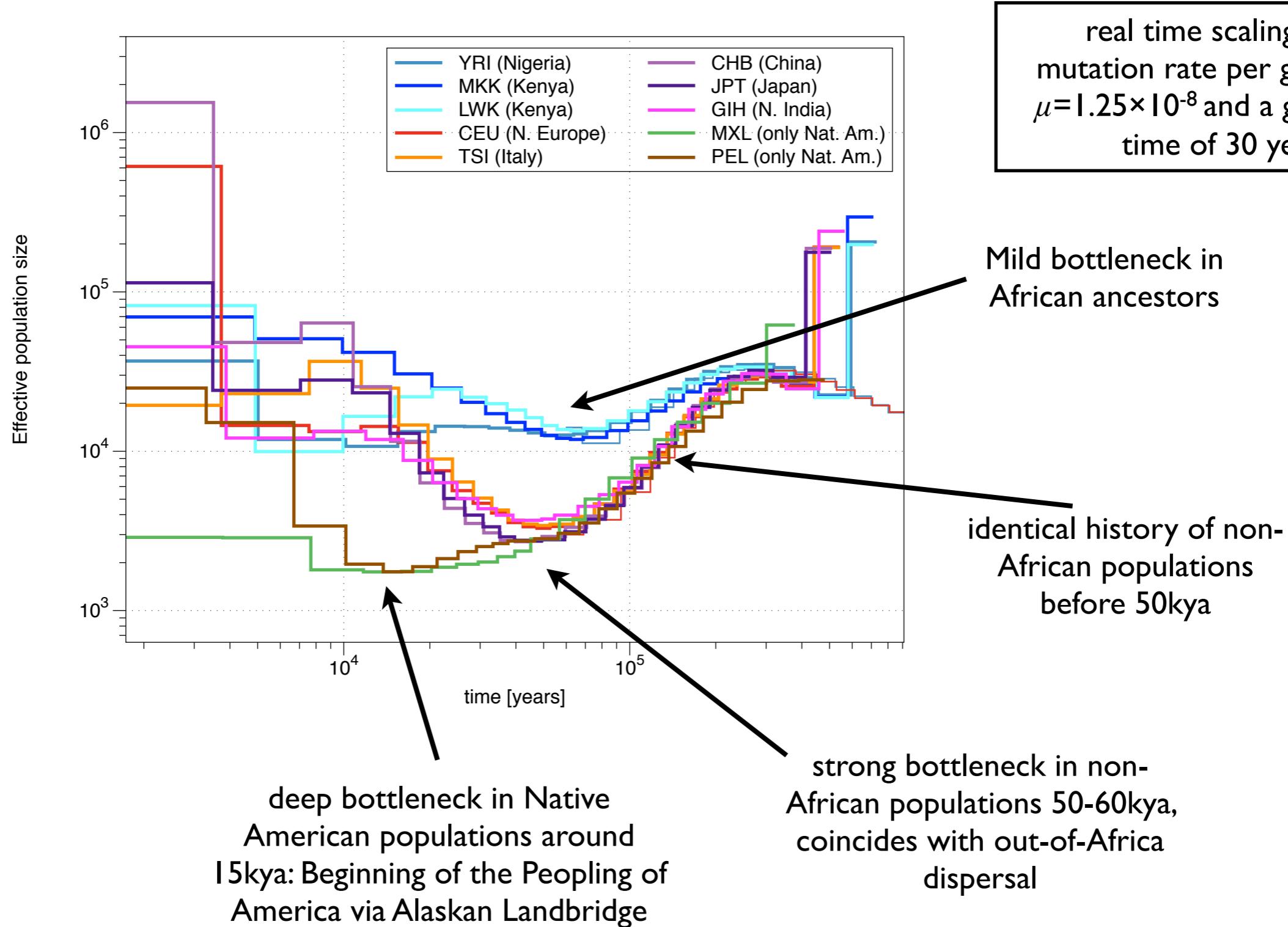
Inferring historical Population Sizes



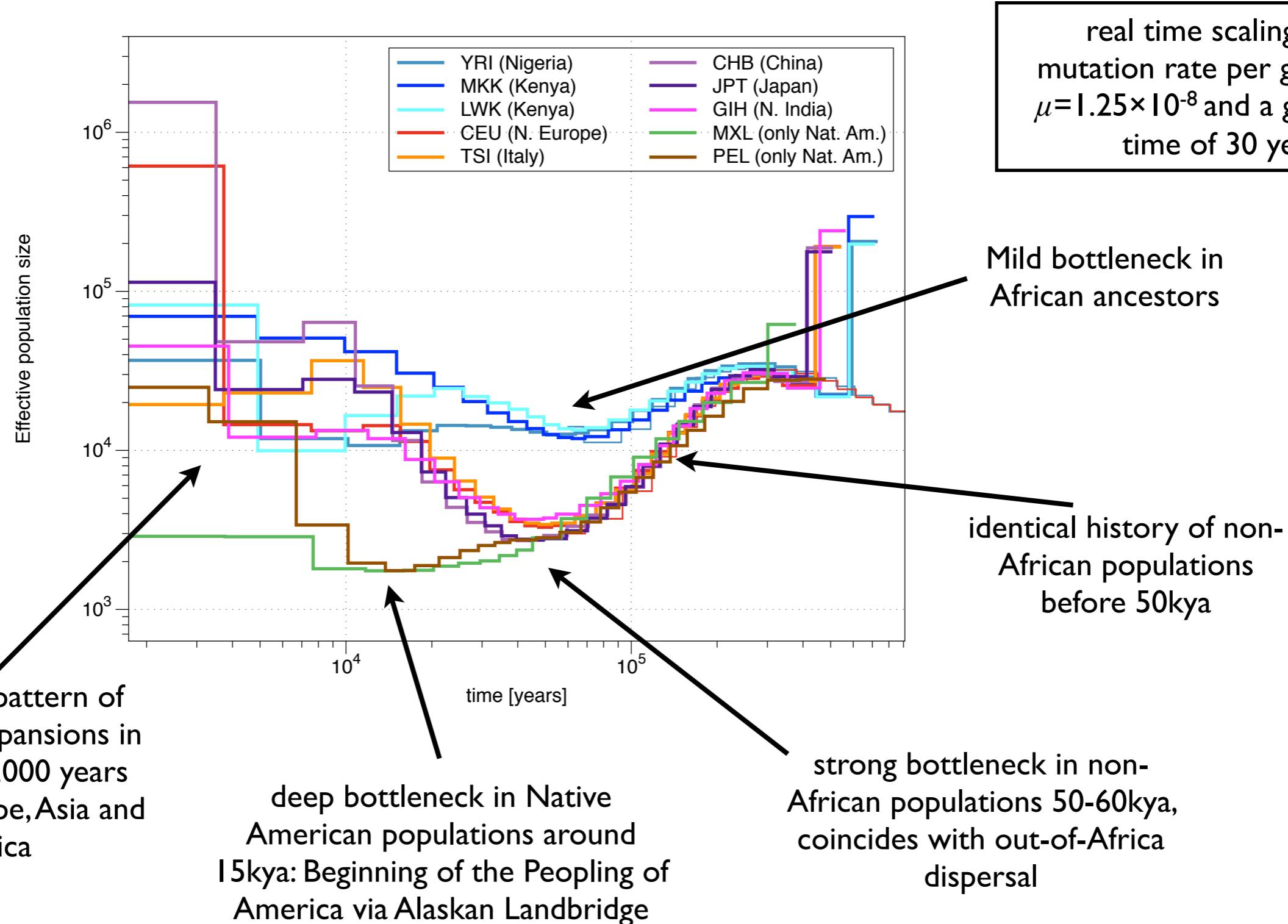
Inferring historical Population Sizes



Inferring historical Population Sizes

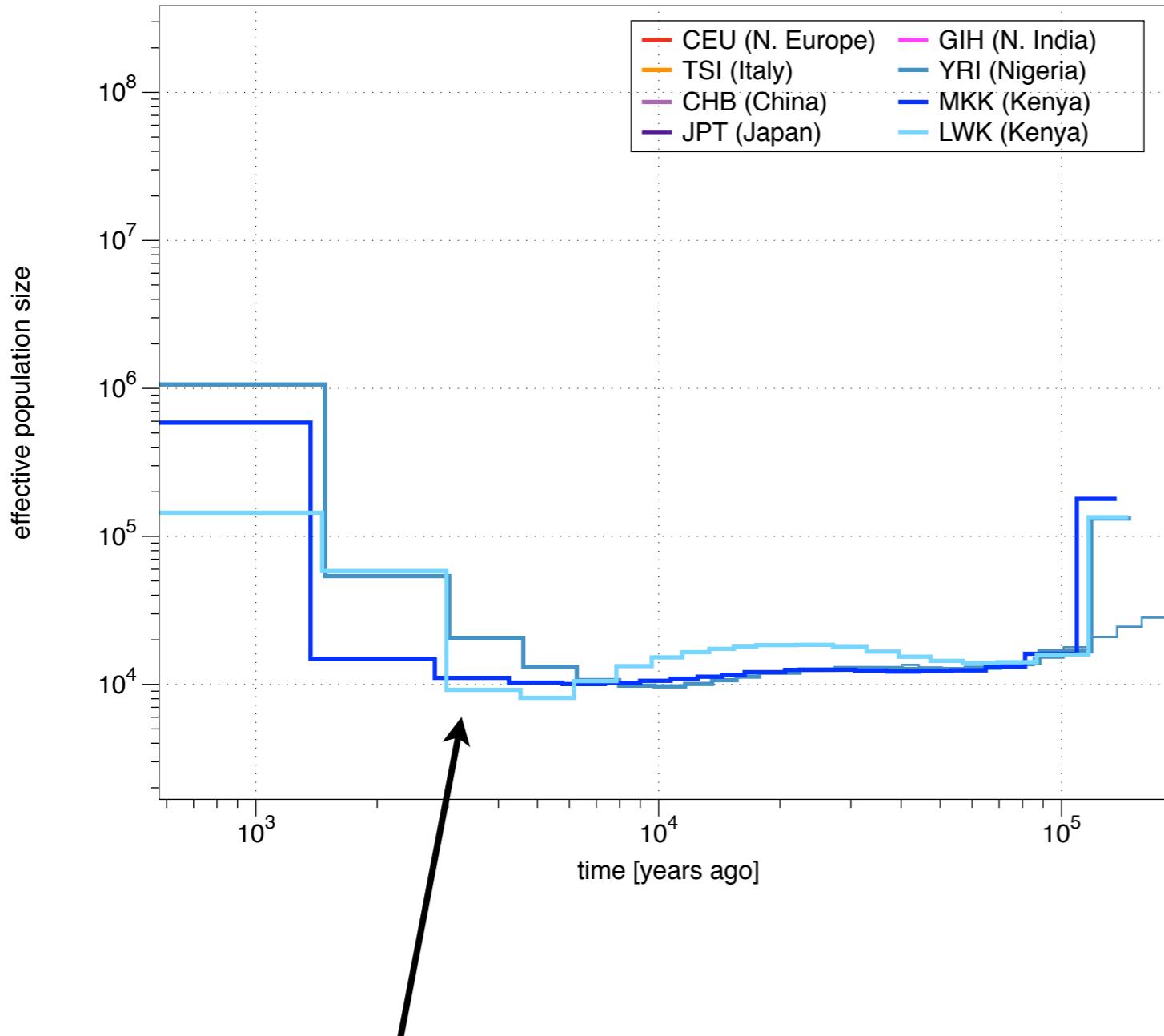


Inferring historical Population Sizes



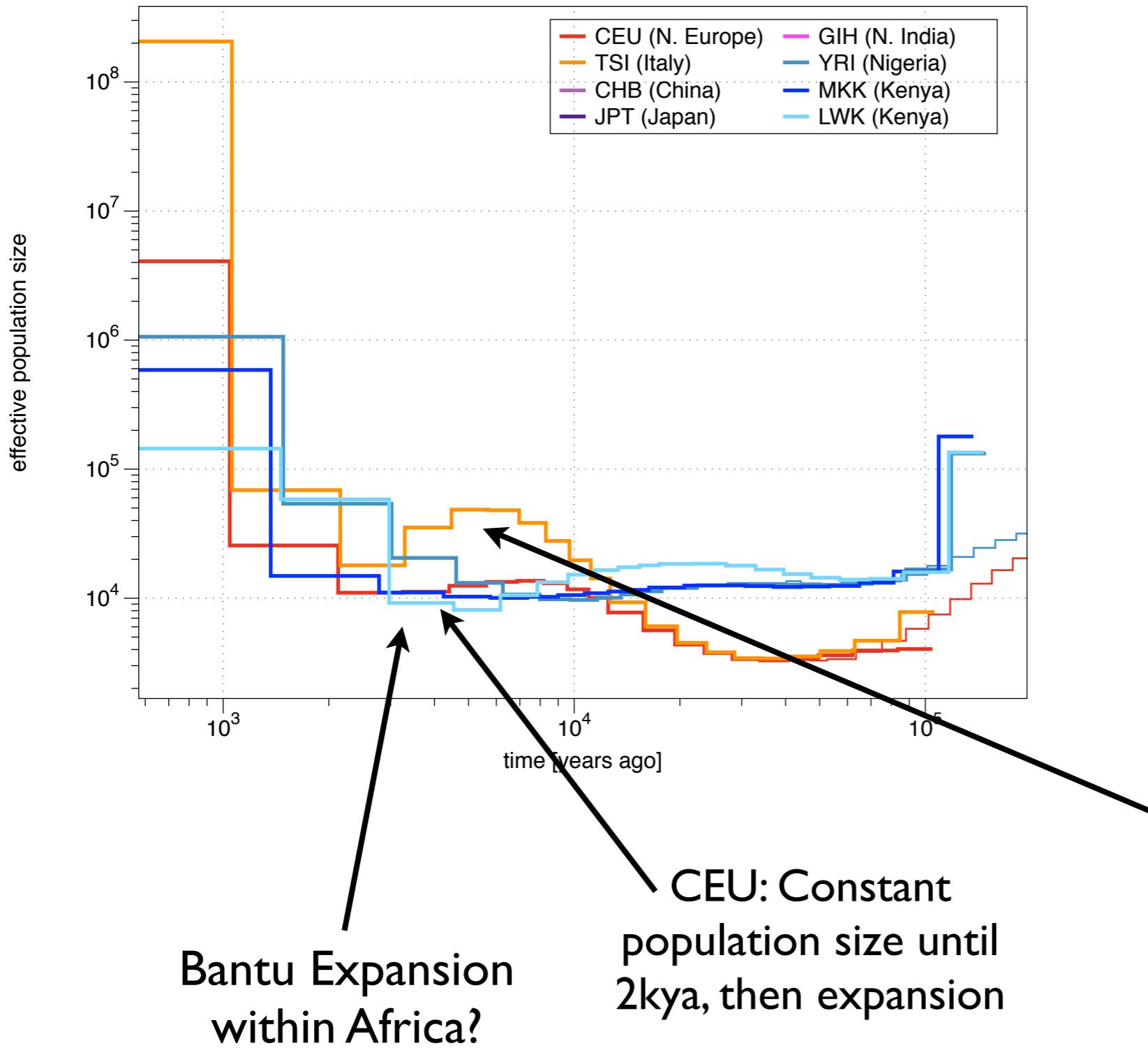
Population size inference from 8 haplotypes

Population size inference from 8 haplotypes



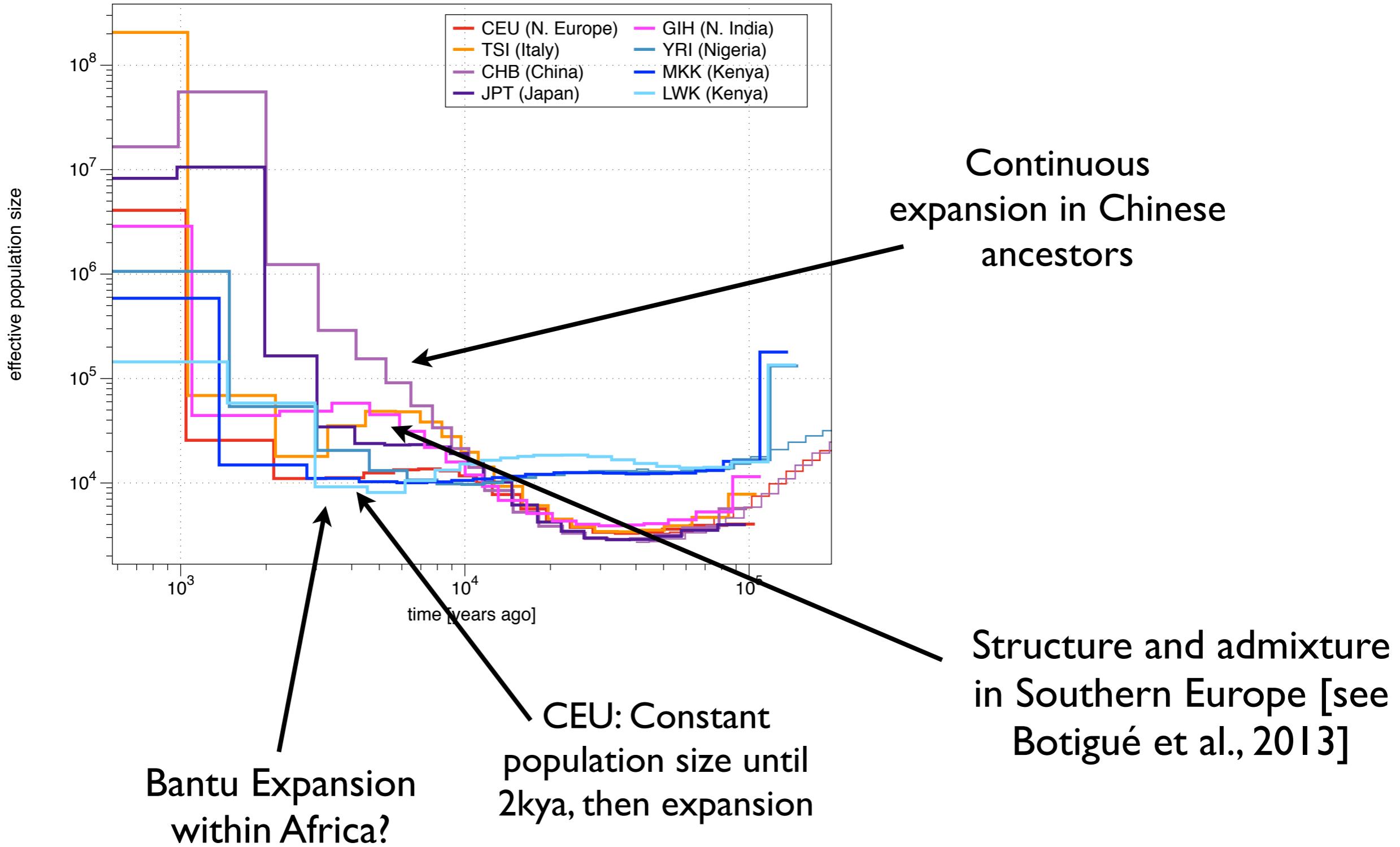
Bantu Expansion
within Africa?

Population size inference from 8 haplotypes



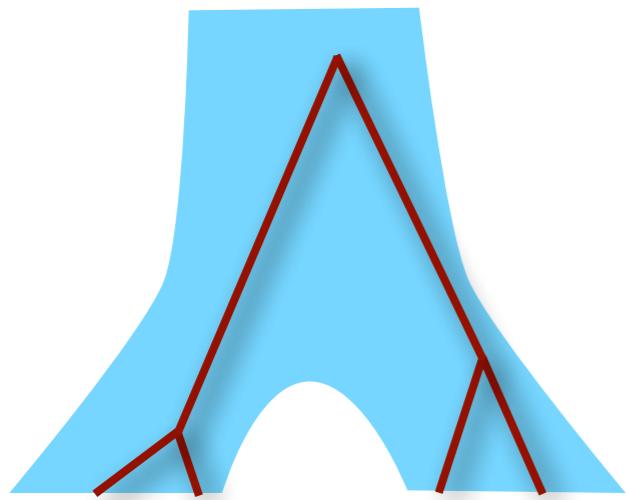
Structure and admixture in Southern Europe [see Botigué et al., 2013]

Population size inference from 8 haplotypes

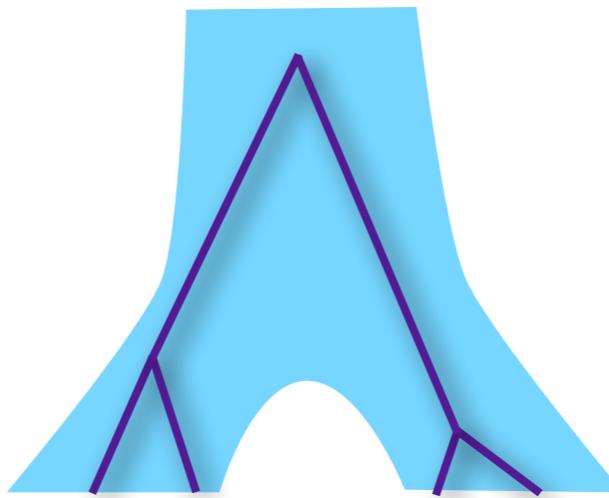


Divergence between populations

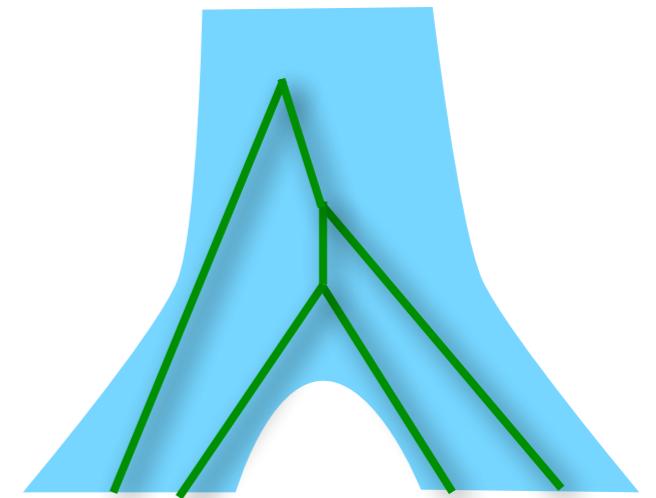
- Idea: Infer separate coalescence rates within and between populations:



First Coalescence
within Population 1



First Coalescence
within Population 2

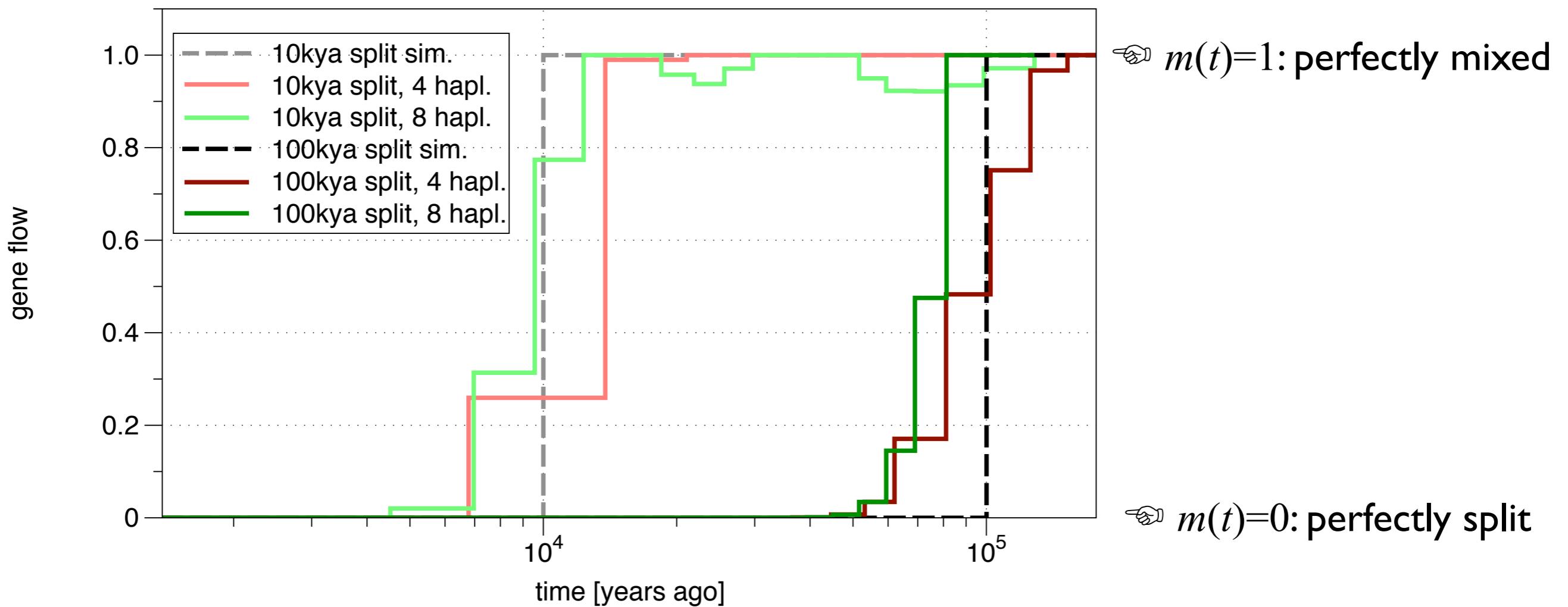


First Coalescence
across both
populations

- MSMC can infer separate coalescence rates within populations,
- Given rates within populations, $\lambda_{11}(t)$ and $\lambda_{22}(t)$, and across populations, $\lambda_{12}(t)$, compute relative gene flow as ratio

$$m(t) = \frac{\lambda_{12}(t)}{[\lambda_{11}(t) + \lambda_{22}(t)] / 2}$$

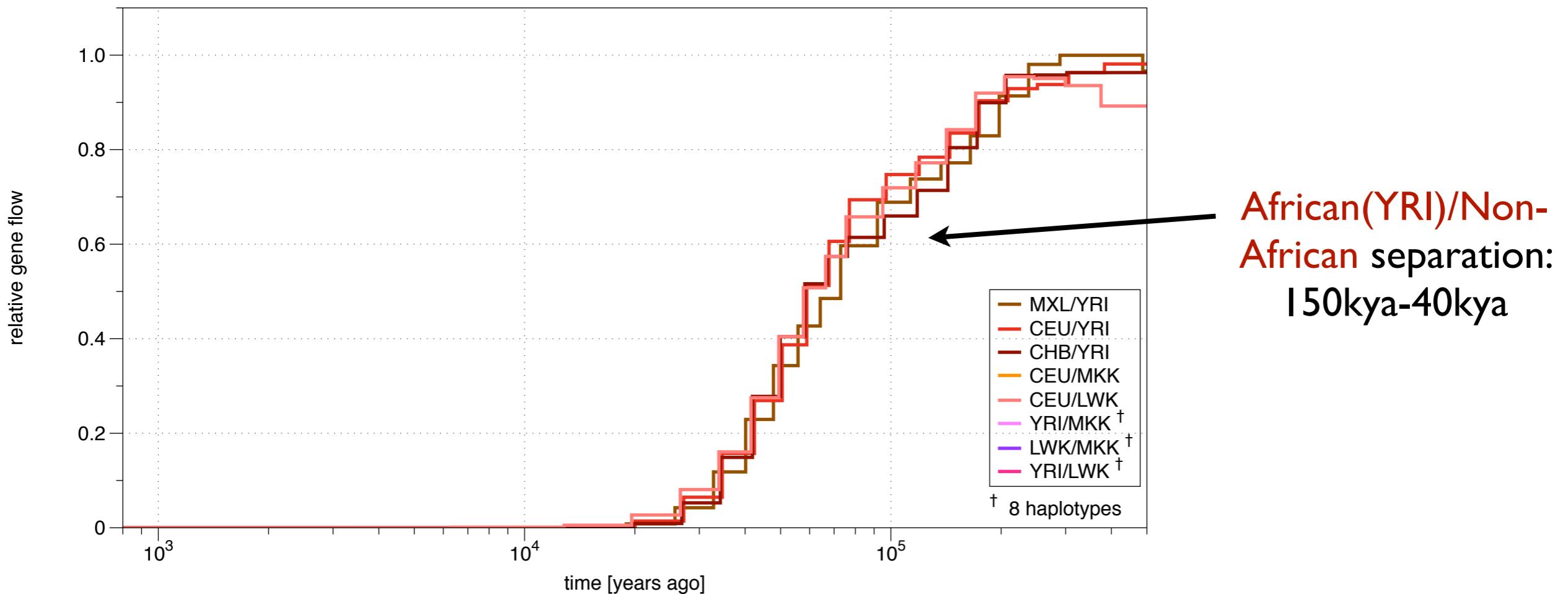
Testing gene flow inference with simulated split



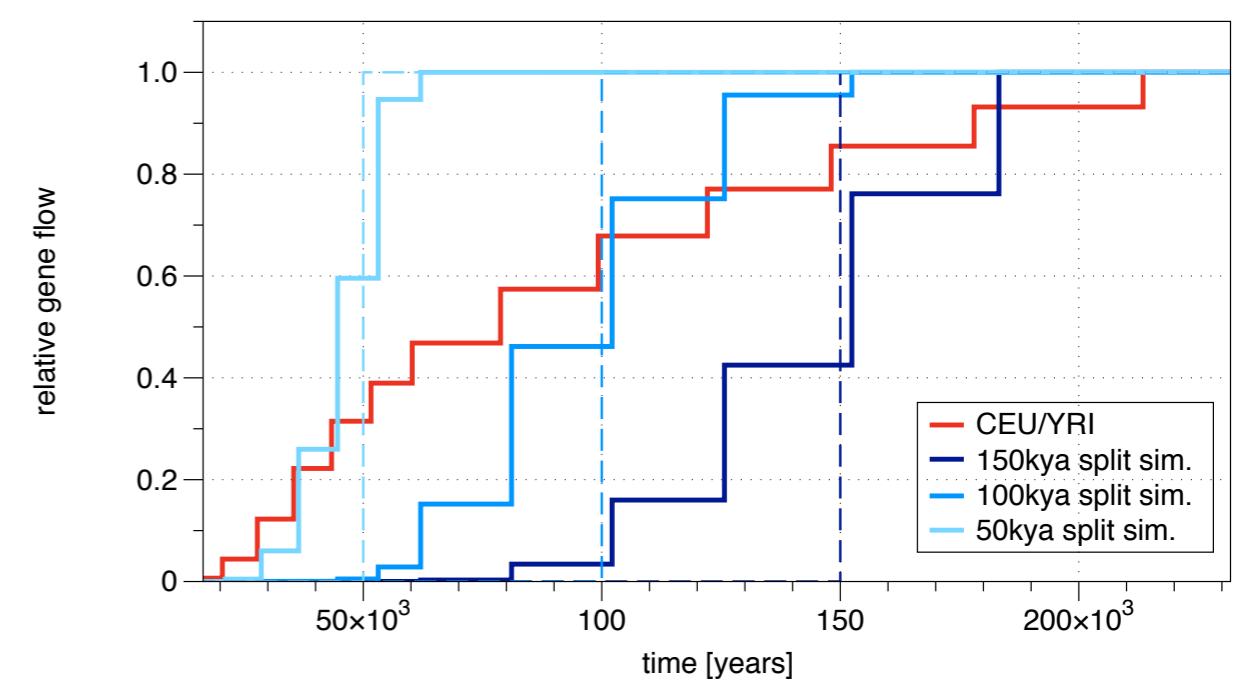
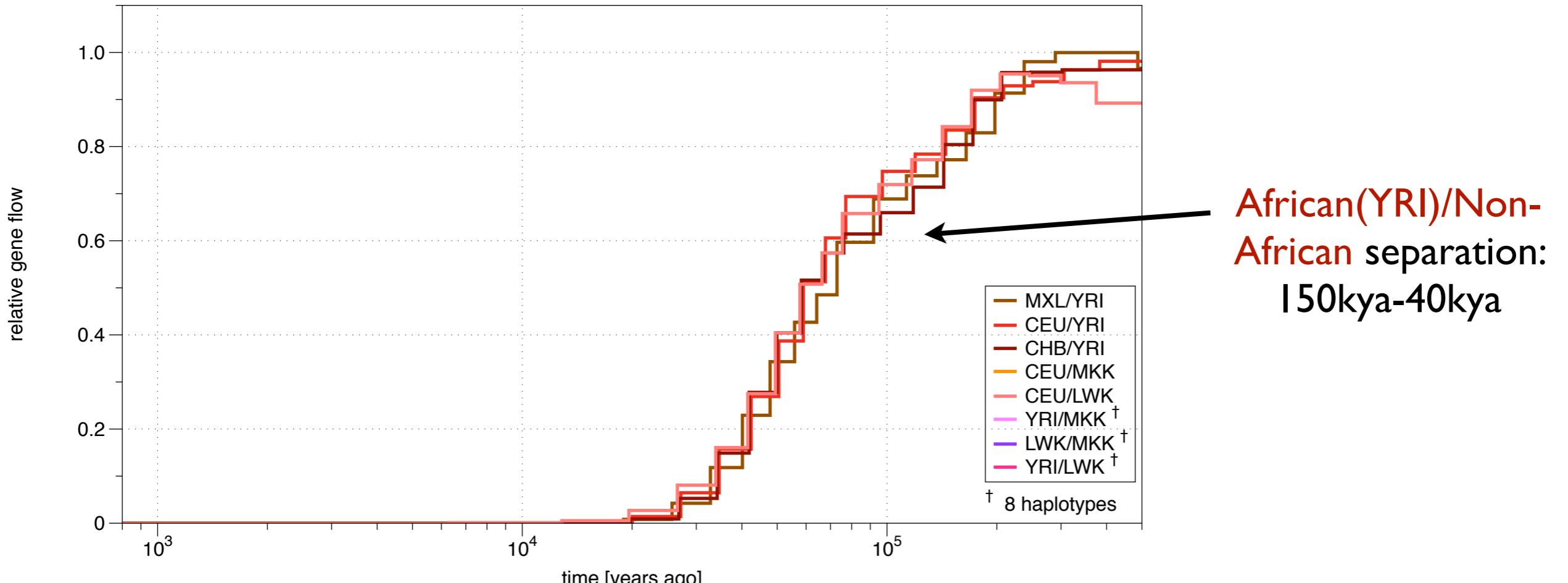
4 haplotypes: good for splits 50-200kya.
8 haplotypes: good for splits 5-50kya.

African Population separations

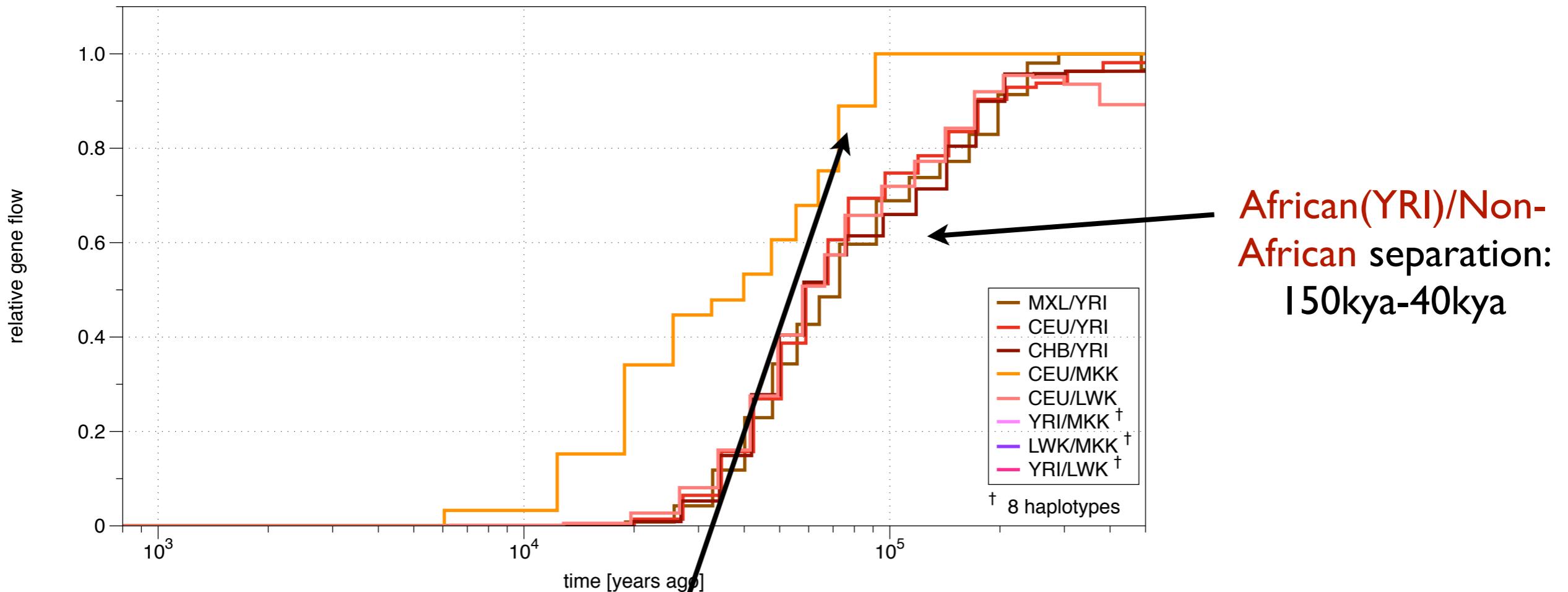
African Population separations



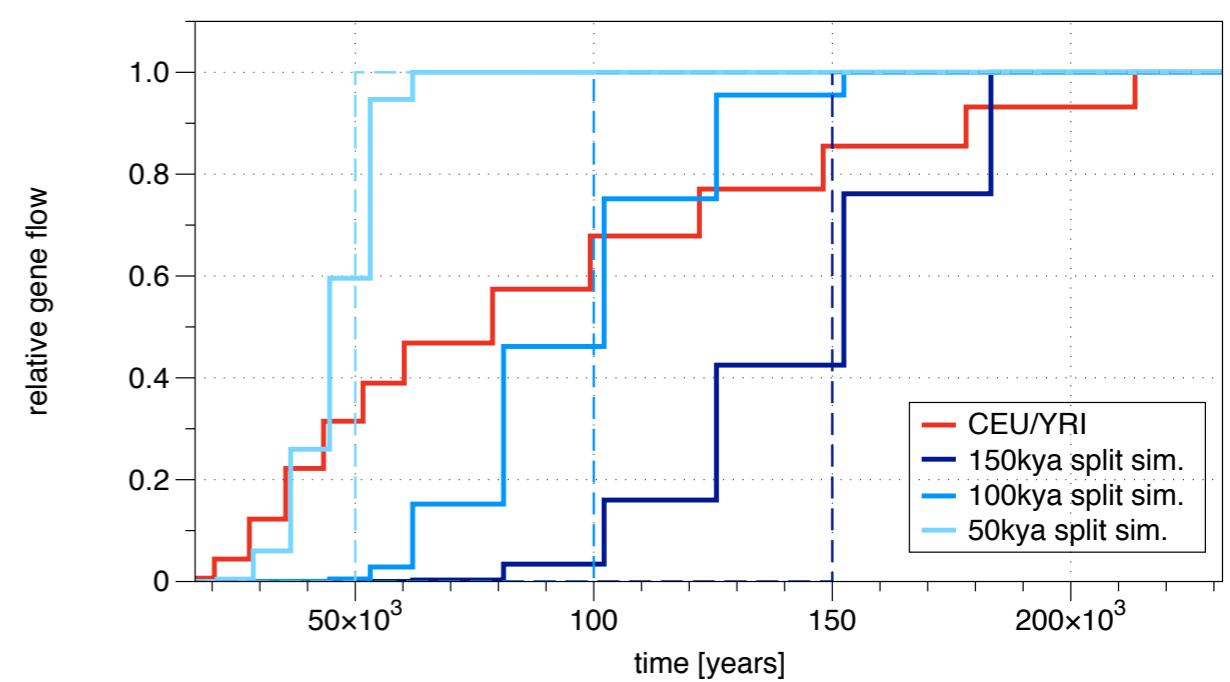
African Population separations



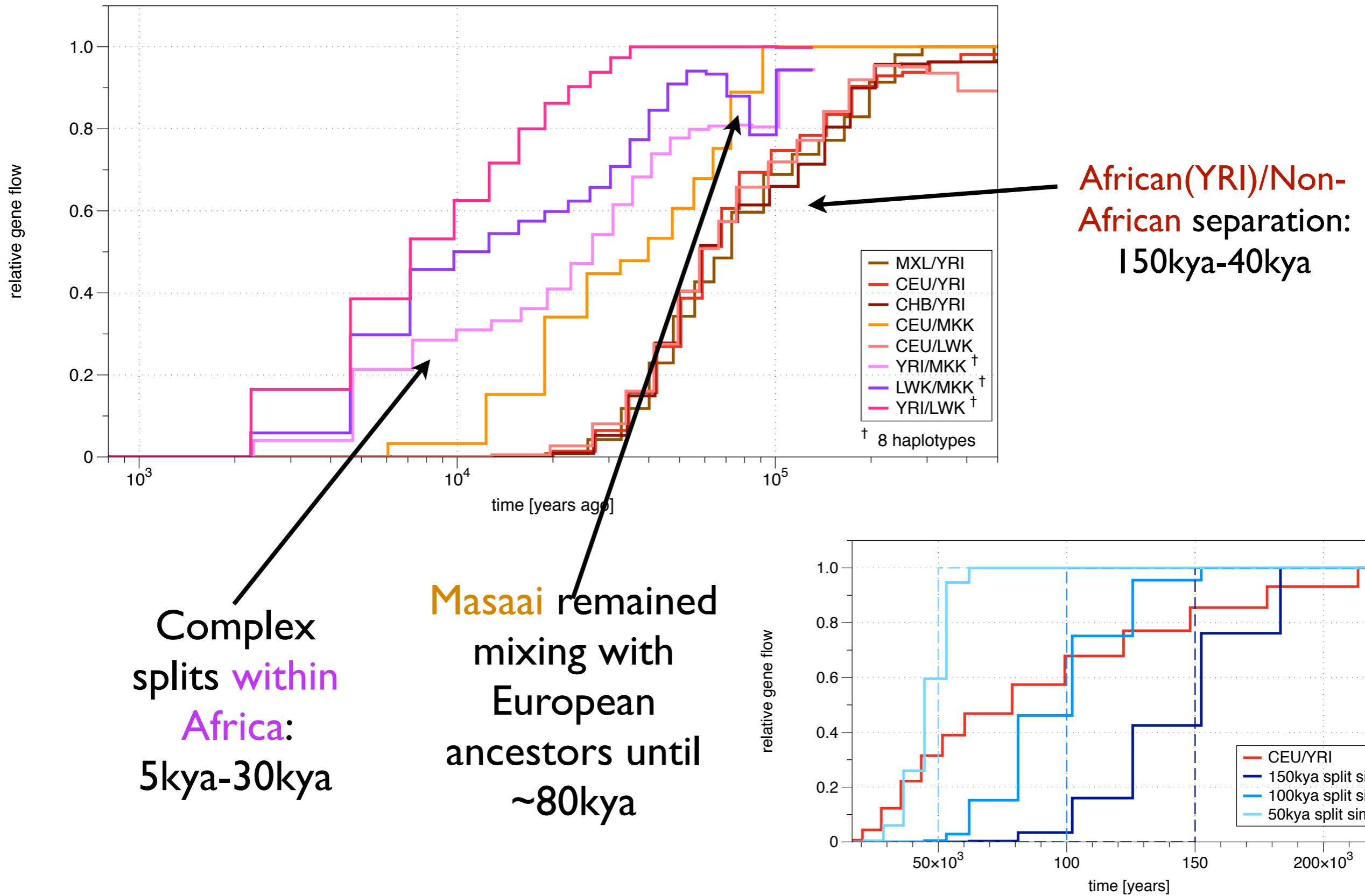
African Population separations



Masai remained mixing with European ancestors until ~80kya

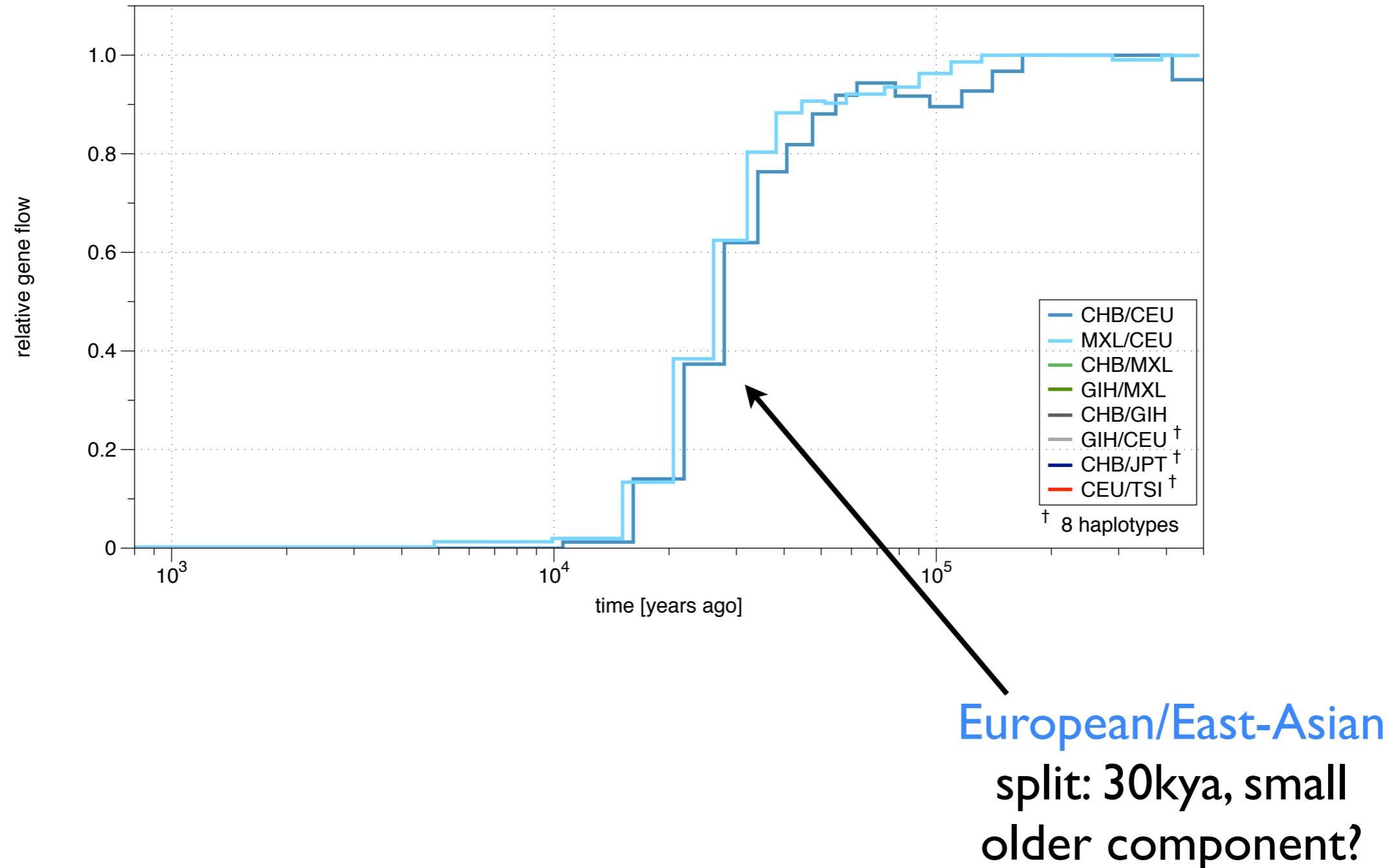


African Population separations

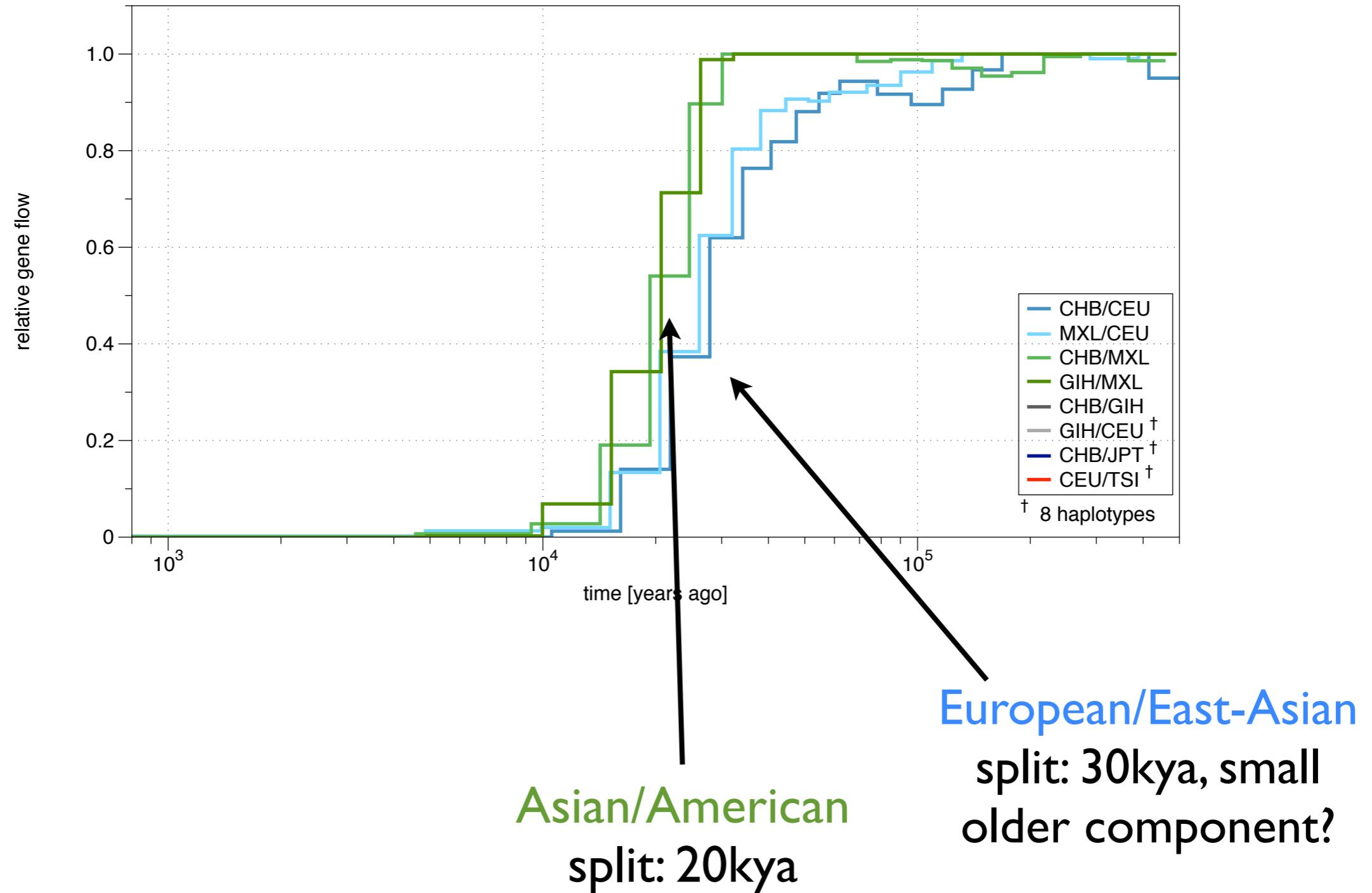


Non-African Population Separations

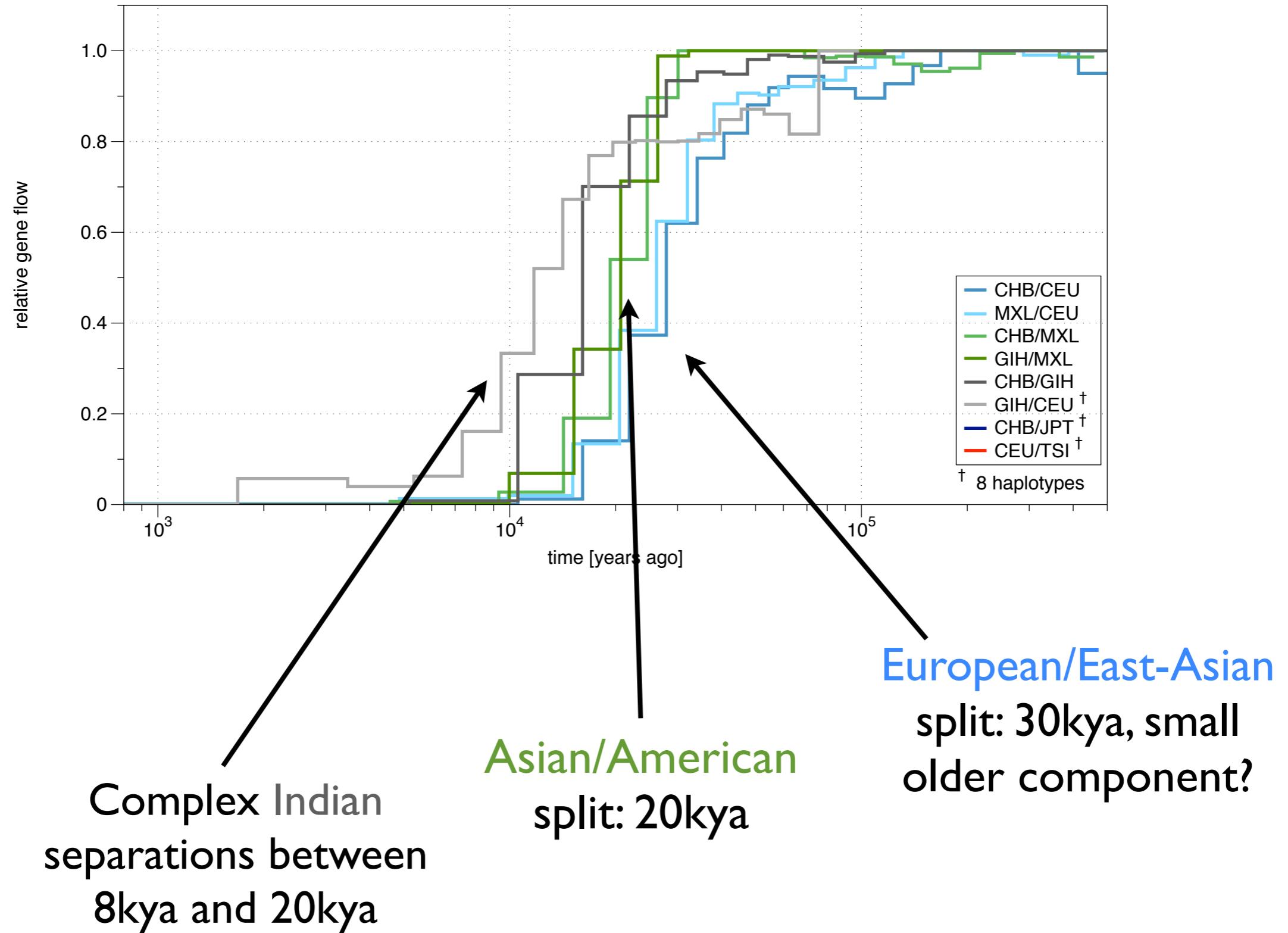
Non-African Population Separations



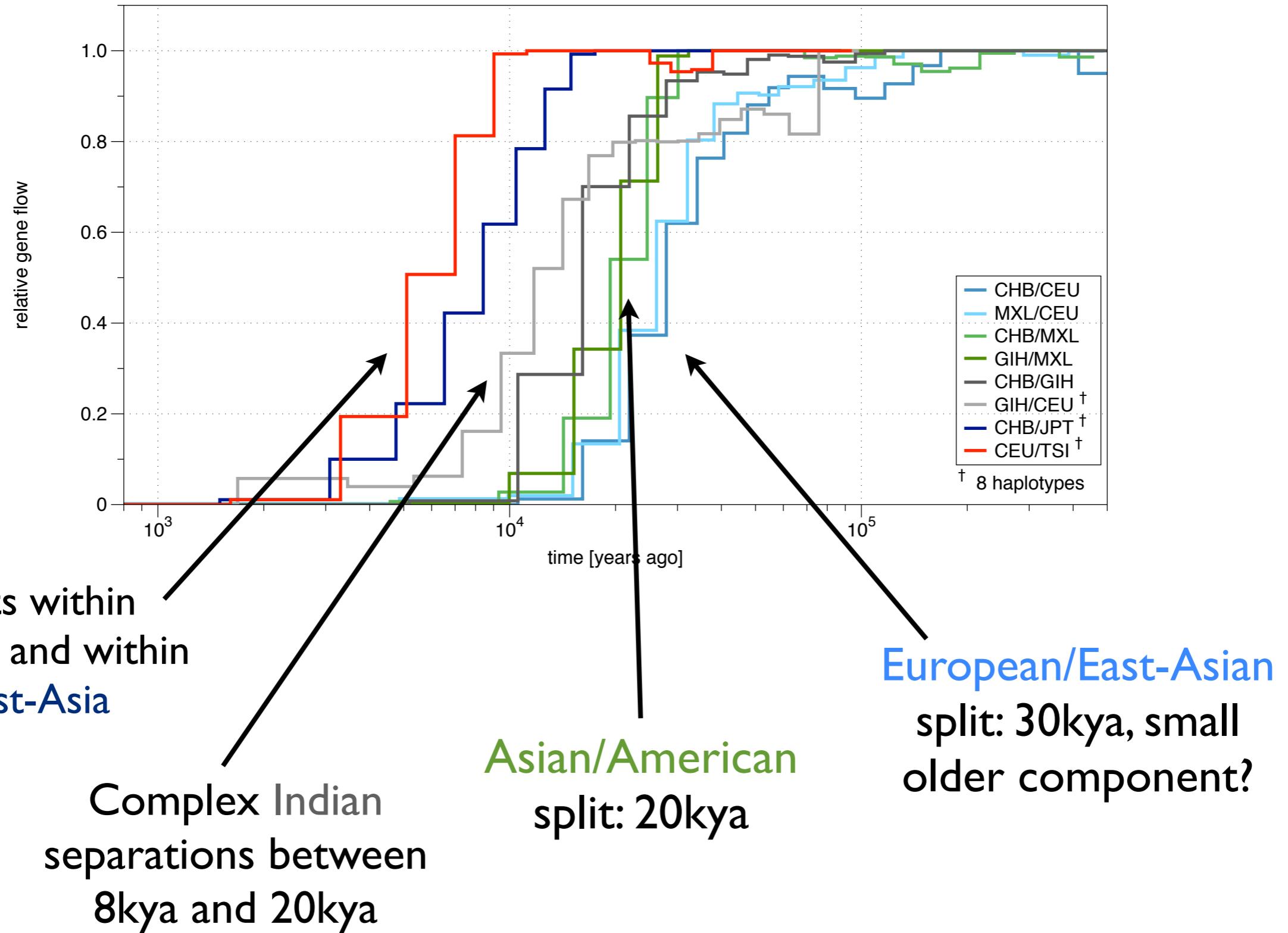
Non-African Population Separations



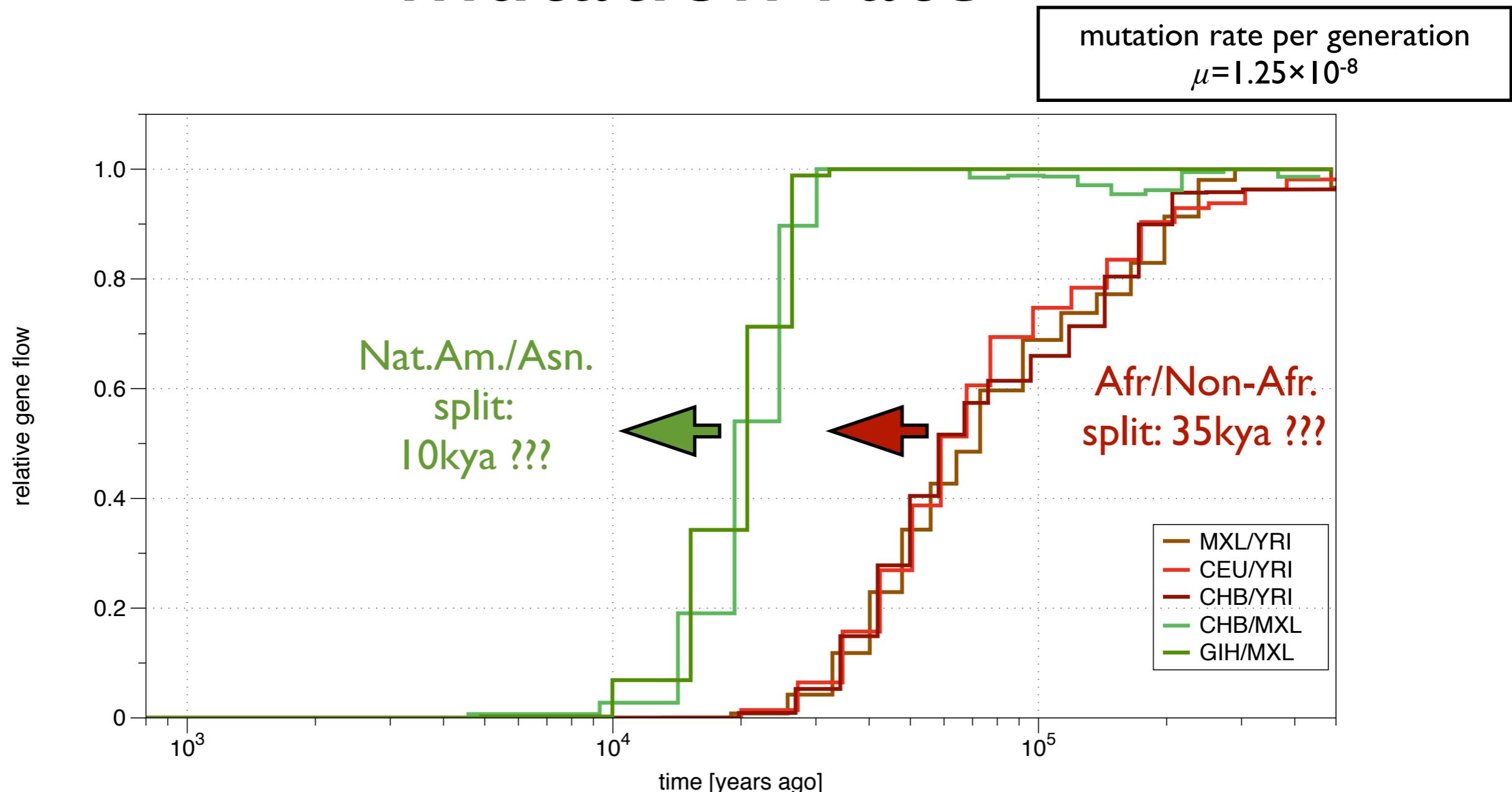
Non-African Population Separations



Non-African Population Separations

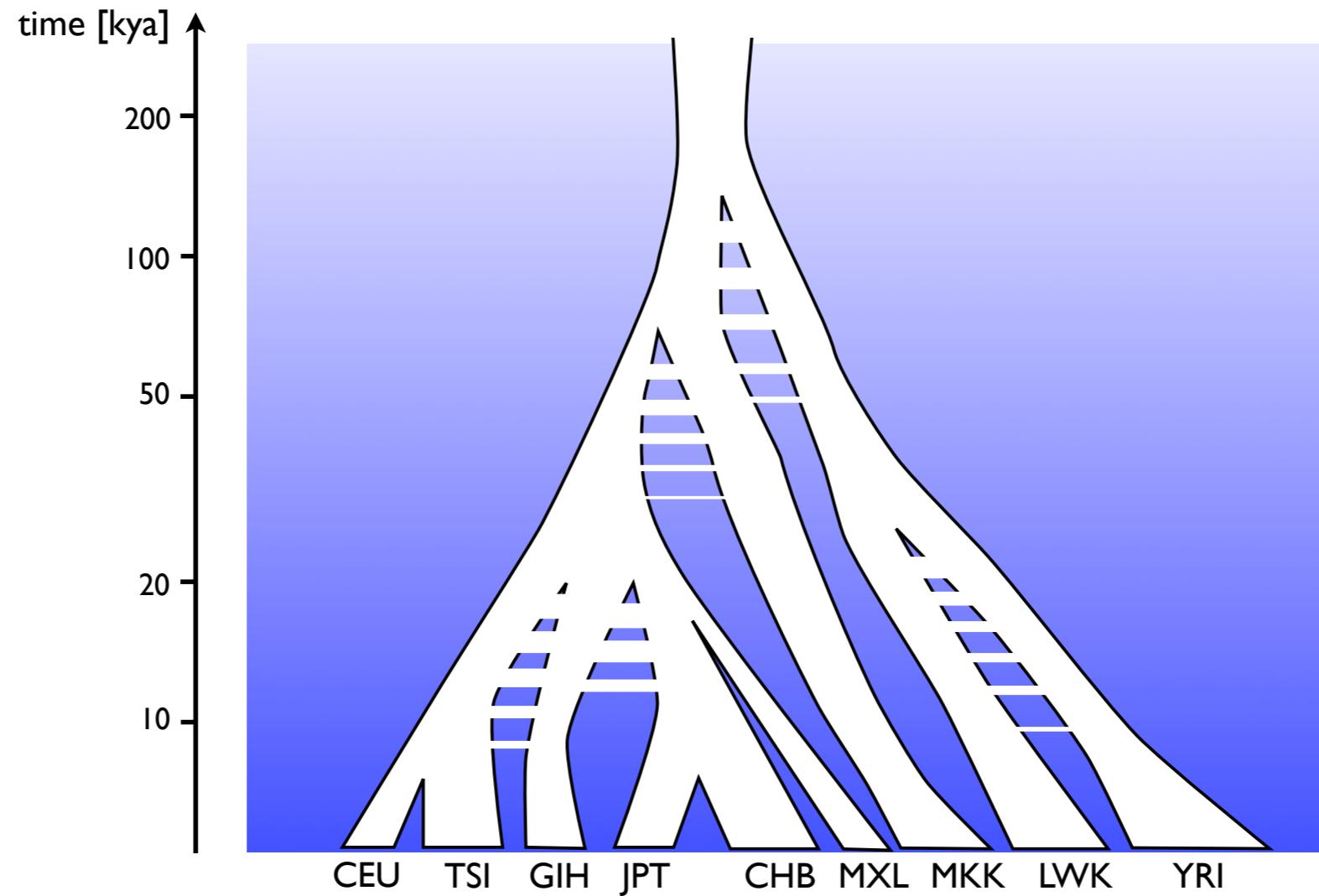


Implications of higher human mutation rate



A higher mutation rate of 2.5×10^{-8} would push these splits towards too recent times

MSMC Summary on separation history



[Schiffels and Durbin, *Nature Genetics*, in press]

Acknowledgements

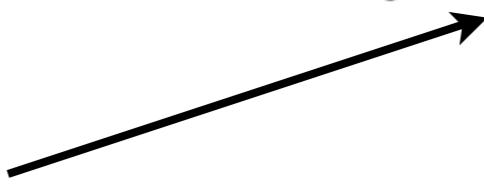
- Richard Durbin
- Aylwyn Scally
- Jeff Kidd, Simon Gravel, Eimear Kenny, Carlos Bustamante for MXL and PEL admixture tracts
- Article in press: Schiffels and Durbin, Nature Genetics
Preprint available on biorxiv.
- Software available on <https://github.com/stschiff/msmc>



Transition probability

$$q(i, j, t | k, l, s) = \delta(t - s) \delta_{i,k} \delta_{j,l} q_1(t) + q_2(t | s)$$

Probability to remain in state (i,j,t)



Probability to change time (and pair) of first coalescence

$$q_1(t) = e^{-M r t} + (1 - e^{-M r t}) \frac{1}{t} \frac{1}{M} \int_0^t \left(1 + (M - 3) \exp \left(-M \int_u^t \lambda(v) dv \right) \right) du$$

$$q_2(t | s) = (1 - e^{-M r s}) \frac{1}{s} \frac{1}{M} 2 \lambda(t) \begin{cases} \int_0^t \exp \left(-M \int_u^t \lambda(v) dv \right) du & \text{if } t < s \\ \exp \left(-\left(\frac{M}{2} \right) \int_s^t \lambda(v) dv \right) \int_0^s \exp \left(-M \int_u^s \lambda(v) dv \right) du & \text{if } t > s. \end{cases}$$

depends on:

- coalescence rates $\lambda(t)$
- recombination rate r
- Number of sequences M