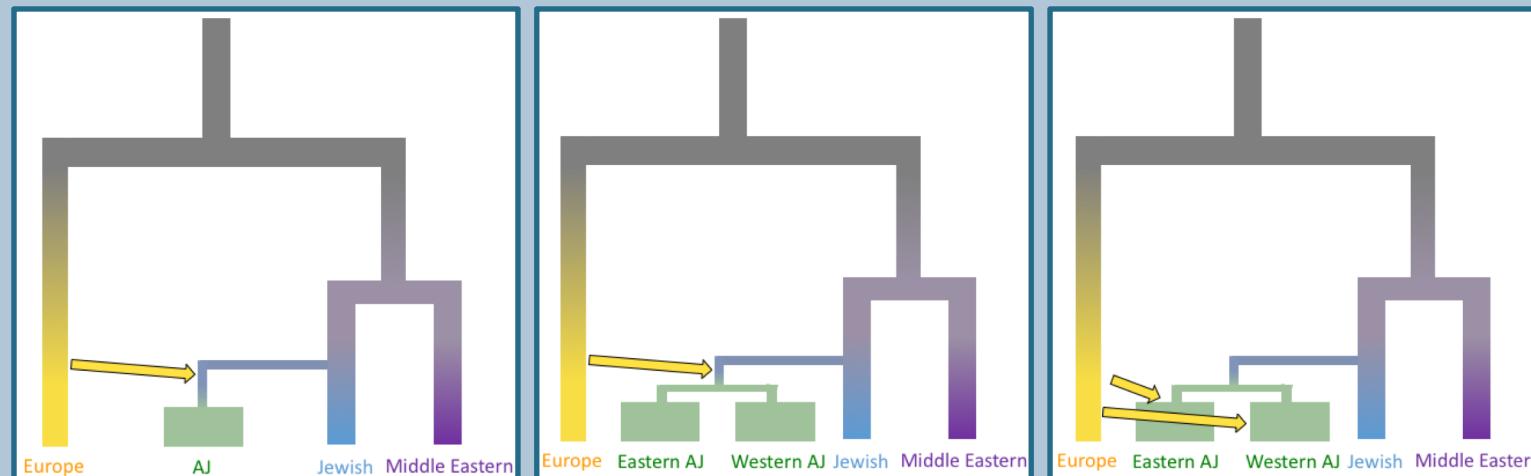


INFERENCE OF RECENT DEMOGRAPHIC HISTORY OF A POPULATION ISOLATE USING SNP ARRAY AND WHOLE GENOME DATA



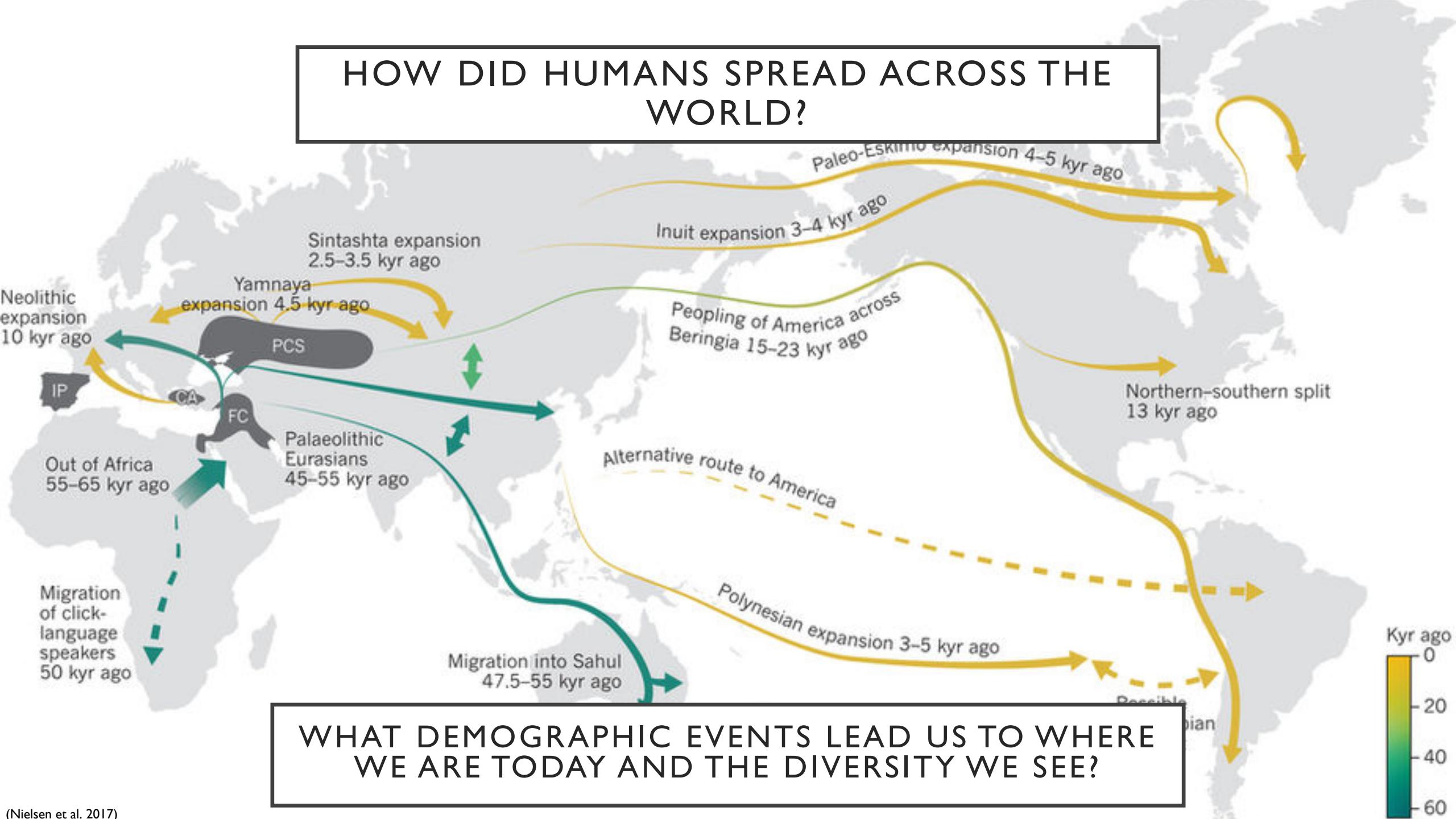
Ariella Gladstein

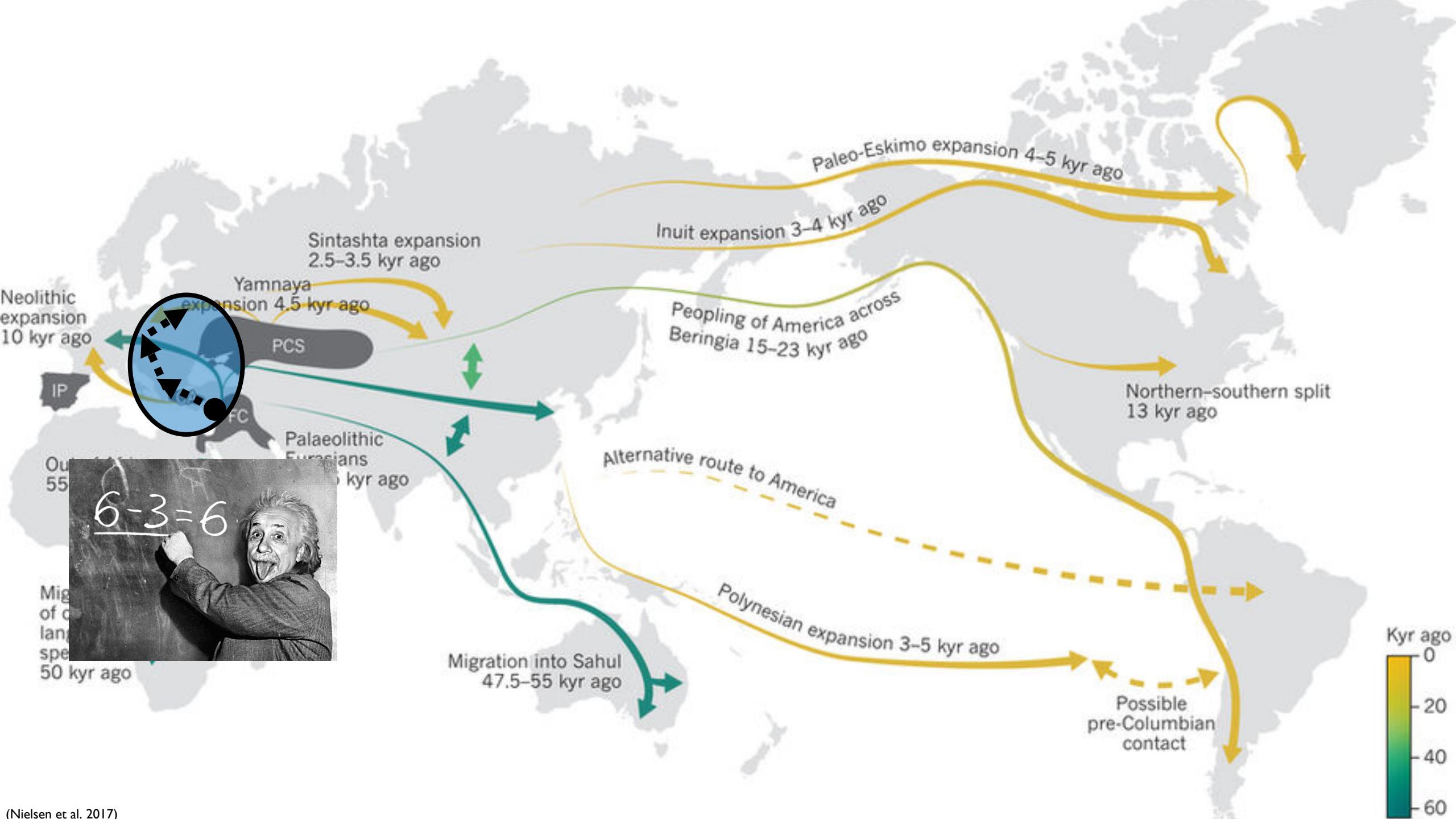
Ecology and Evolutionary Biology





HOW DID HUMANS SPREAD ACROSS THE WORLD?

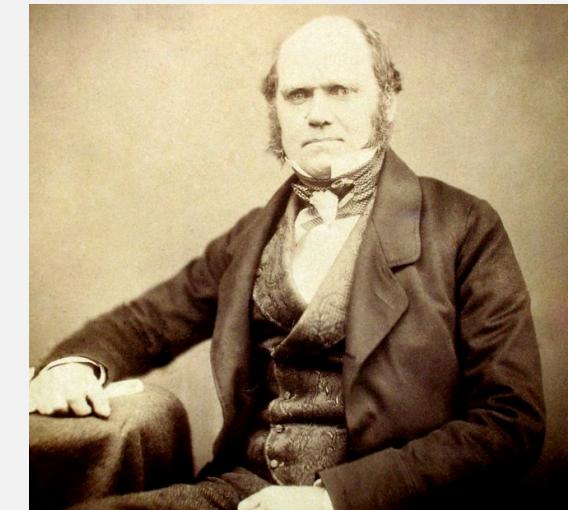
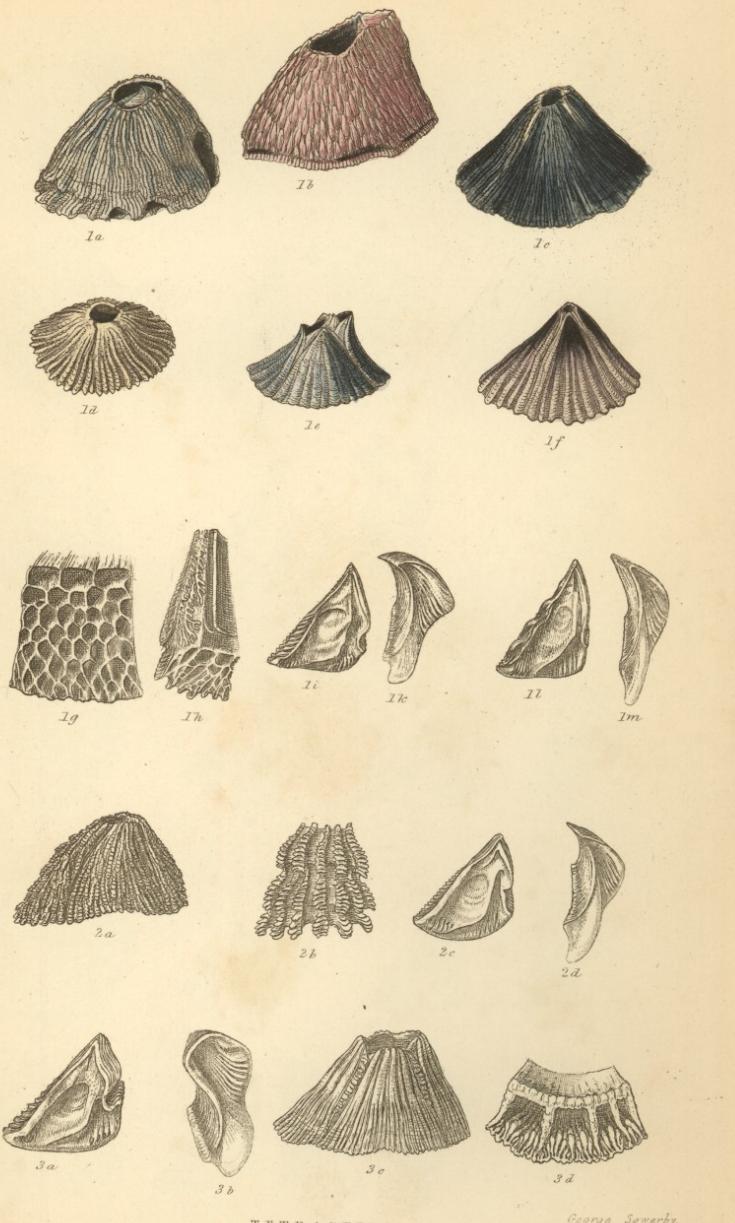




Pl. XX.

BASIC RESEARCH IS IMPORTANT!

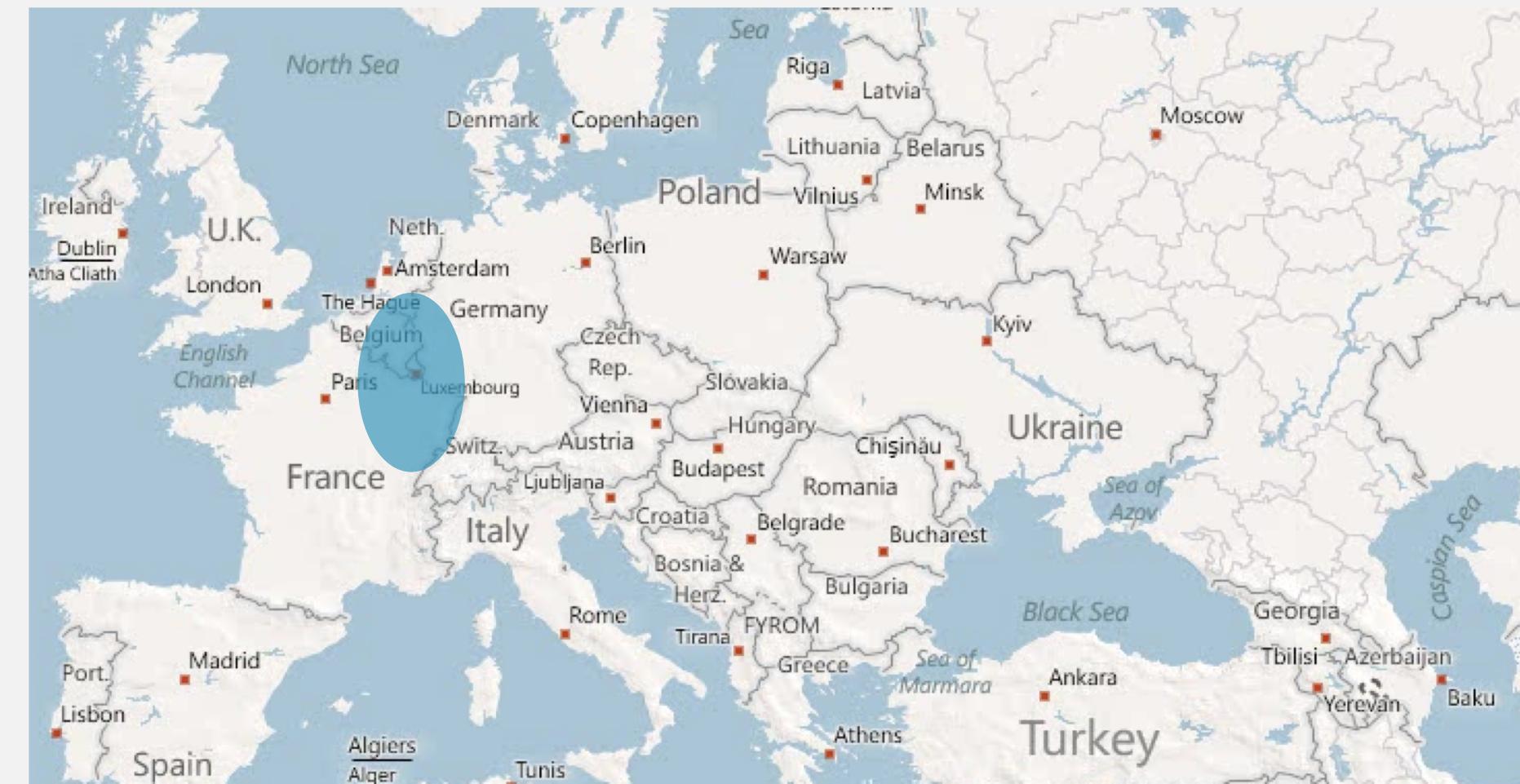
- My research focuses on a small part of basic evolutionary biology questions.
- Huge computational resources and modern techniques to contribute to basic evolution questions



THEMES OF DISSERTATION

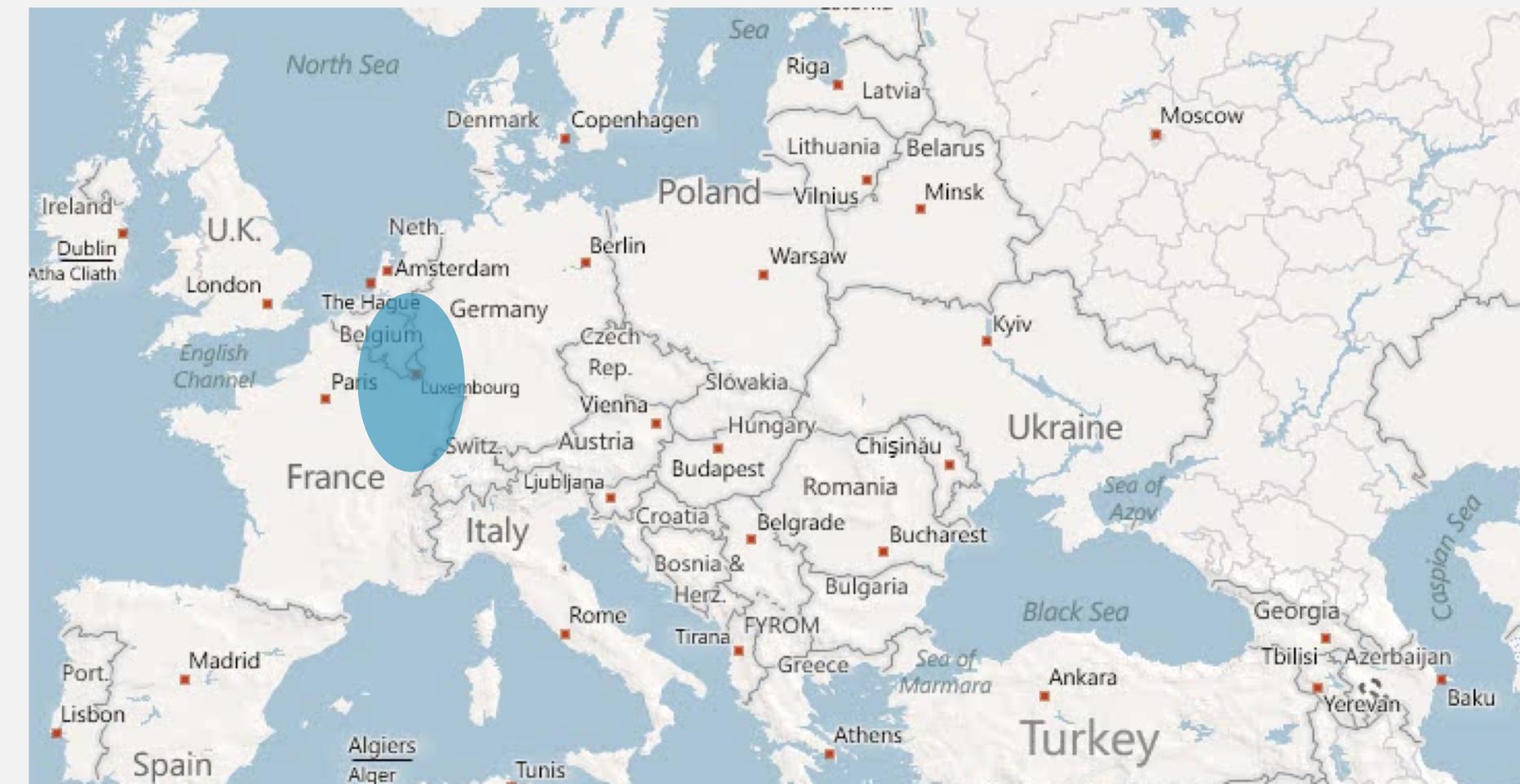
- Detection of runs of homozygosity from SNP arrays
 - Improving identification of runs of homozygosity (Ch. 2)
 - Correcting ascertainment bias in runs of homozygosity (App. C)
- Scaling up Approximate Bayesian Computation for whole chromosomes
 - Create efficient pipeline to simulate demographic models and calculate summary statistics (App. A)
 - Create generalized high throughput workflow (Ch. 4)
- Infer history of the Ashkenazi Jews
 - Substructure in AJ? (Ch. 5)
 - Khazarian origin? (App. B)

WHO ARE THE ASHKENAZI JEWS?



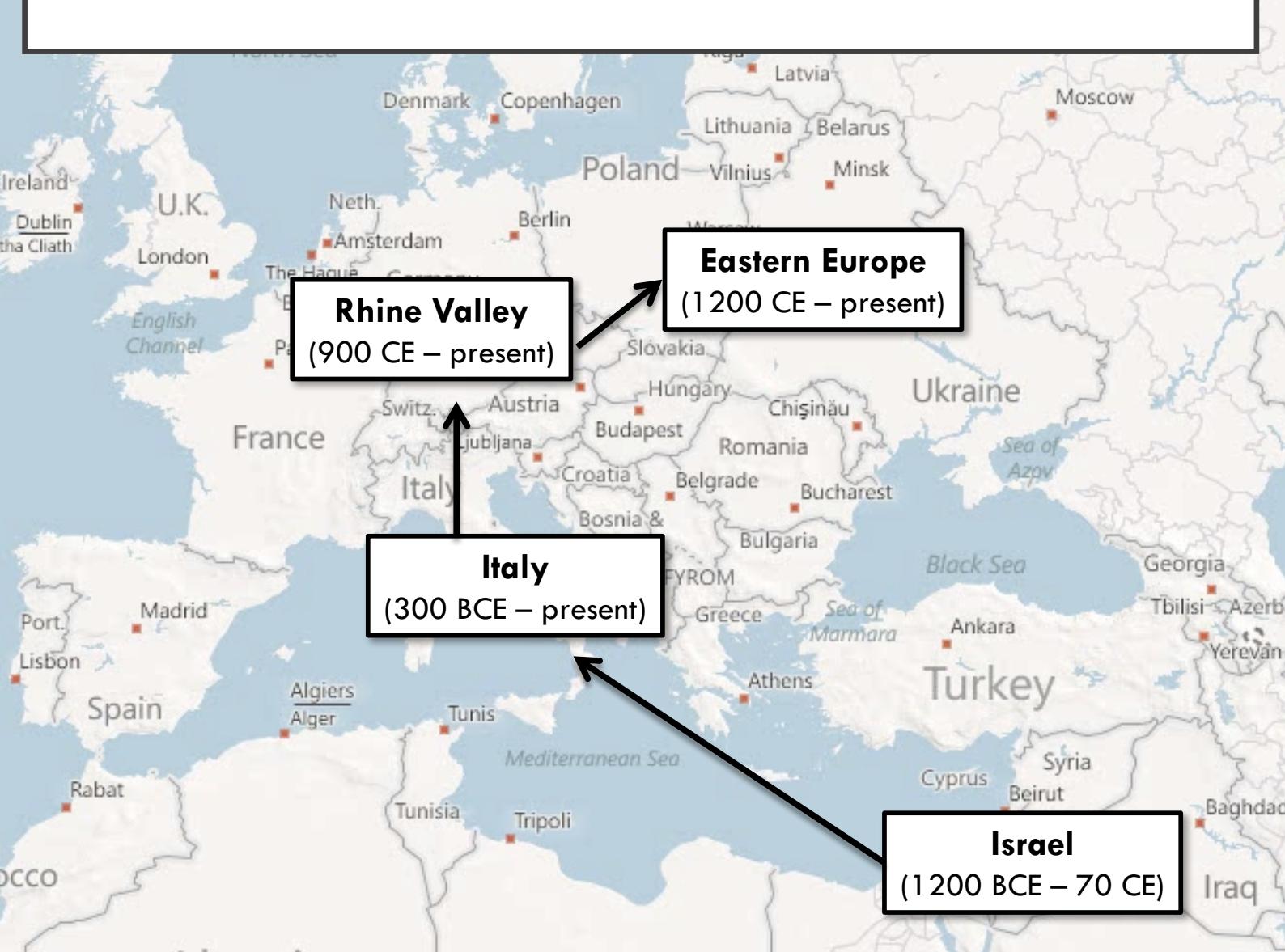
Culturally, religiously, and linguistically identify as Jews whose ancestors came from the Rhine Valley.

ASHKENAZI JEWS: AN INTERESTING STUDY POPULATION

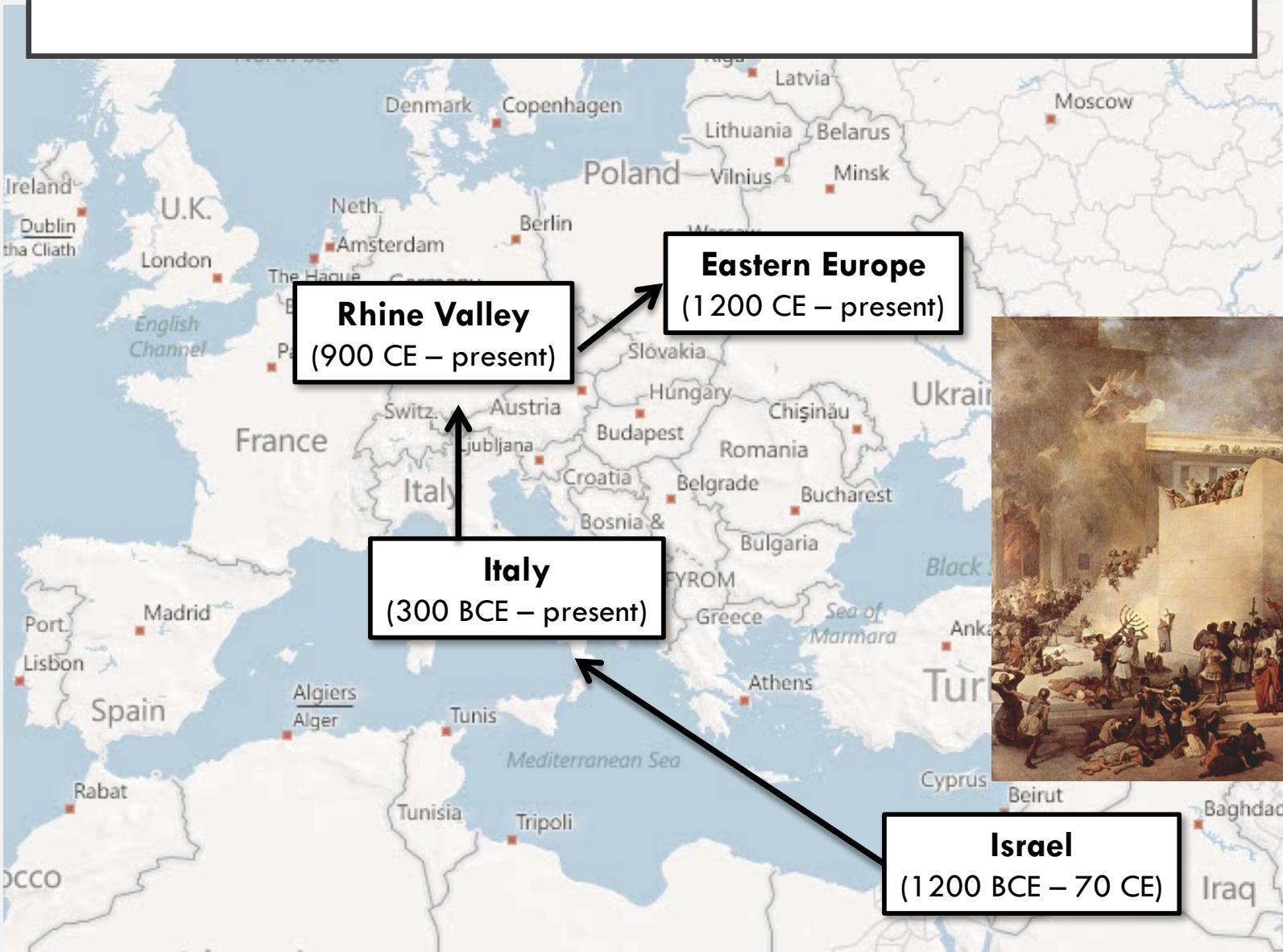


- High frequency of genetic disorders
- Population isolate
- Complex demographic history
- Well documented historical record

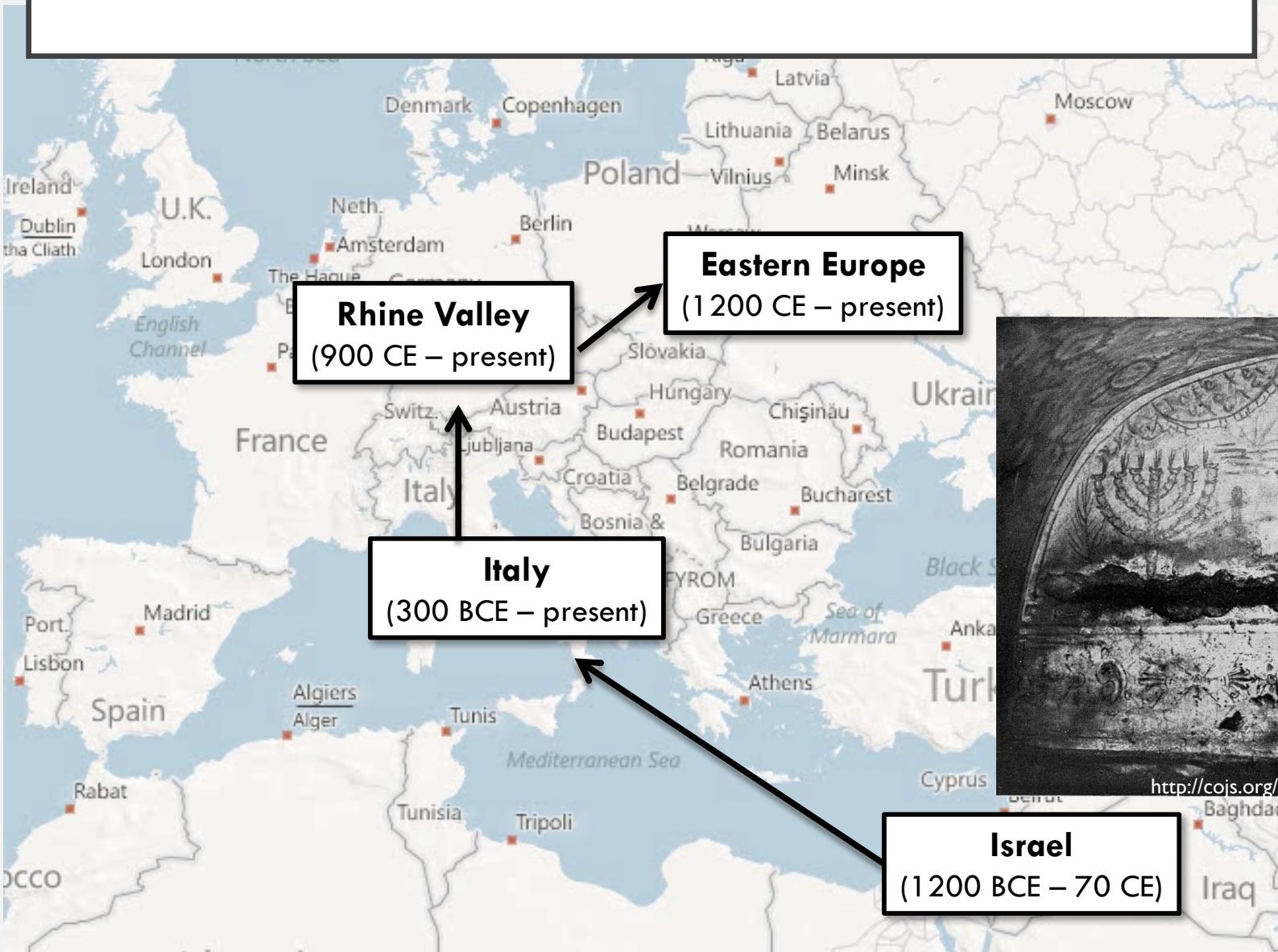
HYPOTHESIS OF ASHKENAZI ORIGINS



HYPOTHESIS OF ASHKENAZI ORIGINS

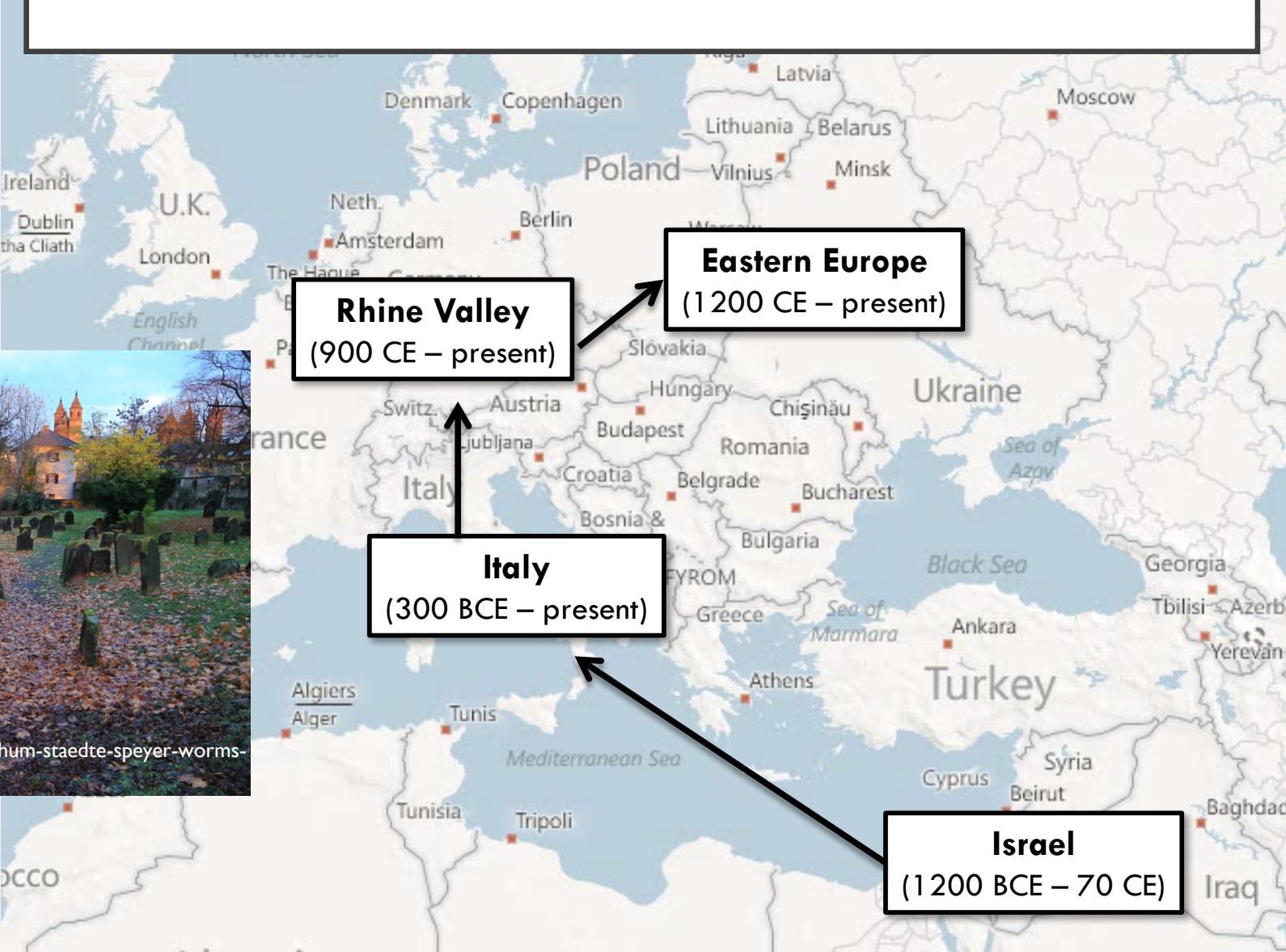


HYPOTHESIS OF ASHKENAZI ORIGINS

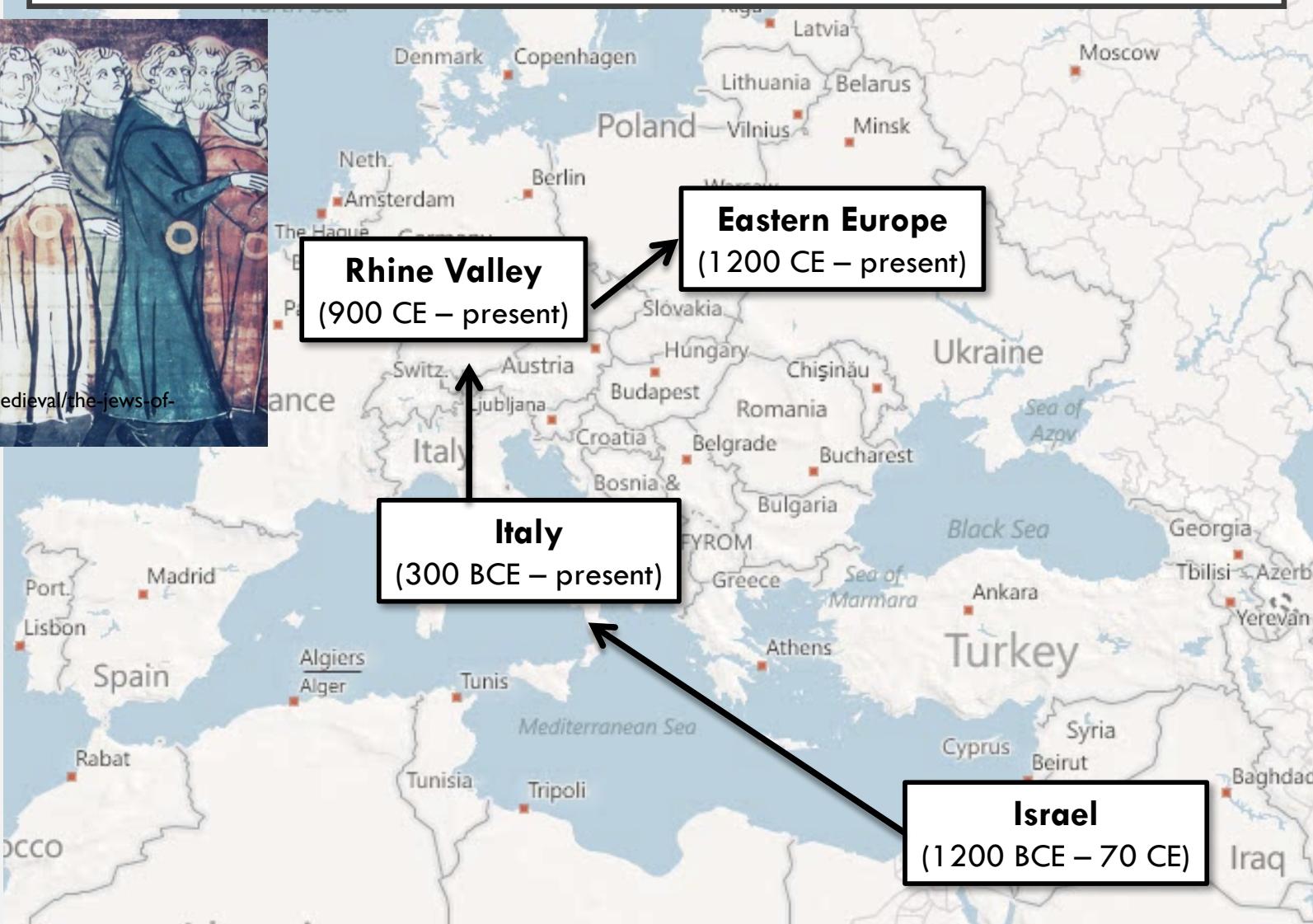


<http://cojs.org/wp-content/uploads/2015/11/Torlonia.jpg>

HYPOTHESIS OF ASHKENAZI ORIGINS



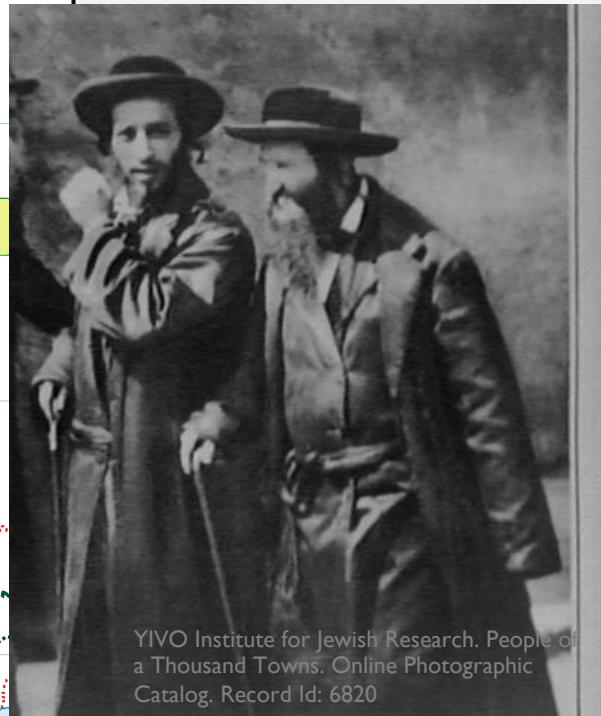
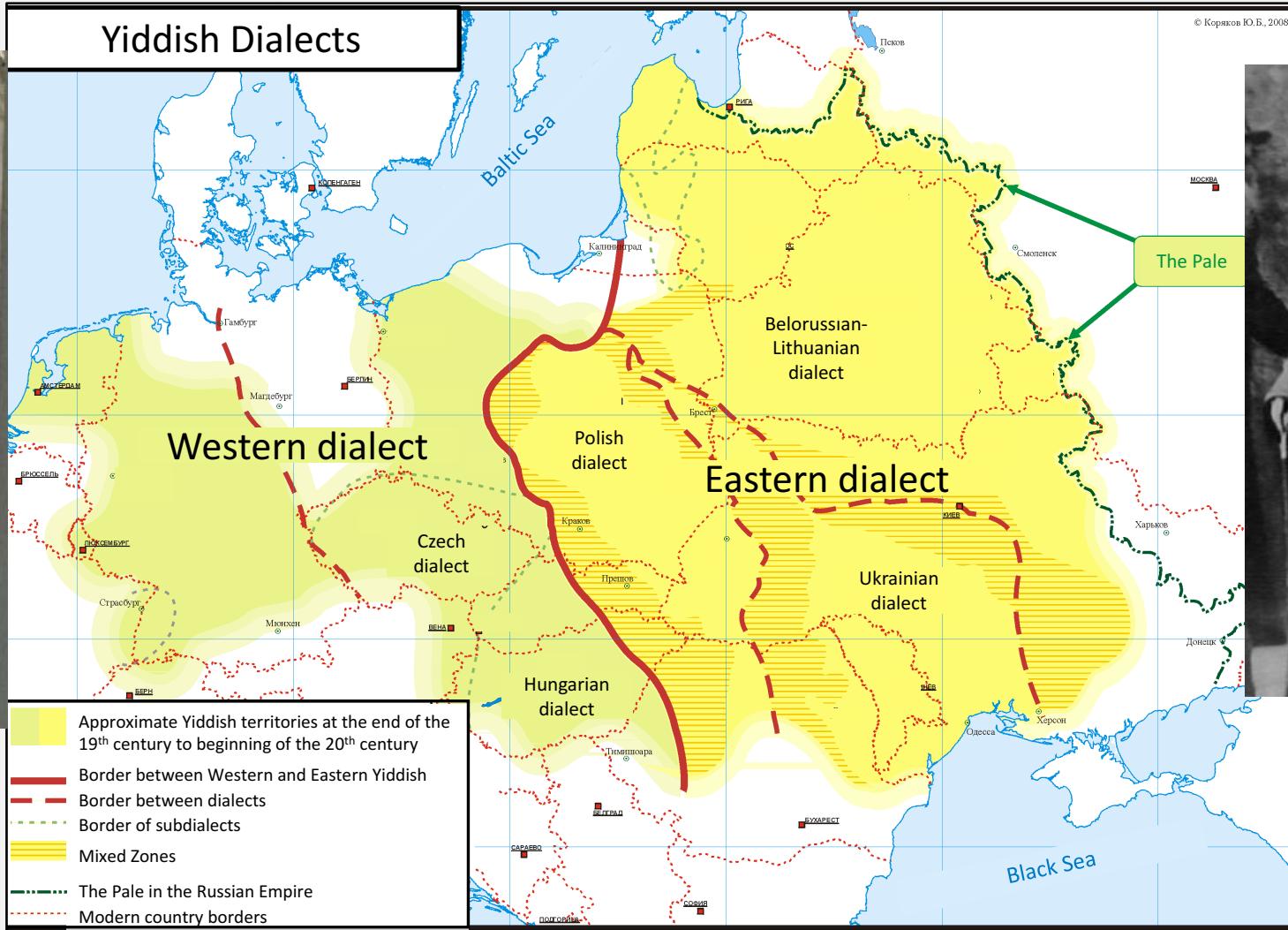
HYPOTHESIS OF ASHKENAZI ORIGINS



WESTERN VS. EASTERN ASHKENAZI JEWS

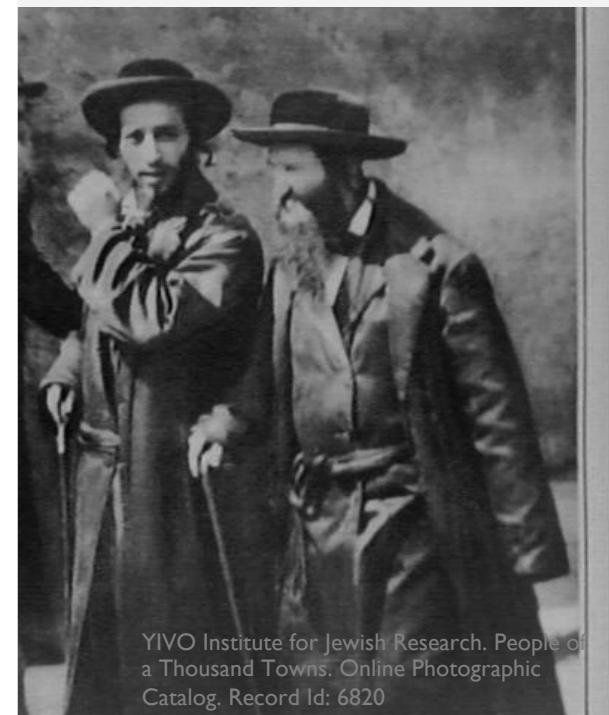
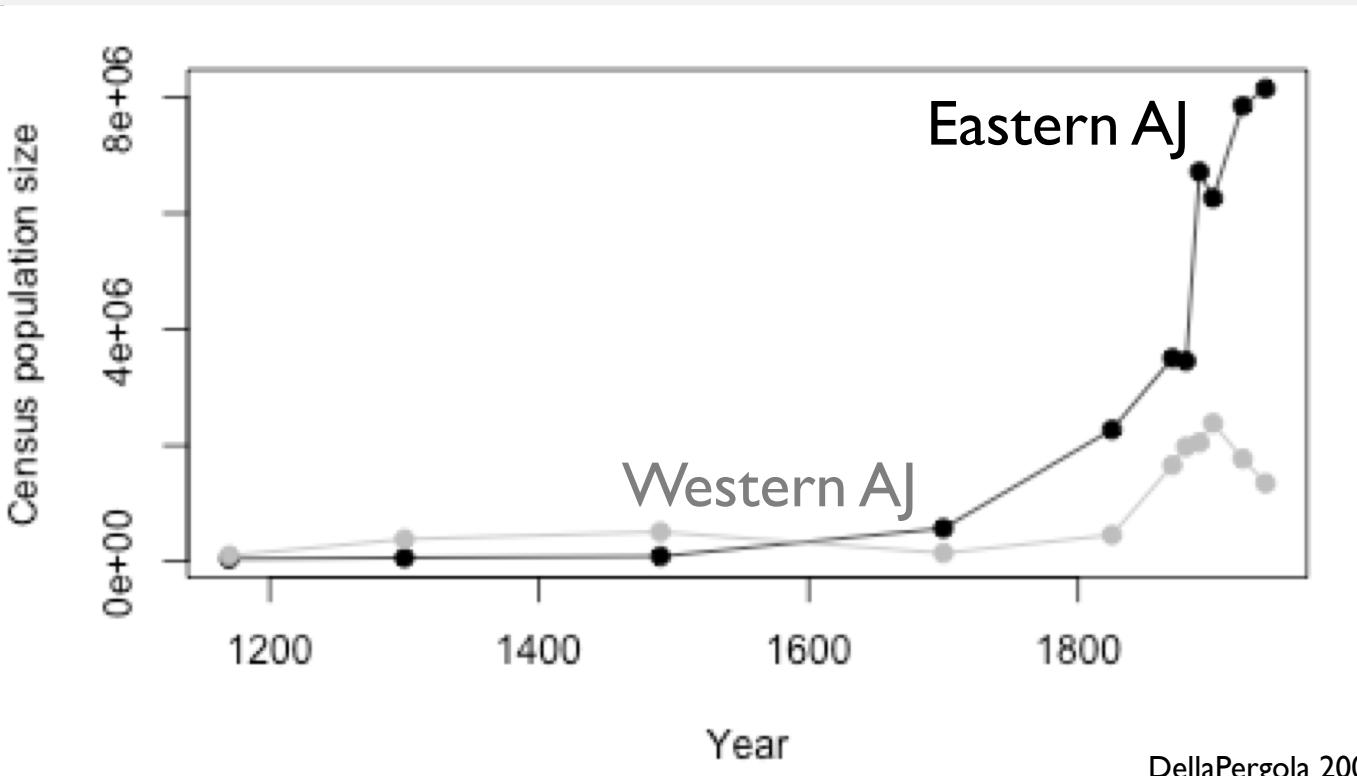


Germany, 1900's



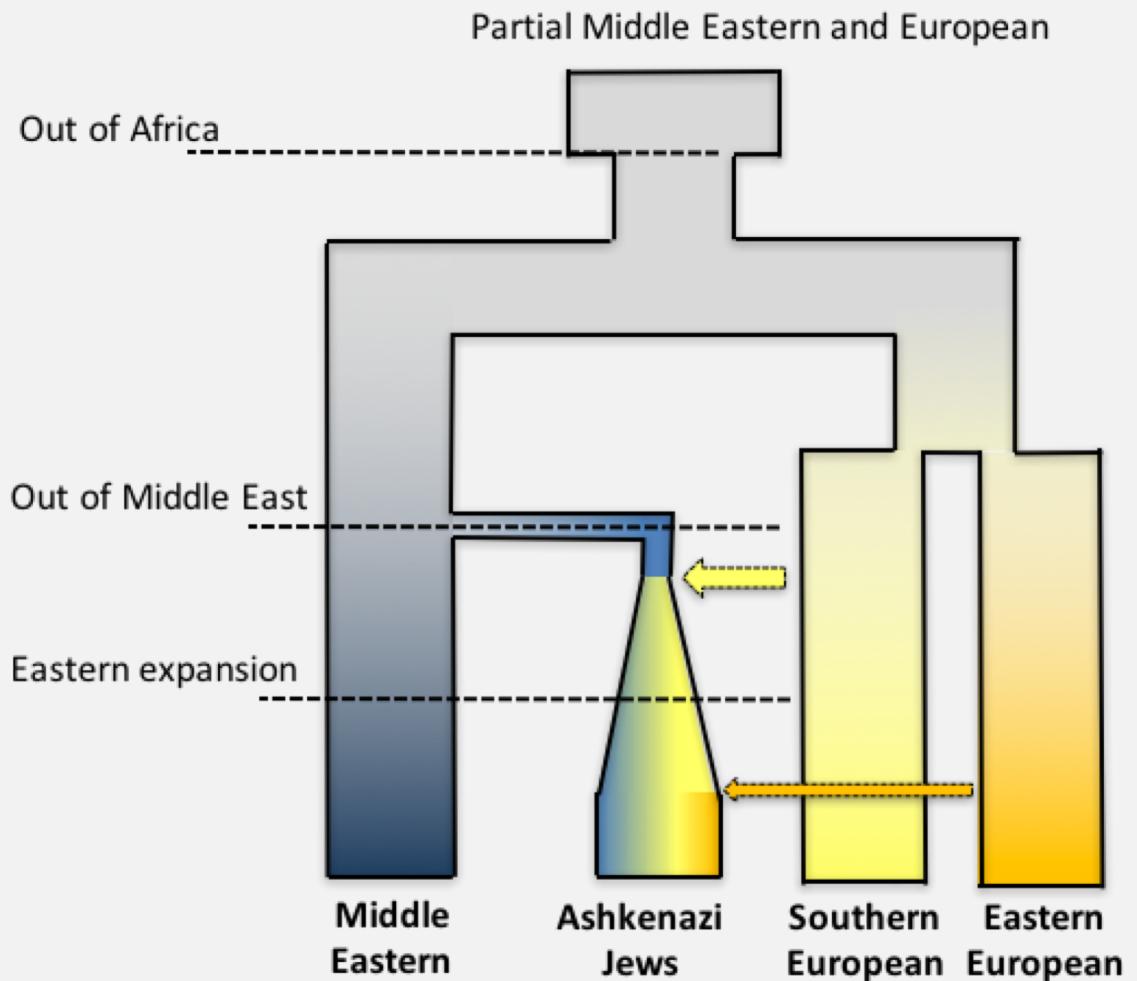
Cracow, Poland. 1932

WESTERN VS. EASTERN ASHKENAZI JEWS



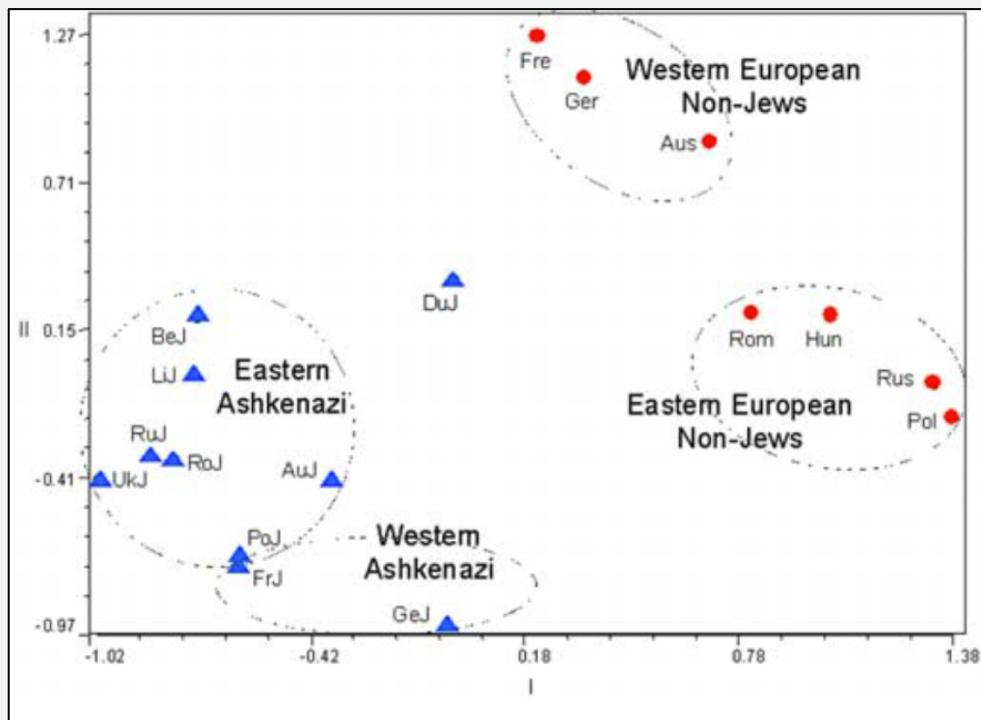
PREVIOUS STUDIES

- Most closely related to other Jewish populations
- Middle Eastern origin
- Bottlenecks
 - Population reduction ~30 gen ago
 - Exponential growth ~30 gen ago
- At least 2 admixture events
 - Southern Europeans ~35-60%, ~25-50 gen ago
 - Eastern Europeans ~15-25%, ~10-20 gen ago

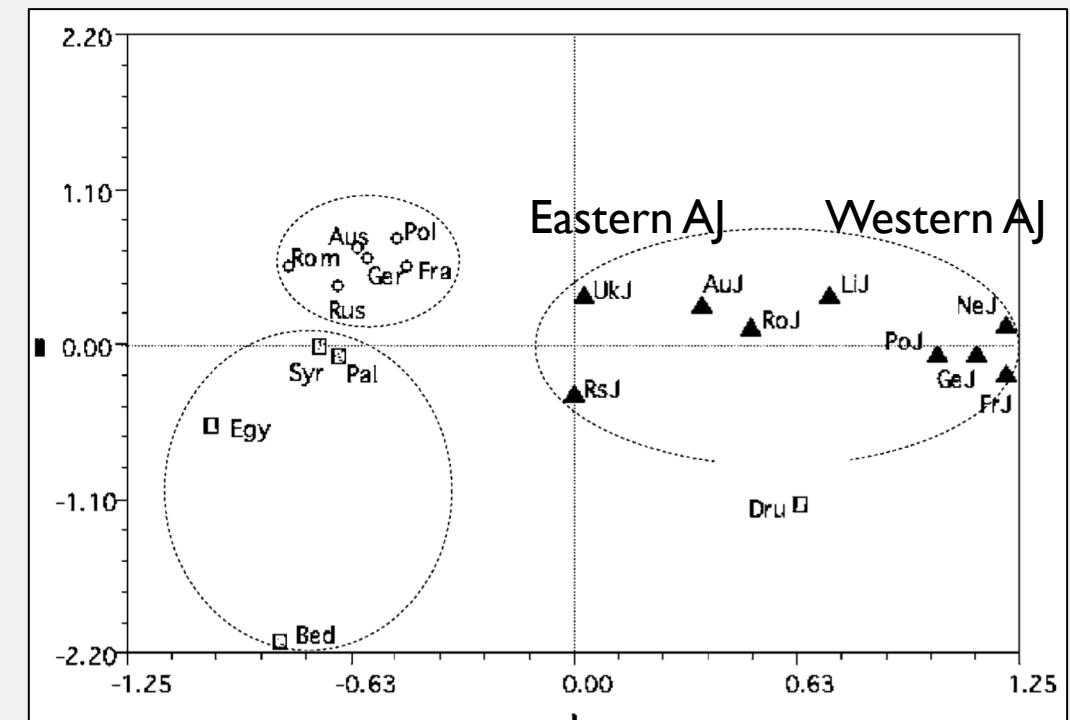


PREVIOUS STUDIES: AJ SUBSTRUCTURE

Y chromosome and mtDNA markers show differences among AJ from different countries



Behar et al. 2004a



Behar et al. 2004b

PREVIOUS STUDIES: AJ SUBSTRUCTURE

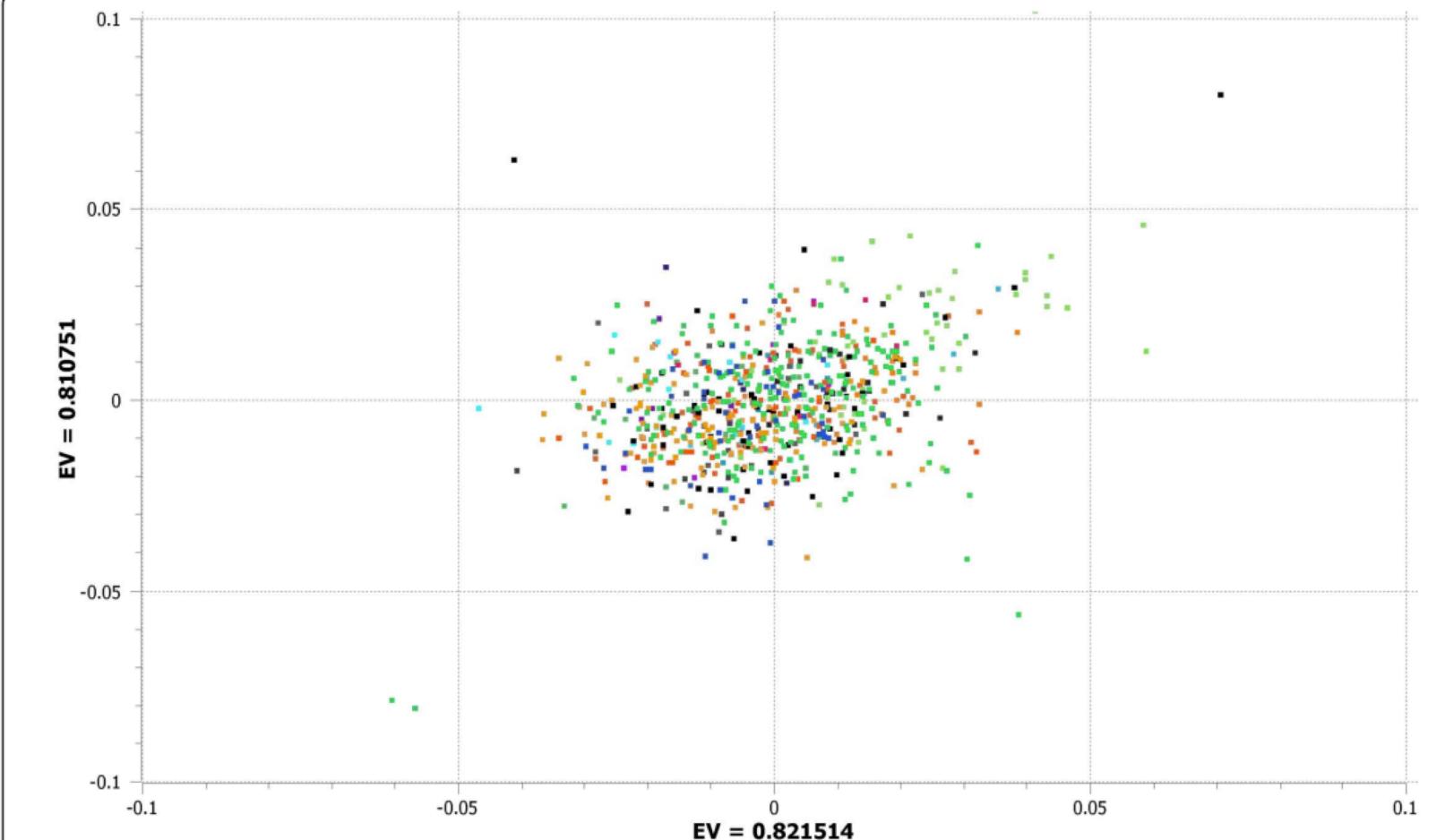
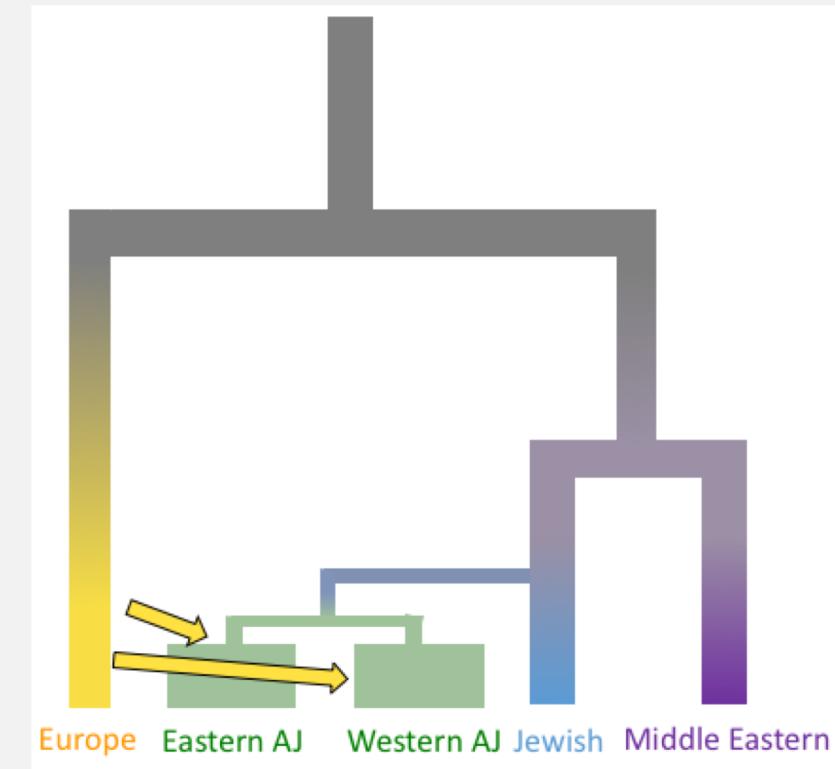
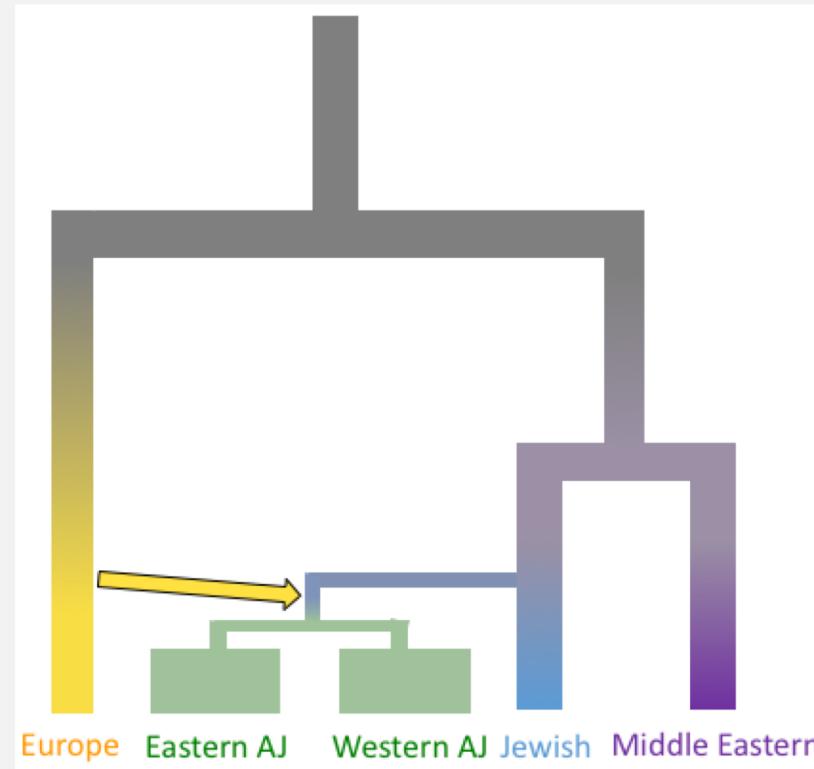
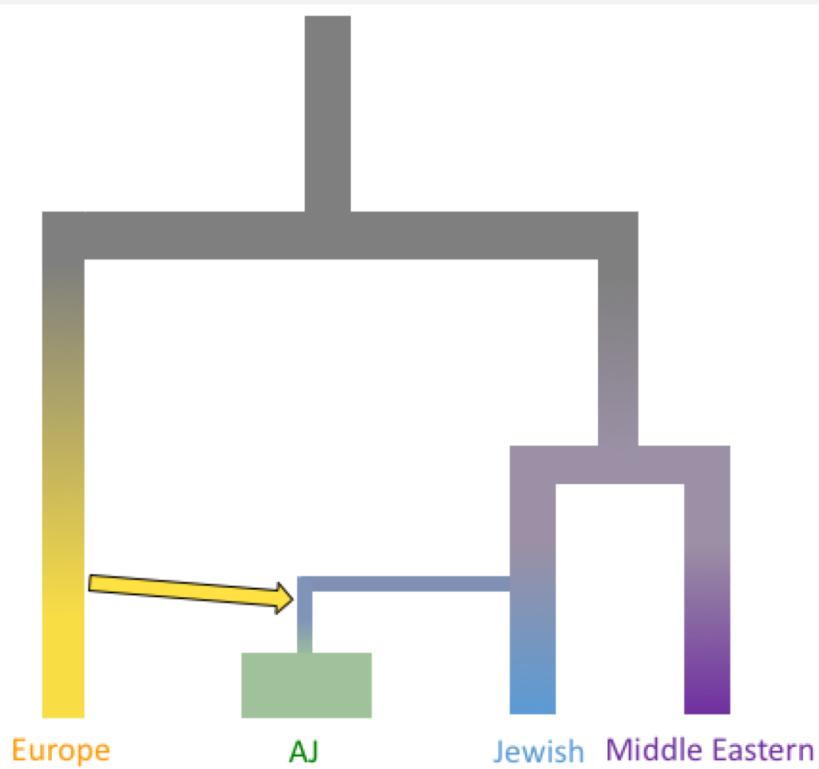


Figure 6 Intra-population principal components analysis of 1,312 Ashkenazi Jewish (AJ) individuals with cluster 3 (C3) scores > 0.6 derived from ADMIXTURE analysis. The x-axis represents the eigenvalue (EV) for principal component 1 (PC1) and the y-axis represents the eigenvalue for principal component 2 (PC2). Different colors represent different geographical origins of AJ individuals.

No intra-population structure found with SNP array data



MODELS OF ASHKENAZI HISTORY

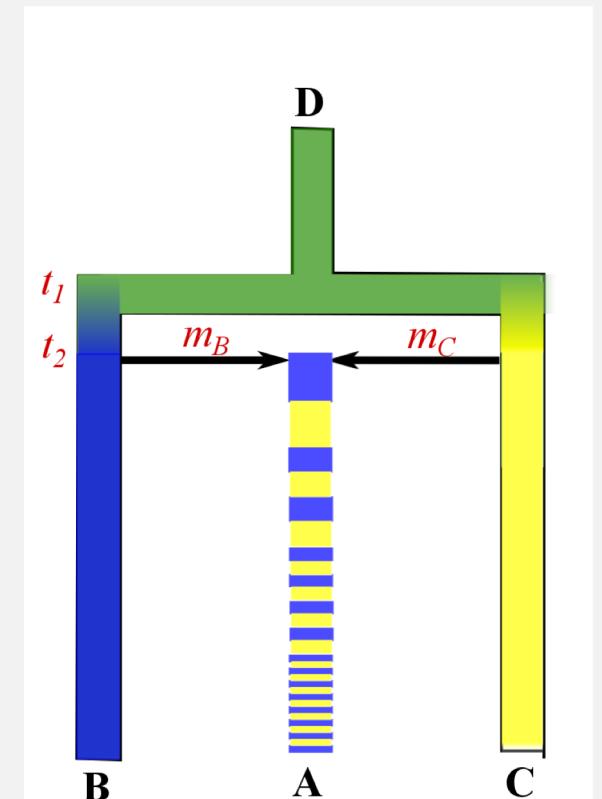
THEMES OF DISSERTATION

- Detection of runs of homozygosity from SNP arrays
 - Improving identification of runs of homozygosity (Ch. 2)
 - Correcting ascertainment bias in runs of homozygosity (App. C)
- Scaling up Approximate Bayesian Computation for whole chromosomes
 - Create efficient pipeline to simulate demographic models and calculate summary statistics (App. A)
 - Create generalized high throughput workflow (Ch. 4)
- Infer history of the Ashkenazi Jews
 - Substructure in AJ? (Ch. 5)
 - Khazarian origin? (App. B)

HAPLOTYPES CAN BE USED TO INFER HISTORY

- Over time, recombination breaks up segments of the genome in predictable ways

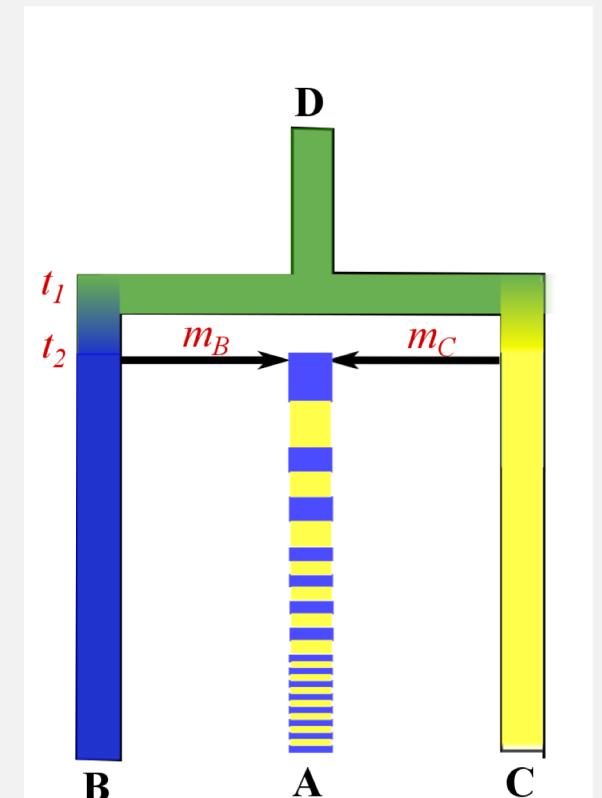
- Segments identical by descent (IBD)
gene flow, effective size, relatedness
- Runs of homozygosity (ROH)
effective size, relatedness, random mating
- Linkage disequilibrium blocks (LD)
gene flow, effective size
- Ancestry blocks
gene flow



HAPLOTYPES CAN BE USED TO INFER HISTORY

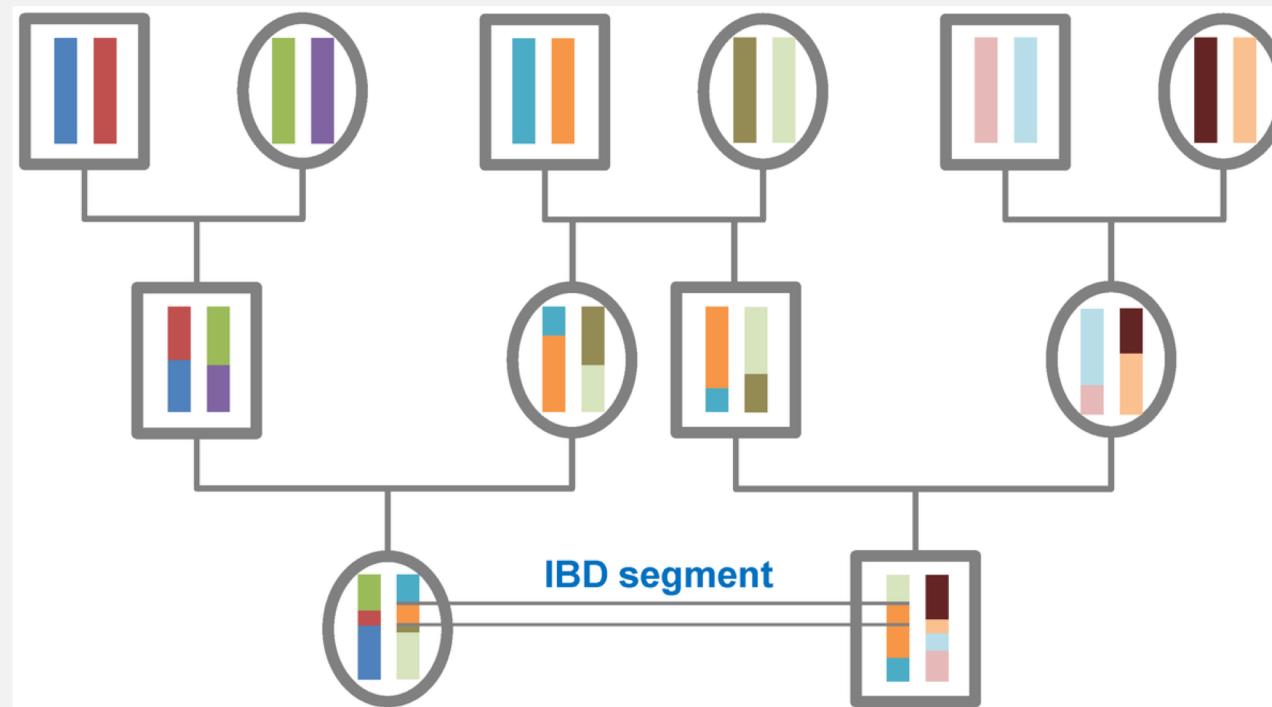
- Over time, recombination breaks up segments of the genome in predictable ways

- Segments identical by descent (IBD)
gene flow, effective size, relatedness
- Runs of homozygosity (ROH)
effective size, relatedness, random mating
- Linkage disequilibrium blocks (LD)
gene flow, effective size
- Ancestry blocks
gene flow



SEGMENTS IDENTICAL BY DESCENT (IBD)

A region on two chromosomes that was inherited from a common ancestor



- Every site is (technically) IBD
- Practically, we define IBD based on a minimum length

RUNS OF HOMOZYGOSITY (ROH)

- A ROH is a genomic segment of continuous homozygous sites.
- ROH are defined based on a minimum length



ROH on a chromosome. Each dotted grey line is a homozygous site and each red line is a heterozygous site. Each shaded grey area is a ROH

RUNS OF HOMOZYGOSITY (ROH)

- A ROH is a genomic segment of continuous homozygous sites.
- ROH are defined based on a minimum length



ROH on a chromosome. Each dotted grey line is a homozygous site and each red line is a heterozygous site. Each shaded grey area is a ROH

ROH reflect relatedness of ancestors
Smaller Ne increases likelihood of creating ROH

SINGLE NUCLEOTIDE POLYMORPHISM (SNP) ARRAYS

- Genome-wide and many SNPs (100 K's – millions)
- Benefits:
 - Inexpensive
 - Low genotyping error rates
 - Easy to work with
- Disadvantages:
 - Ascertainment bias – reduction of represented genetic diversity



HOW DOES BIAS FROM SNP ARRAYS AFFECT HAPLOTYPE STATISTICS?

- Extensive work on the effect of ascertainment bias on the allele frequency spectrum.
- Haplotype statistics considered to be less sensitive to ascertainment bias.

PLINK 1.9 home

plink2-users

GitHub

File formats

PLINK 1.9 index

PLINK 2.0

Introduction, downloads

S: 28 May 2018 (b6.1)

D: 28 May 2018

PLINK 1.90 beta

Google Scholar

"runs of homozygosity"



Articles

About 540 results (0.08 sec)

Any time

Since 2018

Since 2017

[PLINK: a tool set for whole-genome association and population-based linkage analyses](#)

Search within citing articles

Introduction, downloads

S: 28 May 2018 (b6.1)

D: 28 May 2018

Binary file distribution

--indep...

--r/-r2

--show-tags

--blocks

Distance matrices

Identity-by-state/Hamming

(--distance...)

Relationship/covariance

(--make-grm-bin...)

--rel-cutoff

Distance-pheno. analysis

(--ibs-test...)

Identity-by-descent

--genome

--homozyg...

Population stratification

--cluster

--pca

--mds-plot

--neighbour

Association analysis

Basic case/control

(--assoc, --model)

Stratified case/control

PLINK 1.90 beta

Runs of homozygosity

```
--homozyg <group | group-verbose> <consensus-match> <extend> <subtract-1-from-lengths>
--homozyg-snp [min SNP count]
--homozyg-kb [min length]
--homozyg-density [max inverse density (kb/SNP) ]
--homozyg-gap [max internal gap kb length]
--homozyg-het [max hets]
--homozyg-window-snp [scanning window size]
--homozyg-window-het [max hets in scanning window hit]
--homozyg-window-missing [max missing calls in scanning window hit]
--homozyg-window-threshold [min scanning window hit rate]
```

If any of these flags are present, a [set of run-of-homozygosity reports](#) is generated using PLINK 1.07's scanning algorithm. See the [original documentation](#) for more details.

Colorectal cancer risk is not associated with increased levels of homozygosity in Saudi Arabia

Abdul K. Siraj, PhD¹, Hanif G. Khalak, PhD², Mehar Sultana, MSc¹, Maha Prashant Bavi, MD¹, Nasser Al-Sanea, MD³, Fouad Al-Dayel, MD⁴, Shah Fowzan S. Alkuraya, MD, FACMG⁵, Khawla S. Al-Kuraya, MD,

harvard.edu/purcell/plink/) was used, with default parameters except for minimum length of ROH, minimum number of SNPs per ROH, and maximum number of heterozygous SNPs per

The Association of Genotype-Based Inbreeding Coefficient with a Range of Physical Human Traits

Karin J. H. Verweij^{1,2}, Abdel Abdellaoui², Juha Veijola³, Sylvain Sebag-Montefiore⁴, Daniel Visscher⁵, Michael E. Goddard⁵, Matthew C. Keller^{6,7}, Marjo-Riitta Järvelin^{4,5,8,9,10}, Brendan P. Zietsch^{1,2}

In this study we defined ROHs (based on the definition from Howrigan et al. [26], as stretches of at least 1 Mb containing no heterozygous SNPs (not allowing any heterozygous SNP in the window). We set the remaining options to default, thereby ensuring >90% positive predictive value or each of the parameters for “homozyg-snp” option according to our heuristic

SNPs and, hence, to balance the number and density of SNPs in the window. We set the remaining options to default, thereby ensuring >90% positive predictive value or each of the parameters for “homozyg-snp” option according to our heuristic

Response to “Cross-Species Application of SNP Chip Suitable for Identifying Runs of Homozygosity” by S. Miller, and Kardos

Veronika Kharzinova, Alexander A. Sermyagin, Elena A. Gladyr, Gotthard J. Kardos and Natalia A. Zinovieva

We would like to clarify that we used a sliding 100-kb window with a size of 100 SNPs to research ROH. In general, the window size is 10 000 kb, not 100K SNPs as described in Shafer et al. By default, PLINK has a minimum density of 1 SNP/50 kb (Purcell et al. 2007).

Sardinians Genetic Background Explained by Runs of Homozygosity and Genomic Regions under Positive Selection

Francesca Ortu³, Fabio Rosa², Simonetta Guarnera²,
Cristina Barlassina^{4,5}, Chiara Troffa³,
Fresu³, Nicola Glorioso³, Alberto Piazza^{1,2},

Inbreeding and homozygosity in breast cancer survival

Hauke Thomsen¹, Miguel Inacio da Silva Filho¹, Andrea Woltmann¹, Robert Johansson²,
Jorunn E. Eyfjord³, Ute Hamann⁴, Jonas M. Berg⁵,
Roger Henriksson^{2,8}, Stefan Herms^{9,10}, Per L. Jonsson¹¹,
Kari Hemminki^{11,11}, Per Lenner² & Asta Förssell¹²

software (–homozyg option)). The following

Genomic inbreeding estimation in small populations: evaluation of runs of homozygosity in three local dairy cattle breeds

S. Mastrangelo^{1†}, M. Tolone¹, R. Di Gerlando¹, L. Fontanesi², M. T. Sardina¹ and
B. Portolano¹

following criteria were used to define the ROH: (i) the minimum number of SNPs included in the ROH was fixed to 40; (ii) the minimum length that constituted the ROH was set to 4 Mb; (iii) two missing SNPs were allowed in the ROH; (iv) minimum density of one SNP every 100 kb; (v) maximum gap between consecutive SNPs of 1 Mb. Moreover, the number of allowed heterozygous SNPs was set to different values: from one to three. Mean F_{ROH} values obtained

Runs of Homozygosity in European Populations

Ruth McQuillan,¹ Anne-Louise Leutenegger,² Rehab Abdel-Rahman,^{1,7} Christopher Marijana Pericic,³ Lovorka Barac-Lauc,³ Nina Smolej-Narancic,³ Branka Janicijevic,³ Albert Tenesa,⁵ Andrew K. MacLeod,⁶ Susan M. Farrington,⁵ Pavao Rudan,³ Carolin Veronique Vitart,⁷ Igor Rudan,^{1,8,9} Sarah H. Wild,¹ Malcolm G. Dunlop,⁵ Alan F. Walker,¹⁰ Harry Campbell,¹ and James F. Wilson^{1,*}

5000 kb (minimum 50 SNPs) across the genome to detect long contiguous runs of homozygous genotypes. An occasional genotyping error or missing genotype occurring in an otherwise-unbroken homozygous segment could result in the underestimation of ROHs. To address this, the program allows one heterozygous and five missing calls per window.

A threshold was set for the minimum length (kb) needed for a tract to qualify as homozygous. Because strong linkage disequilibrium (LD), typically extending up to about 100 kb, is common throughout the genome,^{48–51} short tracts of homozygosity are very prevalent. For exclusion of these short and very common ROHs that occur in all individuals in all populations, the minimum length for an ROH was set at 500 kb. All empirical studies

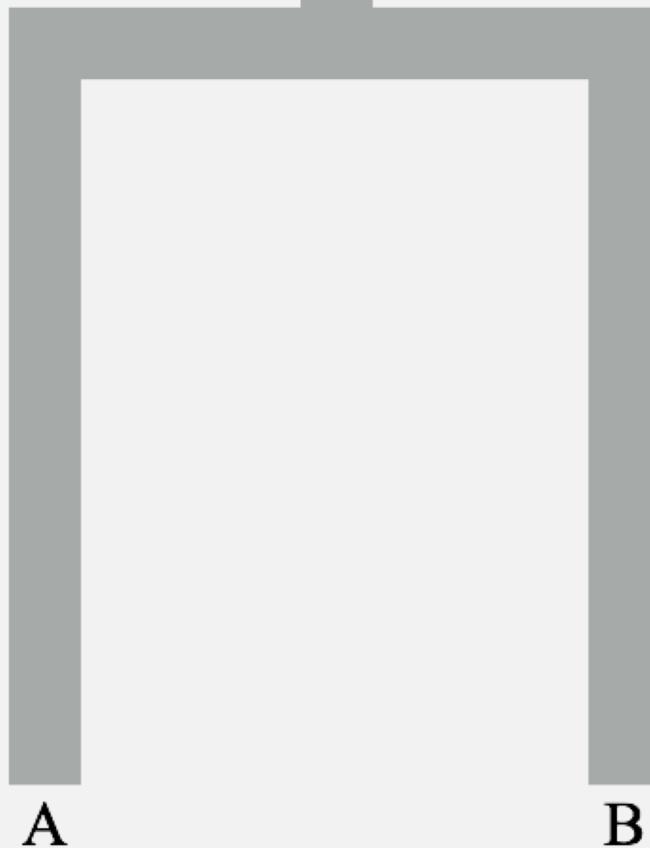
mapping in a family presenting with epilepsy and hearing impairment

Maclean², Muhammad Irfan³, Farooq Naeem⁴, Stephen Cass², Peter J Muir¹, Douglas HR Blackwood¹ and Muhammad Ayub⁵

ins of homozygosity’ analytical tool set. Inspection of homozygous tract lengths was limited to five or more consecutive SNPs, and low SNP densities in centromeric regions were excluded. The length of homozygous regions is taken to be from the most proximal to the most distal homozygous SNP, and the programme allows one heterozygous SNP within this run. Marker positions are

SIMULATE GENOME

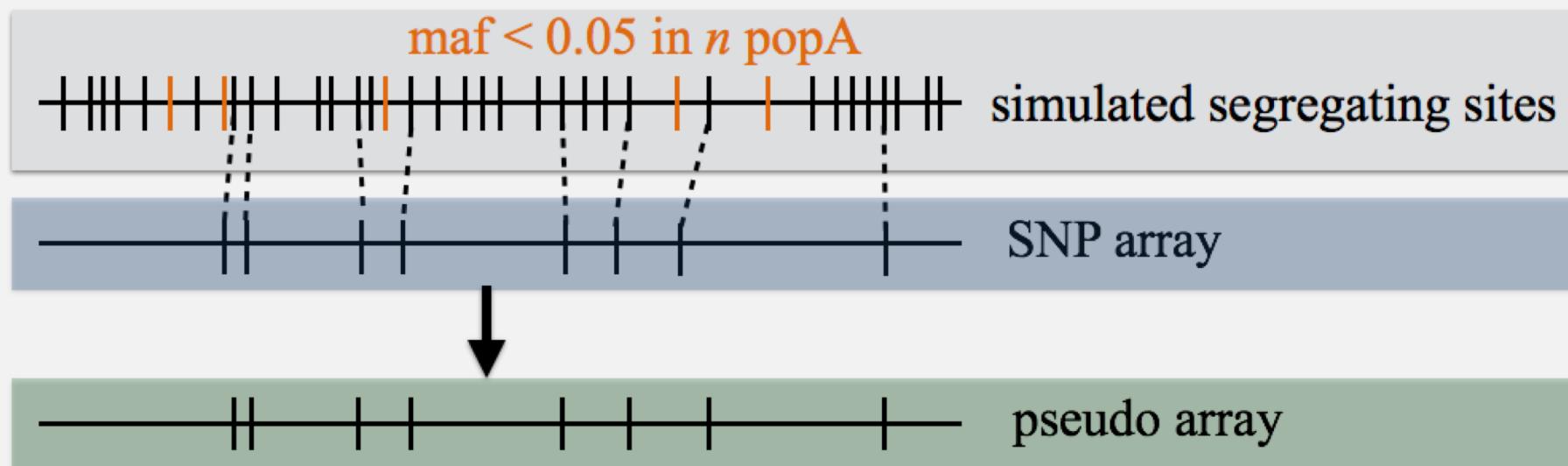
t



- Coalescent simulation
- 100 iterations
- $N_e = 1000$
- Random t , such that $F_{ST} = [0,0.2]$

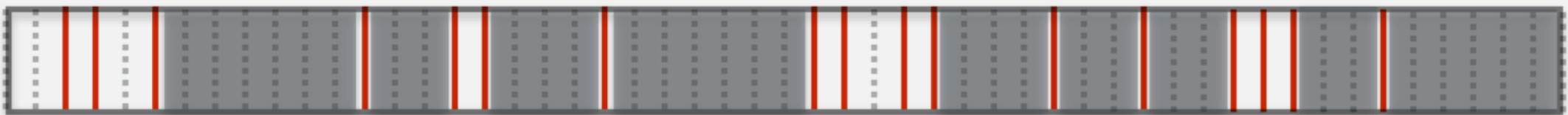
CREATE PSEUDO ARRAY

Use samples from population A to make pseudo array



FIND ROH

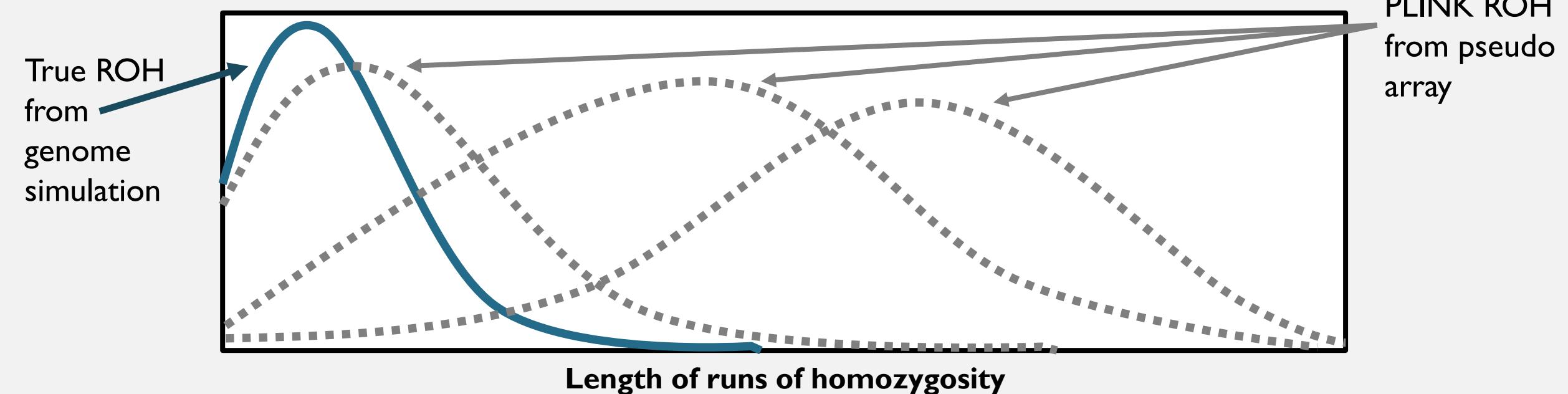
- **Genome ROH:** Script that finds pure ROH longer than k .
- **Pseudo array ROH:**
 - Script that finds pure ROH longer than k .
 - PLINK 1.09 (program used for SNP array data)



ROH on a chromosome. Each dotted grey line is a homozygous site and each red line is a heterozygous site. Each shaded grey area is a ROH

OPTIMIZE PLINK PARAMETERS

- Run PLINK on pseudo array with grid search of parameters (6,561 parameter sets)
- Identify parameter sets that give ROH closest to true ROH



IDENTIFYING CLOSEST ROH TO GENOME ROH

Mean ROH length

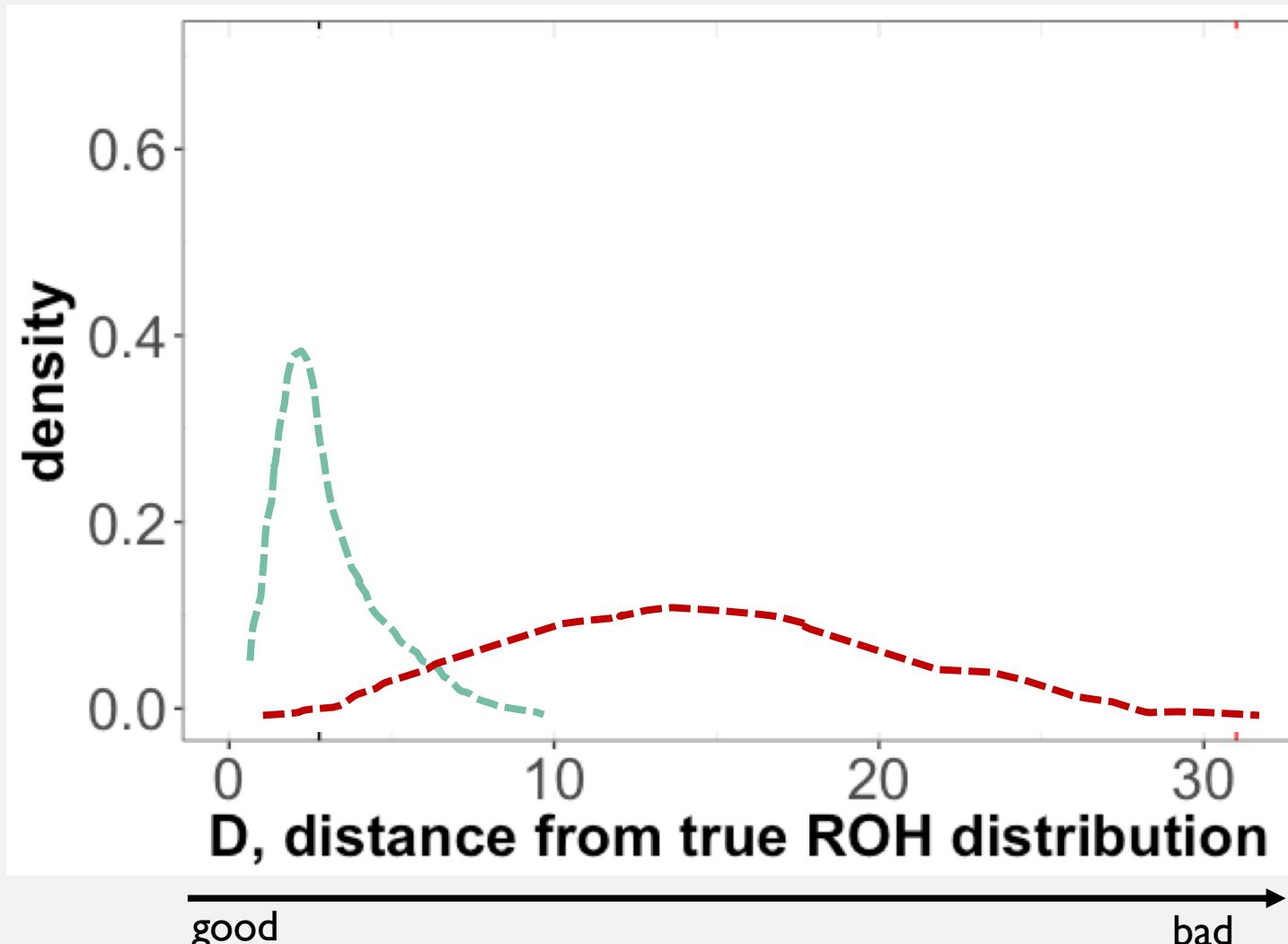
$$D = \left[\ln \left(\frac{\bar{x}}{\bar{y}} \right) \right]^2 + \left[\ln \left(\frac{\text{var}(x)}{\text{var}(y)} \right) \right]^2 + \left[\ln \left(\frac{m}{n} \right) \right]^2,$$

Variance of ROH length

Number of ROH

where x is the set of m ROH found by plink and y is the set of n real ROH.
 $D = 0$ is ideal.

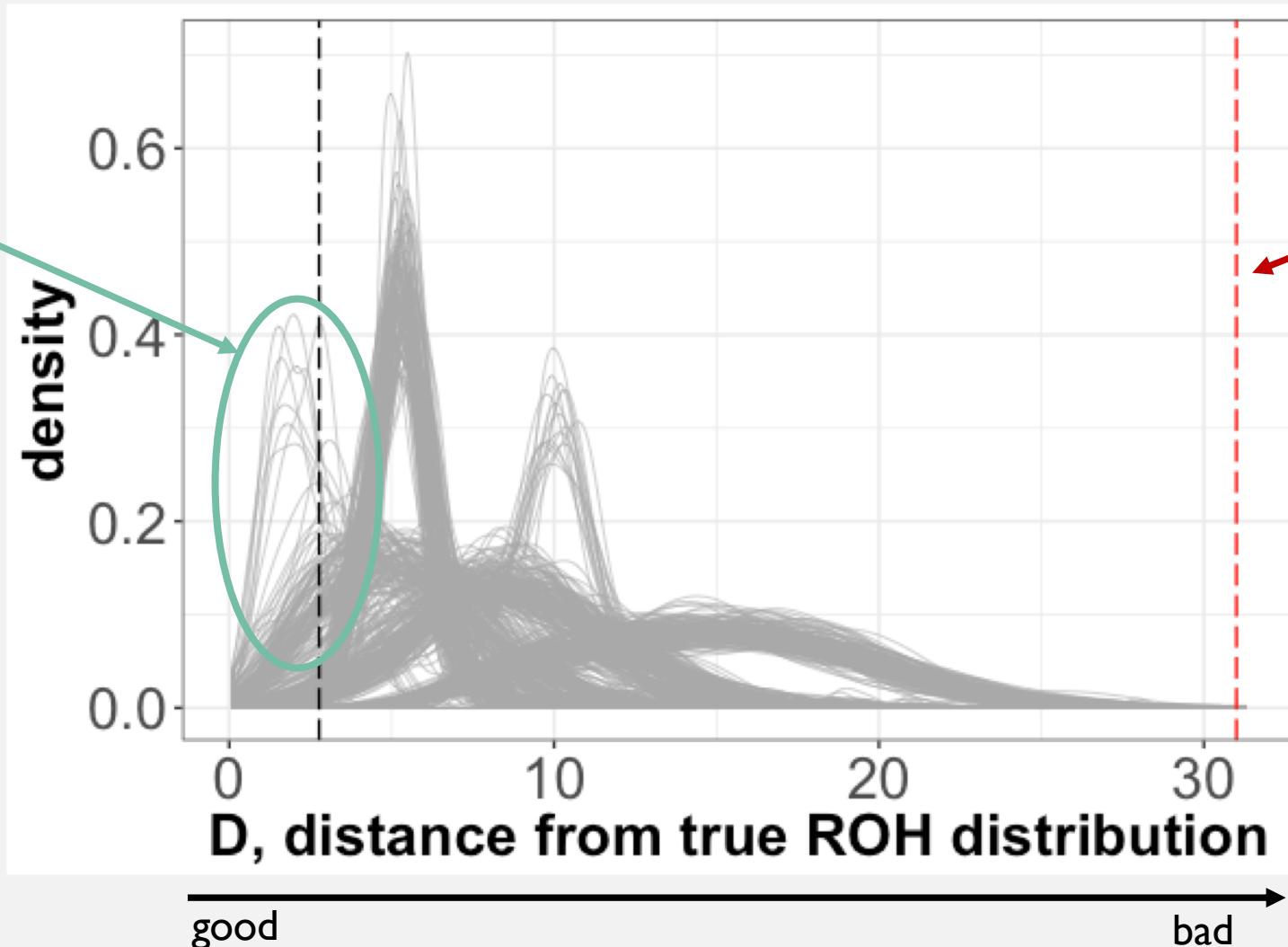
RESULTS: BEST PLINK PARAMETERS



*Not real data

RESULTS: BEST PLINK PARAMETERS

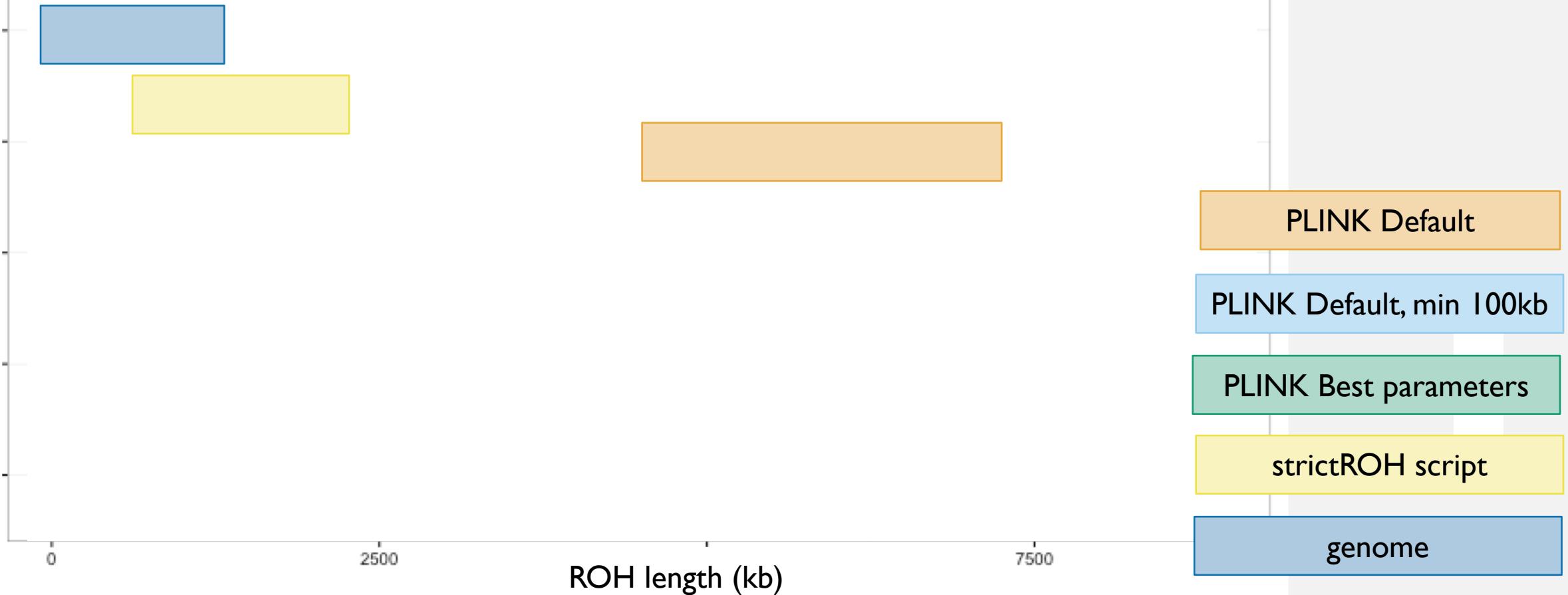
These parameter sets consistently find ROH with similar length distribution to real ROH



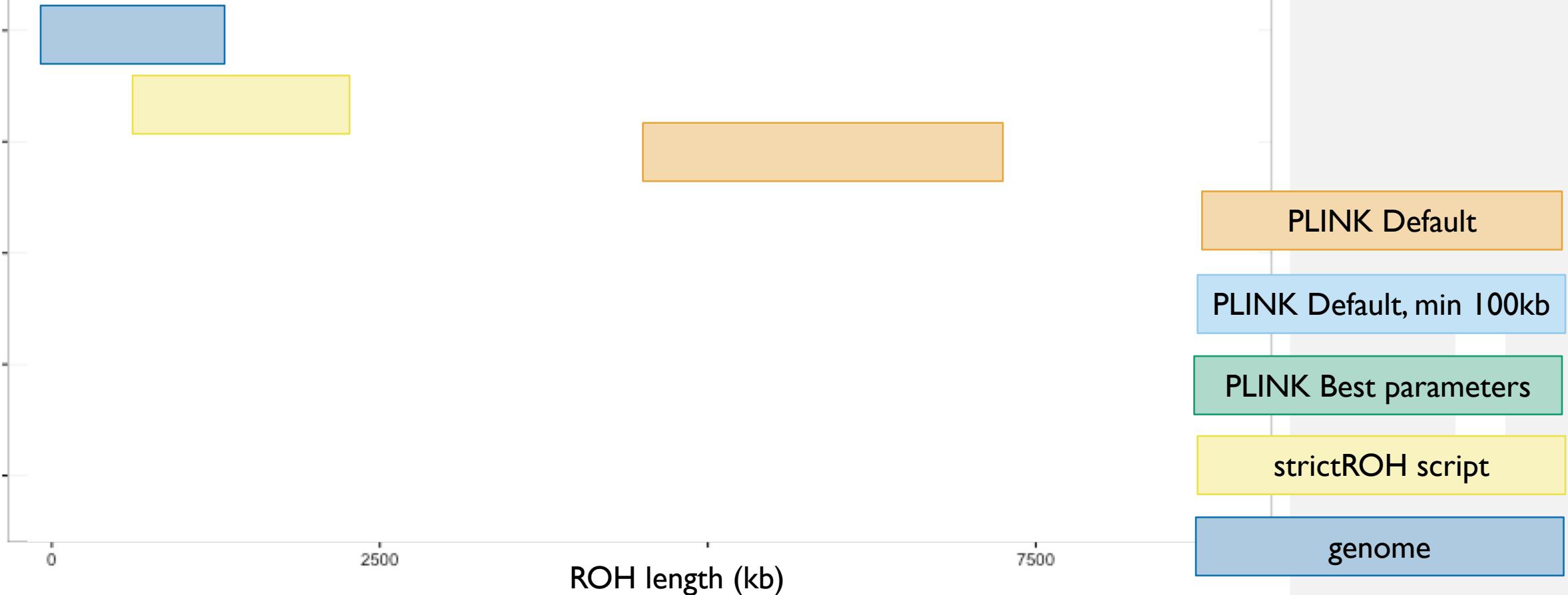
Default parameters

The distribution of d for each parameter set in the top half of the grid search, with 100 simulations. Each line is a different parameter set. The dashed line is the 5% cutoff from the distribution of d with all parameter sets.

DISTRIBUTION OF ROH LENGTHS FROM WHOLE GENOME DATA

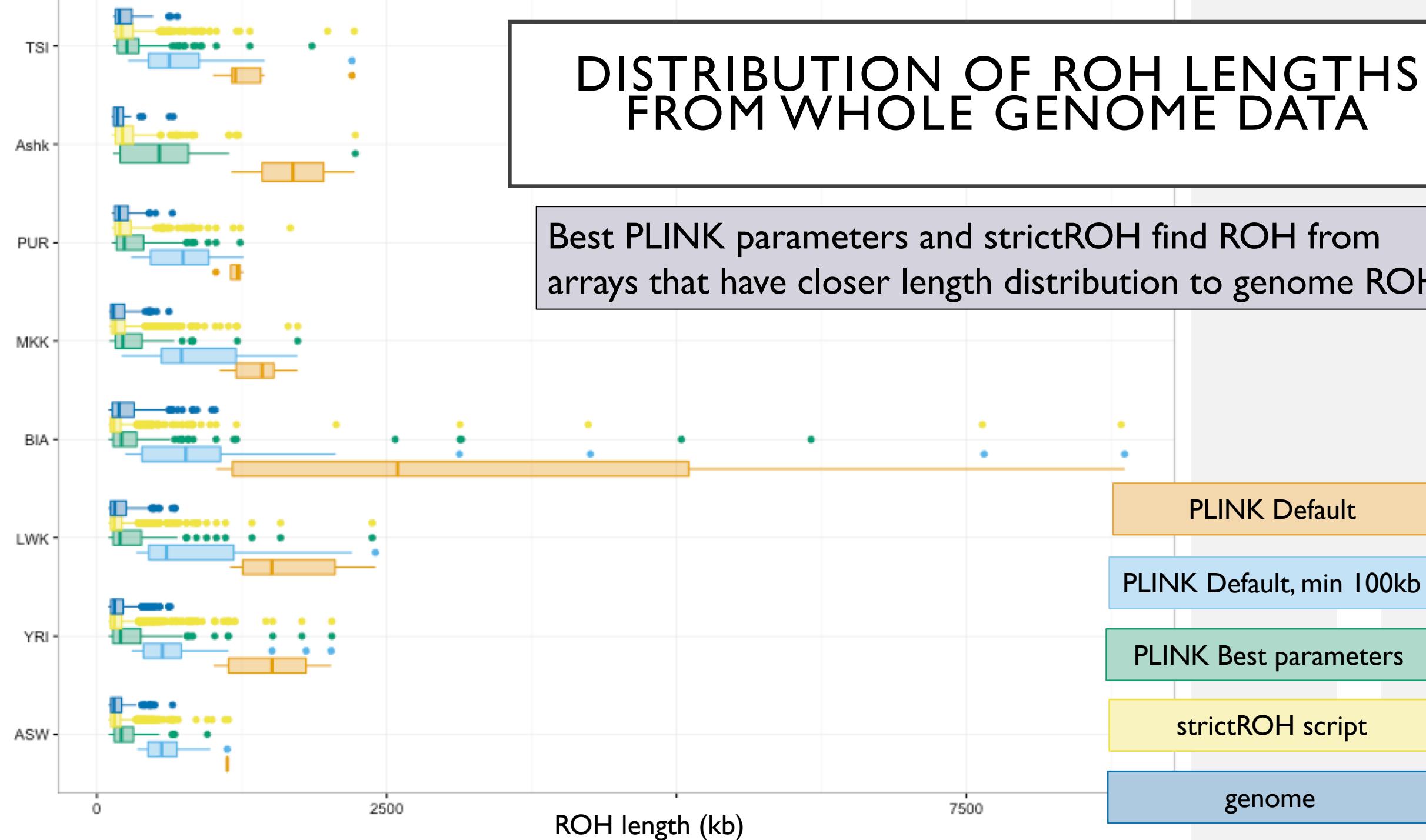


DISTRIBUTION OF ROH LENGTHS FROM WHOLE GENOME DATA



DISTRIBUTION OF ROH LENGTHS FROM WHOLE GENOME DATA

Best PLINK parameters and strictROH find ROH from arrays that have closer length distribution to genome ROH



DEVELOPED CORRECTION FOR ASCERTAINMENT BIAS WITH STRICTROH AND BEST PLINK PARAMETERS

Effect of ascertainment bias on ROH in humans not substantial

- 128 AJ whole genomes published
- Can incorporate SNP array ascertainment into model for ABC

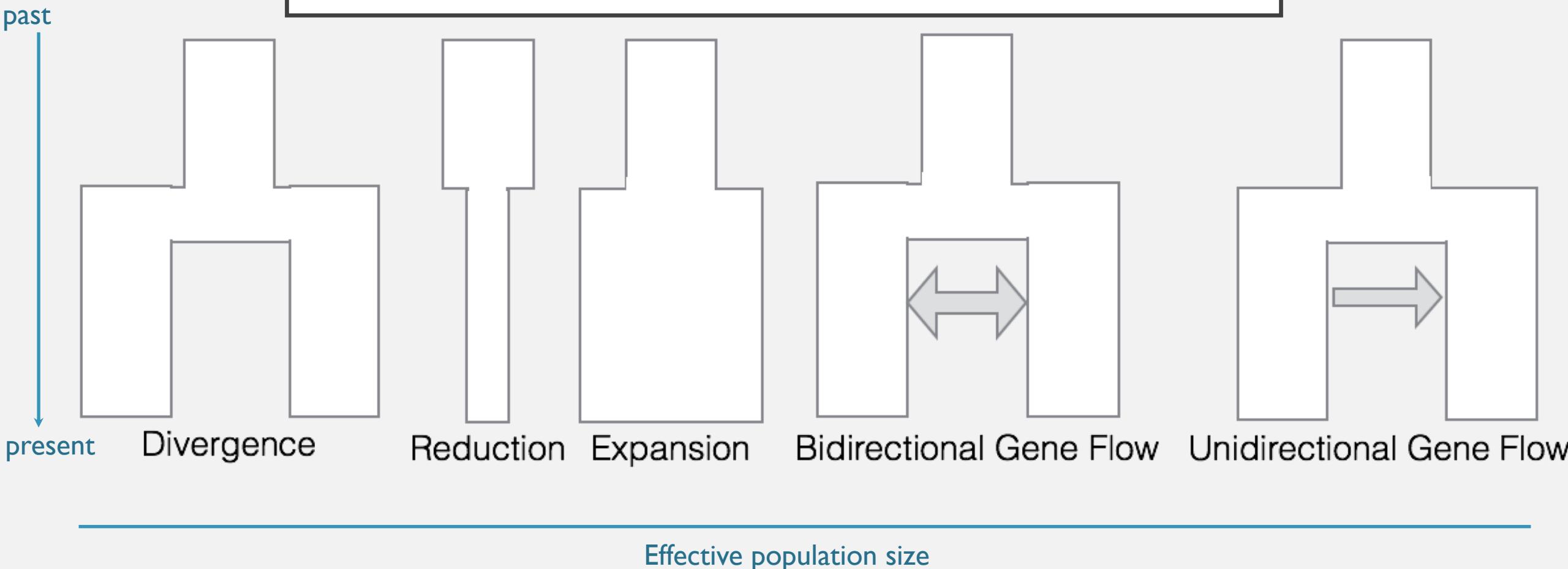
Can check for ascertainment bias in IBD

Effects of ascertainment bias taken care of

THEMES OF DISSERTATION

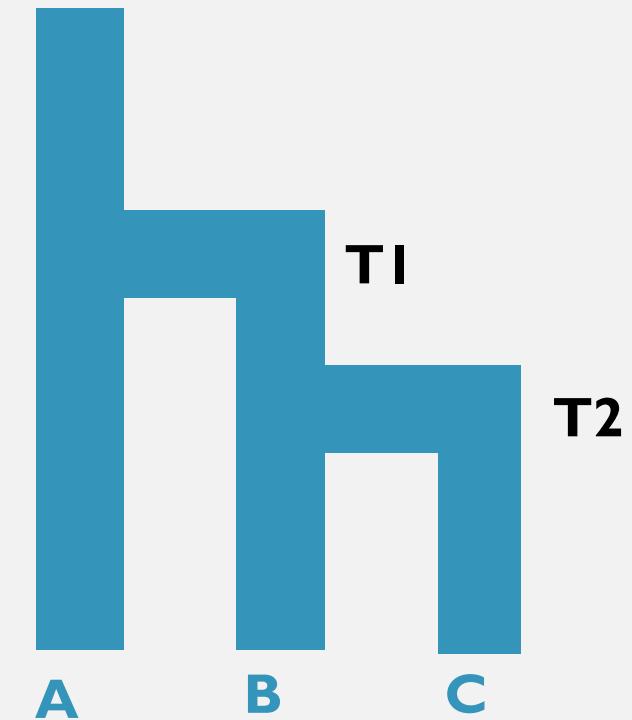
- Detection of runs of homozygosity from SNP arrays
 - Improving identification of runs of homozygosity (Ch. 2)
 - Correcting ascertainment bias in runs of homozygosity (App. C)
- Scaling up Approximate Bayesian Computation for whole chromosomes
 - Create efficient pipeline to simulate demographic models and calculate summary statistics (App. A)
 - Create generalized high throughput workflow (Ch. 4)
- Infer history of the Ashkenazi Jews
 - Substructure in AJ? (Ch. 5)
 - Khazarian origin? (App. B)

EXAMPLES OF DEMOGRAPHIC EVENTS



DEMOGRAPHIC PARAMETERS DEFINE POPULATIONS' HISTORIES

- A demographic model generates data, determined by a set of parameters
- Parameter examples:
 - population sizes,
 - divergence times,
 - admixture proportions, etc.



WHAT IS ABC?

- We want the posterior probability of the parameters given the data (D)

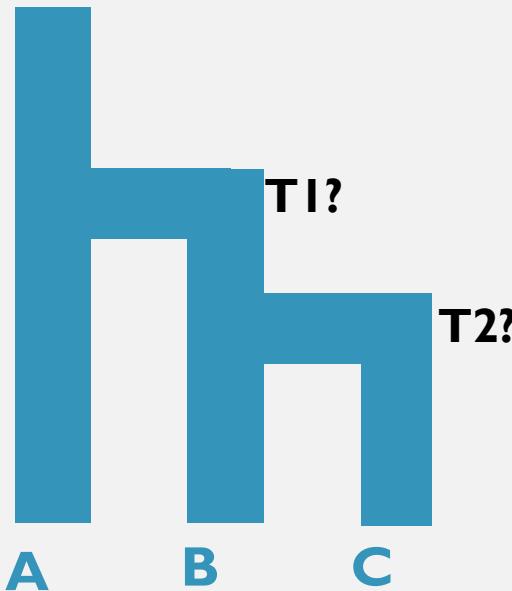
Likelihood of the data

$$P(\theta|D) = \frac{P(D|\theta)P(\theta)}{P(D)}$$

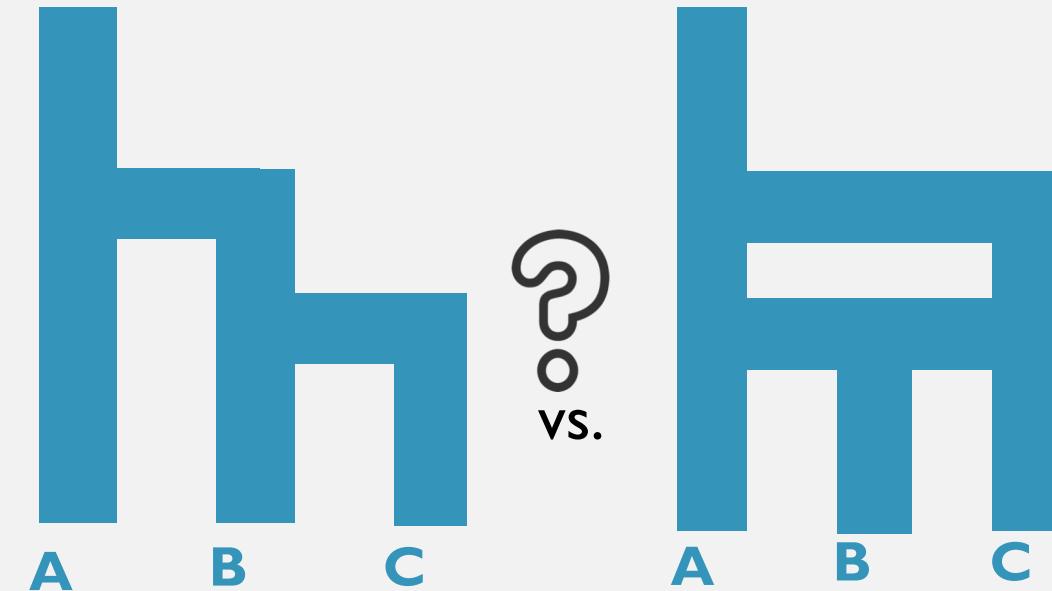
Posterior probability Likelihood of the data Prior probability of the parameter
Marginal likelihood

- Approximate the likelihood function by simulations that are compared to the data (D)

WHAT IS APPROXIMATE BAYESIAN COMPUTATION (ABC) USED FOR?



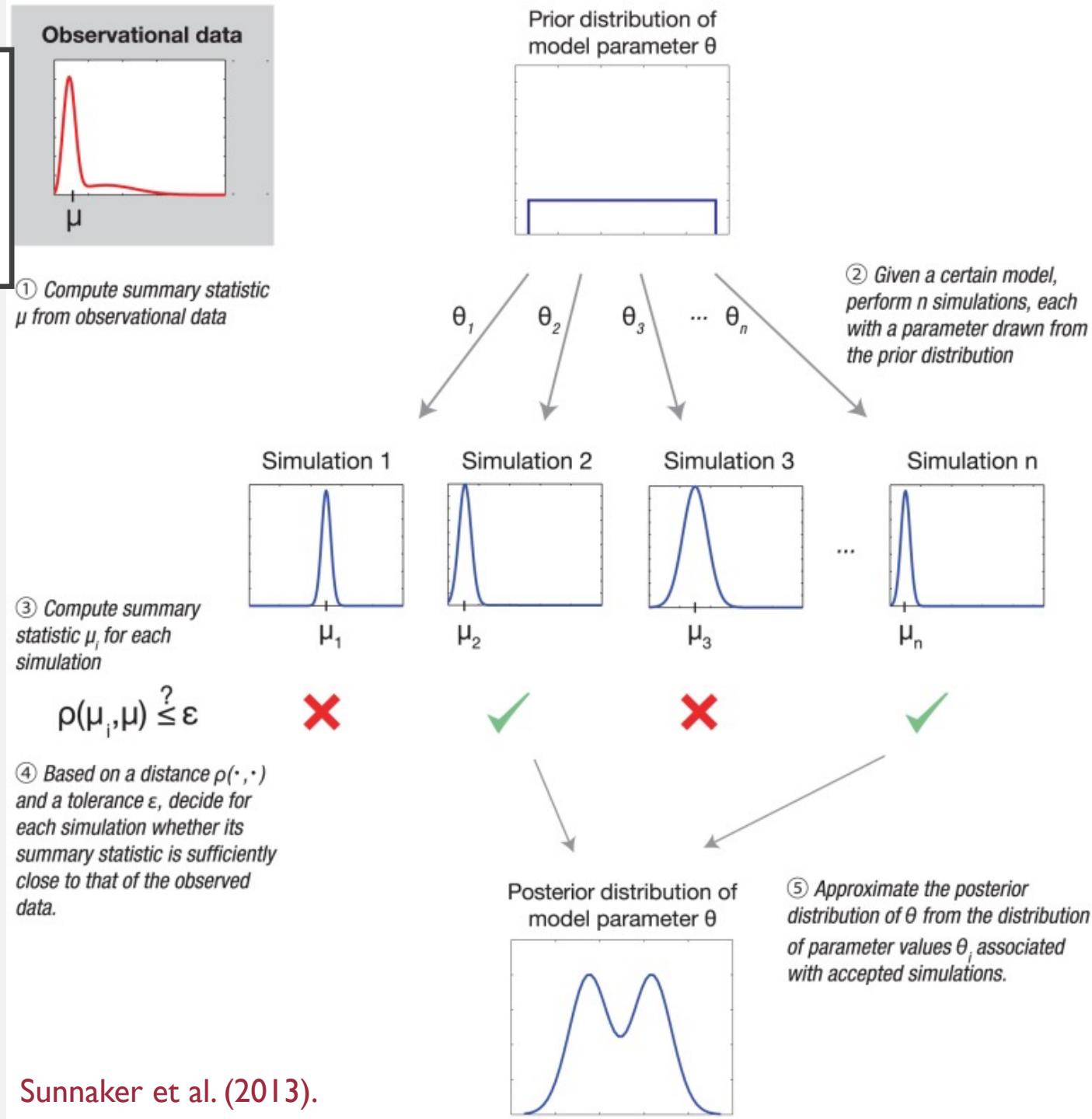
Infer parameter values



Choose among models

OVERVIEW OF ABC STEPS

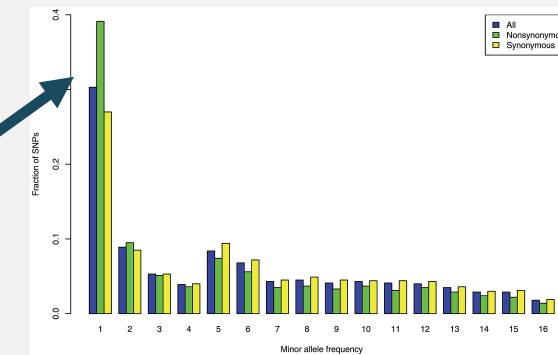
1. Define priors of parameters of model
2. Simulate data many times
3. Summarize genetic data with population genetics statistics
4. Choose model and estimate parameters based on simulations closest to observed data



IMPLEMENTATION OF ABC WITH ABCTOOLBOX: 0. COLLECT DATA AND CALCULATE SUMMARY STATISTICS



```
| rs3094315 0 742429 0 | 0 | 0 0 0 0 | 1 0 0 0 | 0 1 0 0 0 0 0 | 0 0 0 0  
| rs12562034 0 758311 | 1 1 1 1 1 1 | 0 1 1 1 1 1 | 0 1 1 1 1 1 | 0 1 1 1 1 1  
| rs3934834 0 995669 | 1 0 0 1 1 0 0 1 1 1 1 1 | 0 1 0 1 0 1 1 1 1 1 1 1  
| rs9442372 0 1008567 0 | 0 1 0 1 0 0 0 1 1 1 1 0 0 0 0 1 1 1 1 1 1  
| rs3737728 0 1011278 0 0 0 1 0 1 1 0 1 1 1 1 1 0 1 1 1 0 1 1 1 0 1 1 1  
| rs11260588 0 1011521 | 1 1 1 0 1 1 1 0 1 0 1 0 1 0 1 1 1 1 1 1 0 1 0 1  
| rs9442398 0 1011558 0 0 0 1 0 1 0 0 0 1 1 1 1 0 0 0 1 0 1 1 0 1 1 1  
| rs6687776 0 1020428 0 | 0 1 0 1 1 0 1 0 0 0 1 0 1 0 0 0 1 1 1 0 1 0 1  
| rs9651273 0 1021403 | 1 1 1 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1  
| rs4970405 0 1038818 | 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1  
| rs12726255 0 1039813 0 | 0 1 0 1 1 0 1 0 1 0 0 0 0 1 1 0 1 1 1 0 0 0  
| rs7540009 0 1050098 | 1 1 1 0 1 1 1 0 1 1 1 0 1 1 1 1 1 1 1 1 1 0 1 1 1  
| rs11807848 0 1051029 | 1 0 0 0 0 1 1 0 0 1 1 0 0 1 1 1 1 1 0 1 1 1 0 0 1  
| rs9442373 0 1052501 0 | 0 0 0 0 1 1 0 0 1 1 0 0 0 0 1 1 0 1 1 1 0 0 1  
| rs2298217 0 1054842 | 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1  
| rs12145826 0 1055892 | 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1  
| rs4970357 0 1066927 | 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 1 1 1 1 1 1 1  
| rs9442380 0 1077546 | 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1  
| rs7553429 0 1080420 0 | 1 1 0 1 1 1 1 1 1 1 1 1 0 0 0 1 1 0 1 0 0 1  
| rs4970362 0 1084601 | 1 0 1 0 1 0 0 1 1 1 1 1 1 1 1 1 0 1 0 0 0 1 0 1 0 0 0
```

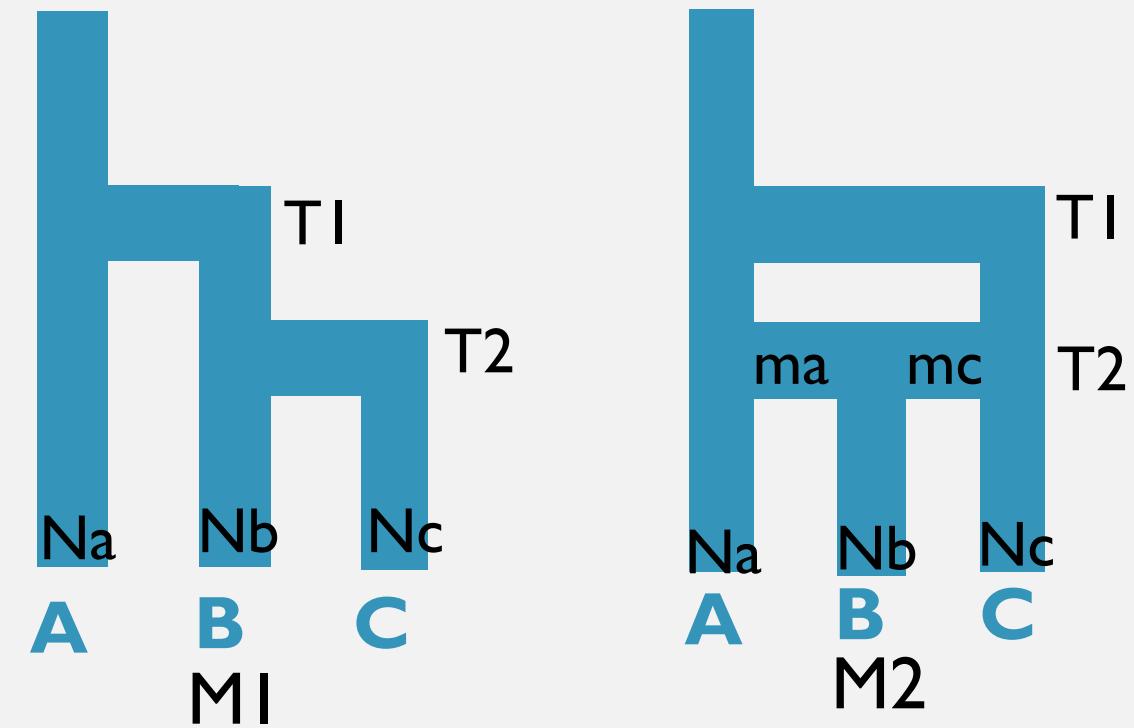
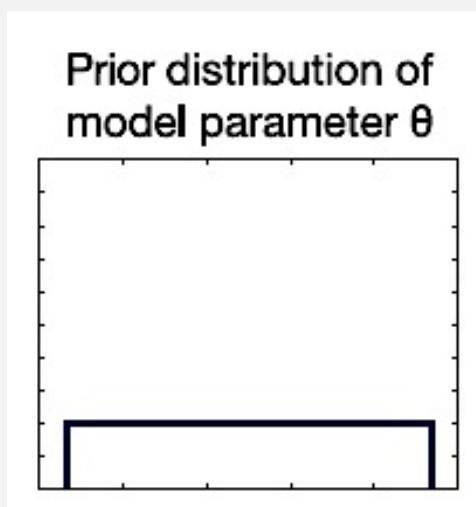


- Number of segregating sites
 - nucleotide diversity
 - Fst
 - Tajima’s D
 - IBD stats, etc.

IMPLEMENTATION OF ABC WITH ABCTOOLBOX:

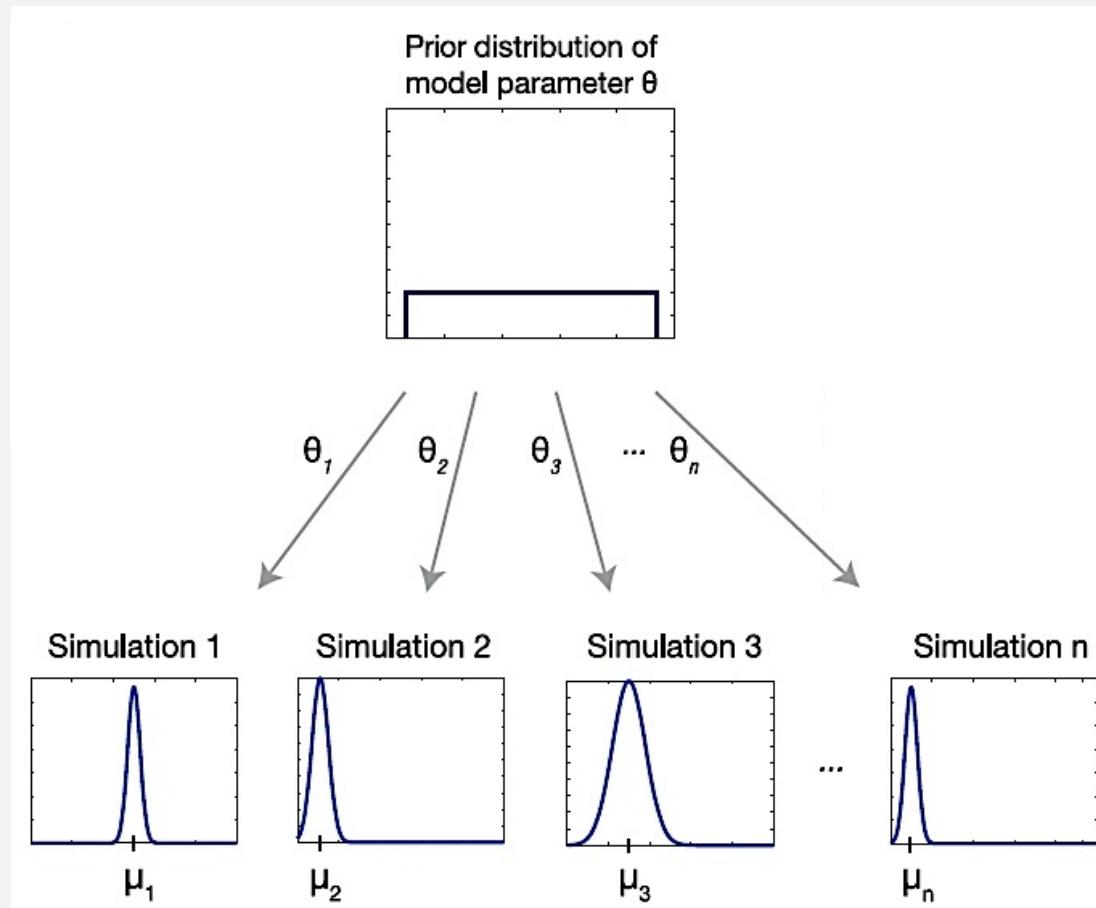
I. PICK MODELS AND PRIORS

- Parameters (Θ):
 - Divergence times (T_i)
 - Population sizes (N_j)
 - Proportion of gene flow (m_j)
 - etc...



IMPLEMENTATION OF ABC WITH ABCTOOLBOX:

2. SIMULATE MODELS ACCORDING TO THE PRIORS AND CALCULATE SUMMARY STATISTICS



IMPLEMENTATION OF ABC WITH ABCTOOLBOX:

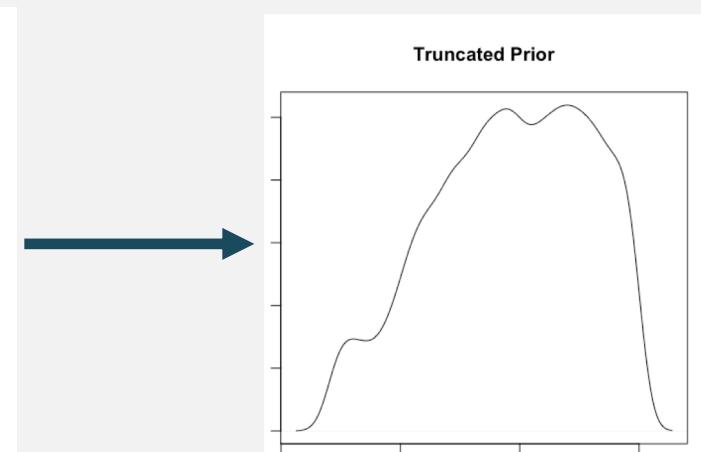
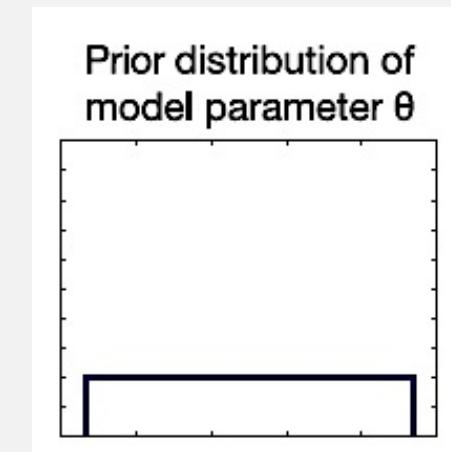
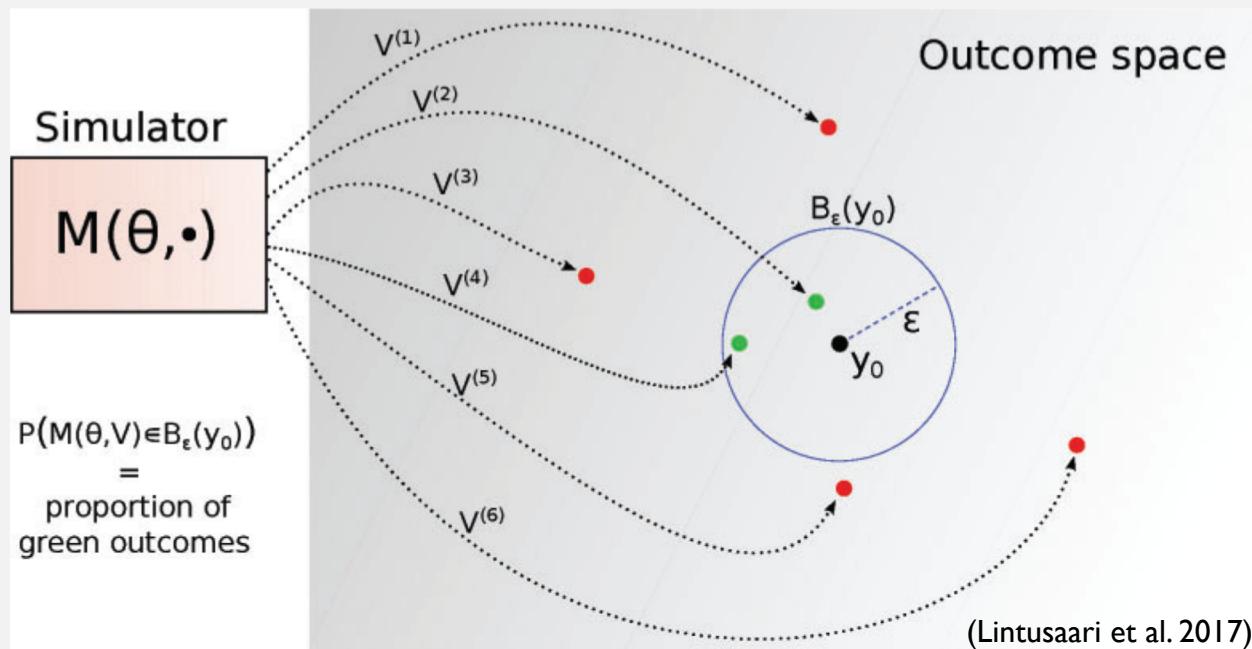
3. ADDRESS CORRELATIONS AMONG STATISTICS

- Prune statistics for high pairwise correlation
 - or
- Transform statistics with Partial Least Squares (PLS)

IMPLEMENTATION OF ABC WITH ABCTOOLBOX:

4. RETAIN N CLOSEST SIMULATIONS TO OBSERVED DATA

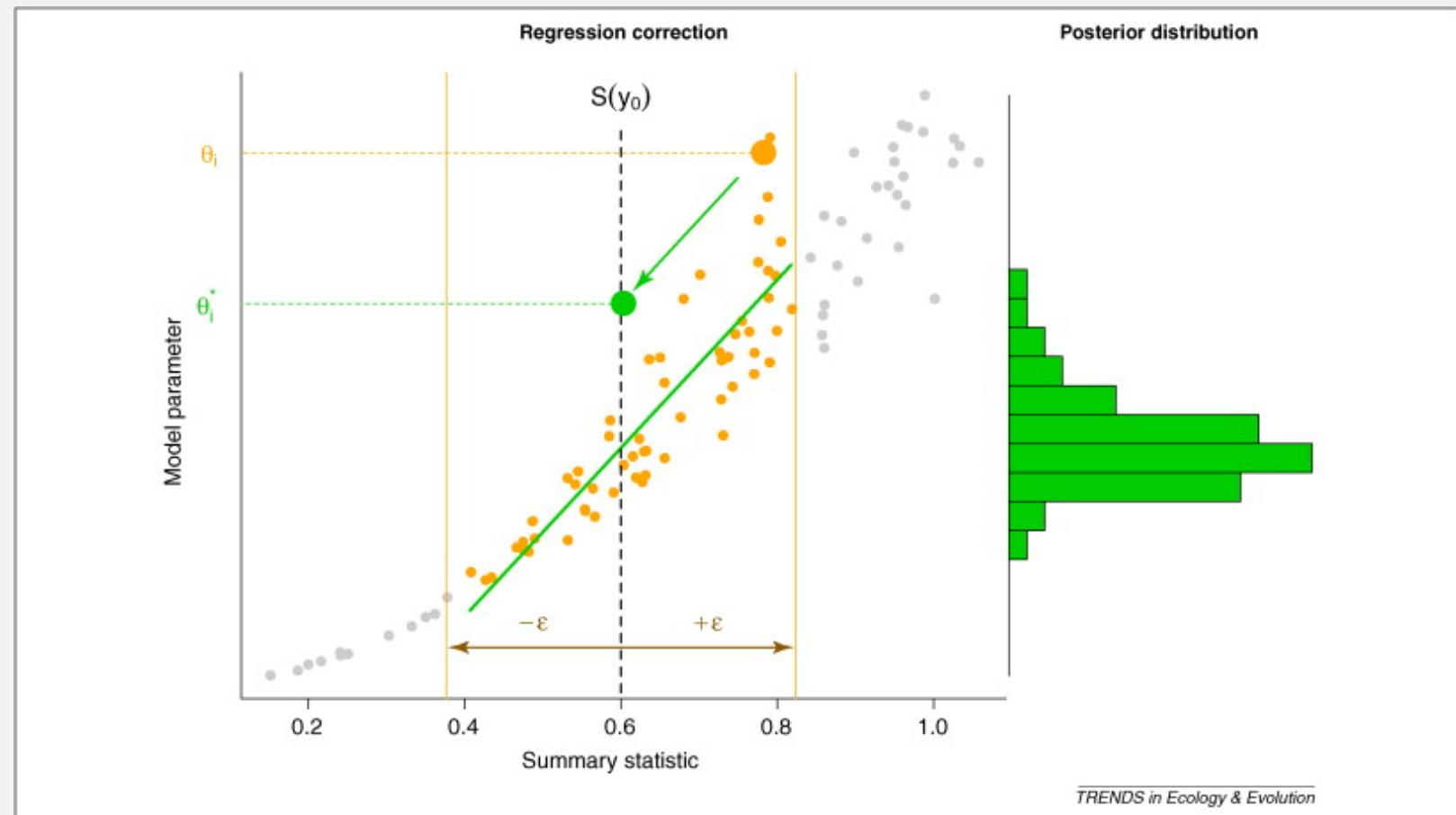
- Creates truncated prior by accepting some proportion of parameters and summary stats pairs closest to observed data
- Closest is defined by Euclidean distance between the simulated and observed summary statistics



IMPLEMENTATION OF ABC WITH ABCTOOLBOX:

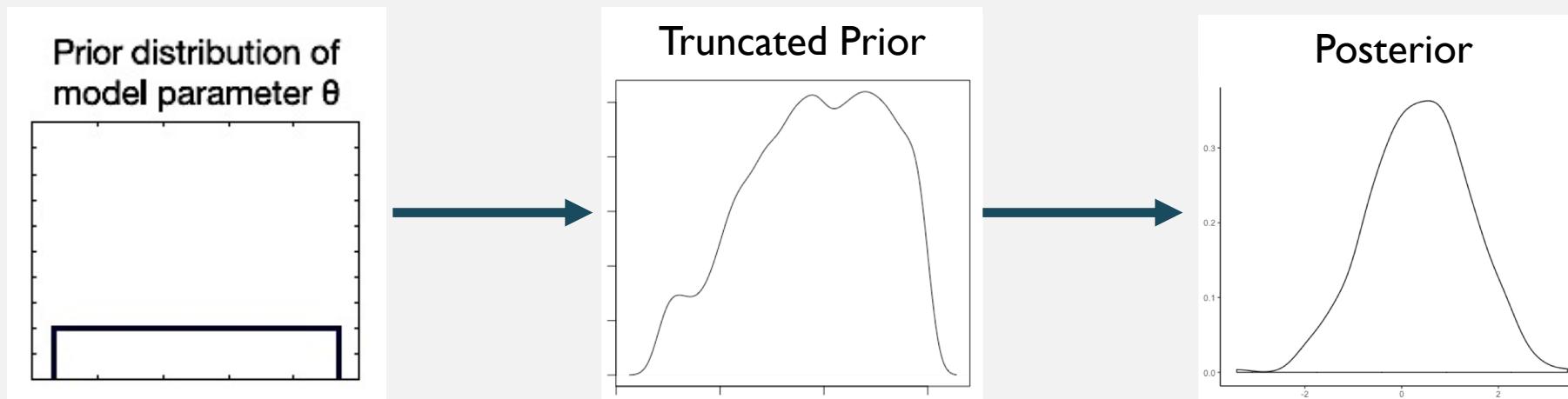
5. LINEAR REGRESSION ON THE SUMMARY STATISTICS AND TRUNCATED PRIOR

- Retained parameter values adjusted according to a linear transformation
- New parameter values form a sample from the posterior

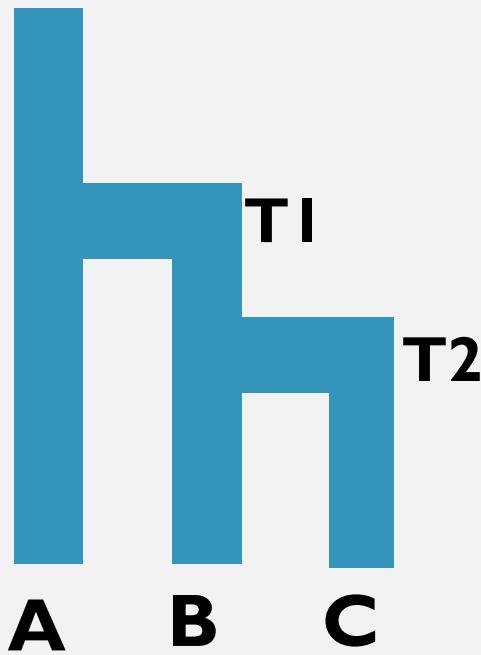


IMPLEMENTATION OF ABC WITH ABCTOOLBOX:

7. BUILD POSTERIOR DISTRIBUTION OF PARAMETERS

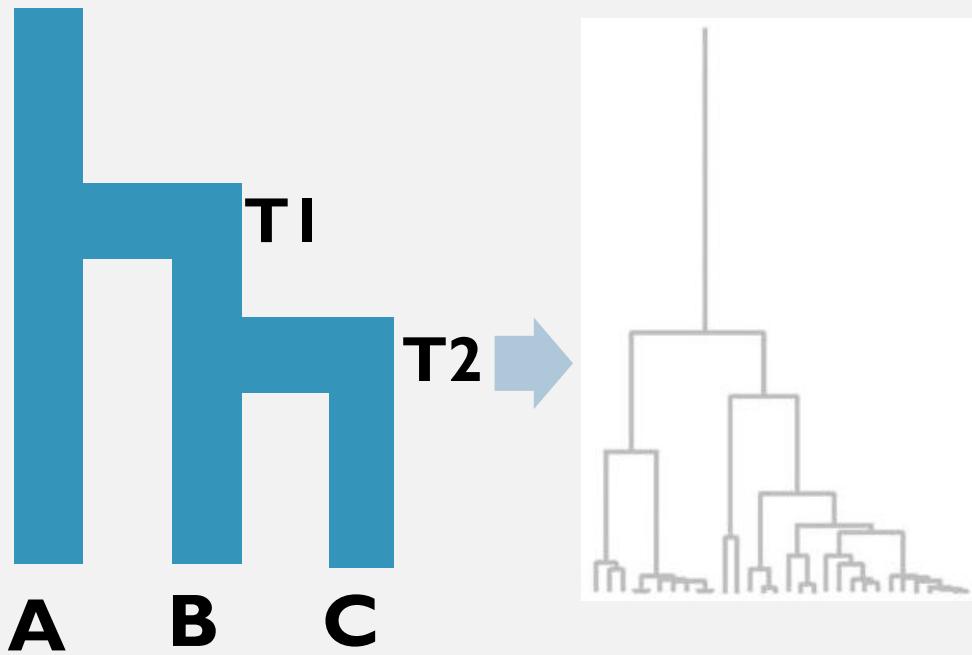


HOW DO WE PERFORM SIMULATIONS AND CALCULATE SUMMARY STATISTICS?



Define
model

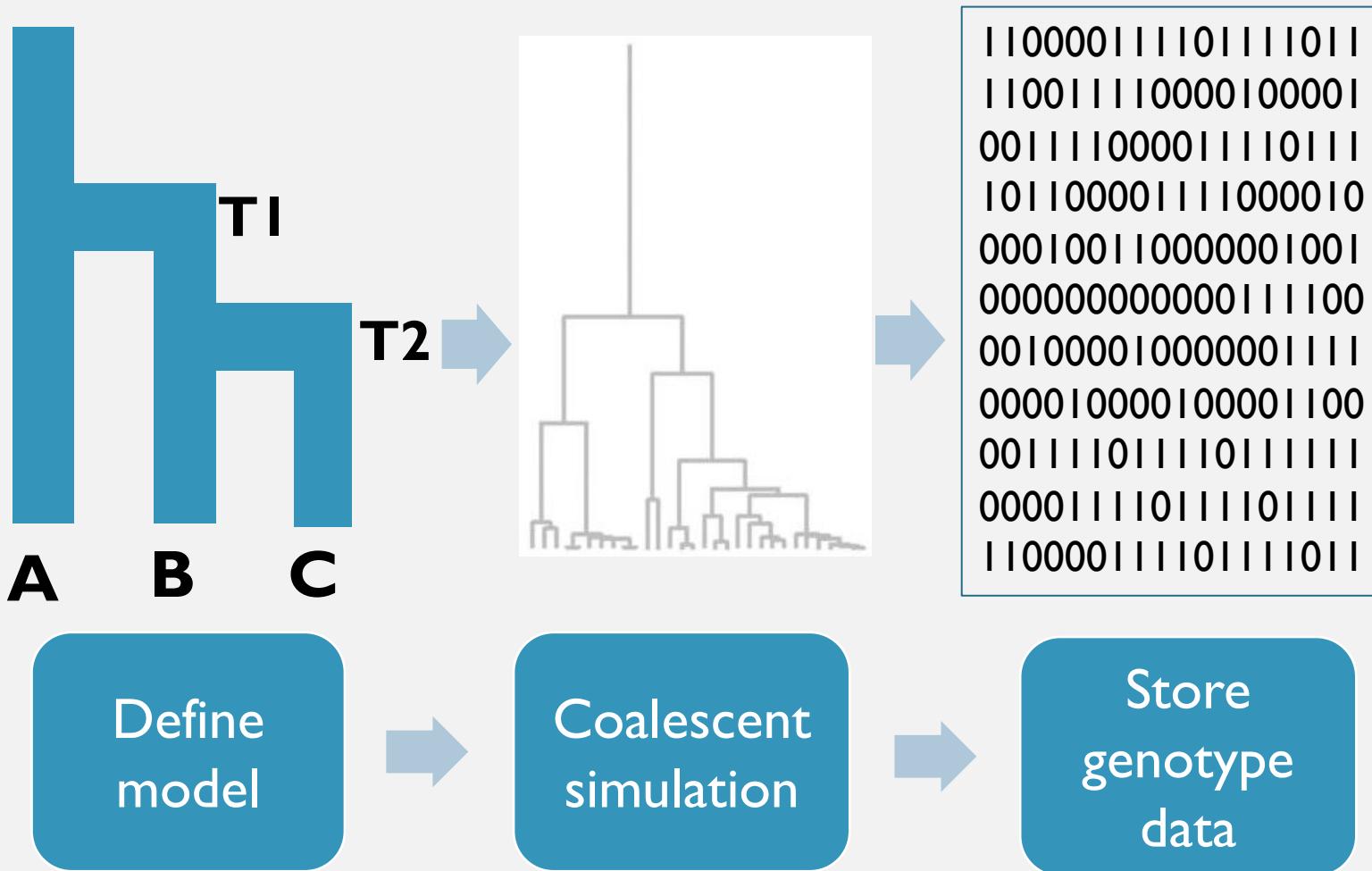
HOW DO WE PERFORM SIMULATIONS AND CALCULATE SUMMARY STATISTICS?



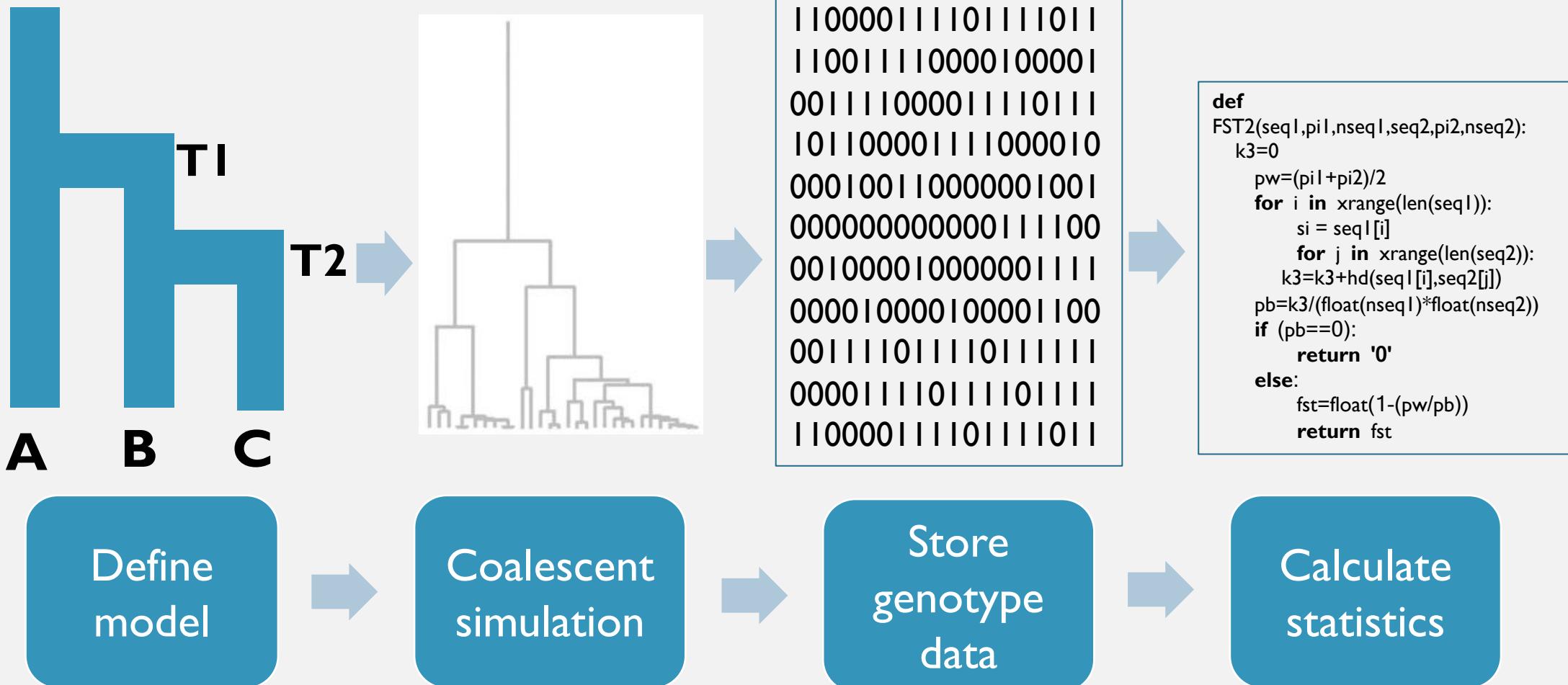
Define
model

Coalescent
simulation

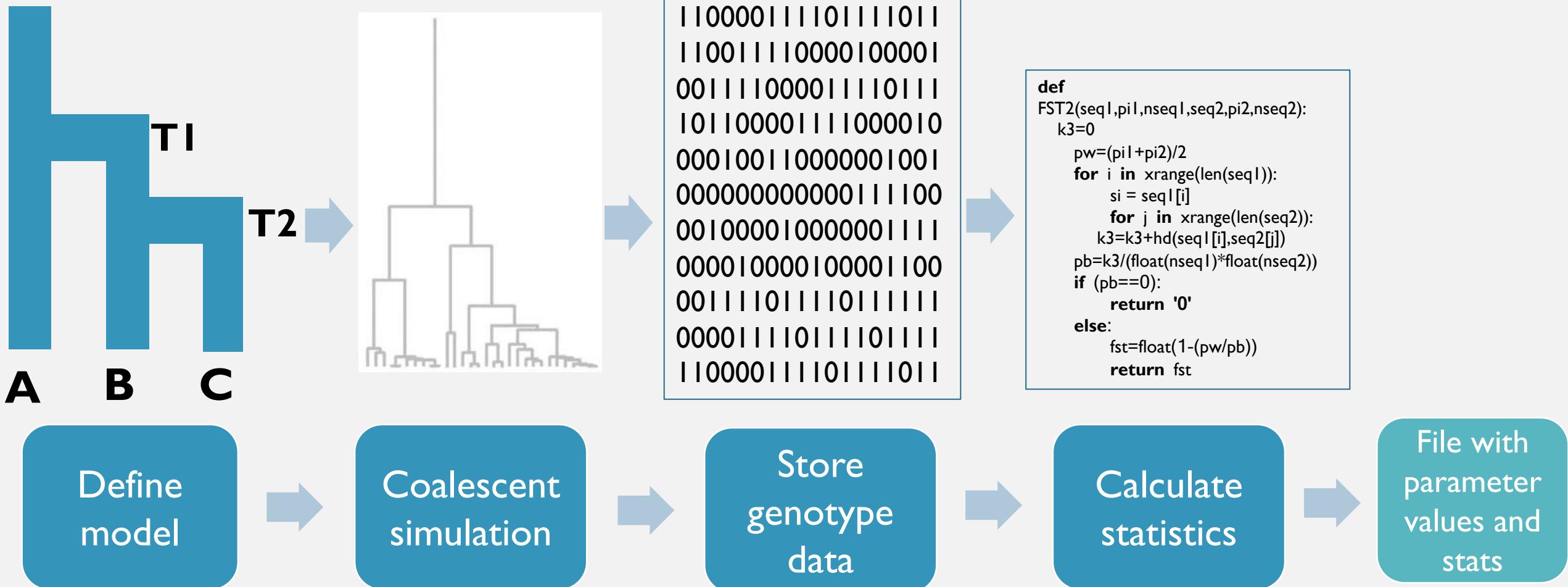
HOW DO WE PERFORM SIMULATIONS AND CALCULATE SUMMARY STATISTICS?

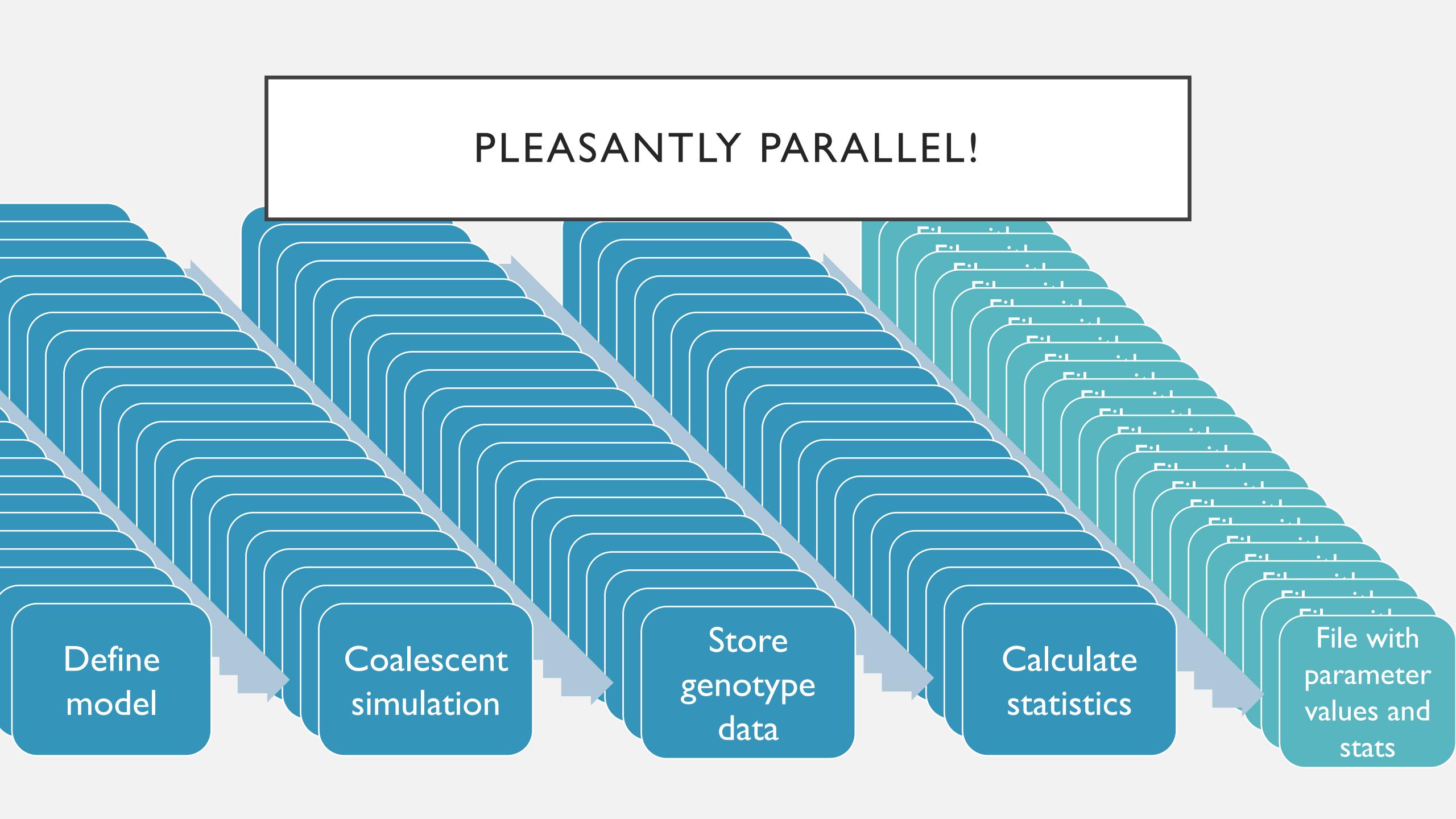


HOW DO WE PERFORM SIMULATIONS AND CALCULATE SUMMARY STATISTICS?



HOW DO WE PERFORM SIMULATIONS AND CALCULATE SUMMARY STATISTICS?





PLEASANTLY PARALLEL!

Define
model

Coalescent
simulation

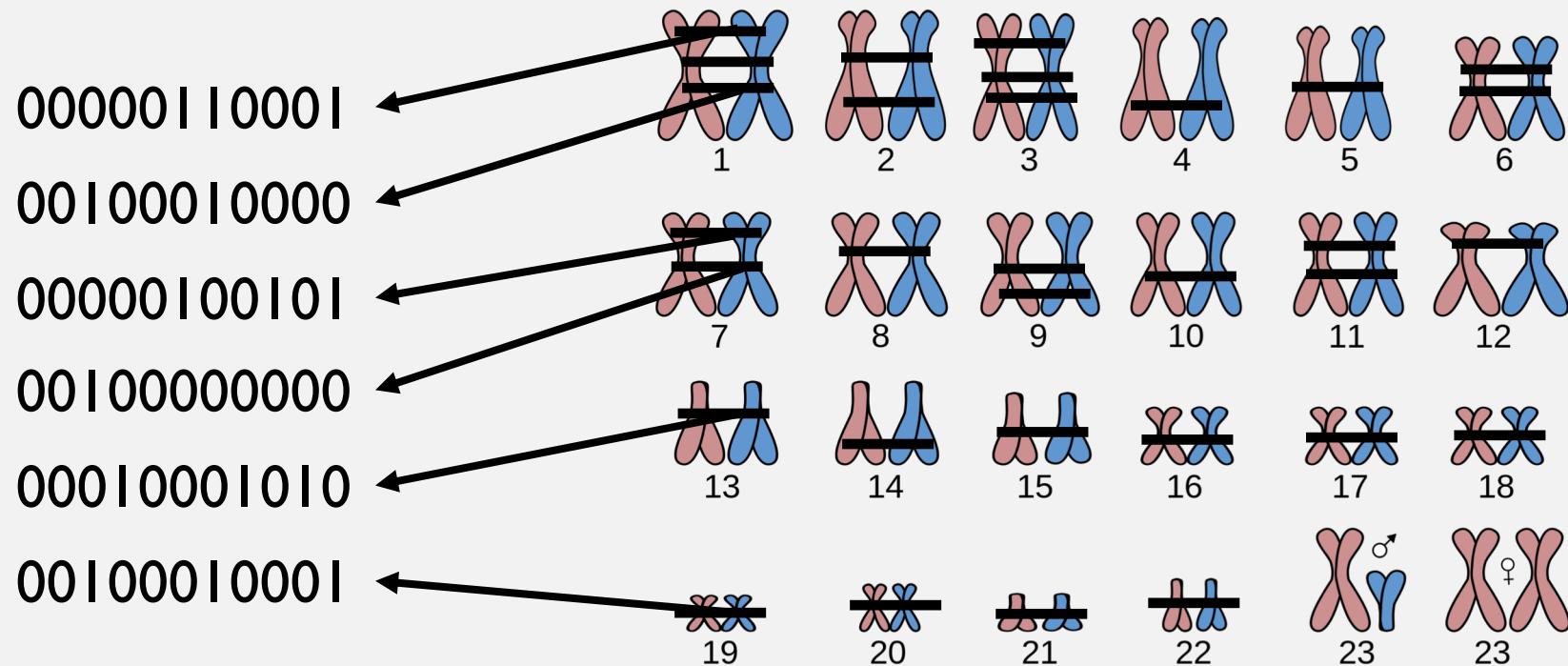
Store
genotype
data

Calculate
statistics

File with
parameter
values and
stats

INHERITED SCRIPT INTENDED FOR SMALL SEQUENCE

- Intended for millions of relatively small simulations
 - 1,389 10kb regions
 - 65 individuals
- Took a few minutes to run one simulation
- Ran parallel on U of A HPC
 - 1 million runs would take approximately 1 month.

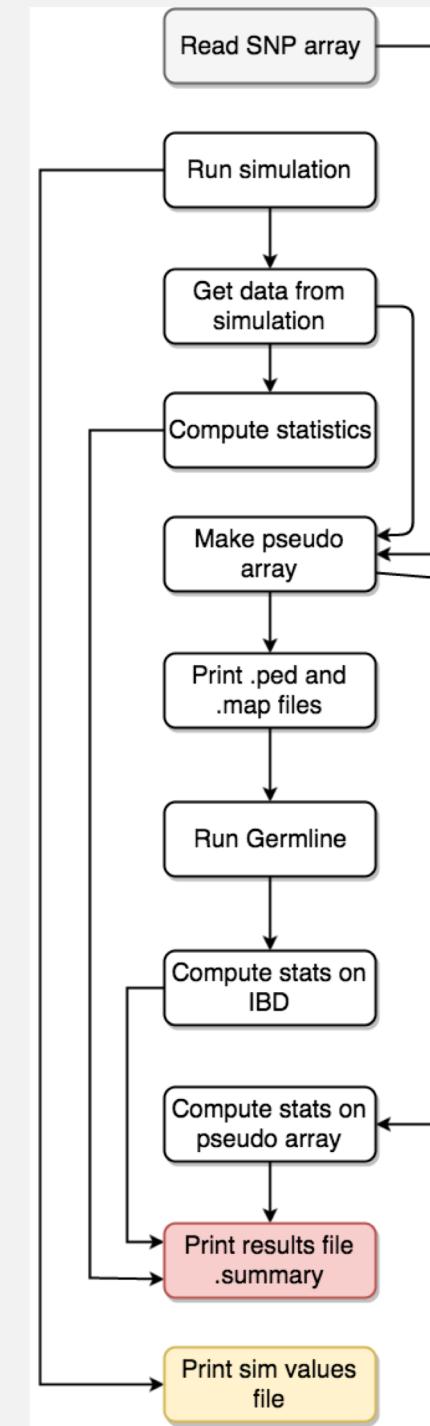


SIMULATE WHOLE CHROMOSOME

~250 million sites on human chromosome I

SIMULATE WHOLE CHROMOSOME

- Modified Python script to
 - Simulate whole chromosome
 - Find IBD segments and calculate IBD stats



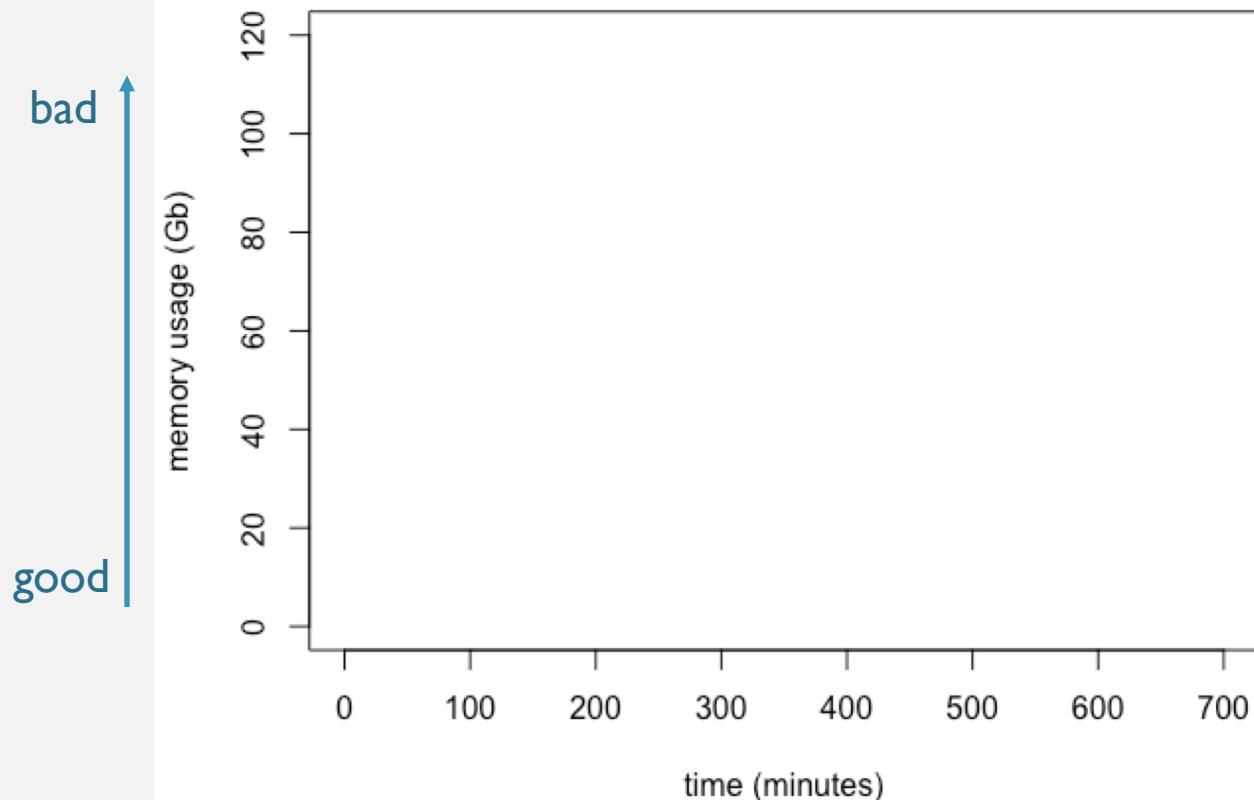
PROBLEM!

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

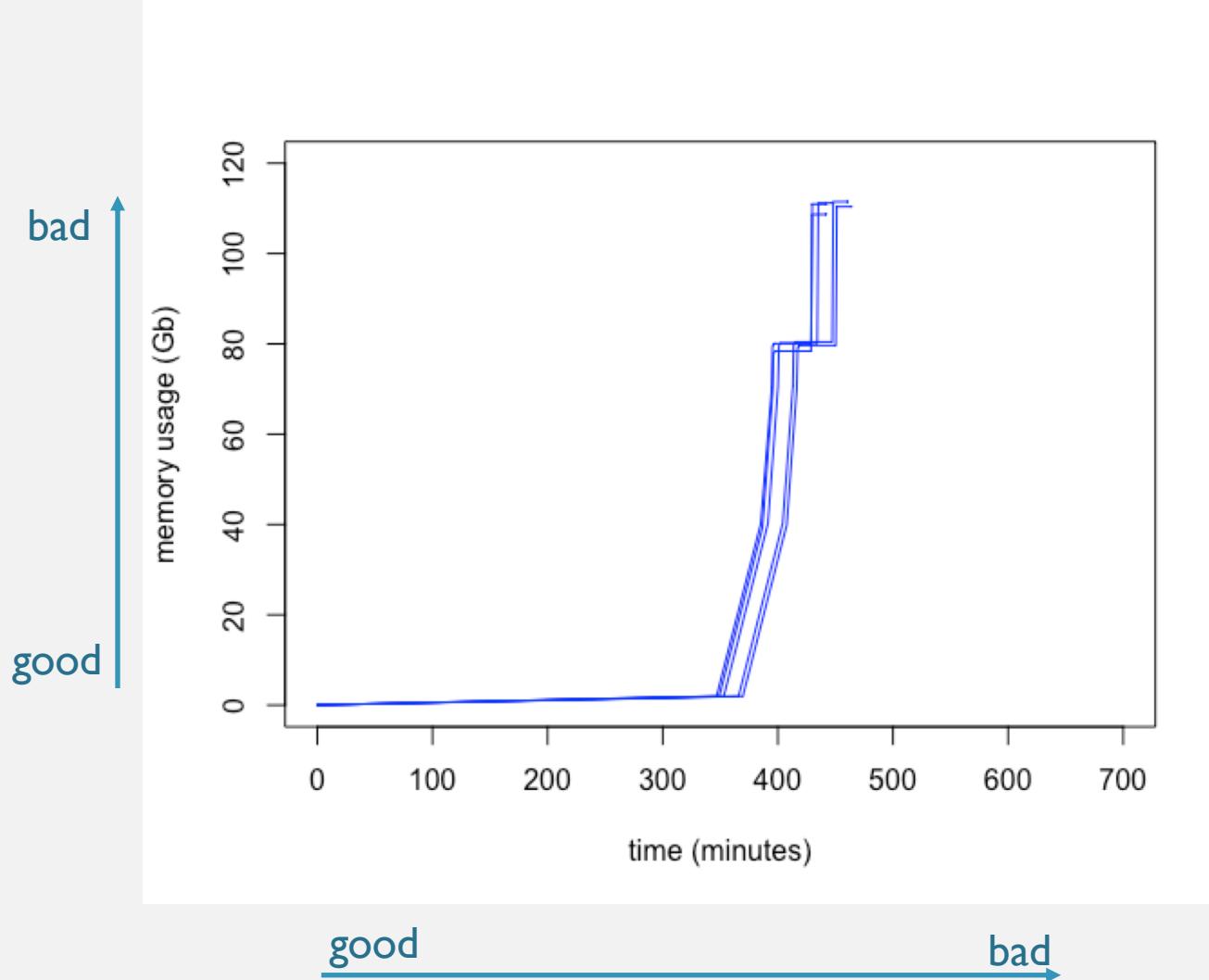
PROFILE OF PYTHON SCRIPT



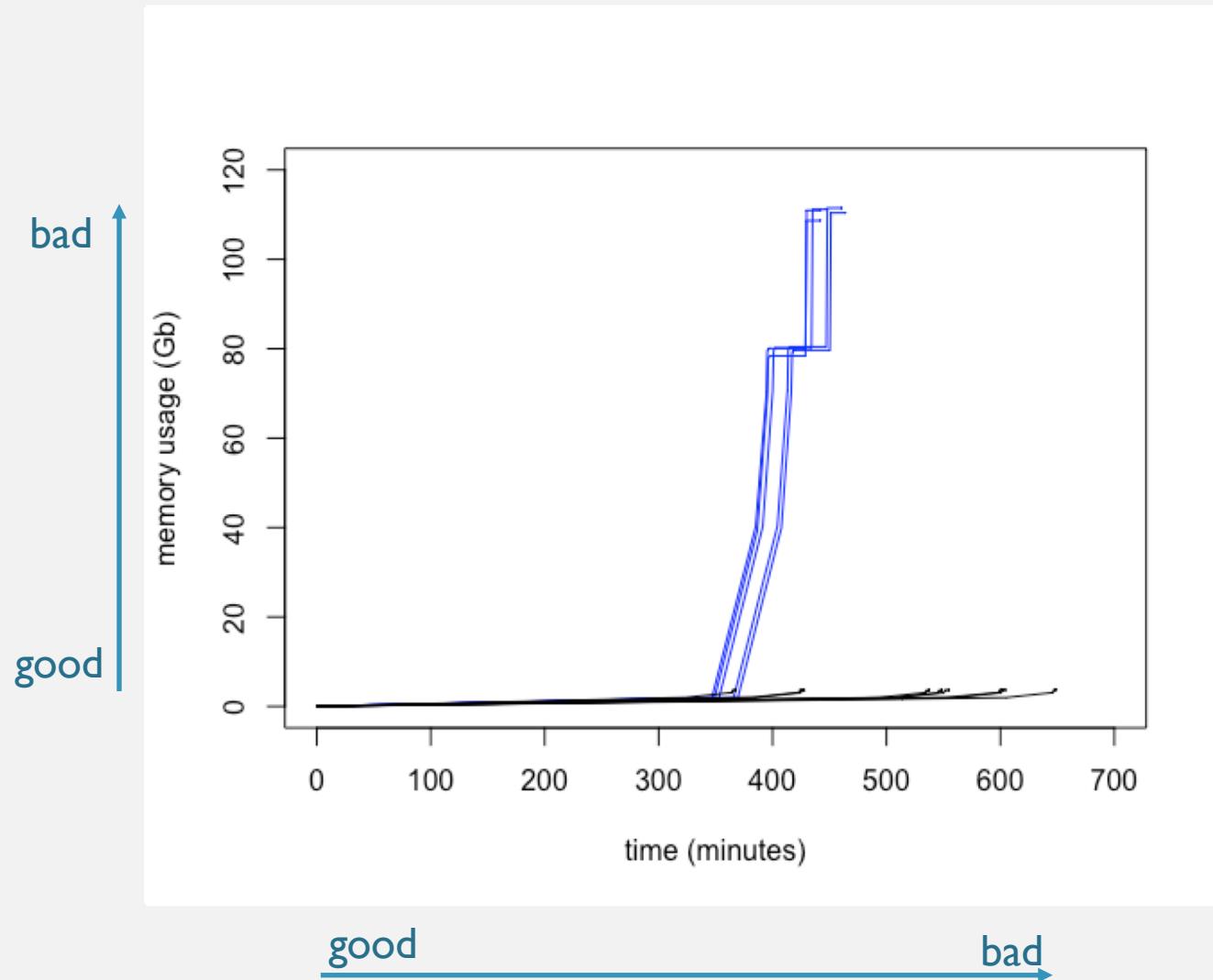
good

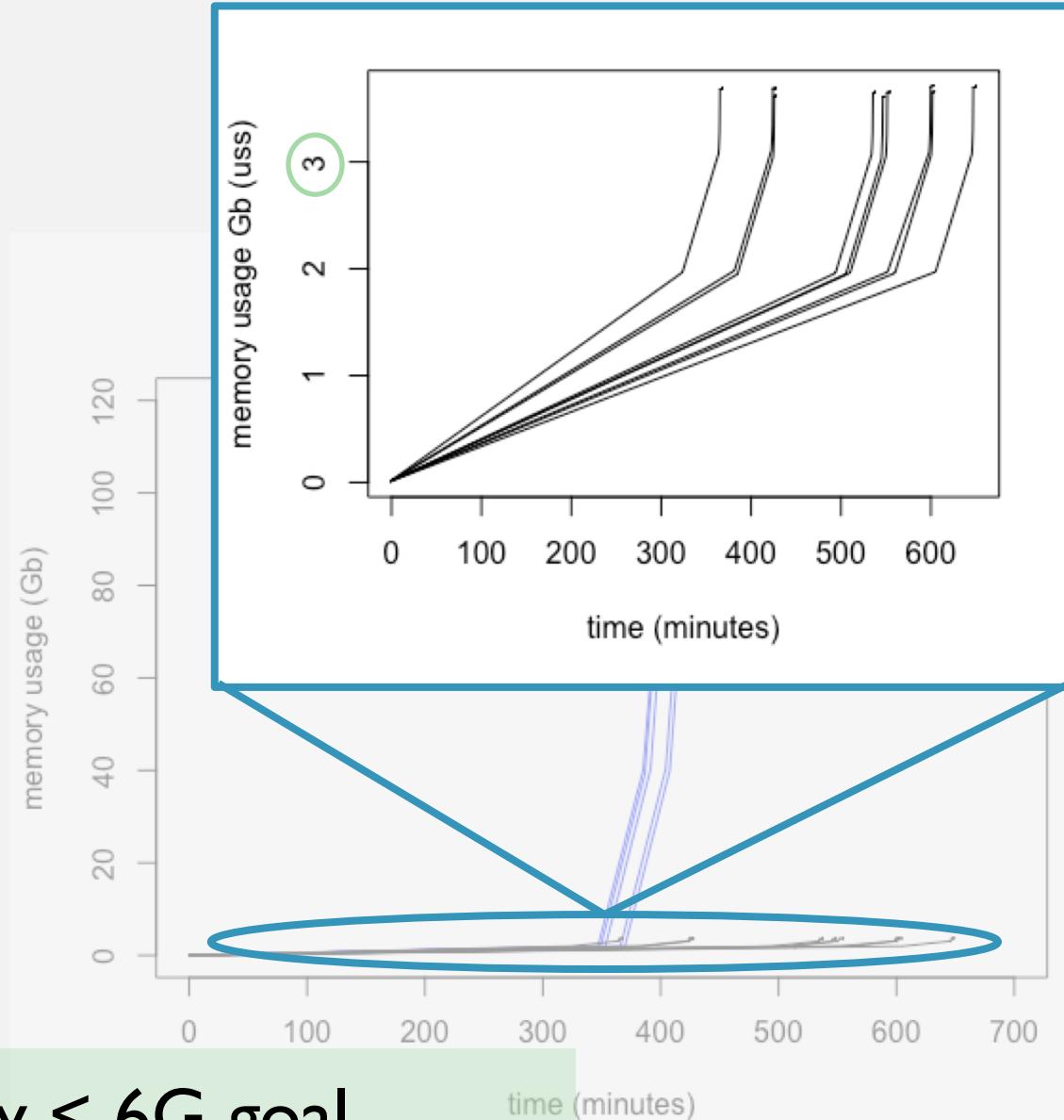
bad

PROFILE OF PYTHON SCRIPT



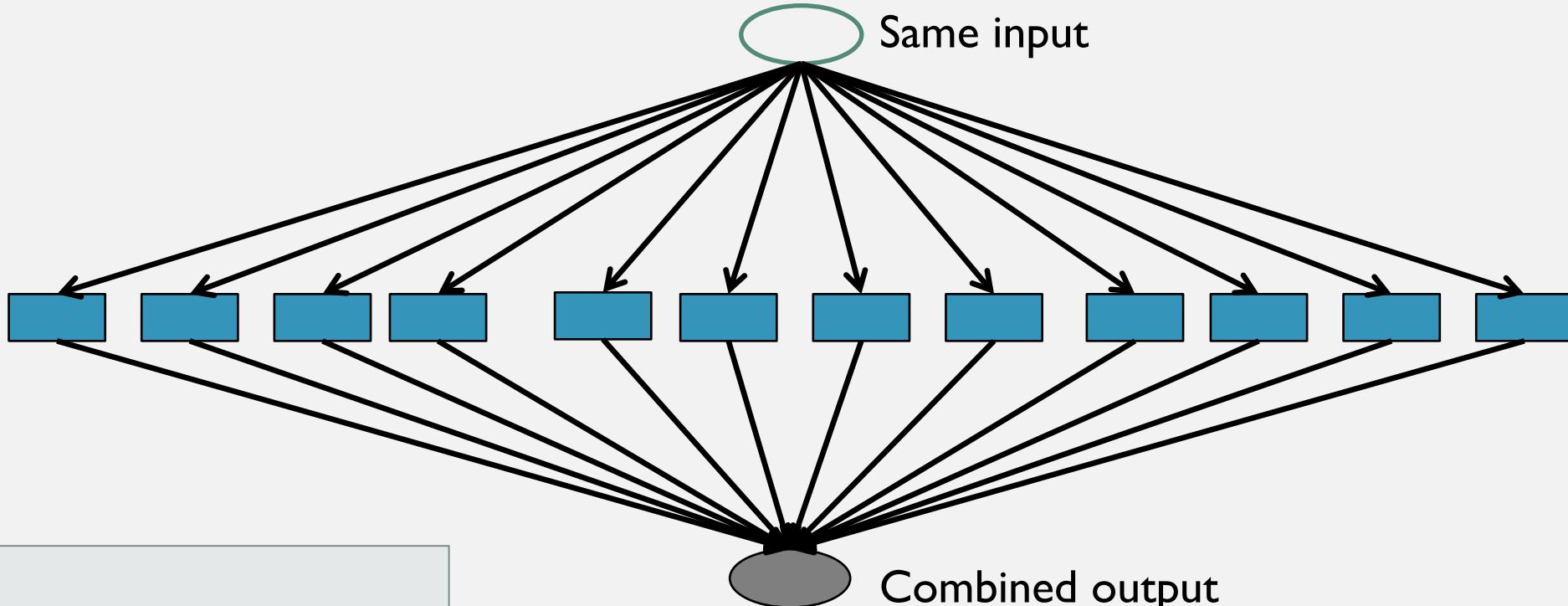
PROFILE OF PYTHON SCRIPT





Max memory < 6G goal
Can now run efficiently in parallel

PLEASANTLY PARALLEL & RESOURCE LIGHT!

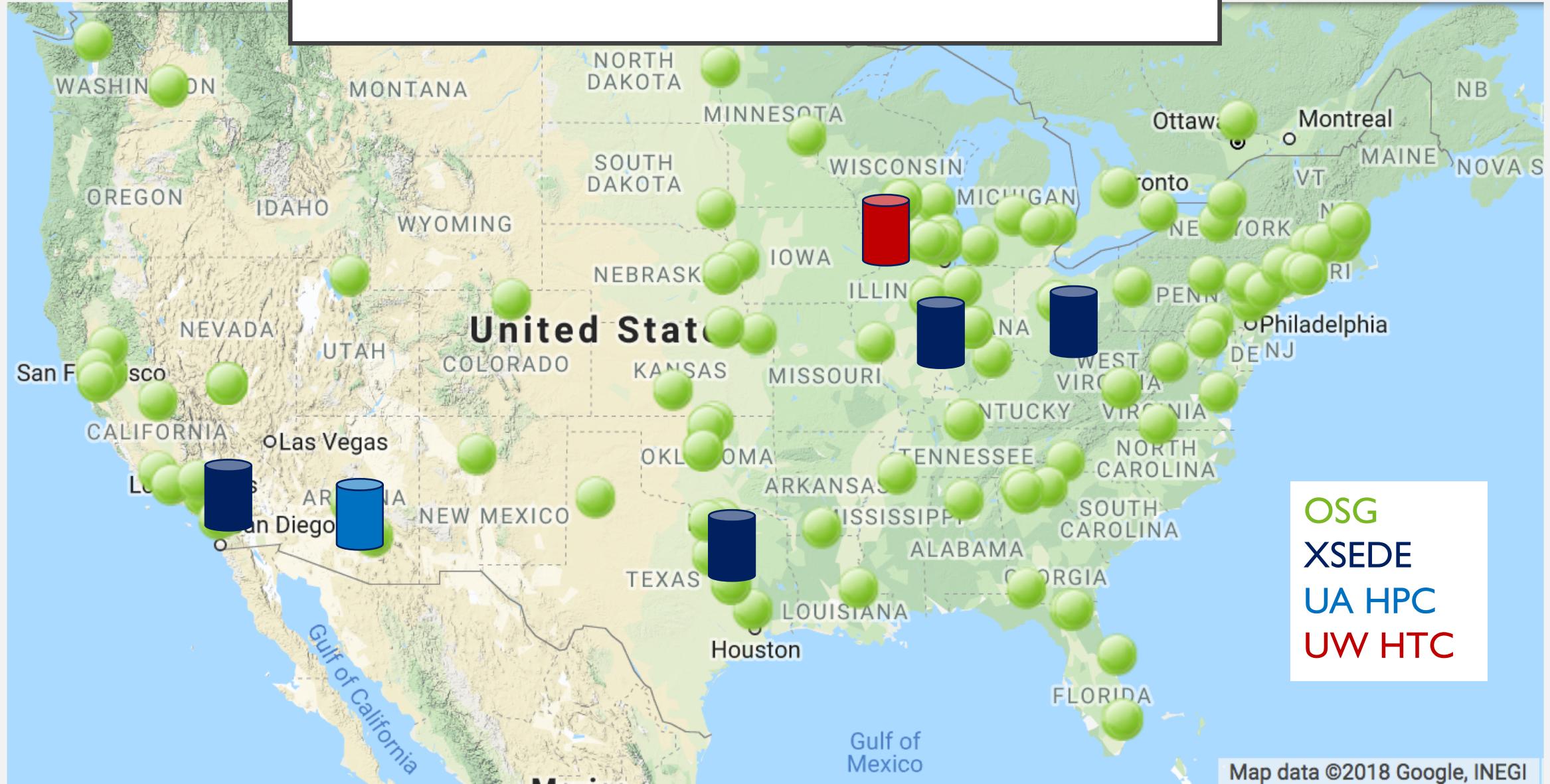


- Each job
 - runs ~40 min, and max 50 hrs
 - Uses ~1G, and max 5G memory
 - Uses ~2M in storage



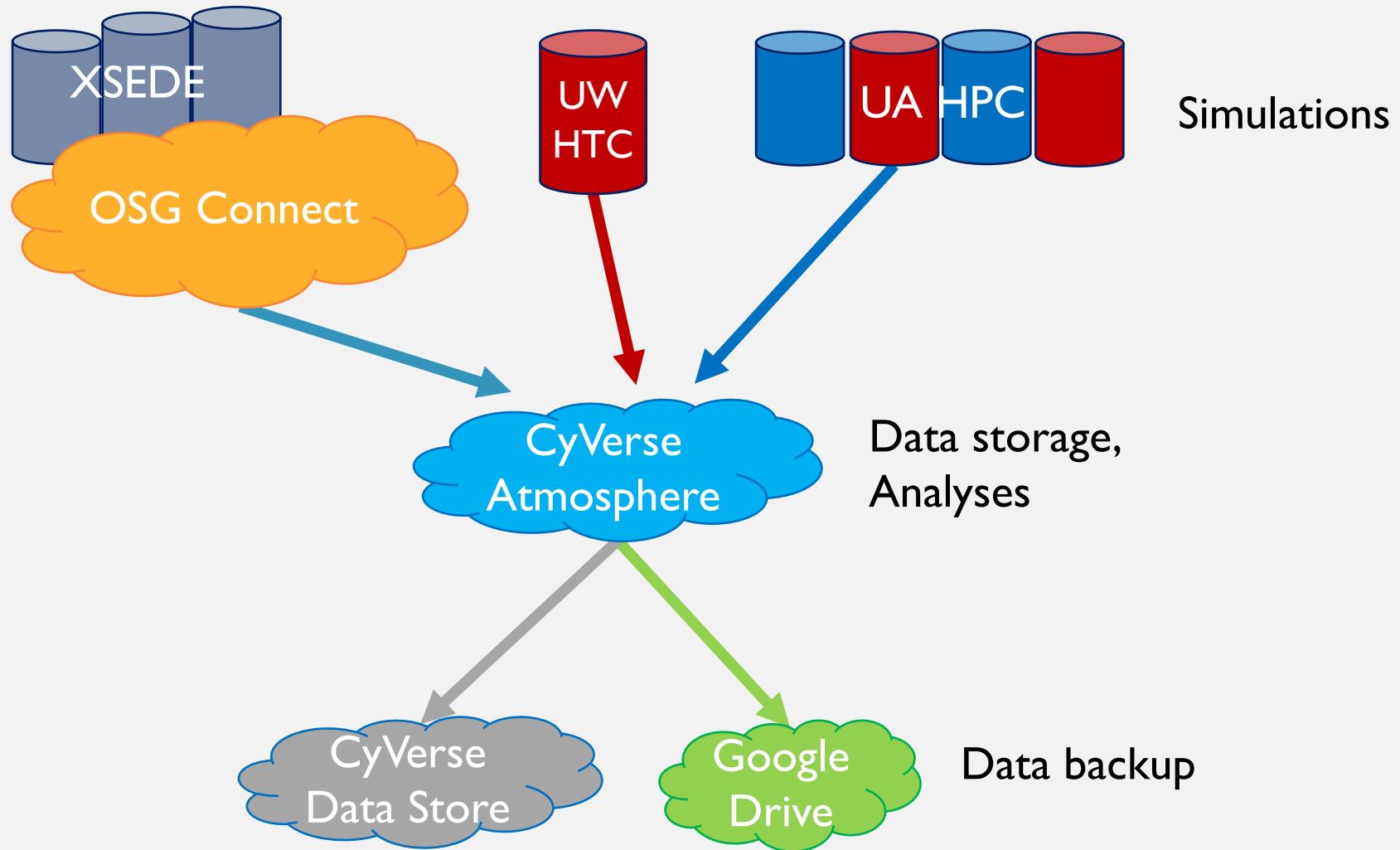
Pegasus

HIGH THROUGHPUT COMPUTING



OSG
XSEDE
UA HPC
UW HTC

SIMULATIONS ON HTC CLUSTERS, ANALYSES ON VM



GENERALIZATION OF CODE AND WORKFLOW

← → ⌂ Secure | <https://agladstein.github.io/SimPrily/>

SimPrily

[Home](#)

[Quick Start](#)

[Documentation](#)

[Citations & Licensing](#)

[Download ZIP](#)

[Download TAR](#)

[View On GitHub](#)

Welcome to SimPrily

SimPrily runs genome simulations with user defined parameters or parameters randomly generated by priors and computes genomic statistics on the simulation output.

- Runs genome simulation with model defined by prior distributions of parameters and demographic model structure.
- Can take into account SNP array ascertainment bias by creating pseudo array based on priors of number of samples of discovery populations and allele frequency cut-off.
- Calculates genomic summary statistics on simulated genomes and pseudo arrays.

This is ideal for use with Approximate Bayesian Computation on whole genome or SNP array data.

Uses c++ programs MaCS and GERMLINE. For more information on these programs, see:

[MaCS Github](#)

[GERMLINE Github](#)

Hosted on [GitHub Pages](#)

Quick Start

Copyright 2018 © SimPrily

To start using right away SimPrily, please visit the [quickstart](#) page.

SIMPRILY HAS UNIQUE FEATURES

Program	Large loci	Priors	Statistics	SNP ascertainment	HTC
SimPrily (2018)	✓	✓	✓	✓	✓
Fastsimcoal2 (2013)	✓	✓			
Msprime (2016)	✓		✓		
BaySICS (2014)		✓	✓		✓
Coala (2016)		✓	✓		
SKELESIM (2017)			✓		

Comparison of SimPrily features with other simulators and wrappers.

POTENTIAL APPLICATIONS OF SIMPRILY

- Simulate genome sequence or SNP array data to
- Test software
- Infer demographic history with Approximate Bayesian Computation
- Use as null model when inferring regions under selection
- Create training and test dataset for machine learning

THEMES OF DISSERTATION

- Detection of runs of homozygosity from SNP arrays
 - Improving identification of runs of homozygosity (Ch. 2)
 - Correcting ascertainment bias in runs of homozygosity (App. C)
- Scaling up Approximate Bayesian Computation for whole chromosomes
 - Create efficient pipeline to simulate demographic models and calculate summary statistics (App. A)
 - Create generalized high throughput workflow (Ch. 4)
- Infer history of the Ashkenazi Jews
 - Substructure in AJ? (Ch. 5)
 - Khazarian origin? (App. B)

DATASET

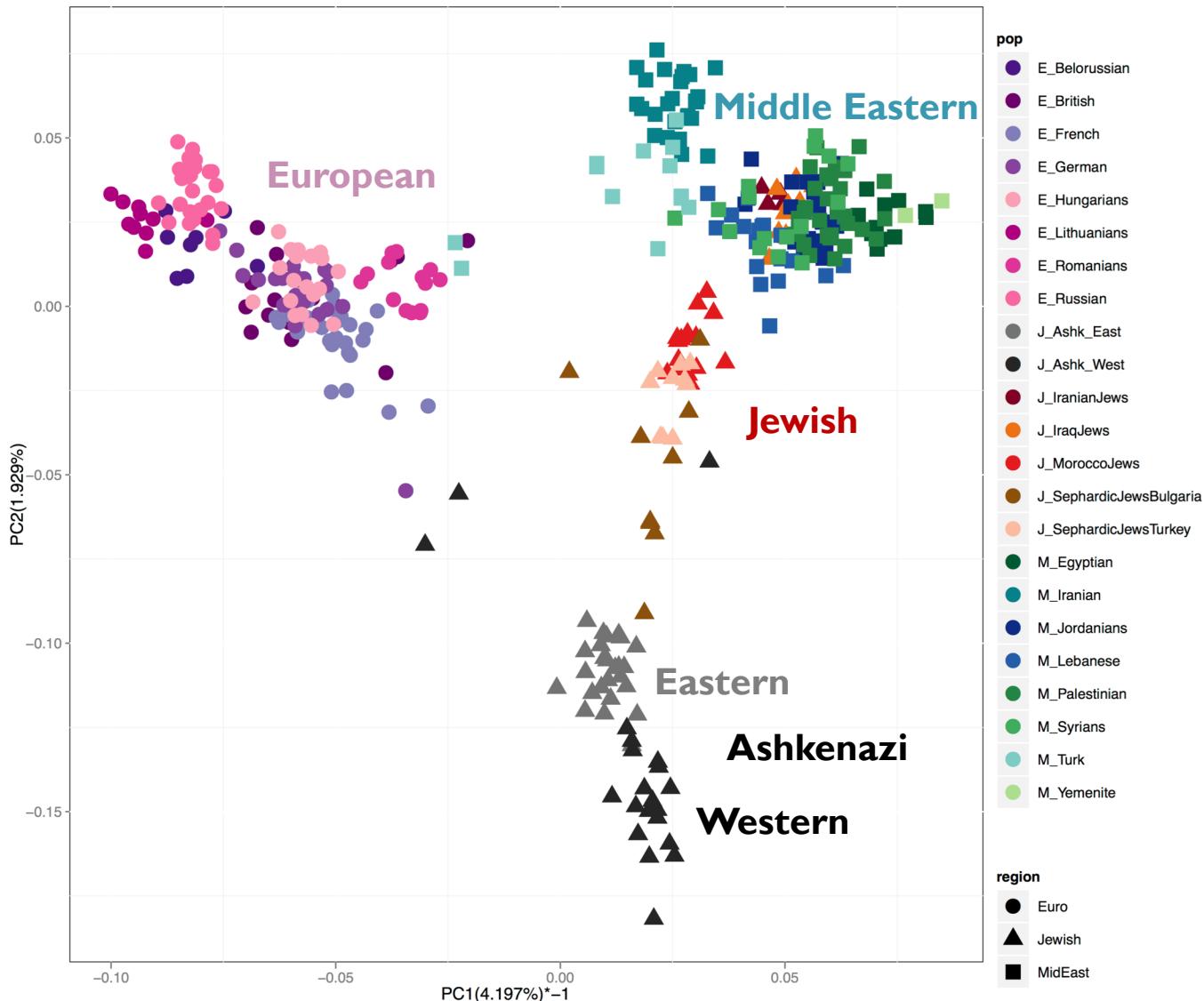
SNP array data	Sample Size	Source
Eastern Ashkenazi	239	Family Tree DNA, Behar et al. 2010
Western Ashkenazi	19	Family Tree DNA, Behar et al. 2010
Jewish (9 pops)	79	Behar et al. 2010
Middle Eastern (11 pops)	211	Behar et al. 2010, Hammer, HGDP
European (8 pops)	139	Behar et al. 2010, Hammer, HGDP

Whole genome data	Sample Size	Source
Ashkenazi	230	Carmi et al. 2014, Hammer Lab
European, African, Asian, American		CGI, 1000 Genomes

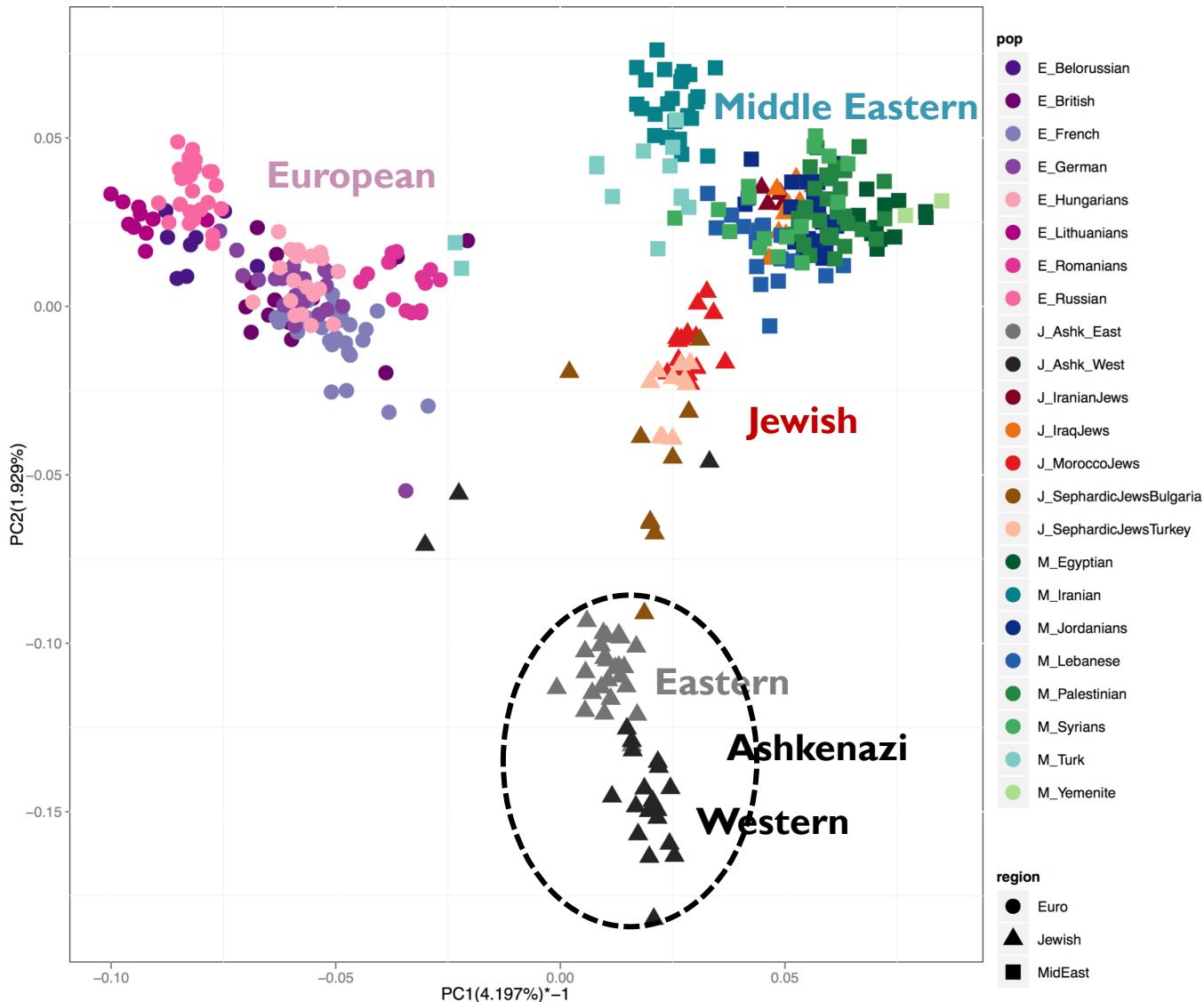
AJ GENETIC RELATIONSHIP TO MIDDLE EASTERN AND EUROPEAN POPULATIONS

- **Principal Component Analysis (PCA)** – a visualization of population genetic structure
- **ADMIXTURE** – visualization of population genetic structure

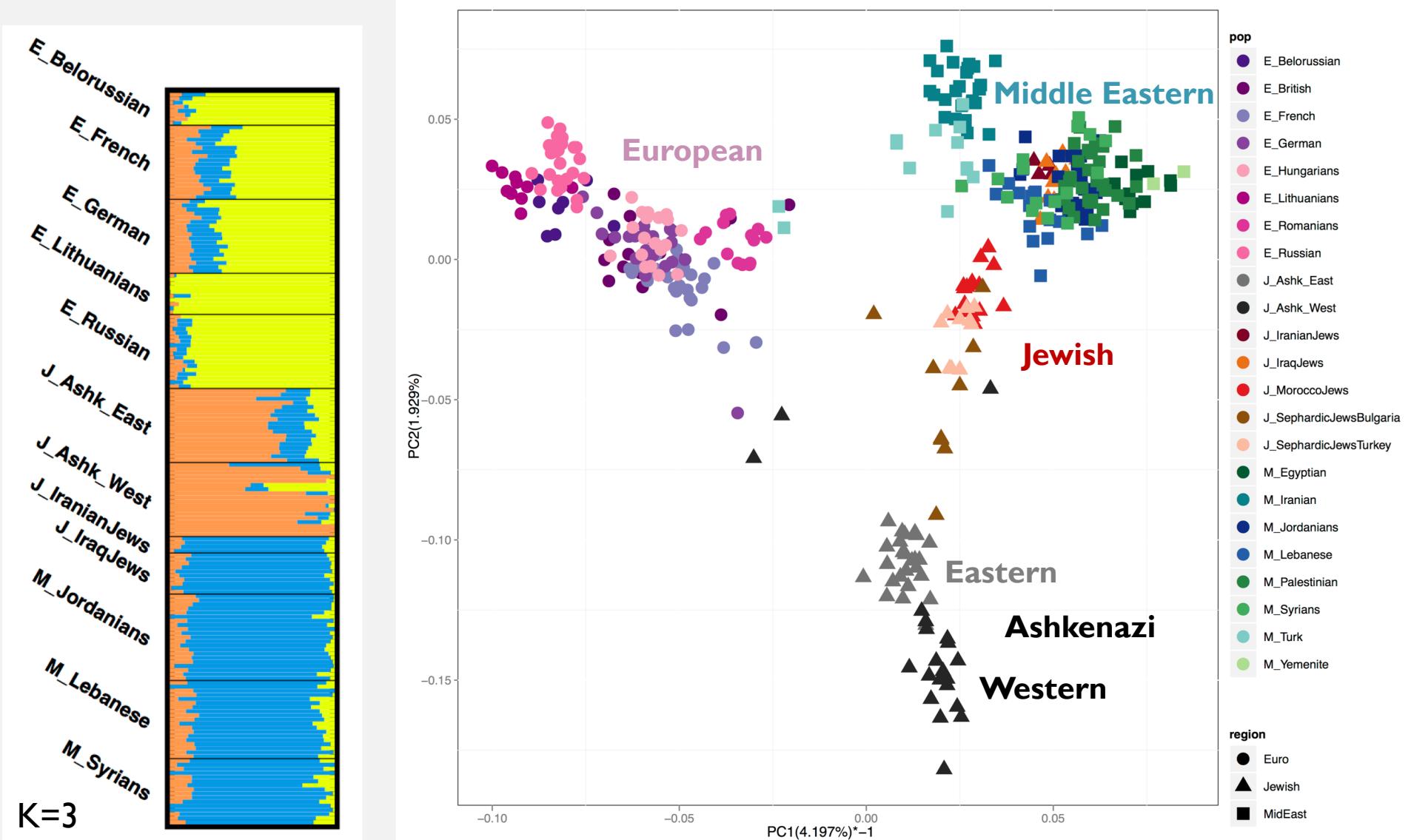
AJ GENETIC RELATIONSHIP TO MIDDLE EASTERN AND EUROPEAN POPULATIONS



AJ GENETIC RELATIONSHIP TO MIDDLE EASTERN AND EUROPEAN POPULATIONS

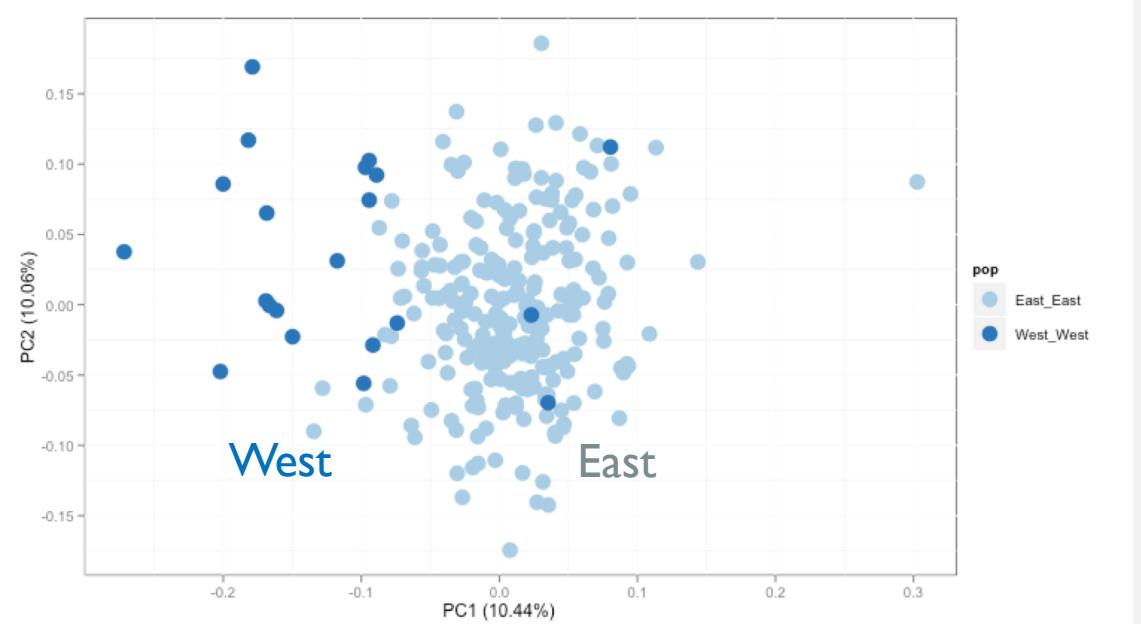


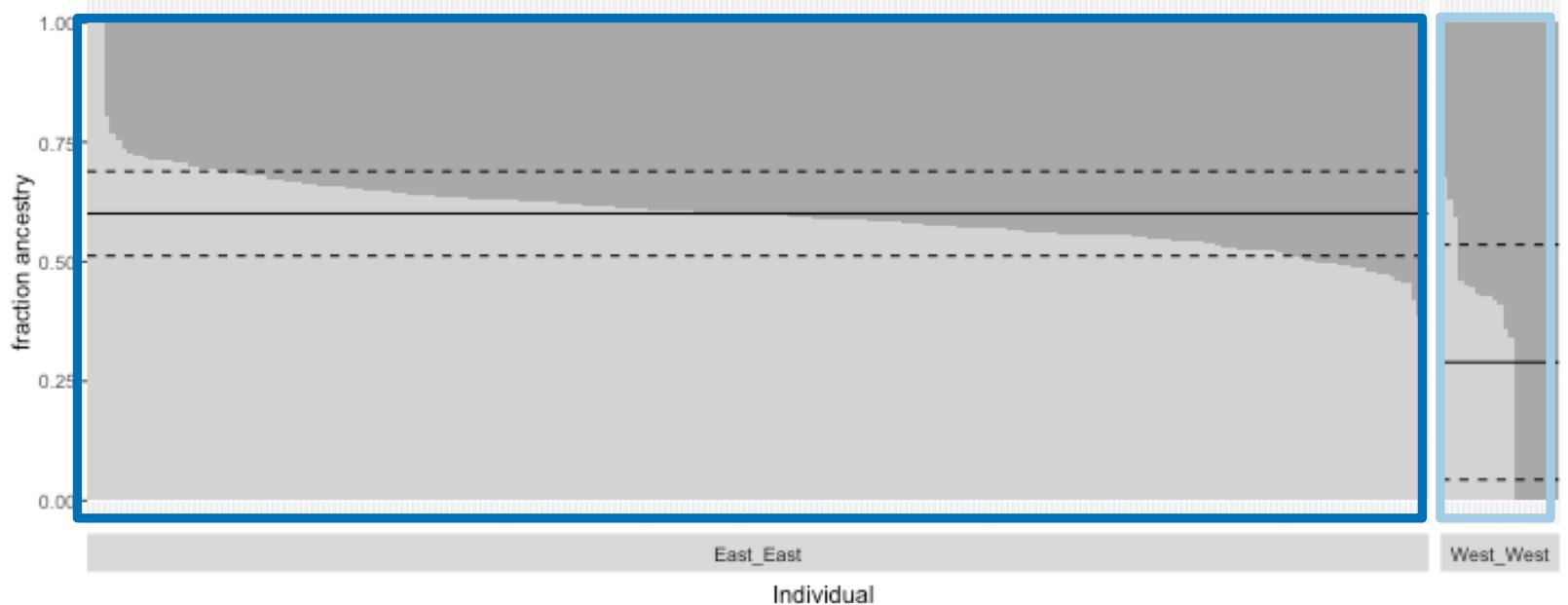
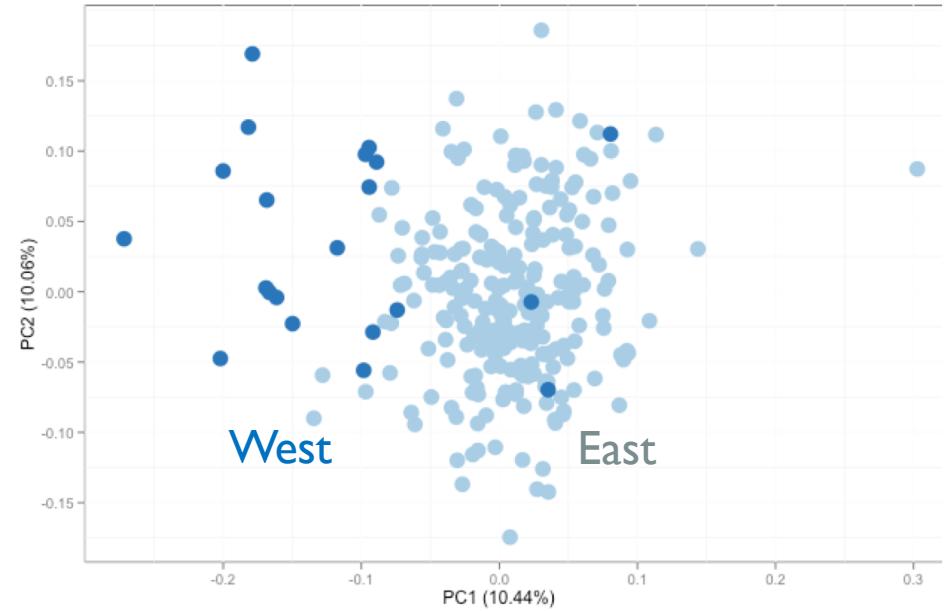
AJ GENETIC RELATIONSHIP TO MIDDLE EASTERN AND EUROPEAN POPULATIONS



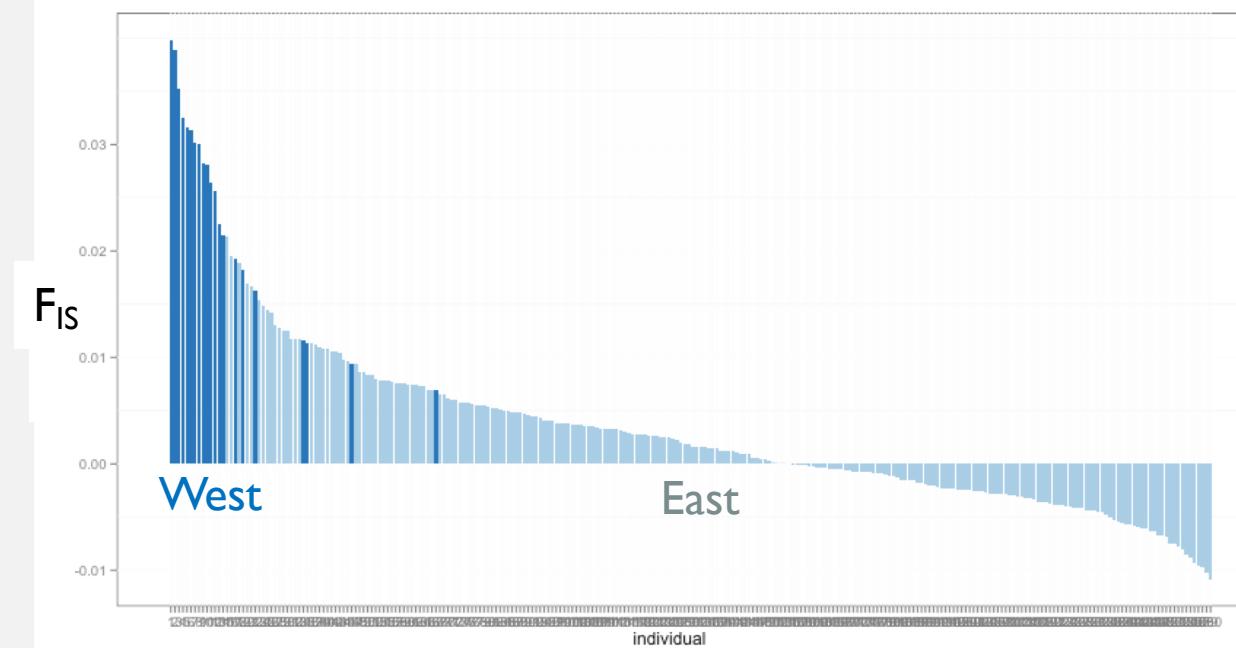
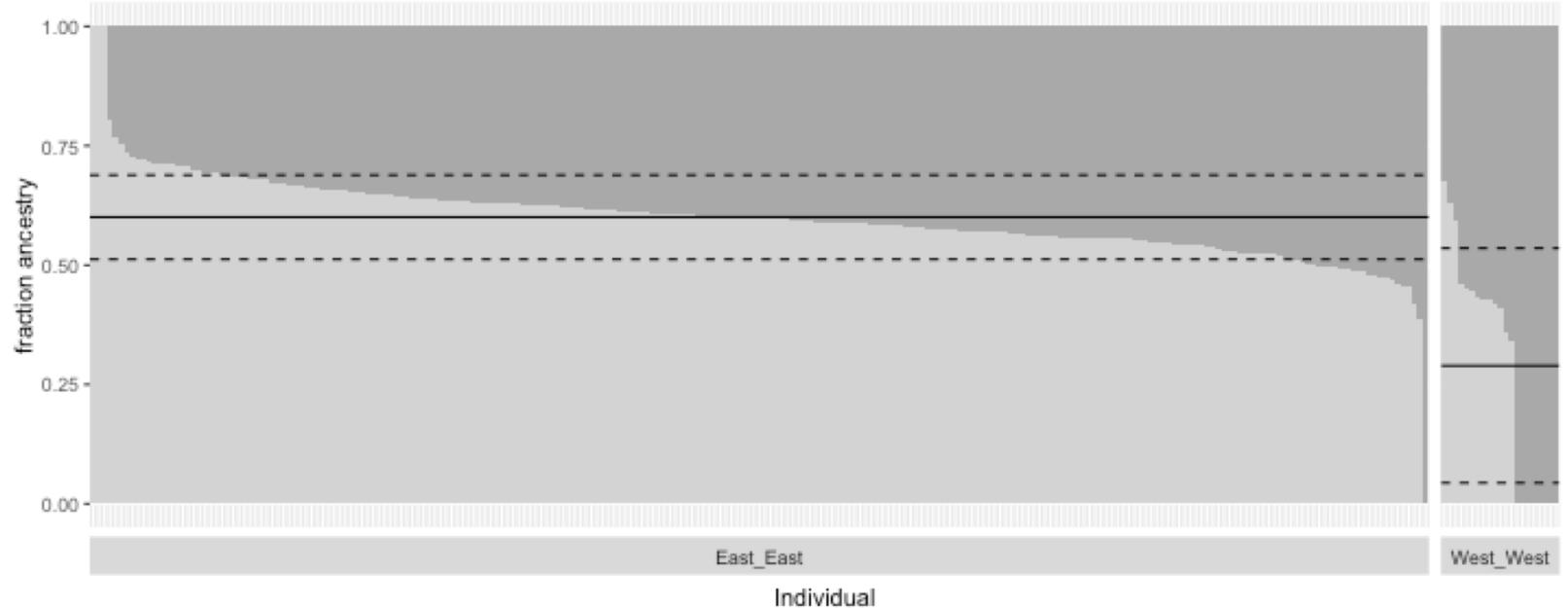
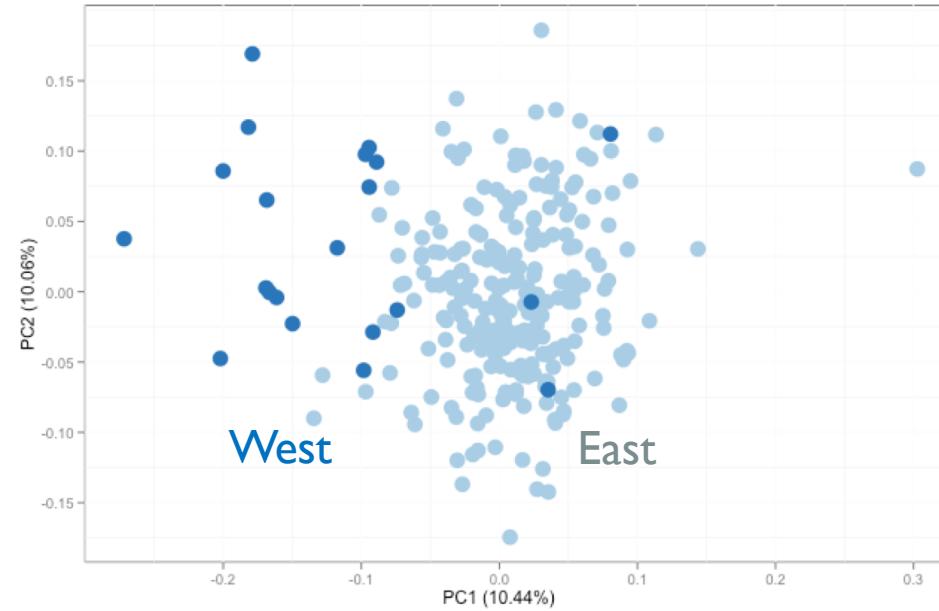
DIFFERENCE BETWEEN EASTERN AND WESTERN?

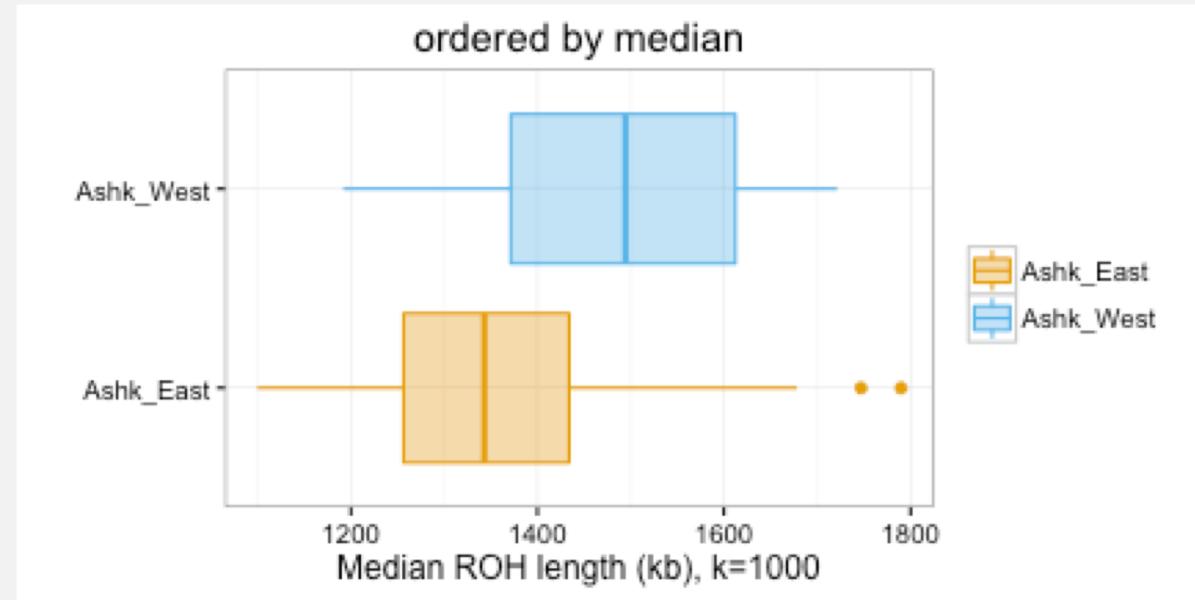
- **Principal Component Analysis (PCA)** – a visualization of population genetic structure
- **ADMIXTURE** – visualization of population genetic structure
- **Runs of homozygosity** – indicates levels of inbreeding or small effective population size
- **Identity by Descent (IBD)** – Indicates shared ancestry between individuals

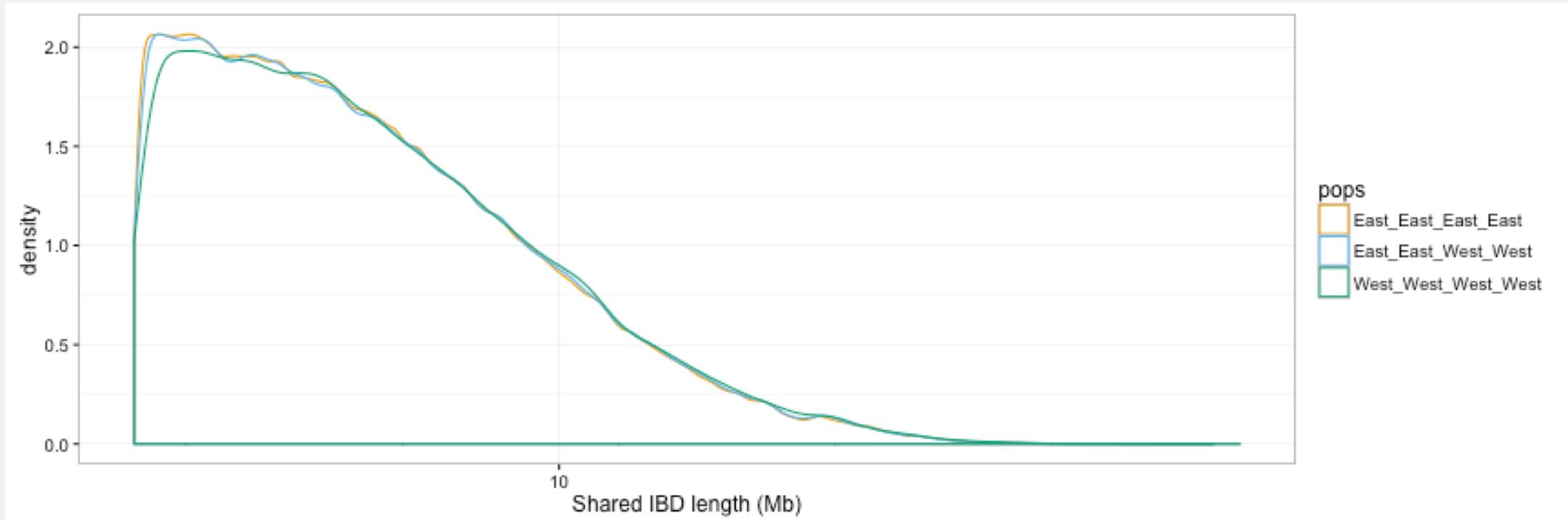
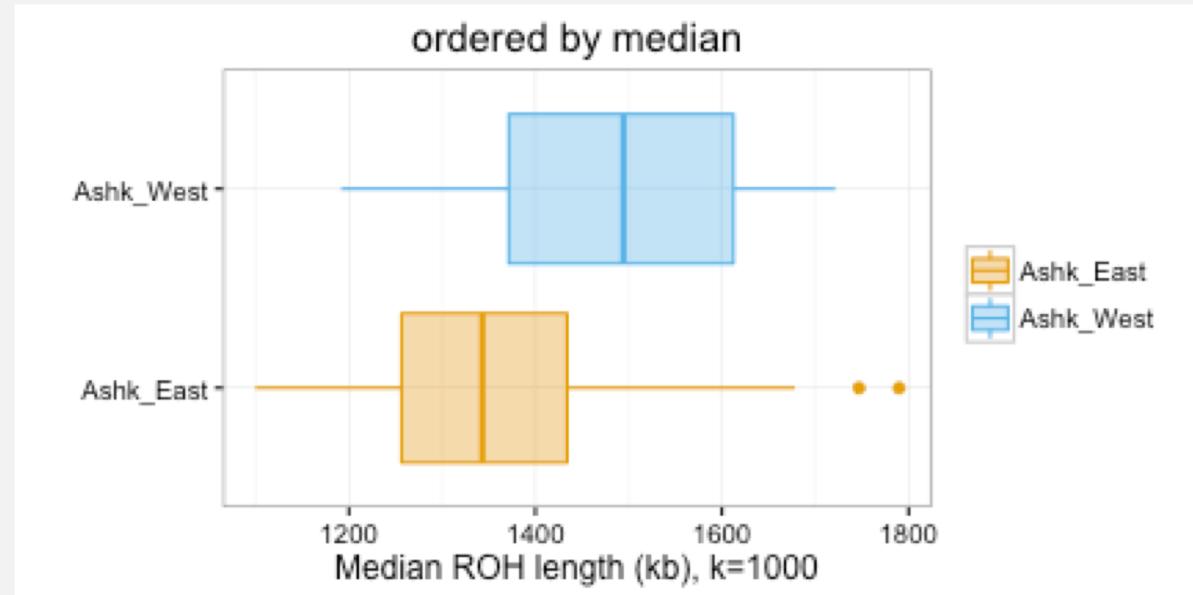


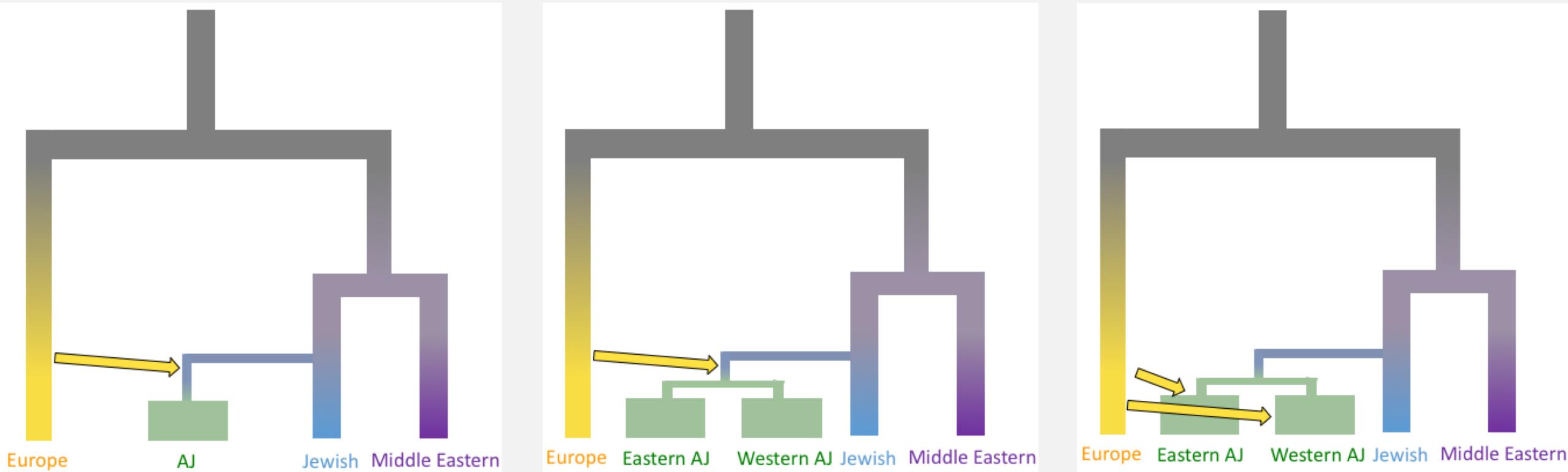


West_West



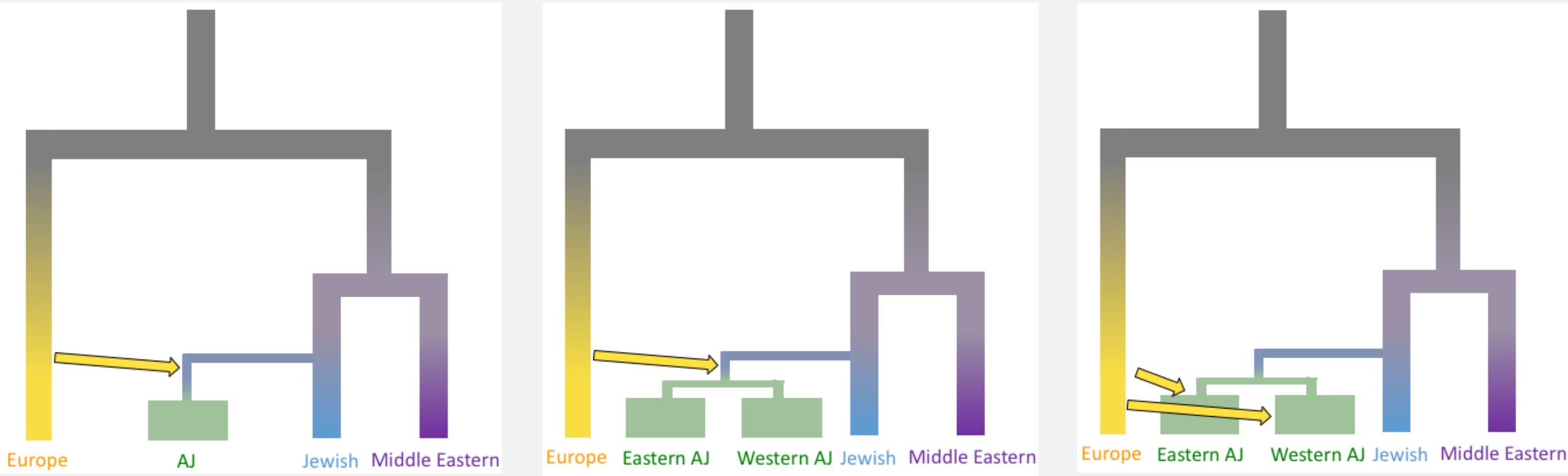






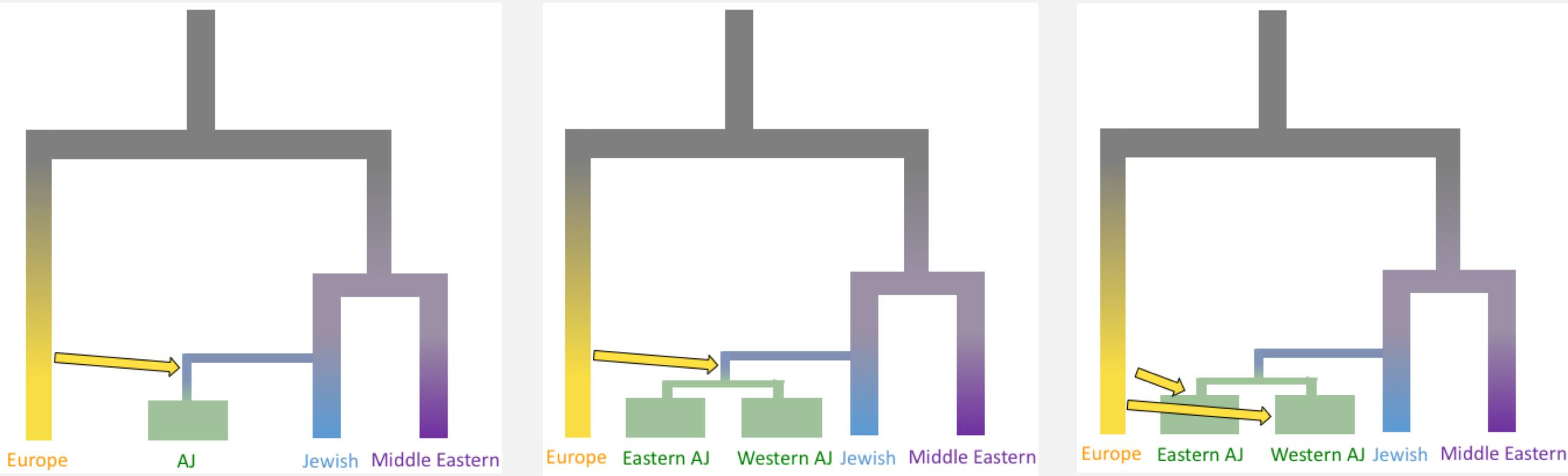
MODELS OF ASHKENAZI HISTORY

APPROXIMATE BAYESIAN COMPUTATION (ABC)



MODELS OF ASHKENAZI HISTORY

APPROXIMATE BAYESIAN COMPUTATION (ABC)



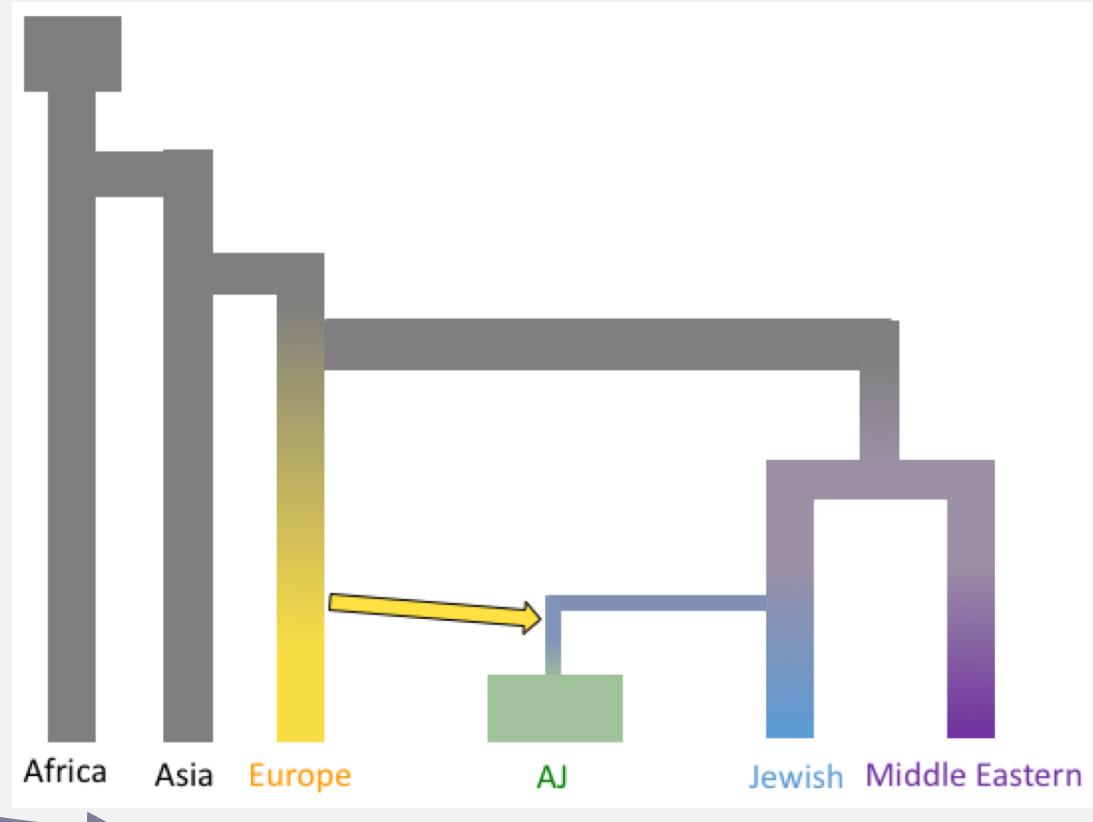
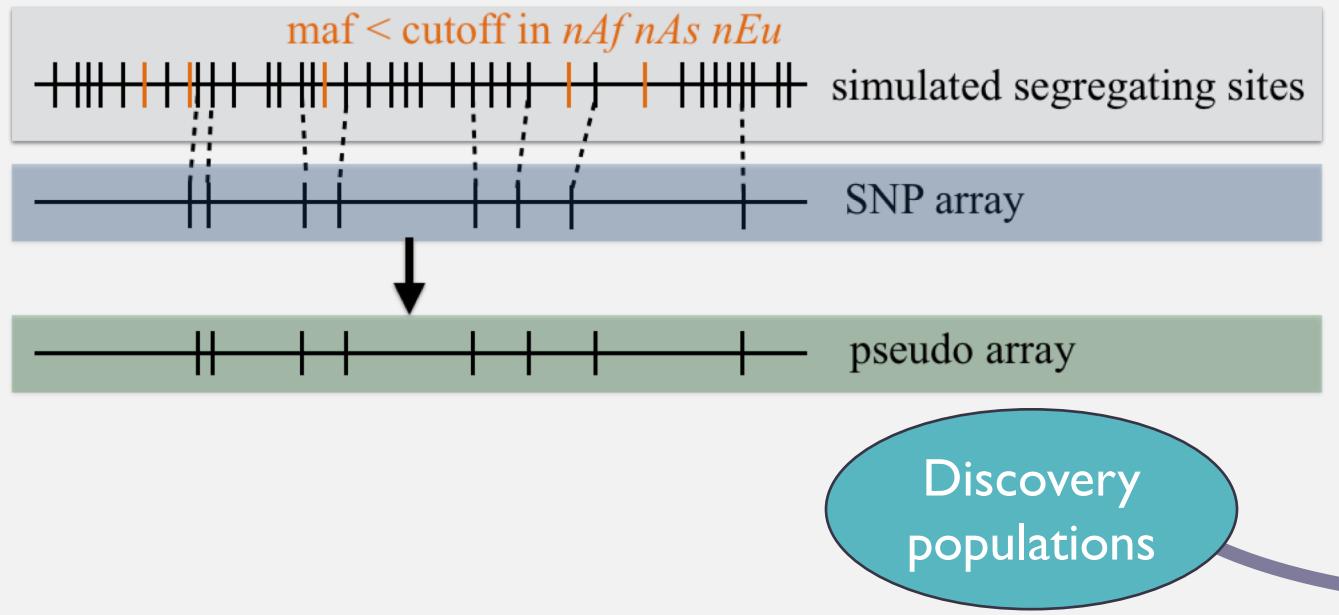
MODELS OF ASHKENAZI HISTORY

APPROXIMATE BAYESIAN COMPUTATION (ABC)

MODEL ASCERTAINMENT BIAS

Ascertainment parameters:

- Sample sizes of discovery populations
- Minor allele frequency cutoff



ABC WORKFLOW

Model Choice

Parameter Estimation

ABC WORKFLOW

Model Choice

Parameter Estimation

Simulate chrl $\sim 1 \times 10^6$
times for each model

ABC WORKFLOW

Model Choice

Simulate chrl $\sim 1 \times 10^6$
times for each model



Parameter Estimation

Find the best stats for
model choice

ABC WORKFLOW

Model Choice

Simulate chr1 $\sim 1 \times 10^6$
times for each model



Find the best stats for
model choice

ABCtoolbox Greedy search algorithm:

1. For all pairs of stats, evaluate the power to distinguish the models, and retain best 10 pairs,
 2. Repeat with triplets,
 3. And so forth until the set of best combinations does not change anymore.
- 100,000 simulations, 1000 retained, 100 cross validations.*
- Keep 77 combinations of stats with power > 0.5,
total of 20 stats

ABC WORKFLOW

Model Choice

Simulate chrl $\sim 1 \times 10^6$
times for each model

Parameter Estimation

Find the best stats for
model choice

Choose best model with
ABC

ABC WORKFLOW

Model Choice

Simulate chrl $\sim 1 \times 10^6$
times for each model

Find the best stats for
model choice

Choose best model with
ABC

Parameter Estimation

- 9, 11, 13 parameters
- $\sim 1 \times 10^6$ simulations
- 1000 retained
- 1000 cross validation

ABC WORKFLOW

Model Choice

Simulate chr1 $\sim 1 \times 10^6$ times for each model



Find the best stats for model choice



Choose best model with ABC

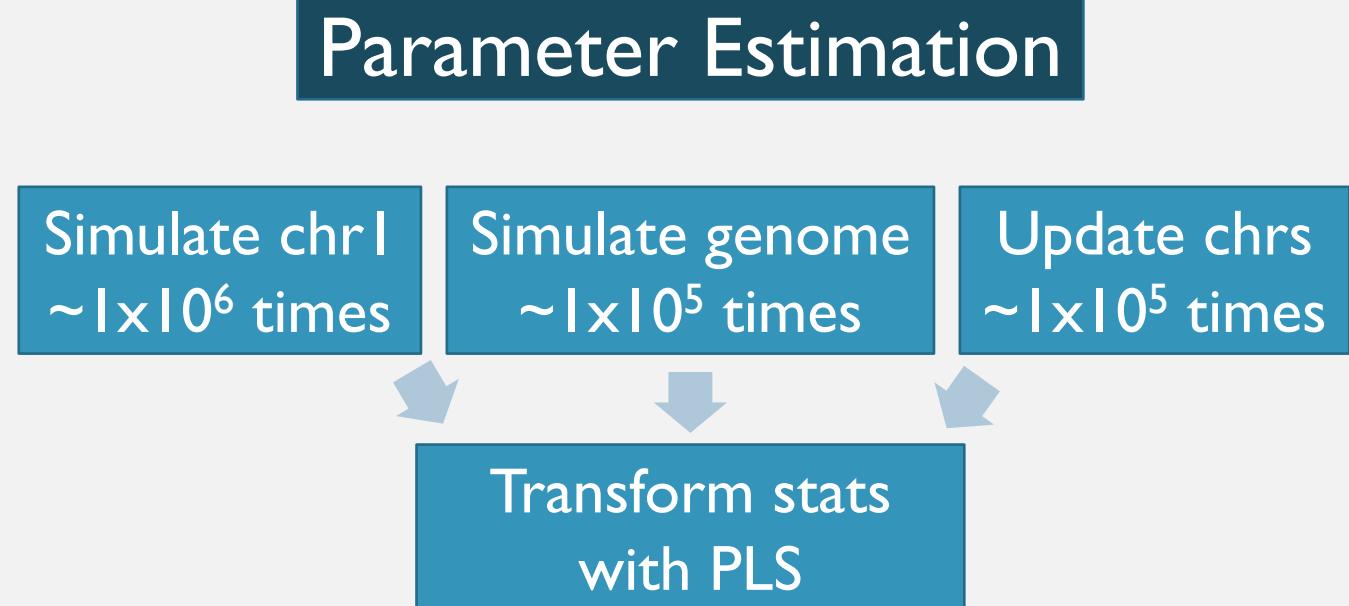
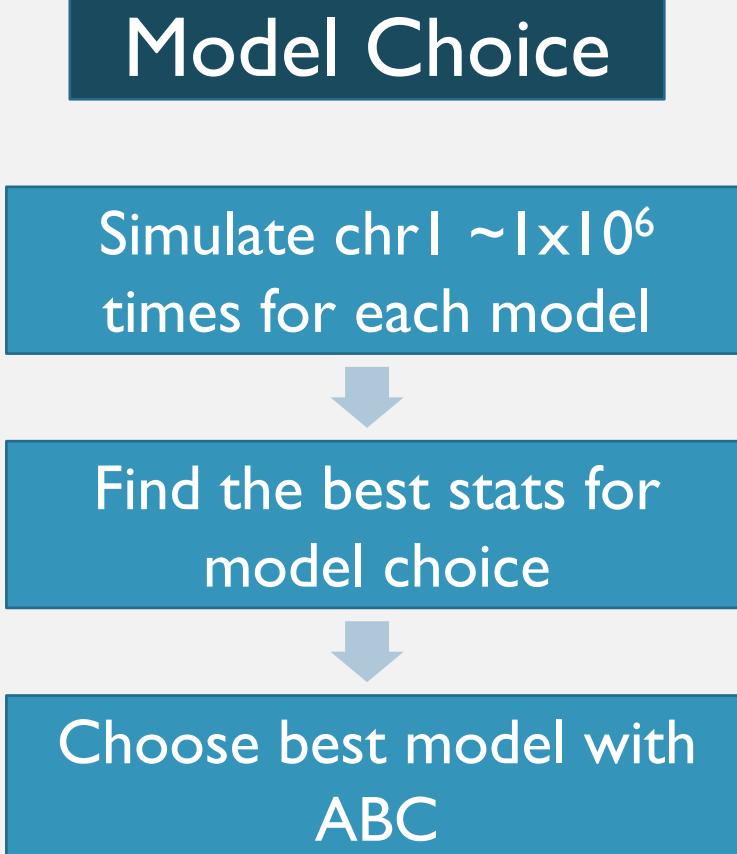
Parameter Estimation

Simulate chr1 $\sim 1 \times 10^6$ times

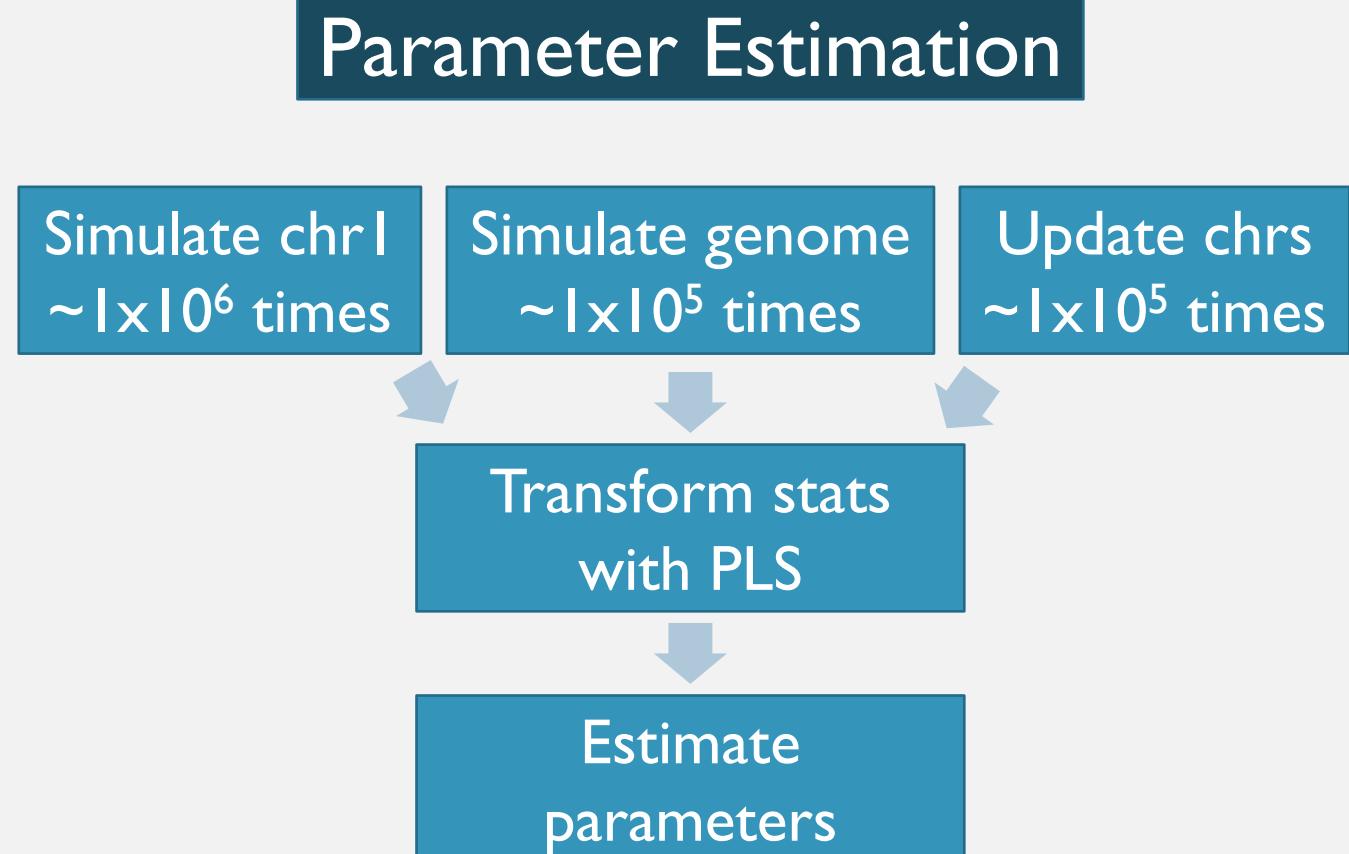
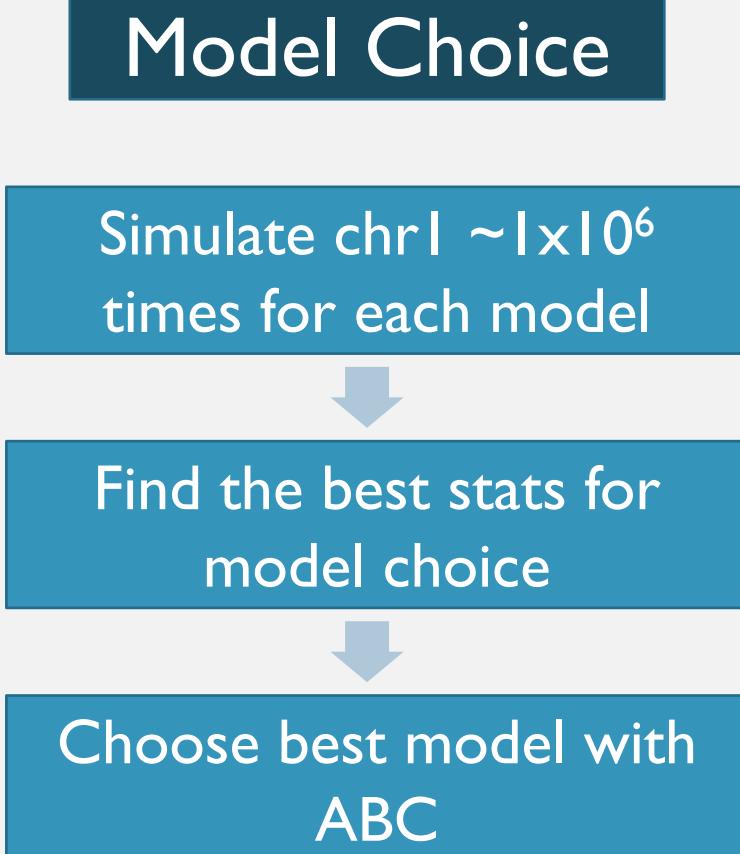
Simulate genome $\sim 1 \times 10^5$ times

Update chrs $\sim 1 \times 10^5$ times

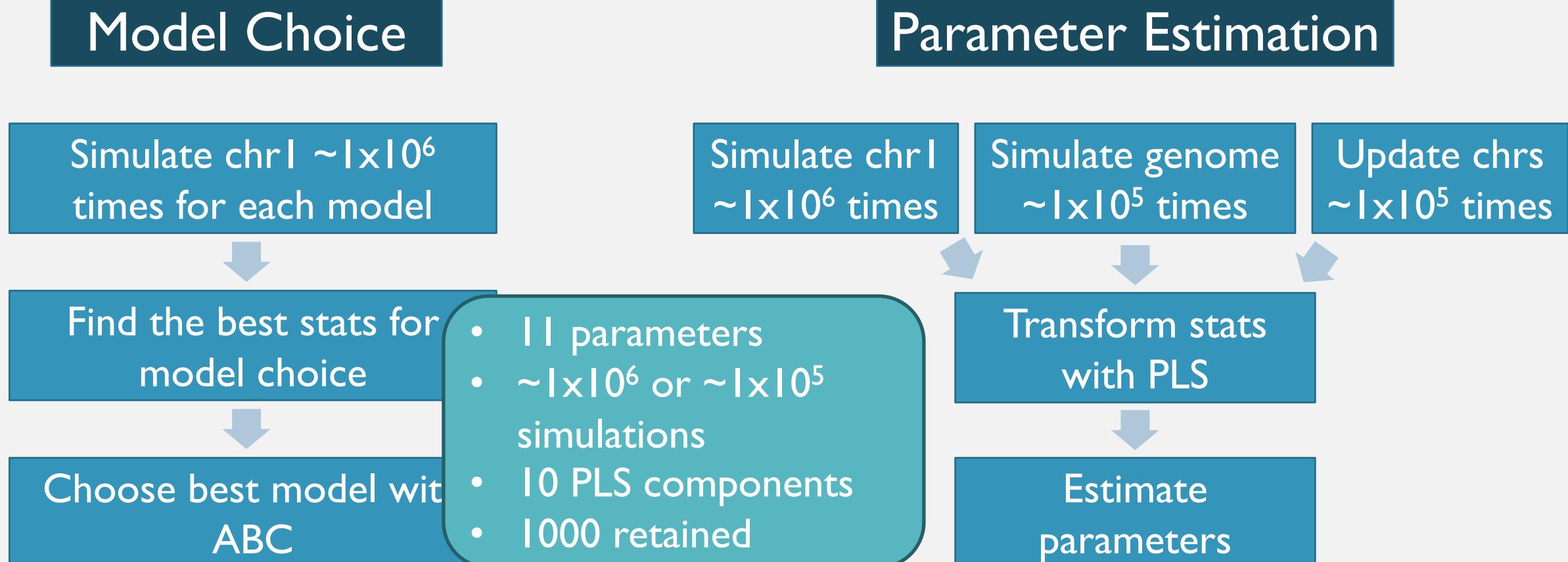
ABC WORKFLOW



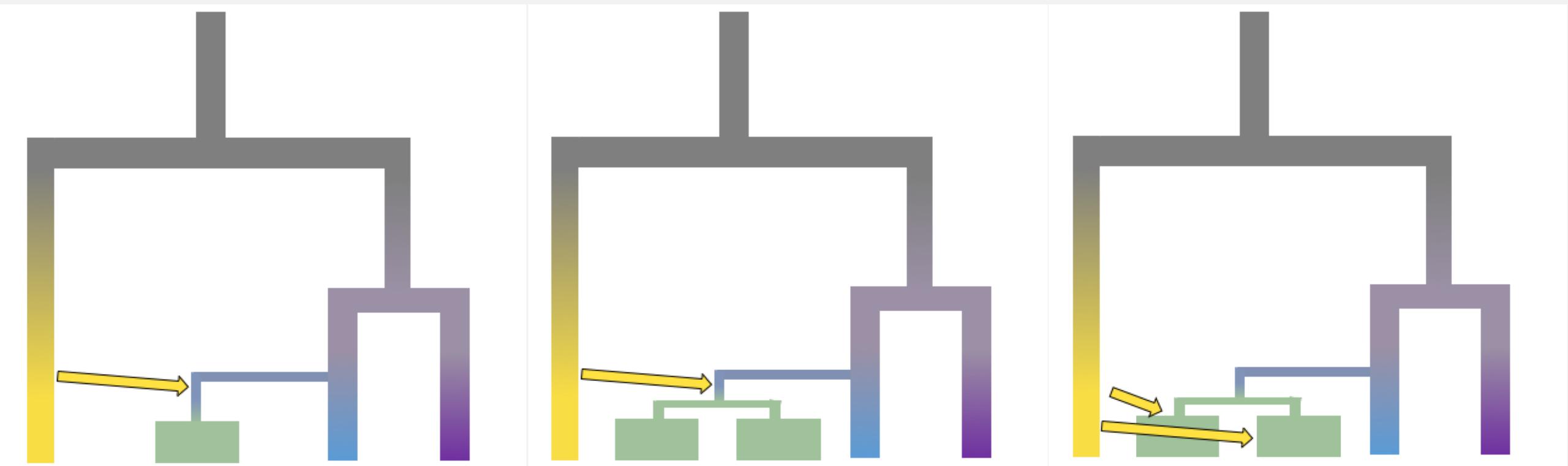
ABC WORKFLOW



ABC WORKFLOW



MODEL CHOICE



Posterior probability: 0.0005

Bayes Factor: 0.0005

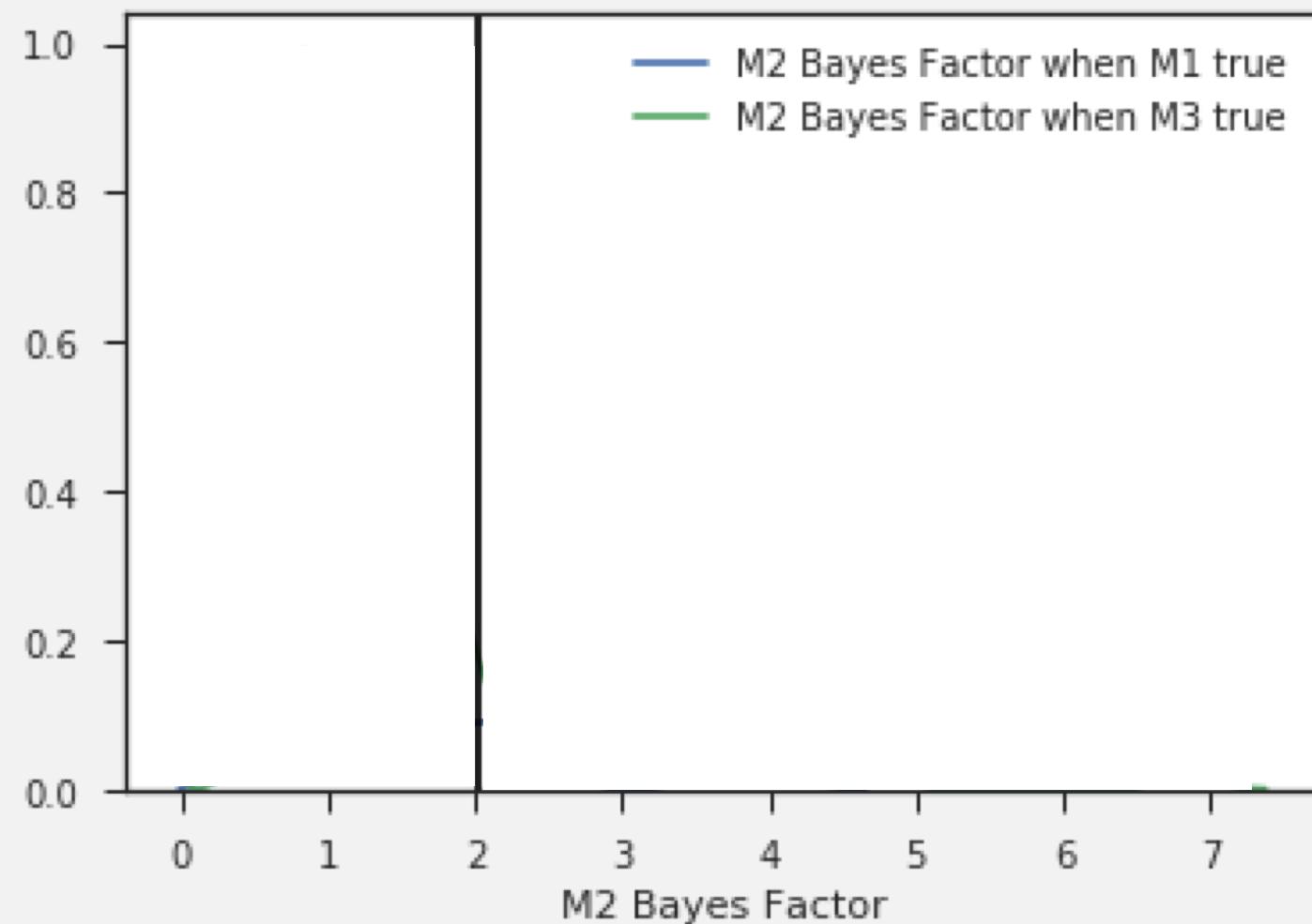
0.67

2.01

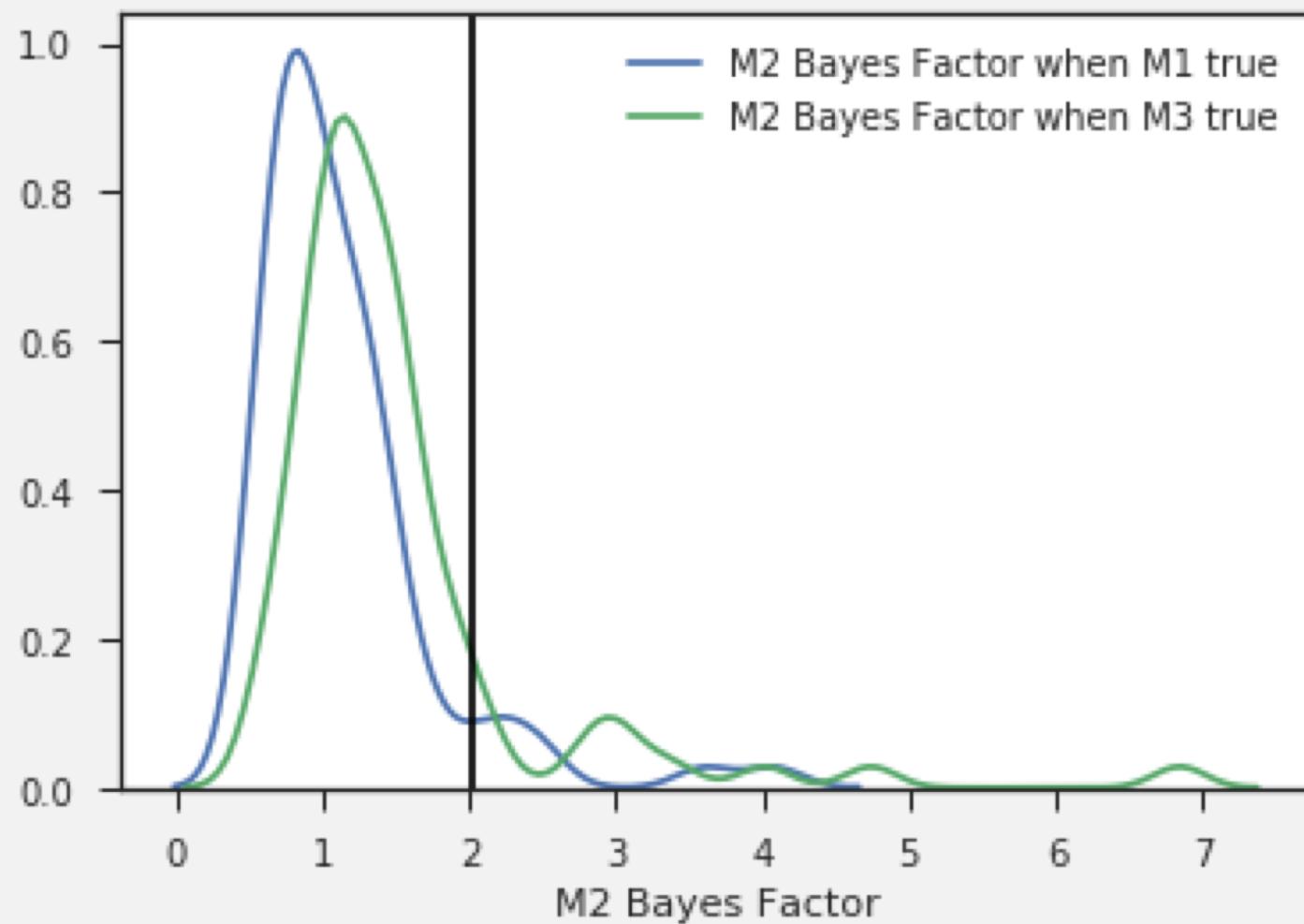
0.33

0.50

CROSS VALIDATION OF MODEL CHOICE



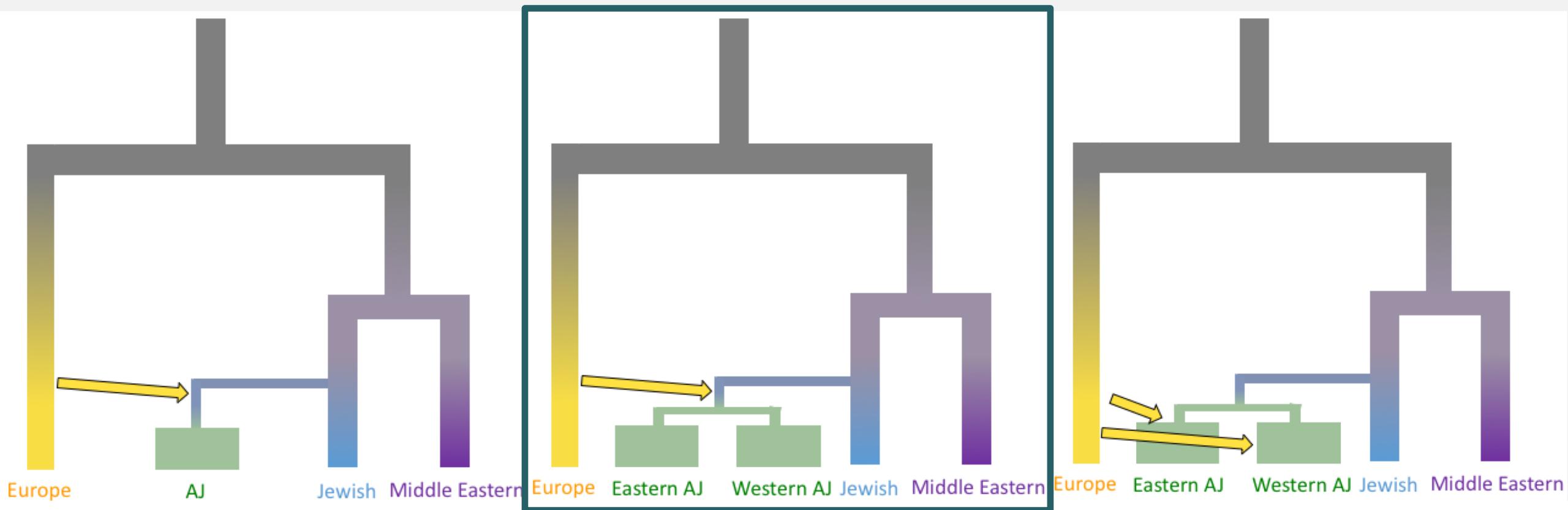
CROSS VALIDATION OF MODEL CHOICE



0.92 Model 2 Bayes factors greater than Model 1 Bayes factors when Model 1 is true model

0.86 Model 2 Bayes factors greater than Model 3 Bayes factors when Model 3 is true model

MODEL CHOICE



Posterior probability: 0.0005

Bayes Factor: 0.0005

Prob. false negative: 0.08

0.67

2.01

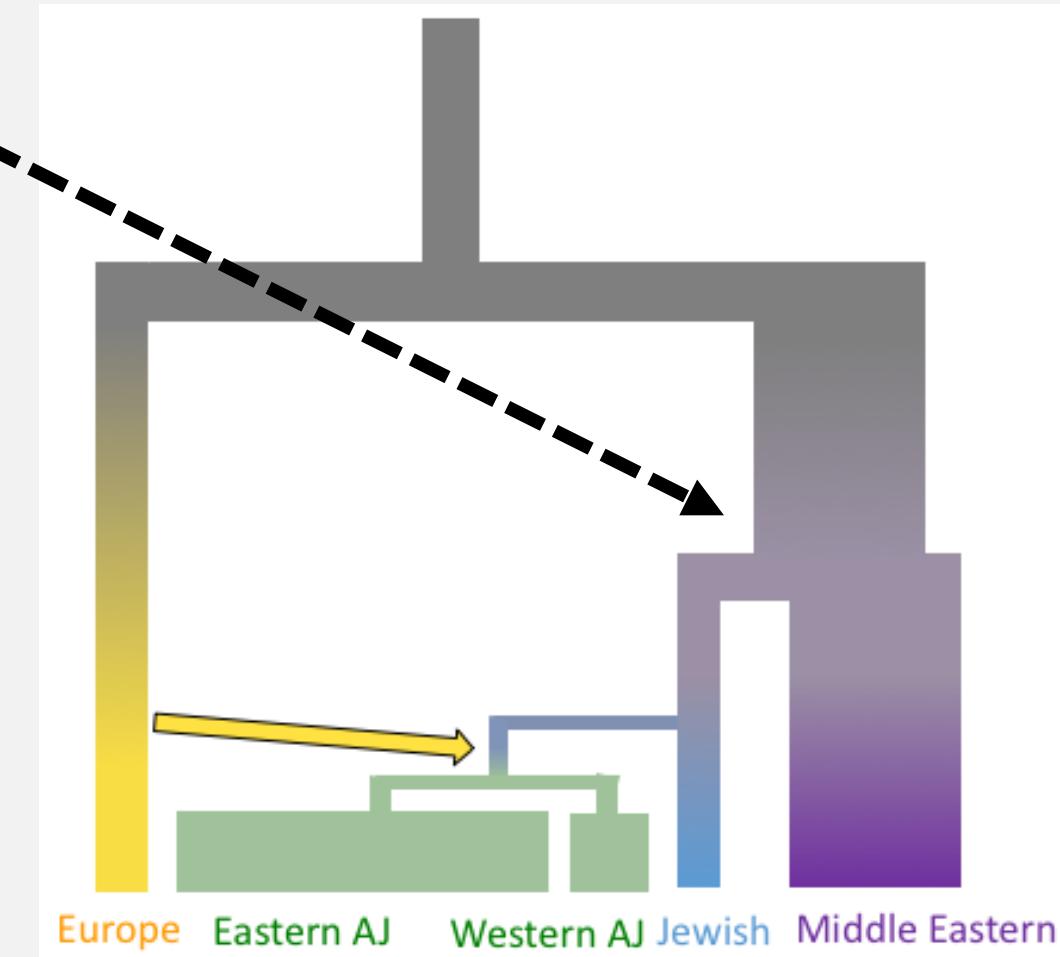
0.33

0.50

0.14

BEST MODEL

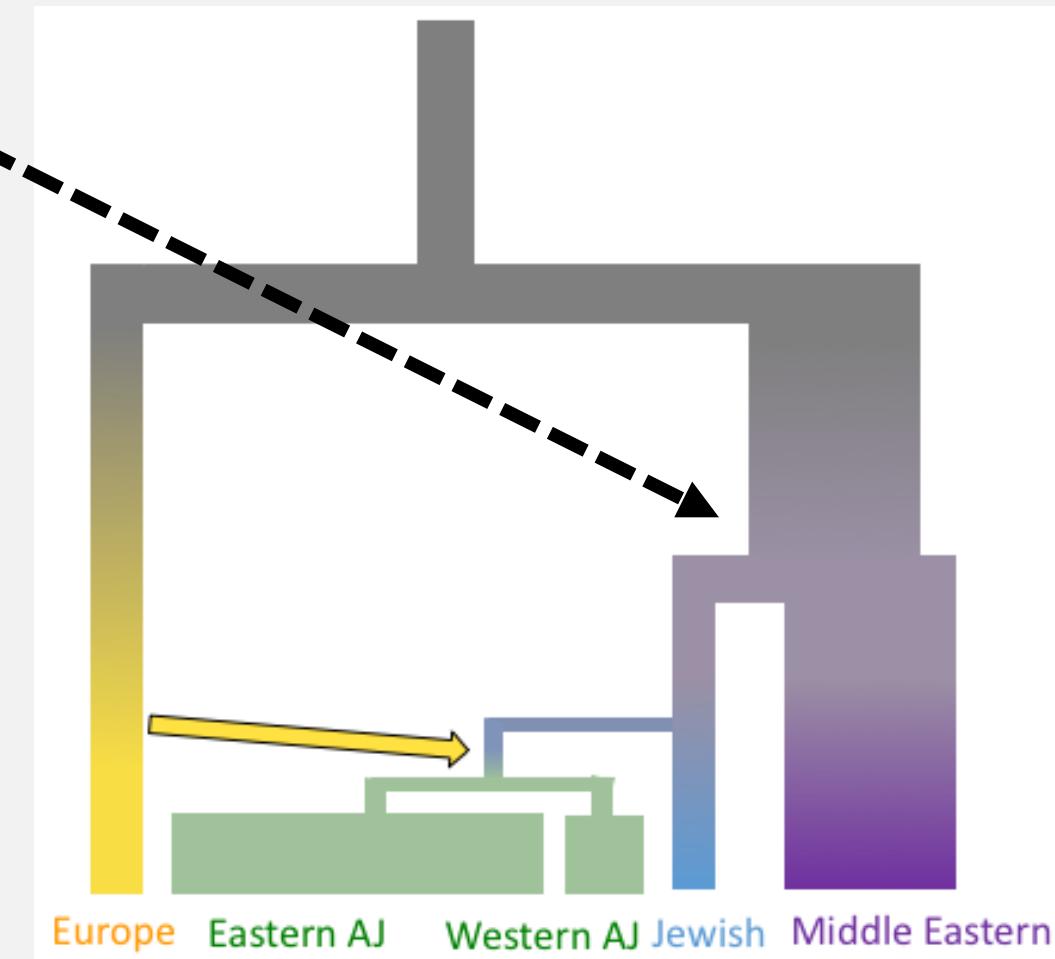
- ~ 3000 BCE ancestors of Jewish populations diverged from other Middle Eastern populations
 - Experienced extreme population size reduction



BEST MODEL

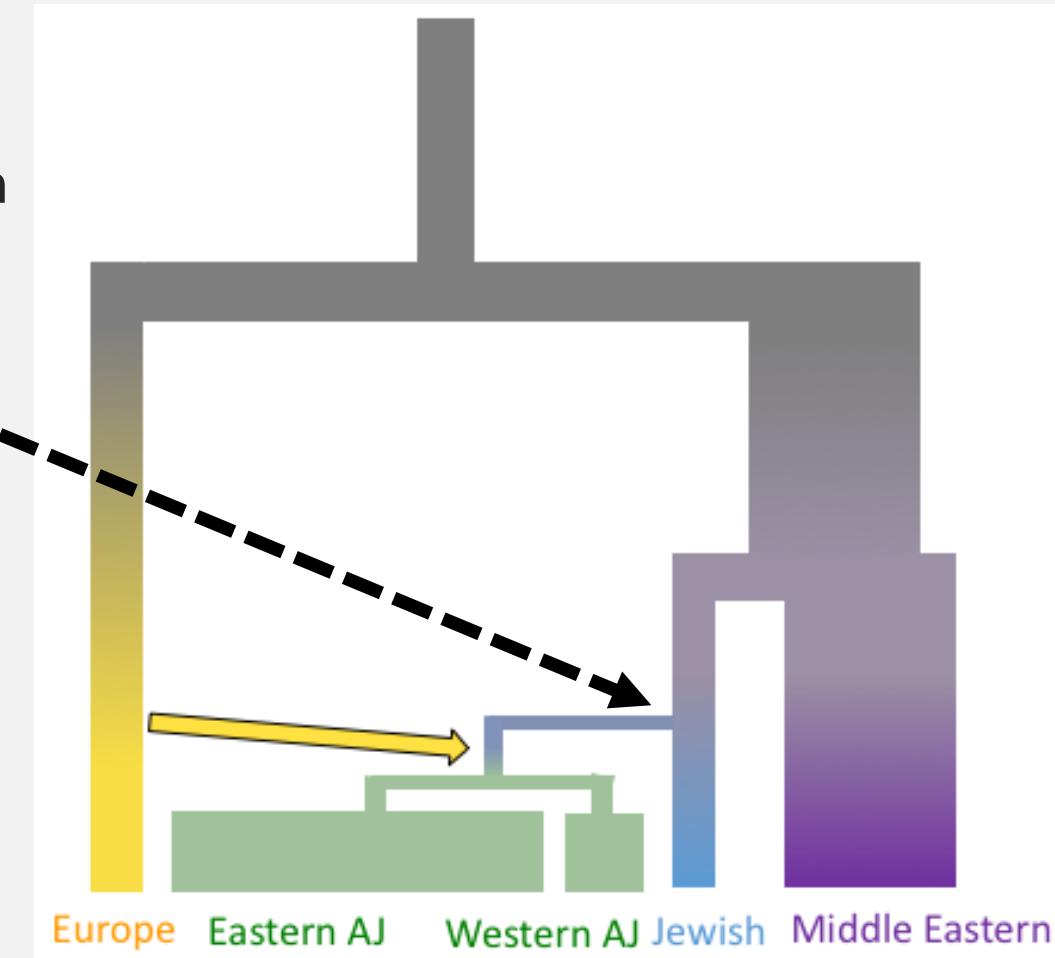
- ~ 3000 BCE ancestors of Jewish populations diverged from other Middle Eastern populations
 - Experienced extreme population size reduction

- First written accounts of “Israel” from Merneptah Stele in 1207 BCE



BEST MODEL

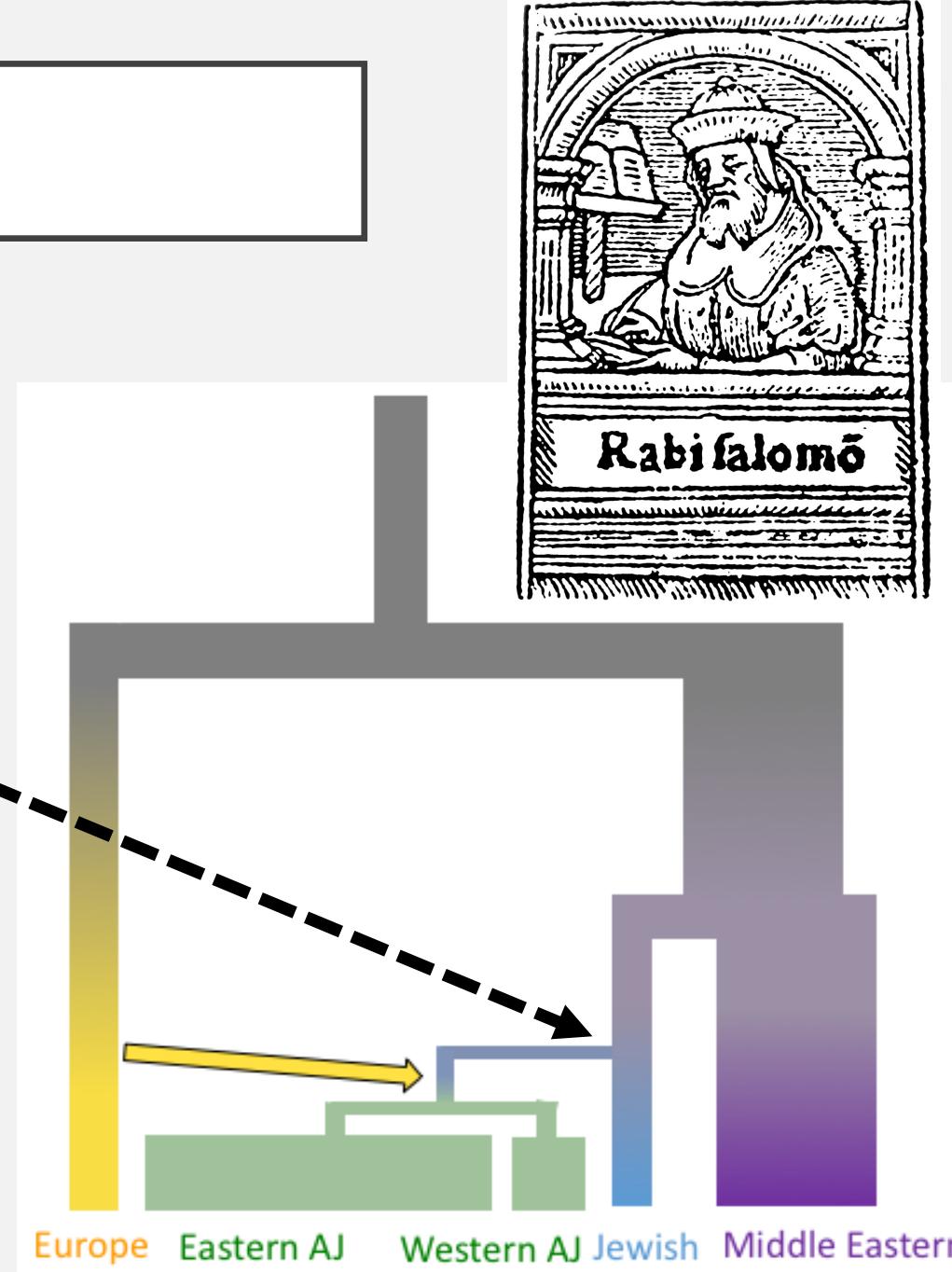
- ~ 3000 BCE ancestors of Jewish populations diverged from other Middle Eastern populations
 - Experienced extreme population size reduction
- 13th century ancestors of Ashkenazi Jews diverged from other Jewish populations
 - Experienced another population size reduction



BEST MODEL

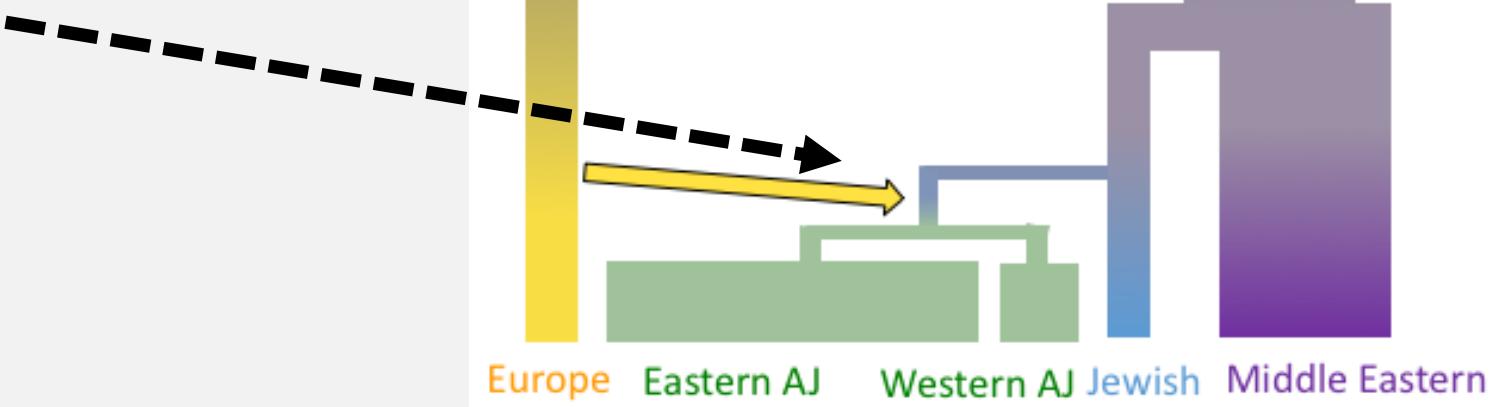
- ~ 3000 BCE ancestors of Jewish populations diverged from other Middle Eastern populations
 - Experienced extreme population size reduction
- 13th century ancestors of Ashkenazi Jews diverged from other Jewish populations
 - Experienced another population size reduction

- Migrations northward from Italy led to AJ community in Rhine Valley by 10th century.
- In the late 10th, 11th, and 12th centuries charters were issued to protect Jews in towns.
- In the 11th and 12th centuries the Ashkenazi rabbinic genres formed.



BEST MODEL

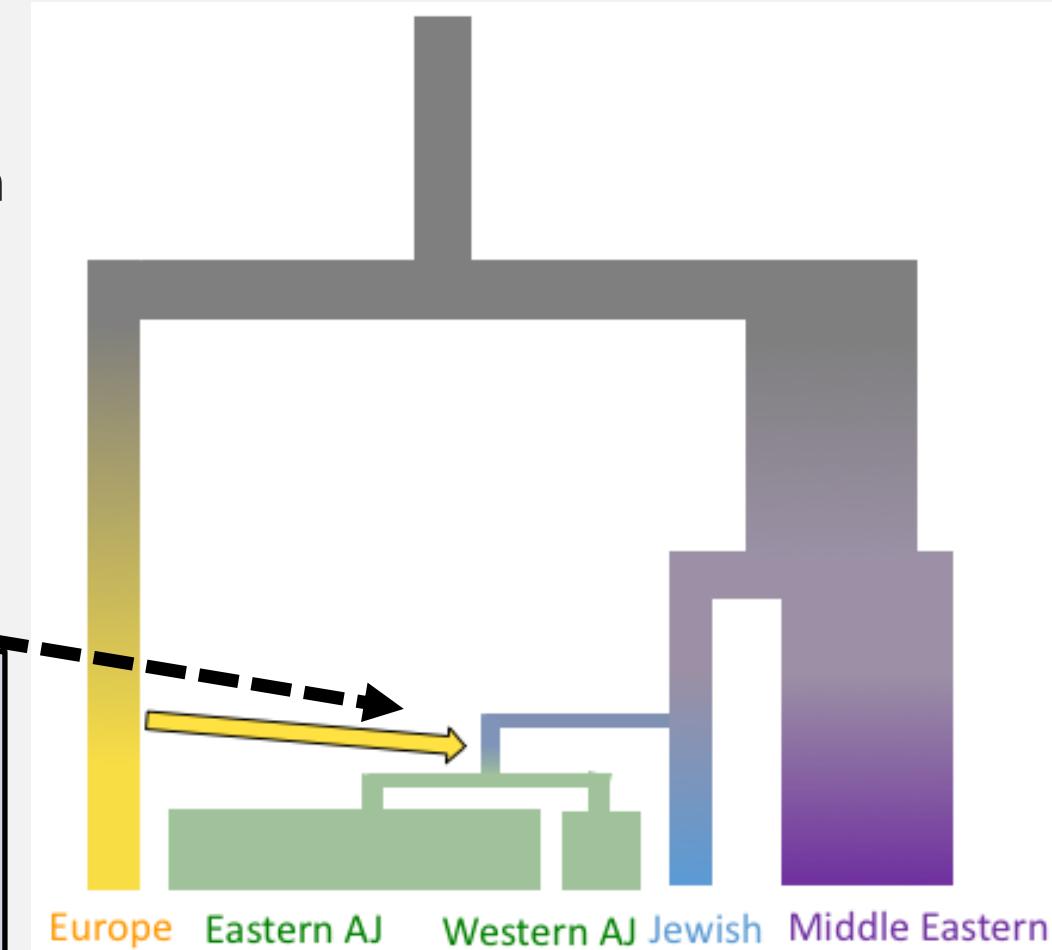
- ~ 3000 BCE ancestors of Jewish populations diverged from other Middle Eastern populations
 - Experienced extreme population size reduction
- 13th century ancestors of Ashkenazi Jews diverged from other Jewish populations
 - Experienced another population size reduction
 - Experienced gene flow from Europeans
(unresolved how much or when)



BEST MODEL

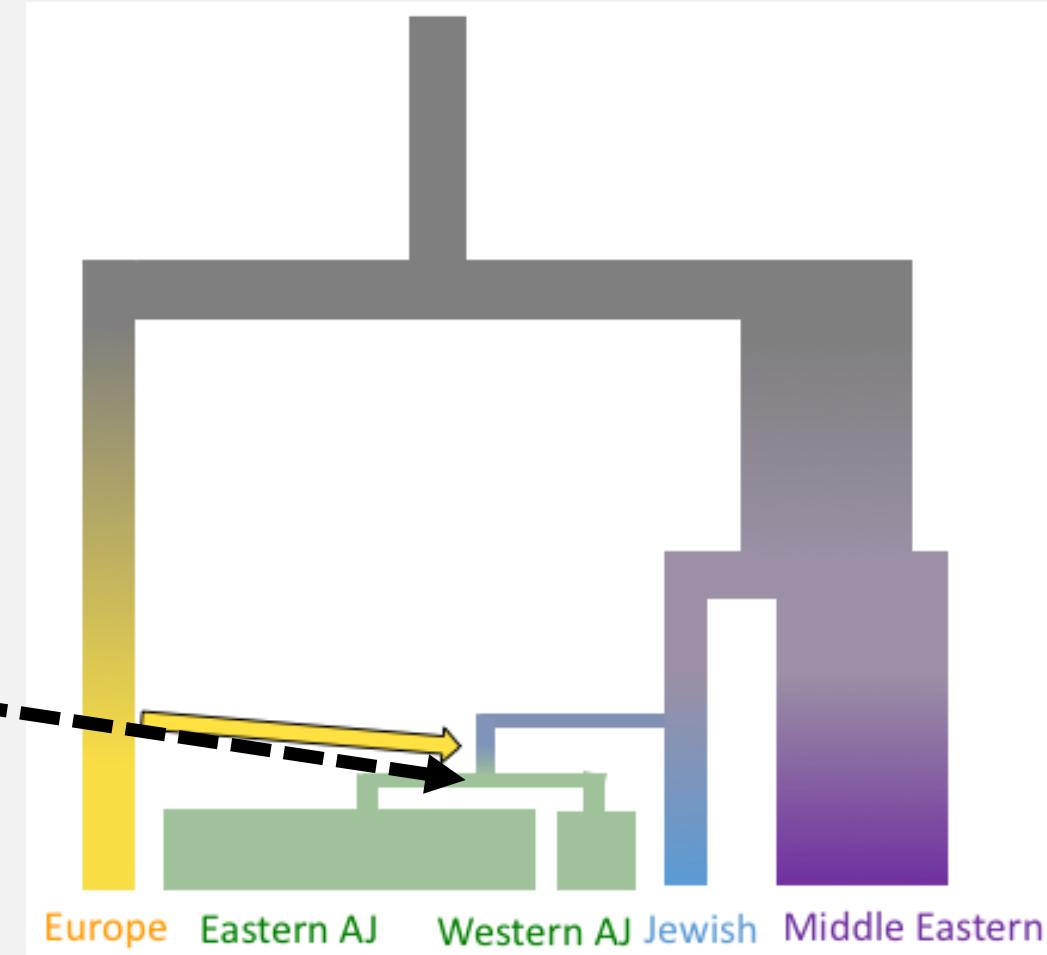
- ~ 3000 BCE ancestors of Jewish populations diverged from other Middle Eastern populations
 - Experienced extreme population size reduction
- 13th century ancestors of Ashkenazi Jews diverged from other Jewish populations
 - Experienced another population size reduction
 - Experienced gene flow from Europeans
(unresolved how much or when)

- Judaism follows matrilineal descent.
- In Central Europe Jews became increasingly integrated into gentile life.
- In Eastern Europe Jews became increasingly isolated.



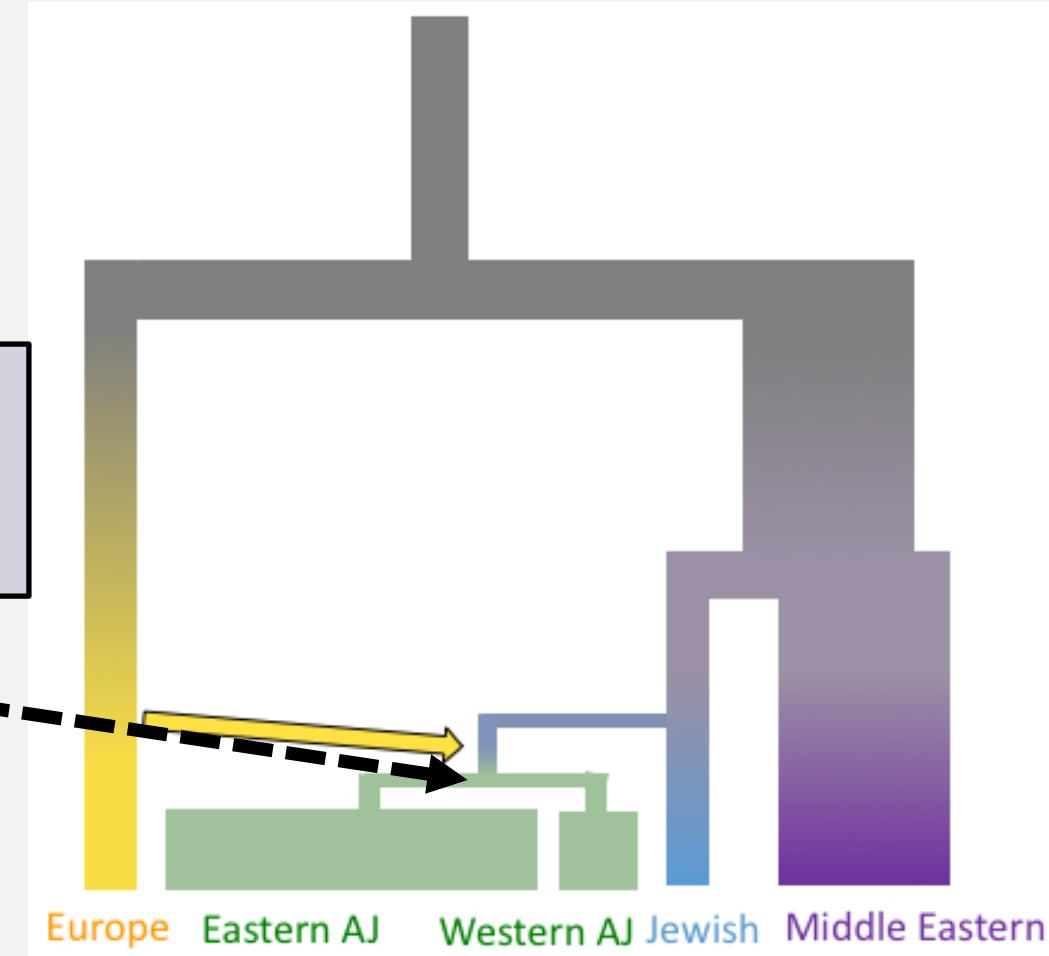
BEST MODEL

- ~ 3000 BCE ancestors of Jewish populations diverged from other Middle Eastern populations
 - Experienced extreme population size reduction
- 13th century ancestors of Ashkenazi Jews diverged from other Jewish populations
 - Experienced another population size reduction
 - Experienced gene flow from Europeans
(unresolved how much or when)
- 16th century Eastern and Western Ashkenazi Jews diverged



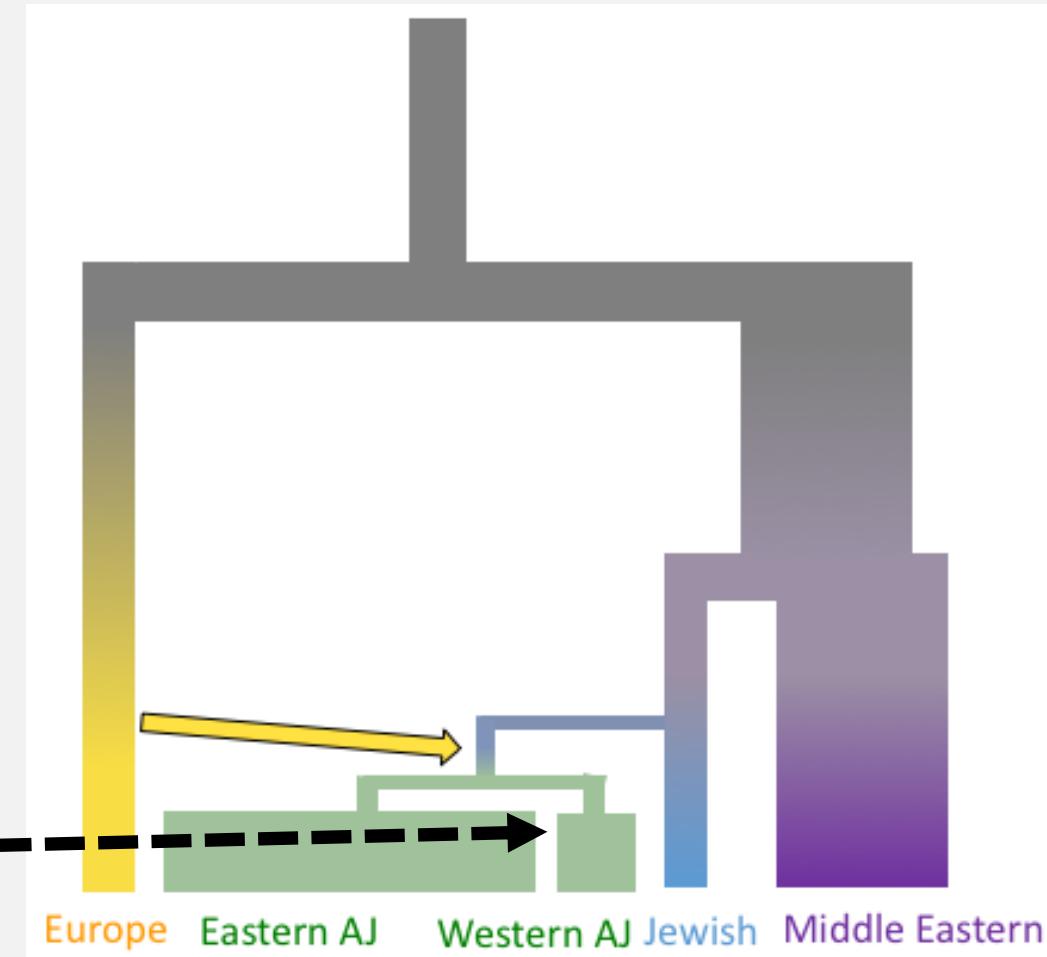
BEST MODEL

- ~ 3000 BCE ancestors of Jewish populations diverged from other Middle Eastern populations
 - Experienced extreme population size reduction
- 13th century ancestors of Ashkenazi Jews diverged from other Jewish populations
 - Migrations from Central Europe to Poland in the 14th, 15th, and 16th centuries.
 - By 16th century Polish Jewry culturally distinct.
- 16th century Eastern and Western Ashkenazi Jews diverged



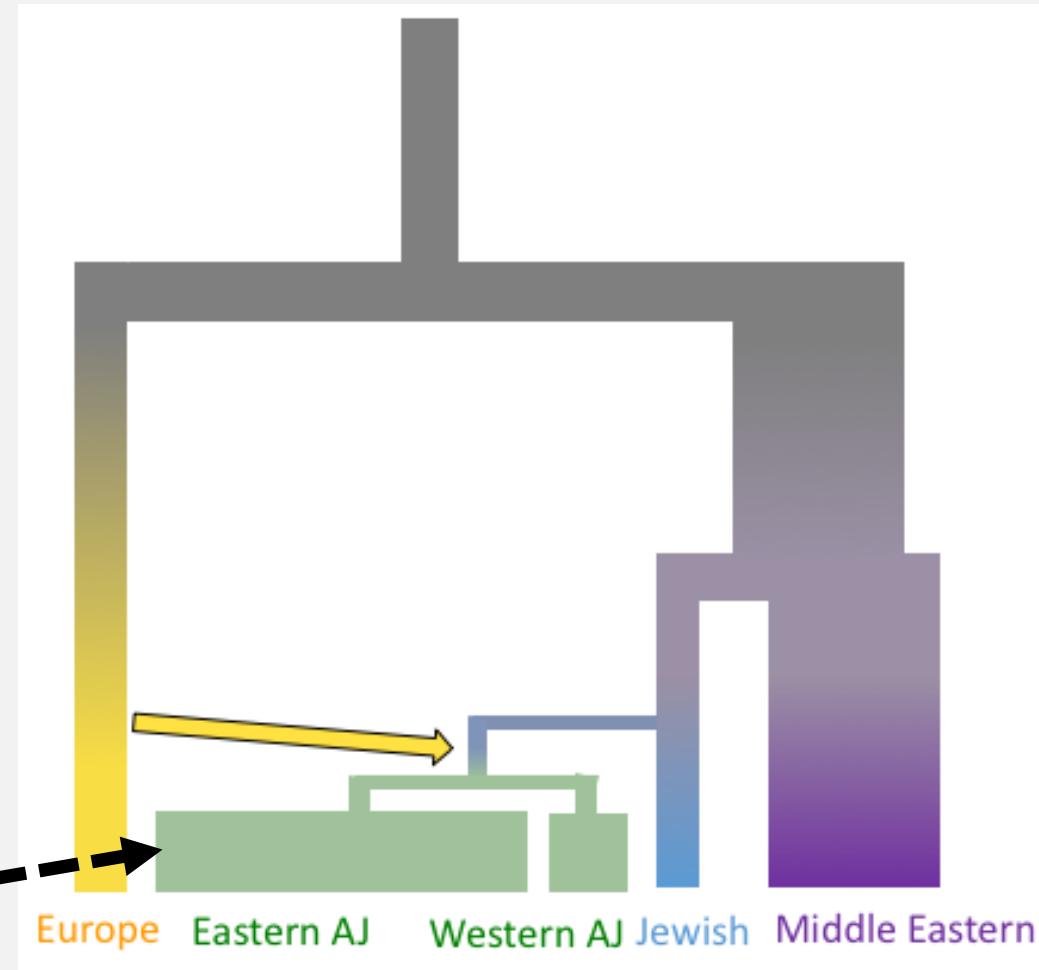
BEST MODEL

- ~ 3000 BCE ancestors of Jewish populations diverged from other Middle Eastern populations
 - Experienced extreme population size reduction
- 13th century ancestors of Ashkenazi Jews diverged from other Jewish populations
 - Experienced another population size reduction
 - Experienced gene flow from Europeans
(unresolved how much or when)
- 16th century Eastern and Western Ashkenazi Jews diverged
 - Western AJ moderately grew in size

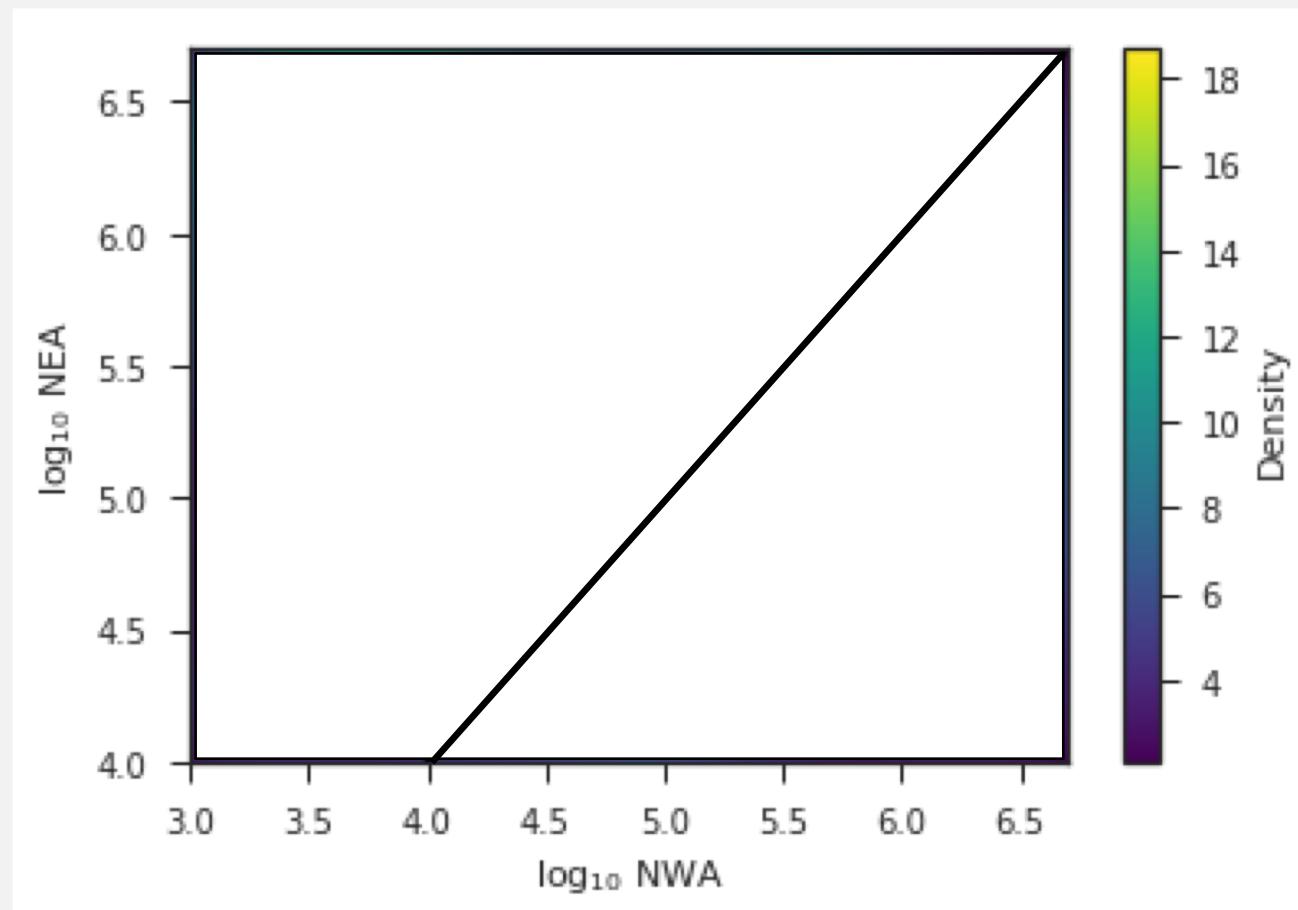


BEST MODEL

- ~ 3000 BCE ancestors of Jewish populations diverged from other Middle Eastern populations
 - Experienced extreme population size reduction
- 13th century ancestors of Ashkenazi Jews diverged from other Jewish populations
 - Experienced another population size reduction
 - Experienced gene flow from Europeans
(unresolved how much or when)
- 16th century Eastern and Western Ashkenazi Jews diverged
 - Western AJ moderately grew in size
 - Eastern AJ massively grew in size

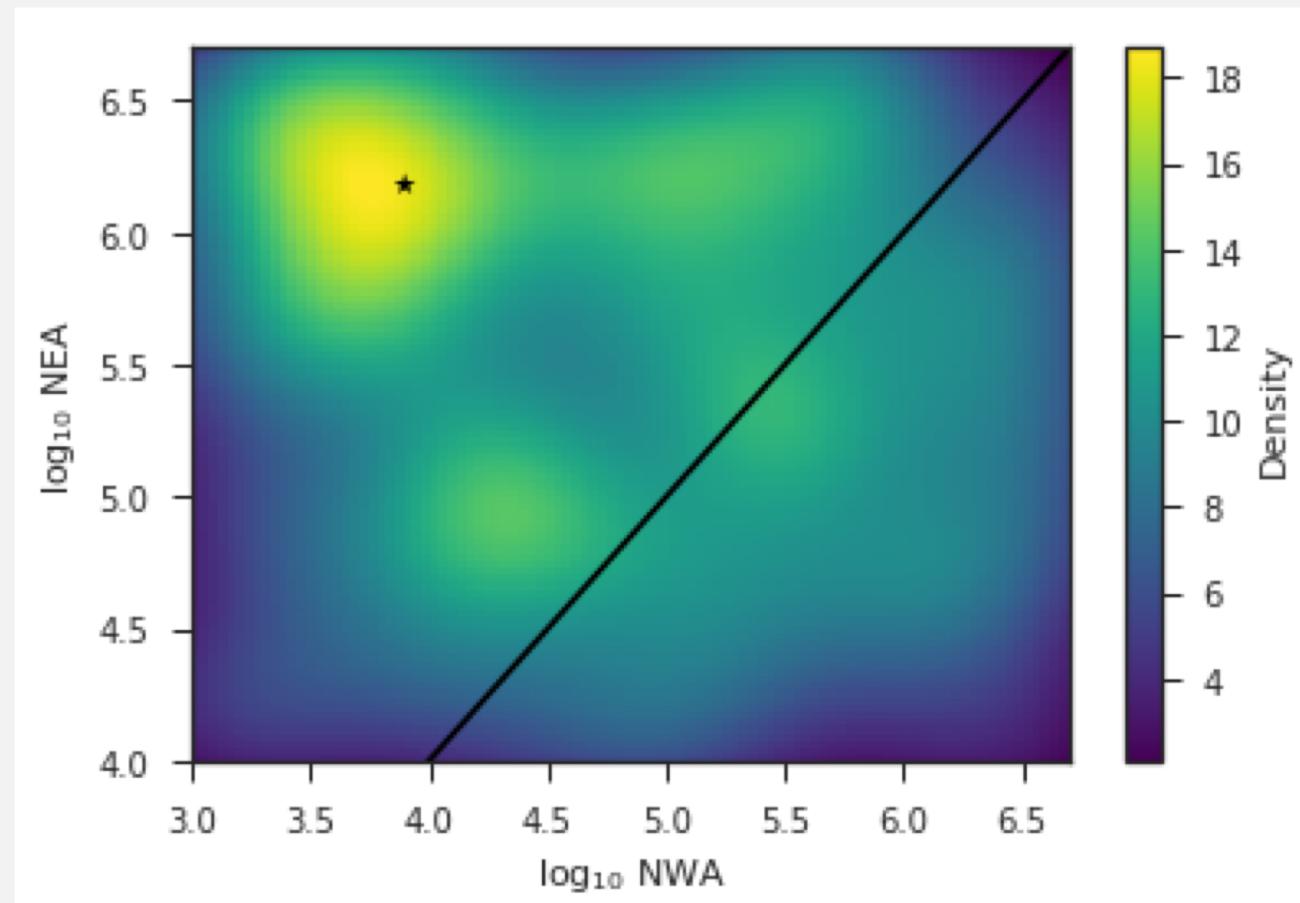


JOINT POSTERIOR OF EFFECTIVE POPULATION SIZE OF EASTERN AND WESTERN AJ



JOINT POSTERIOR OF EFFECTIVE POPULATION SIZE OF EASTERN AND WESTERN AJ

	Probability NEA > NWA
Chr1	0.69
Genome	0.62
Updated chr5	0.70



MORE GROWTH IN EASTERN AJ

Central Europe

- Often expelled from settlements.
- Strict regulations on where Jews could live and what they could do to earn a living.



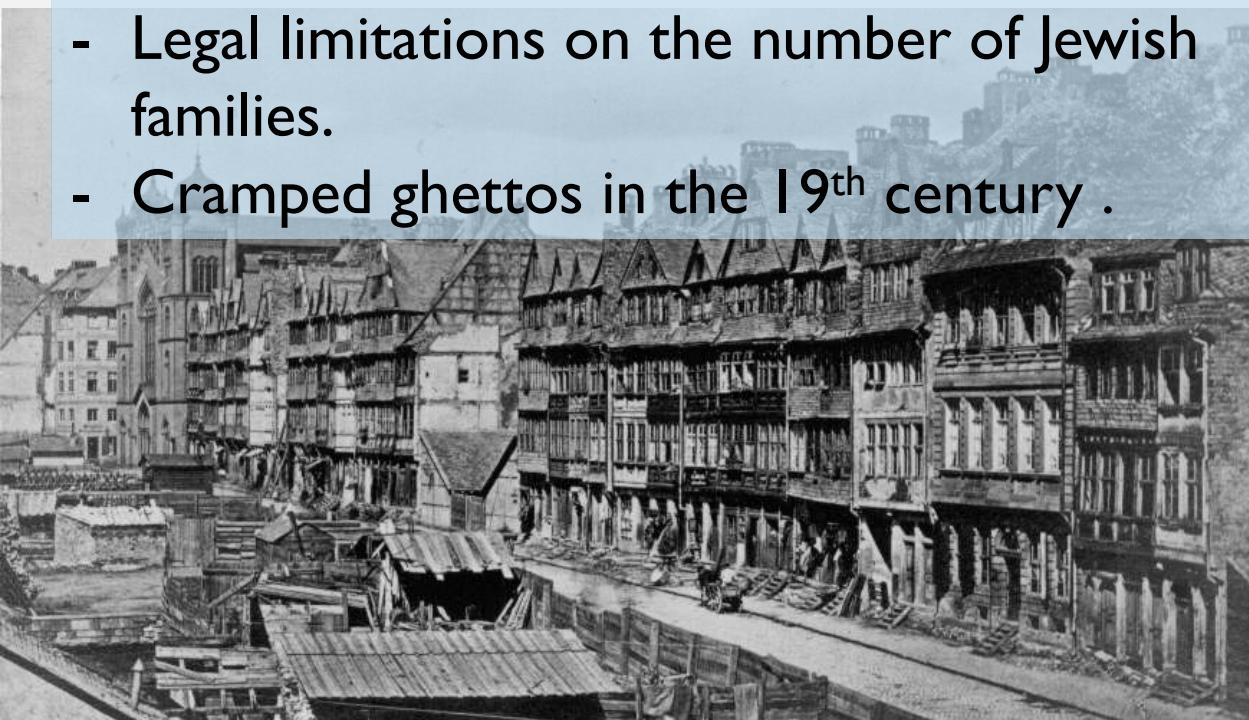
Eastern Europe

- Could generally move freely.
- Protected by nobles.



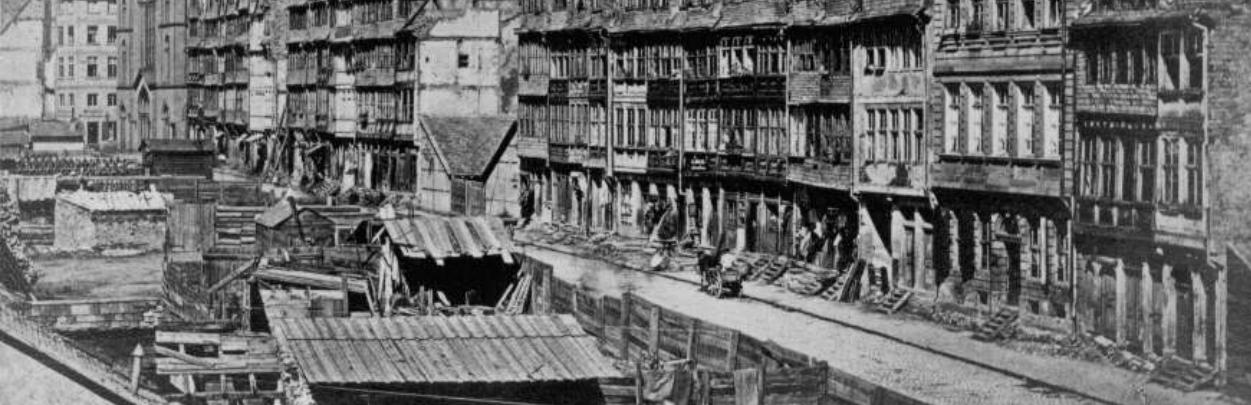
MORE GROWTH IN EASTERN AJ

Central Europe	Eastern Europe
<ul style="list-style-type: none">- Often expelled from settlements.- Strict regulations on where Jews could live and what they could do to earn a living.- Legal limitations on the number of Jewish families.- Cramped ghettos in the 19th century .	<ul style="list-style-type: none">- Could generally move freely.- Protected by nobles.- No limitations from government on number of Jewish marriages



MORE GROWTH IN EASTERN AJ

Central Europe	Eastern Europe
<ul style="list-style-type: none">- Often expelled from settlements.- Strict regulations on where Jews could live and what they could do to earn a living.- Legal limitations on the number of Jewish families.- Cramped ghettos in the 19th century- Integration into non-Jewish society.	<ul style="list-style-type: none">- Could generally move freely.- Protected by nobles.- No limitations from government on number of Jewish marriages- Adherence to religious and traditional norms and economic structures encouraged early marriage and high fertility.



IMPORTANCE OF WORK

Historical / Cultural	Evolution / Population genetics	Medical
Resolved controversial question of Jewish population growth in Eastern Europe.	Demonstration of inference of very recent history.	How do different growth rates in Western and Eastern AJ affect the frequency of deleterious mutations?

FUTURE DIRECTIONS

- Infer demographic history in other populations with histories of population size changes or inbreeding and admixture
- Approximate Bayesian Computation
 - Using other statistics to better infer admixture
- Machine learning
 - Without using genomic statistics

“Big Data”

THANK YOU!

HAMMER LAB (AND FORMER)

- Michael Hammer
- Consuelo Quinto-Cortes
- August Woerner
- Fernando Mendez

UA HPC CONSULTING

- Mike Bruck
- Dima Shyshlov

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

CODING MINIONS

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

CYVERSE

- Blake Joyce
- Julian Pistorius

RESOURCES PROVIDED BY

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