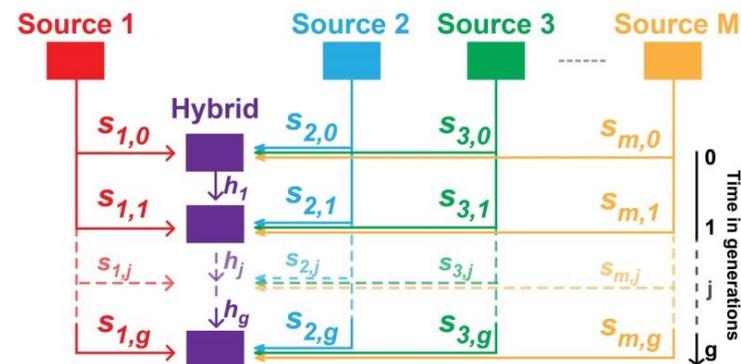# Reconstructing highly complex admixture histories using genetic data

Paul Verdu

Lab: UMR7206 Eco-anthropologie
Institution: CNRS – MNHN - Université Paris Cité
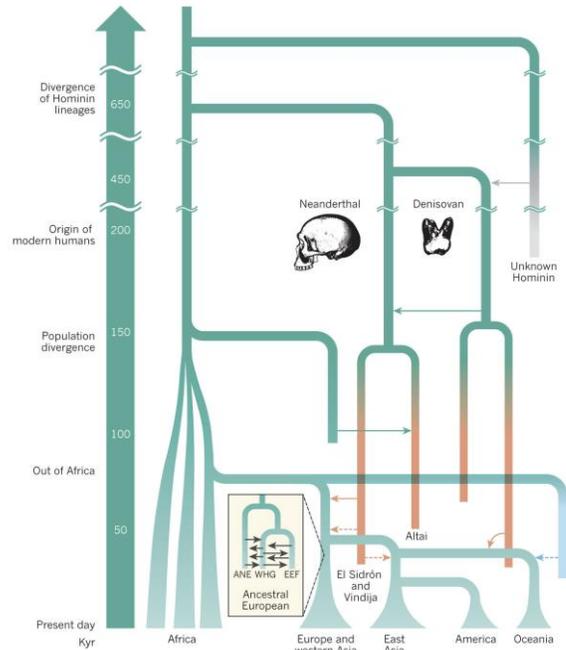
Rencontres de la Chaire MMB, Ecole Polytechnique
25 Avril 2022

# Admixture is ubiquitous in evolutionary history, through time and space

## e.g. *Homo sapiens*



Skoglund et al., *Science* 2012

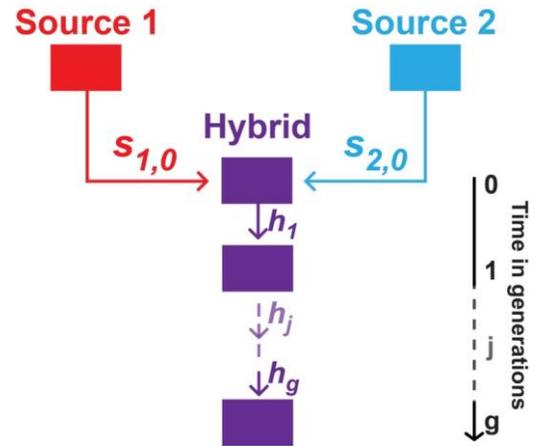

e.g. Nielsen et al., *Nature* 2018



Patin et al., *Science* 2017



Quach and Quintana-Murci, *JEM* 2017

# Inferring admixture history from genetic data

# Inferring admixture history from genetic data



Hellenthal et al., *Nature* 2014

Gravel, *Genetics* 2012

Pritchard et al., *Plos Gen* 2012

**Maximum likelihood approaches relying on**:

► Admixture Linkage-Disequilibrium distributions - e.g. *TRACTS* (Gravel 2012), *GLOBETROTTER* (Hellenthal et al. 2014)

► Moments of allelic frequency spectrum divergences - e.g. *M/ALDER* (Loh et al. 2013), *TreeMix* (Pritchard et al. 2012)

# Inferring admixture history from genetic data



Chromosome painting
Raw painting (chunks):
Cleaned painting:

*Hellenthal et al., Nature 2014*

*Gravel, Genetics 2012*

*Pritchard et al., Plos Gen 2012*

**Maximum likelihood approaches relying on**:

► Admixture Linkage-Disequilibrium distributions - e.g. *TRACTS* (Gravel 2012), *GLOBETROTTER* (Hellenthal et al. 2014)

► Moments of allelic frequency spectrum divergences - e.g. *M/ALDER* (Loh et al. 2013), *TreeMix* (Pritchard et al. 2012)

**Limited by**:

► **"Simple" admixture models**: one or two admixture pulses per source population

► **No formal model-choice**

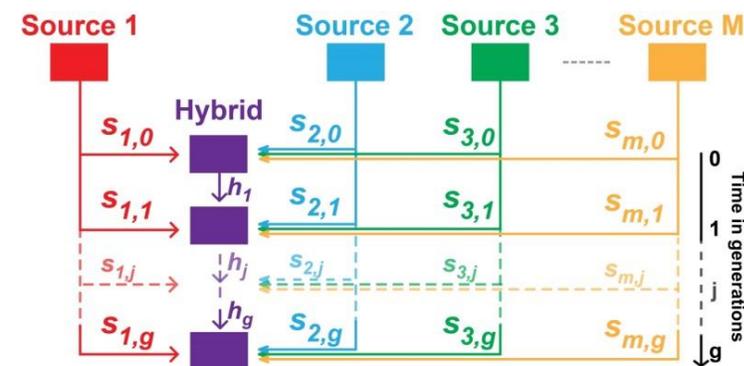► **Massive genomic data** and **accurate phasing** for admixture-LD methods

**ML-inference methods cannot operate**

$$L(ModelParams|Data_{obs}) \propto P(Data_{obs}|ModelParams) \, P(ModelParams)$$

**for highly complex admixture models where likelihoods cannot be written**

**for highly complex admixture models where likelihoods are intractable**



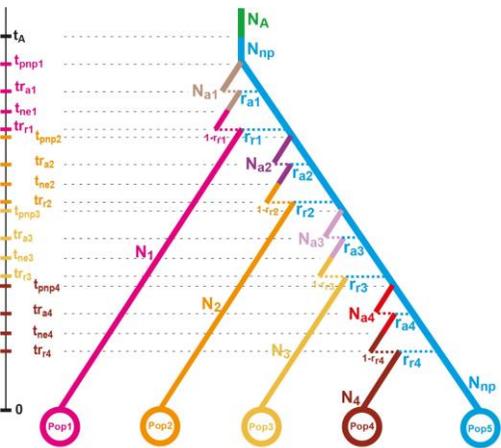**Approximate Bayesian Computation (Tavaré et al. 1997) may represent an alternative**

$$L(ModelParams|Data_{obs}) \propto P(Data_{obs}|ModelParams) \, P(ModelParams)$$

Approximation

$$L(ModelParams|SumStat_{obs}) \propto P(SumStat_{obs}|ModelParams) \, P(ModelParams)$$

# Approximate Bayesian Computation demographic inference



## Parameter prior distributions

| | |
|---|---|
| $N_A$ | Uniform [10 .. 1 000] |
| $N_{np}$ | Uniform [10 .. 100 000] |
| $N_1$ | Uniform [10 .. 10 000] |
| ... | |
| $r_{a1}$ | Normal [0 , 1] |
| $r_{r1}$ | LogUniform [0 , 1] |
| ... | |
| $t_{pnp1}$ | Uniform [1 .. 5 000] |
| $tr_{a1}$ | Uniform [1 .. 250] |
| ... | |

*Tavaré et al., Plos Gen 1997*
*Pritchard et al., Genetics 2000*
*Beaumont et al., Genetics 2002*

# Approximate Bayesian Computation demographic inference

*Tavaré et al., Plos Gen 1997*
*Pritchard et al., Genetics 2000*
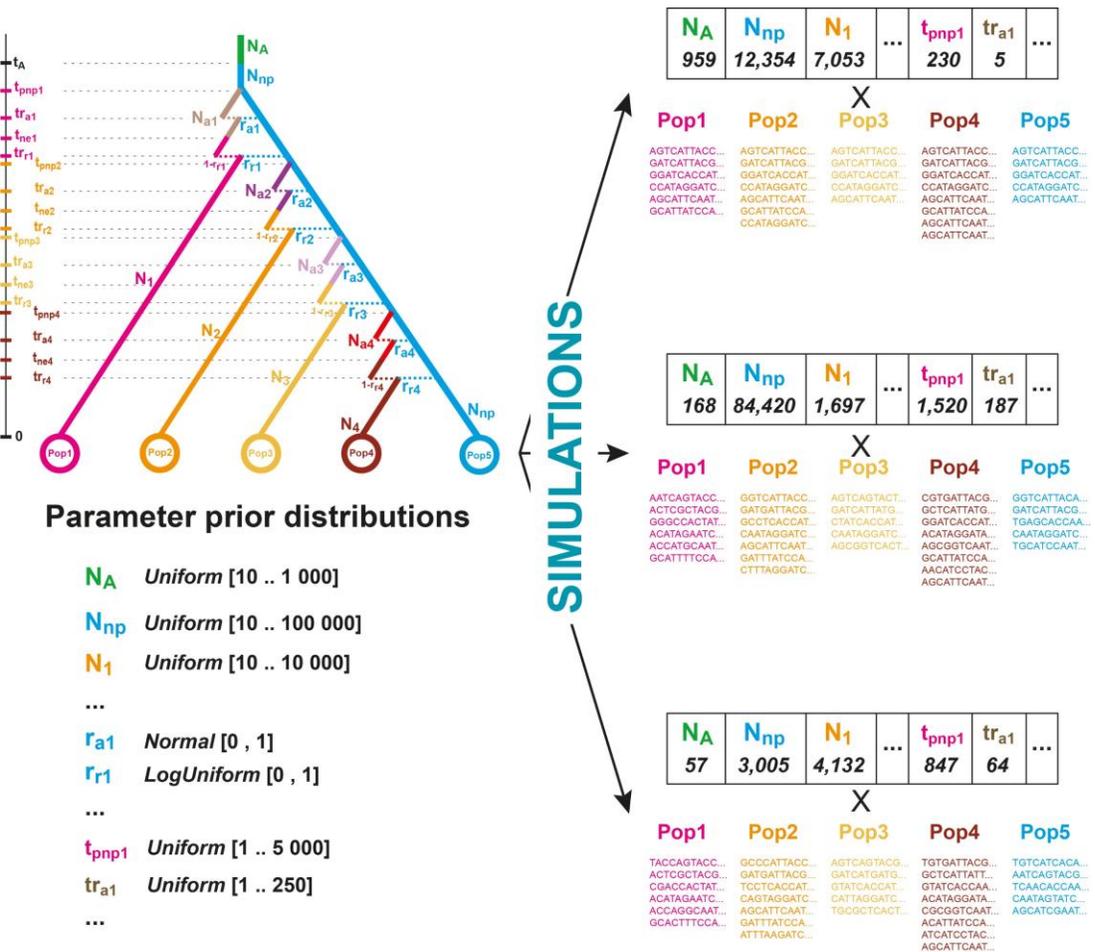*Beaumont et al., Genetics 2002*

# Approximate Bayesian Computation demographic inference



Parameter prior distributions

$N_A$  Uniform [10 .. 1 000]

$N_{np}$  Uniform [10 .. 100 000]

$N_1$  Uniform [10 .. 10 000]

...

$r_{a1}$  Normal [0 , 1]

$r_{r1}$  LogUniform [0 , 1]

...

$t_{pnp1}$  Uniform [1 .. 5 000]

$tr_{a1}$  Uniform [1 .. 250]

...

Tavaré et al., *Plos Gen* 1997
Pritchard et al., *Genetics* 2000
Beaumont et al., *Genetics* 2002

# Approximate Bayesian Computation demographic inference

Tavaré et al., _Plos Gen_ 1997
Pritchard et al., _Genetics_ 2000
Beaumont et al., _Genetics_ 2002

# Approximate Bayesian Computation demographic inference

Tavaré et al., _Plos Gen_ 1997
Pritchard et al., _Genetics_ 2000
Beaumont et al., _Genetics_ 2002

# Approximate Bayesian Computation for complex admixture history

**Approximate Bayesian Computation (Tavaré et al. 1997) represent an alternative ?**

$$L(ModelParams|SumStat_{obs}) \propto P(SumStat_{obs}|ModelParams)\, P(ModelParams)$$



**If <u>simulations are feasible</u>**

**If <u>summary-statistics are informative</u> about model-parameters,**
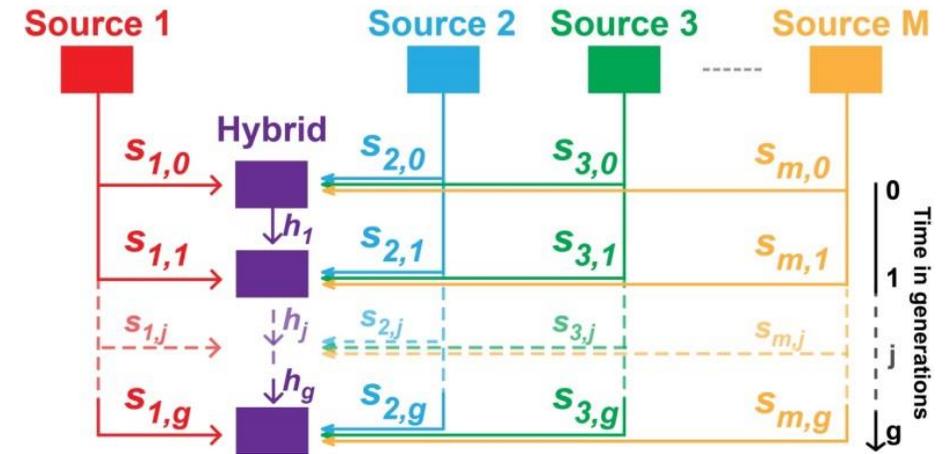
**► ABC inferences may be successful for the reconstruction of complex admixture histories from genetic data**

*Tavaré et al., <u>Plos Gen</u> 1997*
*Pritchard et al., <u>Genetics</u> 2000*
*Beaumont et al., <u>Genetics</u> 2002*

# Complex admixture histories reconstructed with *MetHis*-ABC

**If <u>simulations are feasible…</u>**

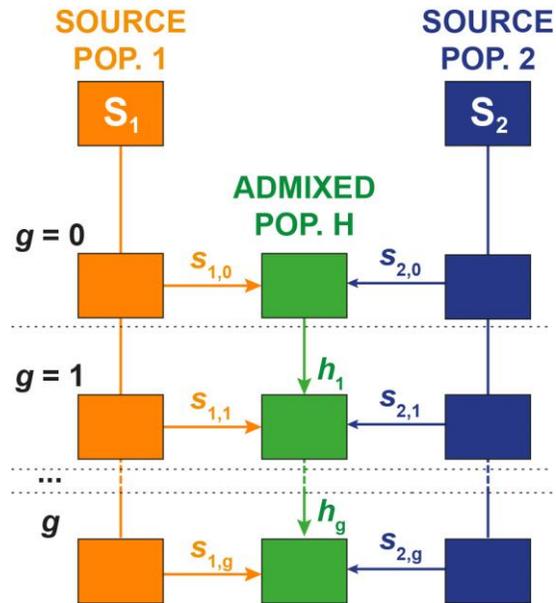Simulating complex admixture histories under the coalescent is often not trivial:

Different pedigree for each independent locus instead of a single pedigree having, in reality, produced all observed gene genealogies (see Wakeley et al. 2012)



*MetHis*

**genetic data simulator under complex admixture models**

# *MetHis*: genetic data simulator under complex admixture histories



*MetHis* simulates a **random-matting admixed population of diploid size $N_g$** at generation $g$, **forward-in-time centered on individuals**

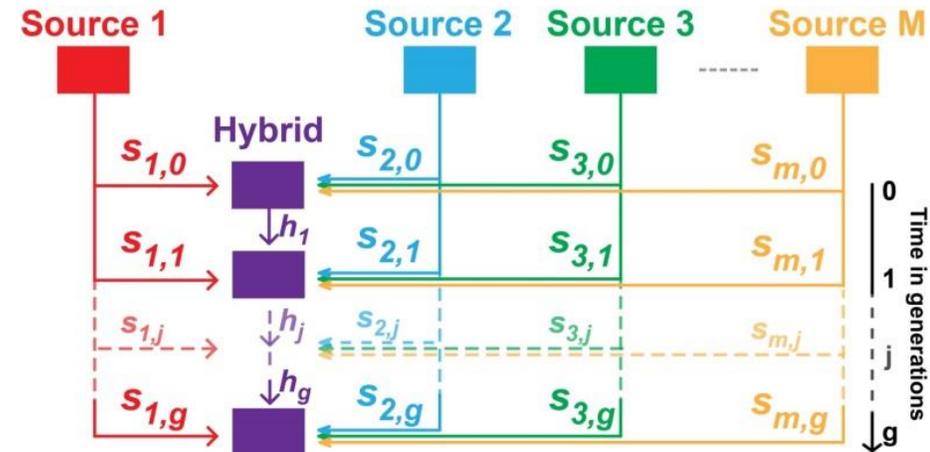*MetHis* simulates any number of **autosomal independent** genetic loci (SNPs or microsatellites)

1. For each $N_{g+1}$ individuals in population H at generation g+1:
Draw randomly parents in populations S1, S2, and H in proportions $s_{1,g}$, $s_{2,g}$ and $h_g$

2. Create haploid gametes for each parent by randomly drawing alleles at each loci

3. Pair gametes to create a new individual in population H, avoiding selfing

*MetHis* simulates a mutation model for microsatellites (GSMM with in/del)

$\forall \, m \in \{1,..,M\}, s_{m,0} \in [0,1]$ such that
$$\sum_{m=1}^{M} s_{m,0} = 1$$

$\forall \, m \in \{1,..,M\}, \forall \, j \in \{1,..,g\},$
$h_j$ and $s_{m,j} \in [0,1]$ such that
$$h_j + \sum_{m=1}^{M} s_{m,j} = 1$$

# Approximate Bayesian Computation for complex admixture history

**Approximate Bayesian Computation (Tavaré et al. 1997) represent an alternative ?**

$$L(ModelParams|SumStat_{obs}) \propto P(SumStat_{obs}|ModelParams)\, P(ModelParams)$$



If <u>simulations are feasible</u>

If <u>summary-statistics are informative</u> about model-parameters,

► ABC inferences may be successful for the reconstruction of complex admixture histories from genetic data

*Tavaré et al., Plos Gen 1997*
*Pritchard et al., Genetics 2000*
*Beaumont et al., Genetics 2002*

# *MetHis*: summary-statistics calculator for ABC inference

**If <u>summary-statistics are informative</u> about model-parameters…**

**"Classical" population genetics summary-statistics**
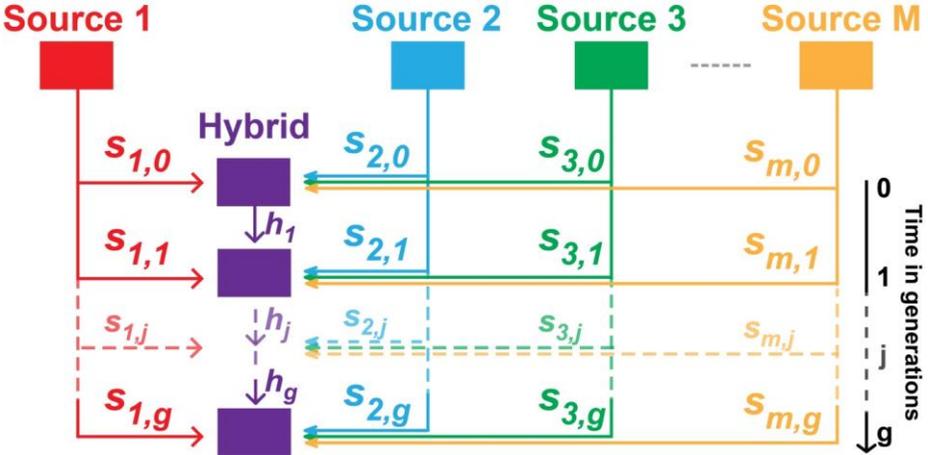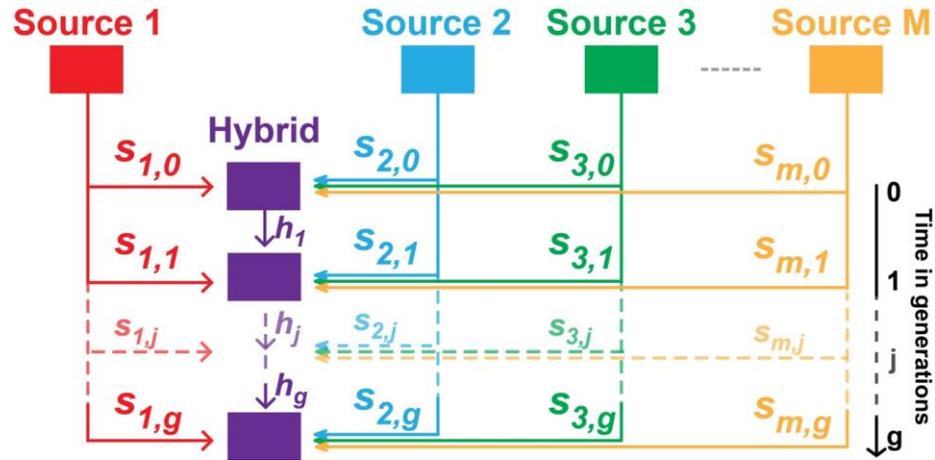
**Expected Heterozygosity** (Nei, 1978)

Inbreeding coefficient $F$ (Danecek et al. 2011)

Multilocus pairwise $F_{ST}$ (Weir and Cockerham, 1984)

$f_3$ **(Admixed; S1, S2)** (Patterson et al., 2012)

Individual pairwise **Allele Sharing Dissimilarities** (Bowcock et al., 1994)

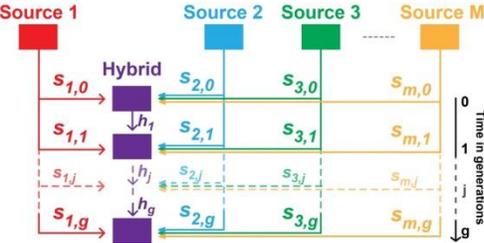# MetHis: Distribution of admixture fractions informative about model-parameters



For a random locus in a random individual in the hybrid population at generation g > 0,
let $H_{i,g}$ be the probability that this locus originally derives from **Source Population _i_** (with _i_ in {1,.., _M_}):

$$H_{i,1} = \begin{cases} 1 & \text{if } Y = S_iS_i, \text{ with } P[Y = S_iS_i] = s_{i,0}^2 \\ \dfrac{1}{2} & \text{if } Y = S_iS_j, \text{ with } P[Y = S_iS_j] = 2s_{i,0}s_{j,0} \\ 0 & \text{if } Y = S_jS_j, \text{ with } P[Y = S_jS_j] = s_{j,0}^2 \\ 0 & \text{if } Y = S_jS_l, \text{ with } P[Y = S_jS_l] = 2s_{j,0}s_{l,0}, \end{cases}$$

# MetHis: Distribution of admixture fractions informative about model-parameters



$\forall\, m \in \{1,..,M\},\, s_{m,0} \in [0,1]$ such that

$$\sum_{m=1}^{M} s_{m,0} = 1$$

$\forall\, m \in \{1,..,M\},\, \forall\, j \in \{1,..,g\},$
$h_j$ and $s_{m,j} \in [0,1]$ such that

$$h_j + \sum_{m=1}^{M} s_{m,j} = 1$$

For a random locus in a random individual in the hybrid population at generation g > 0,
**let $H_{i,g}$ be the probability that this locus originally derives from Source Population *i*** (with *i* in {1,.., *M*}):

$$H_{i,1} = \begin{cases} 1 & \text{if } Y = S_i S_i, \text{ with } P[Y = S_i S_i] = s_{i,0}^2 \\ \dfrac{1}{2} & \text{if } Y = S_i S_j, \text{ with } P[Y = S_i S_j] = 2s_{i,0}s_{j,0} \\ 0 & \text{if } Y = S_j S_j, \text{ with } P[Y = S_j S_j] = s_{j,0}^2 \\ 0 & \text{if } Y = S_j S_l, \text{ with } P[Y = S_j S_l] = 2s_{j,0}s_{l,0}, \end{cases}$$

and

$$H_{i,g} = \begin{cases} 1 & \text{if } Y = S_i S_i, \text{ with } P[Y = S_i S_i] = s_{i,g-1}^2 \\ \dfrac{H_{i,g-1} + 1}{2} & \text{if } Y = S_i H, \text{ with } P[Y = S_i H] = 2s_{i,g-1}h_{g-1} \\ \dfrac{1}{2} & \text{if } Y = S_i S_j, \text{ with } P[Y = S_i S_j] = 2s_{i,g-1}s_{j,g-1} \\ \dfrac{H_{i,g-1}^{(1)} + H_{i,g-1}^{(2)}}{2} & \text{if } Y = HH, \text{ with } P[Y = HH] = h_{g-1}^2 \\ \dfrac{H_{i,g-1}}{2} & \text{if } Y = S_j H, \text{ with } P[Y = S_j H] = 2s_{j,g-1}h_{g-1} \\ 0 & \text{if } Y = S_j S_j, \text{ with } P[Y = S_j S_j] = s_{j,g-1}^2 \\ 0 & \text{if } Y = S_j S_l, \text{ with } P[Y = S_j S_l] = 2s_{j,g-1}s_{l,g-1}. \end{cases}$$

*Verdu and Rosenberg, <u>Genetics</u> 2011*

18

# MetHis: Distribution of admixture fractions informative about model-parameters
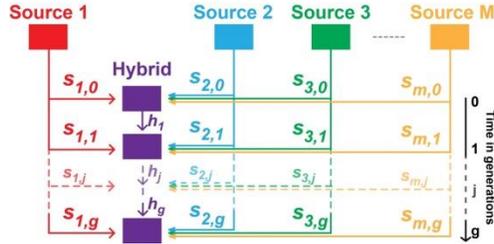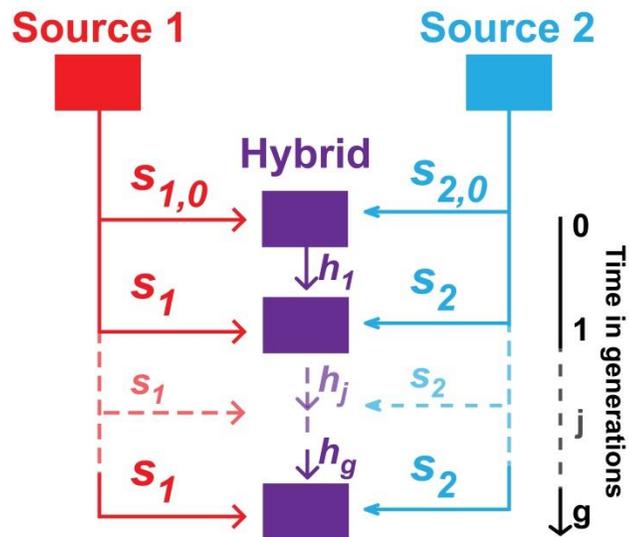


## Distribution of admixture fractions

$$q \in Q_g = \left\{0, \frac{1}{2^g}, \ldots, \frac{2^g - 1}{2^g}, 1\right\}$$

$\forall\, i \geq 1,$

$$P(H_{1,1} = q) = \begin{cases} s_{1,0}^2 & \text{if } q = 1 \\ 2s_{1,0}s_{2,0} & \text{if } q = 1/2 \\ s_{2,0}^2 & \text{if } q = 0. \end{cases}$$

$\forall\, i \geq 1, \forall\, g \geq 2$

$$P(H_{1,g} = q) = h_{g-1}^2 \sum_{r=0}^{2^{g-1}} \left[P\left(H_{1,g-1} = \frac{r}{2^{g-1}}\right)P\left(H_{1,g-1} = \frac{2^g q - r}{2^{g-1}}\right)\right] + 2s_{2,g-1}h_{g-1}P(H_{1,g-1} = 2q) + 2s_{1,g-1}h_{g-1}P(H_{1,g-1} = 2q - 1) + I_g(q),$$

with

$$I_g(q) = \begin{cases} s_{1,g-1}^2 & \text{if } q = 1 \\ 2s_{1,g-1}s_{2,g-1} & \text{if } q = 1/2 \\ s_{2,g-1}^2 & \text{if } q = 0 \\ 0 & \text{otherwise.} \end{cases}$$

*Verdu and Rosenberg, Genetics 2011*

**Law of total expectation**: $E[H_{i,g}] = E_Y\left[E[H_{i,g}|Y]\right] = \sum_{y \in A} P(Y = y)E[H_{i,g}| Y = y].$

***k*th - moment**

$$\forall\, k \geq 1, \forall\, i \geq 1, \qquad E[H_{i,1}^k] = s_{i,0}^2 + \frac{s_{i,0}}{2^{k-1}} \sum_{\substack{j=1 \\ j \neq i}}^{m} s_{j,0}\,.$$

$$\forall\, k \geq 1, \forall\, i \geq 1, \forall\, g \geq 2 \quad E[H_{i,g}^k] = s_{i,g-1}^2 + \frac{s_{i,g-1}h_{g-1}}{2^{k-1}}\left(\sum_{r=0}^{k}\binom{k}{r}E[H_{i,g-1}^r]\right) + \frac{s_{i,g-1}}{2^{k-1}}\sum_{\substack{j=1 \\ j \neq i}}^{m} s_{j,g-1} + \frac{h_{g-1}^2}{2^{k}}\left(\sum_{r=0}^{k}\binom{k}{r}E[H_{i,g-1}^r]E[H_{i,g-1}^{k-r}]\right) + \left(\frac{h_{g-1}}{2^{k-1}}\sum_{\substack{j=1 \\ j \neq i}}^{m} s_{j,g-1}\right)E[H_{i,g-1}^k].$$

*Verdu and Rosenberg, Genetics 2011*

$$E\big[H_{1,g}\big] = \begin{cases} s_{1,0}, & g = 1 \\ s_{1,0}h^{g-1} + s_1\dfrac{1 - h^{g-1}}{1 - h}, & g \ge 2 \end{cases}.$$

$$V[H_{1,g}] = \begin{cases} \dfrac{s_{1,0}(1 - s_{1,0})}{2}, & g = 1 \\ A_1 + A_2 h^{g-1} + A_3\left(\dfrac{h}{2}\right)^{g-1} + A_4\left(\dfrac{h}{2}\right)^{g-1}\sum_{i=1}^{g-1}(2h)^i - \left(s_{1,0}h^{g-1} + s_1\dfrac{1 - h^{g-1}}{1 - h}\right)^2, & g \ge 2. \end{cases}$$

*Verdu and Rosenberg, Genetics 2011*

# MetHis: Distribution of admixture fractions informative about model-parameters



$$E[H_{1,g}] = \begin{cases} s_{1,0}, & g = 1 \\ s_{1,0}h^{g-1} + s_1 \dfrac{1-h^{g-1}}{1-h}, & g \geq 2 \end{cases}.$$

$$V[H_{1,g}] = \begin{cases} \dfrac{s_{1,0}(1-s_{1,0})}{2}, & g = 1 \\ A_1 + A_2 h^{g-1} + A_3 \left(\dfrac{h}{2}\right)^{g-1} + A_4 \left(\dfrac{h}{2}\right)^{g-1}\sum_{i=1}^{g-1}(2h)^i - \left(s_{1,0}h^{g-1} + s_1\dfrac{1-h^{g-1}}{1-h}\right)^2, & g \geq 2. \end{cases}$$

**Use the distribution of admixture fractions as an ABC-informative summary statistics !**

*Verdu and Rosenberg, Genetics 2011*

22

# *MetHis*: summary-statistics calculator for ABC inference

**If <u>summary-statistics are informative</u> about model-parameters…**

**"Classical" population genetics summary-statistics**

      **Expected Heterozygosity** (Nei, 1978)

      Inbreeding coefficient $F$ (Danecek et al. 2011)

      Multilocus pairwise $F_{ST}$ (Weir and Cockerham, 1984)

      $f_3$ **(Admixed; S1, S2)** (Patterson et al., 2012)

      Individual pairwise **Allele Sharing Dissimilarities** (Bowcock et al., 1994)

**Individual admixture fractions**
      **Min** and **Max** admixture fractions
      **10% quantiles** of individual admixture fractions
      **Mode**, **Mean**, **Variance**, **Kurtosis**, **Skewness** of individual admixture fractions

LE CODE NOIR
OU
EDIT DU ROY,
SERVANT DE REGLEMENT

Transatlantic Slave Trade:

► Recent admixture history (~20 generations)

► Variable migration histories to the Americas

► Variable slavery histories in the Americas



*Gravel, Genetics 2012*



*Baharian et al., Plos Gen 2016*

Previous studies tested 2 admixture pulses at most, **with ML methods using >> 1 million SNPs**

They found that 2 admixture pulses from Europe were most likely… but cannot try more complex models.

**21 generations before present ($g = 20$)**

**= Admixed population founded ~ 1500s**



**3 admixture pulses**

**21 generations before present ($g = 20$)**

**= Admixed population founded ~ 1500s**

# Nine competing models for the admixture history of the ASW and the ACB

**21 generations before present ($g = 20$)**

**= Admixed population founded ~ 1500s**

**21 generations before present ($g = 20$)**

**= Admixed population founded ~ 1500s**

*Fortes-Lima et al., Mol Ecol Res 2021*

- 100,000 independent SNPs
- $N_g$ = constant 1,000 individuals (for simplicity here, *MetHis* option)
- 10,000 simulations per model = 90,000 simulations total
- Sample 90 Africans, 89 Europeans, and 50 individuals in the admixed population.

**For each simulated data set, calculate 24 population genetics summary statistics**

**27 cores -> 3 days calculation total.**

**2/3 of this time = summary statistics calculation…**

**Simulation time increases with Ne much more rapidly than with g.**
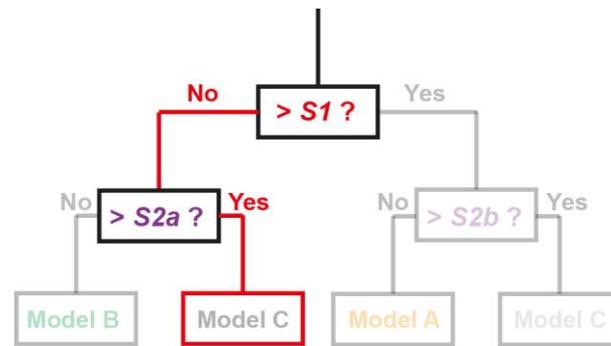
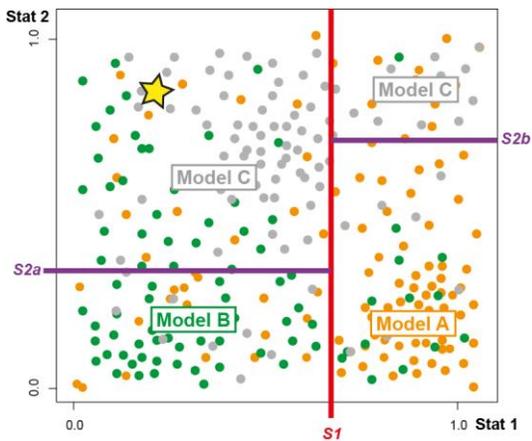# Random Forest algorithm for model choice (Breiman 2001)
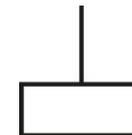
Draw randomly stats and order
them by variance explained

Build Decision Tree 1



Stat 2
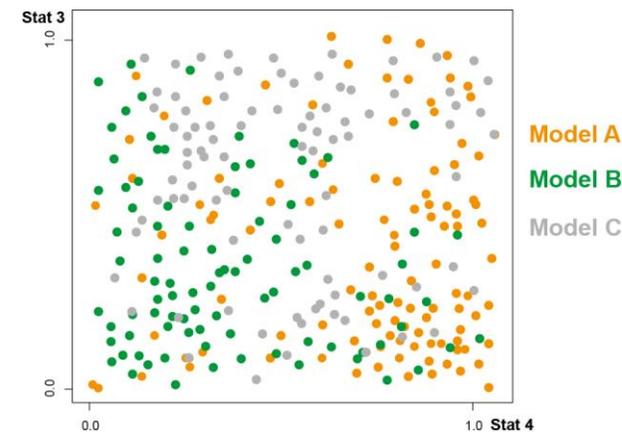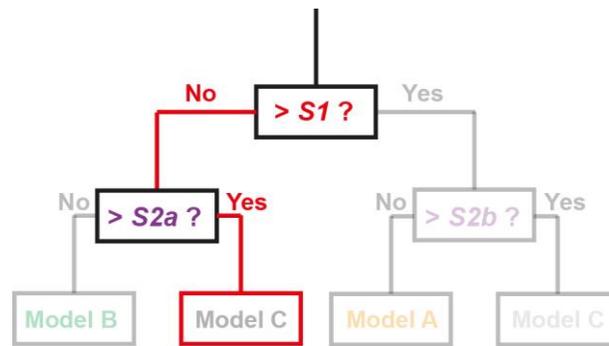
Model A
Model B
Model C

Stat 1

Draw randomly stats and order
them by variance explained

Build Decision Tree 1

Draw randomly stats and order
them by variance explained

Build Decision Tree 1

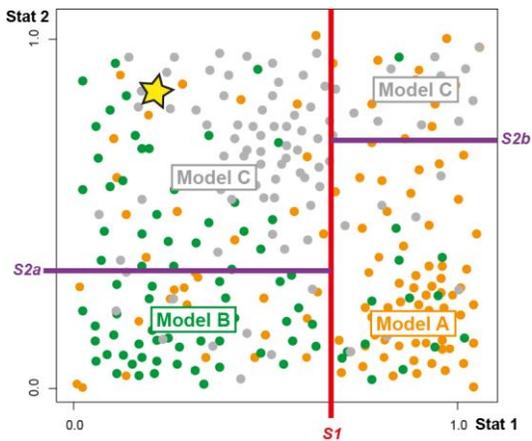# Random Forest algorithm for model choice (Breiman 2001)

Draw randomly stats and order
them by variance explained

Build Decision Tree 1



Decision Tree 1

Draw randomly stats and order
them by variance explained

Build Decision Tree 1



Prediction with decision tree 1 = Model C

# Random Forest algorithm for model choice (Breiman 2001)

Draw randomly stats and order
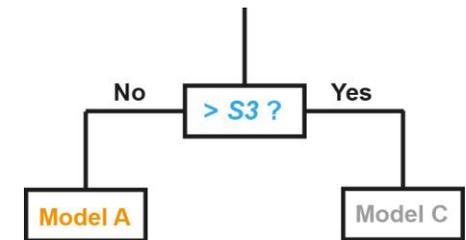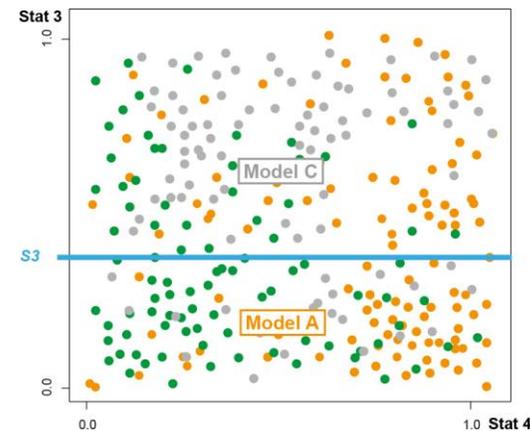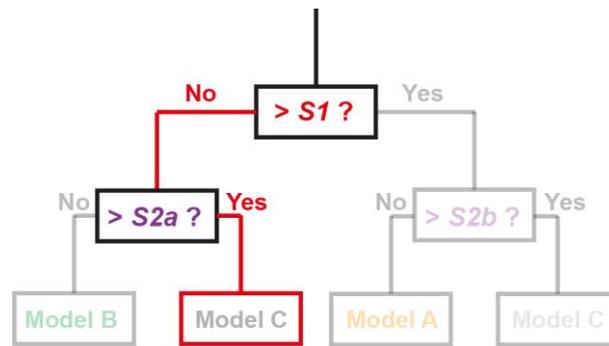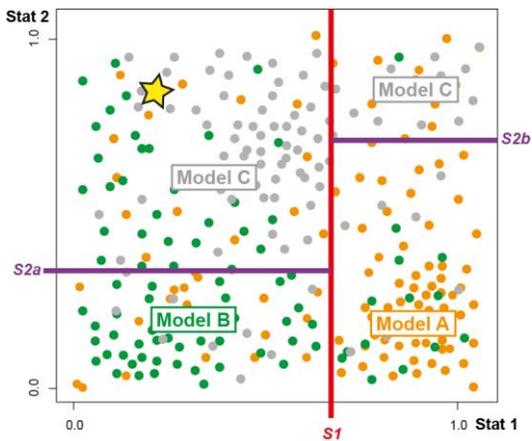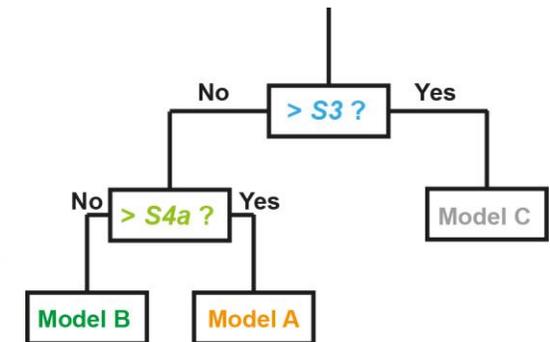them by variance explained

Build Decision Tree 1

Draw randomly other stats and
order them by variance explained

Build Decision Tree 2



Prediction with decision tree 1 = Model C

# Random Forest algorithm for model choice (Breiman 2001)

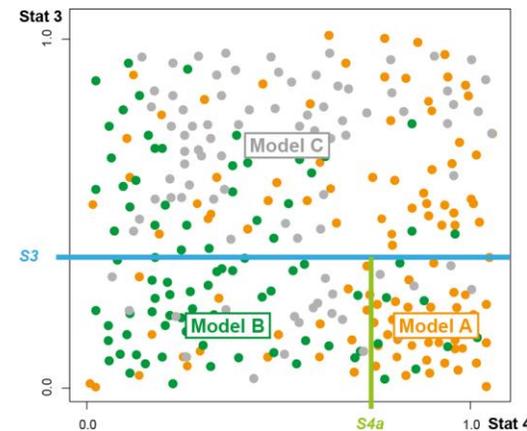Draw randomly stats and order them by variance explained

Build Decision Tree 1

Draw randomly other stats and order them by variance explained

Build Decision Tree 2



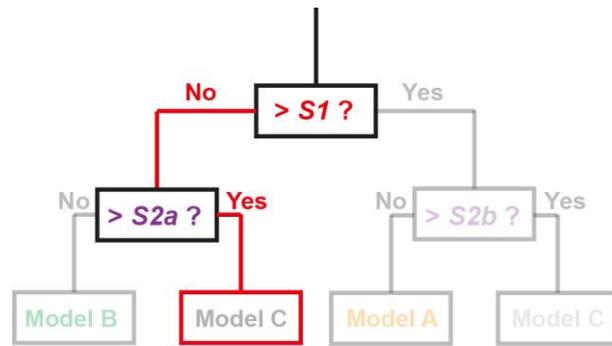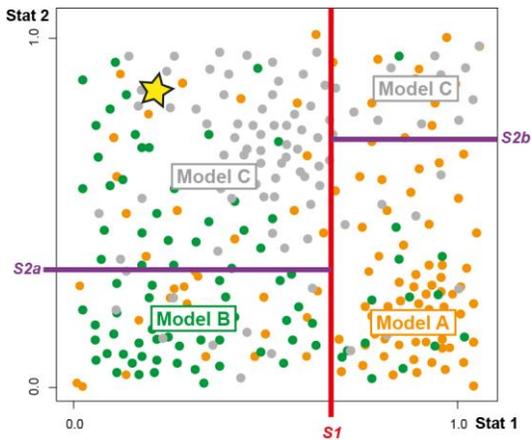Prediction with decision tree 1 = Model C

# Random Forest algorithm for model choice (Breiman 2001)

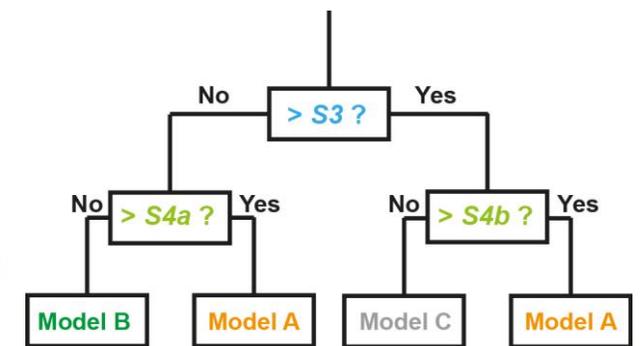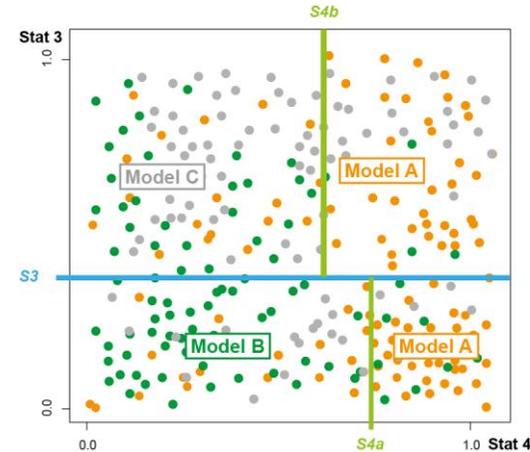Draw randomly stats and order them by variance explained

Build Decision Tree 1

Draw randomly other stats and order them by variance explained

Build Decision Tree 2



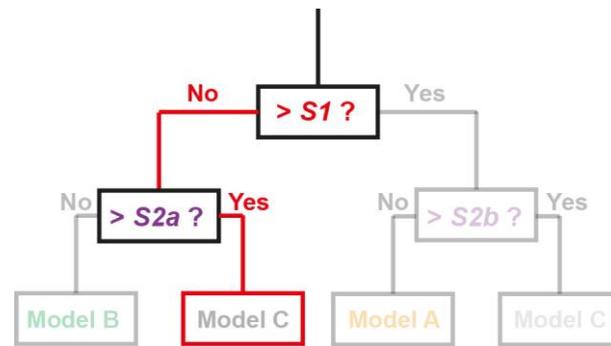Prediction with decision tree 1 = Model C

# Random Forest algorithm for model choice (Breiman 2001)

Draw randomly stats and order them by variance explained

Build Decision Tree 1

Draw randomly other stats and order them by variance explained

Build Decision Tree 2



Prediction with decision tree 1 = Model C
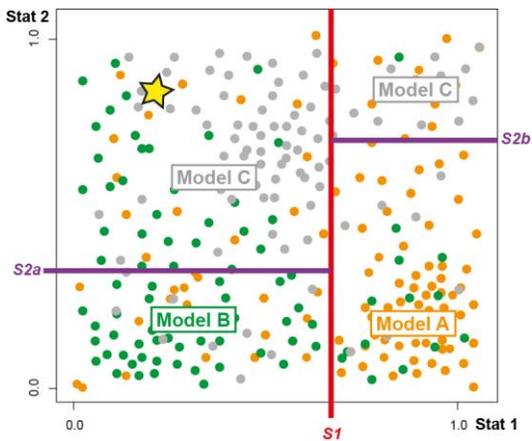


Decision Tree 2

# Random Forest algorithm for model choice (Breiman 2001)

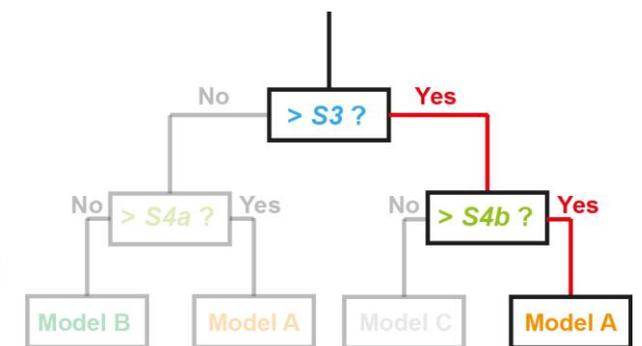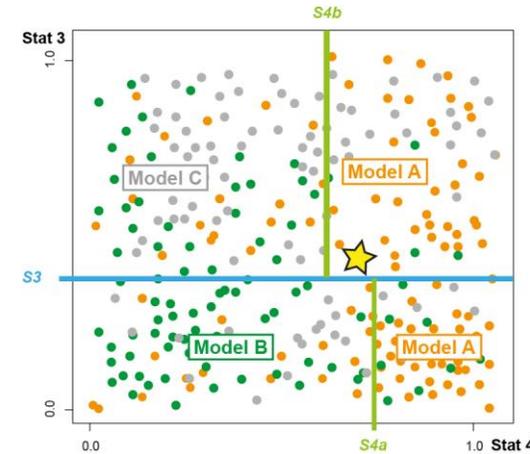Draw randomly stats and order them by variance explained

Build Decision Tree 1

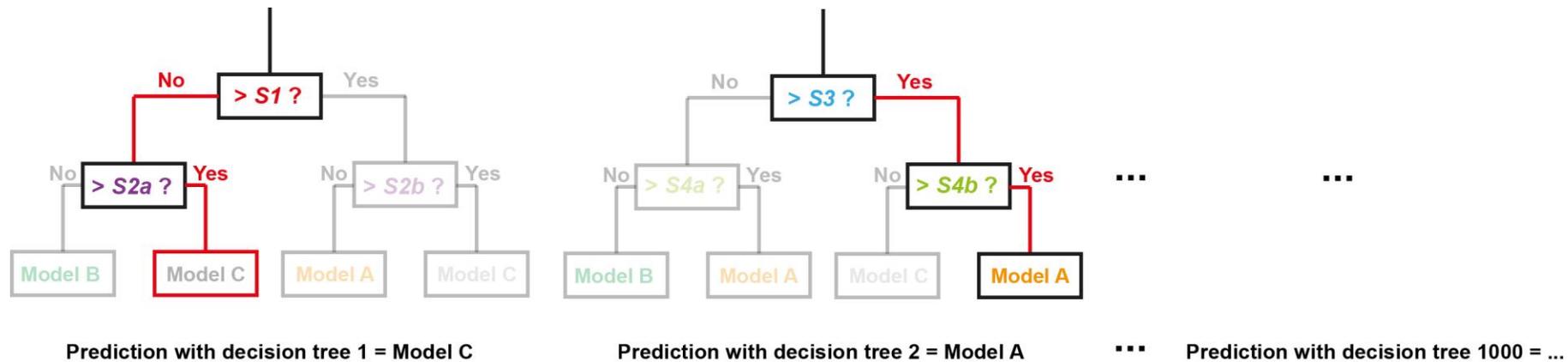Draw randomly other stats and order them by variance explained

Build Decision Tree 2



Prediction with decision tree 1 = Model C

Prediction with decision tree 2 = Model A

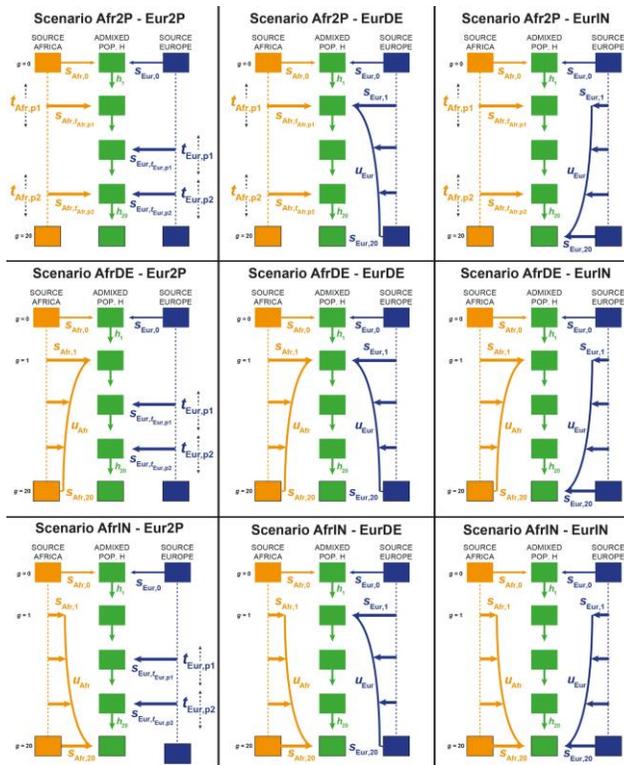# Random Forest algorithm for model choice (Breiman 2001)

**Repeat and build a random-forest of decision trees**



Prediction with decision tree 1 = Model C     Prediction with decision tree 2 = Model A     ⋯     Prediction with decision tree 1000 = ...

**Prediction for the observed data is, e.g., the majority of votes across 1000 trees in the random forest**

# Results: *MetHis*-ABC model-choice with Random-Forest <u>a priori without observed data</u>

- Random Forest ABC *(abcrf* package in *R, Pudlo et al. 2016)* – 10,000 *MetHis* sims/scenario – 1,000 trees
- Cross-validation all simulations
- 100,000 independent SNPs
- Sample sizes:
  - Afr Source : 90 indivs
  - Eur Source : 89 indivs
  - Admixed population H : 50 indivs



| RF-ABC Predicted model | Afr2P - Eur2P | AfrDE - Eur2P | AfrIN - Eur2P | Afr2P - EurDE | AfrDE - EurDE | AfrIN - EurDE | Afr2P - EurIN | AfrDE - EurIN | AfrIN - EurIN |
|---|---|---|---|---|---|---|---|---|---|
| AfrIN - EurIN | 1.2% | 2.7% | 2.9% | 2.9% | 0.1% | 10.8% | 2.8% | 9.7% | 61.4% |
| AfrDE - EurIN | 1.6% | 6.6% | 0.8% | 1.7% | 1.7% | 1.4% | 9.5% | 73.2% | 17.7% |
| Afr2P - EurIN | 5.9% | 2.2% | 0.0% | 5.1% | 0.3% | 0.0% | 76.9% | 8.8% | 0.3% |
| AfrIN - EurDE | 2.0% | 1.5% | 9.9% | 6.0% | 2.0% | 72.2% | 0.9% | 1.5% | 18.4% |
| AfrDE - EurDE | 5.9% | 15.8% | 2.0% | 15.6% | 77.7% | 4.8% | 1.4% | 4.8% | 1.7% |
| Afr2P - EurDE | 11.2% | 1.8% | 0.4% | 58.2% | 7.6% | 0.6% | 6.9% | 1.1% | 0.1% |
| AfrIN - Eur2P | 5.7% | 4.6% | 76.3% | 1.8% | 0.4% | 9.1% | 0.0% | 0.0% | 0.2% |
| AfrDE - Eur2P | 11.2% | 57.3% | 6.7% | 1.8% | 7.6% | 0.9% | 0.6% | 0.6% | 0.2% |
| Afr2P - Eur2P | 55.2% | 7.6% | 0.9% | 7.0% | 2.7% | 0.2% | 0.9% | 0.2% | 0.1% |

**True model**

Prior error rate = 32.41%, Model-choice error a priori = 8/9 = 88.89%

- Random Forest ABC *(abcrf* package in *R, Pudlo et al. 2016)* – 10,000 *MetHis* sims/scenario – 1,000 trees
- Cross-validation all simulations
- 100,000 independent SNPs
- Sample sizes:
  - **Afr Source : 90 indivs**
  - **Eur Source : 89 indivs**
  - **Admixed population H : 50 indivs**

**Model – Nestedness** (Robert et al. 2011)

**• Erroneous model-choice increased among scenarios qualitatively similar**



| RF-ABC Predicted model | Afr2P - Eur2P | AfrDE - Eur2P | AfrIN - Eur2P | Afr2P - EurDE | AfrDE - EurDE | AfrIN - EurDE | Afr2P - EurIN | AfrDE - EurIN | AfrIN - EurIN |
|---|---|---|---|---|---|---|---|---|---|
| AfrIN - EurIN | 1.2% | 2.7% | 2.9% | 2.9% | 0.1% | 10.8% | 2.8% | 9.7% | 61.4% |
| AfrDE - EurIN | 1.6% | 6.6% | 0.8% | 1.7% | 1.7% | 1.4% | 9.5% | 73.2% | 17.7% |
| Afr2P - EurIN | 5.9% | 2.2% | 0.0% | 5.1% | 0.3% | 0.0% | 76.9% | 8.8% | 0.3% |
| AfrIN - EurDE | 2.0% | 1.5% | 9.9% | 6.0% | 2.0% | 72.2% | 0.9% | 1.5% | 18.4% |
| AfrDE - EurDE | 5.9% | 15.8% | 2.0% | 15.6% | 77.7% | 4.8% | 1.4% | 4.8% | 1.7% |
| Afr2P - EurDE | 11.2% | 1.8% | 0.4% | 58.2% | 7.6% | 0.6% | 6.9% | 1.1% | 0.1% |
| AfrIN - Eur2P | 5.7% | 4.6% | 76.3% | 1.8% | 0.4% | 9.1% | 0.0% | 0.0% | 0.2% |
| AfrDE - Eur2P | 11.2% | 57.3% | 6.7% | 1.8% | 7.6% | 0.9% | 0.6% | 0.6% | 0.2% |
| Afr2P - Eur2P | 55.2% | 7.6% | 0.9% | 7.0% | 2.7% | 0.2% | 0.9% | 0.2% | 0.1% |

**True model**

Prior error rate = 32.41%, Model-choice error a priori = 8/9 = 88.89%

- Random Forest ABC *(abcrf* package in *R, Pudlo et al. 2016)* – 10,000 *MetHis* sims/scenario – 1,000 decision trees in the random forest
- Cross-validation all simulations

- 100,000 independent SNPs
- Sample sizes:
  Afr Source : 90 indivs
  Eur Source : 89 indivs
  Admixed population H : 50 indivs



Prior error rate = **32.41%**
Model-choice error a priori = 8/9 = 88.89%

- Random Forest ABC *(abcrf* package in *R, Pudlo et al. 2016)* – 10,000 *MetHis* sims/scenario – 1,000 decision trees in the random forest
- Cross-validation all simulations

- **100,000 independent SNPs**
- **Sample sizes:**
  - Afr Source : 90 indivs
  - Eur Source : 89 indivs
  - Admixed population H : 50 indivs

- **50,000 independent SNPs**
- **Sample sizes:**
  - Afr Source : 90 indivs
  - Eur Source : 89 indivs
  - Admixed population H : 50 indivs



Prior error rate = **32.41%**
Model-choice error a priori = 8/9 = 88.89%

Prior error rate = **33.53%**
Model-choice error a priori = 8/9 = 88.89%

# Results: *MetHis*-ABC model-choice with Random-Forest <u>a priori without observed data</u>

- Random Forest ABC *(abcrf* package in *R, Pudlo et al. 2016)* – 10,000 *MetHis* sims/scenario – 1,000 decision trees in the random forest
- Cross-validation all simulations

- **100,000 independent SNPs**
- **Sample sizes:**
    - Afr Source : 90 indivs
    - Eur Source : 89 indivs
    - Admixed population H : 50 indivs

- **50,000 independent SNPs**
- **Sample sizes:**
    - Afr Source : 90 indivs
    - Eur Source : 89 indivs
    - Admixed population H : 50 indivs

- **10,000 independent SNPs**
- **Sample sizes:**
    - Afr Source : 90 indivs
    - Eur Source : 89 indivs
    - Admixed population H : 50 indivs



Prior error rate = **32.41%**
Model-choice error a priori = 8/9 = 88.89%

Prior error rate = **33.53%**
Model-choice error a priori = 8/9 = 88.89%

Prior error rate = **37.93%**
Model-choice error a priori = 8/9 = 88.89%

# Results: *MetHis*-ABC model-choice with Random-Forest <u>a priori without observed data</u>
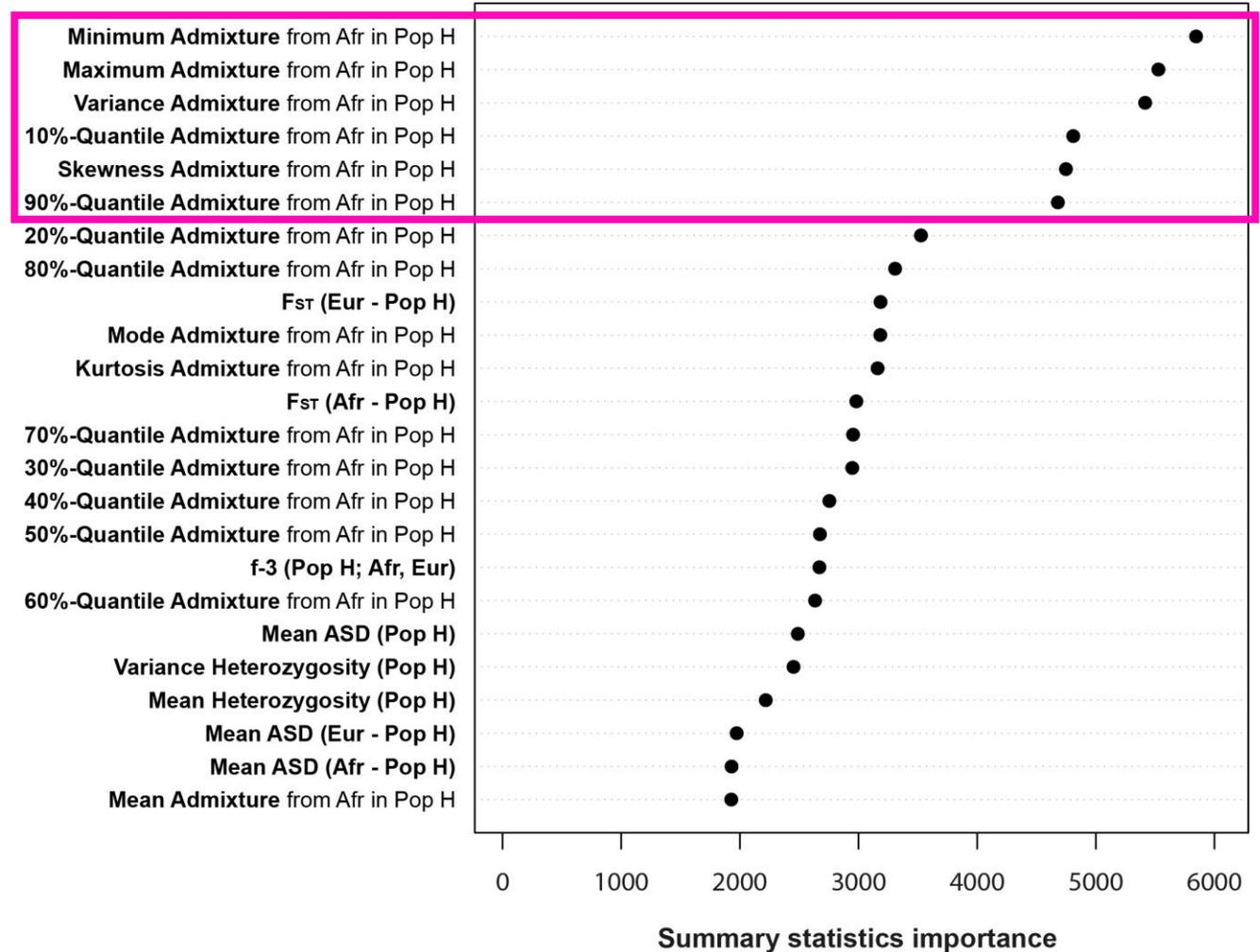
- Random Forest ABC *(abcrf* package in *R, Pudlo et al. 2016)* – 10,000 *MetHis* sims/scenario – 1,000 decision trees in the random forest
- Cross-validation all simulations

- 100,000 independent SNPs
- Sample sizes:
   - Afr Source : 90 indivs
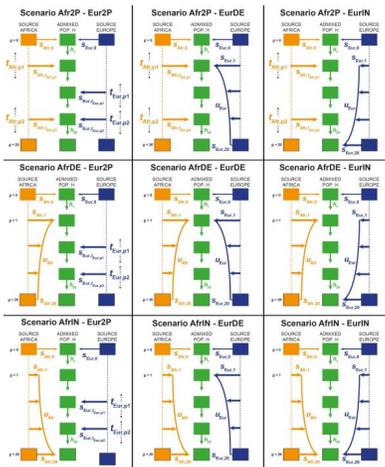   - Eur Source : 89 indivs
   - Admixed population H : 50 indivs

- 100,000 independent SNPs
- Sample sizes:
   - Afr Source : 18 indivs
   - Eur Source : 18 indivs
   - Admixed population H : 10 indivs



Prior error rate = **32.41%**
Model-choice error a priori = 8/9 = 88.89%

Prior error rate = **48.39%**
Model-choice error a priori = 8/9 = 88.89%

**<u>ABC relies on summary-statistics informativeness rather than absolute amount of data</u>**

*Fortes-Lima et al., <u>Mol Ecol Res</u> 2021*

# Results: *MetHis*-ABC model-choice with Random-Forest <u>a priori without observed data</u>

- Random Forest ABC *(abcrf* package in *R, Pudlo et al. 2016)* – 10,000 *MetHis* sims/scenario – 1,000 decision trees in the random forest
- Cross-validation all simulations

  - 100,000 independent SNPs
  - Sample sizes:
    - Afr Source : 90 indivs
    - Eur Source : 89 indivs
    - Admixed population H : 50 indivs



Prior error rate = **32.41%**
Model-choice error a priori = 8/9 = 88.89%

- *MetHis*-ABC Random-Forest model-choice a priori powerful to distinguish highly complex historical admixture models

- Errors are made in the parameter-space where models are highly nested and thus biologically similar

- *MetHis*-ABC Random-Forest model-choice performances are robust to reduced SNP and Sample sets    -> ABC relies on summary-statistics informativeness rather than absolute amount of data

- Distribution of admixture fractions is highly informative for admixture history inference
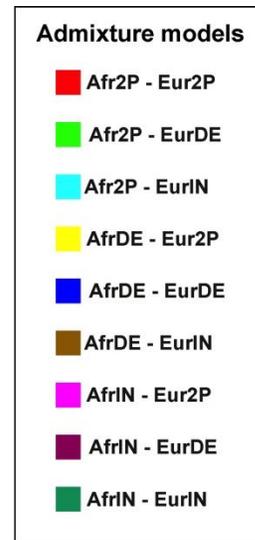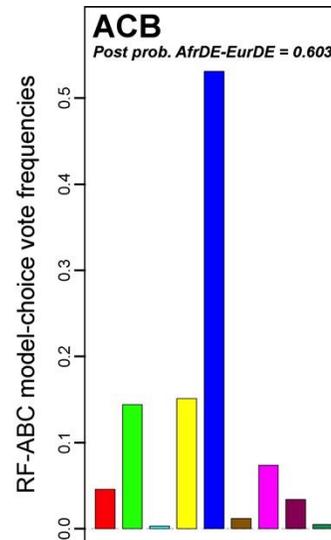
**Model-choice**
**Random Forest ABC**
*abcrf* package in *R* (Pudlo et al. 2016)

**1,000 decision trees in the forest**

**9 competing scenarios**
**10,000 *MetHis* sims/scenario**
**100,000 independent SNPs**
**24 summary-statistics**

**Parameter-inference**
**Neural Network ABC**
*abc* package in *R (*Csilléry et al. 2012)

**Tolerance 1% (1,000 sims closest to obs.)**
**4 neurons in the hidden layer**

**1 best scenario**
**100,000 *MetHis* sims**
**100,000 independent SNPs**
**24 summary-statistics**

**Parameter-inference**
**Neural Network ABC**
*abc* package in *R (*Csilléry et al. 2012)

**Tolerance 1% (1,000 sims closest to obs.)**
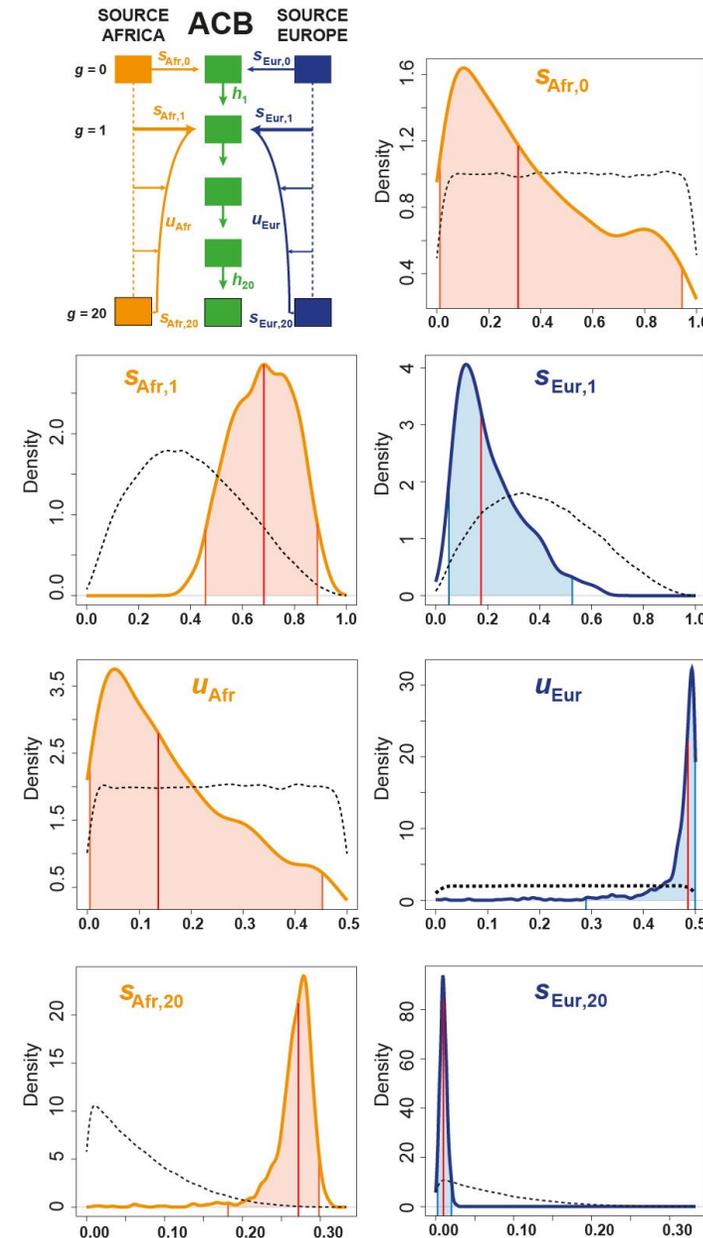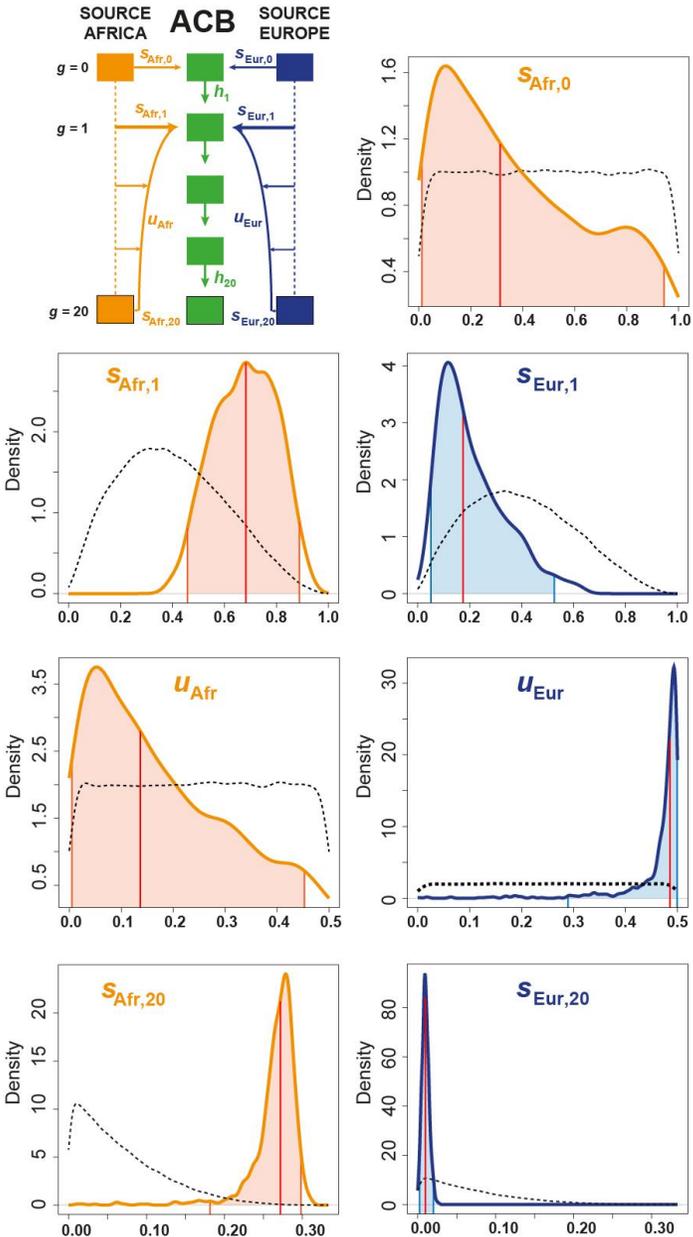**4 neurons in the hidden layer**

**1 best scenario**
**100,000 *MetHis* sims**
**100,000 independent SNPs**
**24 summary-statistics**

**Cross-validation post param error**
**1,000 closest simulations in turn used as controlled**
**pseudo-observed data for NN-ABC param inference**

| *AfrDE-EurDE* parameters | ACB | |
| --- | --- | --- |
| | **Av. absolute Error** | Mean-square Error / Var. |
| $s_{Afr,0}$ | 0.2530 | 1.0070 |
| $s_{Afr,1}$ | **0.1206** | 0.8533 |
| $s_{Afr,20}$ | **0.0274** | 0.4162 |
| $u_{Afr}$ | 0.1166 | 0.9974 |
| $s_{Eur,1}$ | **0.0952** | 1.0526 |
| $s_{Eur,20}$ | **0.0044** | 0.6452 |
| $u_{Eur}$ | 0.1084 | 0.9431 |

*Fortes-Lima et al., Mol Ecol Res 2021*

# Conclusions

***MetHis* – RF-ABC is successful in distinguishing a priori among competing highly complex admixture models**

**Admixture distribution is highly informative for ABC model choice as expected theoretically**

***MetHis* – NN-ABC produces accurate posterior parameter estimation and relatively conservative 95% CI inference**

*Fortes-Lima et al., Mol Ecol Res 2021*

# Conclusions

## Complex genetic admixture histories reconstructed with Approximate Bayesian Computation
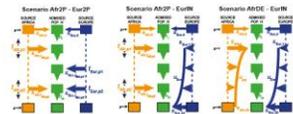
Cesar A. Fortes-Lima,* Romain Laurent,* Valentin Thouzeau, Bruno Toupance, Paul Verdu ✉

Special issue *Machine Learning in Molecular Ecology*

https://github.com/romain-laurent/MetHis



Step 1. Design competing admixture models under the two source populations version of the Verdu and Rosenberg (2011) general admixture model, thought *a priori* to explain the observed data

etc.

Step 2. Simulate genetic data from source populations using prefered methods and tools

Step 3. Build model parameter vectors reference table by drawing parameter values in prior distributions
— *MetHis parameter generator tool*
— Other tools

Step 4. Simulate genetic data with *MetHis* under the models designed in Step 1 with source populations data produced in Step2 and model parameters produced in Step 3.
**Output = Simulated genetic data files**

Step 5. Calculate summary statistics on each simulated dataset and the observed dataset
— *MetHis summary statistics calculator*
— Other tools
**Output = Summary statistics reference table corresponding to model parameter vectors produced in Step 3 and used in Step 4.**

Step 6. Approximate Bayesian Computation inference

Step 6a. ABC model-choice
— Random-Forest ABC with *R* package *abcrf* (Pudlo et al. 2016; Raynal et al. 2019)
— Other ABC model-choice tools
**Output = Model designed in Step 1 best explaining the observed data with associated posterior-probabilities**

Step 6b. ABC posterior parameter estimation for all parameters from a given model designed in Step 1, and simulated in Step 4 with model parameters from Step 3
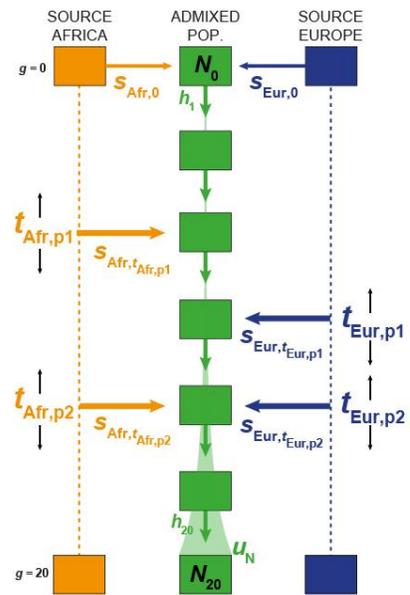— Neural-Network ABC with *R* package *abc* (Csilléry et al. 2012)
— Other ABC posterior parameter estimation tools
**Output = Posterior distribution of model parameters best explaining the observed data**
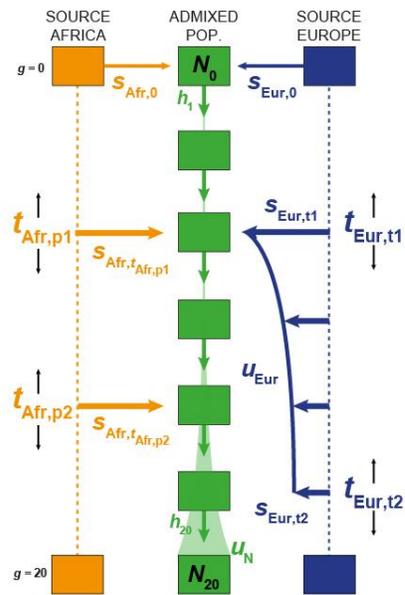
# The admixture histories of Cabo Verde
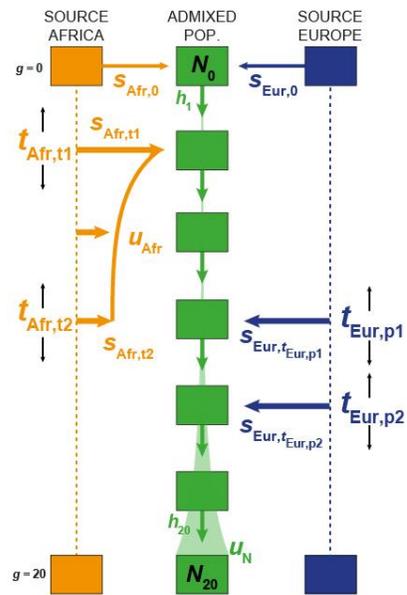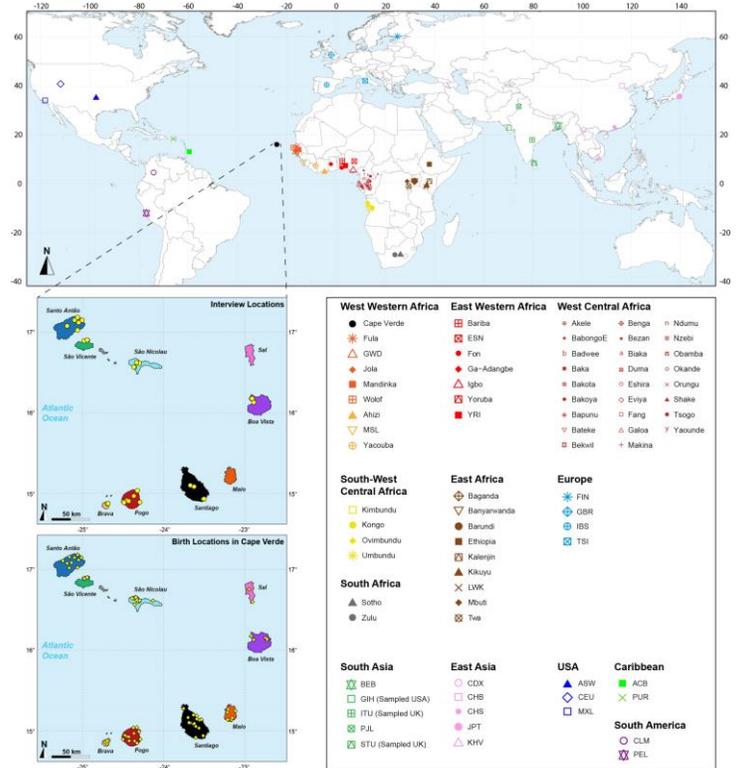## Laurent et al. *Nature Communications* (under revision), BioRxiv https://doi.org/10.1101/2022.04.11.487833

**Reconstructing the admixture history of Cabo Verde with Romain Laurent (UMR7206), Noah Rosenberg (Stanford University) and Marlyse Baptista (University of Michigan)**

# The admixture histories of Cabo Verde
## Laurent et al. *Nature Communications* (under revision), BioRxiv https://doi.org/10.1101/2022.04.11.487833

# Ongoing !

**Microsatellites and *Terminalia superba* hybridization zone in the « Dahomey Gap » (West Africa)**
**With Romain Laurent (UMR7206) and Olivier Hardy (Univ. Bruxelles)**

# Great Many Thanks !