

# INFERENCE OF EVOLUTIONARY HISTORY WITH APPROXIMATE BAYESIAN COMPUTATION

Ariella Gladstein

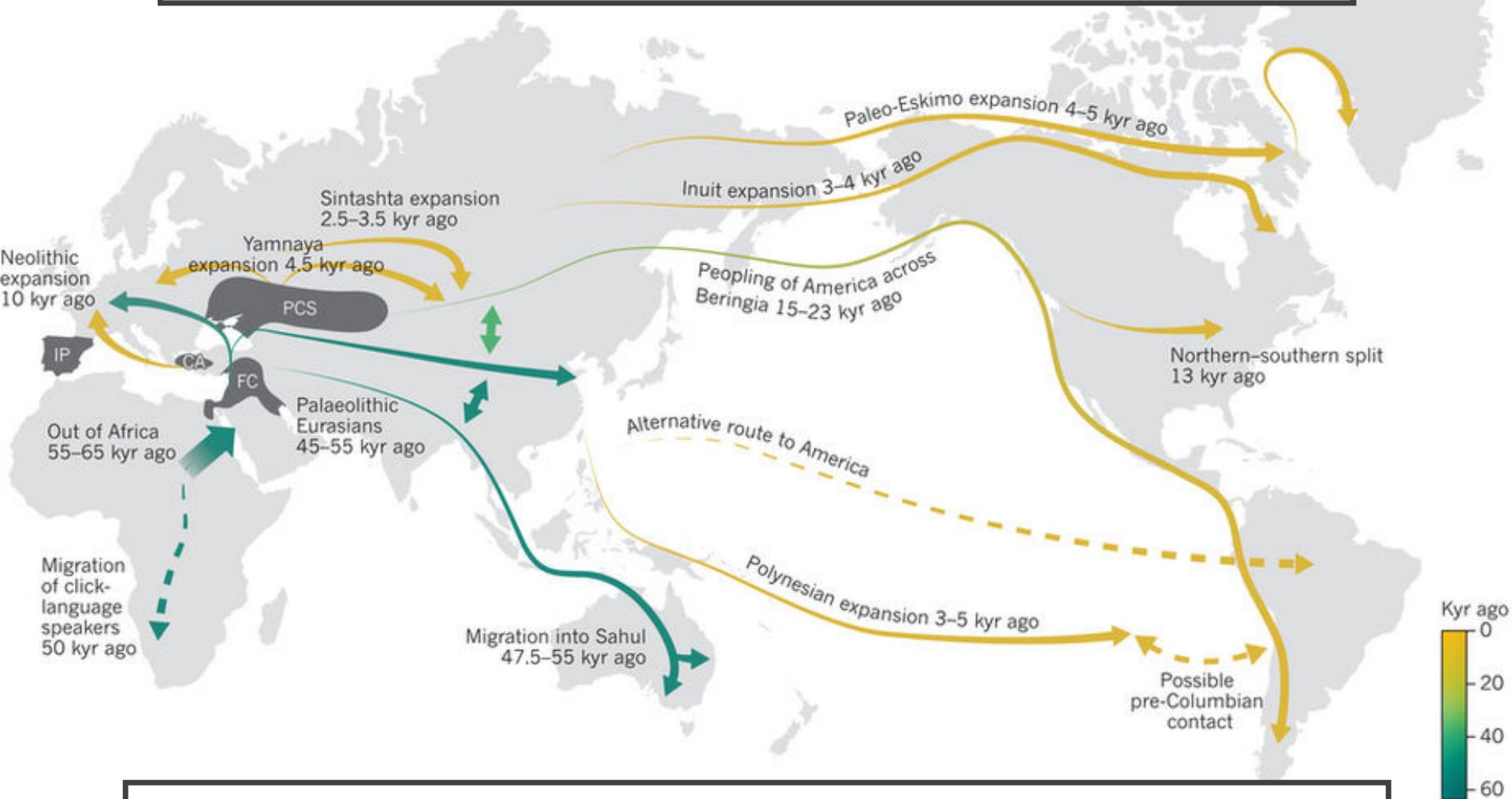
Ecology and Evolutionary Biology

University of Arizona



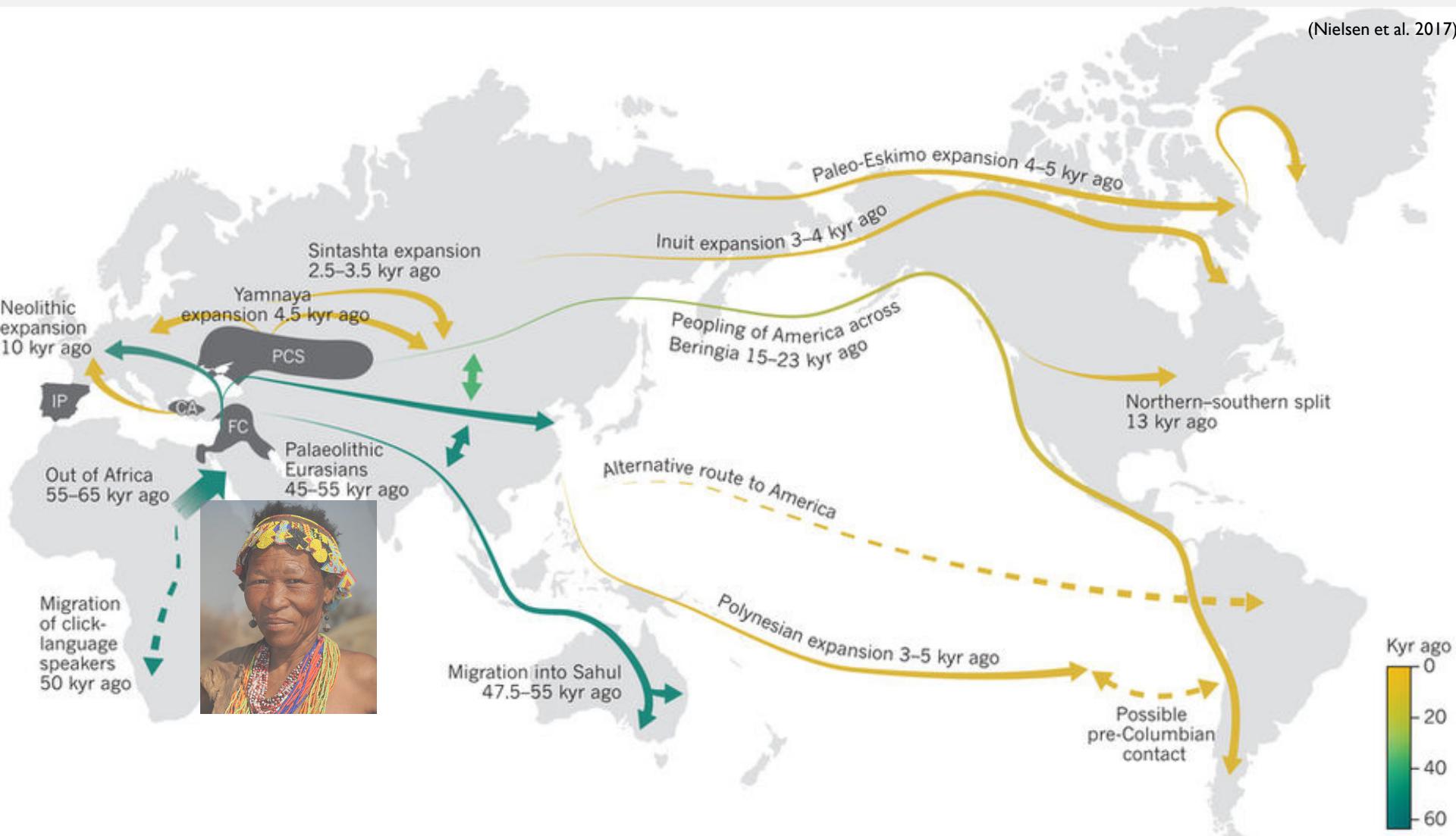
# HOW DID HUMANS SPREAD ACROSS THE WORLD?

(Nielsen et al. 2017)

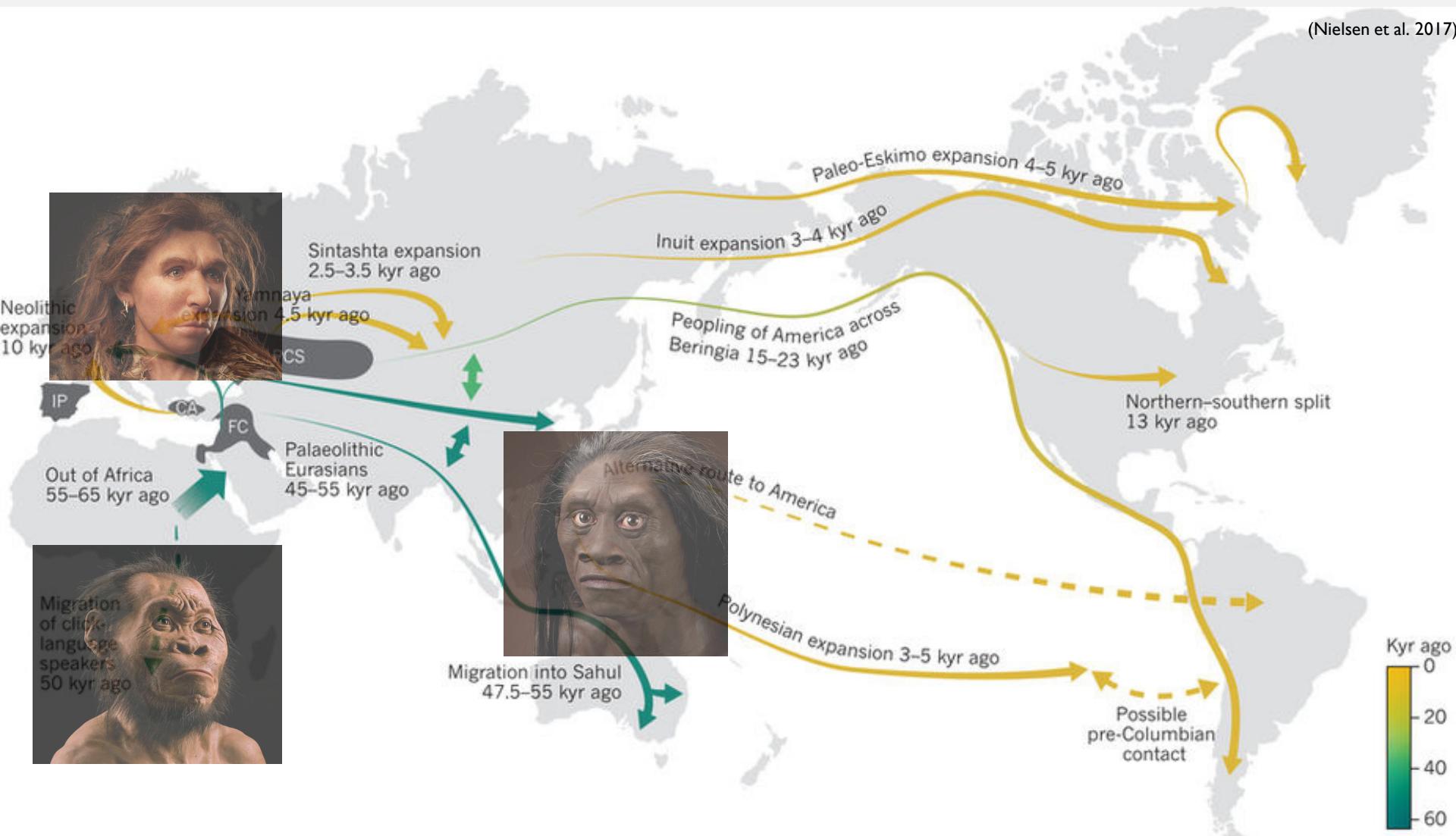


WHAT DEMOGRAPHIC EVENTS LEAD US TO WHERE WE ARE TODAY AND THE DIVERSITY WE SEE?

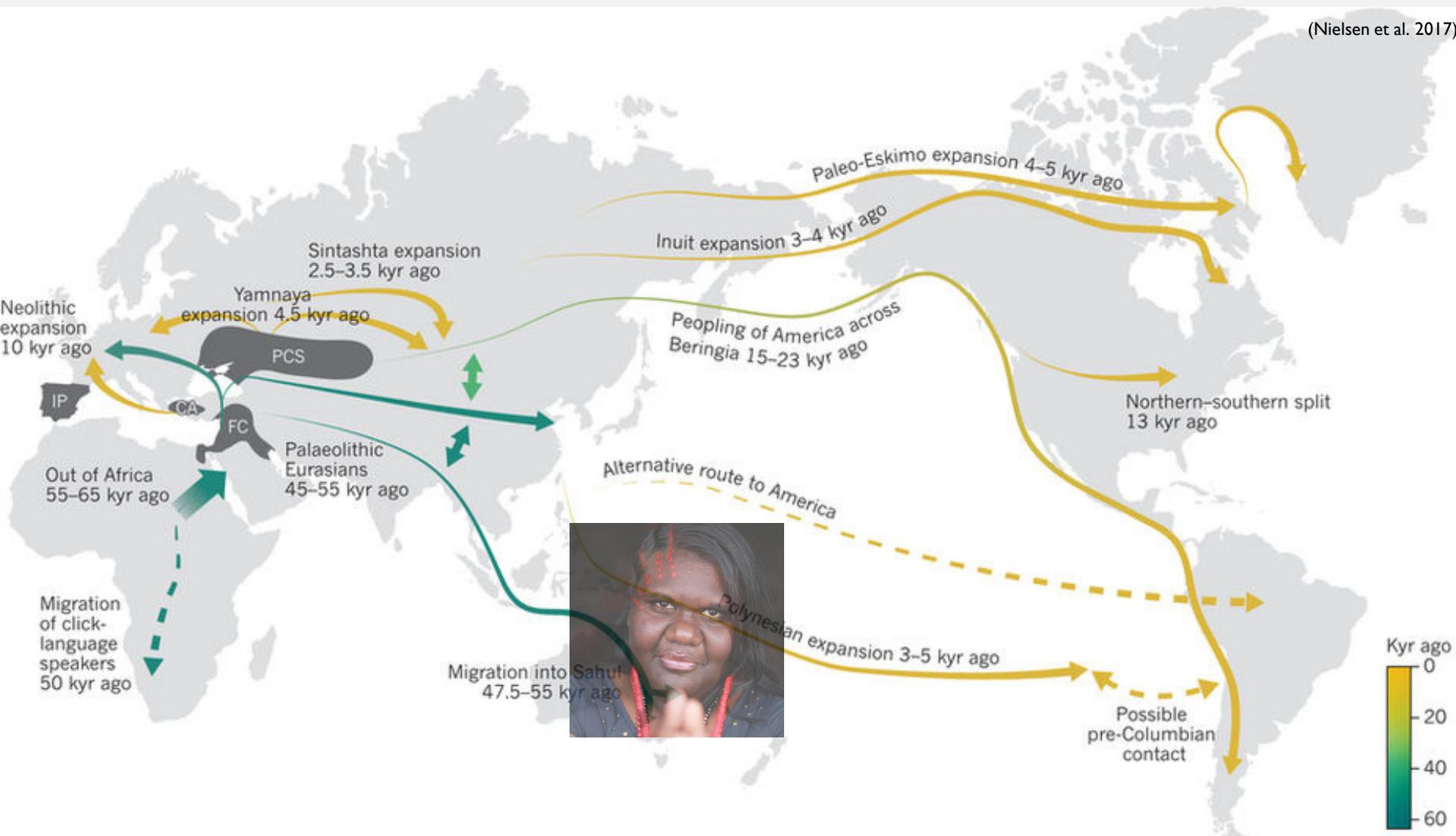
(Nielsen et al. 2017)



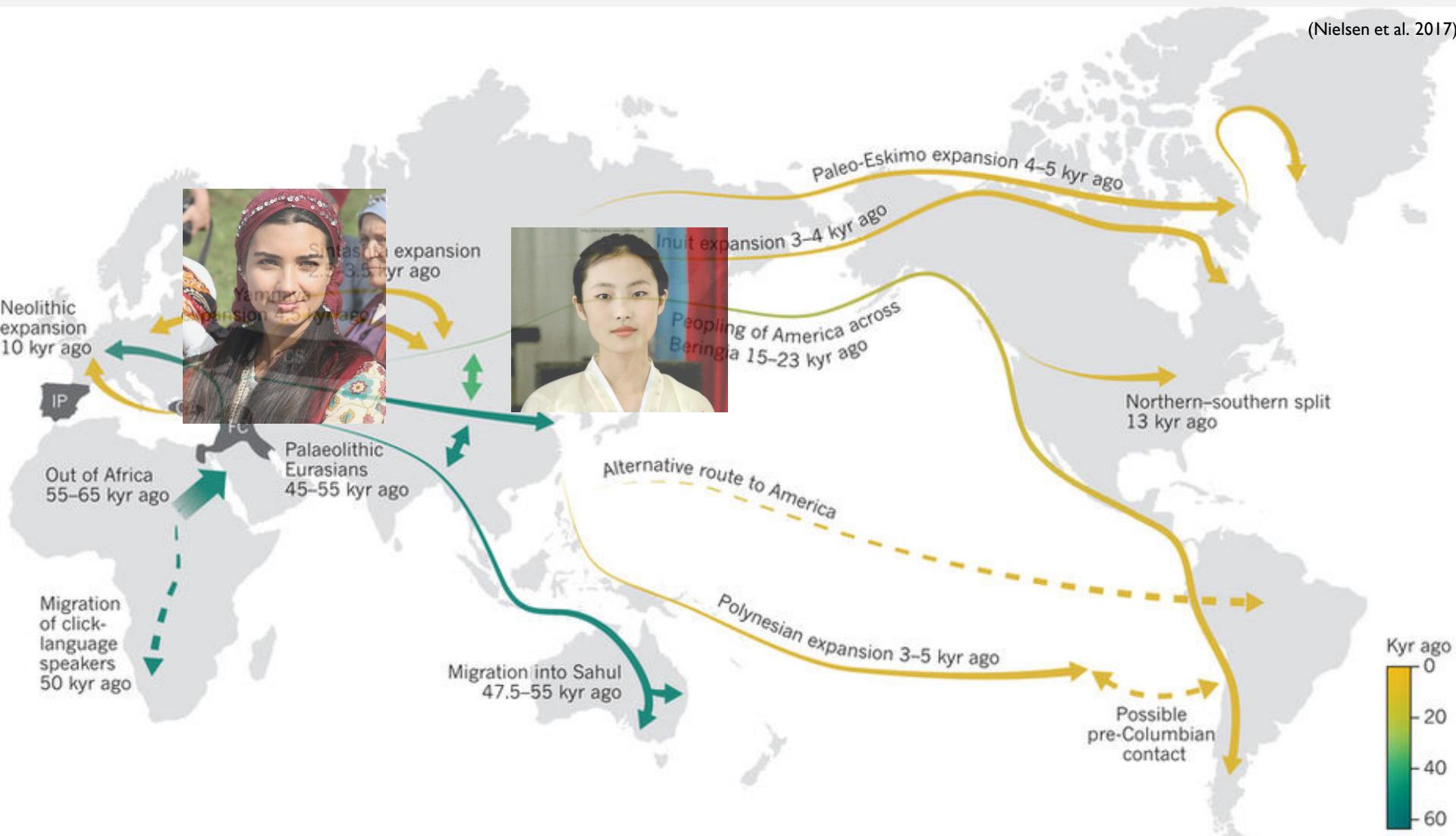
(Nielsen et al. 2017)

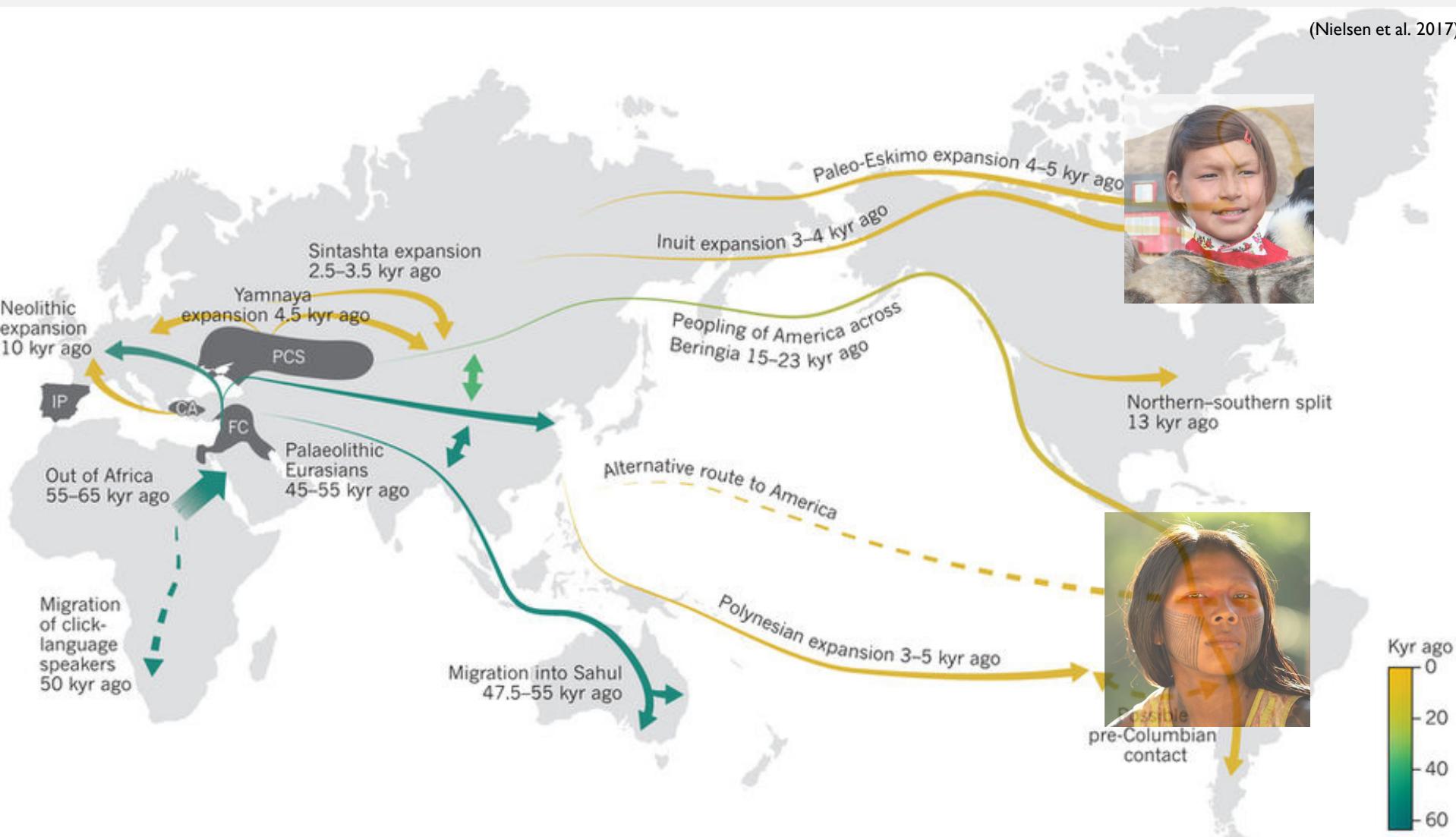


(Nielsen et al. 2017)

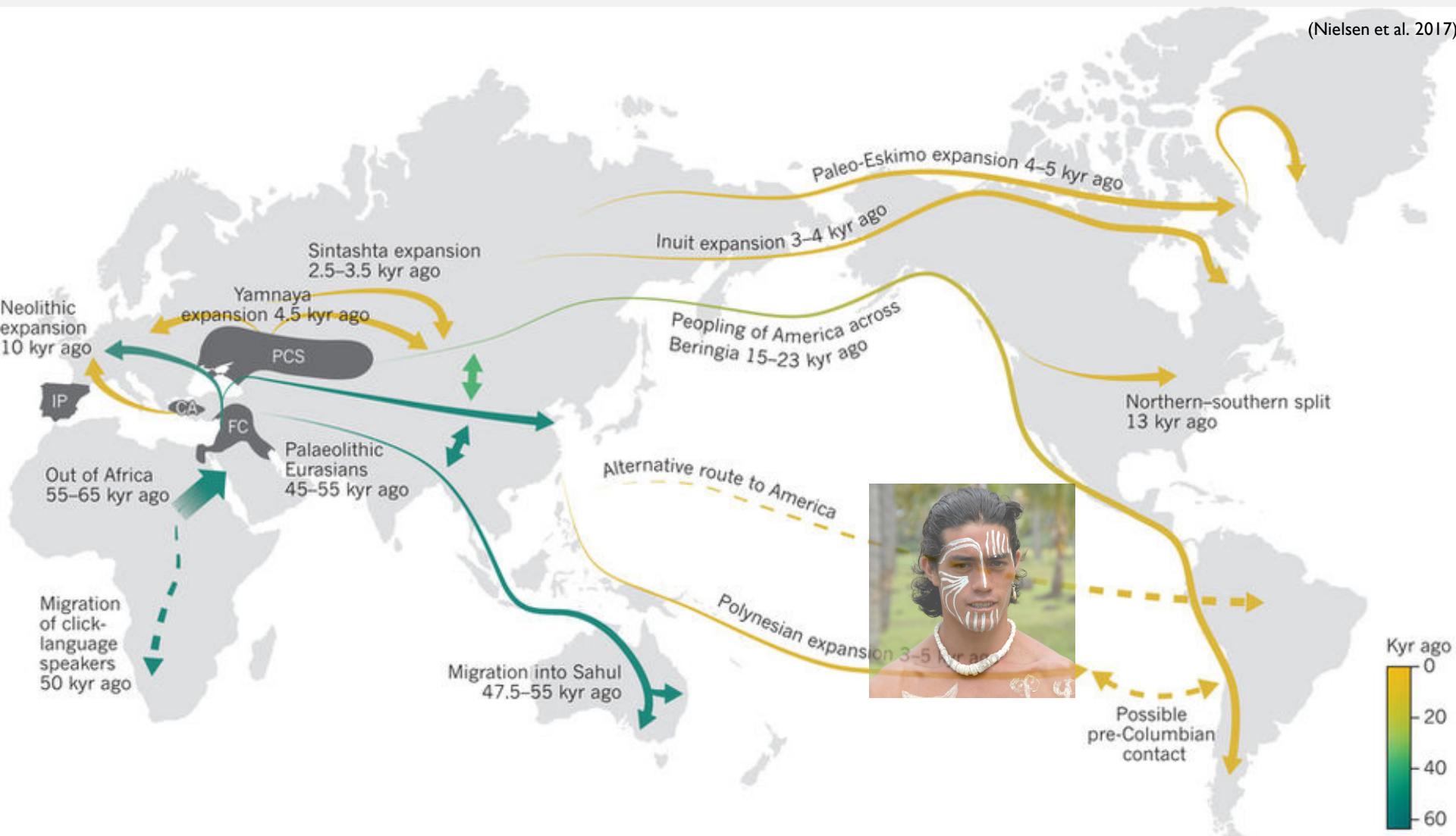


(Nielsen et al. 2017)

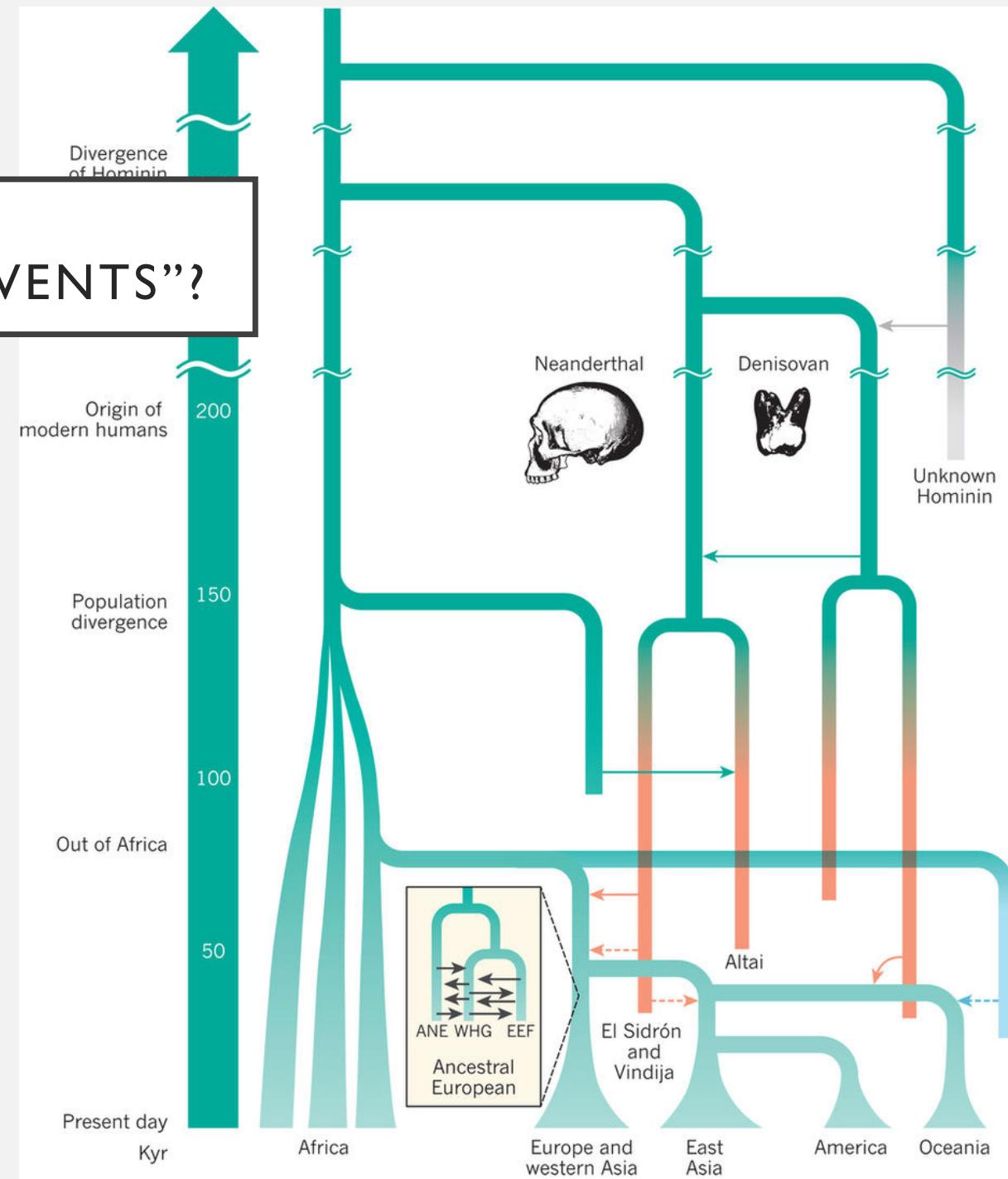




(Nielsen et al. 2017)

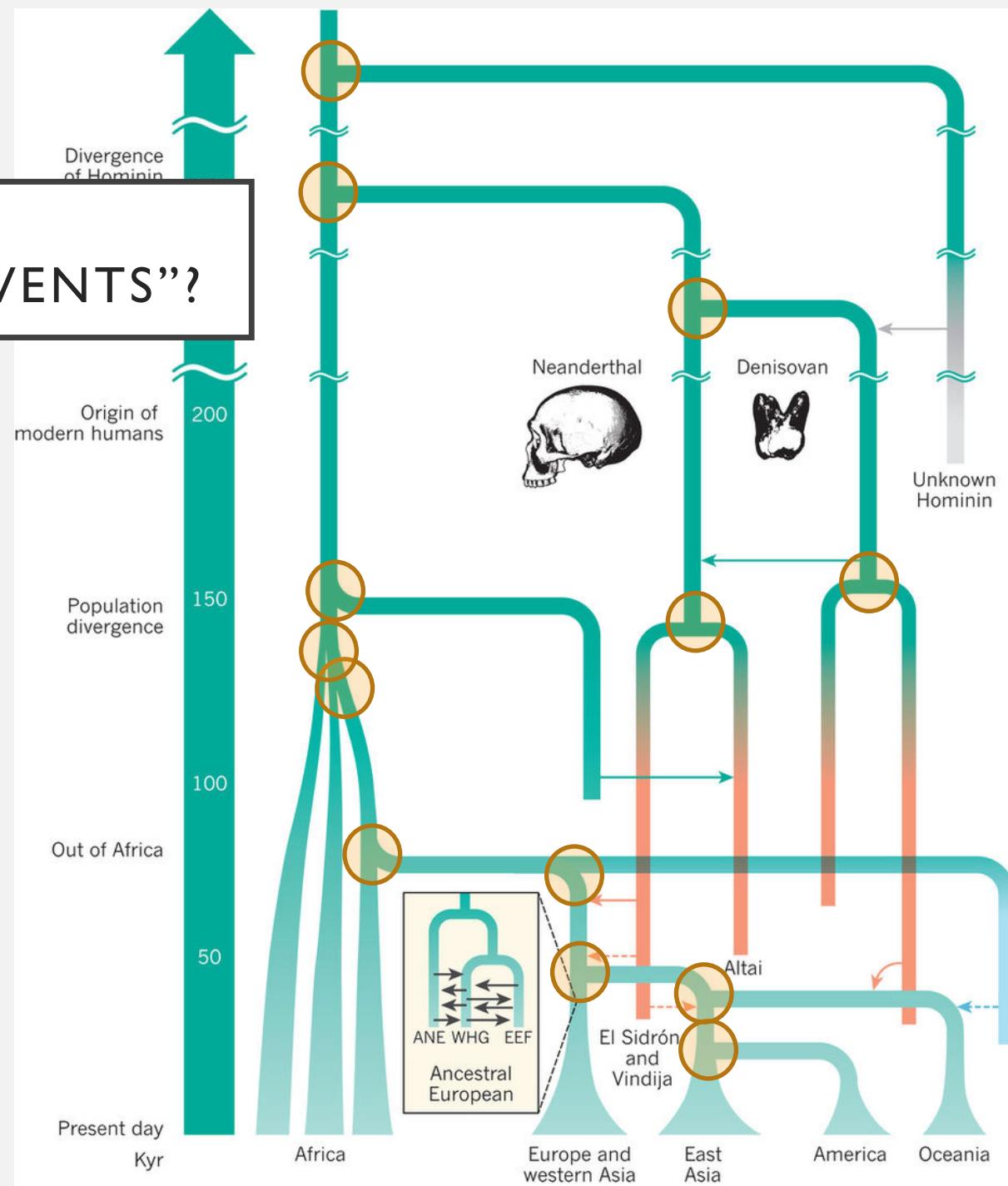


# WHAT ARE “DEMOGRAPHIC EVENTS”?



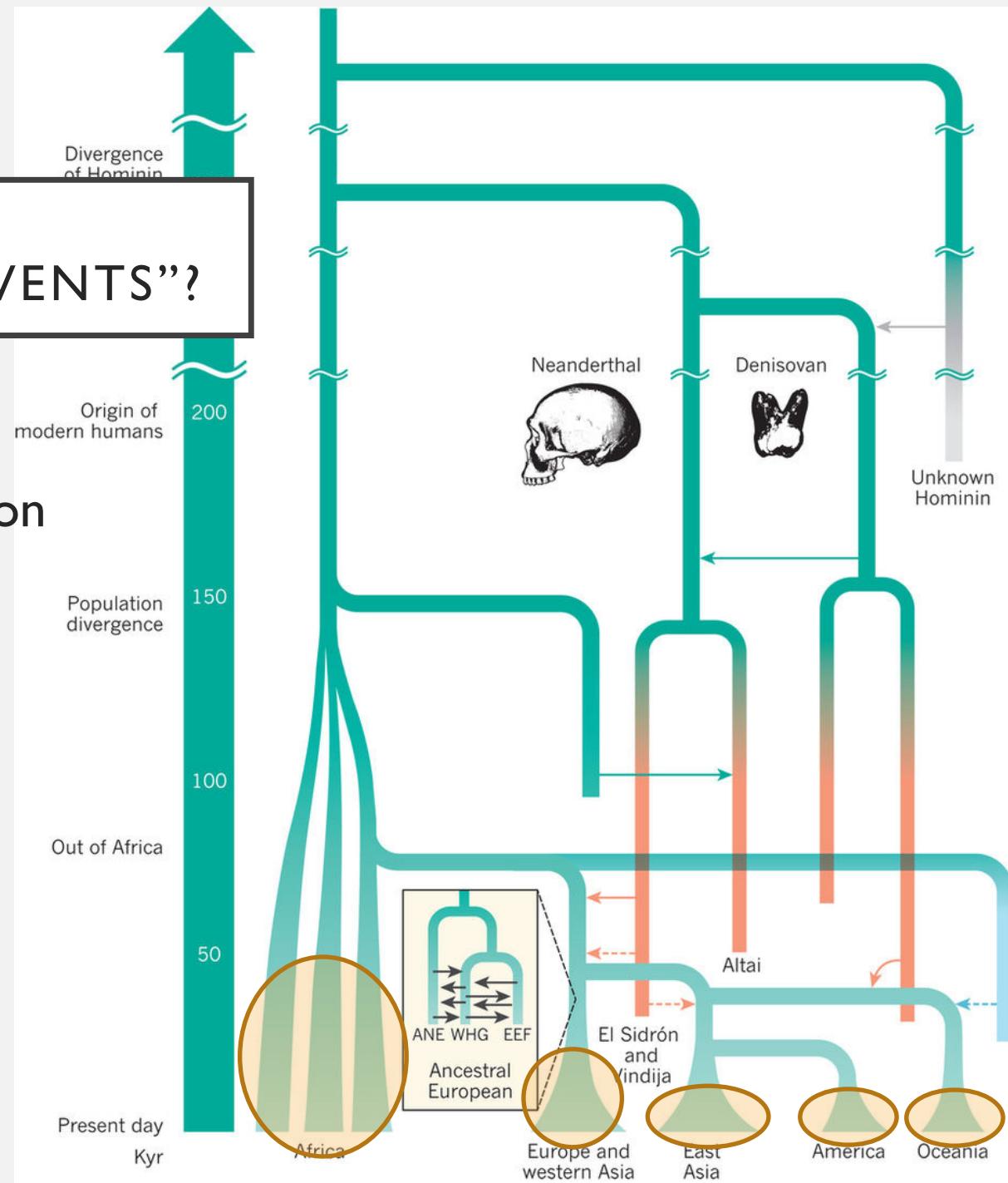
# WHAT ARE “DEMOGRAPHIC EVENTS”?

- Divergence



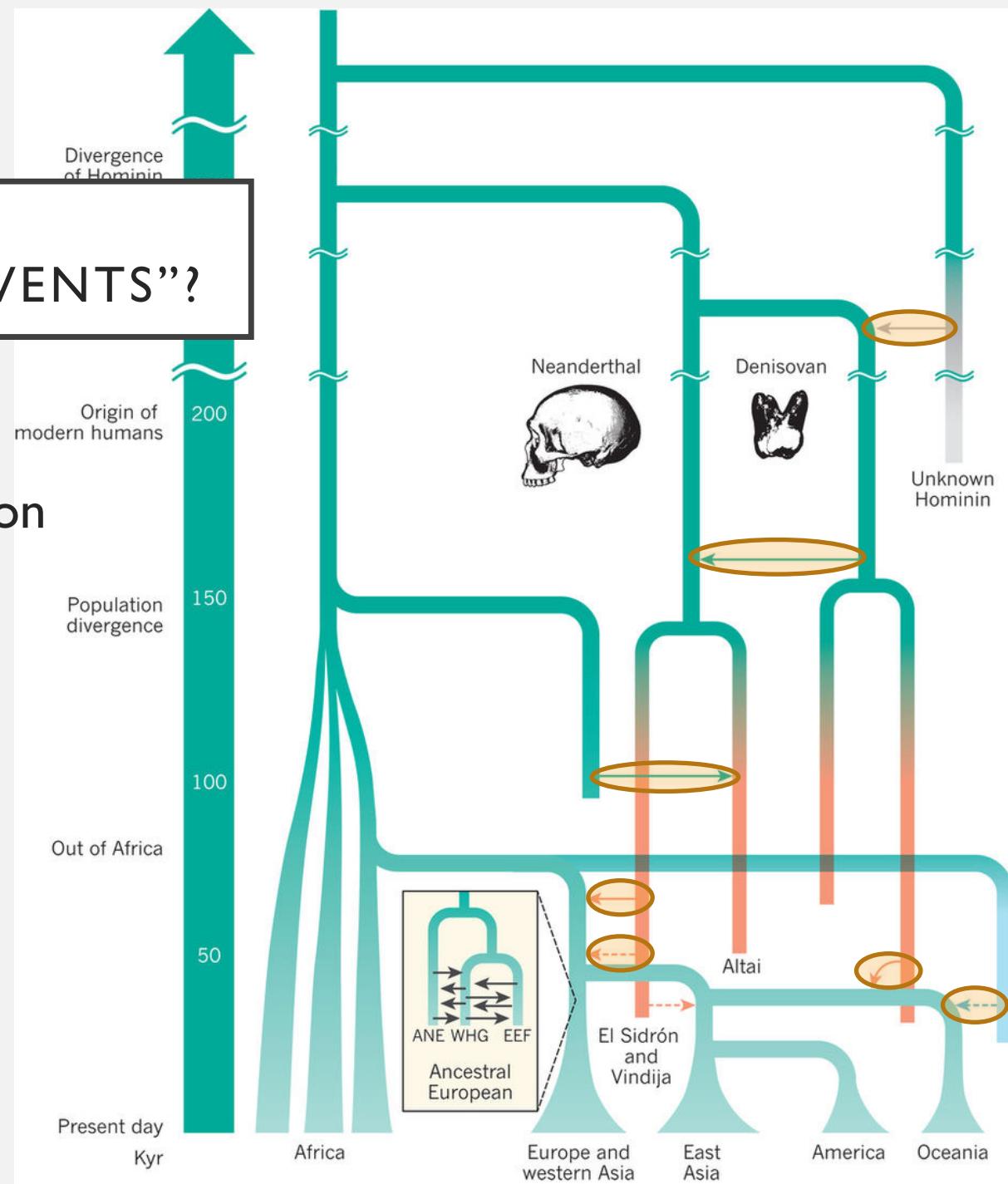
# WHAT ARE “DEMOGRAPHIC EVENTS”?

- Divergence
- Expansion or reduction

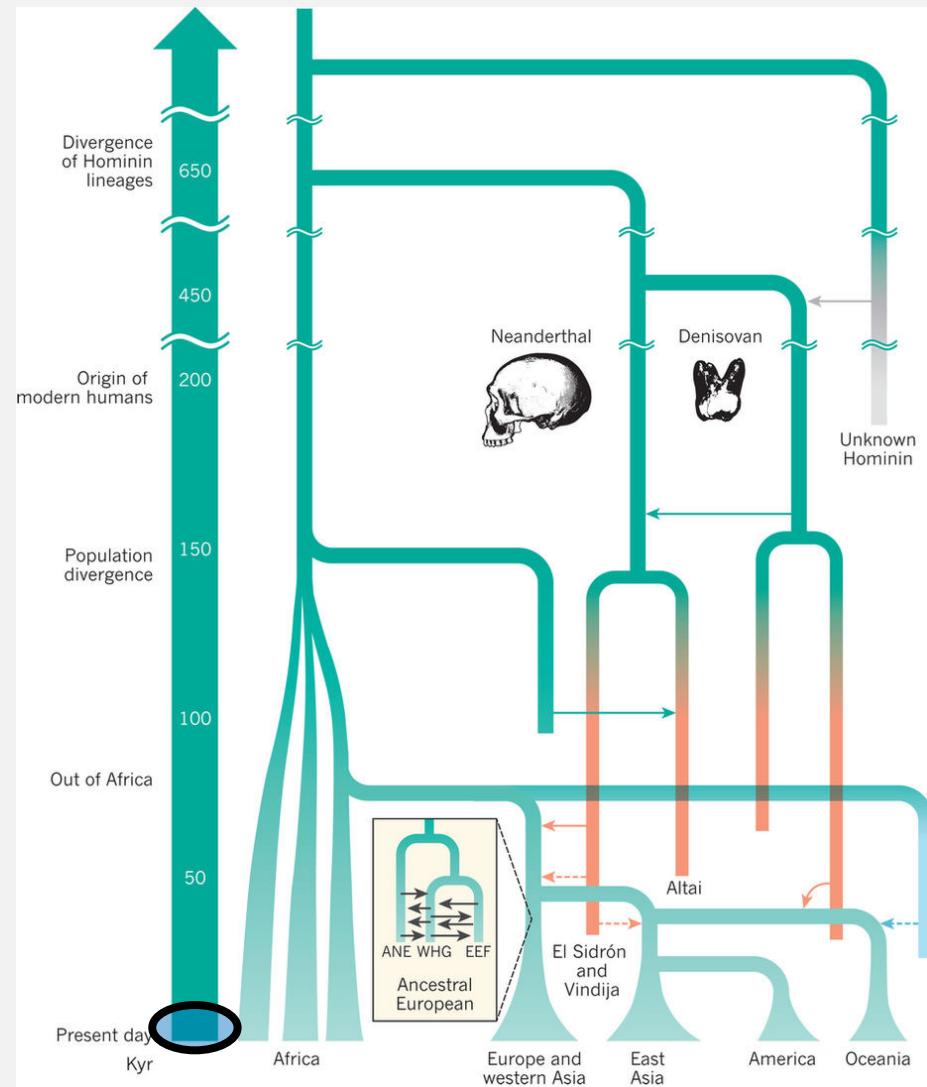
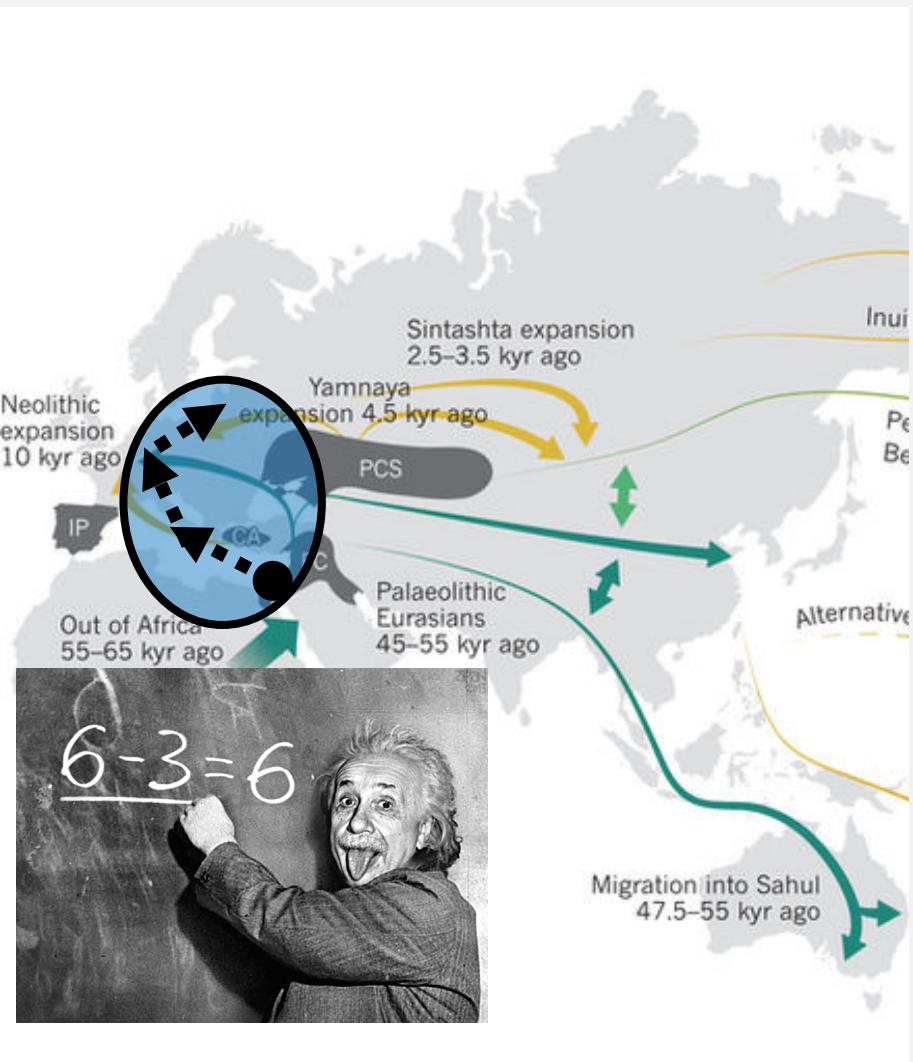


## WHAT ARE “DEMOGRAPHIC EVENTS”?

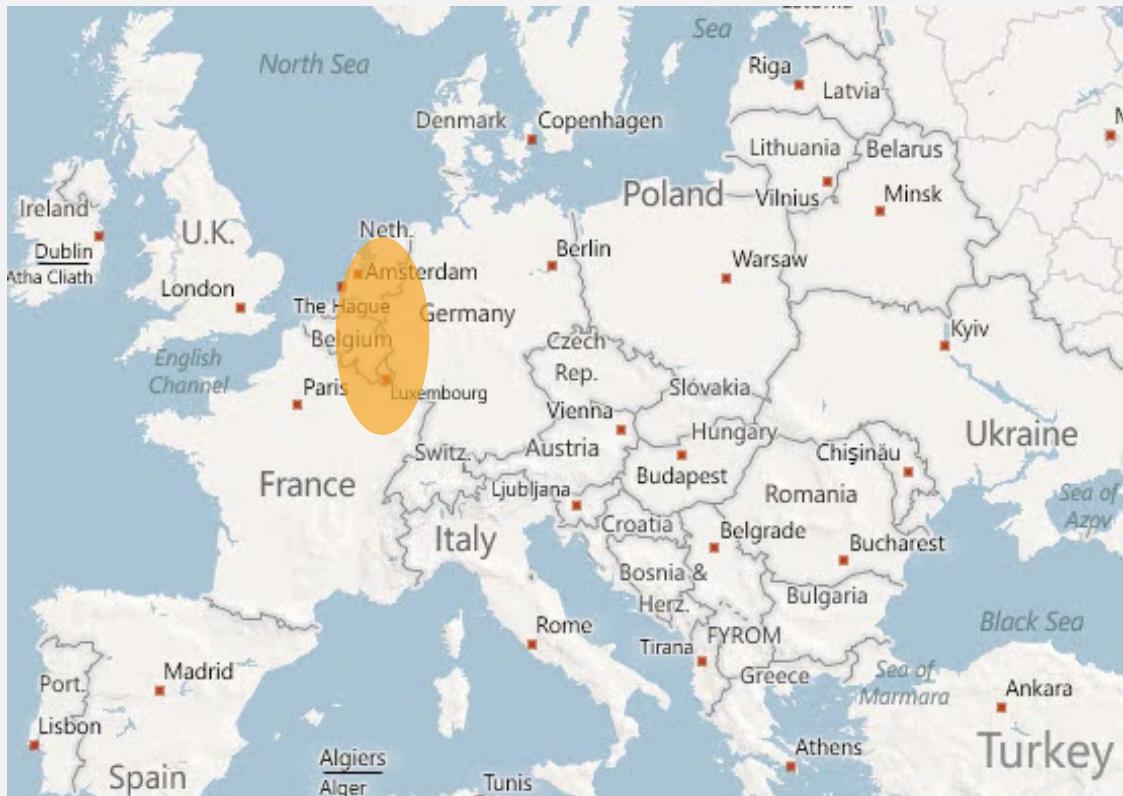
- Divergence
- Expansion or reduction
- Gene flow



# AIM: INFER THE DEMOGRAPHIC HISTORY OF THE ASHKENAZI JEWS.

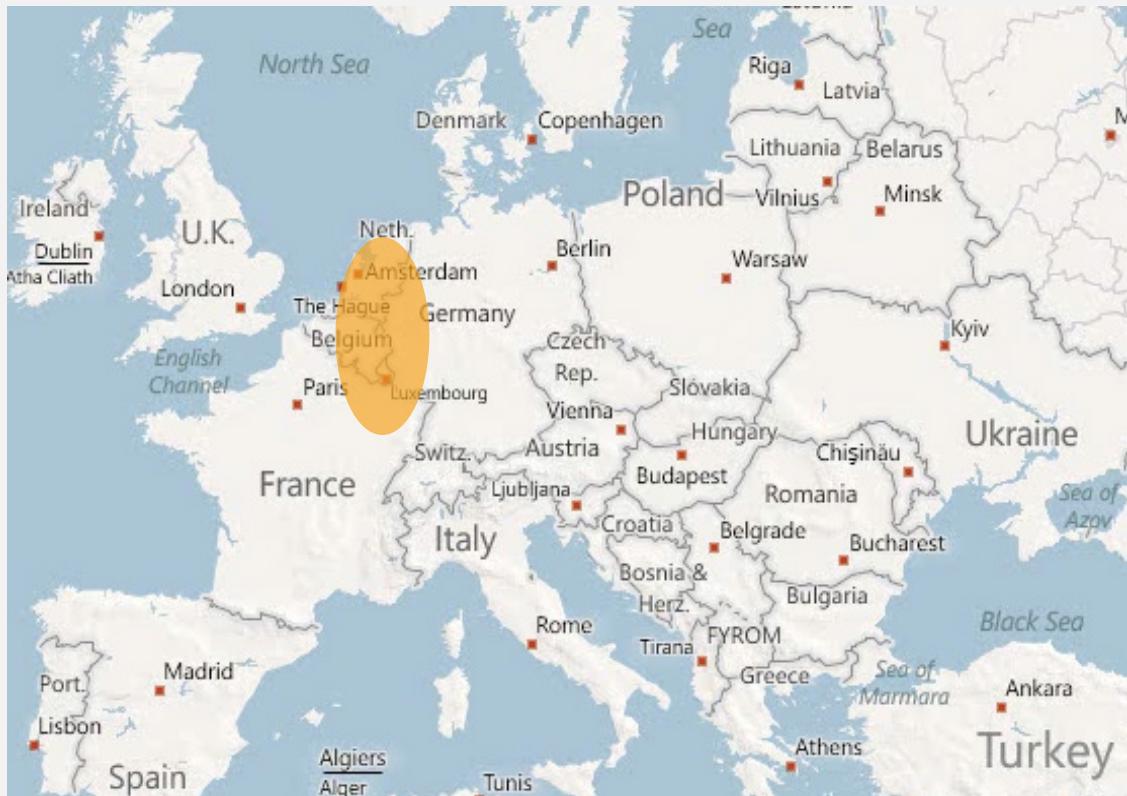


# ASHKENAZI JEWS: AN INTERESTING STUDY POPULATION



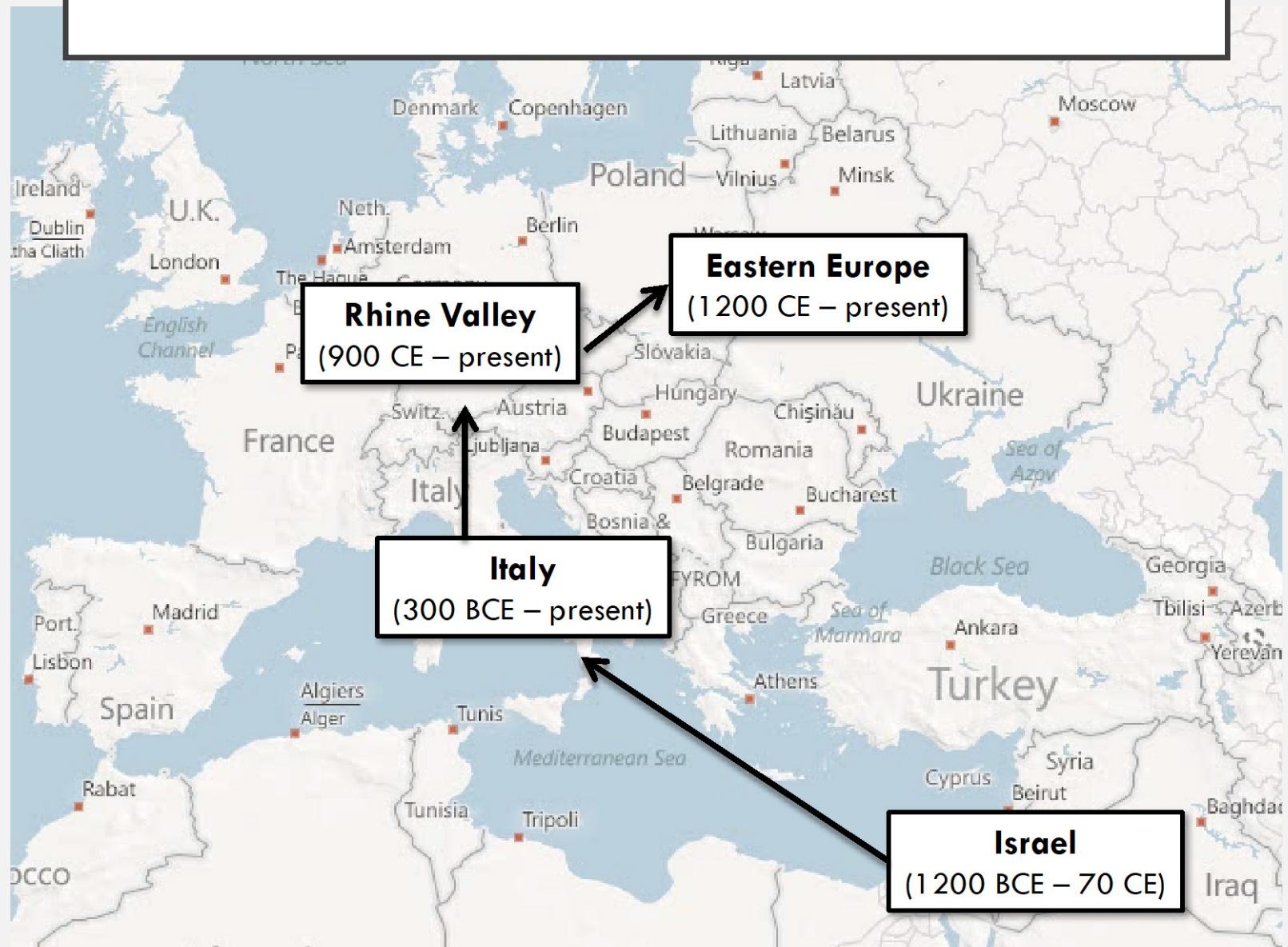
- High frequency of genetic disorders
- Population isolate
- Complex demographic history
- Well documented historical record

# ASHKENAZI JEWS: AN INTERESTING STUDY POPULATION



- High frequency of genetic disorders
- Population isolate
- Complex demographic history
- Well documented historical record

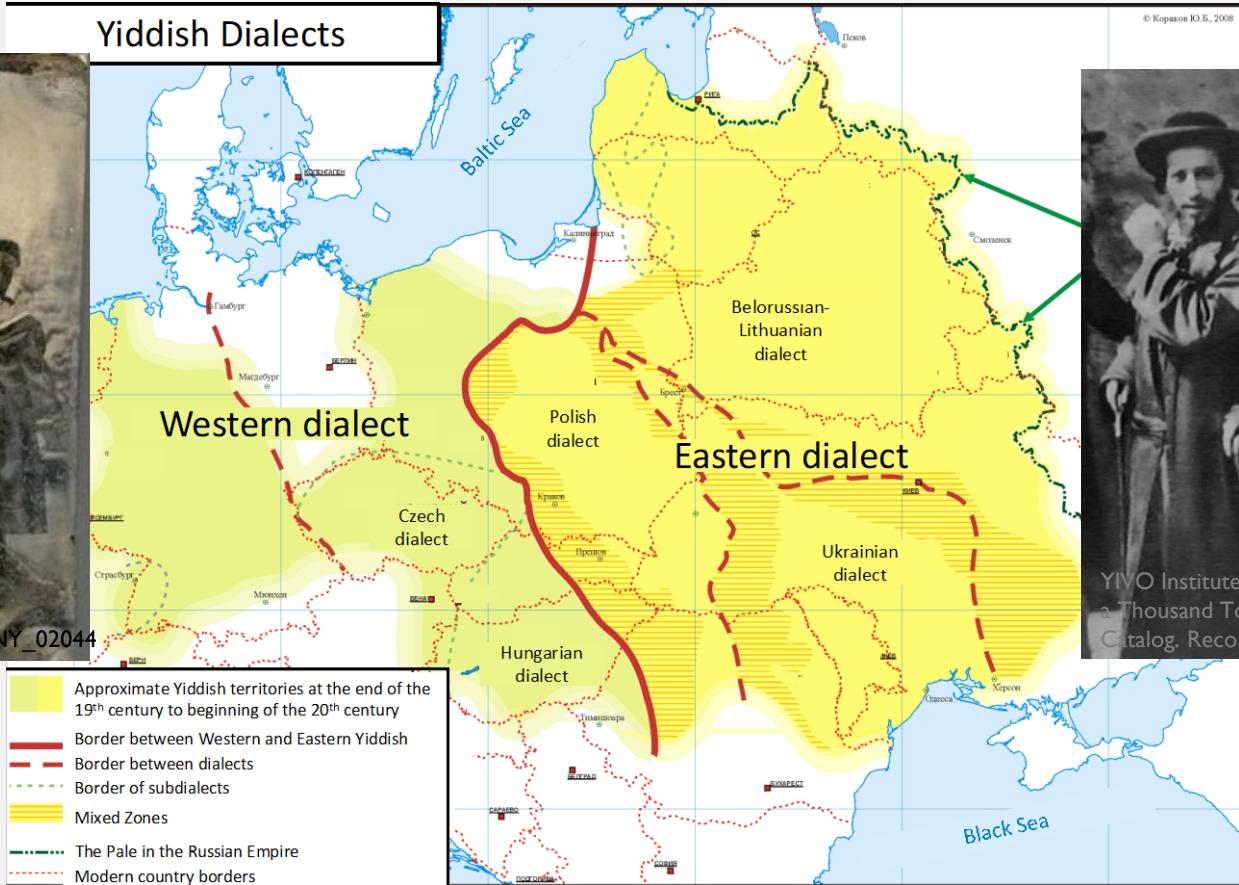
# HYPOTHESIS OF ASHKENAZI ORIGINS



# WESTERN VS. EASTERN ASHKENAZI JEWS



Yiddish Dialects

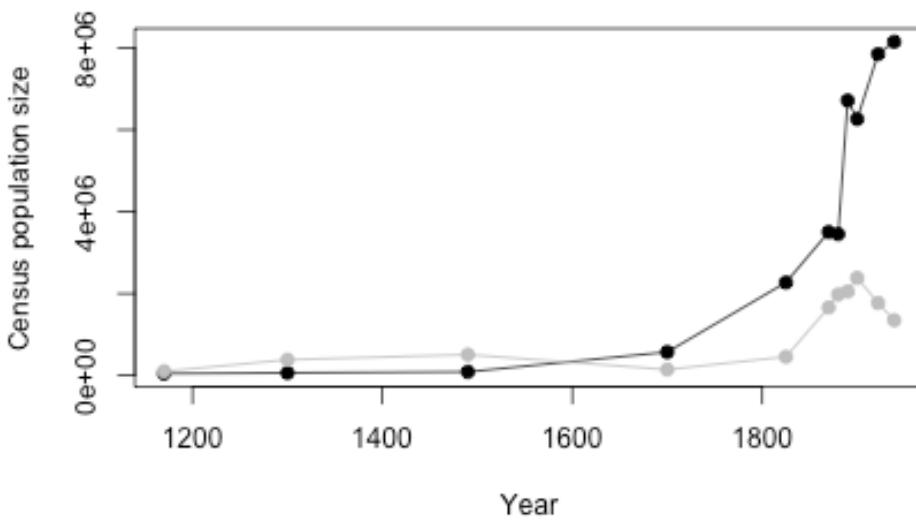


# WESTERN VS. EASTERN ASHKENAZI JEWS



JDC Archives. Reference Code: NY\_0204

Germany, 1900's



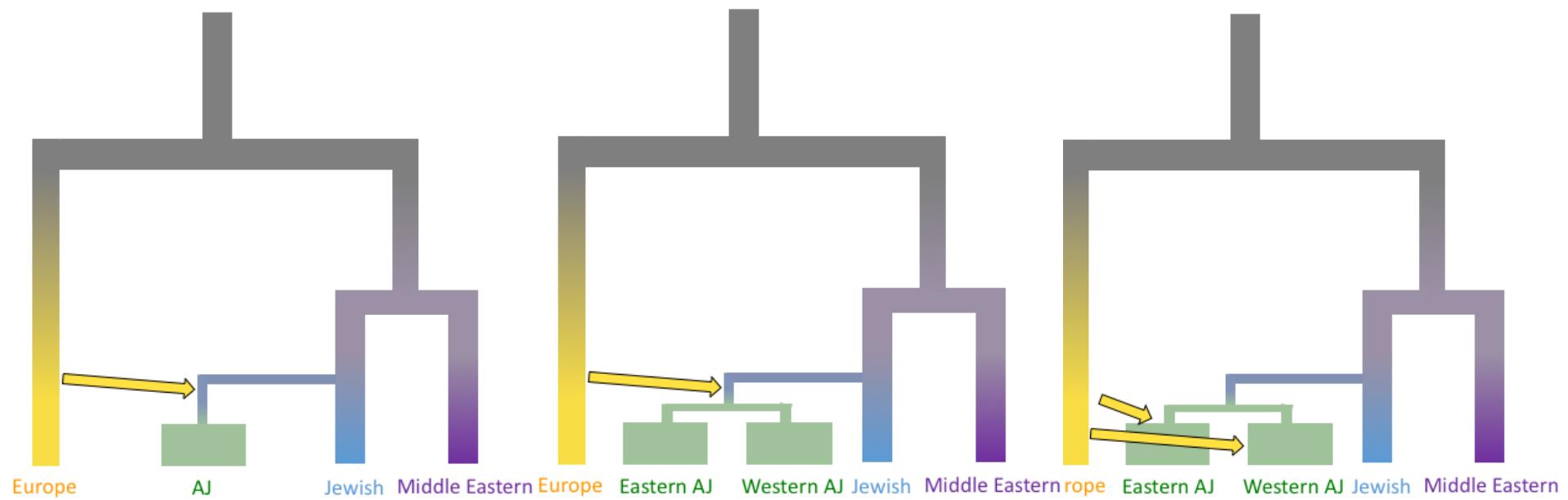
YIVO Institute for Jewish Research. People of a Thousand Towns. Online Photographic Catalog. Record Id: 6820

Cracow, Poland. 1932

Reference census data

## MOTIVATION

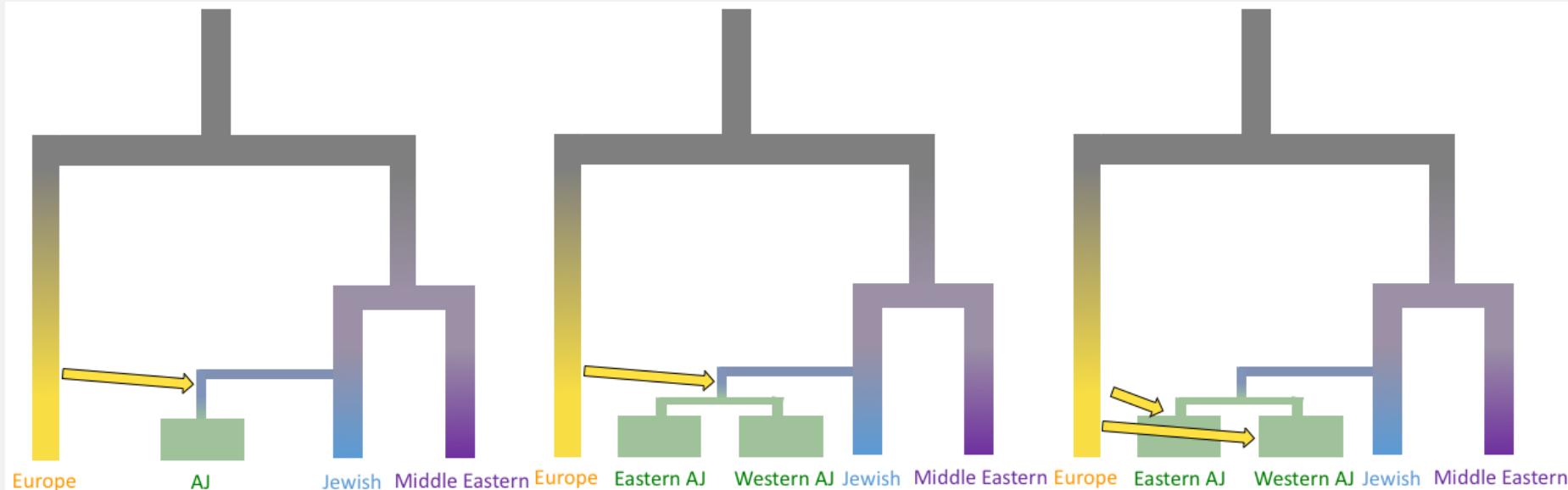
- Numerous genetic studies on the Ashkenazi Jews.
  - All genome-wide studies treat Ashkenazi Jews as one population.
- Preliminary work consistent with genetic differentiation.
  - Not informative of cause of differentiation.



## MODELS OF ASHKENAZI HISTORY

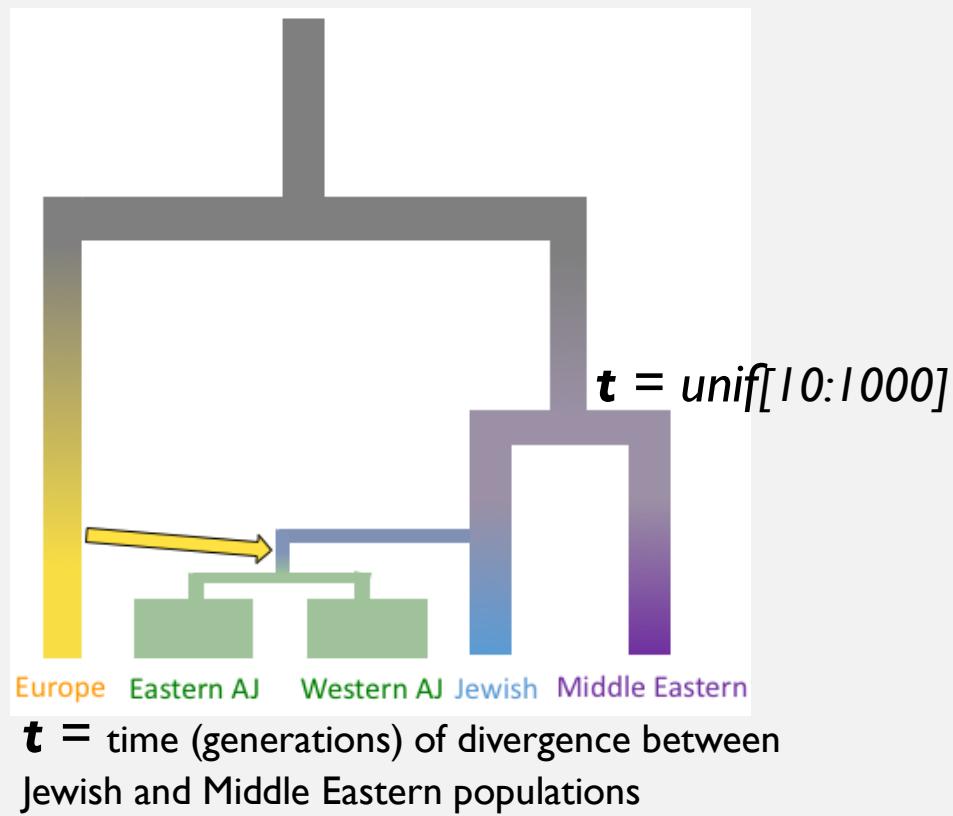
# APPROXIMATE BAYESIAN COMPUTATION

- Infer parameter values
- Choose among models



# APPROXIMATE BAYESIAN COMPUTATION

## I. Define priors of parameters of model



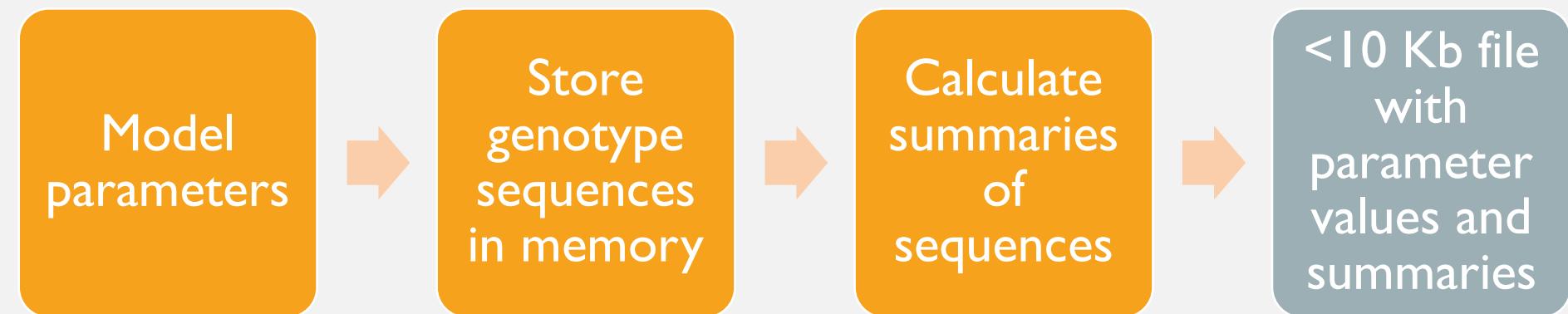
## APPROXIMATE BAYESIAN COMPUTATION

1. Define priors of parameters of model
2. Simulate data many times

## APPROXIMATE BAYESIAN COMPUTATION

1. Define priors of parameters of model
2. Simulate data many times
3. Choose model and estimate parameters based on simulations closest to real data

# SIMULATION

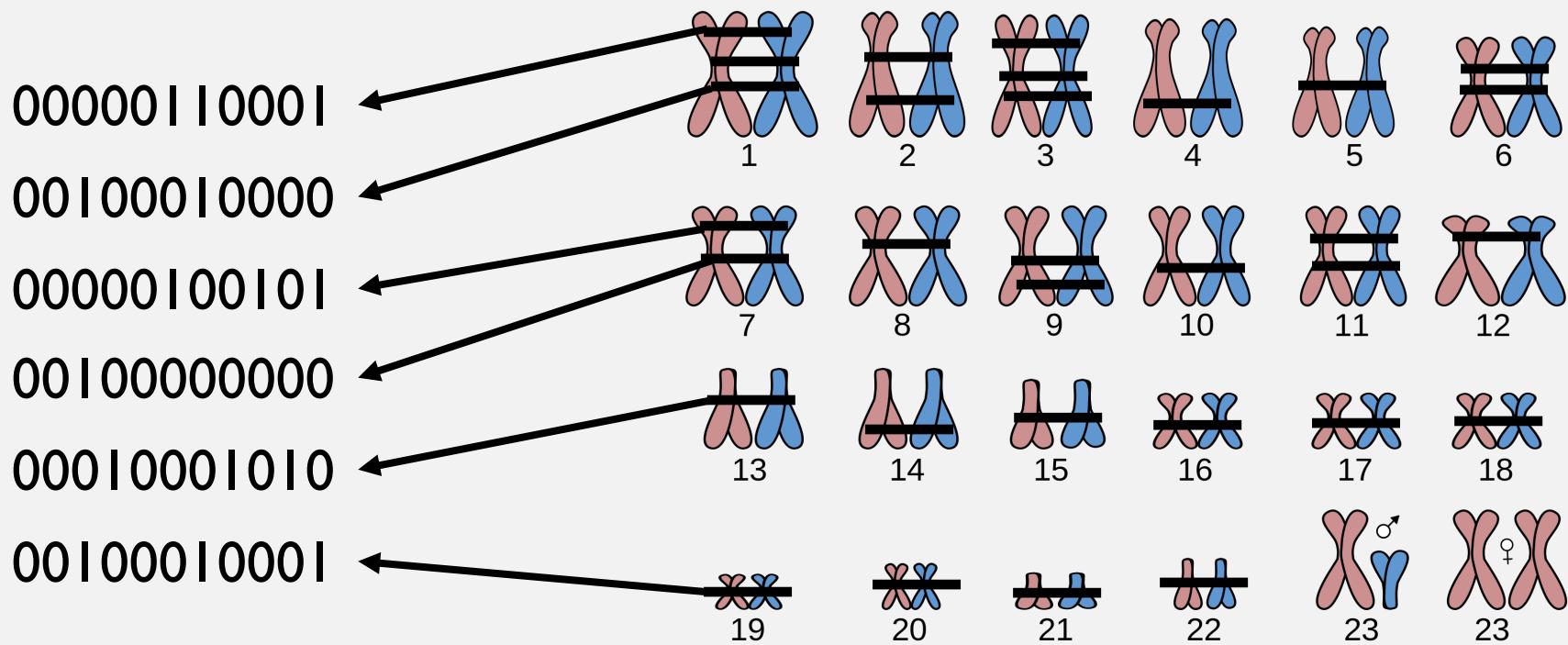


# EMBARRASSINGLY PARALLEL!



# INHERITED SCRIPT INTENDED FOR SMALL SEQUENCE

1,389 10kb regions



# SIMULATE WHOLE CHROMOSOME

~250 million sites on human chromosome 1

# PROBLEM!

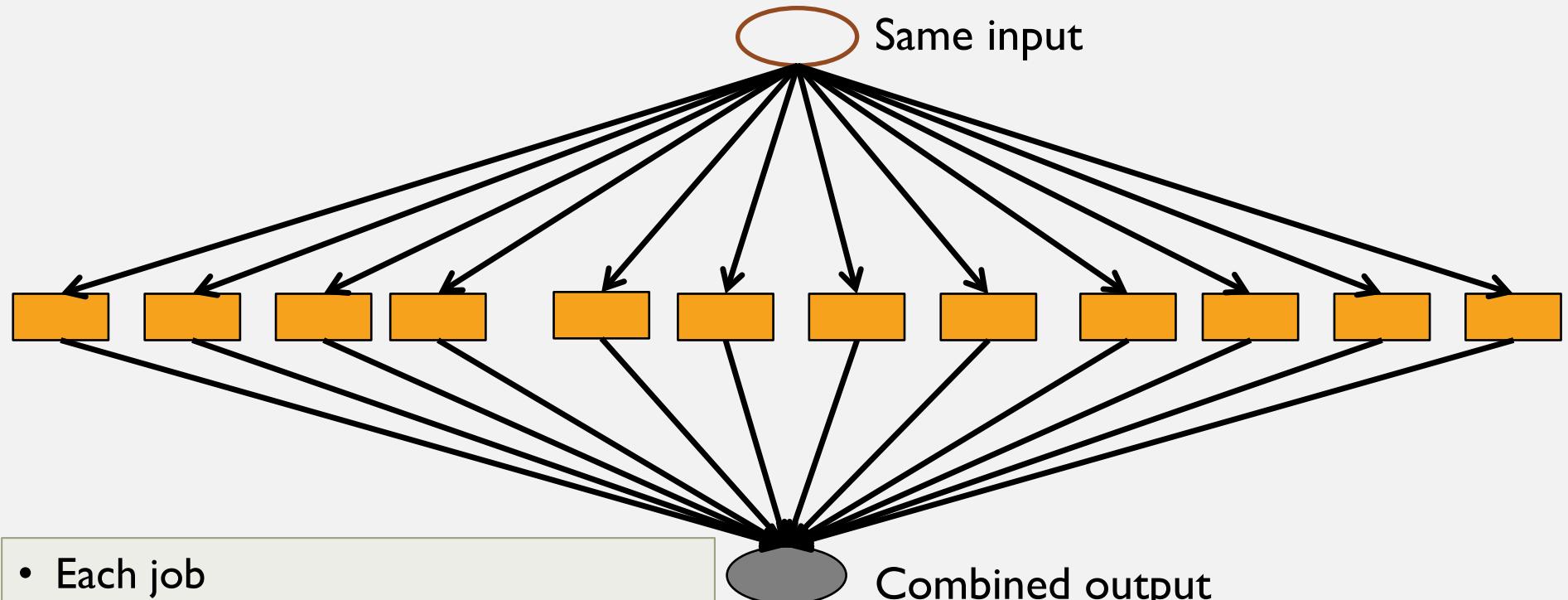
Parameters	Average Walltime	Average Memory
Minimum	00:21:00	2.7 Gb
Random	00:55:11	20 Gb
Maximum	08:02:11	117 Gb

Too much memory!

Over a decade to complete  
6000 runs/month w/ UA resources

Each core on UA HPC has 6G - **Need memory < 6G** for each run

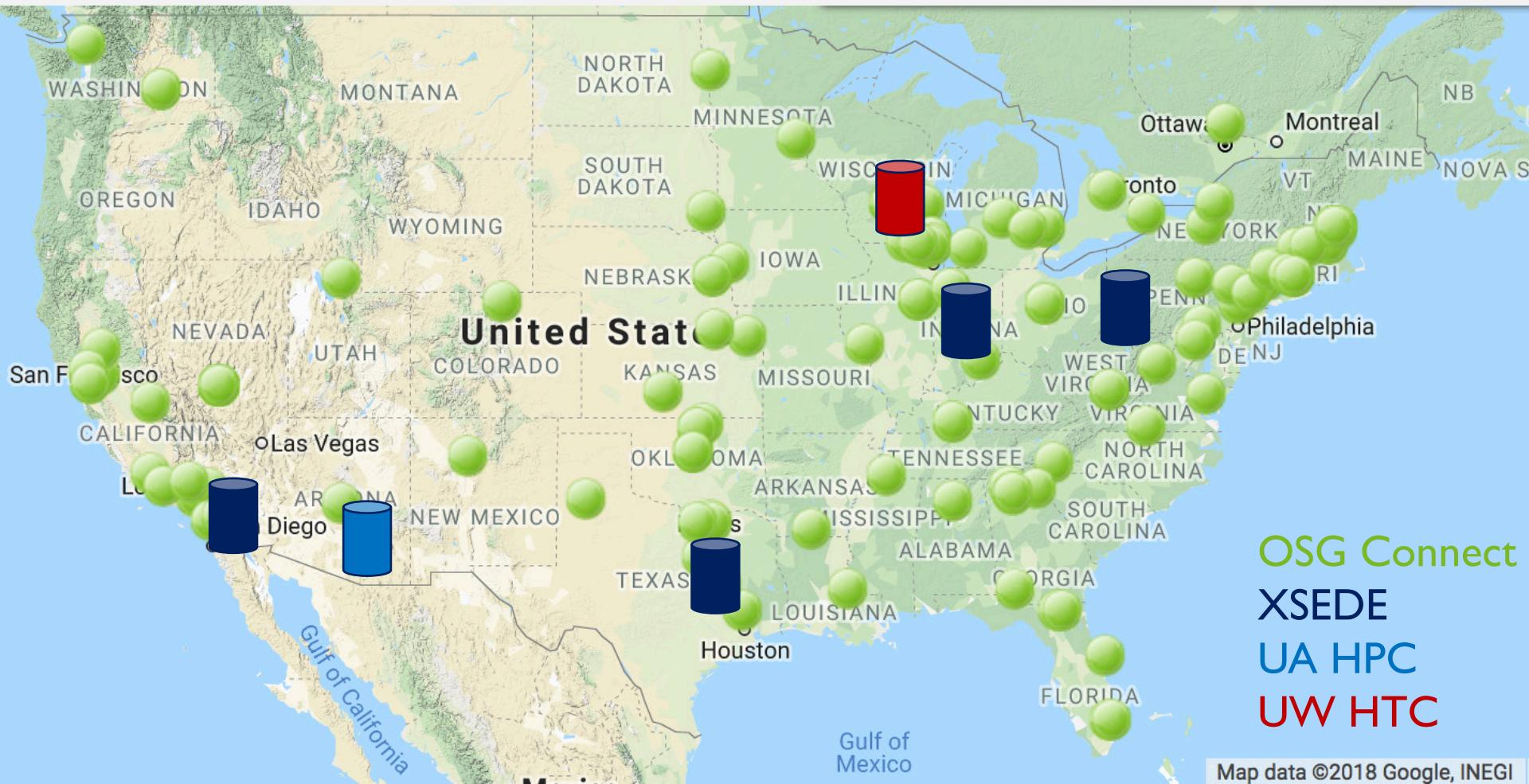
# EMBARRASSINGLY PARALLEL & RESOURCE LIGHT!



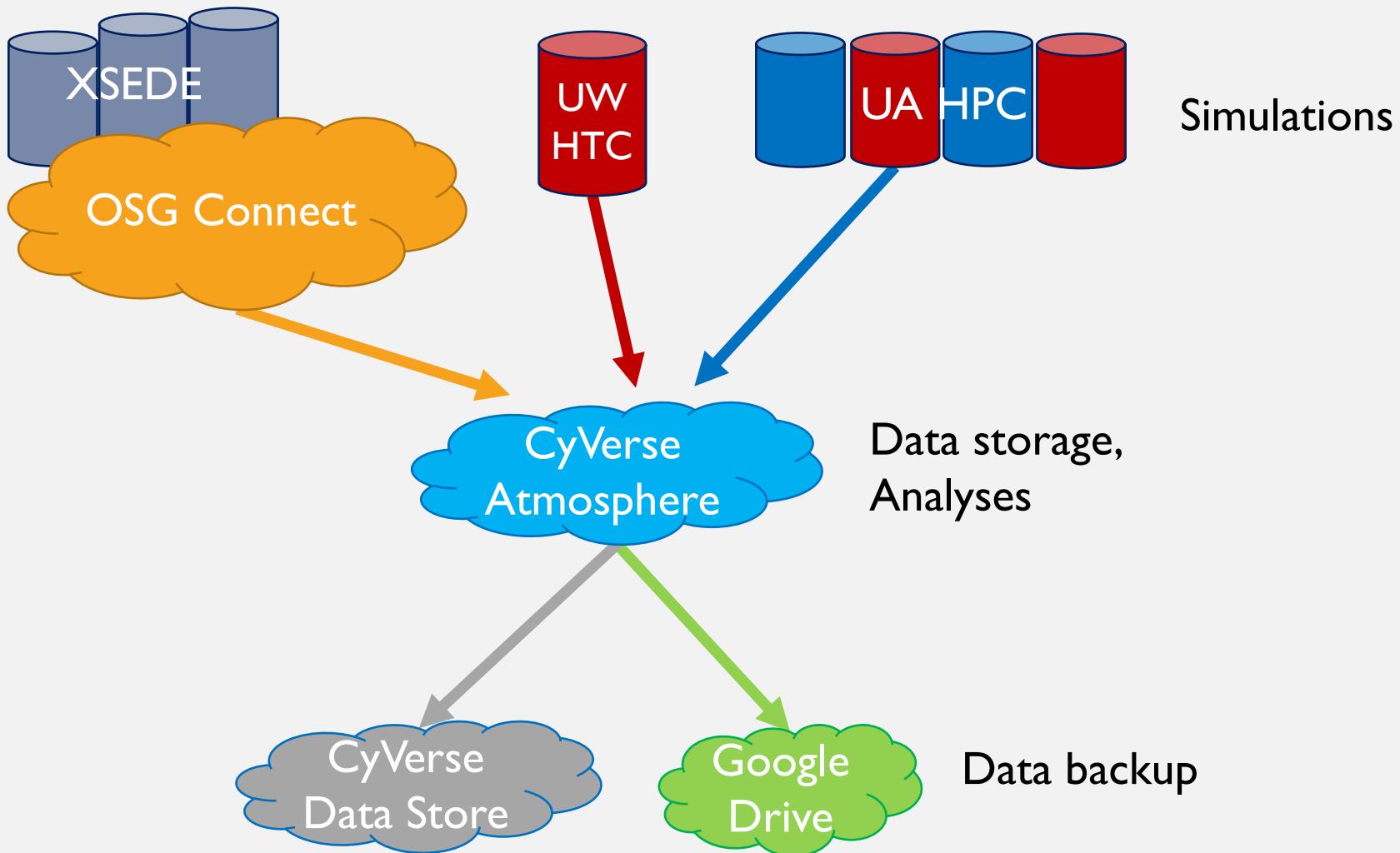
- Each job
  - runs ~40 min, and max 50 hrs
  - Uses ~1G, and max 5G memory
  - Uses ~2M in storage



# HIGH THROUGHPUT COMPUTING



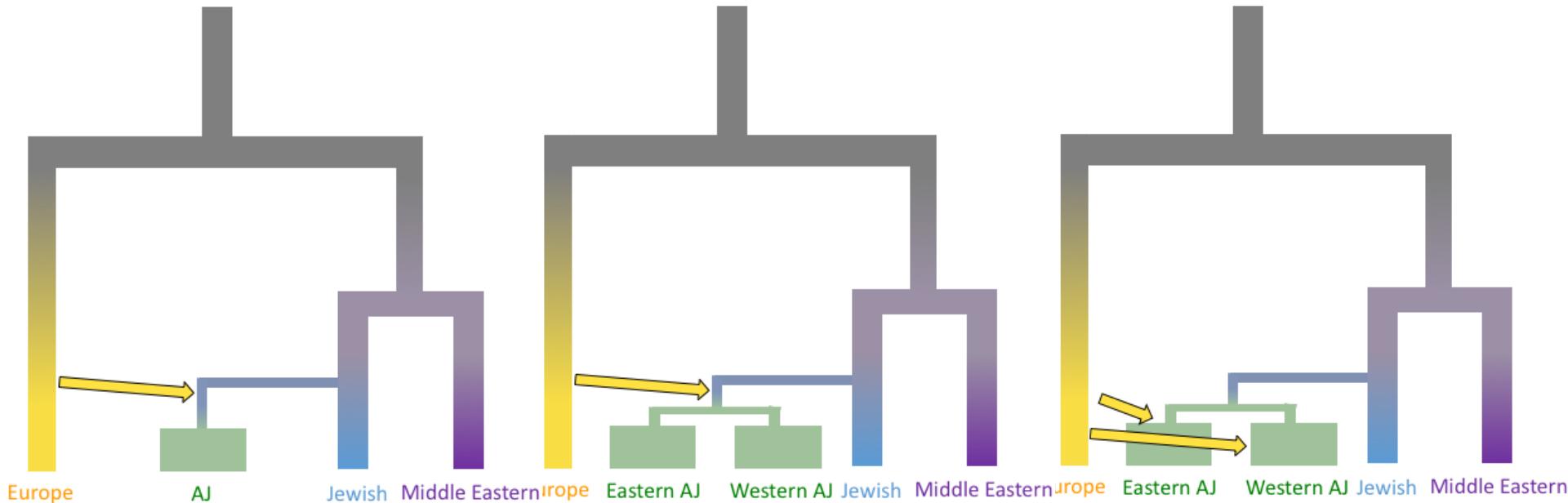
## SIMULATIONS ON HTC CLUSTERS, ANALYSES ON VM



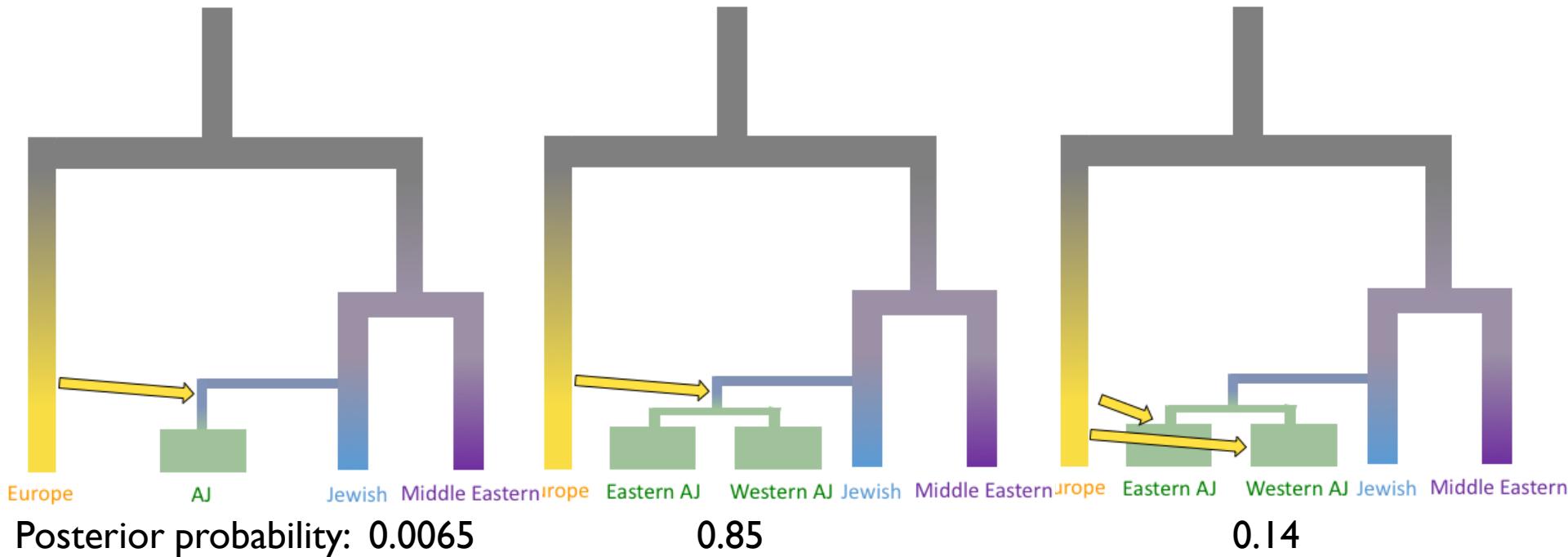
## CHALLENGES: TECHNICAL

- How to handle millions of files?
  - UA HPC has file number limit
  - If there are too many files in a directory simple things take a long time
- How to not overload UA HPC system?
- How to reliably backup data?
- Why do jobs fail?

>1 MILLION SIMULATIONS OF EACH MODEL

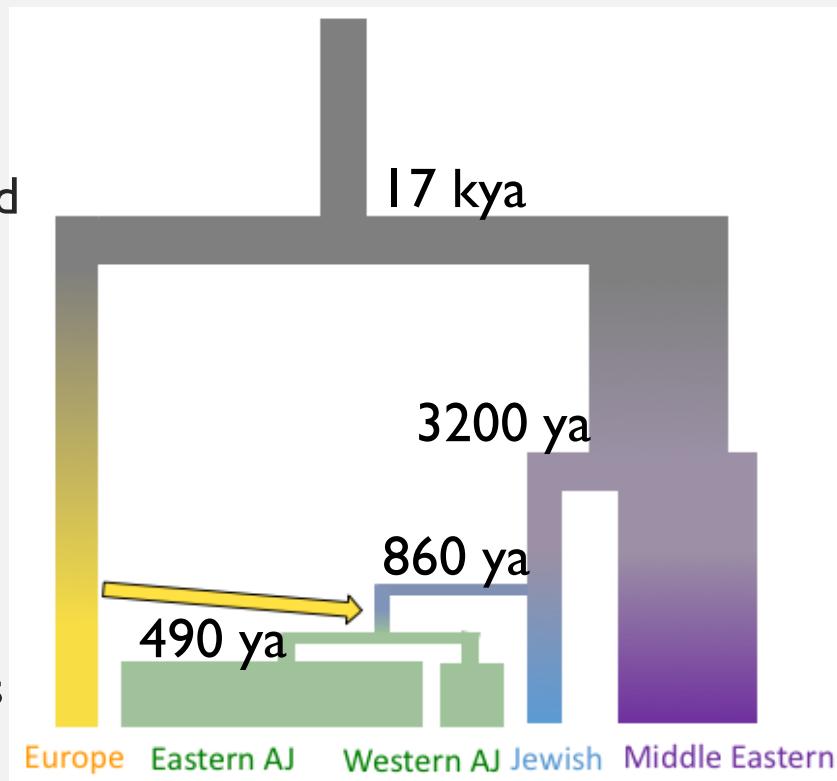


## MODEL CHOICE

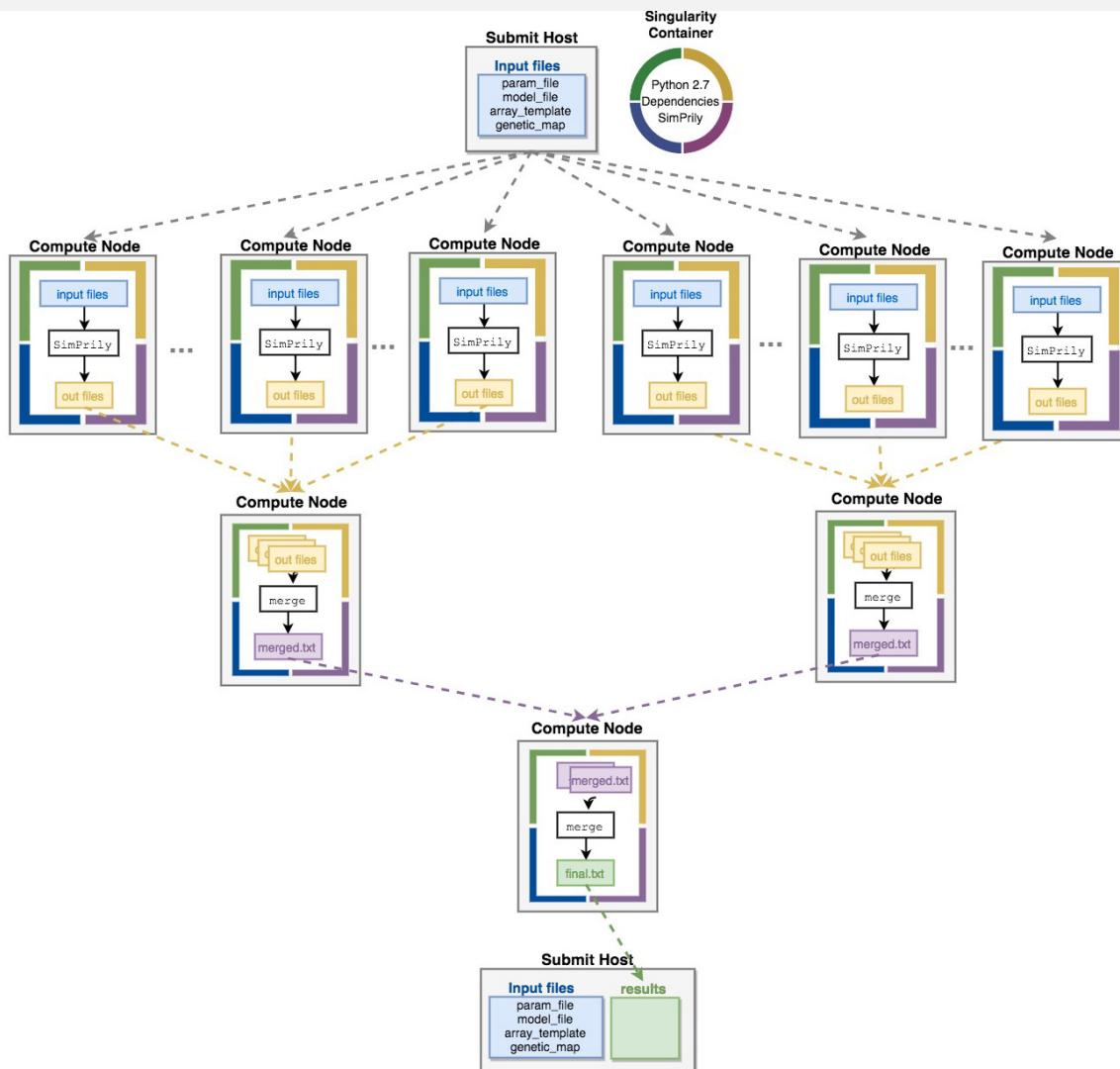


## BEST MODEL

- ~1200 BCE ancestors of Jewish populations diverged from other Middle Eastern populations
  - Experienced extreme population size reduction
- ~1100 CE ancestors of Ashkenazi Jews diverged from other Jewish populations
  - Experienced another population size reduction
  - Experienced gene flow from Europeans (unresolved how much or when)
- ~1500 CE Eastern and Western Ashkenazi Jews diverged
  - Western AJ moderately grew in size
  - Eastern AJ massively grew in size



# SIMPRILY: GENERALIZATION OF CODE AND WORKFLOW



- Developed program to simulate any demographic model
  - Memory & space efficient
- Use Singularity container
- Pegasus workflow for OSG

<https://agladstein.github.io/SimPrily/>

# THANK YOU!

## HAMMER LAB

- Michael Hammer
- Consuelo Quinto-Cortes

## CYVERSE

- Blake Joyce
- Julian Pistorius

## UA HPC CONSULTING

- Mike Bruck
- Dima Shyshlov



THE UNIVERSITY  
OF ARIZONA®



WISCONSIN  
UNIVERSITY OF WISCONSIN-MADISON

## OPEN SCIENCE GRID & PEGASUS

- Mats Rynge
- ## UW CENTER FOR HTC

- Lauren Michael
- Christina Koch

## OPEN SCIENCE GRID USER SCHOOL

- Tim Cartwright
- Lauren Michael
- Christina Koch



Pegasus

## CODING MINIONS

- David Christy
- Logan Gantner
- Mack Skodiak
- Daniel Olson
- Rafael Lopez
- Kayleen Gurrola
- Katie McCready

## RESOURCES PROVIDED BY

- University of Arizona HPC
- University of Wisconsin HTC
- CyVerse
- Open Science Grid
- XSEDE
  - Bridges
  - Comet
  - Jetstream

XSEDE

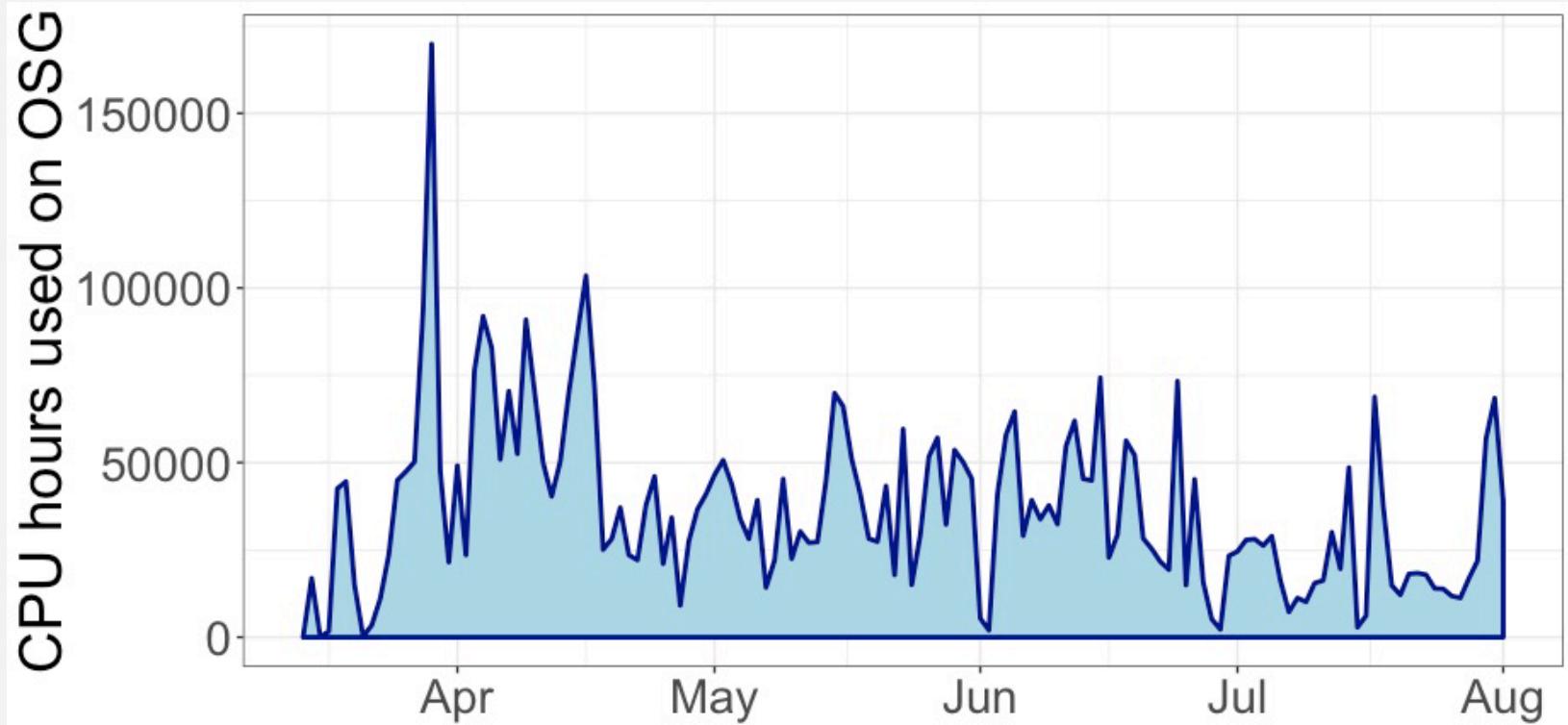
Extreme Science and Engineering  
Discovery Environment



CYVERSE™

osg connect

## CPU HOURS ON THE OPEN SCIENCE GRID



## DNA SEQUENCE

Indiv I

AATCATTTCGGTTTAATGCTTGGGCTGCATTGGGAAA  
AATCATATCGGTCTTAATGCTTGCGCTGCCTTGGTAAA

## DNA SEQUENCE, SEGREGATING SITES

Indiv I

AATCATTTCGGTTTAATGCTTGGGCTGCATTGGGAAA  
AATCATATCGGTCTTAATGCTTGCCTGCTGCCTTGGTAAA

## DNA SEQUENCE, SEGREGATING SITES

Indiv 1

AATCATTTCGGTTTAATGCTTGGGCTGCATTGGGAAA  
AATCATATCGGTCTTAATGCTTGCCTGCGCTGCCTTGGTAAA

Indiv 2

AATCATTTCGGTTTAATGCTTGGGCTGCCTTGGTAAA  
AAACATTTCGGTCTTATGGTTGCCTGCGCTGCATTGGGGAA

## DNA SEQUENCE, GENOTYPES ENCODED 0/1

# Indiv I

AATCATTTCGGTTTAATGCTTGGGCTGCATTGGGAAA  
AATCATATCGGTCTTAATGCTTGGCTGCCTTGGTAAA

# Indiv I

000000000000000000000000000010000000000001000  
0000001000000100000000000000000000001000000000

## Indiv 2

AATCATTTCGGTTAATGCTTGGGCTGCCTTGGTA  
AAACATTTCCGTCTTATGGTTGCGCTGCATTGGGGAA

## Indiv 2

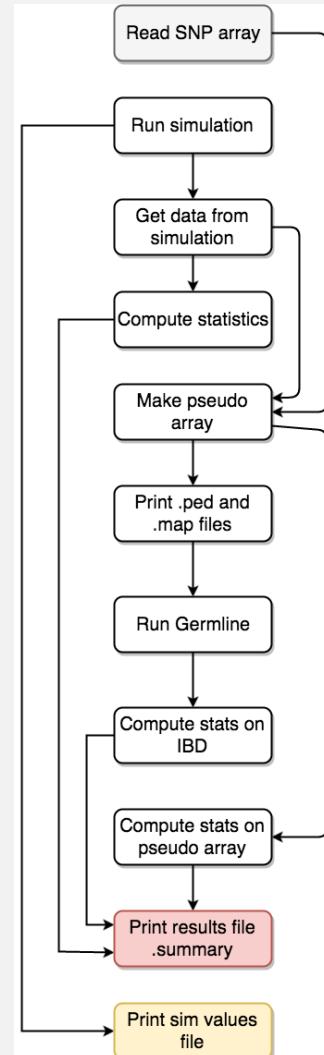
## SEQUENCE OF GENOTYPES, ONLY SEGREGATING SITES

**Indiv 1** 0000001010  
0101000100

**Indiv 2** 0000000100  
1011111011

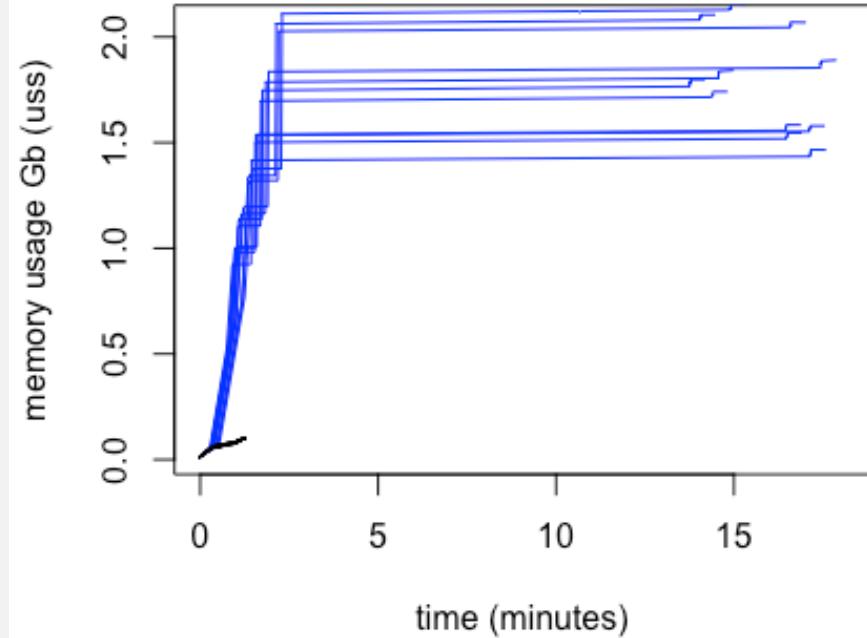
# PYTHON SCRIPT: GENOME SIMULATIONS AND COMPUTE SUMMARY STATISTICS

- Inherited from lab mates
- Intended for millions of relatively small simulations
  - 1,389 10kb regions
  - 65 individuals
- Originally took a few minutes to run
- Originally ran parallel on U of A HPC
  - 1 million runs would take approximately 1 month.

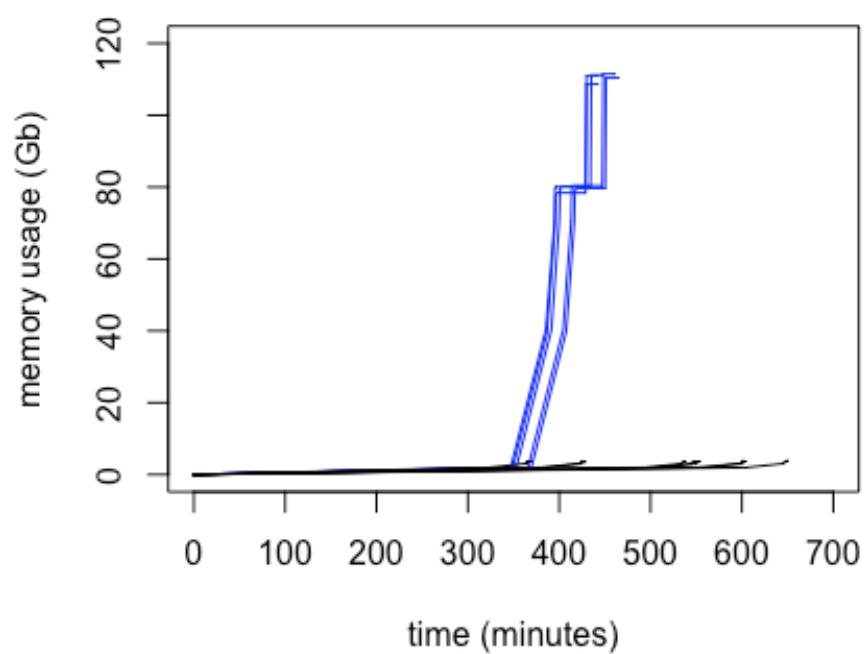


# PROFILE OF PYTHON SCRIPT

Minimum Simulation Parameters



Maximum Simulation Parameters



\*Note different scales

