

Demographic inference based on Site frequency spectrum (SFS) – Part II

Vitor Sousa

CE3C – center for ecology, evolution and environmental changes

2018 WSPG Cesky Krumlov
22 Jan 2020

vmsousa@fc.ul.pt



Outline part II

Example of Applications:

- Human dispersal out of Africa (high quality whole-genome) – lessons on choice of models
- Human colonization of Siberia and America (ancient whole-genome data) - lessons on dealing with sequencing errors
- Deer mice colonization of Nebraska Sand Hills (targeted re-capture data) – lessons on effects of filtering
- Inferring divergence times and gene flow in sawflies (ddRAD-seq data) – lessons from comparing models



Nourlangie, Kakadu National Park, NT, Australia

A genomic history of Aboriginal Australia

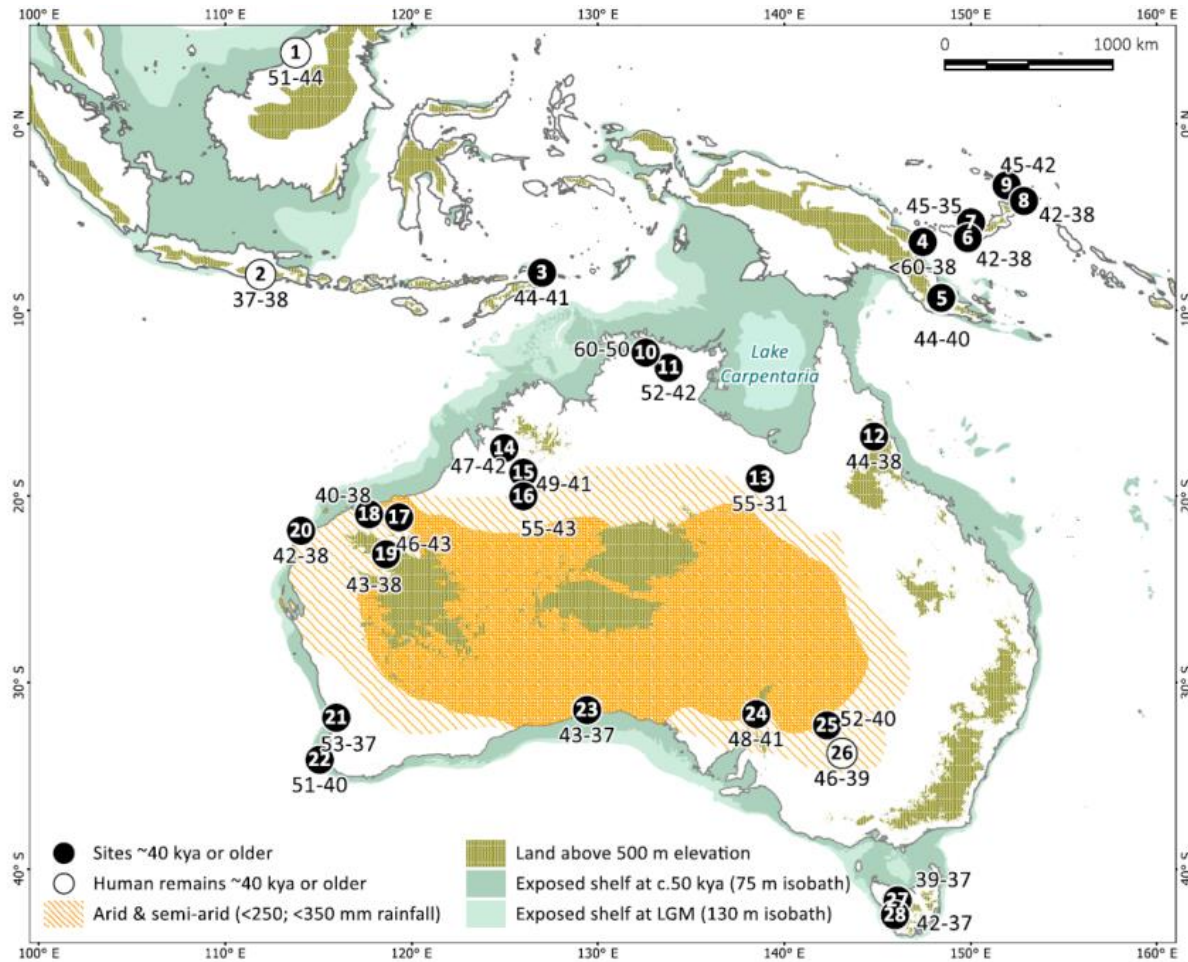
Anna-Sapfo Malaspinas^{1,2,3*}, Michael C. Westaway^{4*}, Craig Muller^{1*}, Vitor C. Sousa^{2,3*}, Oscar Lao^{5,6*}, Isabel Alves^{2,3,7*}, Anders Bergström^{8*}, Georgios Athanasiadis⁹, Jade Y. Cheng^{9,10}, Jacob E. Crawford^{10,11}, Tim H. Heupink⁴, Enrico Macholdt¹², Stephan Peischl^{3,13}, Simon Rasmussen¹⁴, Stephan Schiffels¹⁵, Sankar Subramanian⁴, Joanne L. Wright⁴, Anders Albrechtsen¹⁶, Chiara Barbieri^{12,17}, Isabelle Dupanloup^{2,3}, Anders Eriksson^{18,19}, Ashot Margaryan¹, Ida Moltke¹⁶, Irina Pugach¹², Thorfinn S. Korneliussen¹, Ivan P. Levkivskyi²⁰, J. Víctor Moreno-Mayar¹, Shengyu Ni¹², Fernando Racimo¹⁰, Martin Sikora¹, Yali Xue⁸, Farhang A. Aghakhanian²¹, Nicolas Brucato²², Søren Brunak²³, Paula F. Campos^{1,24}, Warren Clark²⁵, Sturla Ellingvåg²⁶, Gudjugudju Fourmile²⁷, Pascale Gerbault^{28,29}, Darren Injie³⁰, George Koki³¹, Matthew Leavesley³², Betty Logan³³, Aubrey Lynch³⁴, Elizabeth A. Matisoo-Smith³⁵, Peter J. McAllister³⁶, Alexander J. Mentzer³⁷, Mait Metspalu³⁸, Andrea B. Migliano²⁹, Les Murgha³⁹, Maude E. Phipps²¹, William Pomat³¹, Doc Reynolds⁴⁰, Francois-Xavier Ricaut²², Peter Siba³¹, Mark G. Thomas²⁸, Thomas Wales⁴¹, Colleen Ma'run Wall⁴², Stephen J. Oppenheimer⁴³, Chris Tyler-Smith⁸, Richard Durbin⁸, Joe Dortch⁴⁴, Andrea Manica¹⁸, Mikkel H. Schierup⁹, Robert A. Foley^{1,45}, Marta Mirazón Lahr^{1,45}, Claire Bowern⁴⁶, Jeffrey D. Wall⁴⁷, Thomas Mailund⁹, Mark Stoneking¹², Rasmus Nielsen^{1,48}, Manjinder S. Sandhu⁸, Laurent Excoffier^{2,3}, David M. Lambert⁴ & Eske Willerslev^{1,8,18}

Nature(2016)



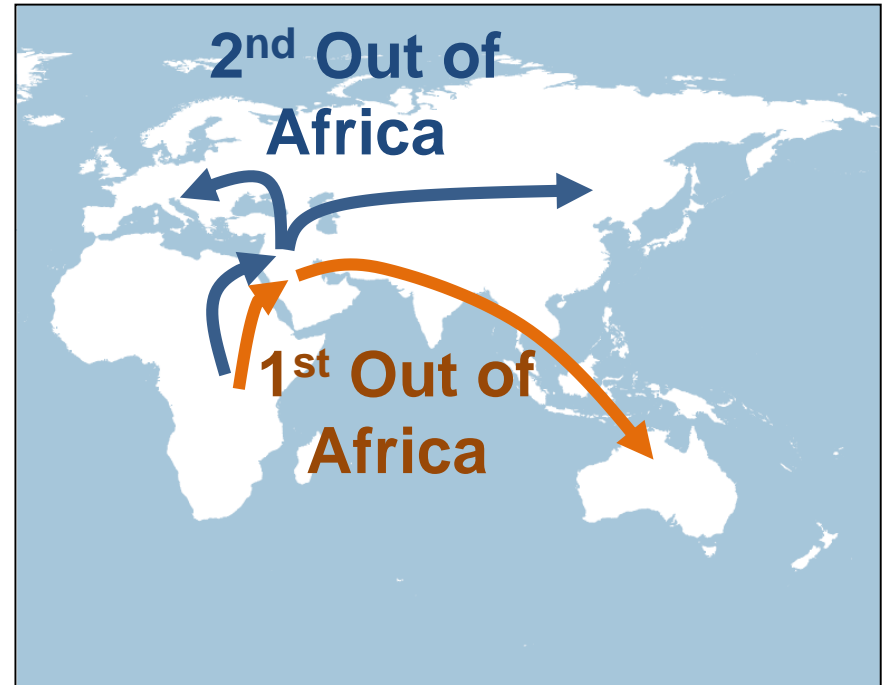
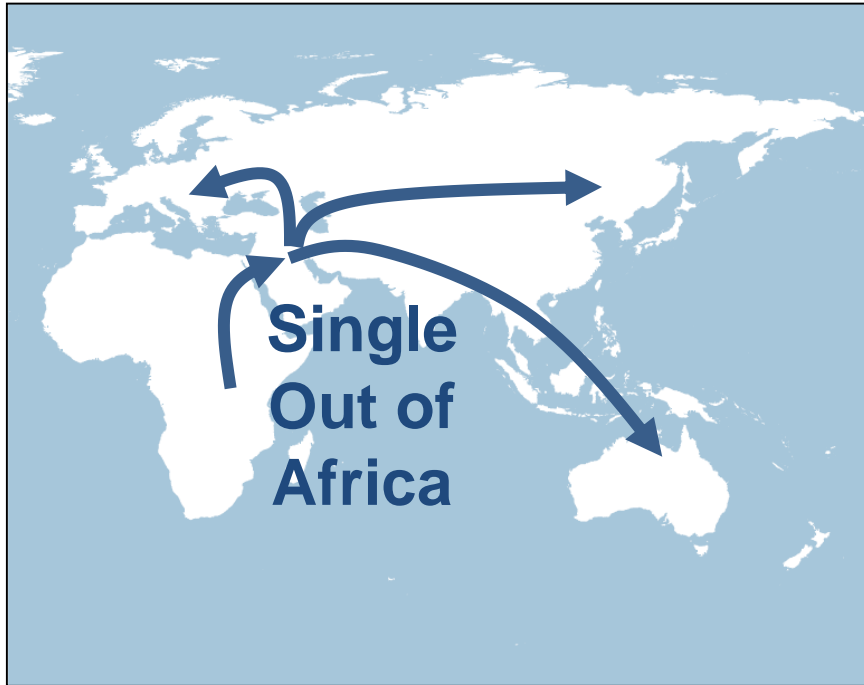
Ewaninga Rock Carvings Conservation Reserve, NT, Australia

Australia harbors some of the oldest modern human remains outside Africa

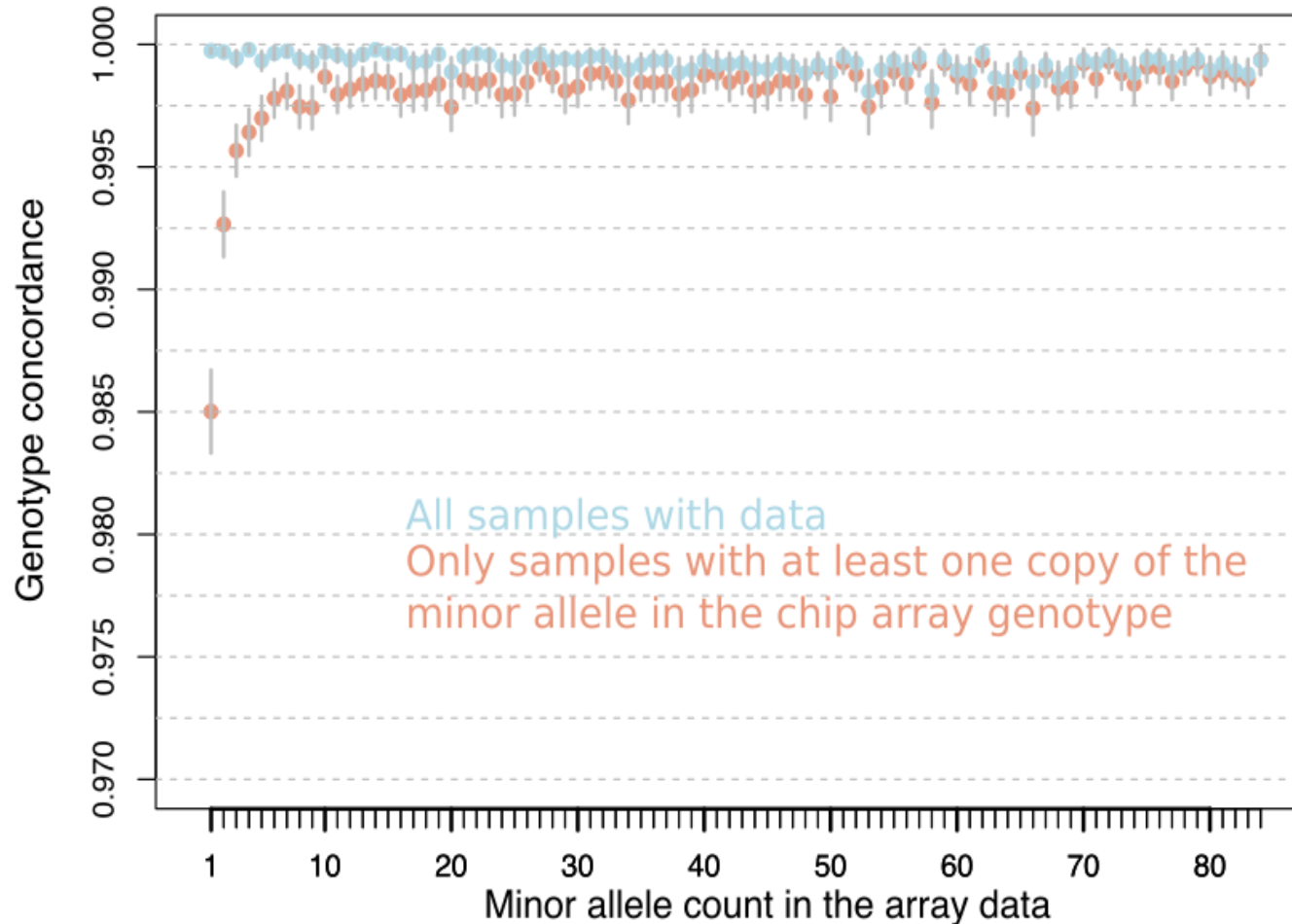


Many sites and remains dated to be older than 40 kya, suggesting a human settlement 47.5-55 kya

One wave out of Africa vs Two waves out of Africa



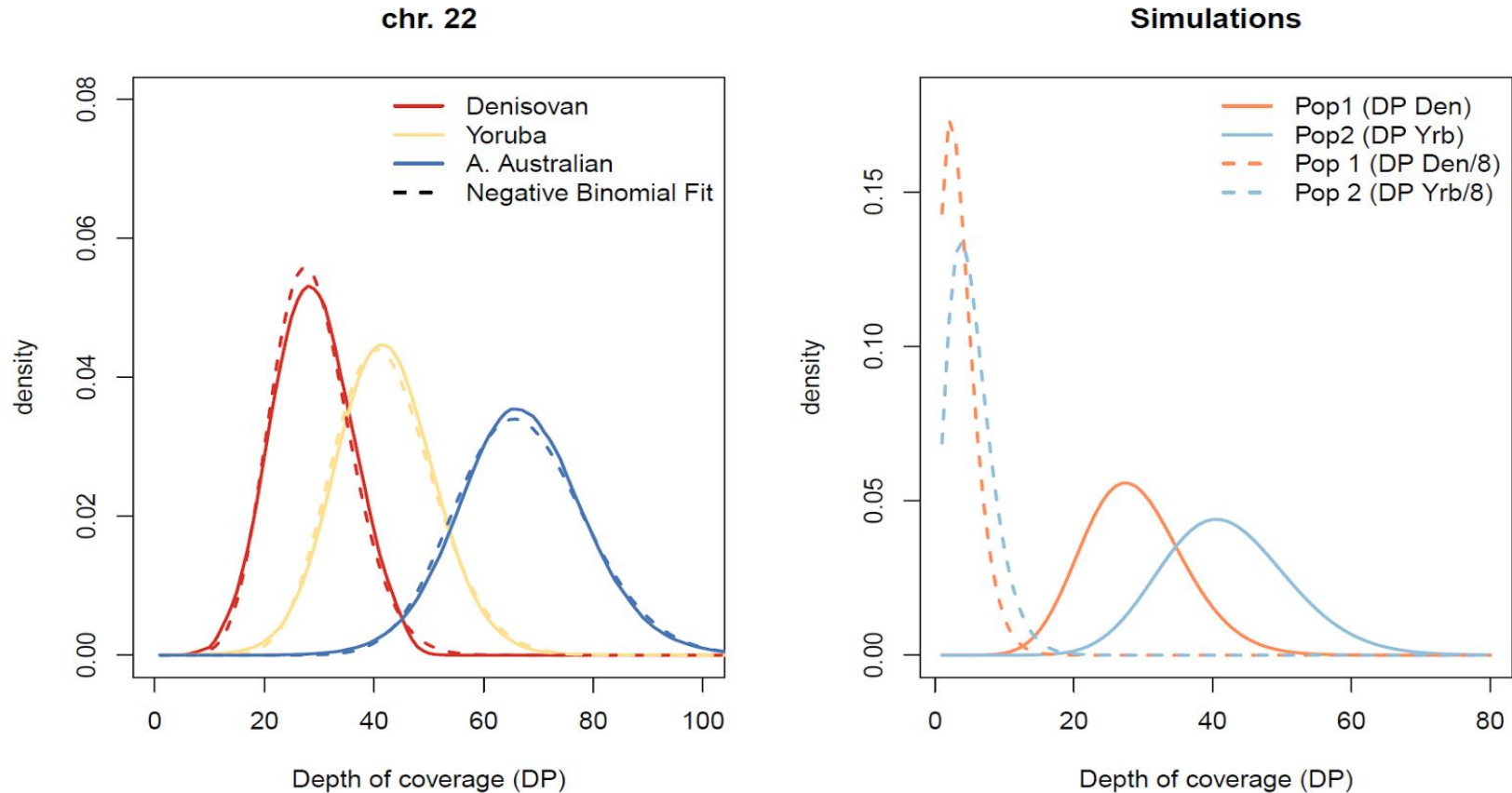
83 high-coverage Aboriginal Australians genomes



Average depth of coverage: 65x

Very good quality of genotype calls

Effect of depth of coverage on SFS

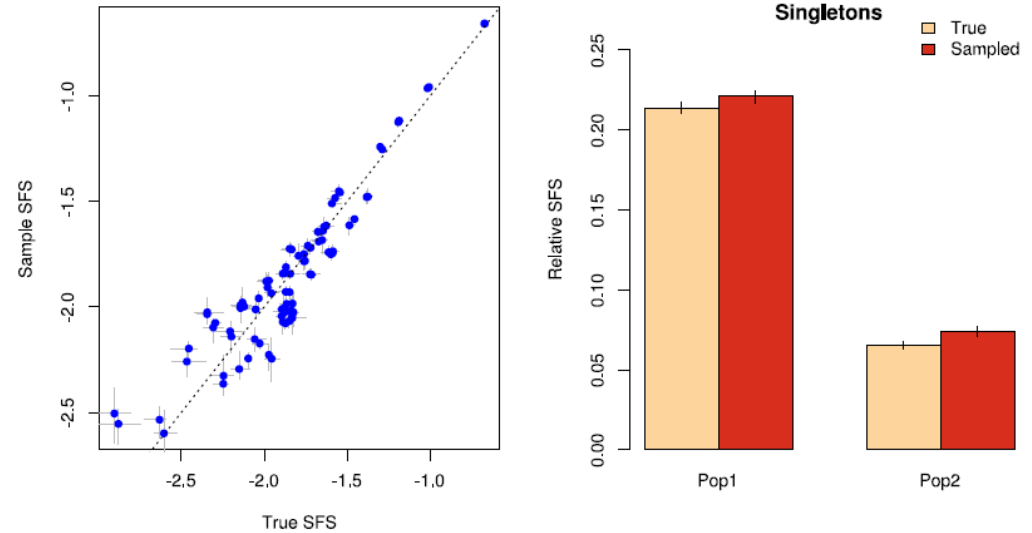


- Compared 2D SFS based on depth of coverage of observed data (mean larger than $>20x$), with a distribution 8 times smaller.

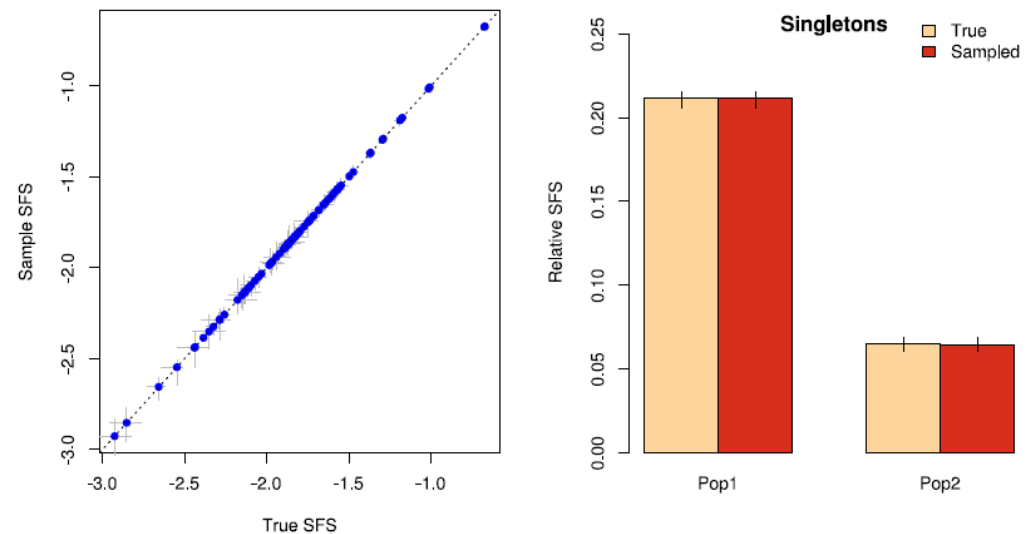
A note on recovering the SFS from genomic data

- Simulation study
- Low depth of coverage and missing data lead to biased SFS towards rare variants

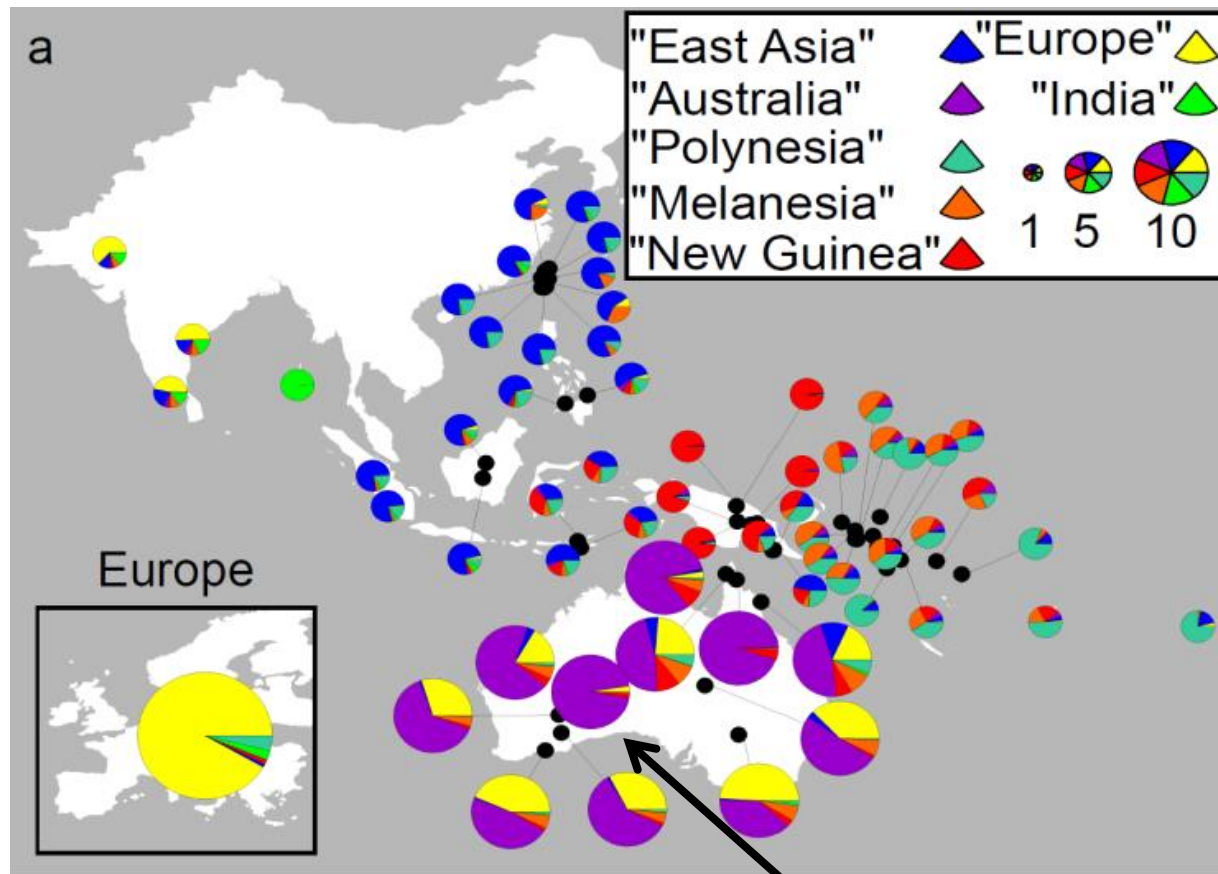
a) Low depth of coverage, no GQ filter, allowing missing data



b) Depth of coverage similar to observed data, GQ>30 filter, no missing data



83 high-coverage Aboriginal Australians genomes



Western Central Desert (WCD)

Average depth of coverage: 65x



- ★ Archaic human genomes:
- 1 Neanderthal (~66 kya)
 - 1 Denisovan (~52 kya)

Mutation rate assumed

1.25×10^{-8} /site/gen

Scally and Durbin (2012) *Nat. Rev. Genet.*

Generation time

29 years/gen

Fenner (2005) *Am. J. Phys. Anthropol.*

Since we want to infer demography we tried to minimize the number of sites affected by selection:

- 985 1Mb blocks outside genic regions and CpG islands (~4.3 Million SNPs)
- 5 dimensional SFS (16,875 entries)
- Confidence intervals obtained using block-bootstrap

Towards a model to test the hypotheses: One vs Two waves Out of Africa

- Data (SFS)



- (Re-)Define model
(hypotheses to test)



- Run fastsimcoal2



- Estimates!
 - Assess the fit to the data



Do you have an outgroup?

- **Yes** – use the derived (unfolded) SFS
- **No** – use the minor allele frequency spectrum (folded)

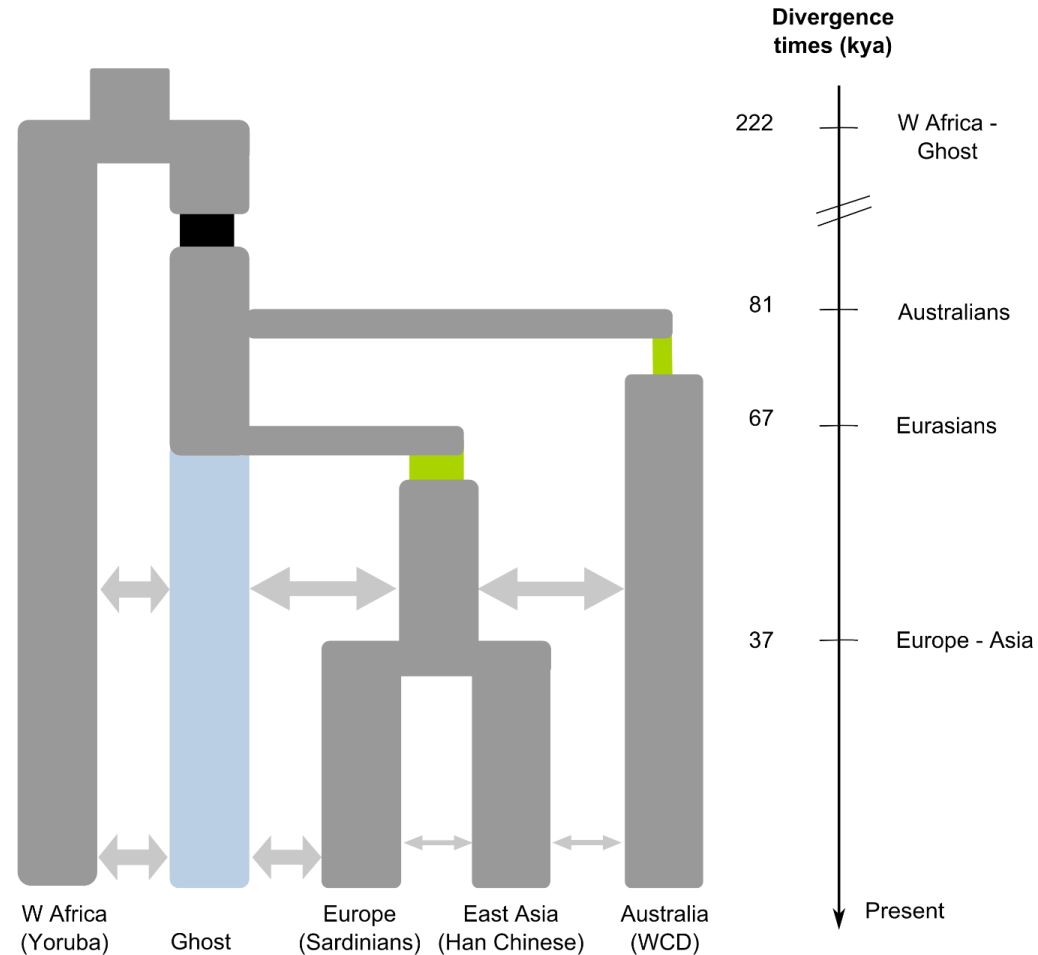
Do you have monomorphic sites?

- **Yes** - then, given a mutation rate you can infer the absolute times and effective sizes
- **No** – then all your estimates need to be relative to a fixed parameter (fixed N_e or fixed time)

We always get results...

Evidence of two waves Out of Africa:

- Old split leading to colonization of Australia (81kya)
- More recent split leading to colonization of Eurasia (67 kya)



Legend:

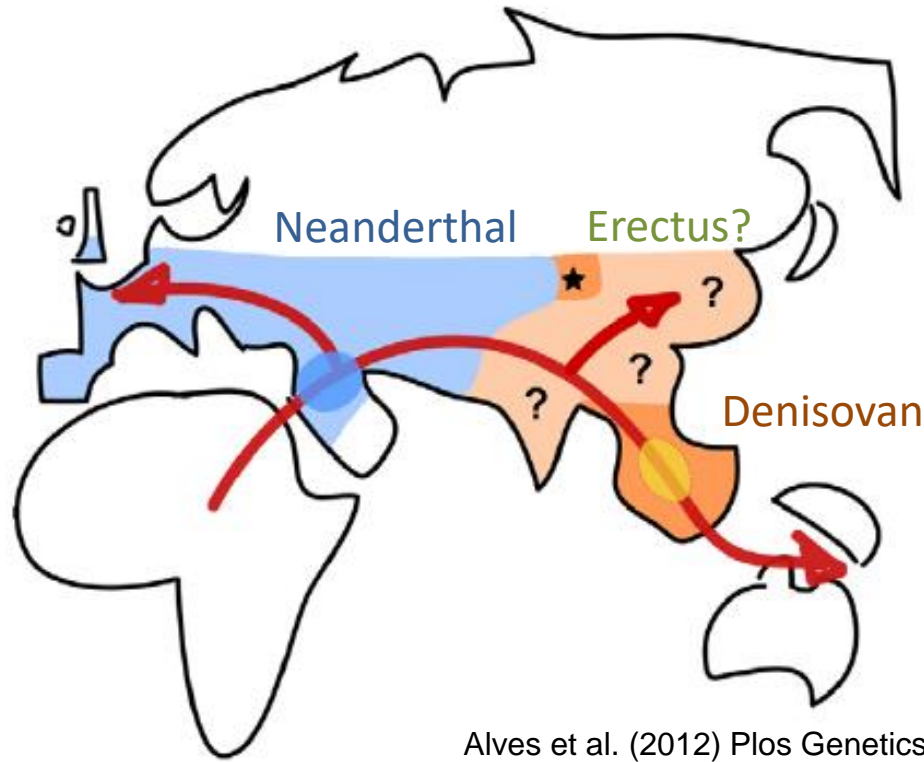
↔ Migration, $2Nm > 1$

■ Ancestral bottleneck

↔ Migration, $2Nm < 1$

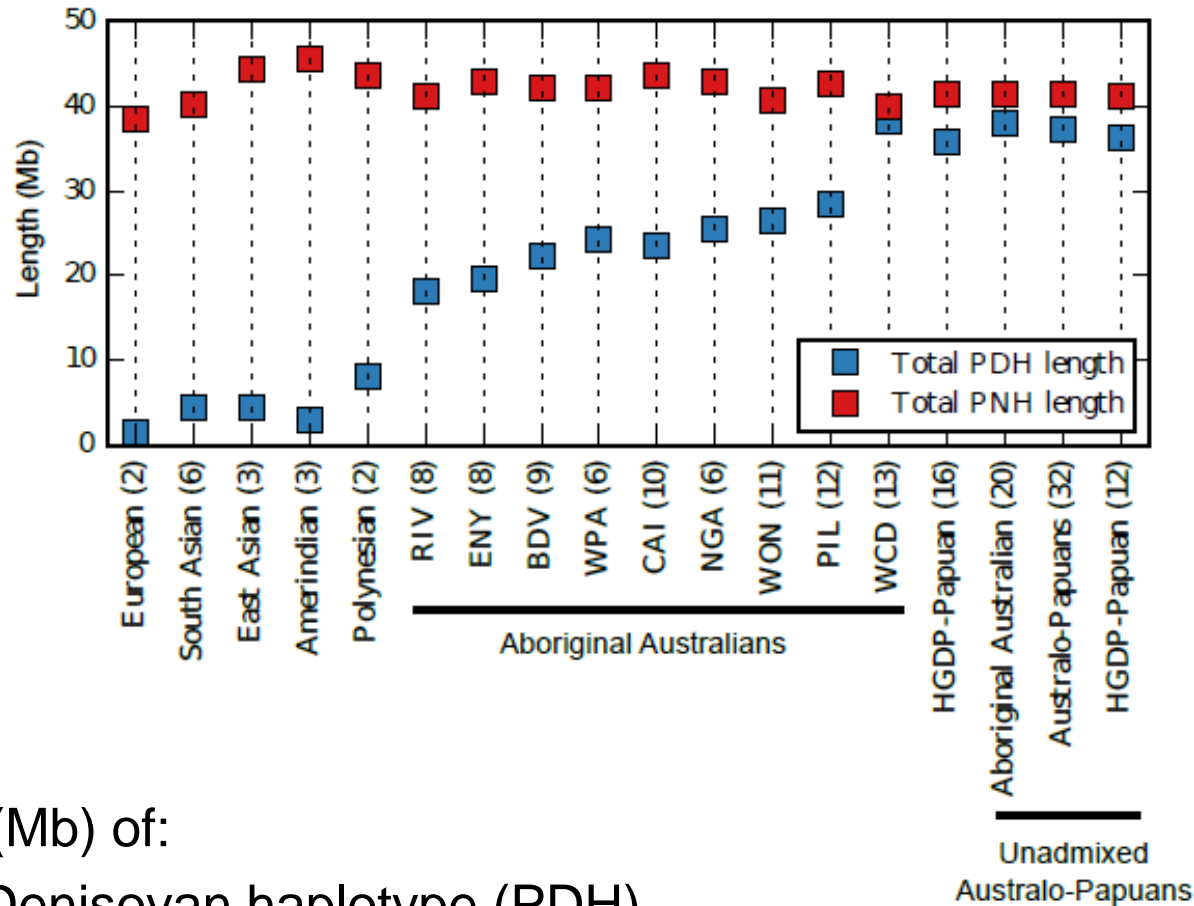
■ Continent-specific bottlenecks

Towards a model incorporating Neanderthal and Denisovan admixture



- Non-African populations: 1-4% estimated Neanderthal admixture
- Aboriginal Australians and New Guineans: 3-6% estimated Denisovan admixture
- Archaic admixture can affect times of split estimates

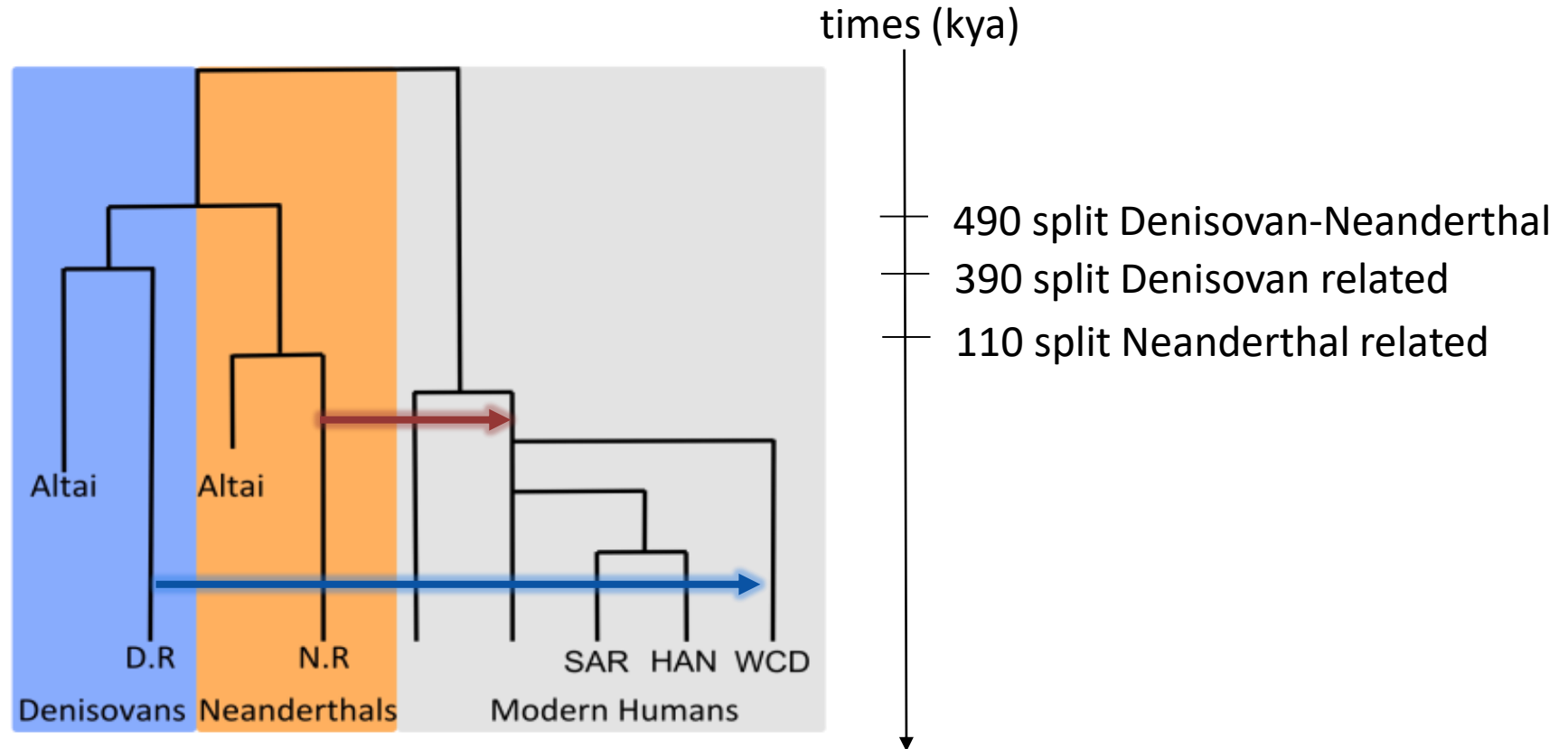
Evidence of archaic introgression



Total length (Mb) of:

- Putative Denisovan haplotype (PDH)
- Putative Neanderthal haplotypes (PNH)

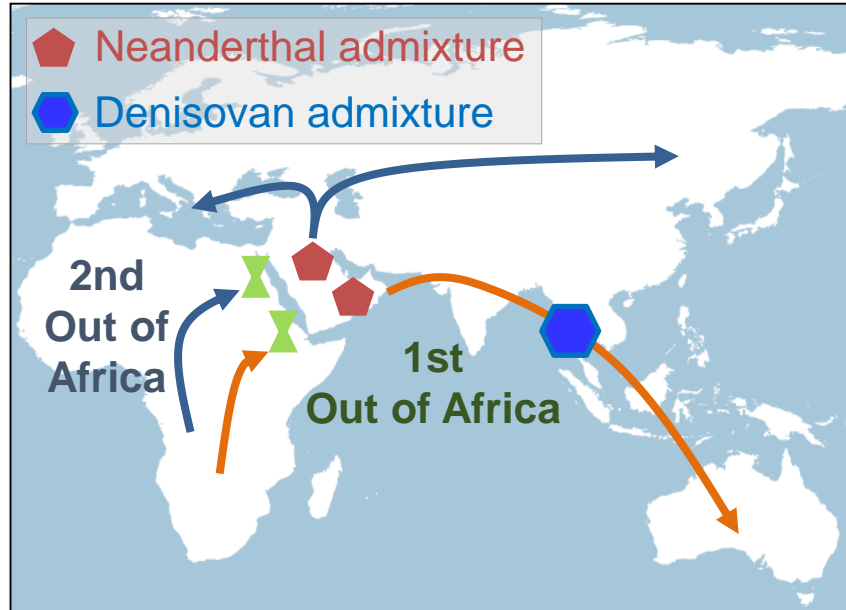
Accounting for shared ancestry of Neanderthal and Denisovan



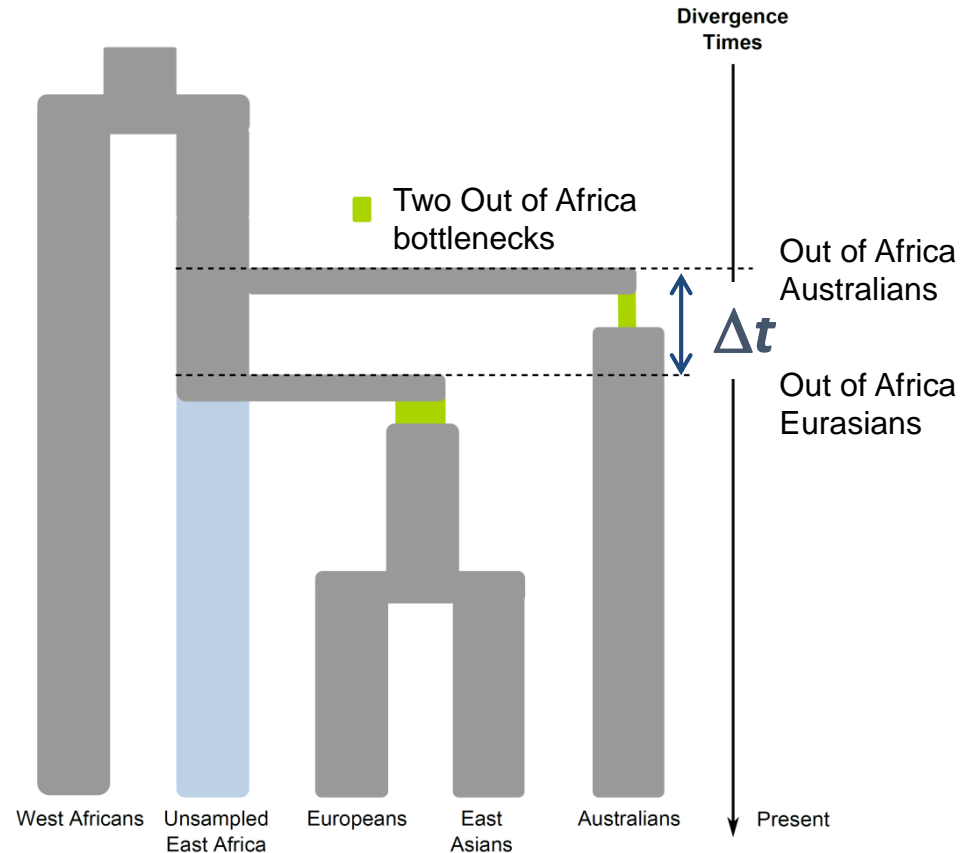
Admixture occurs between modern humans and:

- Denisovan-related (D.R.) population
- Neanderthal-related (N.R.) population

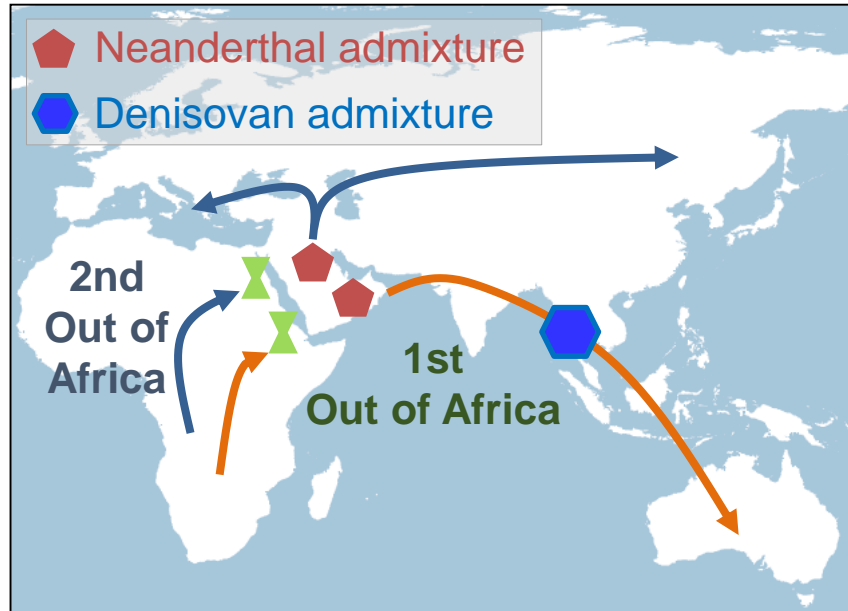
Two-waves out of Africa



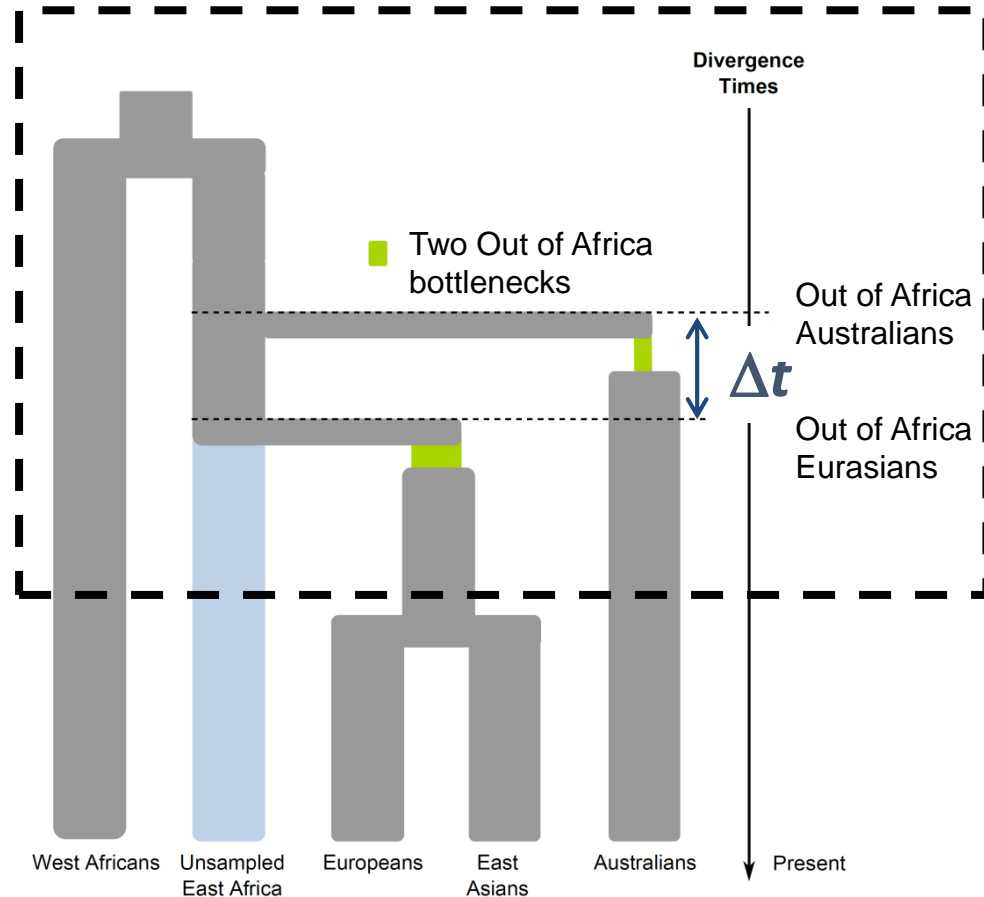
- Two different divergence times ($\Delta t \gg 0$)
- Two independent bottlenecks associated with the two Out of Africa events



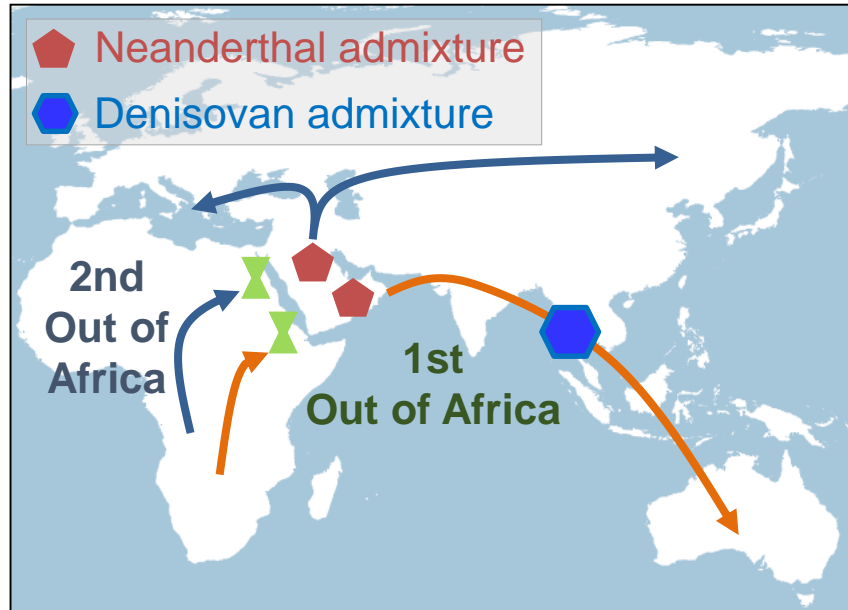
Two-waves out of Africa



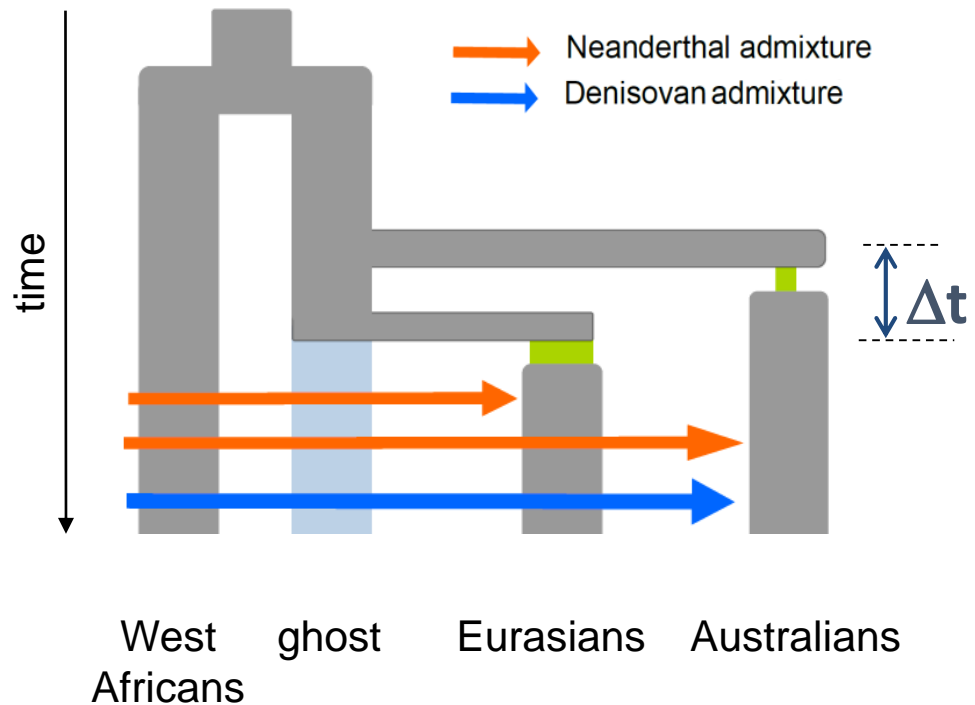
- Two different divergence times ($\Delta t \gg 0$)
- Two independent bottlenecks associated with the two Out of Africa events



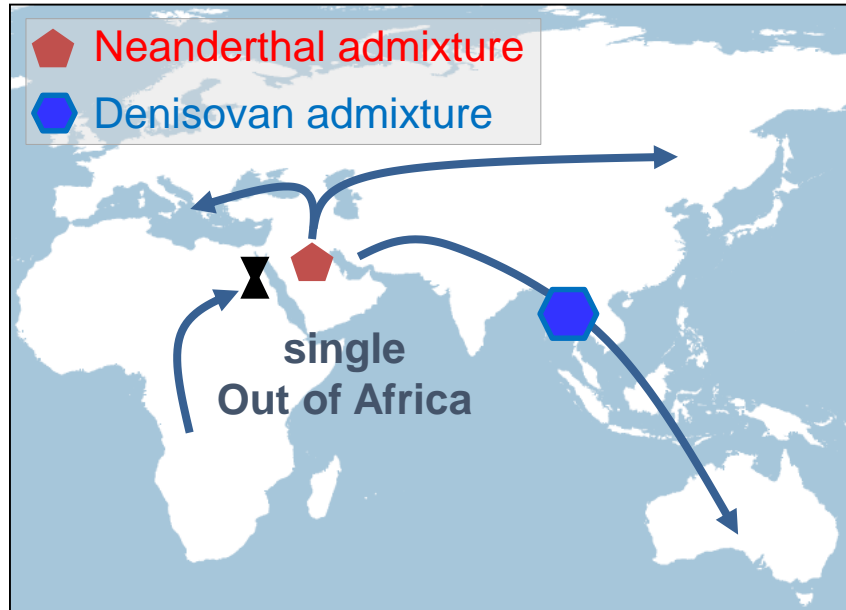
Two-waves out of Africa



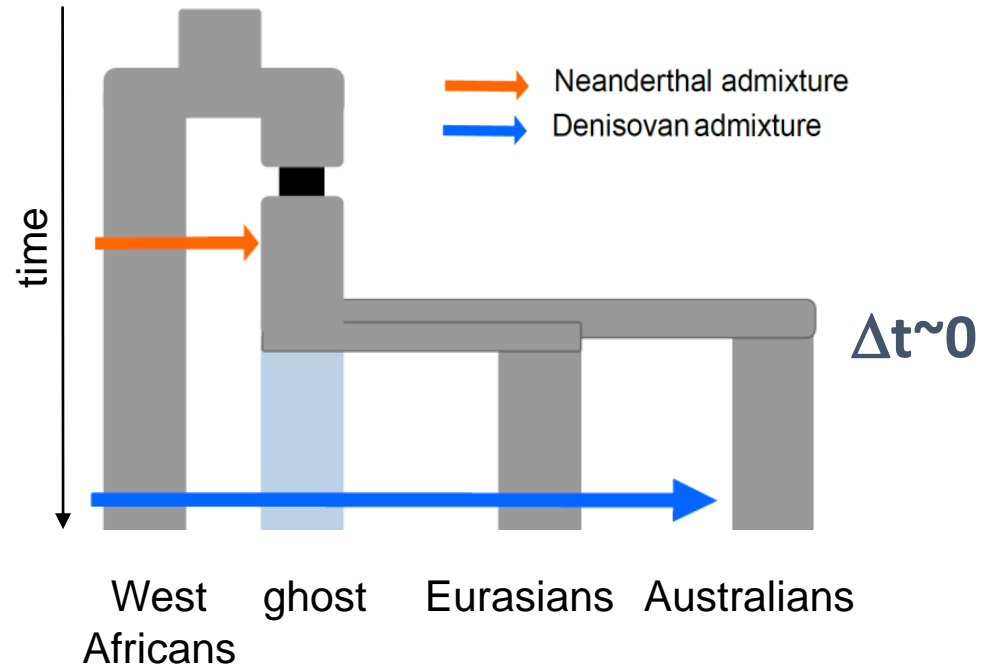
- Two different divergence times ($\Delta t \gg 0$)
- Two independent bottlenecks associated with the two Out of Africa events



One wave out of Africa

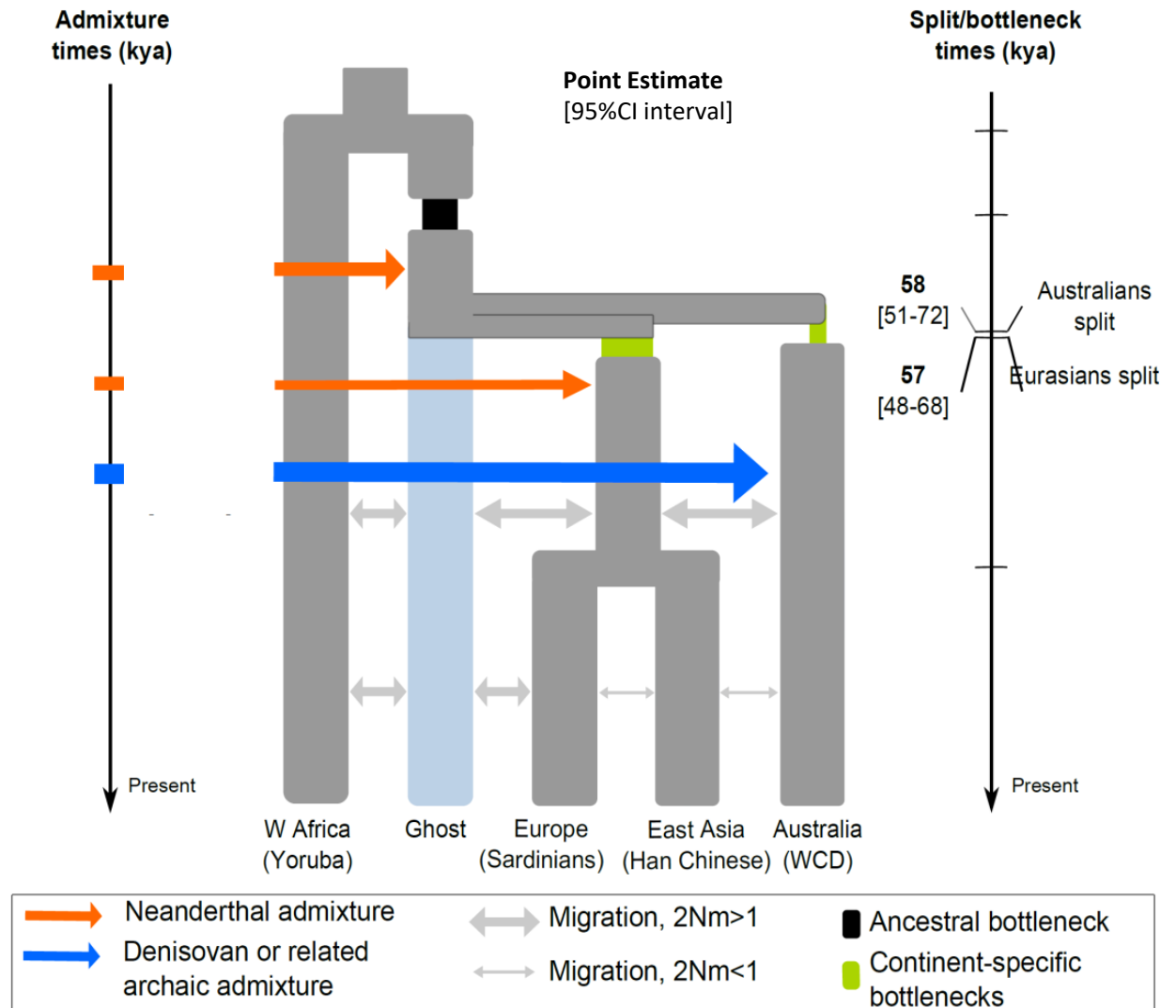


- Similar divergence times (Δt close to zero)
- One single bottlenecks associated with the Out of Africa events
- A major admixture pulse with Neanderthal



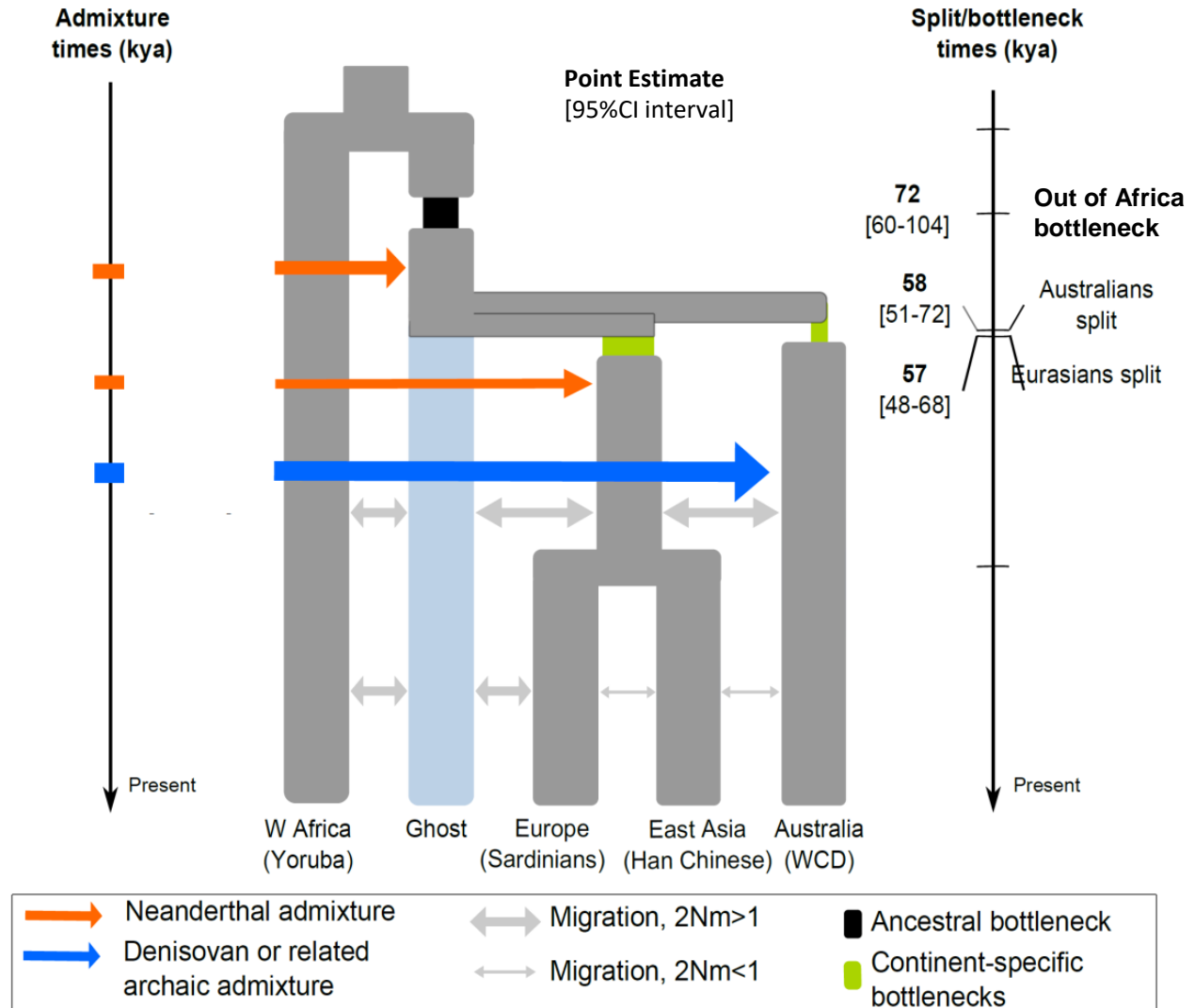
A single wave Out of Africa is consistent with our estimates when accounting for archaic admixture

- Similar divergence time (Δt close to zero)



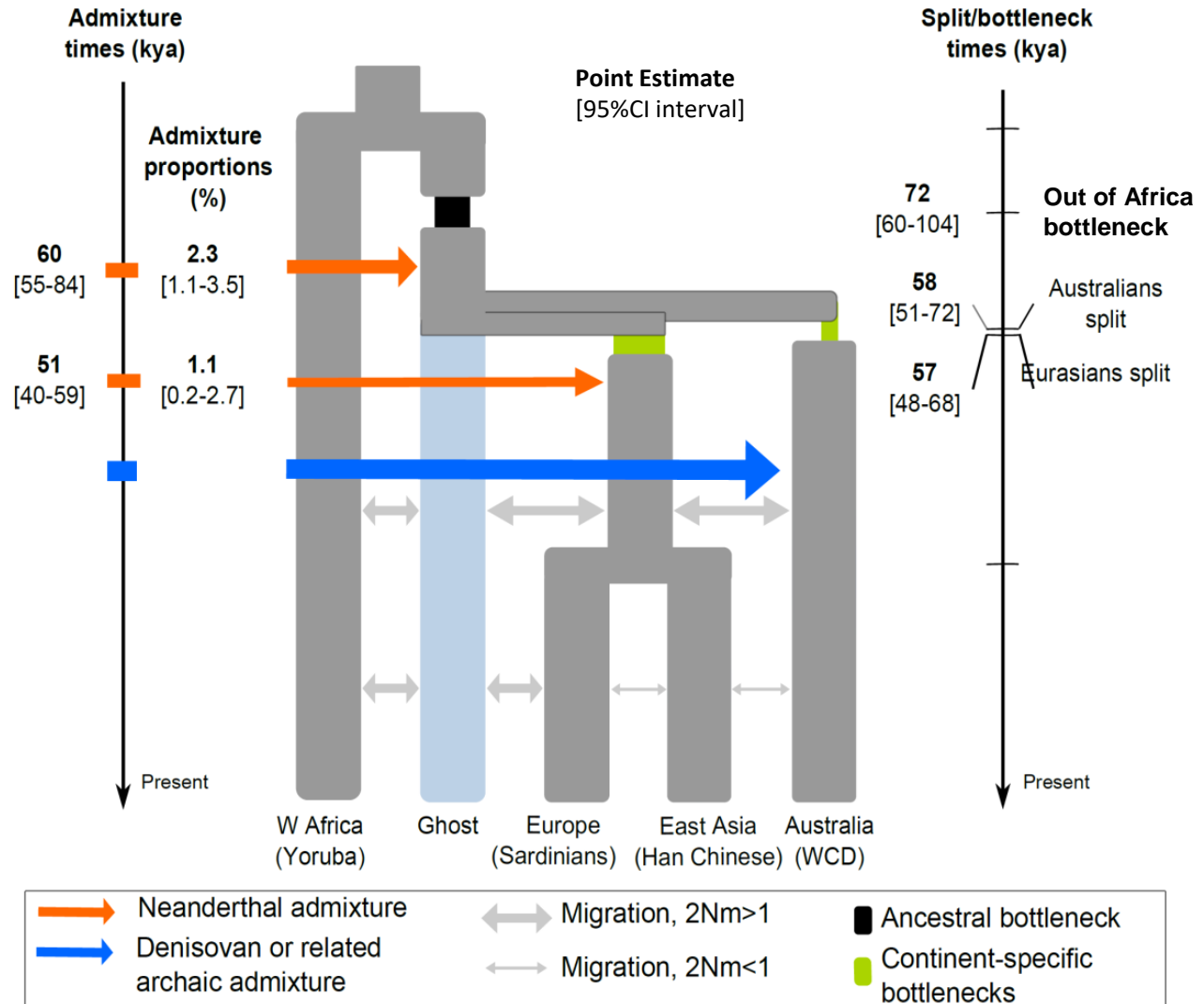
A single wave Out of Africa is consistent with our estimates when accounting for archaic admixture

- Similar divergence time (Δt close to zero)
- Bottleneck associated with the Out of Africa event



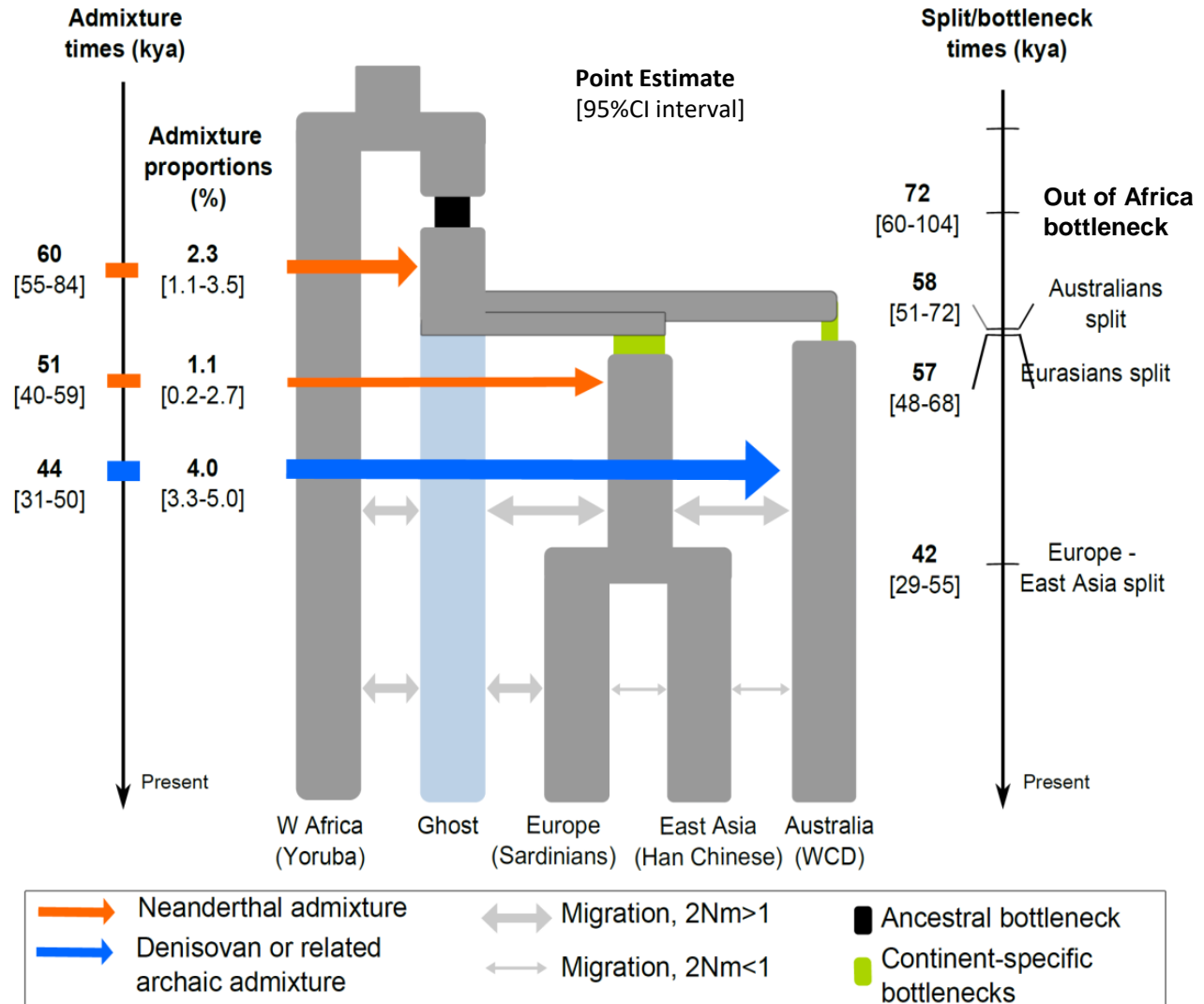
A single wave Out of Africa is consistent with our estimates when accounting for archaic admixture

- Similar divergence time (Δt close to zero)
- Bottleneck associated with the Out of Africa event
- A major admixture pulse with Neanderthal in ancestors of all non-Africans



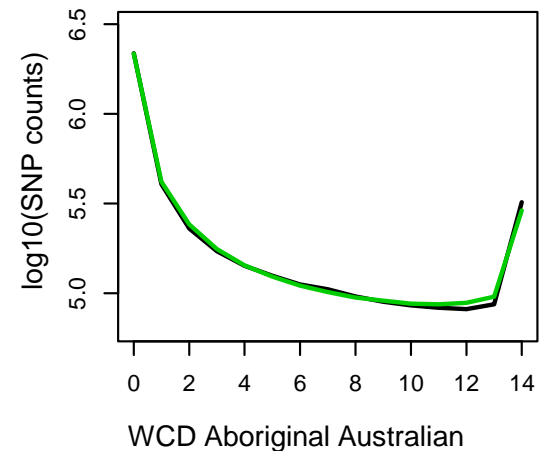
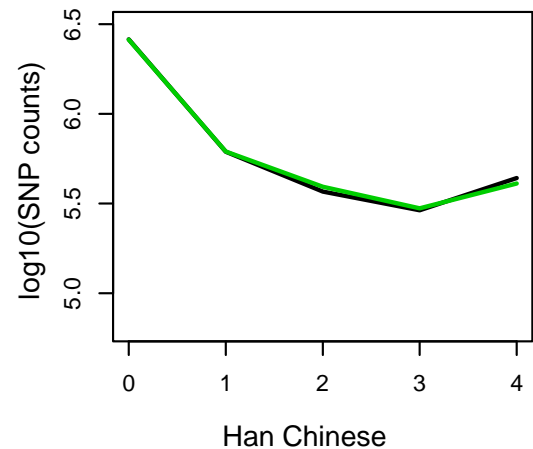
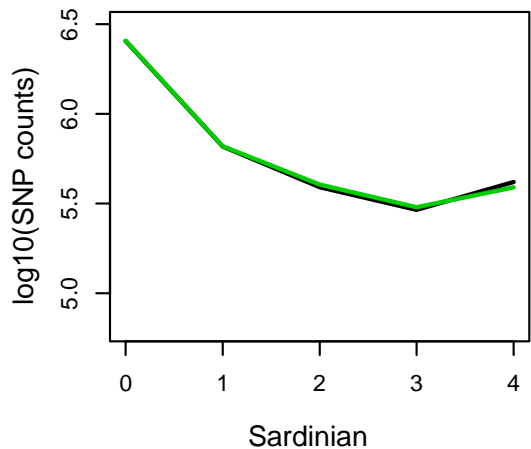
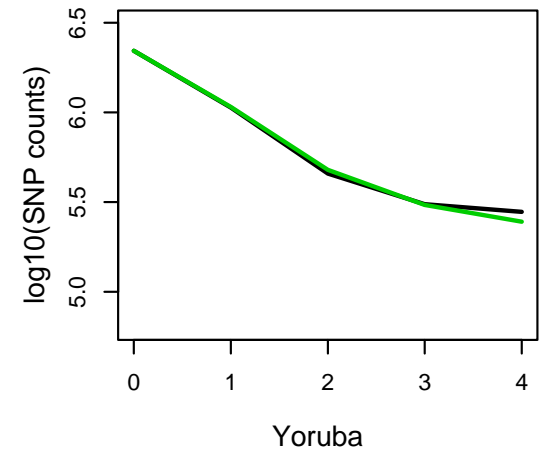
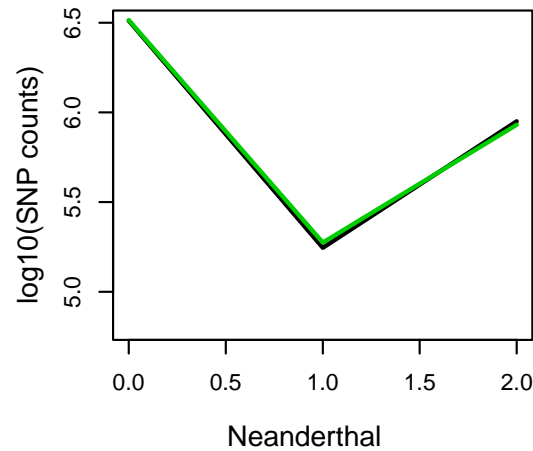
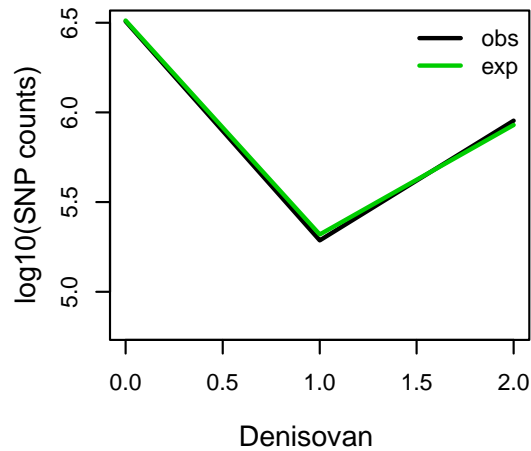
A single wave Out of Africa is consistent with our estimates when accounting for archaic admixture

- Similar divergence time (Δt close to zero)
- Bottleneck associated with the Out of Africa event
- A major admixture pulse with Neanderthal in ancestors of all non-Africans

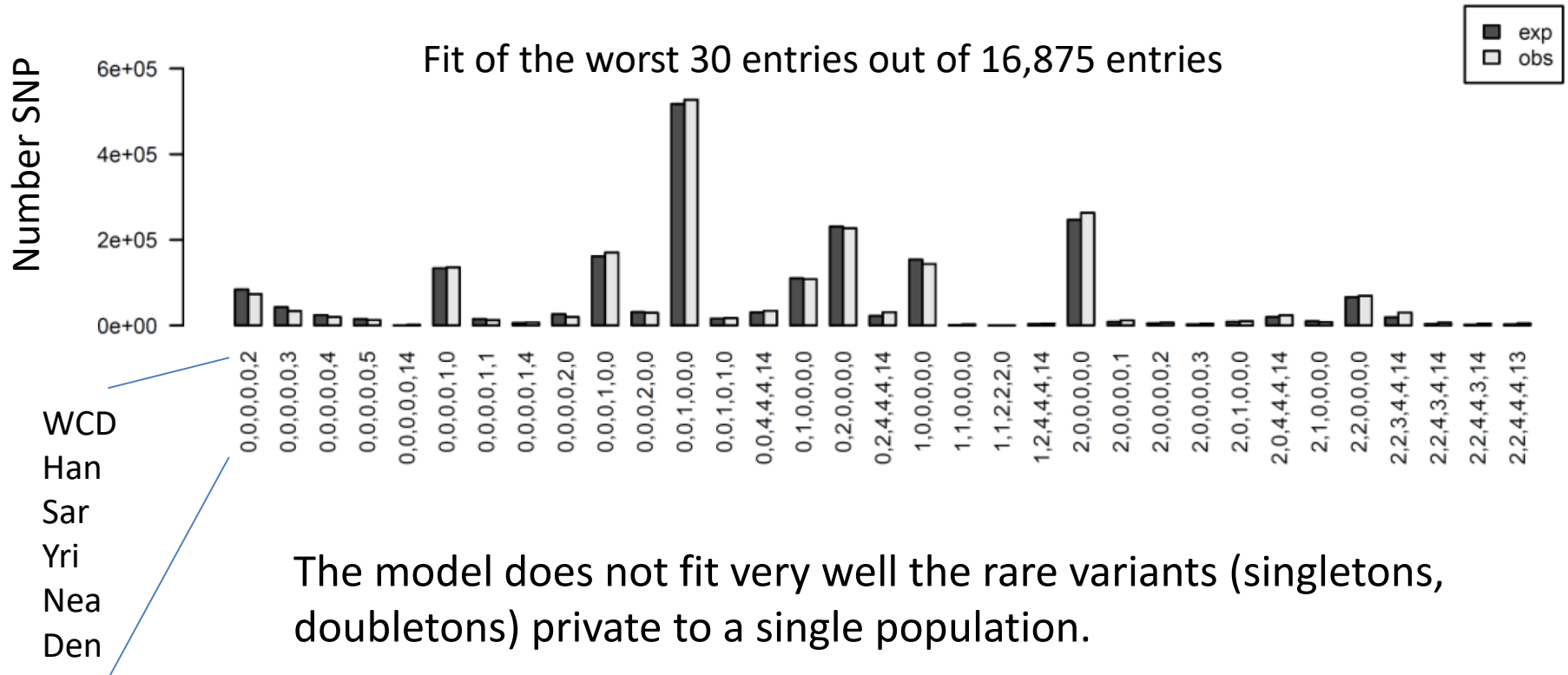


Model captures aspects about the observed data

Good fit to the marginal 1D site frequency spectrum

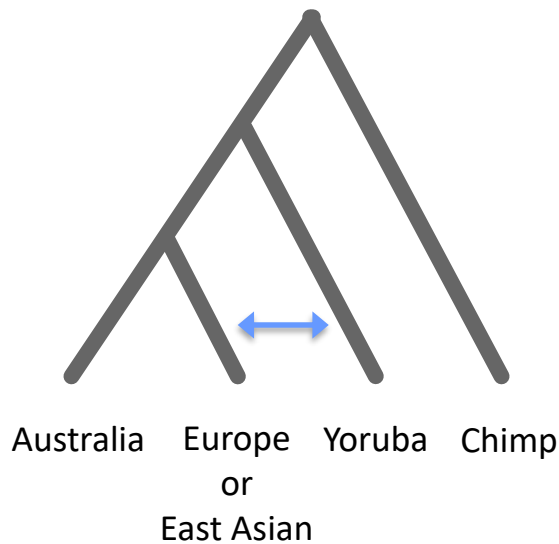


What entries are not well fitted?

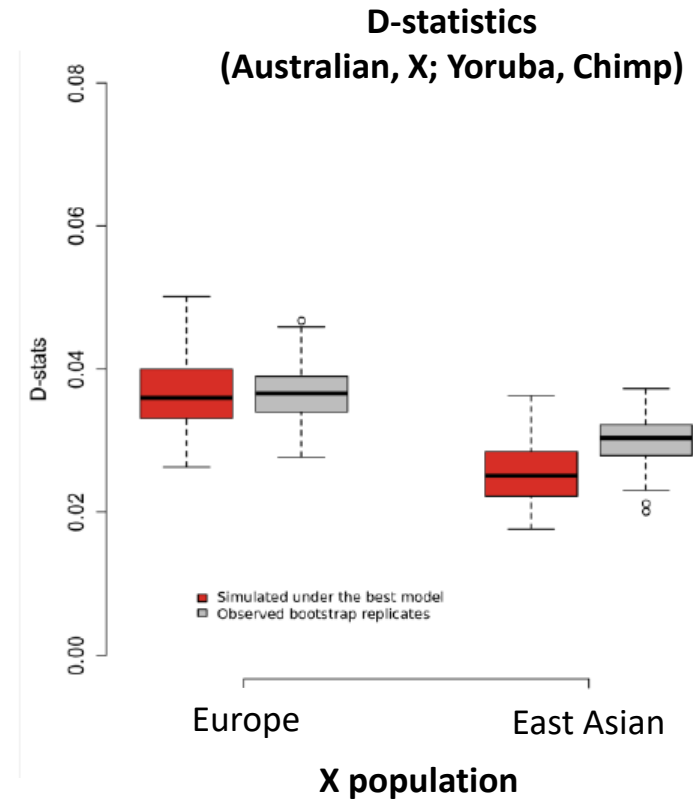


Pagani et al (2016) suggests two waves: Papuan genomes with signature of admixture with humans from first wave (at least 2% of their genome).

Model captures the higher derived allele sharing between Eurasians and Yoruba



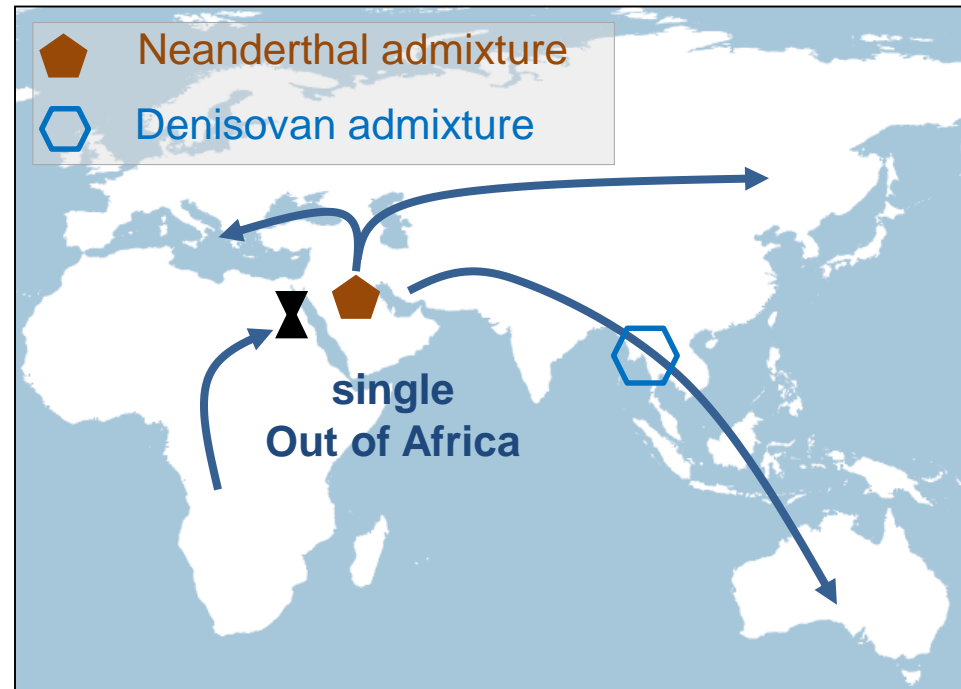
D-statistics suggest that Yoruba and Eurasians share more derived alleles than Yoruba and Australians



Summary

Aboriginal Australians genomes support a single major wave out of Africa

- Accounting for archaic admixture with Neanderthal and Denisovan was crucial to understand population divergence
- Genomic data consistent with a single major dispersal event out of Africa (60-104 kya)
- Two major dispersal waves into Asia: Aboriginal Australians diverged 51-72 kya from Eurasians



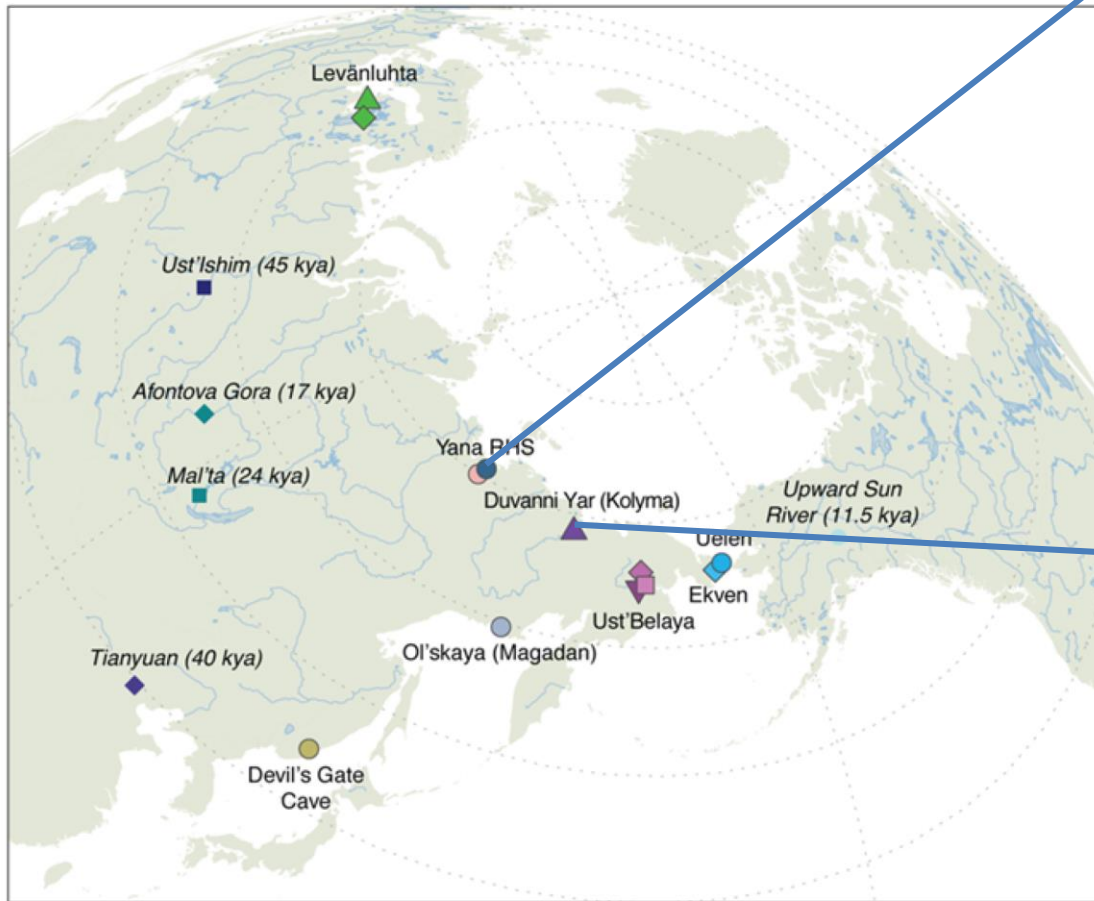
The population history of northeastern Siberia since the Pleistocene

Martin Sikora^{1,43*}, Vladimir V. Pitulko^{2,43*}, Vitor C. Sousa^{3,4,5,43}, Morten E. Allentoft^{1,43}, Lasse Vinner¹, Simon Rasmussen^{6,41}, Ashot Margaryan¹, Peter de Barros Damgaard¹, Constanza de la Fuente^{1,42}, Gabriel Renaud¹, Melinda A. Yang⁷, Qiaomei Fu⁷, Isabelle Dupanloup⁸, Konstantinos Giampoudakis⁹, David Nogués-Bravo⁹, Carsten Rahbek⁹, Guus Kroonen^{10,11}, Michaël Peyrot¹¹, Hugh McColl¹, Sergey V. Vasilyev¹², Elizaveta Veselovskaya^{12,13}, Margarita Gerasimova¹², Elena Y. Pavlova^{2,14}, Vyacheslav G. Chasnyk¹⁵, Pavel A. Nikolskiy^{2,16}, Andrei V. Gromov¹⁷, Valeriy I. Khartanovich¹⁷, Vyacheslav Moiseyev¹⁷, Pavel S. Grebenyuk^{18,19}, Alexander Yu. Fedorchenko²⁰, Alexander I. Lebedintsev¹⁸, Sergey B. Slobodin¹⁸, Boris A. Malyarchuk²¹, Rui Martiniano²², Morten Meldgaard^{1,23}, Laura Arppe²⁴, Jukka U. Palo^{25,26}, Tarja Sundell^{27,28}, Kristiina Mannermaa²⁷, Mikko Putkonen²⁵, Verner Alexandersen²⁹, Charlotte Primeau²⁹, Nurbol Baimukhanov³⁰, Ripan S. Malhi^{31,32}, Karl-Göran Sjögren³³, Kristian Kristiansen³³, Anna Wessman^{27,34}, Antti Sajantila²⁵, Marta Mirazon Lahr^{1,35}, Richard Durbin^{22,36}, Rasmus Nielsen^{1,37}, David J. Meltzer^{1,38}, Laurent Excoffier^{4,5*} & Eske Willerslev^{1,36,39,40*}

Nature (2019)



Colonization of Siberia



Yana RHS (31,600 years ago)
Whole-genome depth of coverage 25x



Kolyma (9,800 years ago)
Whole-genome depth of coverage 14x



Hypothesis: Continuity vs Replacement of populations

Data: Ancient and present-day samples; 625 blocks of 1Mb (~1.5 Million SNP), far from genic regions and CpG islands

Method: Composite likelihood - *fastsimcoal2*
(Excoffier et al, 2013 Plos Genetics)

Europe (Sardinia)	Ancient North Siberians (Yana)	Ancient Paleo- siberian (Kolyma)	Neo- siberian (Even)	East Asia (Han)
----------------------	---	---	----------------------------	-----------------------

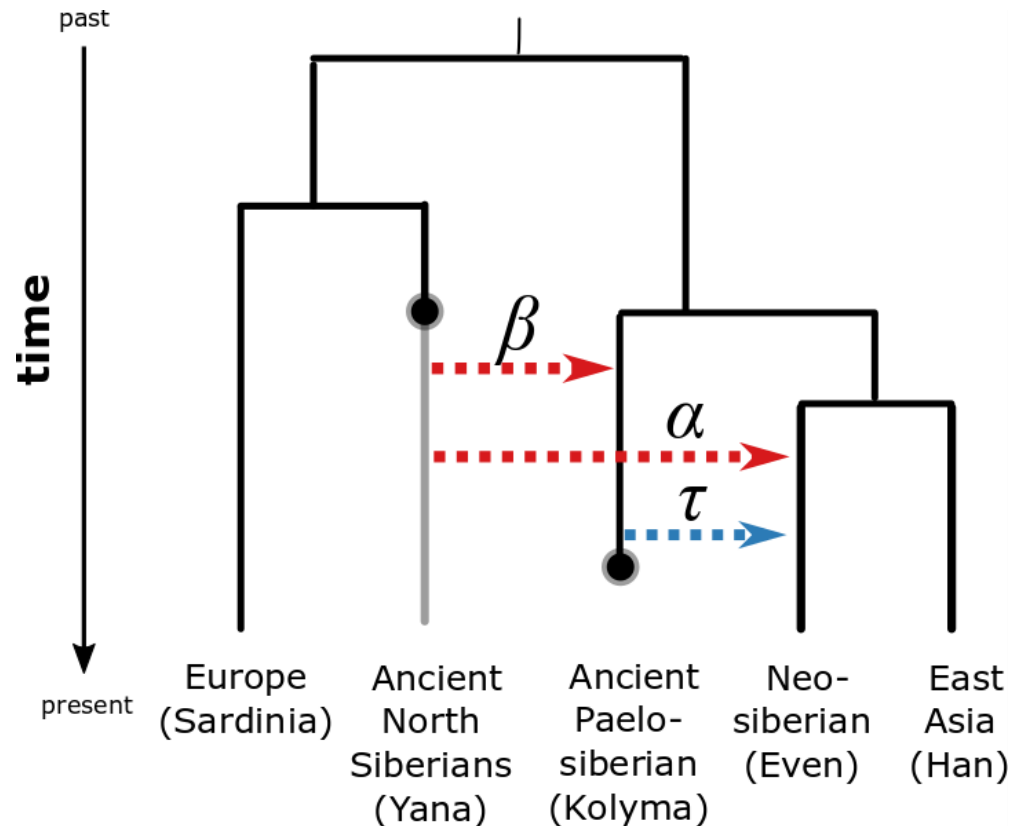


Hypothesis: Continuity vs Replacement of populations

For instance:

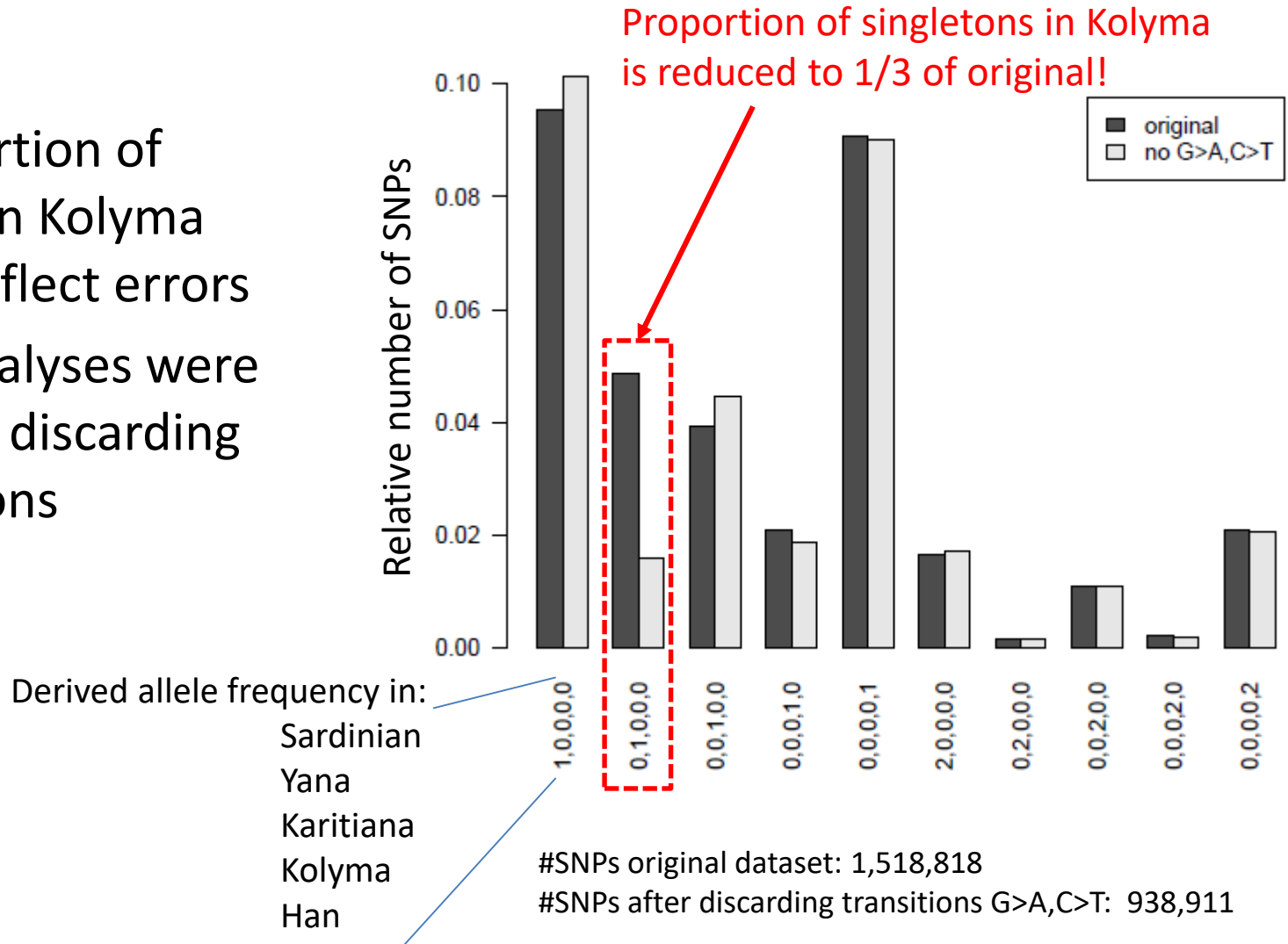
$\beta = 1$ indicates continuity:
Kolyma descends from Yana

$\beta = 0$ indicates replacement
of Yana by Kolyma

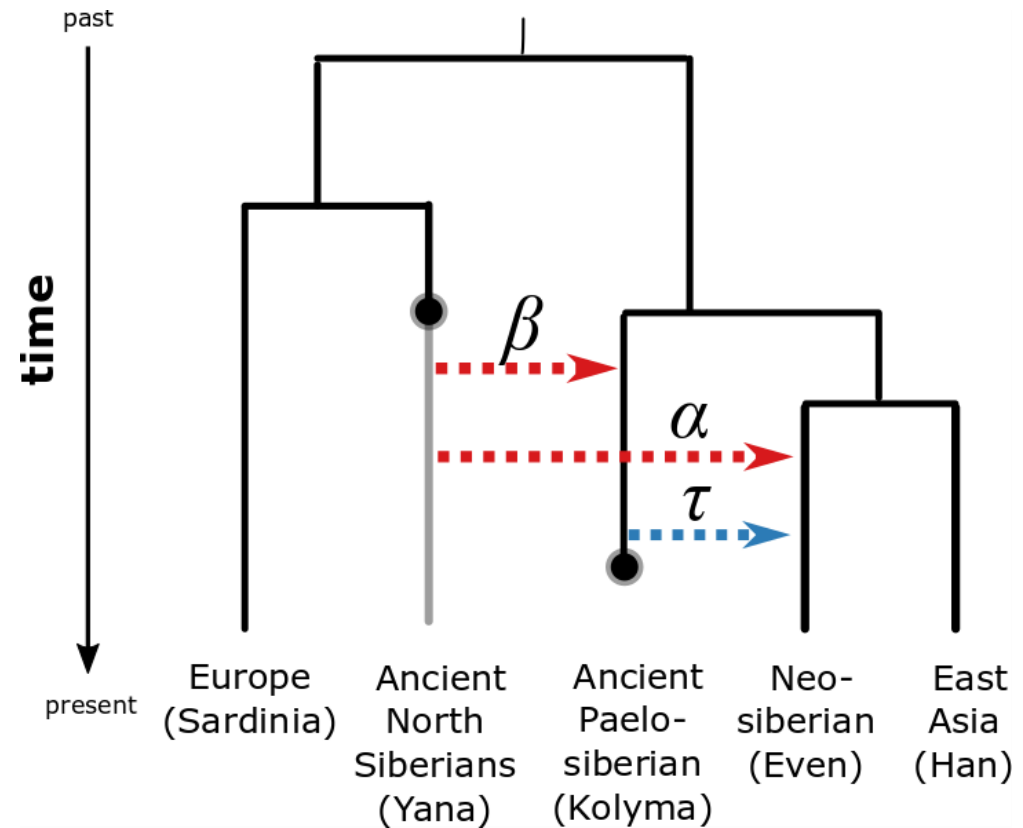
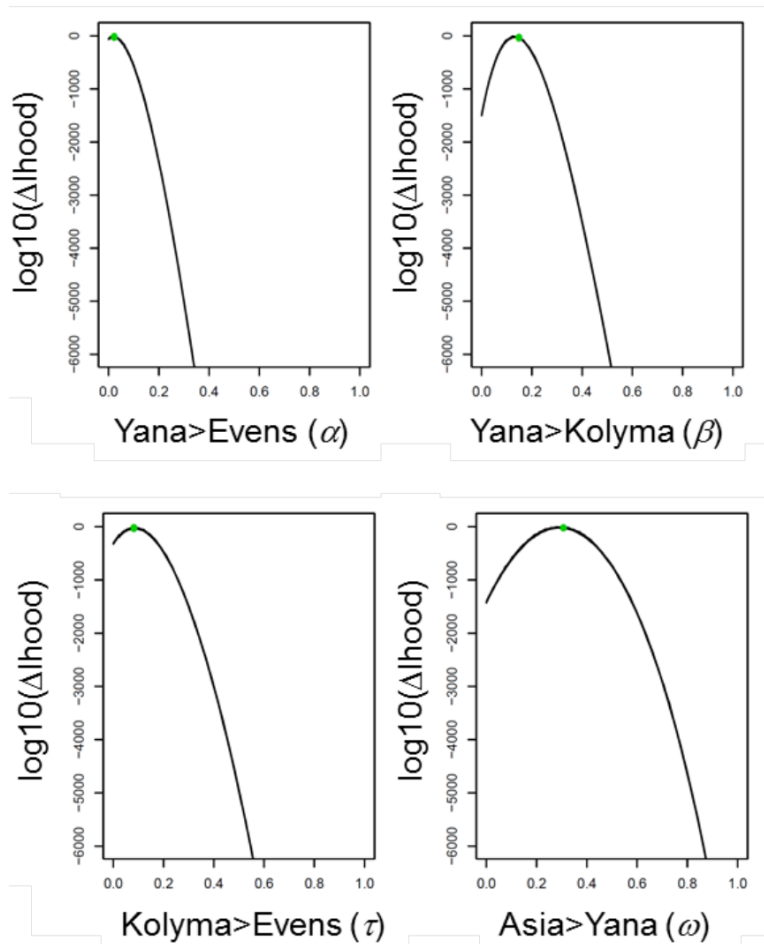


Site frequency spectrum is affected by damage patterns in ancient DNA

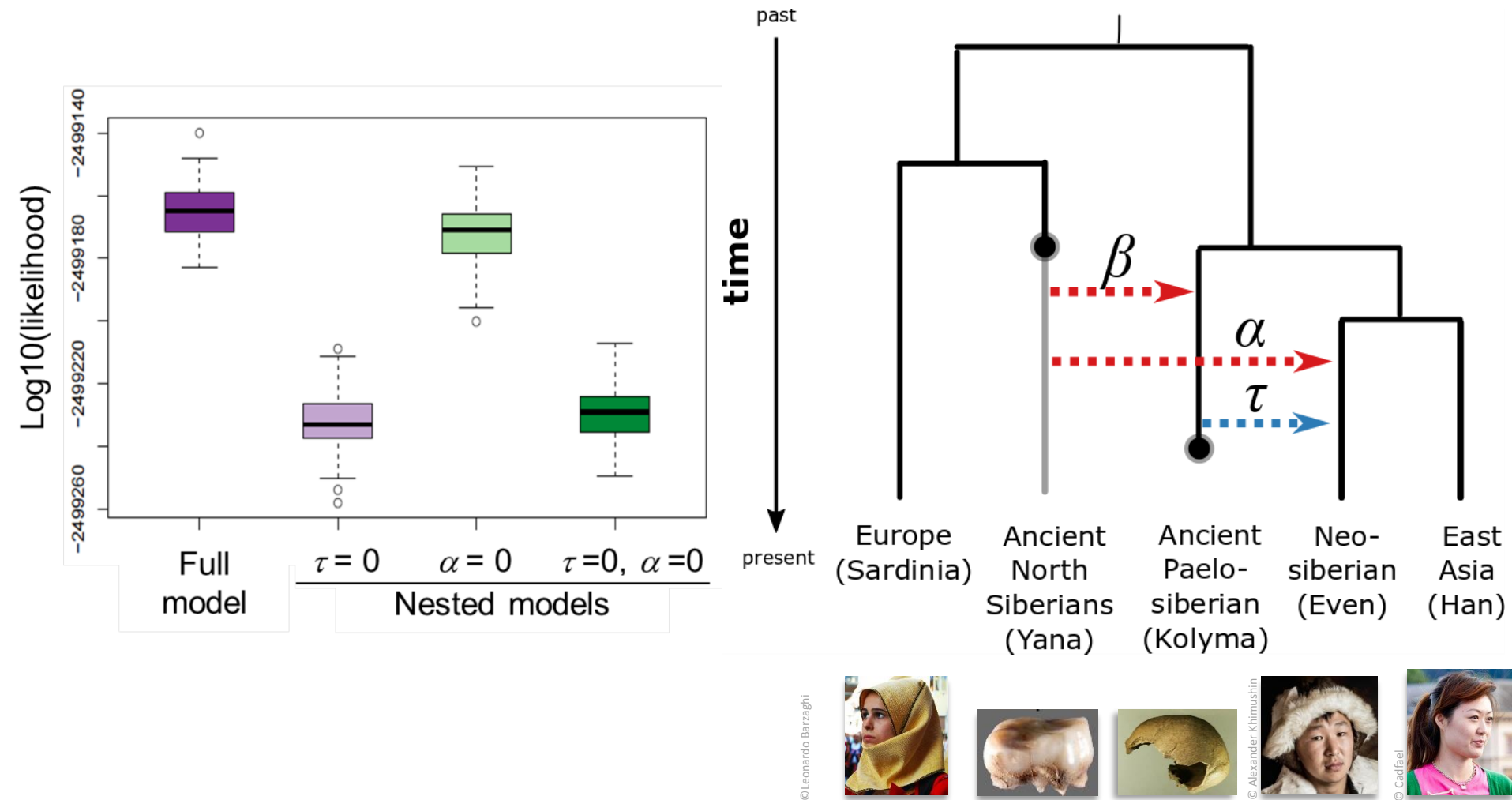
- High proportion of singletons in Kolyma probably reflect errors
- Thus, all analyses were performed discarding the singletons



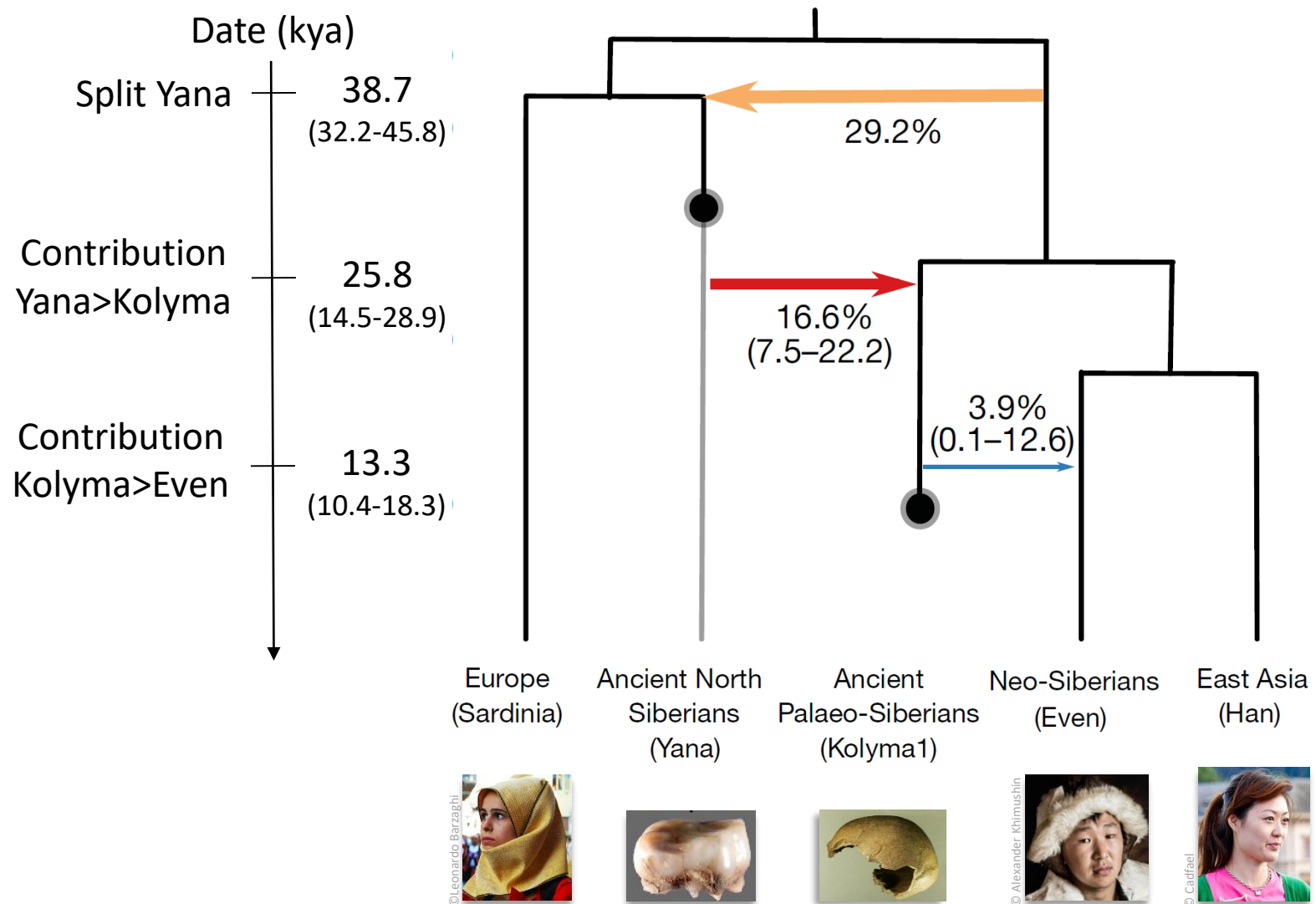
Model comparison and likelihood profiles consistent with replacement with gene flow



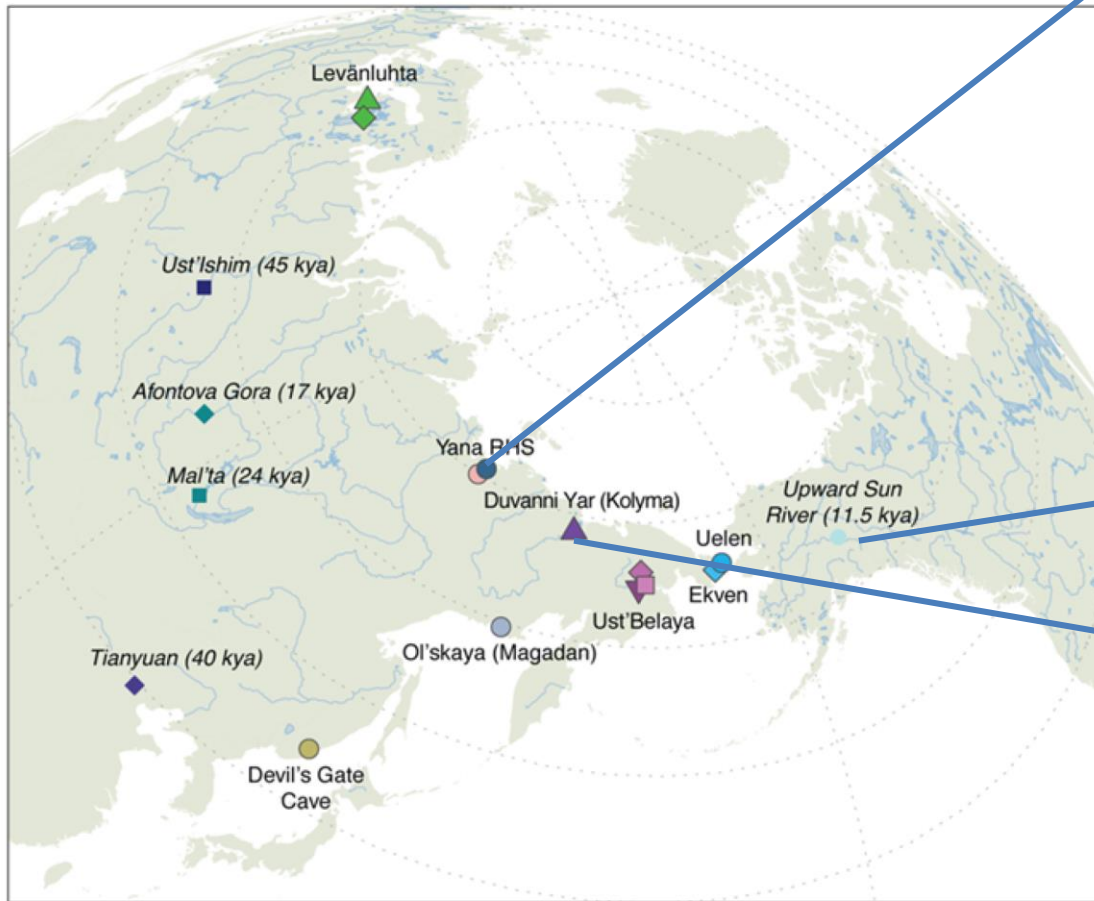
Model comparison and likelihood profiles consistent with replacement with gene flow



Estimates of best nested model indicate replacement with gene flow



Siberia and colonization of the Americas



Yana RHS (31,600 years ago)
Whole-genome depth of coverage 25x

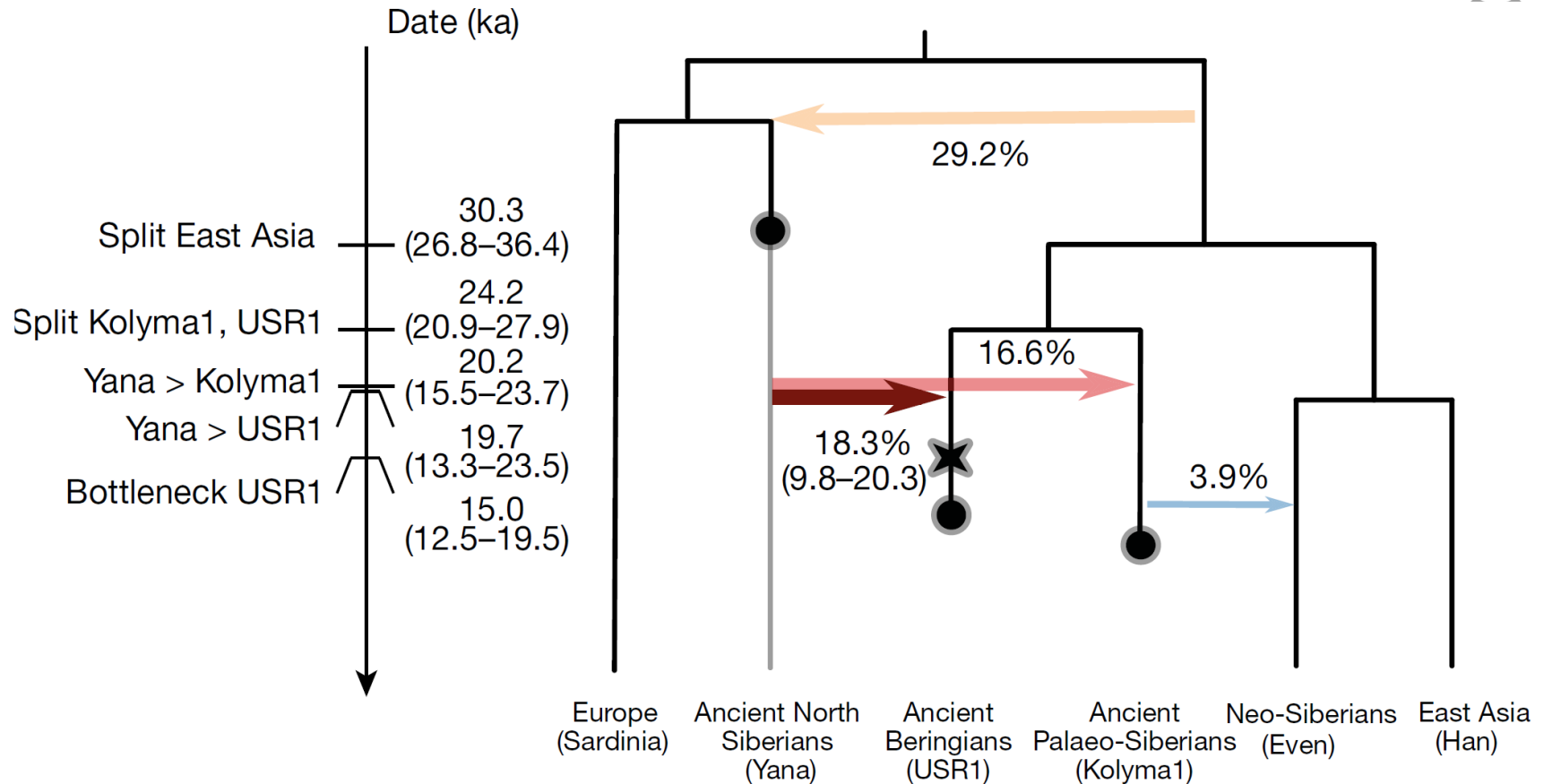


USR1 (11,500 years ago) Alaska

Kolyma (9,800 years ago)
Whole-genome depth of coverage 14x

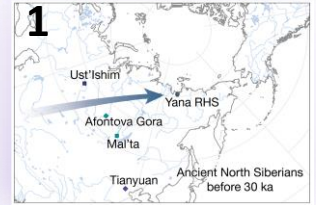


Estimates consistent with replacement with gene flow

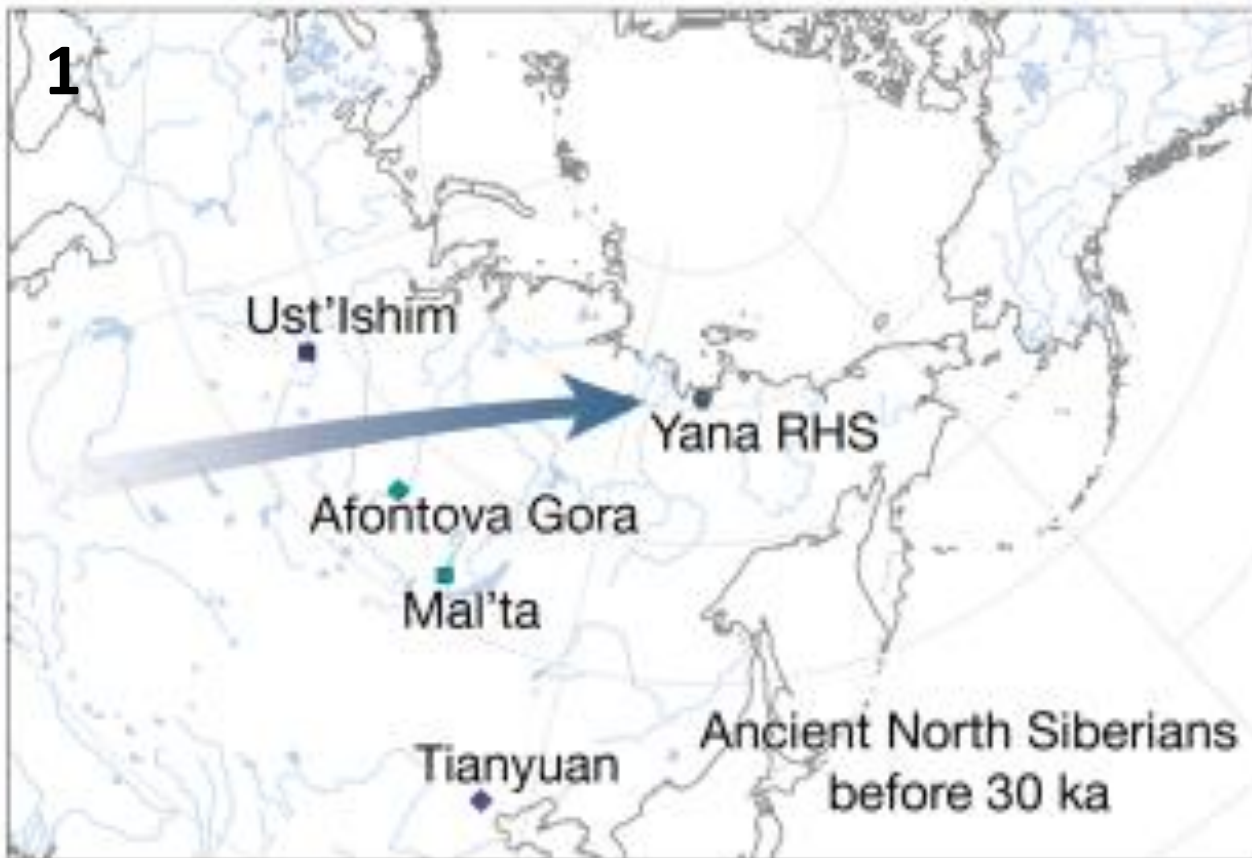


- Kolyma is the closest population to Native Americans (USR1 and Karitiana)
- Native Americans with a contribution of up to 20% from Yana

Summary: 3 migration waves



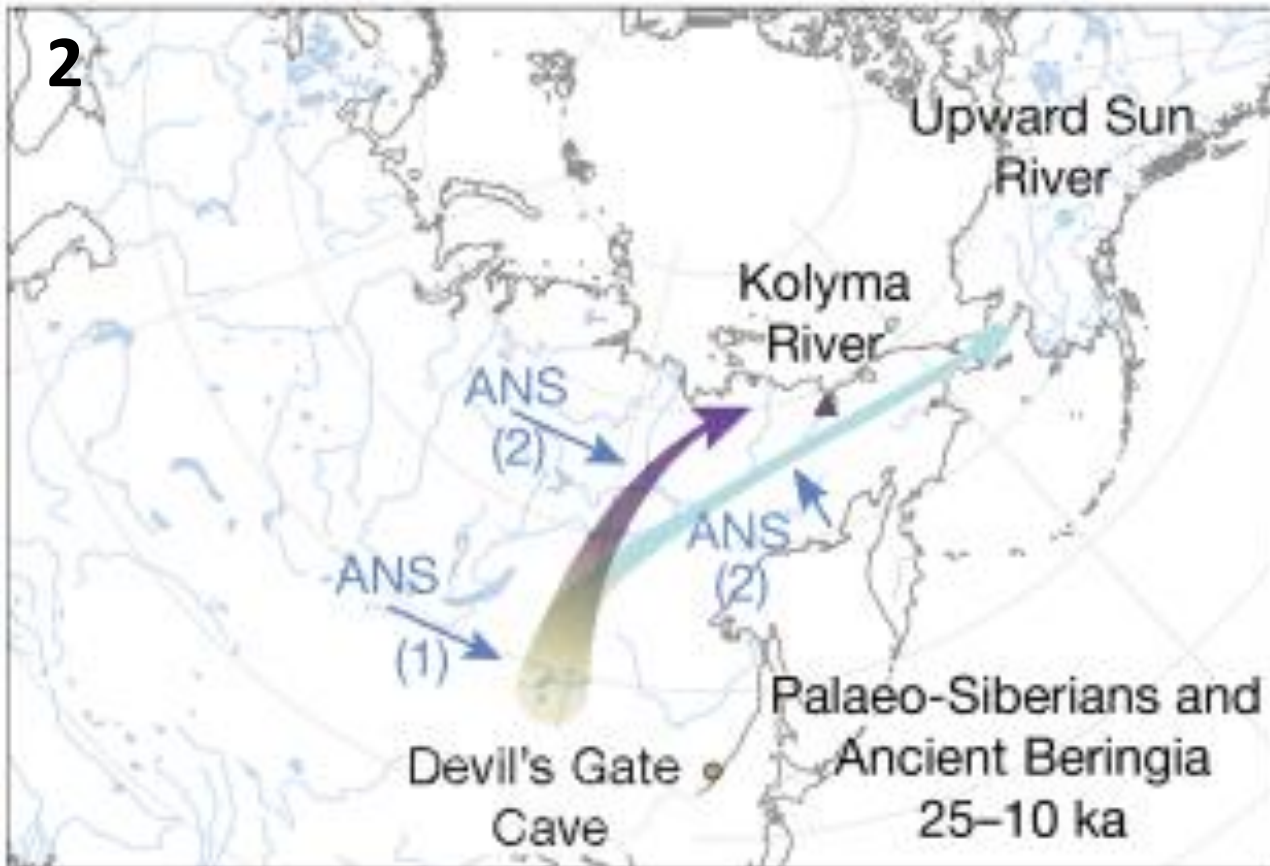
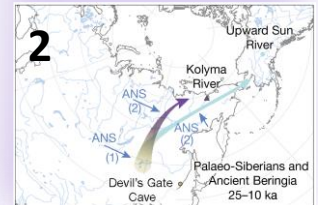
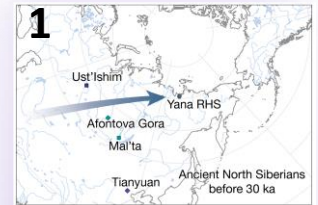
- Ancient North Siberians (Yana) reached Siberia before 30 ka (thousand-years ago)



1st migration wave

Summary: 3 migration waves

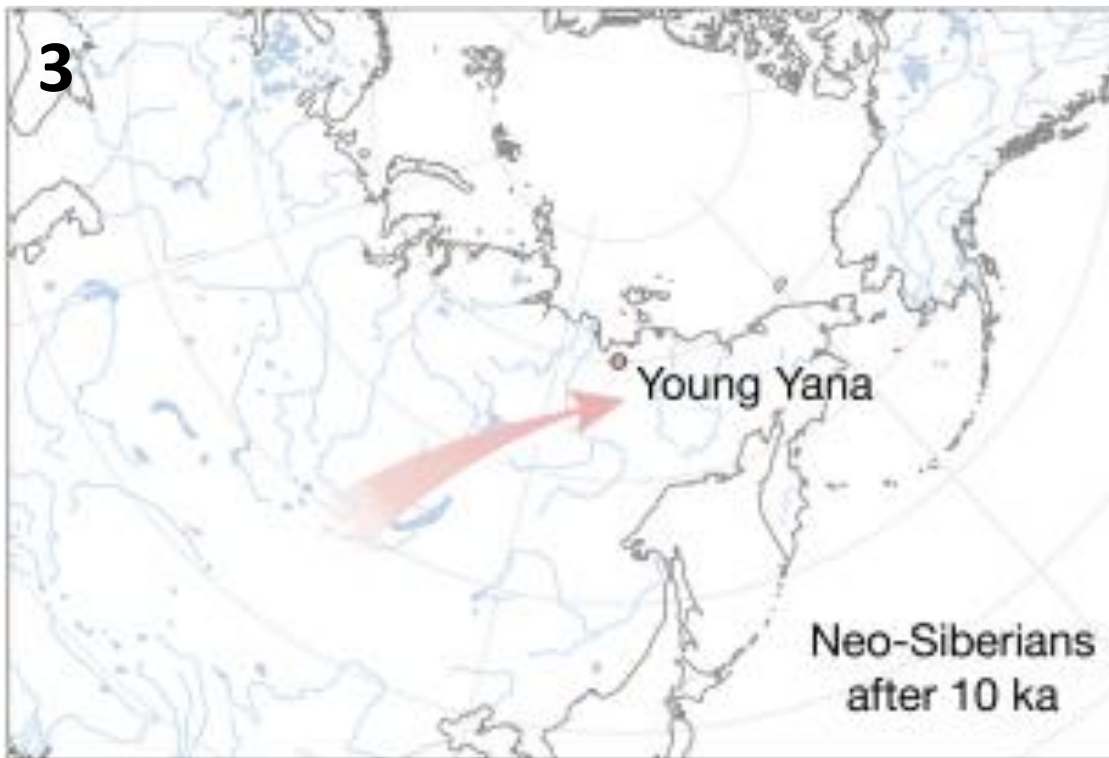
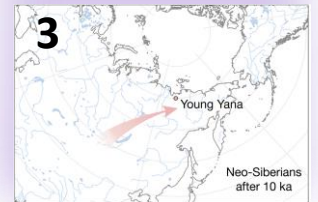
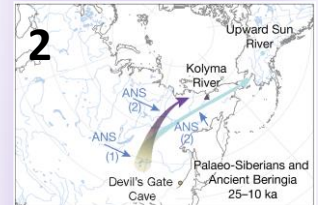
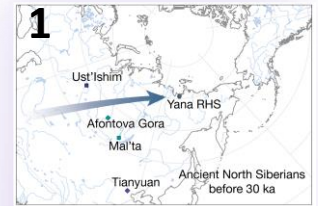
- Ancient North Siberians (Yana) reached Siberia before 30 kya
- Paleo-Siberians (Kolyma) migrated after Last Glacial Maximum (26.5 ka)
- Native-Americans are closer to Kolyma, with 20% of Yana contribution



2nd migration wave

Summary: 3 migration waves

- Ancient North Siberians (Yana) reached Siberia before 30 ka
- Paleo-Siberians (Kolyma) likely migrated after Last Glacial Maxima
- Native-Americans are closer to Kolyma, with 20% of Yana contribution
- Paleo-Siberians (Kolyma) were replaced by Neo-Siberians, likely associated with the cooler period “Younger Dryas” (12.8-11.5 ka)



3rd migration wave

Deer mice from Nebraska Sand Hills



S. Pfeifer, S. Laurent, V. Sousa, C. Linnen, H. Hoekstra, L. Excoffier, J. Jensen

Coat color adaptation in deer mice *Peromyscus maniculatus*

- Habitat (soil color) correlated with coat phenotype
- Field experiments suggest that light color confers selective advantage against visually hunting predators
- Nebraska Sand Hills were formed 8000 to 15,000 years ago



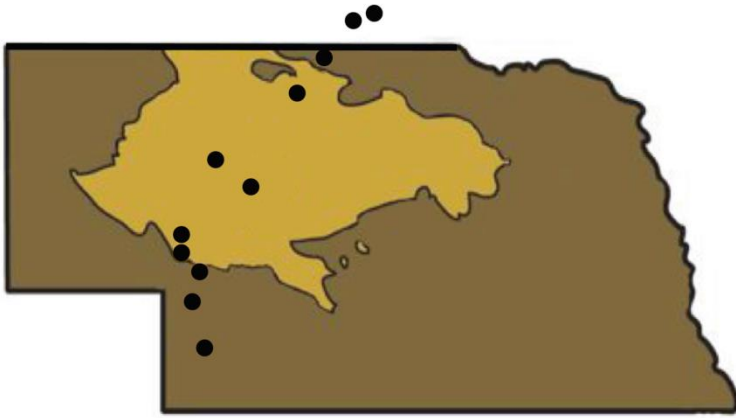
Linnen et al (2013) Science

Pfeifer*, Laurent*, Sousa* et al (in press) MBE

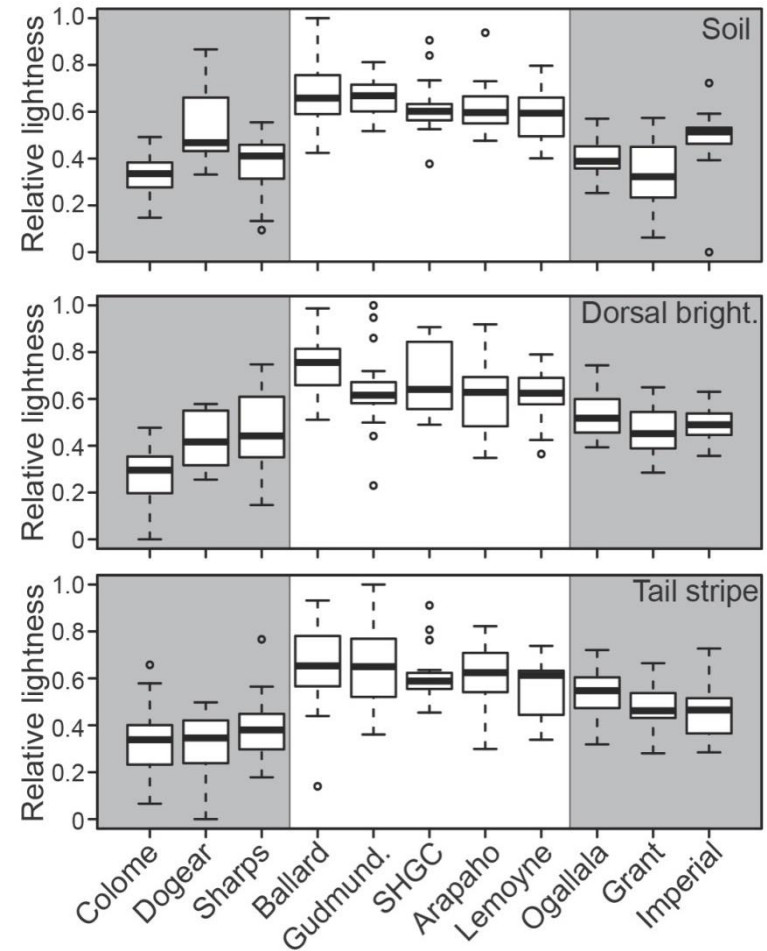
A transect across the Sand Hills (ON and OFF)

Sample locations “off” and “on” the Sand Hills

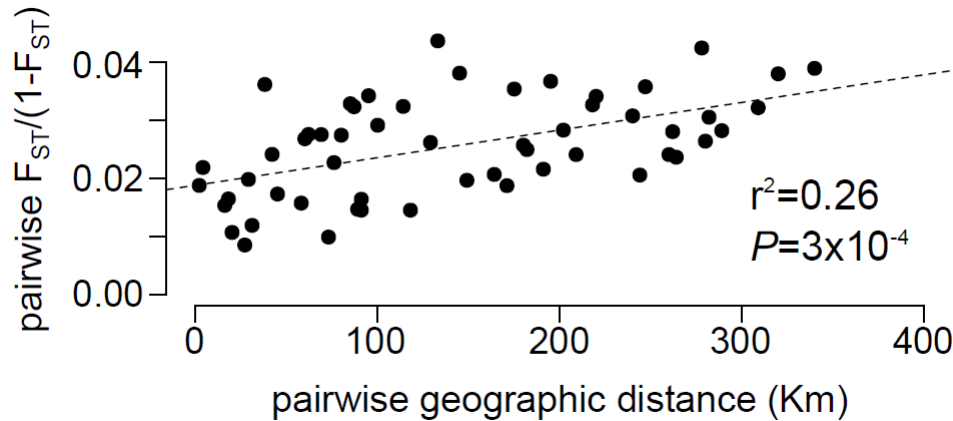
- 11 populations
- 330 individuals



- Genomic data (NGS) data
 - Target 10,000 random 1.5kb regions
 - 185kbp region comprising the *Agouti* gene
- Phenotypic data for each individual

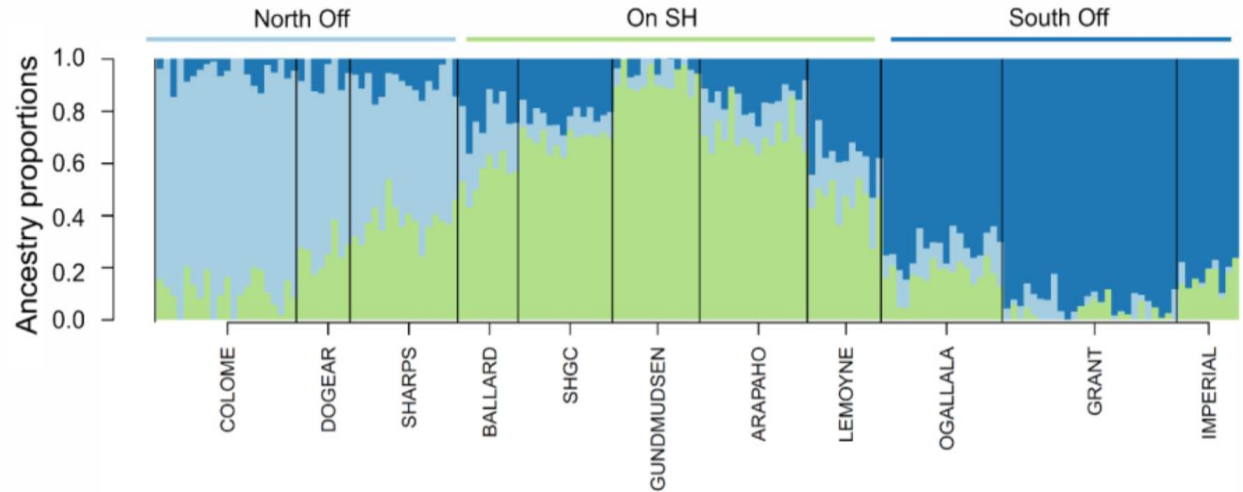
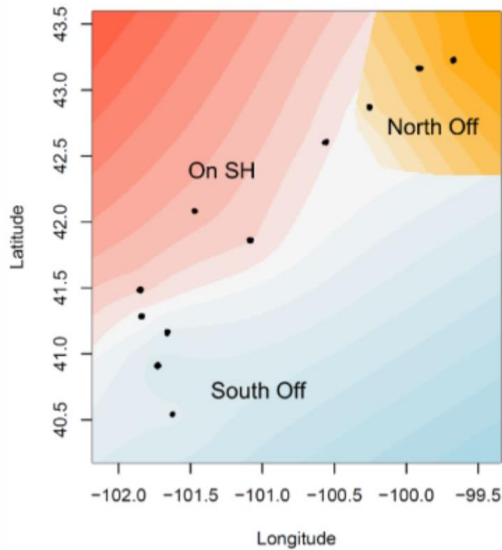


Evidence for isolation by distance but three groups



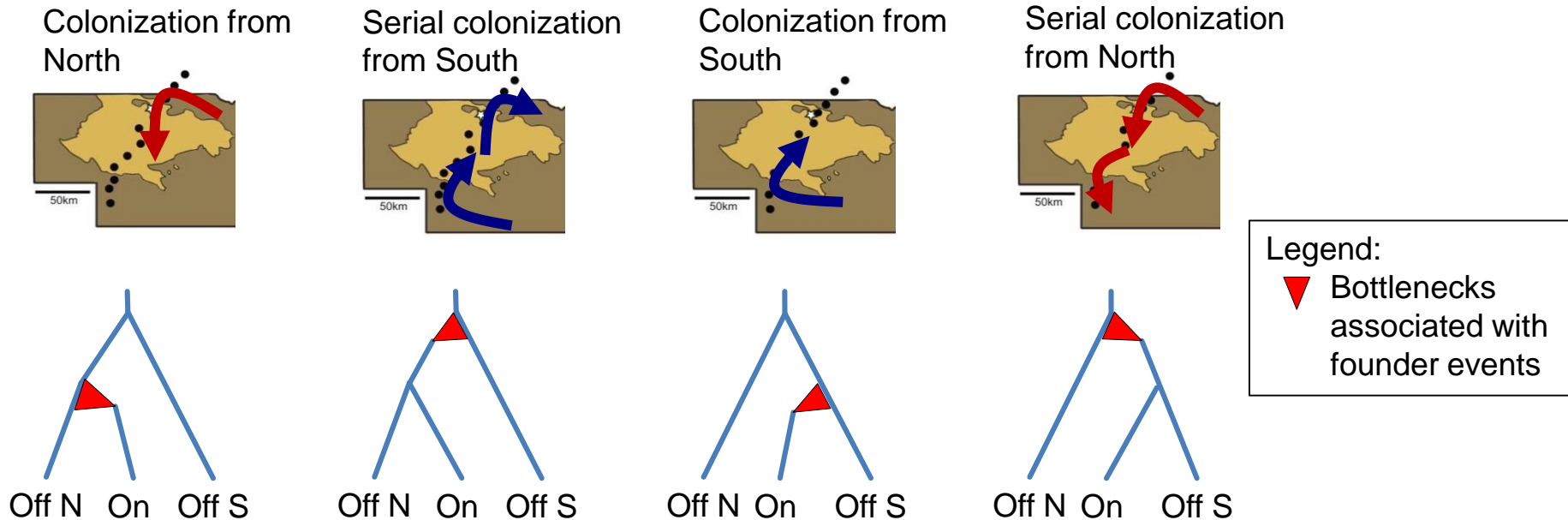
Geographically closer samples are genetically more similar

Ancestry coefficients



Model-based inference

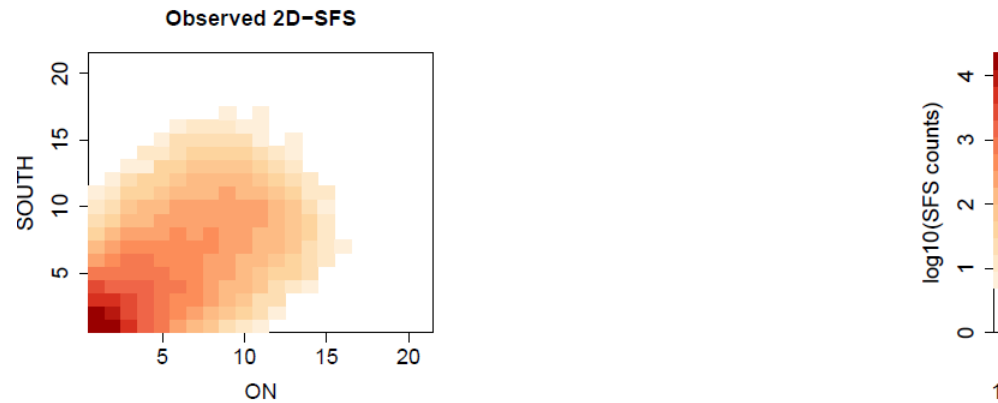
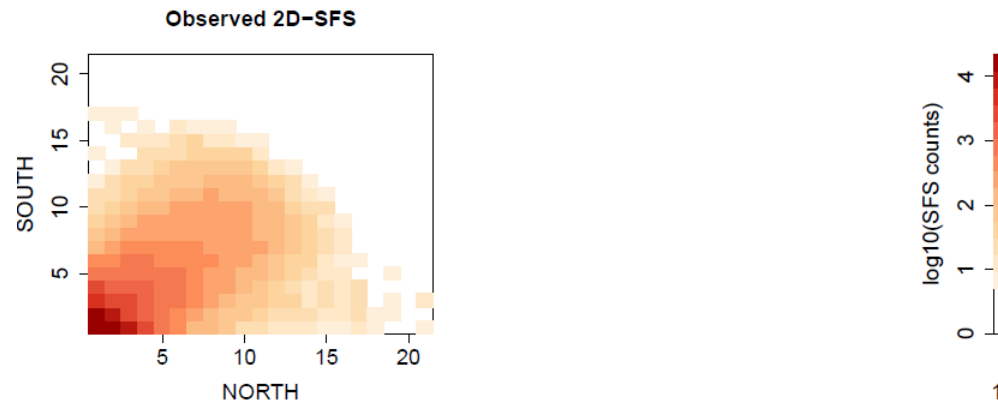
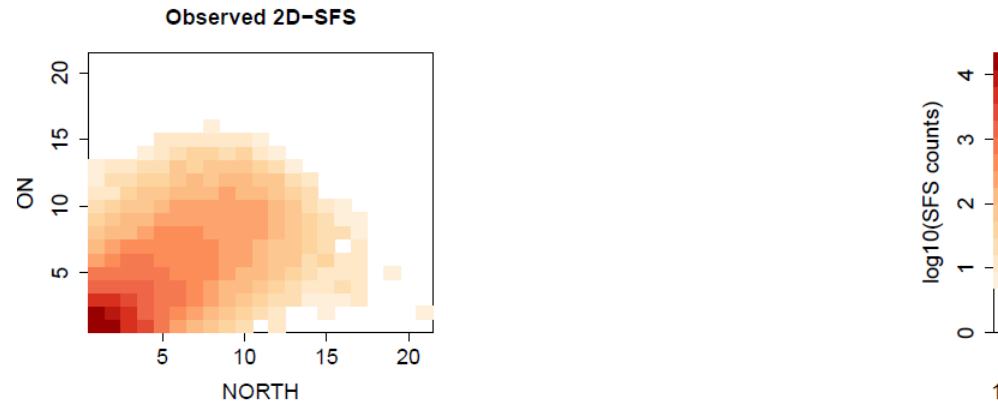
Is there evidence of gene flow between Off and On the Sand Hills?



Estimates based on the joint **3D site frequency spectrum (SFS)**:
- folded SFS with 140,358 SNPs

Deer mice: Pairwise marginal 2D SFS

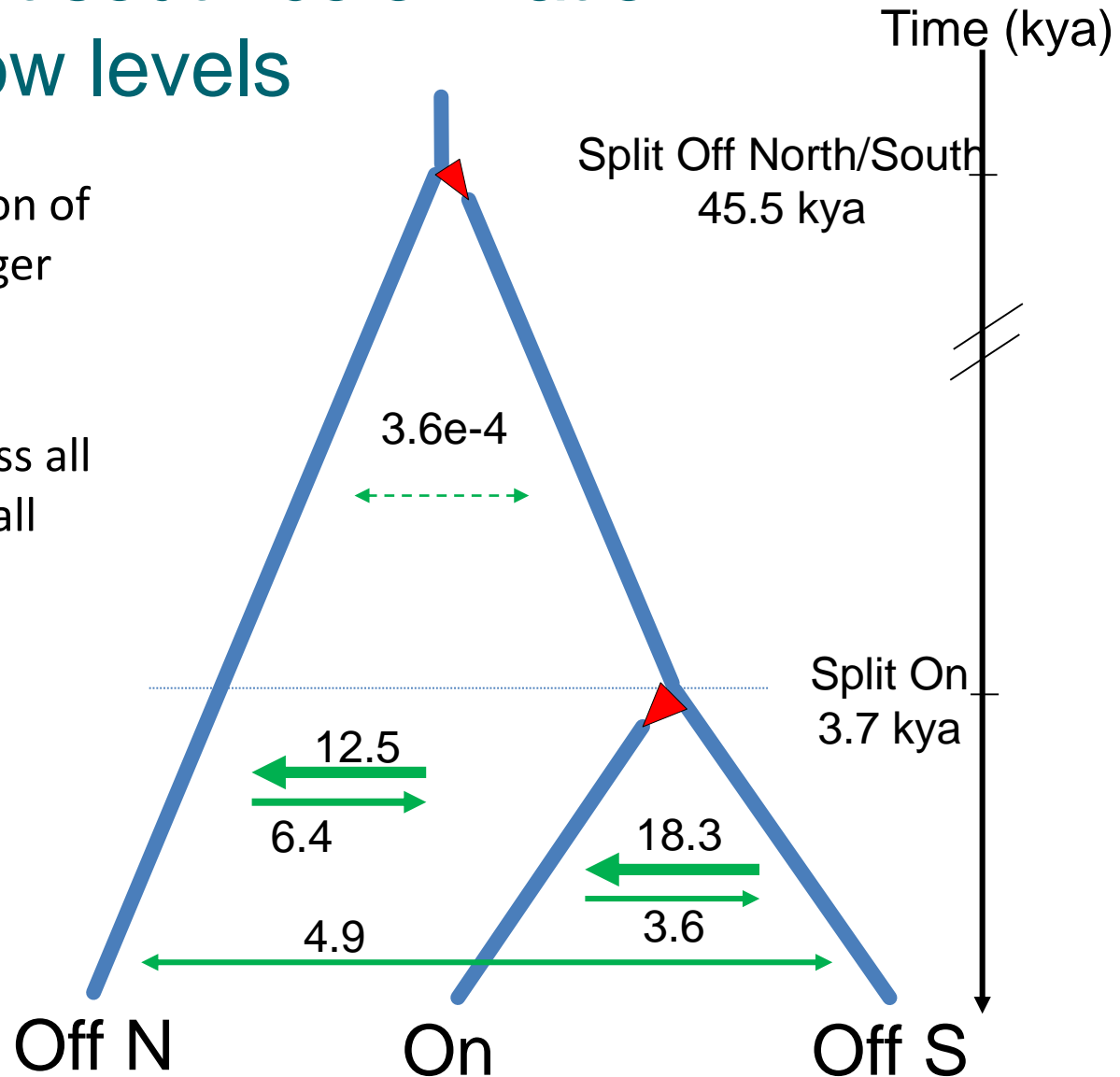
Since we did not have an outgroup we used the folded SFS



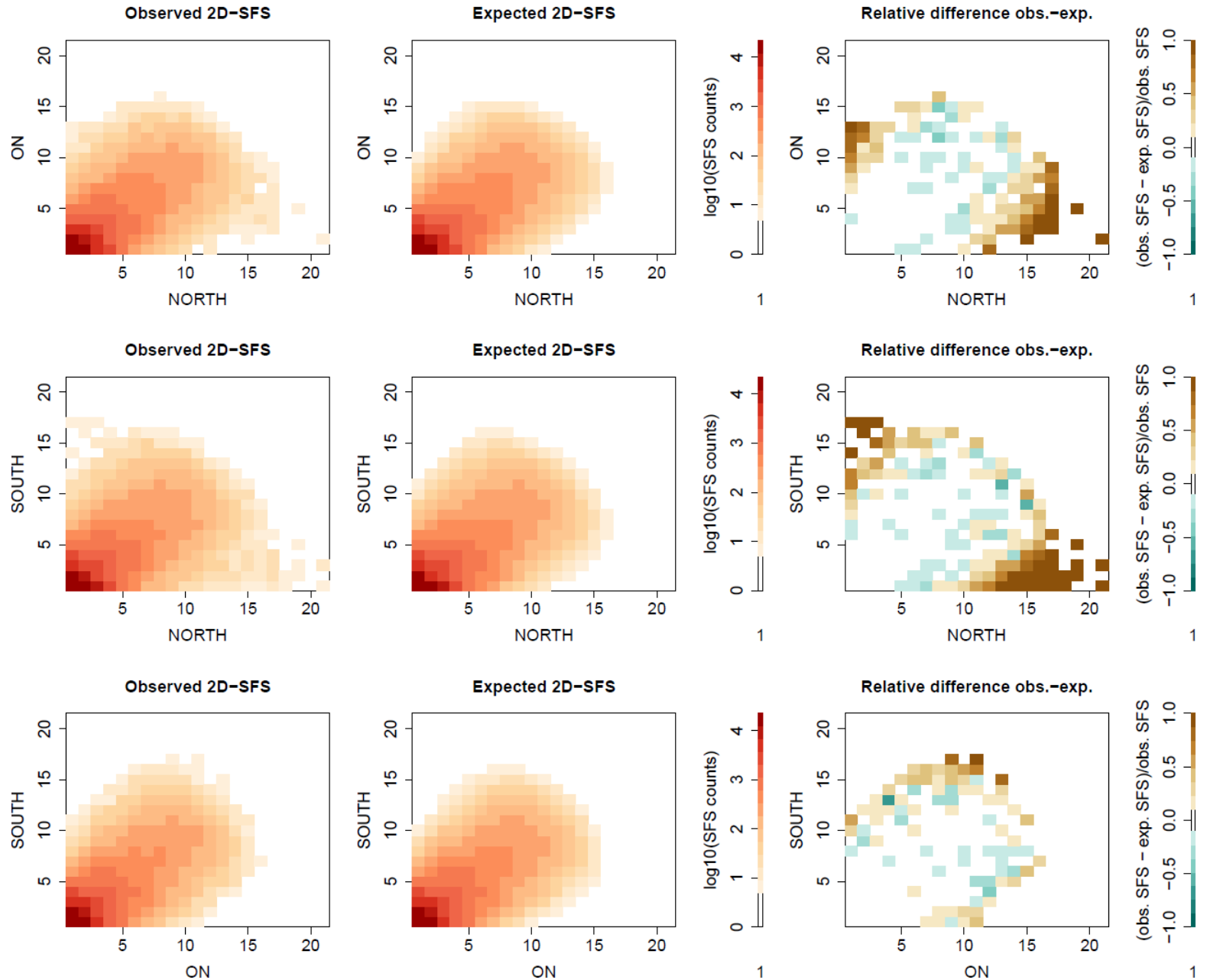
Estimates support south colonization and high gene flow levels

- Recent time of colonization of Sand Hills ~3-5 kya, younger than formation of Sand Hills 8-15 kya
- High migration rates across all populations, inferred for all models

Migration rates above/below arrows in units of $2Nm$, i.e. average number of immigrants per generation.



Deer mice: Model fit to marginal SFS

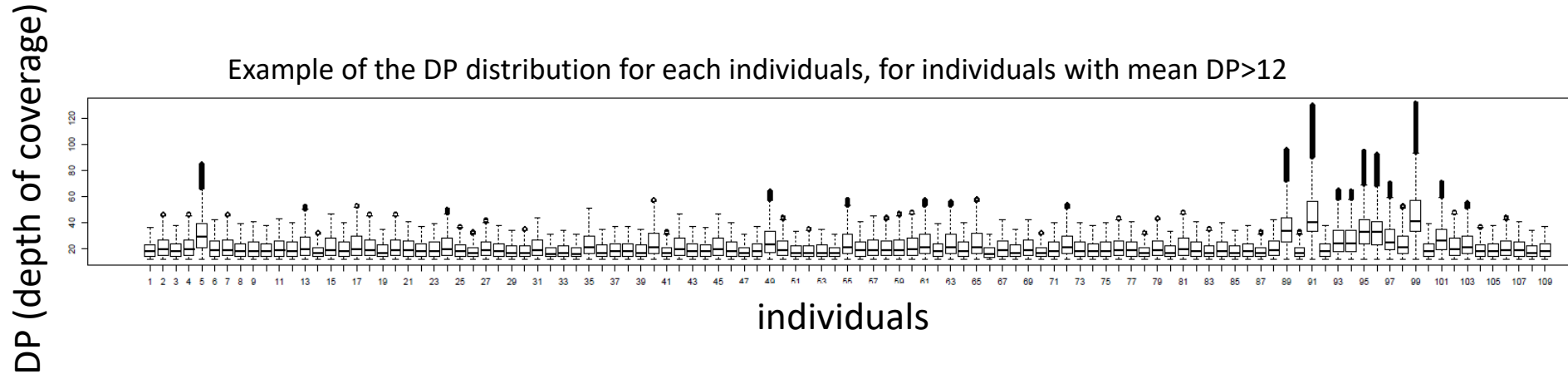


Some lessons I learned working with the deer mice data

- Be carefull when applying Hardy-Weinberg filters to your data
- Be carefull when filtering on depth of coverage applying the same thresholds for all individuals

The depth of coverage varied considerably across individuals

Example of the DP distribution for each individuals, for individuals with mean DP>12



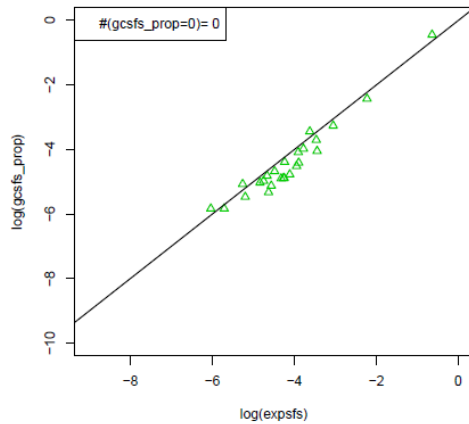
- Applying the same threshold for all individuals can lead to biases
- Apply a filter on DP for each individual

Effect of DP filters on the SFS

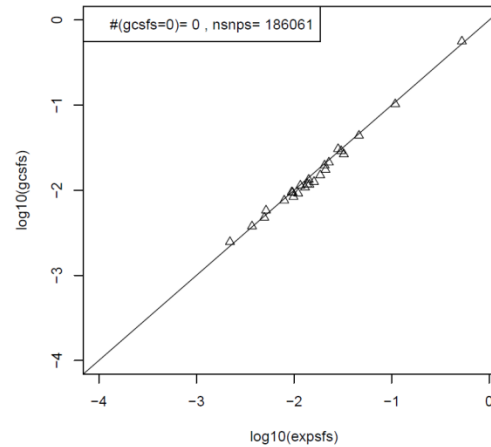
Simulation study

SFS based on
called
genotypes

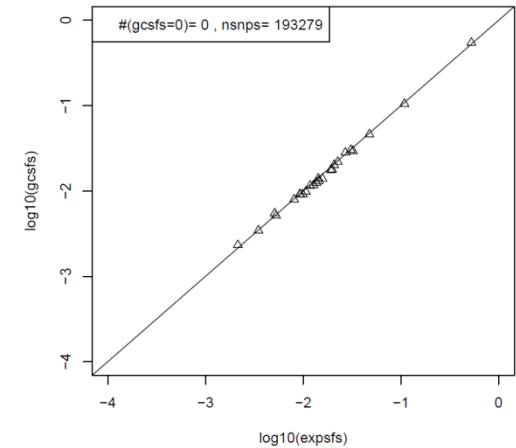
DP > 10



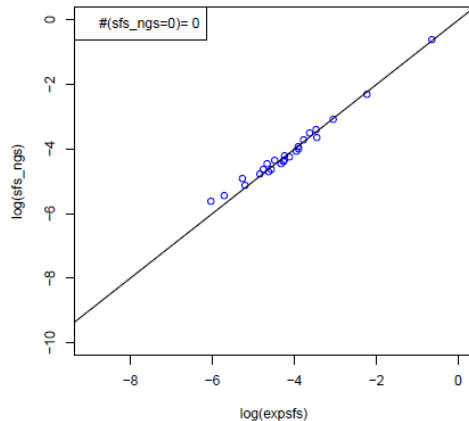
DP > 15



DP > 20



SFS accounting
for genotype
uncertainty
(ANGSD)

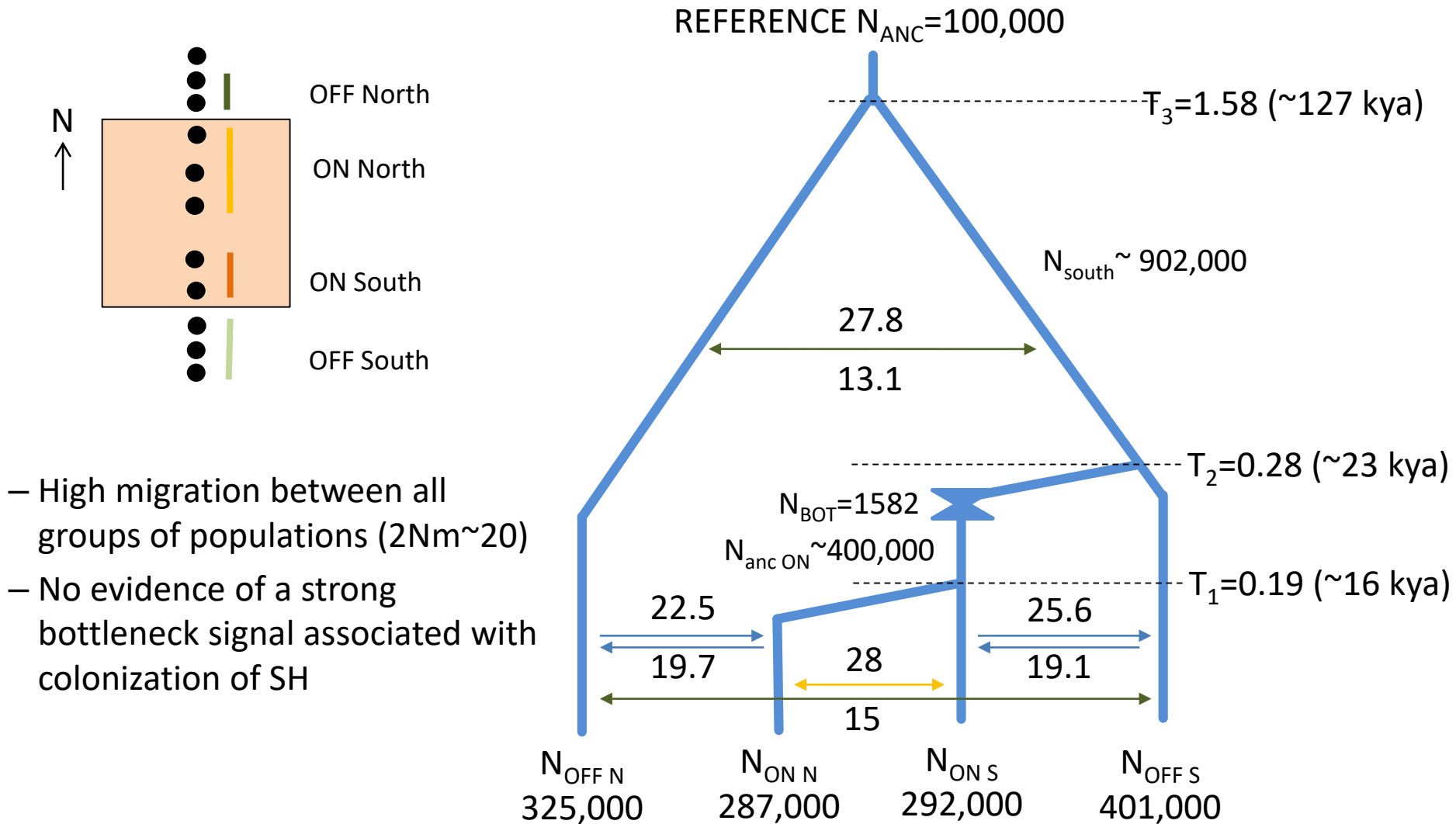


Simulated 2 pops SFS sampling 4 diploids from each pop, 200000 SNPs, mean coverage=**10x**, error rate=0.01. Simulated with correlated allele frequencies model ($F_{ST}=(0.275, 0.01)$)

With DP>15 we have a very good approximation to the correct SFS, even when using the called genotypes

Effect of HW filtering on demographic estimates

Removing sites with HWE excess and deficit leads to different estimates



Sawflies and RAD data

MOLECULAR ECOLOGY

Molecular Ecology (2016)

doi: 10.1111/mec.13972

History, geography and host use shape genomewide patterns of genetic variation in the redheaded pine sawfly (*Neodiprion lecontei*)

ROBIN K. BAGLEY,* VITOR C. SOUSA,† MATTHEW L. NIEMILLER‡ and CATHERINE R. LINNEN*

*Department of Biology, University of Kentucky, Lexington, KY 40506, USA, †cE3c - Centre for Ecology, Evolution and Environmental Changes, Faculdade de Ciências, Universidade de Lisboa, 1749-016 Lisboa, Portugal, ‡Illinois Natural History Survey, Prairie Research Institute, University of Illinois Urbana-Champaign, Champaign, IL 61820, USA

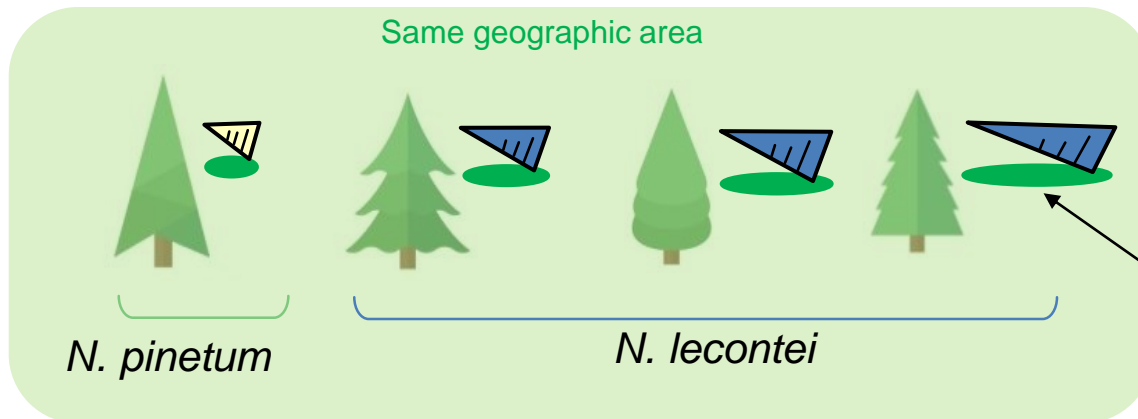
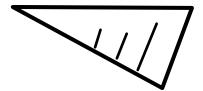


Sawflies *Neodiprion lecontei*

- Hymenoptera
- Plant-feeding insects
- Pine tree specialists



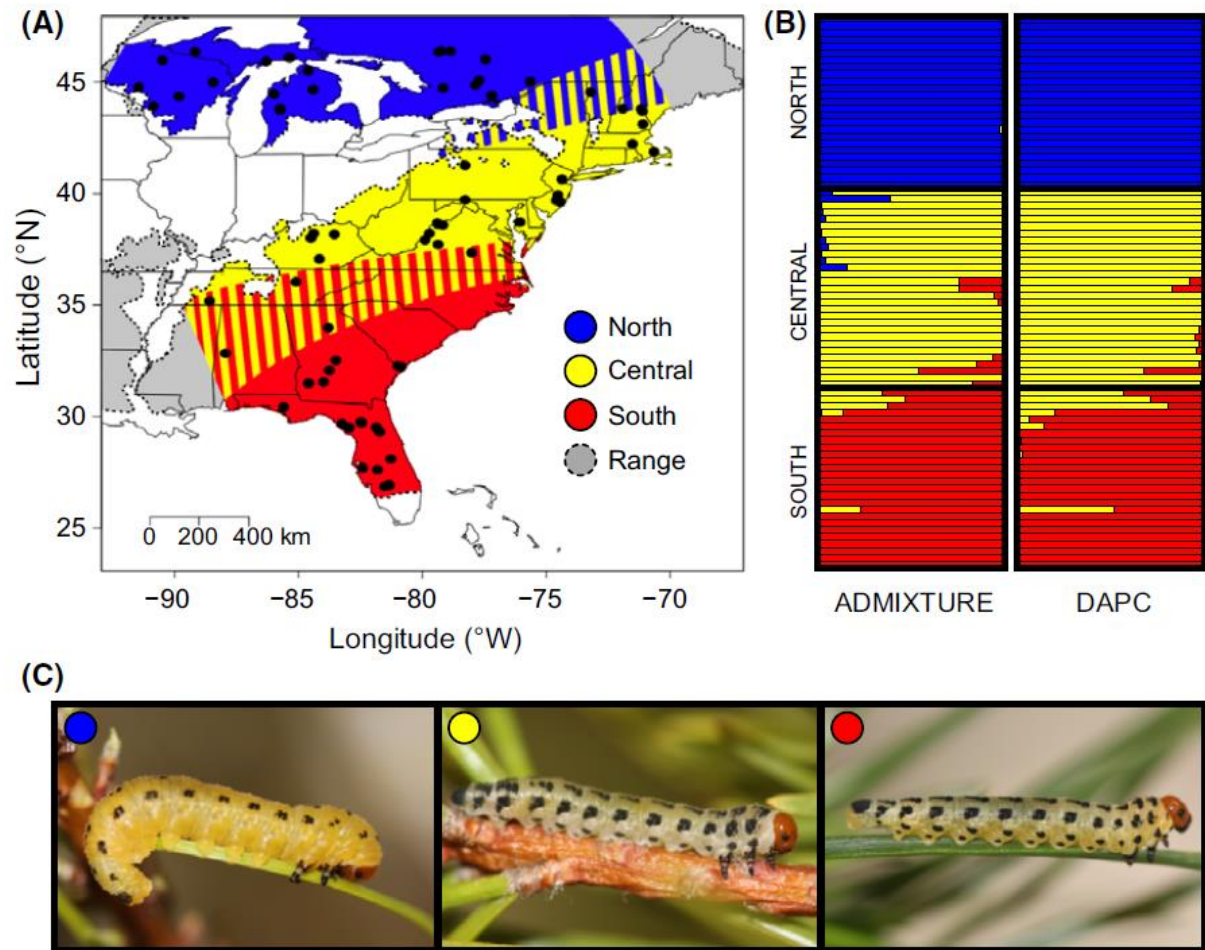
Ovipositor
(saw)



needle
width

ddRAD seq data

- 80 individuals from 77 localities and 13 host species
- 100 bp paired-end reads, mapped to reference genome of *N. lencontei*
- Depth of coverage filter DP>10

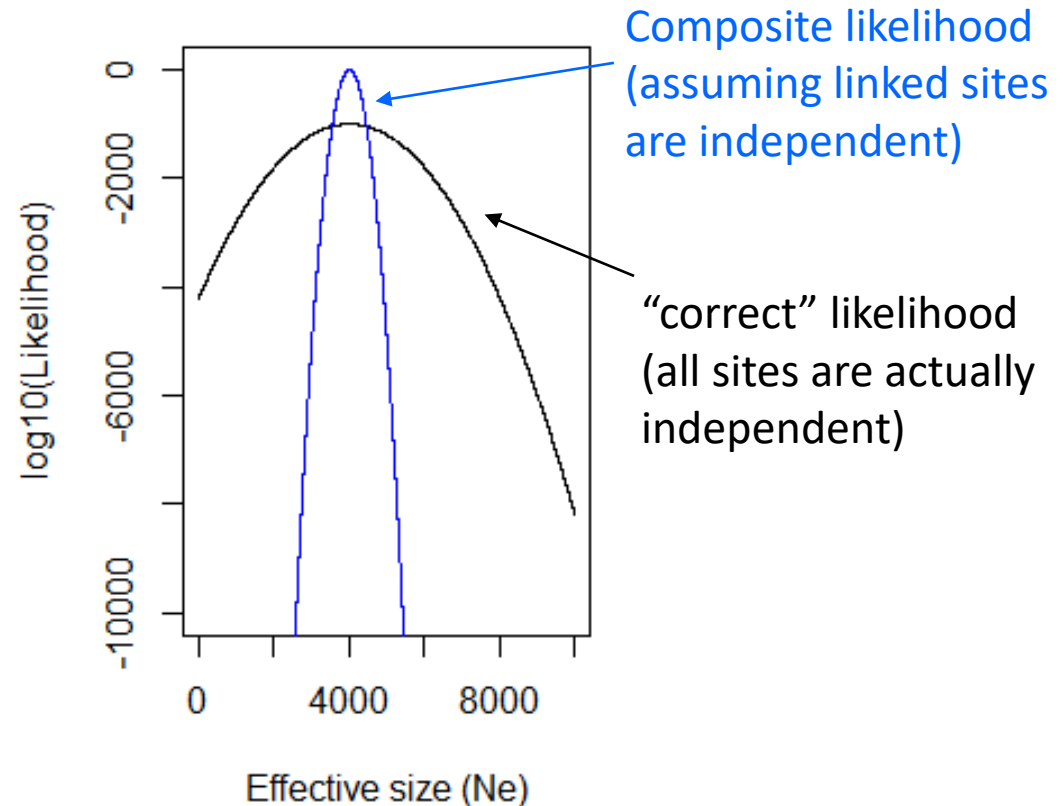


Given the detected three groups (North, Central, South):

- What is the the population tree topology?
- What are the split times?
- What are the migration levels among groups?

Comparing models with composite likelihoods

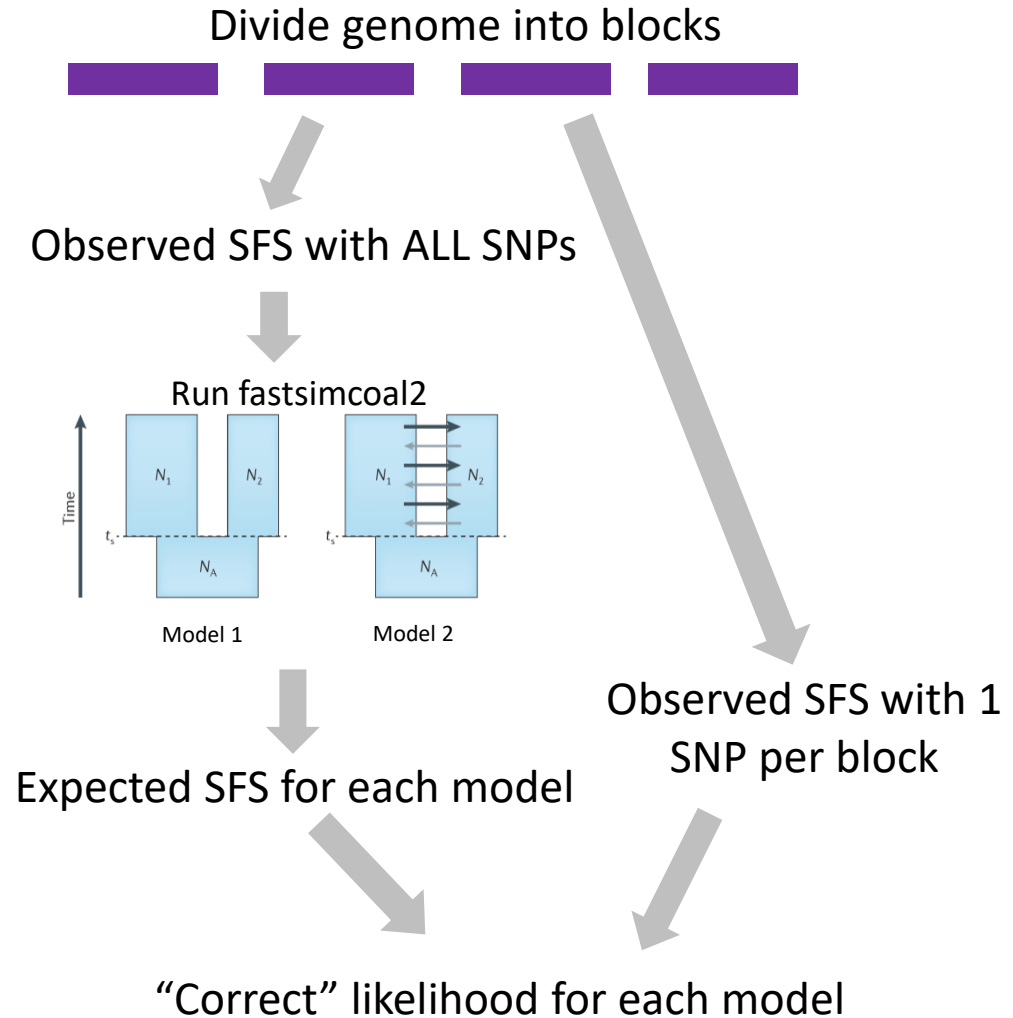
- Fastsimcoal2 likelihood is “correct” if all SNPs are independent
- We can then compare the model likelihoods using Akaike Information Criterion (AIC)



Composite likelihood provide unbiased maximum likelihood parameter estimates, but the likelihoods are inflated

A strategy to compare models

1. Divide the dataset into LD blocks.
2. Create a dataset with all SNPs (including linked SNPs)
3. For each model, obtain the parameters that maximize the likelihood (this is ok even with linked sites!) and the corresponding expected SFS
4. Create a dataset with “independent” SNPs (1 SNP per RAD tag)
5. Given the expected SFS of each model, compute the “correct” likelihood for each model with the dataset with independent SNPs
6. Compare models with AIC

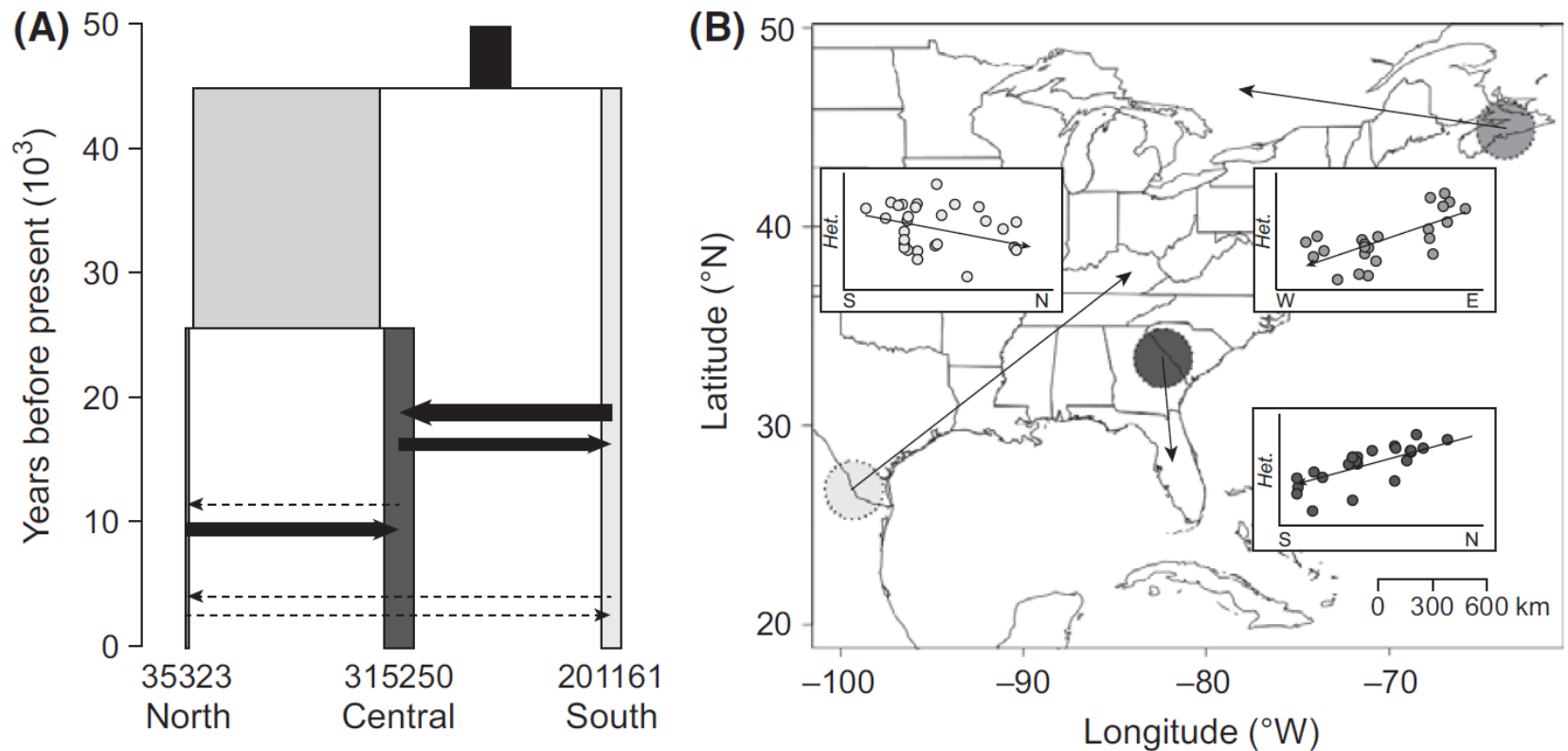


Comparing alternative models

Table 2 Summary of the likelihoods for the sixteen demographic models tested. Lhood (ALL SNPs) and Lhood (1 SNP) correspond to the mean likelihood computed with the data sets containing ‘all SNPs’ (including monomorphic sites) and a ‘single SNP’ (without monomorphic sites) per RAD locus, respectively. Mean likelihoods were computed based on 100 expected site frequency spectra simulated according to the parameters that maximized the likelihood of each model. Topology names for each model are as indicated in Fig. S1 (Supporting information). AIC scores and relative likelihoods (Akaike’s weight of evidence) were calculated based on the ‘single SNP’ data set following Excoffier *et al.* 2013.

Topology	Migration allowed?	Exponential growth?	North bottleneck?	log ₁₀ (Lhood) ALL SNPs	log ₁₀ (Lhood) 1 SNP	# Parameters	AIC	ΔAIC	Relative likelihood
North–South	No	No	No	−46502.02	−7381.4	7	34006.70	75.69	0.000
North–Central	No	No	No	−46475.82	−7369.0	7	33949.44	18.43	0.000
South–Central	No	No	No	−46502.18	−7381.6	7	34007.60	76.59	0.000
Trifurcation	No	No	No	−46501.54	−7380.4	5	33998.07	67.06	0.000
North–South	Yes	No	No	−46470.49	−7365.0	15	33947.25	16.24	~0.000
North–Central	Yes	No	No	−46462.24	−7361.5	15	33931.01	0.00	0.851
South–Central	Yes	No	No	−46467.69	−7363.8	15	33941.57	10.56	0.004
Trifurcation	Yes	No	No	−46470.28	−7364.7	11	33937.93	6.91	0.027
North–South	Yes	Yes	No	−46469.48	−7362.8	18	33942.91	11.90	0.002
North–Central	Yes	Yes	No	−46461.17	−7361.7	18	33937.82	6.80	0.028
South–Central	Yes	Yes	No	−46463.73	−7363.9	18	33948.15	17.13	~0.000
Trifurcation	Yes	Yes	No	−46467.72	−7363.3	14	33937.39	6.37	0.035
North–South	Yes	Yes	Yes	−46467.45	−7361.5	20	33940.86	9.85	0.006
North–Central	Yes	Yes	Yes	−46461.25	−7362.1	20	33943.82	12.81	0.001
South–Central	Yes	Yes	Yes	−46463.58	−7364.1	20	33953.08	22.07	0.000
Trifurcation	Yes	Yes	Yes	−46466.06	−7362.4	16	33936.93	5.92	0.044

Estimates favors a scenario where
North and Central diverged more recently with asymmetric gene flow



The inferred population tree topology and divergence times are consistent with divergence and range expansion from different refugia after LGM

Summary

- Fastsimcoal2 can be applied to RAD seq data
- We used a strategy to obtain (as close as possible) the “correct” likelihood by dividing the data into blocks, inferring the expected SFS for each model with ALL SNPs, and then re-computing the “true” likelihood with independent SNPs (1 SNP per block)
- Despite the reduced number of SNPs we were able to discriminate models based on their likelihoods

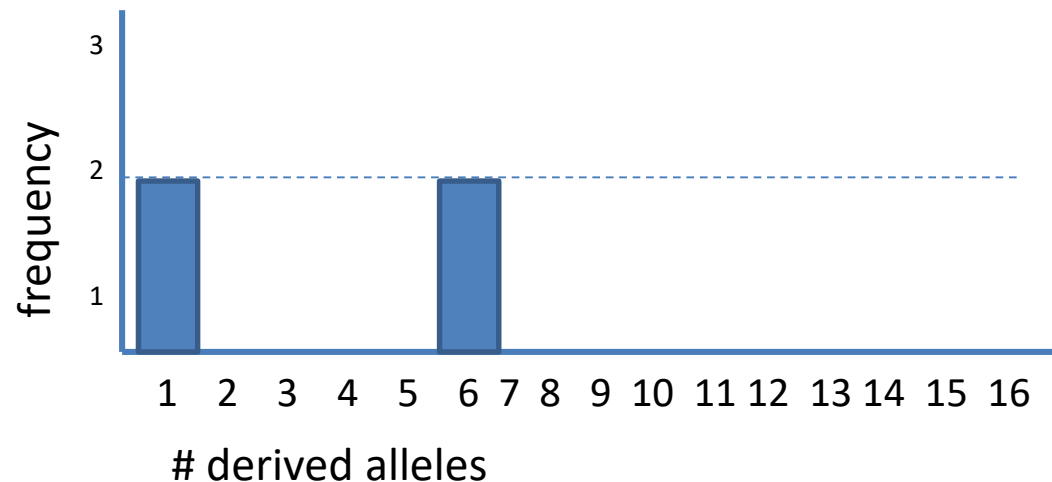
Protocol for model comparison based on AIC when we have independent SNPs

- Get the observed SFS
- Define the alternative models
- Perform 50-100 runs under each model
- Select the runs with maximum likelihood under each model
- Compute the AIC (Akaike information criteria) for each model
- Select the model with minimum AIC

Estimating SFS from observed data

- The sample size can vary across SNPs due to missing data
- How to deal with missing data?

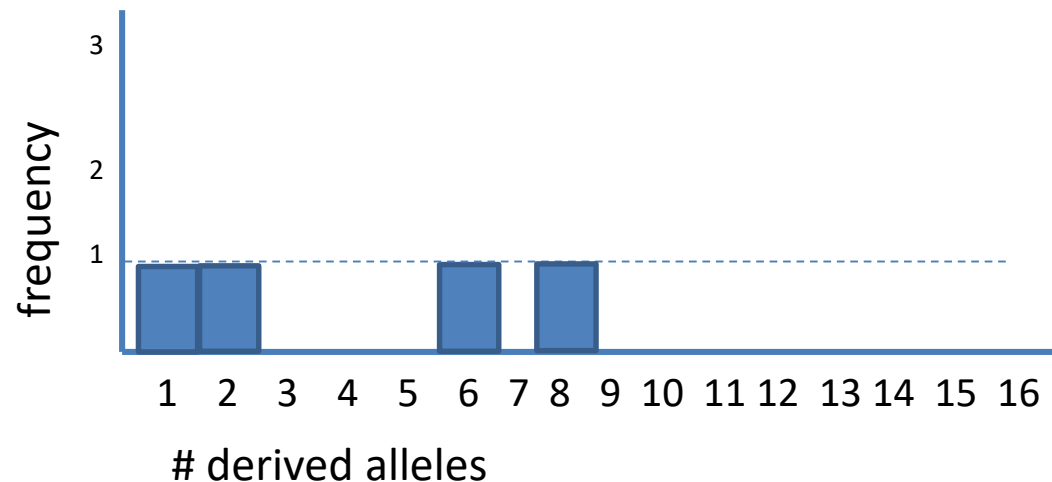
	Freq. derived	Sample size	Rel. freq
SNP1	1	16	1/16
SNP2	6	12	1/2
SNP3	1	12	1/12
SNP4	6	16	3/8



Estimating SFS from observed data

- The sample size can vary across SNPs due to missing data
- How to deal with missing data?

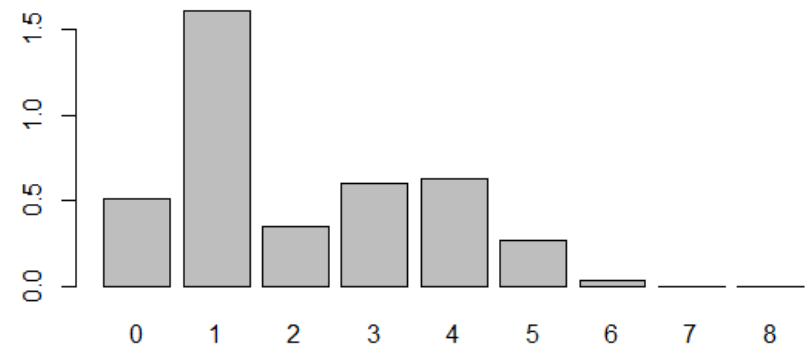
	Freq. derived	Sample size	Rel. freq
SNP1	1	16	1/16
SNP2	6	12	1/2
SNP3	1	8	1/12
SNP4	6	16	3/8



Estimating SFS from observed data

- The sample size can vary across SNPs due to missing data
- How to deal with missing data?
- Solution:
 - Find minimum sample size
 - Resample without replacement

	Freq. derived	Sample size	Rel. freq
SNP1	1	16	1/16
SNP2	6	12	1/2
SNP3	1	8	1/12
SNP4	6	16	3/8



FASTSIMCOAL2 INPUT FILES

Vitor Sousa

vmsousa@fc.ul.pt

Cesky Krumlov 2020

Examples of observed SFS

1PopExpInst20Mb_DAFpop0.obs

```
1 observations
d0_0      d0_1      d0_2      d0_3      d0_4      d0_5      d0_6      d0_7      d0_8      d0_9      d0_10
19973842  24630      810       173       145       111       88        84        61        56        0
```

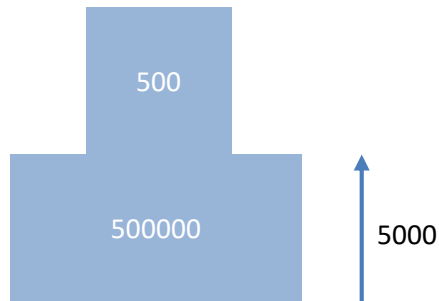
2PopDivMigr20Mb_jointDAFpop1_0.obs

```
1 observations
          d0_0      d0_1      d0_2      d0_3      d0_4      d0_5
d1_0      19985747  8350     1628     360      62       8
d1_1      9660      0        0        0        0
d1_2      4790      0        0        0        0
d1_3      3280      0        0        0        0
d1_4      2490      0        0        0        0
d1_5      1760      13       18       13       19       0
```

2PopDiv20Mb_jointDAFpop1_0.obs

```
1 observations
          d0_0      d0_1      d0_2      d0_3      d0_4      d0_5
d1_0      19985547  8211     1415     316      55      10
d1_1      1266      101      37       16       5       1
d1_2      61142     20       8        2        0
d1_3      48631     12       5        0        0
d1_4      47915     9        2        3        1
d1_5      1189      46       22       19       18       0
```

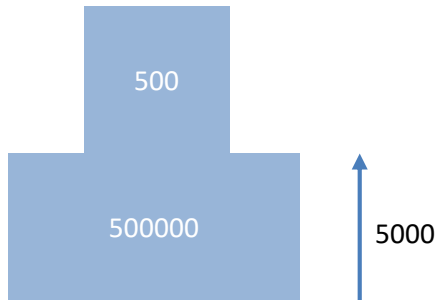
Parameter estimation settings files



1PopExpInst20Mb

- Additional files necessary to estimate parameters:
- Template file (TPL) defining the model
 - Estimation file (EST) with search range for parameters

Parameter estimation settings files



1PopExpInst20Mb

- Additional files necessary to estimate parameters:
- **Template file (TPL) defining the model**
 - Estimation file (EST) with search range for parameters

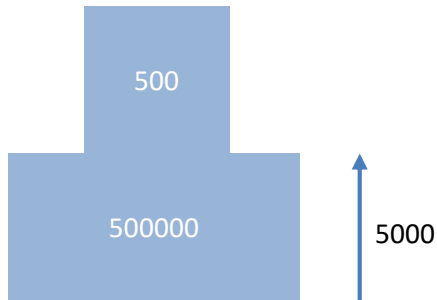
Tags for parameter we want to estimate:
\$NPOP\$, \$TEXP\$, \$RESIZE\$

Template file (filename.tpl)

1PopExpInst20Mb/1PopExpInst20Mb.tpl

```
//Parameters for the coalescence simulation program : fsimcoal2.exe
1 samples to simulate :
//Population effective sizes (number of genes)
$NPOP$
//Samples sizes and samples age
10
//Growth rates: negative growth implies population expansion
0
//Number of migration matrices : 0 implies no migration between demes
0
//historical event: time, source, sink, migrants, new deme size, new growth rate,
migration matrix index
1 historical event
$TEXP$ 0 0 0 $RESIZE$ 0 0
//Number of independent loci [chromosome]
1 0
//Per chromosome: Number of contiguous linkage Block: a block is a set of contiguous loci
1
//per Block:data type, number of loci, per generation recombination and mutation rates
and optional parameters
FREQ 1 0 2.5e-8 OUTEXP
```

Parameter estimation settings files



- Additional files necessary to estimate parameters:
- Template file (TPL) defining the model
 - **Estimation file (EST) with search range for parameters**

1PopExpInst20Mb

Tags for parameter we want to estimate:
\$NPOP\$, \$TEXP\$, \$RESIZE\$

Estimation file (filename.est)

1PopExpInst20Mb/1PopExpInst20Mb.est

```
// Search ranges and rules file
// *****

[PARAMETERS]
// #isInt? #name      #search #min  #max
// all Ns are in number of haploid individuals
1  $NPOP$           logunif  1000   1e7   output
1  $NANC$           logunif   10    1e5   output
1  $TEXP$           unif     10    1e5   output

[RULES]

[COMPLEX PARAMETERS]

0  $RESIZE$         = NANC/NPOP      hide
```

INPUT files for fastsimcoal2:

Defining an evolutionary model with TPL file

Number of samples
to simulate

2PopDivMig.tpl

```
//Parameters for the coalescence simulation program : fsimcoal2.exe  
2 samples to simulate :  
//Population effective sizes (number of genes)
```

\$NPOP1\$

\$NPOP2\$

```
//Samples sizes and samples age
```

5

5

```
//Growth rates: negative growth implies population expansion
```

0

0

```
//Number of migration matrices : 0 implies no migration between demes
```

2

```
//Migration matrix 0
```

0 0

\$MIG10\$ 0

```
//Migration matrix 1: No migration
```

0 0

0 0

```
//historical event: time, source, sink, migrants, new deme size, new growth rate, migration matrix  
index
```

```
2 historical event
```

\$TMIG_STOP\$ 0 0 0 1 0 1

\$TDIV_POP01\$ 1 0 1 **\$RESIZES\$** 0 1

```
//Number of independent loci [chromosome]
```

1 0

```
//Per chromosome: Number of contiguous linkage Block: a block is a set of contiguous loci
```

1

```
//per Block: data type, number of loci, per generation recomb. and mut. rates and optional parameters
```

FREQ 1 0 2.5e-8 OUTEXP

Deme sizes (2N)

Sample sizes

Growth rates

Migration
matrices

Historical events

Keep these as default
(not used for SFS)

Definition of genetic data type to simulate.

For SFS inference use:

FREQ 1 0 fixedMutationRate OUTFREQ

NOTE: for the SFS you cannot jointly infer the effective sizes and mutation rates! You need to **give a fixed mutation rate** if you have the number of monomorphic sites.

Otherwise, with "-0" option the mutation rate is ignored.

FREQ indicates you will use the SFS

OUTFREQ means the expected SFS will be output

TPL files

These files are very important! Check carefully all the definitions. Errors in the TPL file are difficult to detect and imply the model specification is incorrect! This means that all inferences will be wrong, and also that all parameter estimates will be incorrect!

Defining population sizes and sample sizes

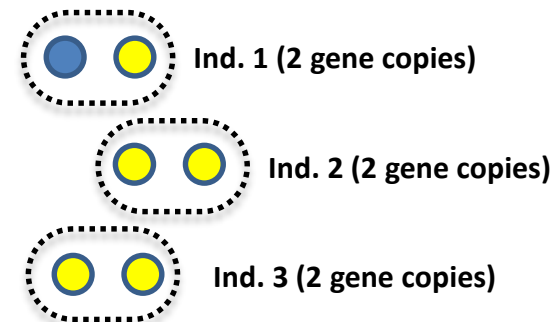
2PopDivMigr10Loci.par

```
//Parameters for the coalescence simulation program : fsimcoal2.exe
2 samples to simulate :
//Population effective sizes (number of genes)
$NPOP1$
$NPOP2$
//Samples sizes and samples age
6
6
//Growth rates: negative growth implies population expansion
0
0
```

Parameter tags

Population effective sizes are given in number of gene copies. For a diploid species with $N=500$ individuals, this corresponds to a $2N=1000$ gene copies, as each individual carries two gene copies at any given site.

The sample size is also given in gene copies. The value of 6 means that we sampled 3 diploid individuals.



TPL files

MIGRATION

```
//Number of migration matrices : 0 implies no migration between demes
1
//migration matrix
0.000 $MIG_01$
$MIG_10$ 0.000
```

Parameter tags

The migration matrix can be asymmetric, and in the case the entry m_{ij} list the **migration rates backward in time** from population in row i to population in column j . The above-mentioned matrix states that, for each generation (backward in time), any gene from population 0 has probability MIG_01 to be sent to population 1, and that a gene from population 1 has a probability MIG_10 to move to population 0.

If no migration matrix is defined, no migration is assumed between populations.

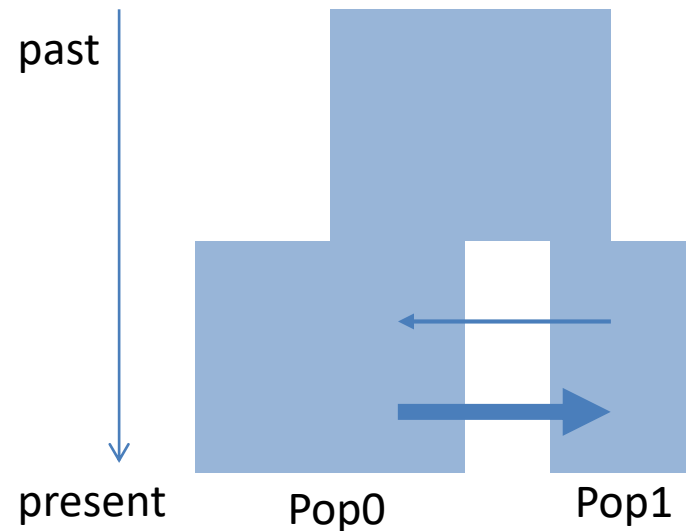
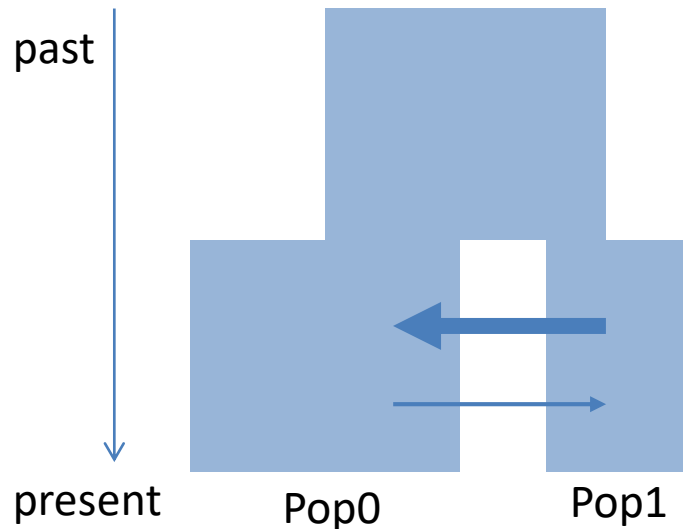
1PopStationary10Loci.par

```
//Number of migration matrices : 0 implies no migration between demes
0
```

A note on looking backward in time

Assuming that we look **forward in time** and that the size of the arrows are proportional to the migration rate, to what model does the following migration matrix corresponds to?

```
//Number of migration matrices : 0 implies no migration between demes
1
//migration matrix
0.000 0.005
0.001 0.000
```



A note on looking backward in time

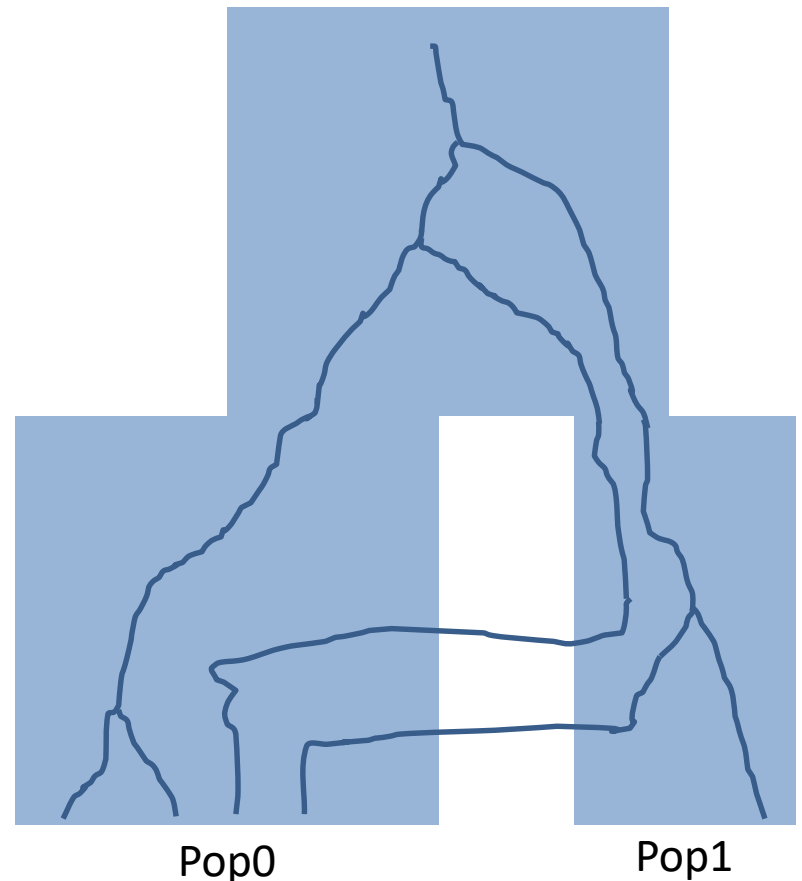
Assuming that we look **forward in time** and that the size of the arrows are proportional to the migration rate, to what model does the following migration matrix corresponds to?

```
//Number of migration matrices  
1  
//migration matrix  
0.000 0.005  
0.001 0.000
```

This means that there are more lineages migrating ("jumping") from pop0 to pop1 backward in time.

Thus, in Pop0 there are many individuals whose ancestors were migrants from Pop1 into Pop0.

past
↓
present

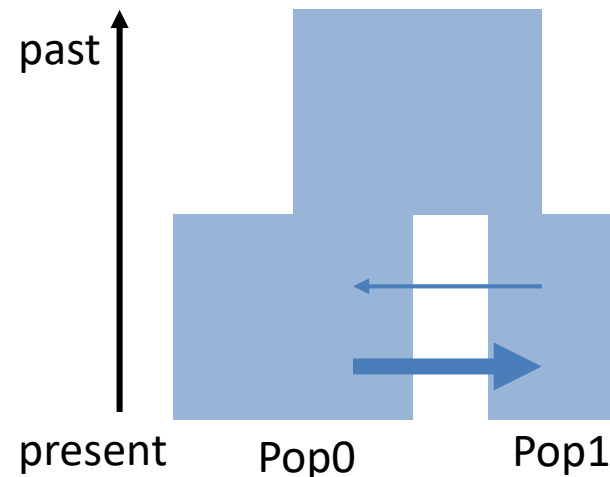
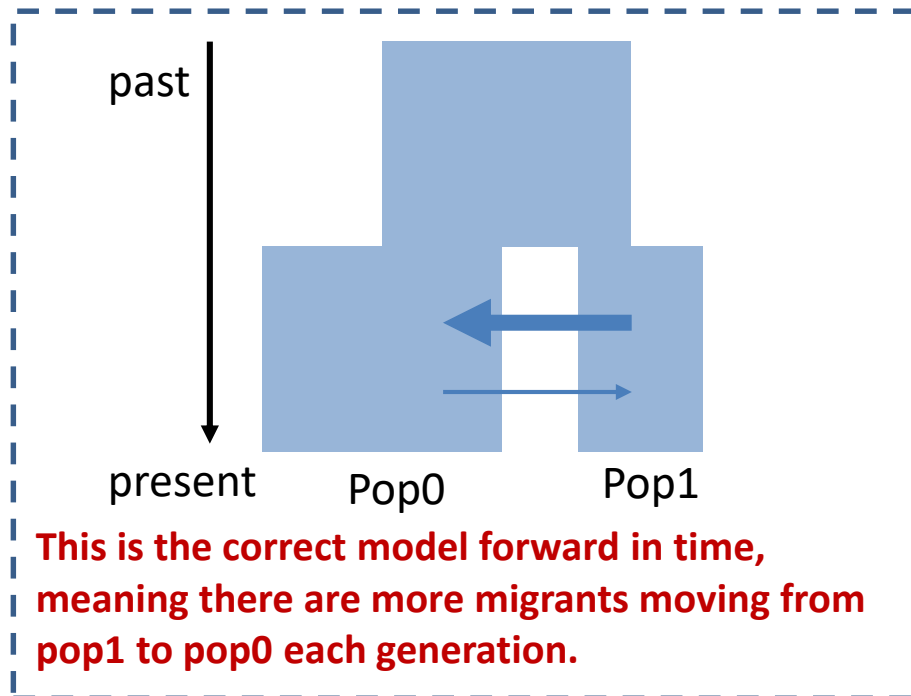


A note on looking backward in time

Assuming that we look forward in time and that the size of the arrows are proportional to the migration rate, to what model does the following migration matrix corresponds to?

```
//Number of migration matrices : 0 implies no migration between demes
1
//migration matrix
0.000 0.005
0.001 0.000
```

Note that in the PAR and TPL files everything is backward in time!!



Historical events in fastsimcoal2

Historical events can be used to:

- Change the size of a given population
- Change the growth rate of a given population
- Change the migration matrix to be used between populations
- Move a fraction of the genes of a given population to another population. This amounts to implementing a (stochastic) admixture or introgression event.
- Move all genes from a population to another population. This amounts to fusing two populations into one looking backward in time.
- One or more of these events at the same time

Defining the historical events is crucial to have a correct model!

Historical events (backward in time)

Each historical event is coded with a line with the following arguments

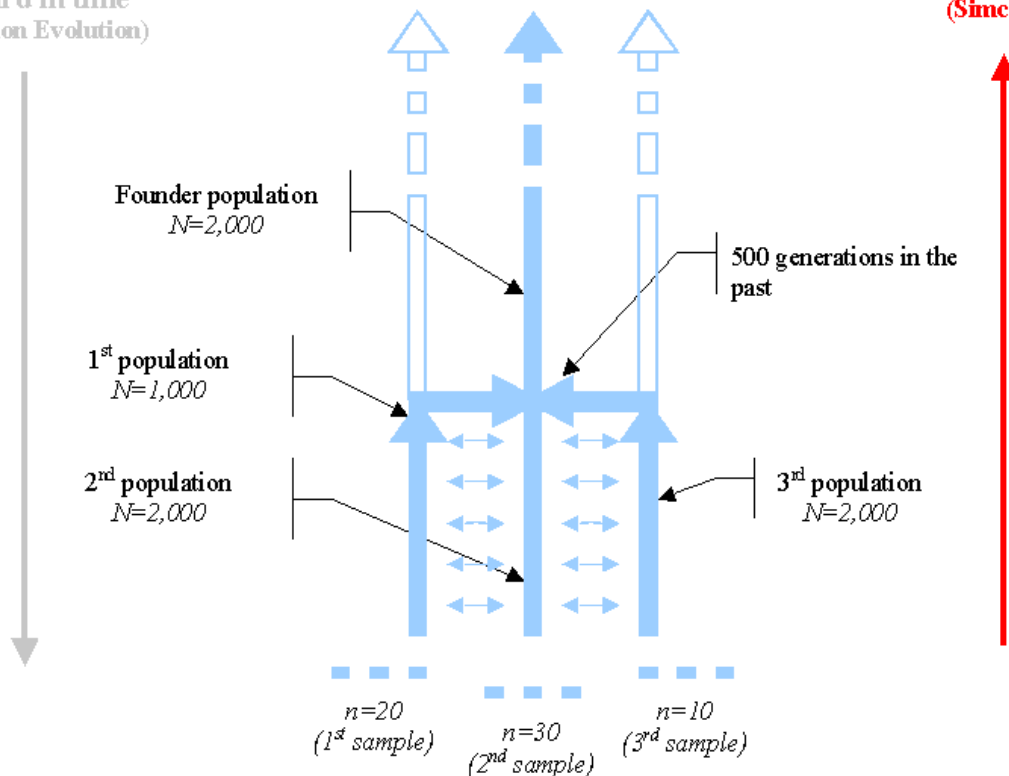
time, **source**, **sink**, **migrants**, **new deme size**, **new growth rate**, migration matrix index

500 0 1 1 1 0 1
500 2 1 1 1 0 1

500 generations ago, 100% (**migrants=1.0**) of lineages in **pop0** (**source =0**) migrated to **pop1** (**sink=1**). The size of the sink (pop1) remained the same (**new deme size=1.0**, i.e. $N_2=2000$). The new growth rate is zero. The migration rate that is active after the event is given in the migration matrix 1.

Forward in time
(Population Evolution)

Backward in time
(Simcoal2)



Historical events (backward in time)

Each historical event is coded with a line with the following arguments

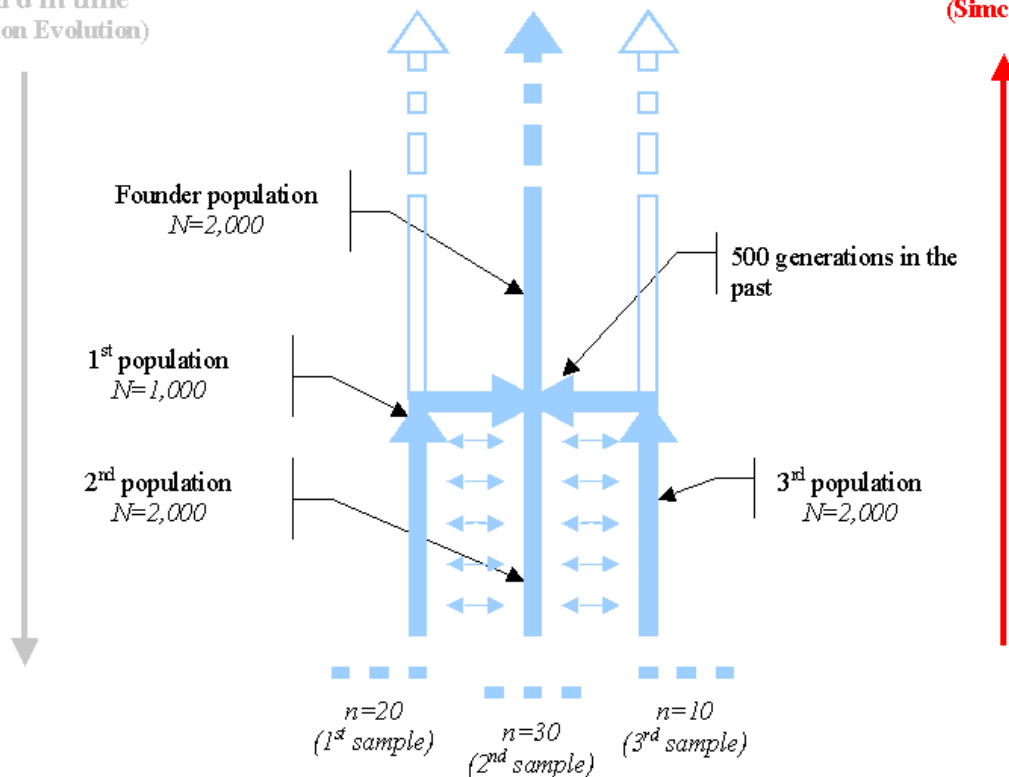
time, **source**, **sink**, **migrants**, **new deme size**, **new growth rate**, migration matrix index

500 0 1 1 1 0 1

500 2 1 1 1 0 1

Forward in time
(Population Evolution)

Backward in time
(Simcoal2)



500 generations ago, 100% of lineages (**migrants=1.0**) in **pop2** (**source =2**) migrated to **pop1** (**sink=1**). The size of the sink (pop1) remained the same (**new deme size=1.0**, i.e. $N_2=2000$). The new growth rate is zero. The migration rate that is active after the event is given in the migration matrix 1.

Historical events in fastsimcoal2

Change the size of a given population

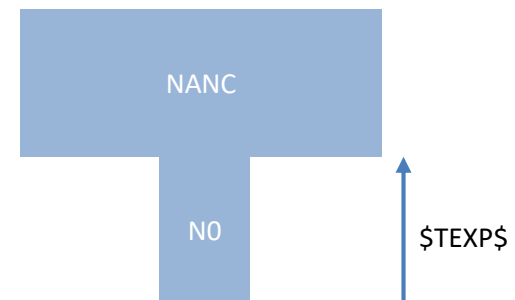
1PopContrInst10Loci.par

```
//Parameters for the coalescence simulation program : fsimcoal2.exe
1 samples to simulate :
//Population effective sizes (number of genes)
1000
//Samples sizes and samples age
10
//Growth rates: negative growth implies population expansion
0
//Number of migration matrices : 0 implies no migration between demes
0
//historical event: time, source, sink, migrants, new deme size, new growth rate, migration matrix index
1 historical event
$TEXP$ 0 0 0 $RESIZE$ 0 0
```



- **\$TEXP\$** generations ago, 0% (migrants=0) of lineages in pop0 (source) migrated to pop1 (sink). This means that 100% of lineages remained in pop0.
- The sink population (pop0) has a size **\$RESIZE\$** times larger after the event (**\$RESIZE\$=\$NANC\$/\$NO\$**). Given **NO** diploids at time zero, it implies that **NANC=NO*RESIZE** diploids.
- The migration matrix valid after the event is the migration rate 0. Since it is not defined it implies no migration.

Recent instantaneous
demographic contraction



1PopContrInst10loci.par

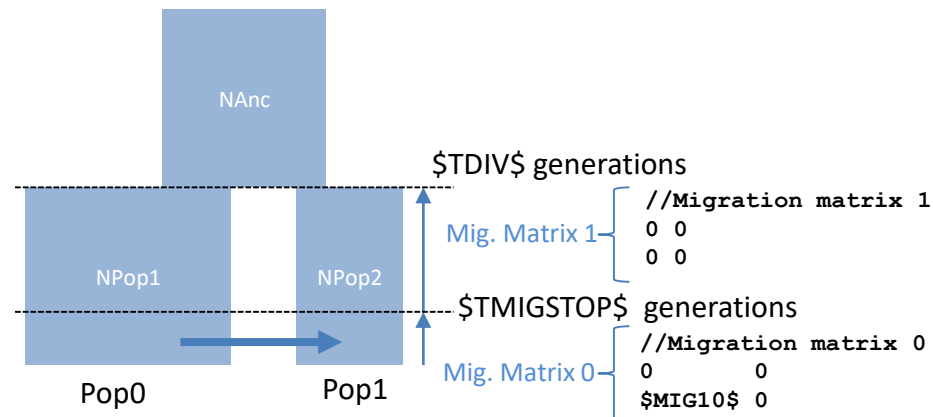
Historical events in fastsimcoal2

Change the migration matrix to be used between populations

2PopDivMigr10Loci.par

```
//Number of migration matrices : 0 implies no migration between demes
2
//Migration matrix 0
0 0
$MIG10$ 0
//Migration matrix 1: No migration
0 0
0 0
//historical event: time, source, sink, migrants, new deme size, new growth rate, migration matrix
index
2 historical event
$TMIGSTOP$ 0 0 0 1 0 1
$TDIV$ 1 0 1 $RESIZES$ 0 1
```

- At generation `$TMIGSTOP$` in the past, 0% (migrants=0) of lineages migrated from pop0 (source=0) to pop1 (sink=0).
- After the historical event, the deme size of the sink population (pop1) remained the same (new deme size=1).
- After the historical event the growth rate was set to zero.
- After the historical event the migration rate matrix was set to matrix 1, i.e. no migration between populations.



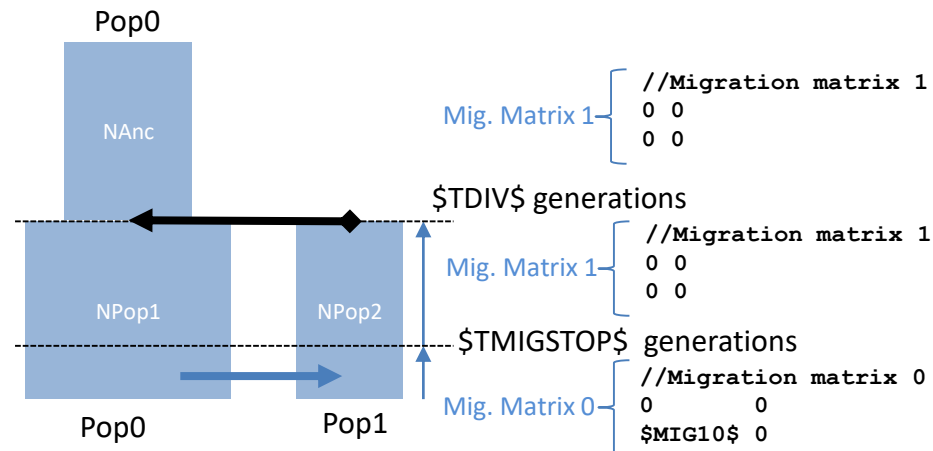
Historical events in fastsimcoal2

Population split (merge populations going backwards in time)

2PopDivMigr10Loc1.par

```
//Number of migration matrices : 0 implies no migration between demes
2
//Migration matrix 0
0 0
$MIG10$ 0
//Migration matrix 1: No migration
0 0
0 0
//historical event: time, source, sink, migrants, new deme size, new growth rate, migration matrix
index
2 historical event
$TMIGSTOP$ 0 0 0 1 0 1
$TDIV$ 1 0 1 $RESIZE$ 0 1
```

- At generation $\$TDIV\$$ in the past, 100% (migrants=1) of lineages migrated from pop1 (source=1) to pop0 (sink=0).
- After the population split, the deme size of the sink population (pop0) is $\$NANC\$$, and hence $\$RESIZE\$ = \$NANC\$ / \$NPOPO\$$.
- After the historical event the growth rate of the sink population pop0 is zero.
- After the historical event the migration rate matrix was set to matrix 1, i.e. no migration between populations.



Estimation file (.est)

Estimation file ("NoMigSan_Maya.est")

NoMigSan_Maya.est

```
// Search ranges and rules file
// *****

[PARAMETERS]
//#isInt? #name      #search.#min  #max
//all Ns are in number of haploid individuals
1  $NPOP1$          unif  10    1e5    output
1  $NPOP2$          unif  10    1e5    output
1  $NANC$           unif  10    1e5    output
1  $NBOTP1$         unif  1     1e3    output  bounded

1  $TDIV$           unif  100   1e4    output
0  $RELTBOT$        unif  1e-5  1      hide    bounded

[RULES]

[COMPLEX PARAMETERS]

0  $RES_BOT_START$ = $NBOTP1$/$NPOP1$      hide
0  $RES_BOT_END$   = $NPOP1$/$NBOTP1$      hide

1  $TBOT_START$    = $TDIV$ * $RELTBOT$     output
1  $TBOT_END$      = $TBOT_START$ + 10      hide

0  $RESIZE0$       = $NANC$/$NPOP1$         hide
```

Each line must contain the following:

#isInt? 0 for continuous, 1 for integers

#name Parameter tag name

#search "unif" for uniform scale

"logunif" for log10 scale

#min minimum search range (lower bound)

#max maximum search range. If the keyword bounded is not used, then if likelihood is higher near maximum value, fastsimcoal2 will keep increasing the maximum value. The **bounded** keyword prevents this.

Complex parameters depend on the values of other parameters. Only one operation per line can be done. Thus, you cannot have something with many operations in a single line:

$\$BLA\$ = (\$BL\$ * \$A\$) + (\$BLA\$ / \$LA\$) - \text{WRONG!}$

Estimation file (.est)

Estimation file ("NoMigSan_Maya.est")

NoMigSan_Maya.est

```
// Search ranges and rules file
// *****

[PARAMETERS]
// #isInt? #name      #search.#min  #max
// all Ns are in number of haploid individuals
1  $NPOP1$           unif  10    1e5    output
1  $NPOP2$           unif  10    1e5    output
1  $NANC$            unif  10    1e5    output
1  $NBOTP1$          unif  1     1e3    output  bounded

1  $TDIV$            unif  100   1e4    output
0  $RELTBOT$         unif  1e-5  1      hide    bounded

[RULES]

[COMPLEX PARAMETERS]

0  $RES_BOT_START$ = $NBOTP1$/$NPOP1$      hide
0  $RES_BOT_END$   = $NPOP1$/$NBOTP1$      hide

1  $TBOT_START$ = $TDIV$ * $RELTBOT$      output
1  $TBOT_END$   = $TBOT_START$ + 10      hide

0  $RESIZE0$      = $NANC$/$NPOP1$        hide
```

Note that complex parameters can be used to define the order of events.

By using a \$RELTBOT\$ between 1e-5 and 1, and then specifying that

$$\text{\$TBOT_START\$} = \text{\$TDIV\$} * \text{\$RELTBOT\$}$$

This means that the TBOT_START is always more recent than the time of divergence.

If this is not well specified you can get errors, because events need to happen in a specific order. Another solution is to actually estimate the time between time events.

Estimation file (.est)

Estimation file ("NoMigSan_Maya.est")

NoMigSan_Maya.est

```
// Search ranges and rules file
// *****

[PARAMETERS]
//#isInt? #name      #search.#min  #max
//all Ns are in number of haploid individuals
1  $NPOP1$          unif  10    1e5    output
1  $NPOP2$          unif  10    1e5    output
1  $NANC$           unif  10    1e5    output
1  $NBOTP1$         unif   1    1e3    output  bounded

1  $TBOT_END$       unif  100   1e4    output
0  $TDIV_TBOT_INT$  unif   10    1e3    hide

[RULES]

[COMPLEX PARAMETERS]

0  $RES_BOT_START$ = $NBOTP1$/$NPOP1$      hide
0  $RES_BOT_END$   = $NPOP1$/$NBOTP1$      hide

1  $TBOT_START$ = $TBOT_END$ - 10          output
1  $TDIV$ = $TBOT_END$ + $TDIV_TBOT_INT$  output

0  $RESIZE0$       = $NANC$/$NPOP1$        hide
```

Another solution is to actually estimate the time between time events, as shown on the left.

In this case, we would estimate the parameter \$TDIV_TBOT_INT\$

And then in complex parameters:

$TDIV = TBOT_END + TDIV_TBOT_INT$

Estimation file (.est)

Estimation file ("NoMigSan_Maya.est")

NoMigSan_Maya.est

```
// Search ranges and rules file
// *****

[PARAMETERS]
// #isInt? #name      #search.#min  #max
// all Ns are in number of haploid individuals
1  $NPOP1$           unif  10    1e5    output
1  $NPOP2$           unif  10    1e5    output
1  $NANC$            unif  10    1e5    output
1  $NBOTP1$          unif  1     1e3    output  bounded

1  $TBOT_END$         unif  100   1e4    output
0  $TDIV_TBOT_INT$   unif  10    1e3    hide

[RULES]

[COMPLEX PARAMETERS]

0  $RES_BOT_START$ = $NBOTP1$/$NPOP1$      hide
0  $RES_BOT_END$   = $NPOP1$/$NBOTP1$      hide

1  $TBOT_START$ = $TBOT_END$ - 10          output
1  $TDIV$ = $TBOT_END$ + $TDIV_TBOT_INT$ output

0  $RESIZE0$      = $NANC$/$NPOP1$         hide
```

Finally, a note about inferring bottlenecks associated with founder events.

It is difficult to jointly infer the duration and Effective population size of a bottleneck.

Instead, we can infer the bottleneck intensity, which is given by

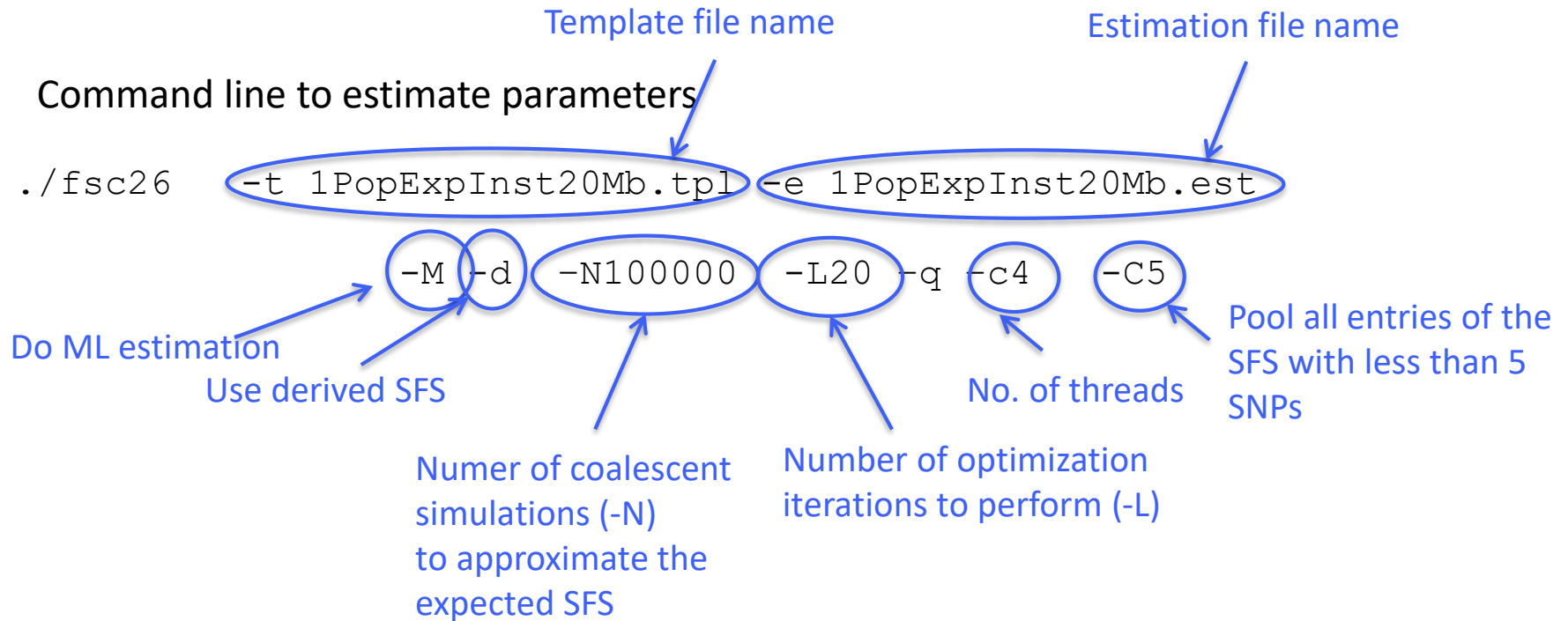
$$I_B = \frac{\text{Time Duration Bottleneck (generations)}}{(2 * \text{Effective size during bottleneck})}$$

Thus, we usually fix the duration of the bottleneck and infer the effective size.

In this case, we fix the duration of the bottleneck to be 10 generations.

If \$NBOTP1\$ is larger than 500, then actually there was no bottleneck, as $I_B < 0.01$ ($10/(2*500)$).

Launching parameter estimations



Observed SFS file must have the same name as template file and extension
_DAFpop0.obs. e.g. `1PopExpInst20Mb_DAFpop0.obs`