# Population genetics: Inference using trees of individuals
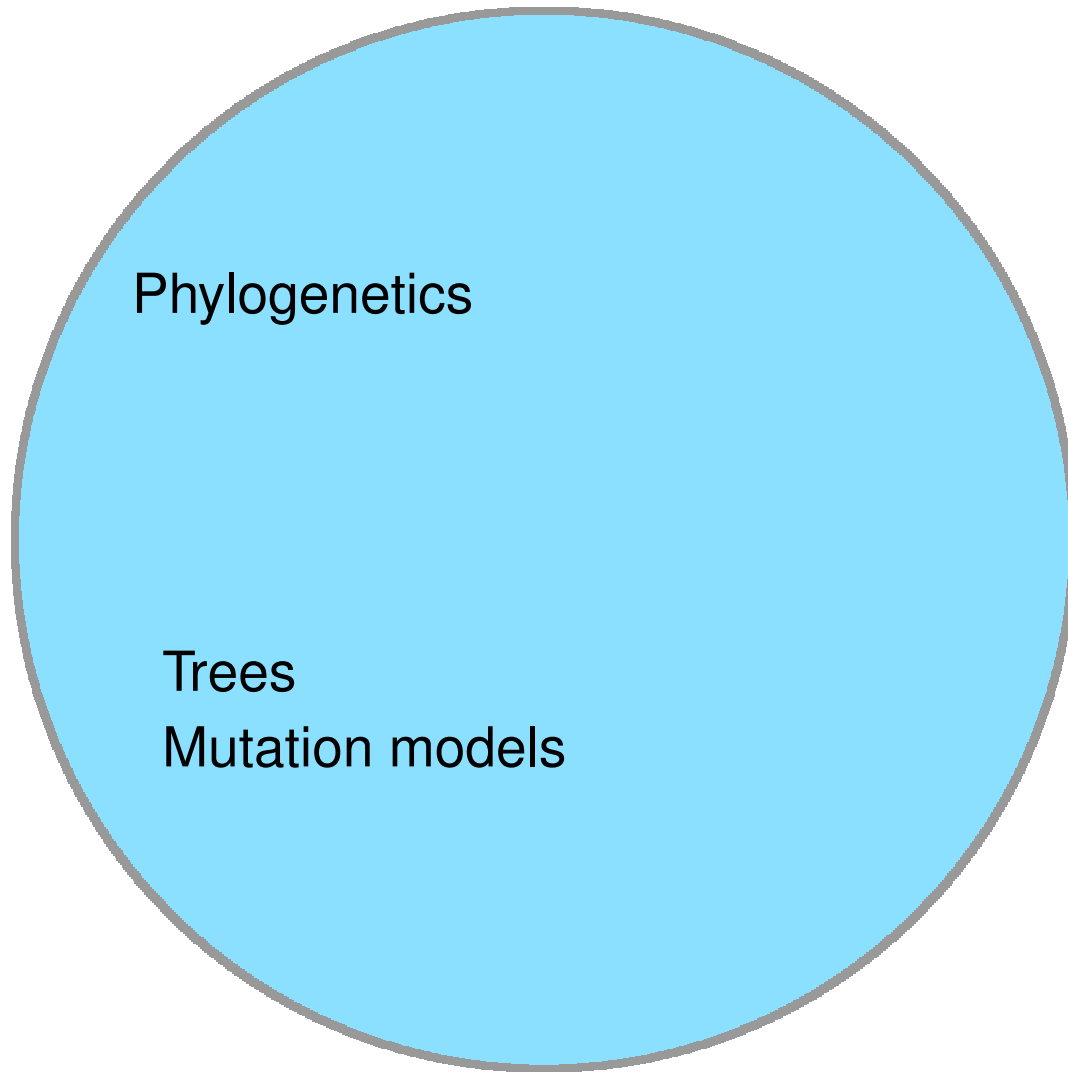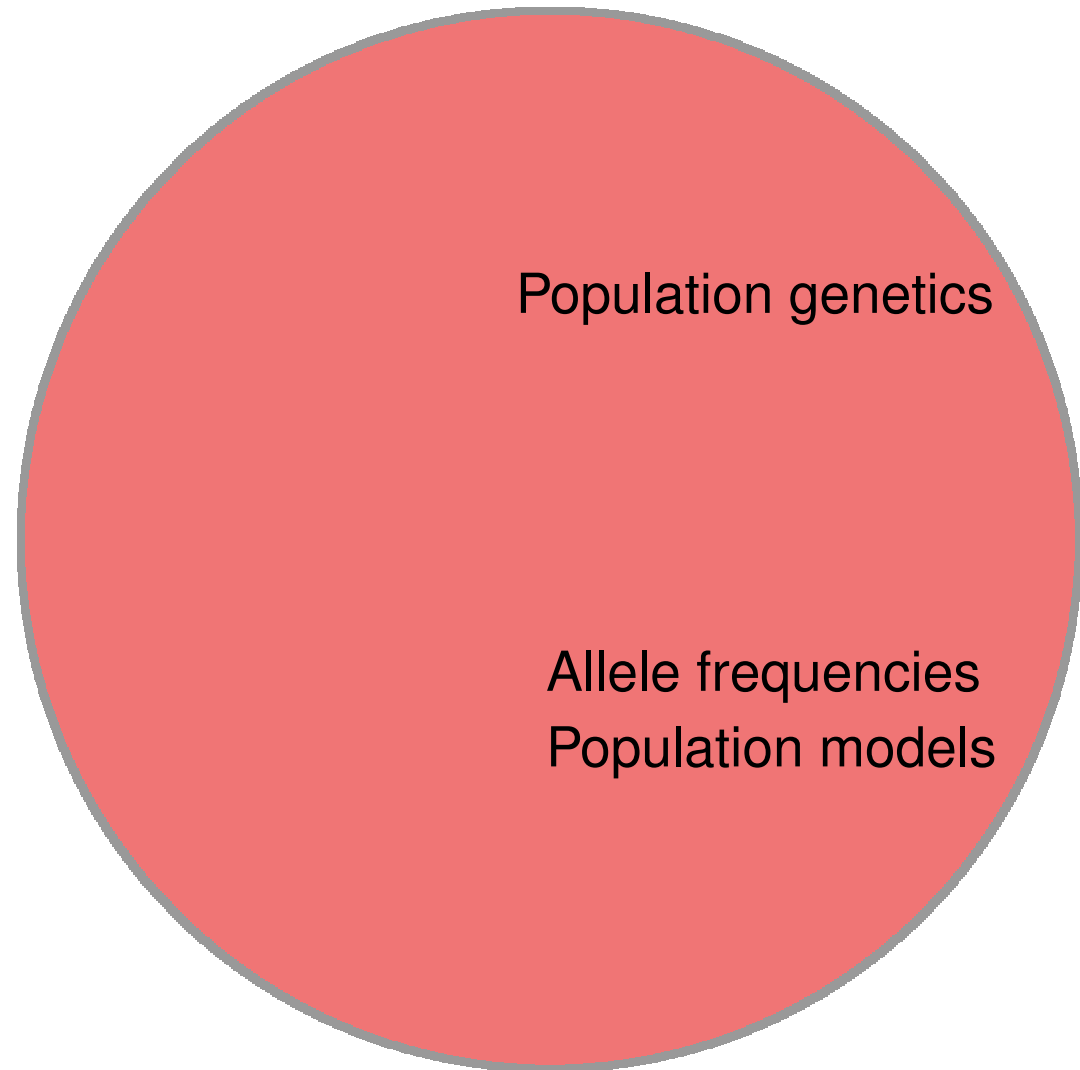
Peter Beerli
Florida State University
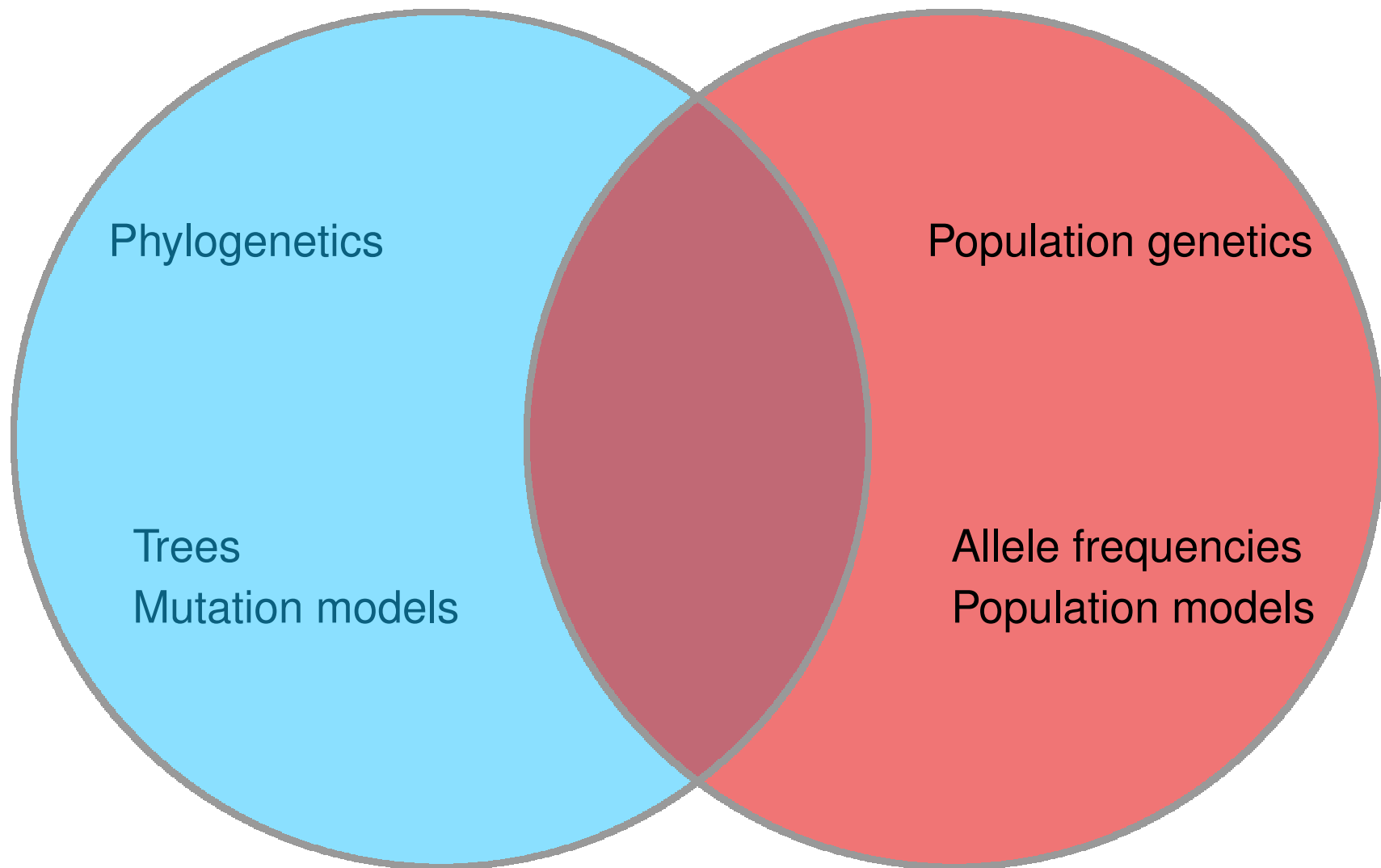
#MolEvol2018 MBL Woods Hole

Phylogenetics

Trees
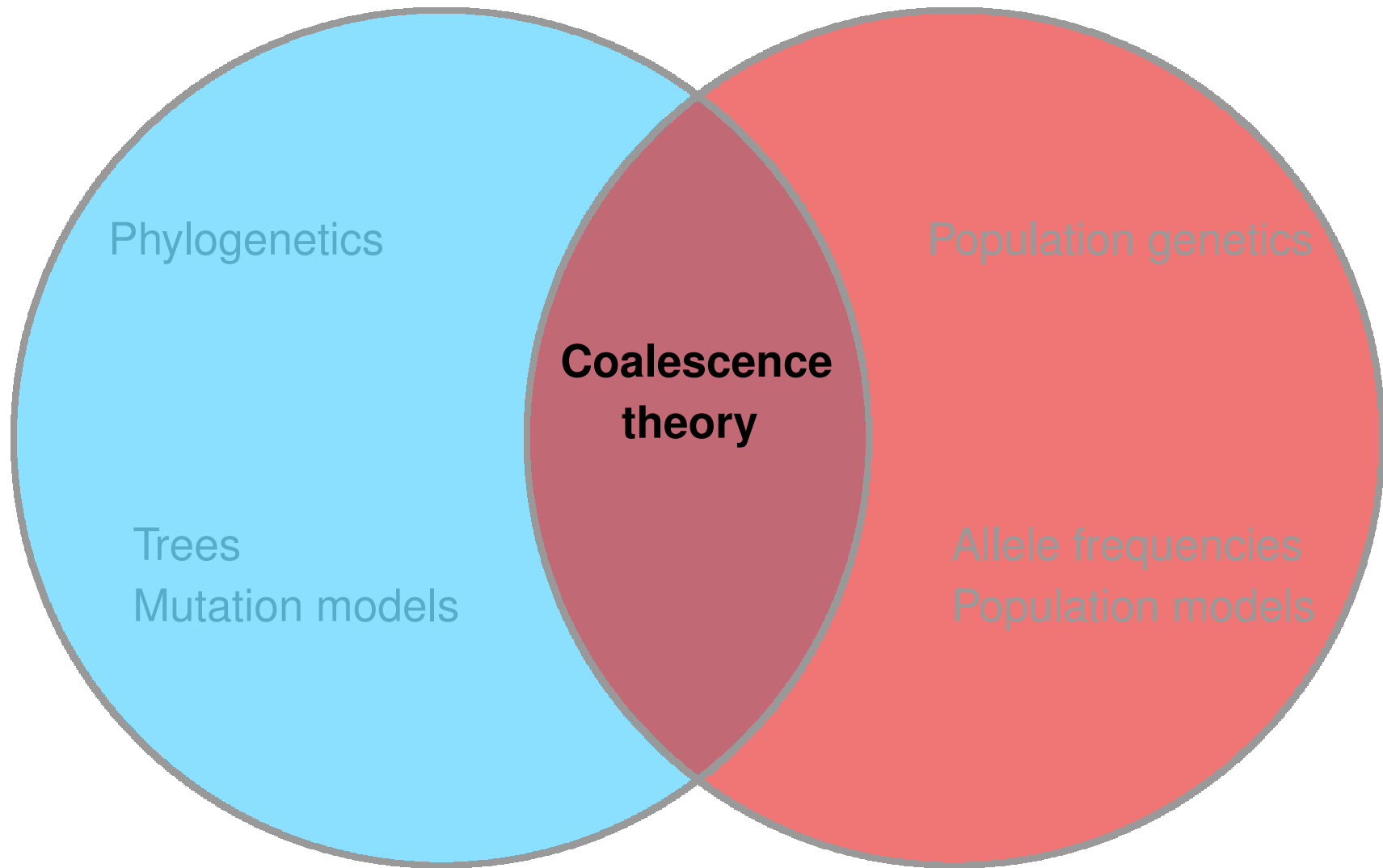Mutation models

Population genetics

Allele frequencies
Population models

# The big overview



Phylogenetics

Population genetics

**Coalescence theory**

Trees
Mutation models

Allele frequencies
Population models

## co•a•lesce |ˌkōəˈles|

verb [ intrans. ]

come together and form one mass or whole : *the puddles had* ***coalesced into*** *shallow streams* | *the separate details coalesce to form a single body of scientific thought.*
  • [ trans. ] combine (elements) in a mass or whole : *to help coalesce the community, they established an office.*
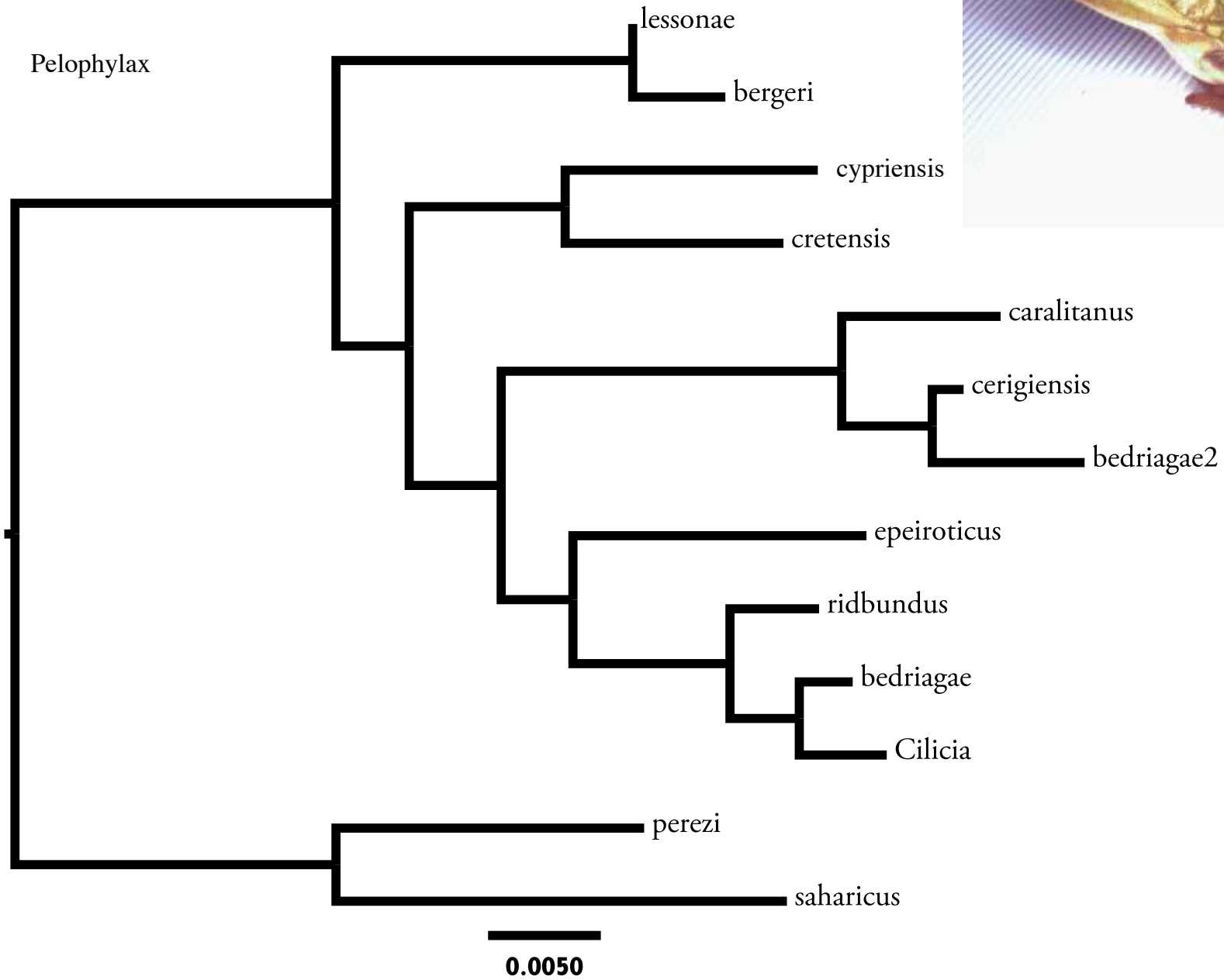
DERIVATIVES
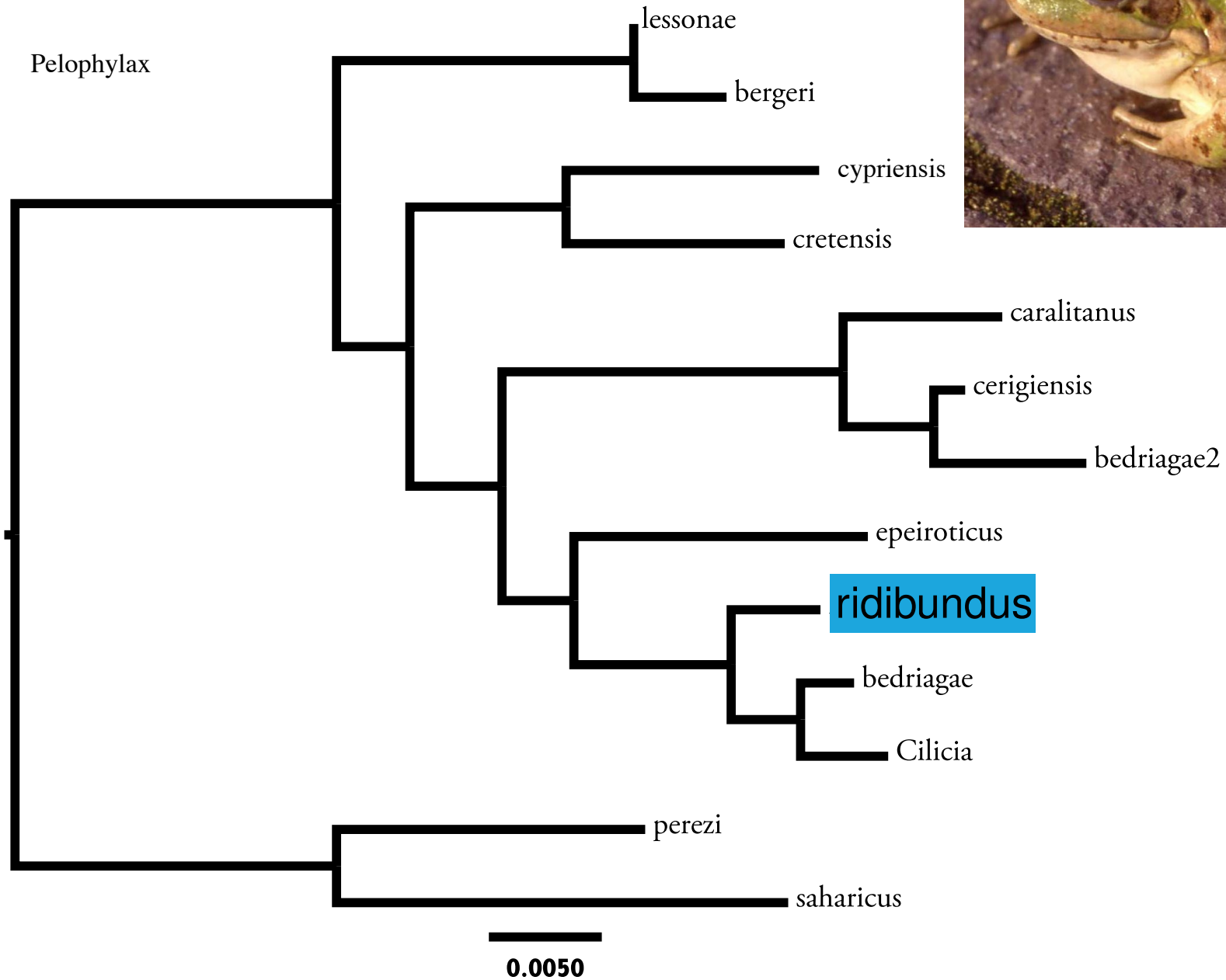
**co•a•les•cence** |-ˈlesəns| noun

**co•a•les•cent** |-ˈlesənt| adjective

ORIGIN mid 16th cent. (in the sense [bring together, unite] ): from Latin ***coalescere***, from ***co-*** (from ***cum 'with'*** ) + ***alescere 'grow up'*** (from ***alere 'nourish'*** ).
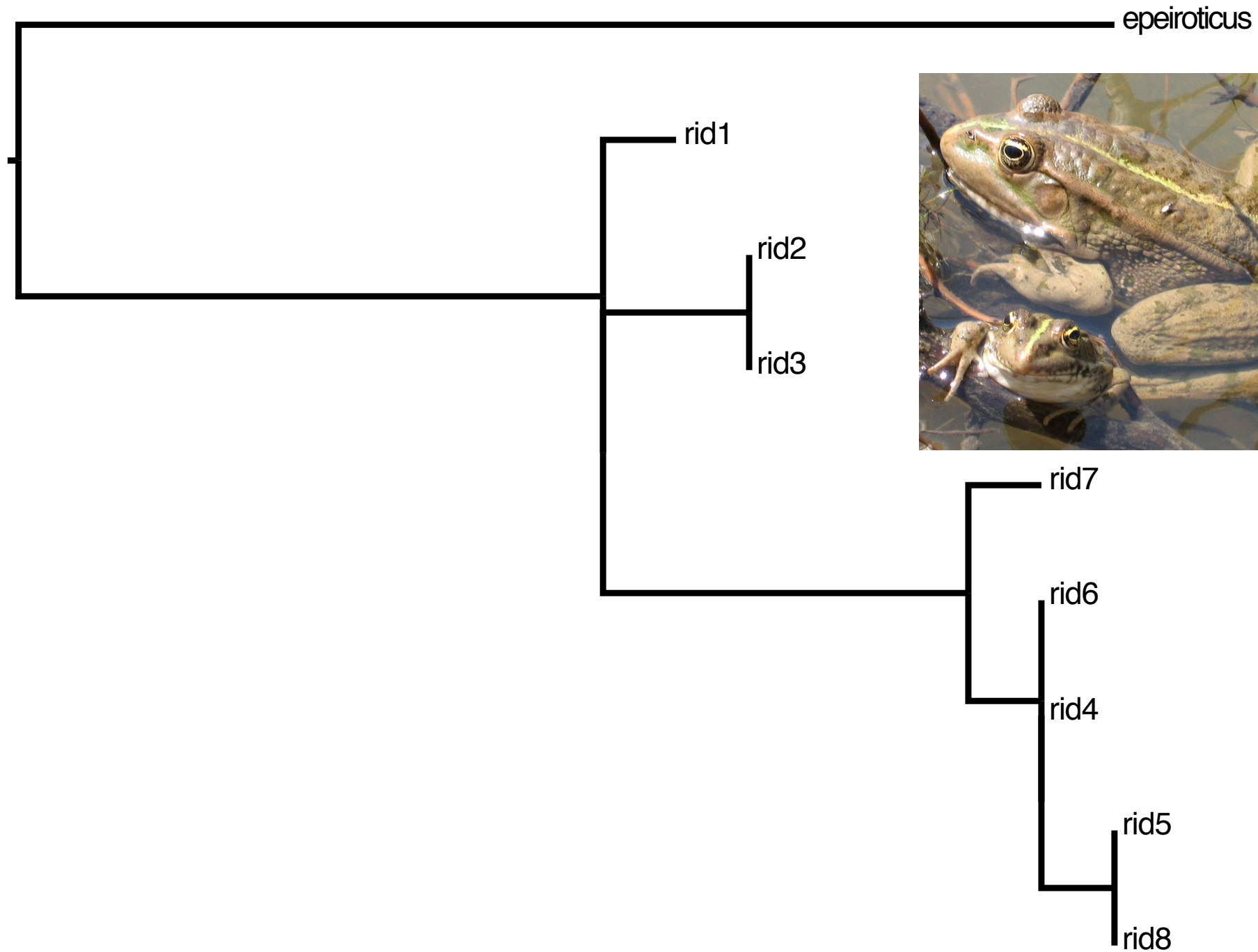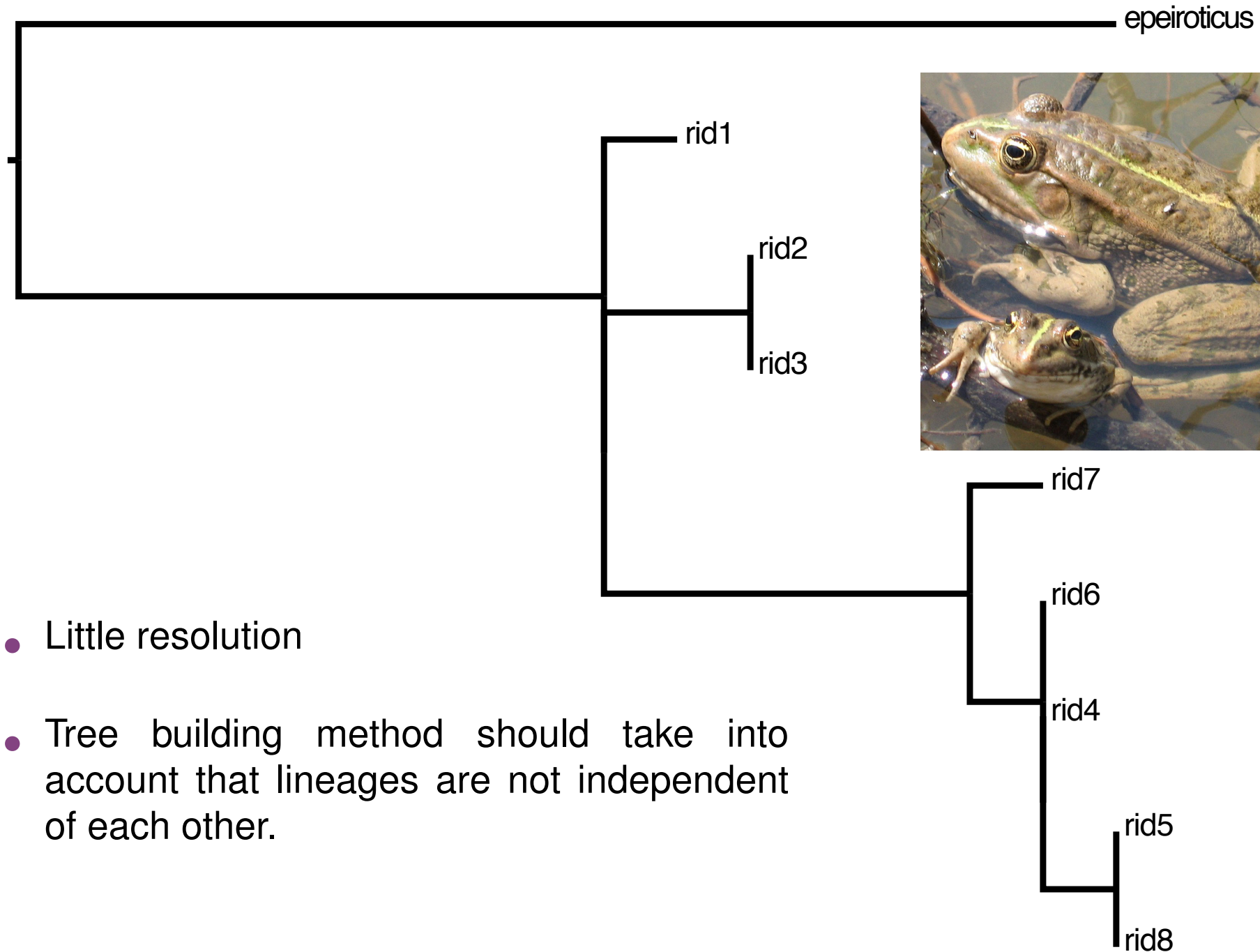
# Species trees



Pelophylax

```
                                    ┌─── lessonae
                              ┌─────┤
                              │     └─── bergeri
                        ┌─────┤
                        │     │     ┌─── cypriensis
                        │     └─────┤
                        │           └─── cretensis
                  ┌─────┤
                  │     │                 ┌─── caralitanus
                  │     │           ┌─────┤
                  │     │           │     │  ┌─── cerigiensis
                  │     │     ┌─────┤      └─┤
                  │     │     │     │        └─── bedriagae2
                  │     └─────┤     
                  │           │           ┌─── epeiroticus
                  │           └───────────┤
                  │                       │     ┌─── ridbundus
                  │                       └─────┤
                  │                             │  ┌─── bedriagae
                  │                              └─┤
                  │                                └─── Cilicia
                  │
                  │           ┌─── perezi
                  └───────────┤
                              └─── saharicus
```

0.0050

# Species trees



Pelophylax

- lessonae
- bergeri
- cypriensis
- cretensis
- caralitanus
- cerigiensis
- bedriagae2
- epeiroticus
- ridibundus
- bedriagae
- Cilicia
- perezi
- saharicus

0.0050

epeiroticus

rid1

rid2

rid3
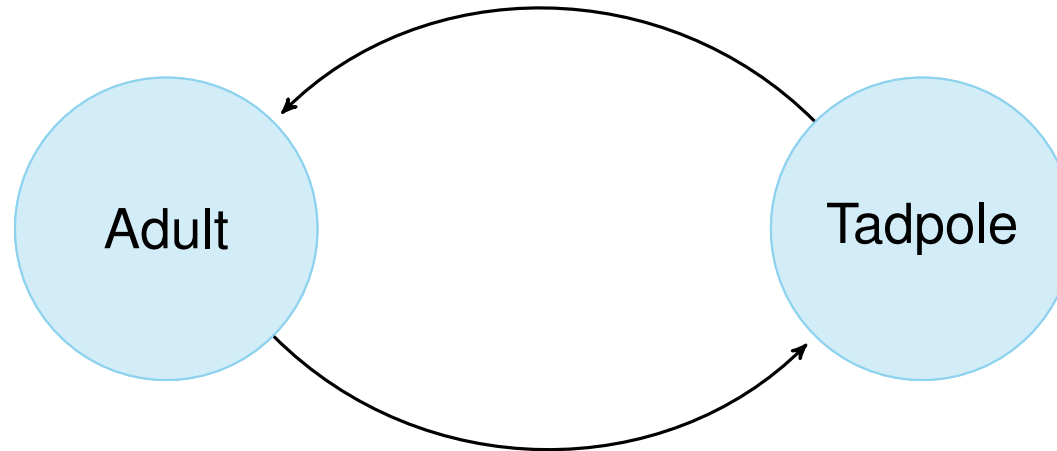
rid7

rid6

rid4

rid5

rid8

epeiroticus

rid1

rid2

rid3

rid7

rid6

rid4

rid5

rid8

- Little resolution

- Tree building method should take into account that lineages are not independent of each other.

Adult ⟷ Tadpole

Wright-Fisher population model

- All individuals live one generation and get replaced by their offspring

- All have same chance to reproduce, all are equally fit

- The number of individuals in the population is constant
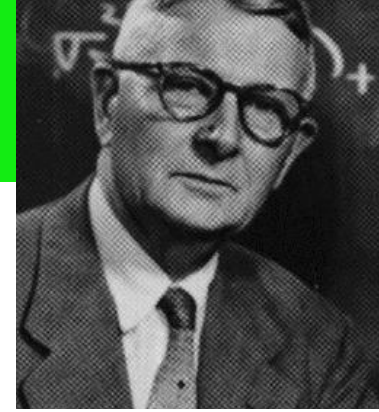
Past

Present

Past
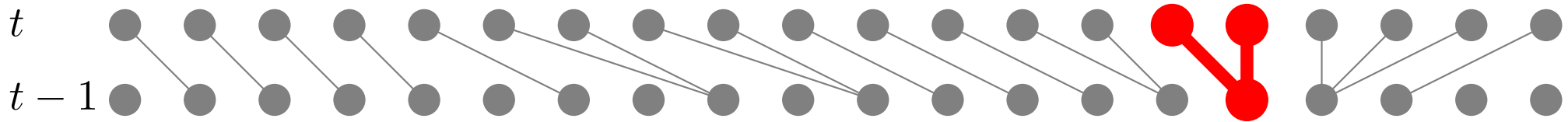
Present

Past

Present

Past

Present

Sewall Wright evaluated the probability that two randomly chosen individuals in generation $t$ have a common ancestor in generation $t - 1$. If we assume that there are $2N$ chromosomes then the probability of sharing a common ancestor in the last generation is
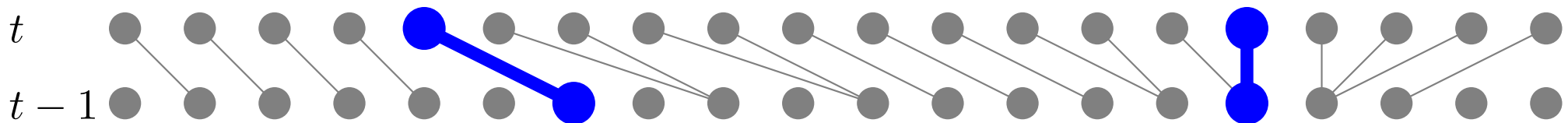
# Population model

Sewall Wright evaluated the probability that two randomly chosen individuals in generation $t$ have a common ancestor in generation $t - 1$. If we assume that there are $2N$ chromosomes then the probability of sharing a common ancestor in the last generation is

$$1.0$$

**Wright**

Sewall Wright evaluated the probability that two randomly chosen individuals in generation $t$ have a common ancestor in generation $t-1$. If we assume that there are $2N$ chromosomes then the probability of sharing a common ancestor in the last generation is

$$1.0 \times \frac{1}{2N}$$

$t$

$t-1$

Sewall Wright evaluated the probability that two randomly chosen individuals in generation $t$ have a common ancestor in generation $t-1$. If we assume that there are $2N$ chromosomes then the probability of sharing a common ancestor in last generation is

$$\frac{1}{2N}$$



The probability that two randomly picked chromosome do not have a common ancestor is

$$1 - \frac{1}{2N}$$

If we know the genealogy of the two individuals then we can calculate the probability as

$$\mathrm{P}(\tau|N) = \left(1 - \frac{1}{2N}\right)^{\tau} \left(\frac{1}{2N}\right)$$

where $\tau$ is the number of generations with no coalescence. This formula is the Geometric Distribution and we can calculate the expectation of the waiting time until two random individuals coalesce:

$$\mathbb{E}(\tau) = 2N$$

Present

Past

Present

Past

Present



Past

Present

Past

Present

Past

10000 random draw from a population with size $2N = 20$ leads to this distribution of times until two randomly chosen individuals have a common ancestor. The observed mean waiting time of 2N=20.34

# Observations

- For the time of coalescence in a sample of $\textcolor{red}{\text{TWO}}$ , we will wait on average $\textcolor{red}{2N}$ generations assuming it is a Wright-Fisher population

- The model assumes that the generations are discrete and non-overlapping

- Real populations do not necessarily behave like a Wright-Fisher (the *'ideal' population*)

- *We assume that calculation using Wright-Fisher populations can be extrapolated to real populations.*

# Other population models

| Wright-Fisher | Canning | Moran |
|---|---|---|



$$\sigma^2_{\text{offspring}} \simeq 1 \qquad\qquad \sigma^2_{\text{offspring}} = x \qquad\qquad \sigma^2_{\text{offspring}} = \frac{2}{2N}$$

$$\mathbb{E}(\tau) = 2N \qquad\qquad \mathbb{E}(\tau) = 2N/x \qquad\qquad \mathbb{E}(\tau) = \frac{1}{2}(2N)^2$$

generation time $g = 1$ $\qquad\qquad g = 1 \qquad\qquad g = 2N$

You can generate graphs like this using the python program *popsim* (check out my faculty page for the link)

Past

Present

Past

Present

Past

Present

Past

Present

Past

Present

# Samples larger than two

Sir J. F. C. Kingman described in 1982 the $n$-coalecent. He showed the behavior of a sample of size $n$, and its probability structure looking backwards in time.
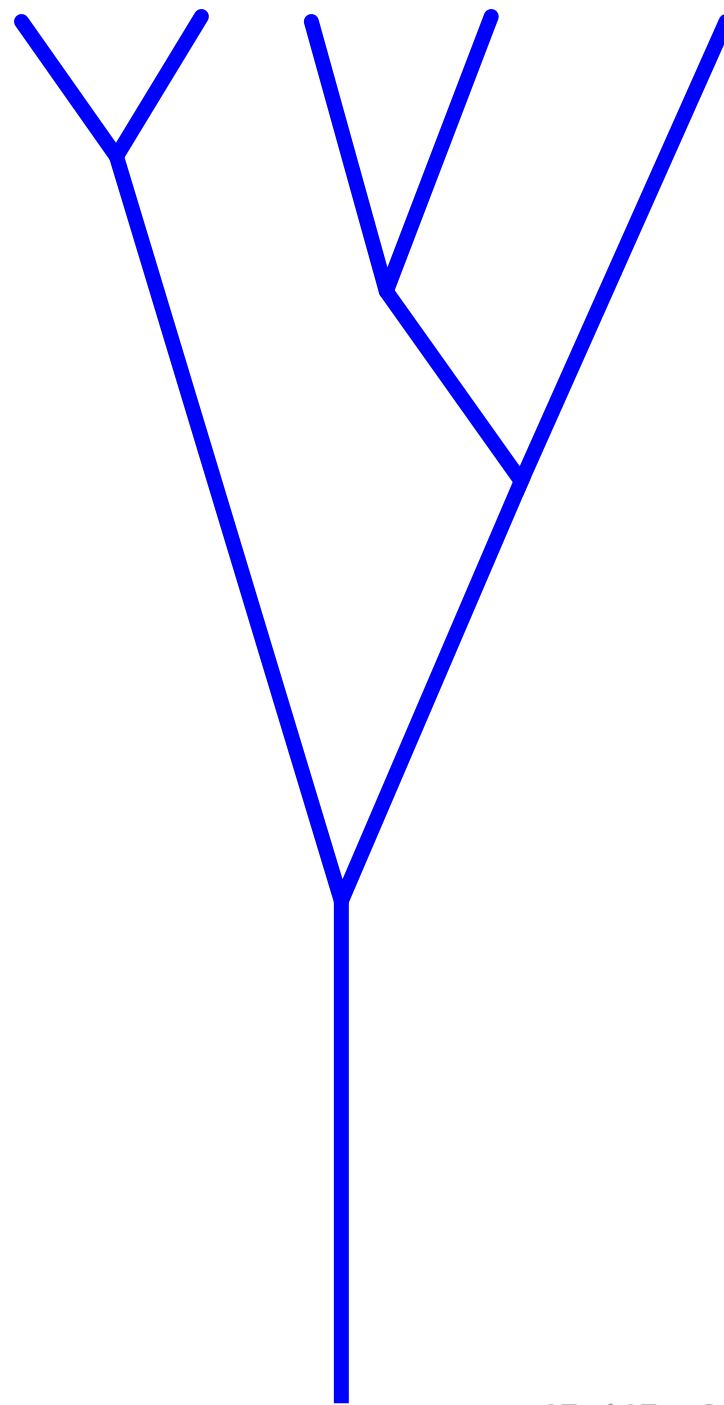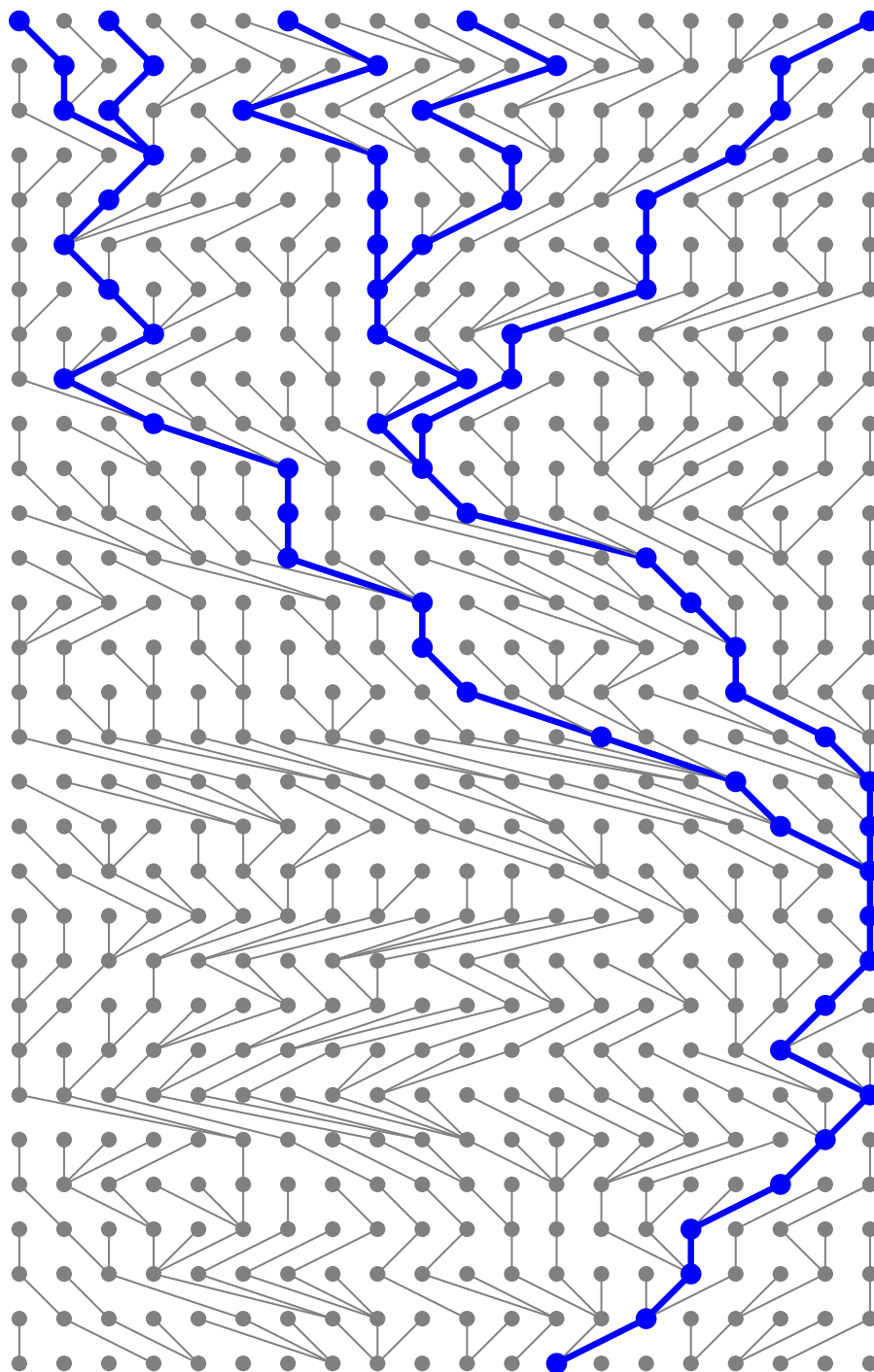
General findings:

$$\text{coalescence rate} = \binom{n}{2} = \frac{n(n-1)}{2}$$

Once a coalescence happened $n$ is reduced to $n-1$ because two lineages merged into one. He then imposed a continuous approximation of the Canning's exchangeable model to get results.
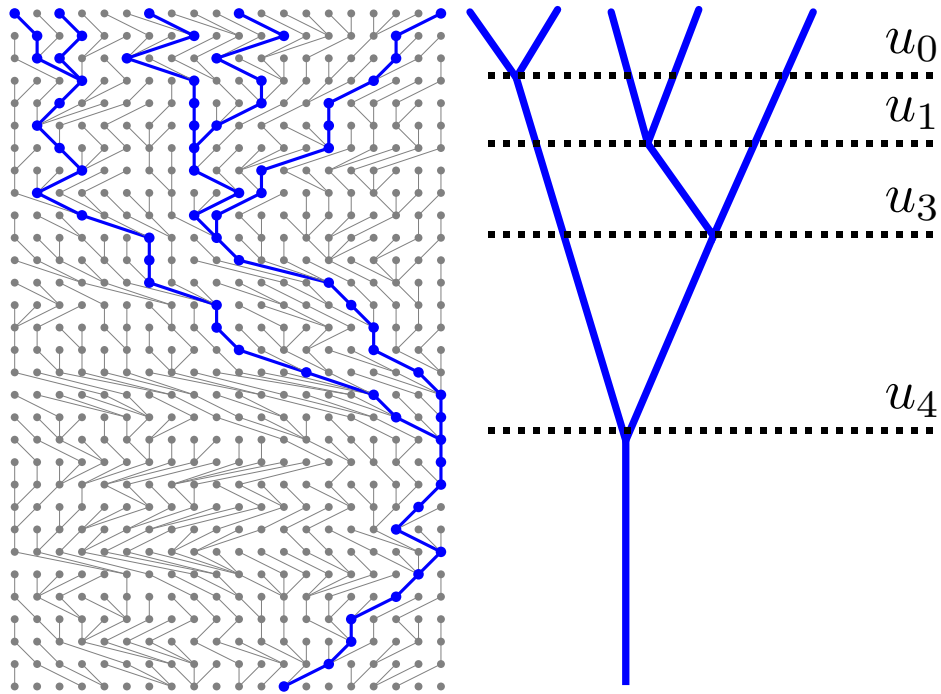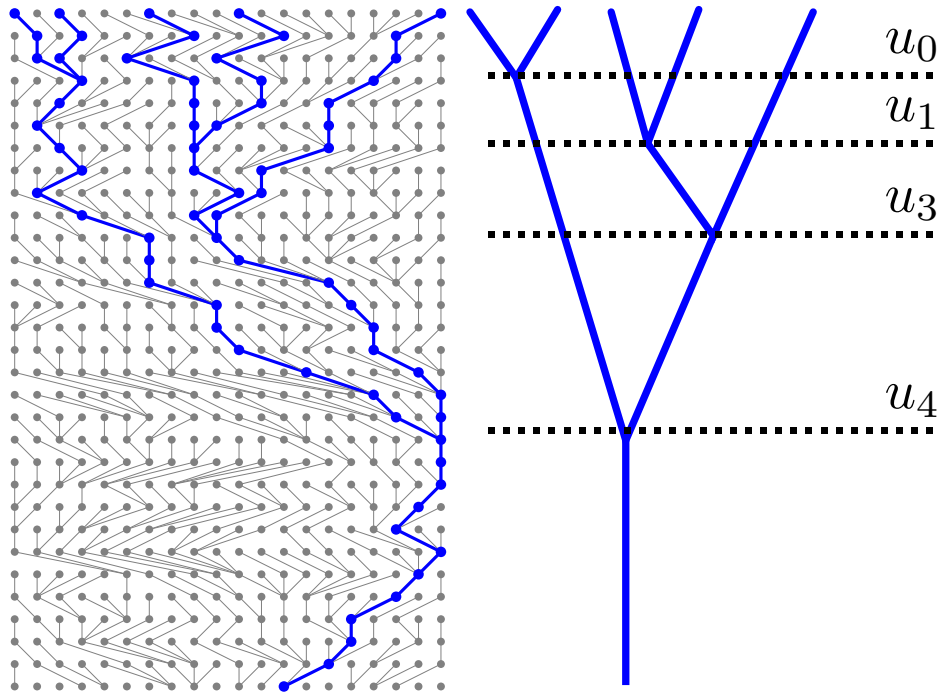
$u_0$

$u_1$

$u_3$

$u_4$

Looking backward in time, the first coalescence between two random individuals is the result of a waiting process that depends on the sample $n$ and the total population size $N$.
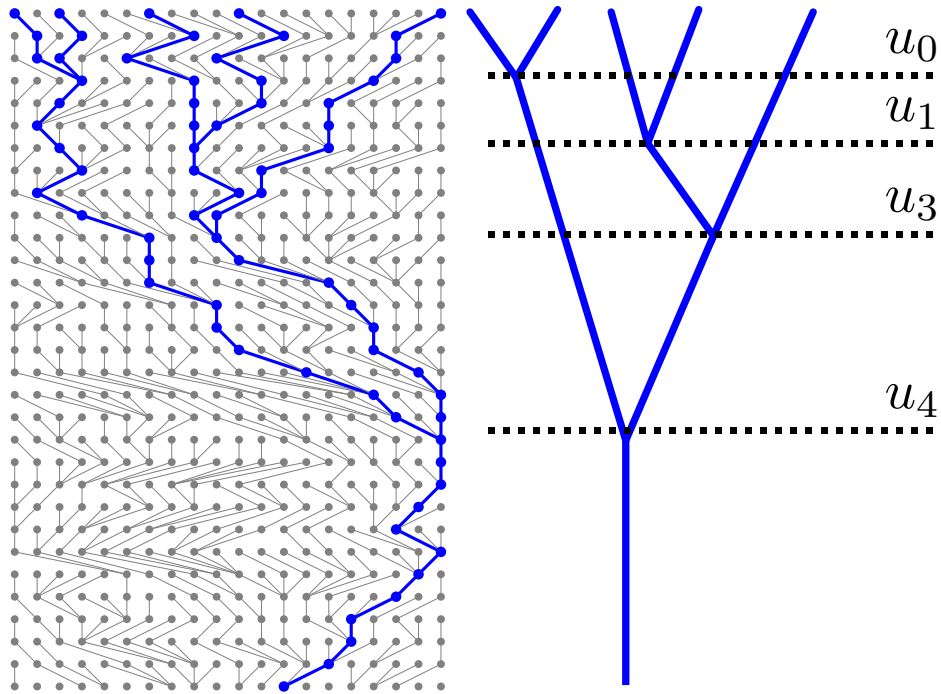
# Samples larger than two



$u_0$
$u_1$
$u_3$
$u_4$

Looking backward in time, the first coalescence between two random individuals is the result of a waiting process that depends on the sample $n$ and the total population size $N$.
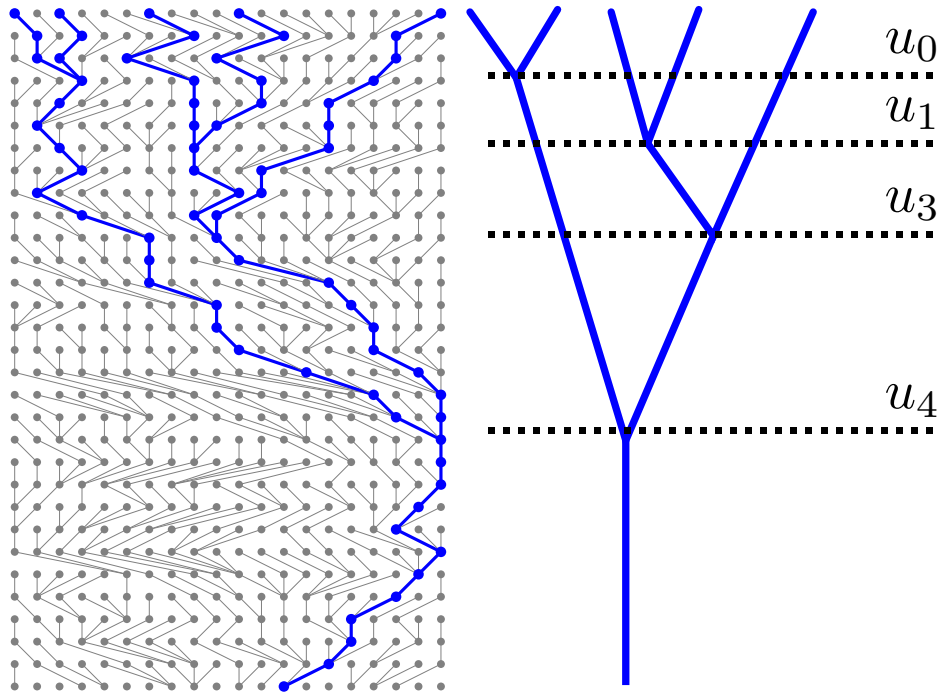
Using Kingman's coalescence rate and imposing a time scale we can approximate the process with an exponential distribution:

# Samples larger than two
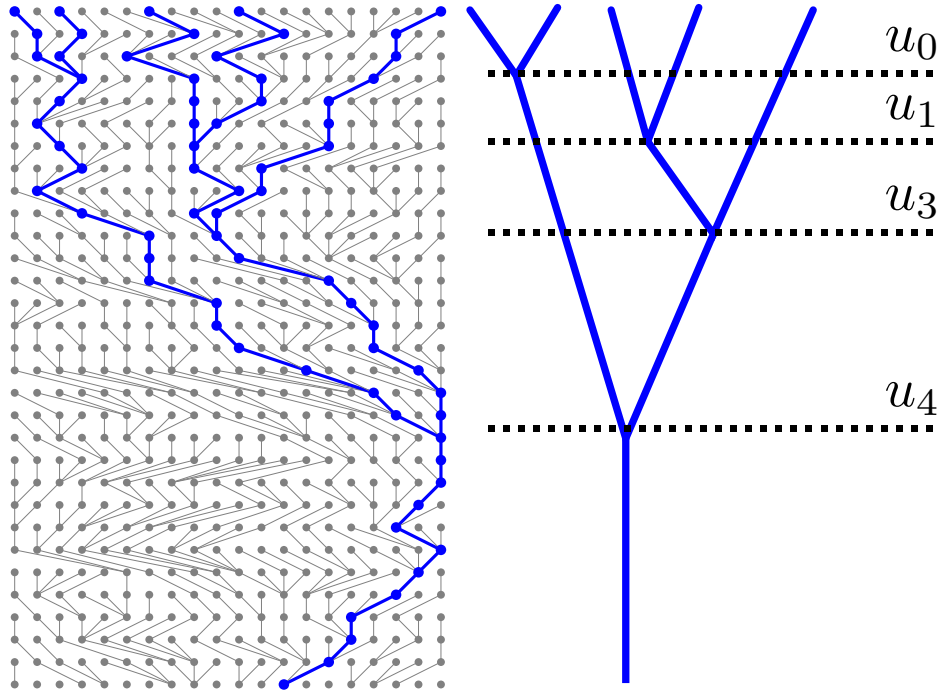


$u_0$
$u_1$
$u_3$
$u_4$

Looking backward in time, the first coalescence between two random individuals is the result of a waiting process that depends on the sample $n$ and the total population size $N$.

Using Kingman's coalescence rate and imposing a time scale we can approximate the process with an exponential distribution:
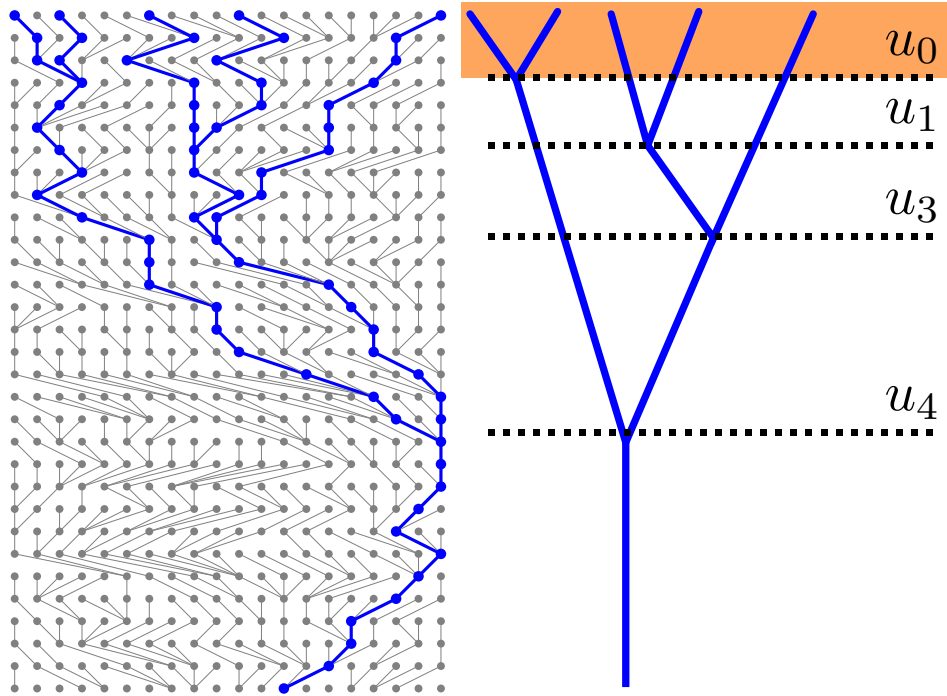
$$\mathrm{P}(u_j|N) = e^{-u_j\lambda}\lambda$$

with the scaled coalescence rate

$$\lambda = \binom{k}{2}\frac{1}{2N} \times \mathrm{Prob}(\text{others do not coalesce})$$

# Samples larger than two



$u_0$
$u_1$
$u_3$
$u_4$

Looking backward in time, the first coalescence between two random individuals is the result of a waiting process that depends on the sample $n$ and the total population size $N$.

Using Kingman's coalescence rate and imposing a time scale we can approximate the process with a exponential distribution:

$$\mathrm{P}(u_j|N) = e^{-u_j\lambda}\lambda$$

with the scaled coalescence rate

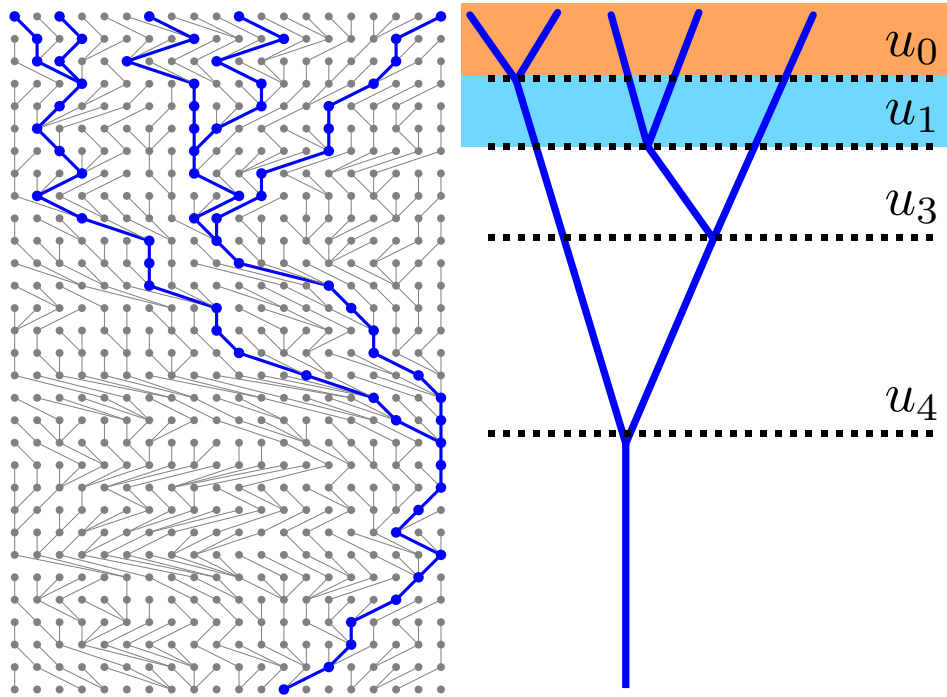$$\lambda = \binom{k}{2}\frac{1}{2N} = \frac{k(k-1)}{2(2N)} = \frac{k(k-1)}{4N}$$
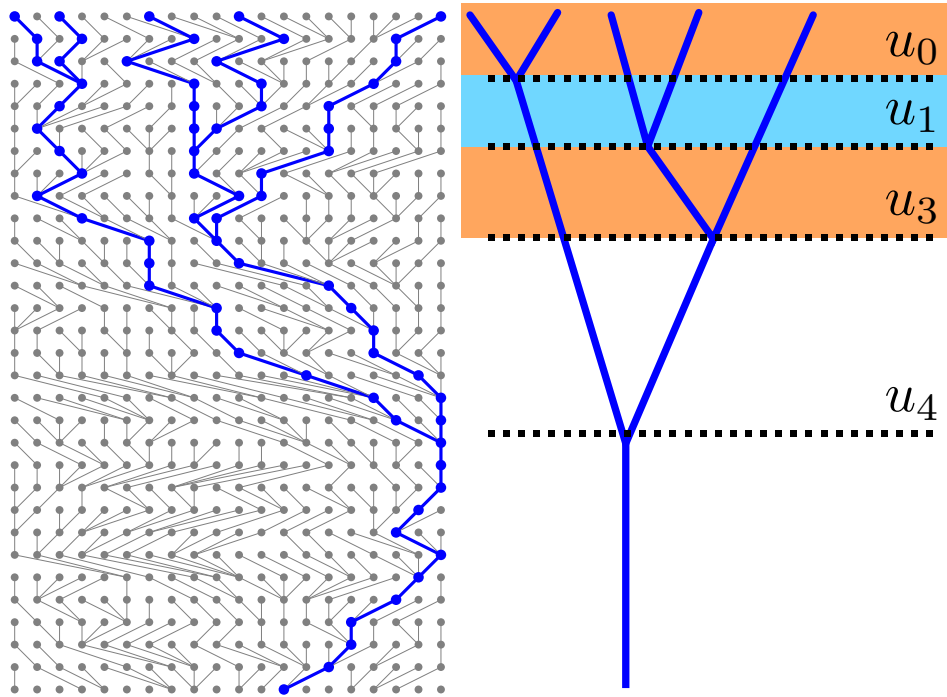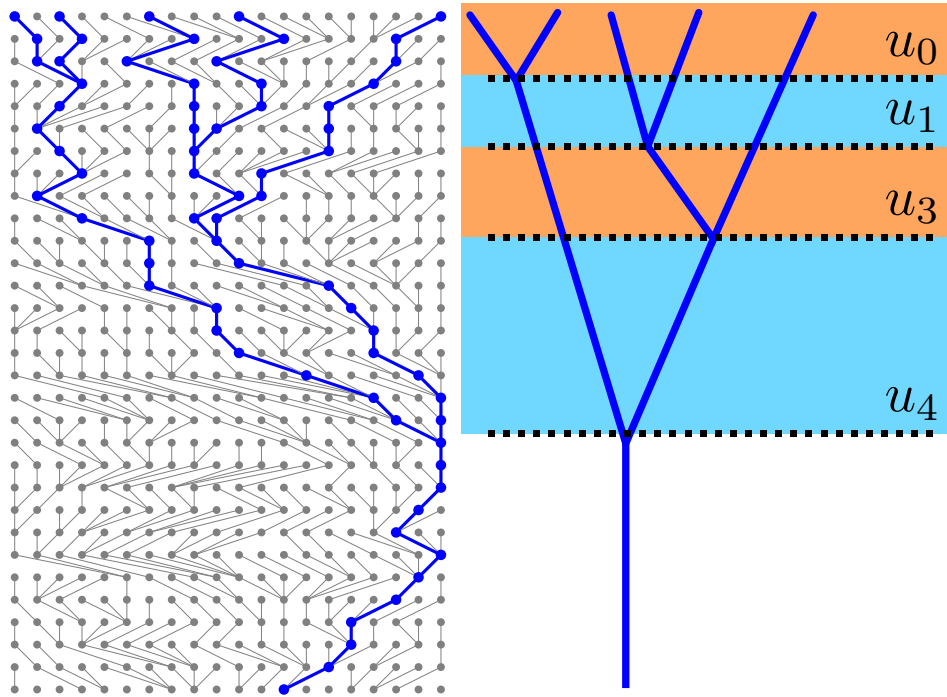
$u_0$

$u_1$

$u_3$

$u_4$

We are now able to calculate the probability of a whole relationship tree (Genealogy $G$). We assume that each coalescence is independent from any other:

$$\mathrm{P}(G|N)$$
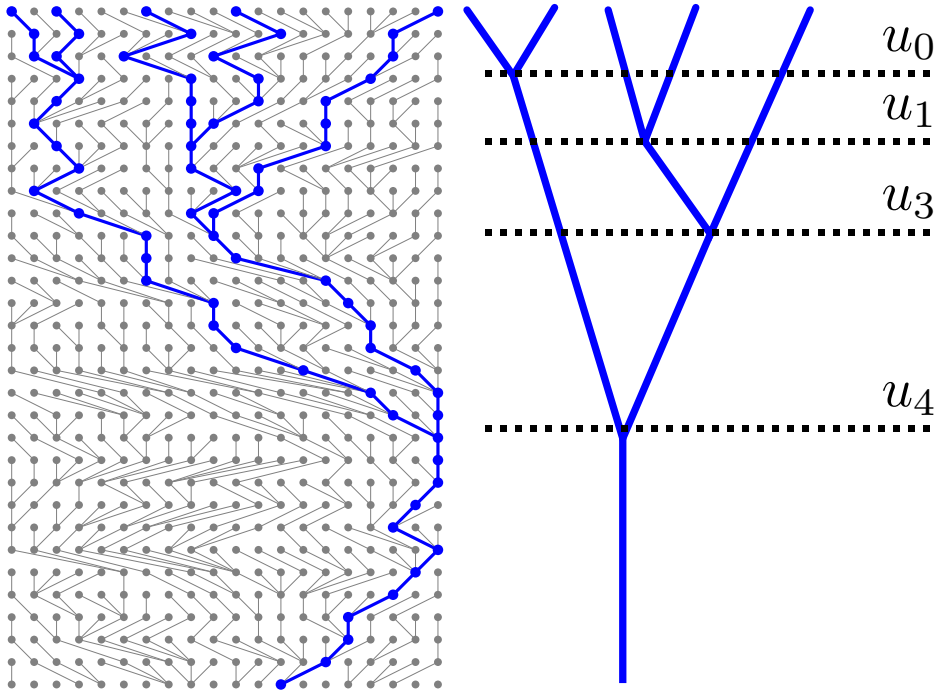
We are now able to calculate the probability of a whole relationship tree (Genealogy $G$). We assume that each coalescence is independent from any other:

$$\mathrm{P}(G|N) = \mathrm{P}(u_0|N, i_1, i_2)$$

$$\times$$

We are now able to calculate the probability of a whole relationship tree (Genealogy $G$). We assume that each coalescence is independent from any other:
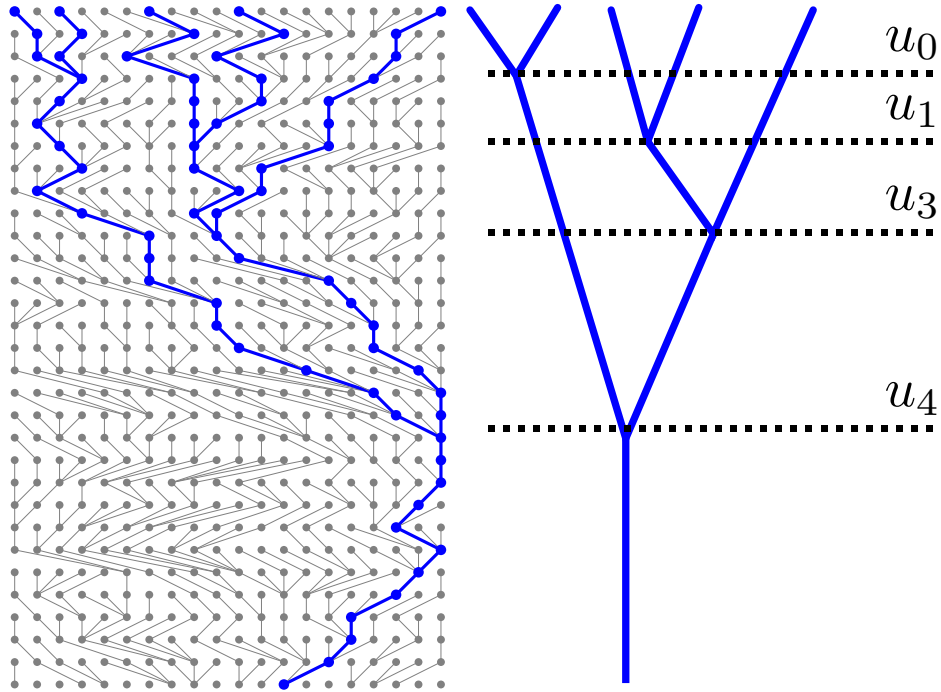
$$\mathrm{P}(G|N) = \mathrm{P}(u_0|N, i_1, i_2)$$
$$\times \mathrm{P}(u_1|N, i_3, i_4)$$

We are now able to calculate the probability of a whole relationship tree (Genealogy $G$). We assume that each coalescence is independent from any other:

$$\mathrm{P}(G|N) = \mathrm{P}(u_0|N, i_1, i_2)$$

$$\times \mathrm{P}(u_1|N, i_3, i_4)$$

$$\times \mathrm{P}(u_3|N, i_{3,4}, i_5)$$

We are now able to calculate the probability of a whole relationship tree (Genealogy $G$). We assume that each coalescence is independent from any other:

$$\mathrm{P}(G|N) = \mathrm{P}(u_0|N, i_1, i_2)$$
$$\times \mathrm{P}(u_1|N, i_3, i_4)$$
$$\times \mathrm{P}(u_3|N, i_{3,4}, i_5)$$
$$\times \mathrm{P}(u_4|N, i_{1,2}, i_{3,4,5})$$

We are now able to calculate the probability of a whole relationship tree (Genealogy $G$). We assume that each coalescence is independent from any other:

$$\mathrm{P}(G|N) = \quad \mathrm{P}(u_0|N, i_1, i_2)$$

$$\times \mathrm{P}(u_1|N, i_3, i_4)$$

$$\times \mathrm{P}(u_3|N, i_{3,4}, i_5)$$

$$\times \mathrm{P}(u_4|N, i_{1,2}, i_{3,4,5})$$

$$\mathrm{P}(G|N) = \prod_{j=0}^{T} e^{-u_j \frac{k_j(k_j-1)}{4N}} \frac{2}{4N}$$

Each interval $u_j$ is independent of the others, the expected length of the interval is the inverse of the coalescent rate. Thus we can sum these expectations to get to expectation of the depth of the genealogy.

$$\mathbb{E}(\tau_{\text{MRCA}}) = \text{Sum of the expectation of each time interval} = \sum_{j=0}^{J} \frac{4N}{k_j(k_j - 1)}$$

$$\lim_{k \to \infty} \mathbb{E}(\tau_{\text{MRCA}}) = 2N + \frac{2}{3}N + \frac{1}{3}N + \frac{1}{5}N + \frac{2}{15}N + ... = 4N \qquad \lim_{k \to \infty} \sigma(\tau_{\text{MRCA}}) = 4N$$

If we know the genealogy $G$ with certainty then we can calculate the population size $N$. Finding the maximum probability $\mathrm{P}(G|N,k)$ is simple, we evaluate all possible values for $N$ and pick the value with the highest probability.
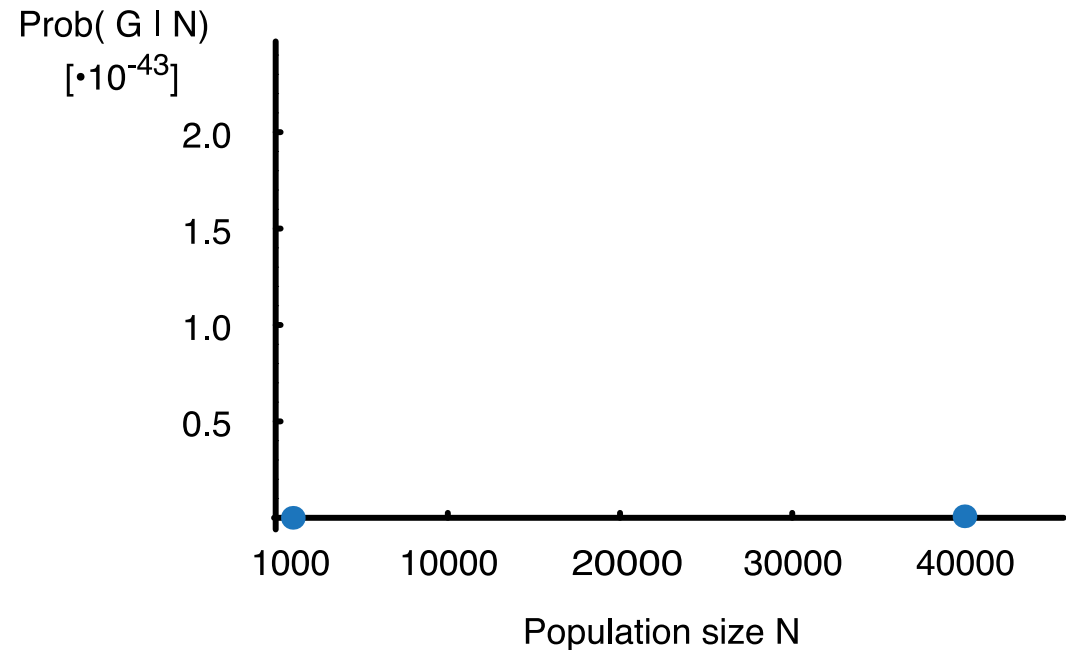
If we know the genealogy $G$ with certainty then we can calculate the population size $N$. Finding the maximum probability $\mathrm{P}(G|N,k)$ is simple, we evaluate all possible values for $N$ and pick the value with the highest probability.

If we know the genealogy $G$ with certainty then we can calculate the population size $N$. Finding the maximum probability $\mathrm{P}(G|N, k)$ is simple, we evaluate all possible values for $N$ and pick the value with the highest probability.

Prob( G | N)
[$\cdot 10^{-43}$]

2.0

1.5

1.0

0.5

1000    10000    20000    30000    40000
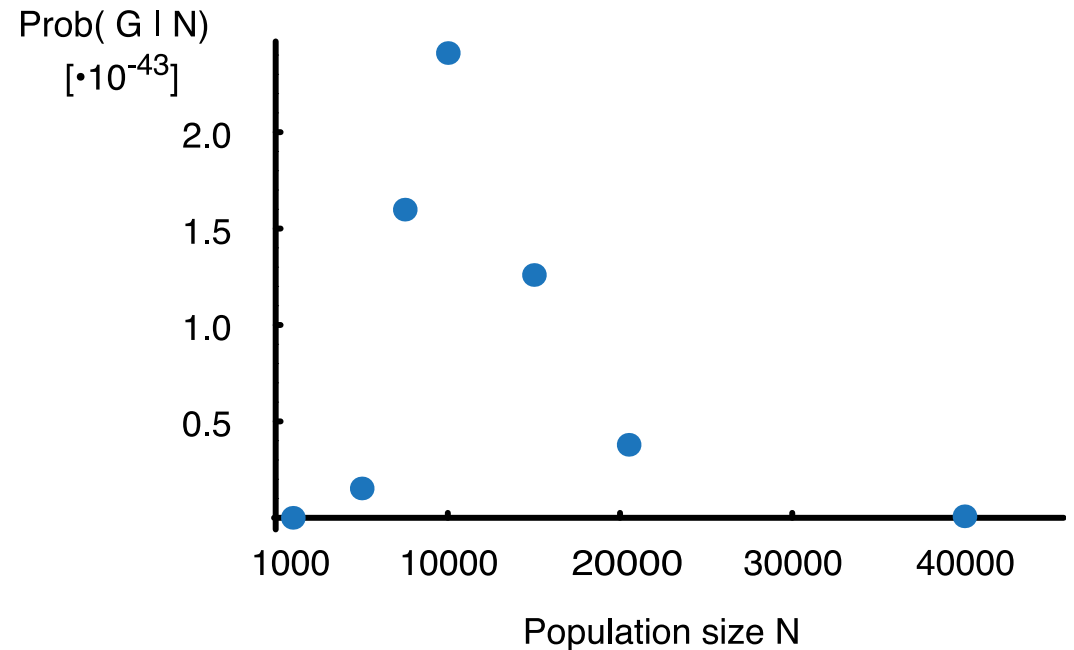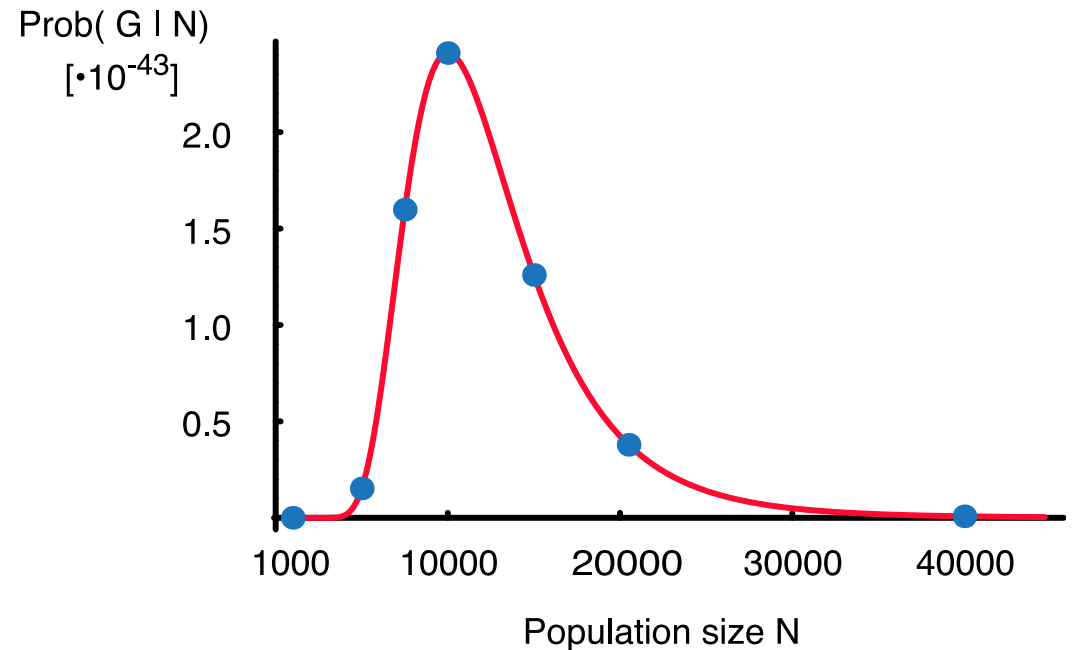
Population size N

If an oracle gives us the true relationship tree $G$ then we can calculate the population size $N$.

$$p(G|N, n) = \prod_{k=2}^{n} \exp\left(-u_k \frac{k(k-1)}{4N}\right) \frac{2}{4N}$$

Prob( G | N)
$[\cdot 10^{-43}]$

2.0

1.5

1.0

0.5

1000   10000   20000   30000   40000
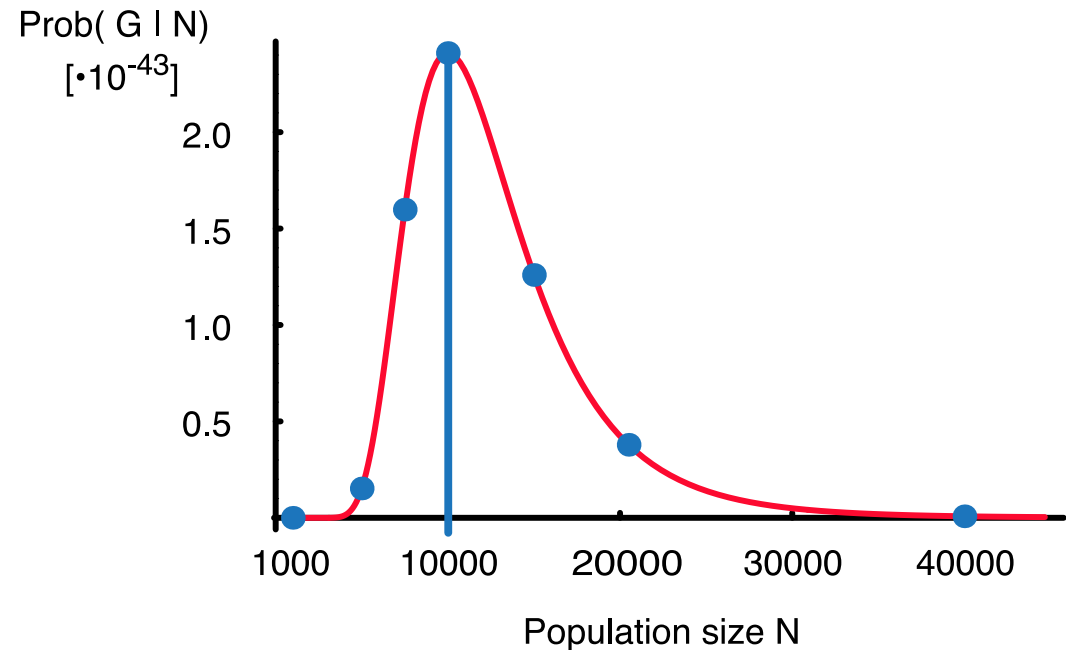
Population size N

If an oracle gives us the true relationship tree $G$ then we can calculate the population size $N$.

$$p(G|N, n) = \prod_{k=2}^{n} \exp\left(-u_k \frac{k(k-1)}{4N}\right) \frac{2}{4N}$$

Prob( G | N)
$[\cdot 10^{-43}]$

2.0

1.5

1.0

0.5
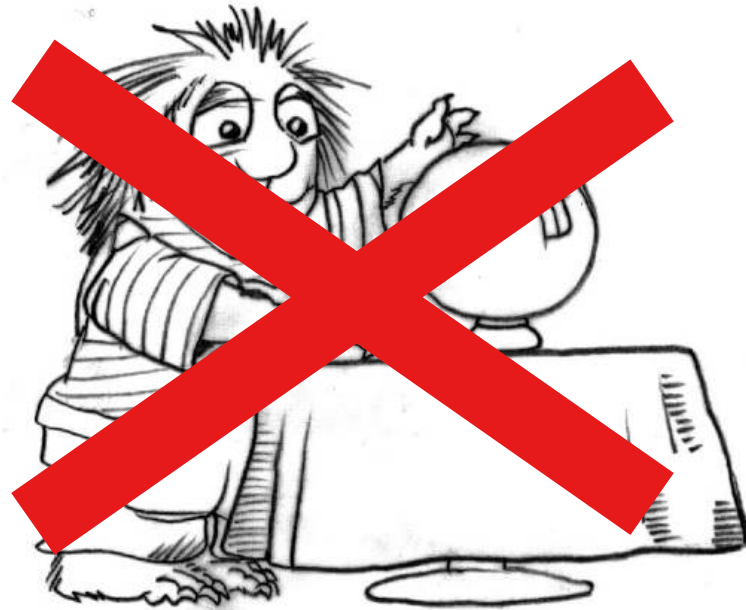
1000  10000  20000  30000  40000

Population size N

If an oracle gives us the true relationship tree $G$ then we can calculate the population size $N$.

$$p(G|N, n) = \prod_{k=2}^{n} \exp\left(-u_k \frac{k(k-1)}{4N}\right) \frac{2}{4N}$$

Prob( G | N)
$[\cdot 10^{-43}]$

2.0

1.5

1.0

0.5

1000   10000   20000   30000   40000

Population size N

If an oracle gives us the true relationship tree $G$ then we can calculate the population size $N$.

$$p(G|N, n) = \prod_{k=2}^{n} \exp\left(-u_k \frac{k(k-1)}{4N}\right) \frac{2}{4N}$$

Prob( G | N) [•10⁻⁴³] vs Population size N

If an oracle gives us the true relationship tree $G$ then we can calculate the population size $N$.

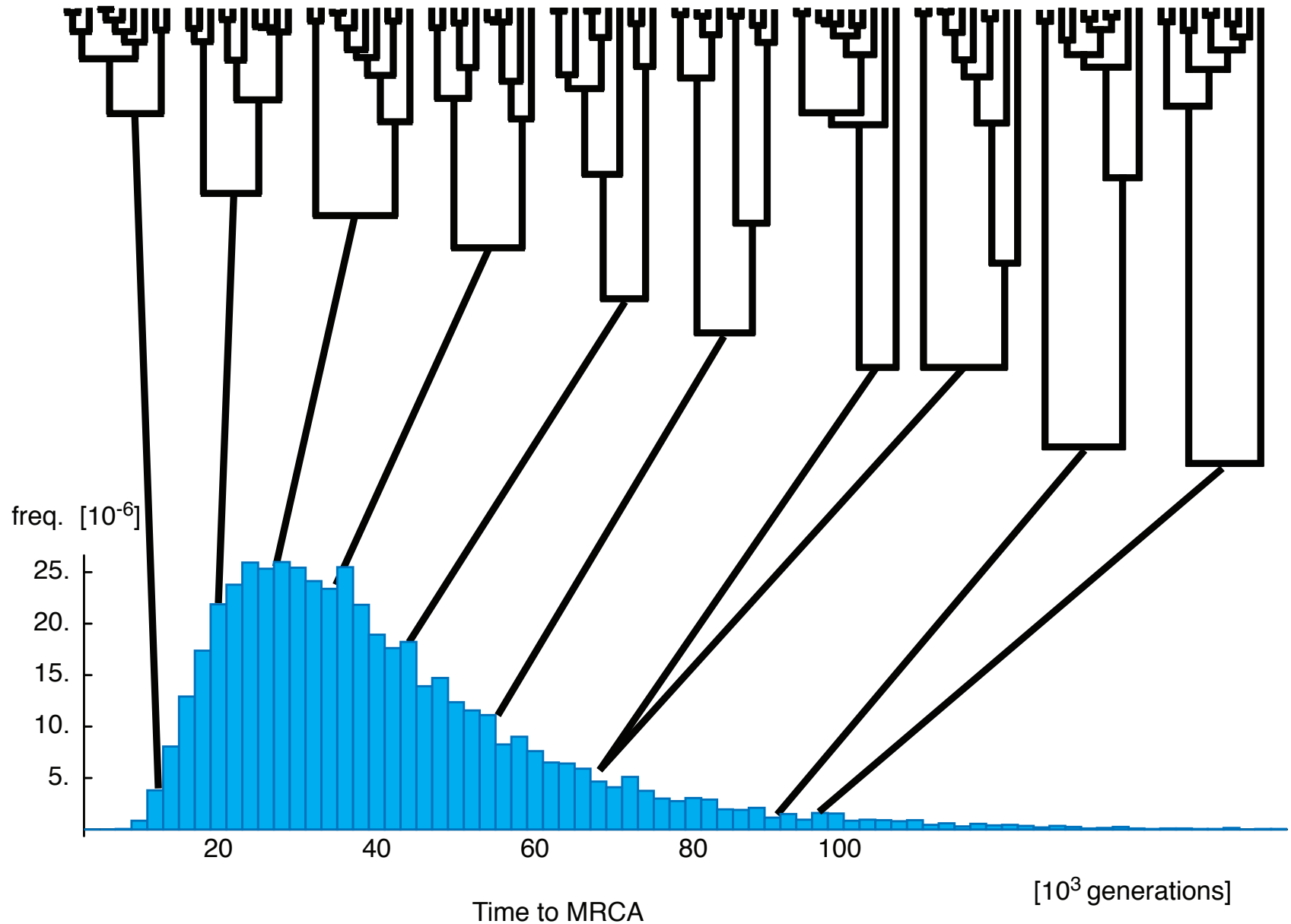$$p(G|N, n) = \prod_{k=2}^{n} \exp\left(-u_k \frac{k(k-1)}{4N}\right) \frac{2}{4N}$$

Prob( G | N) $[\cdot 10^{-43}]$

If an oracle gives us the true relationship tree $G$ then we can calculate the population size $N$.

$$p(G|N, n) = \prod_{k=2}^{n} \exp\left(-u_k \frac{k(k-1)}{4N}\right) \frac{2}{4N}$$

Prob( G I N) [$\cdot 10^{-43}$]

If an oracle gives us the true relationship tree $G$ then we can calculate the population size $N$.

$$p(G|N, n) = \prod_{k=2}^{n} \exp\left(-u_k \frac{k(k-1)}{4N}\right) \frac{2}{4N}$$

Prob( G I N)
[$\cdot 10^{-43}$]

2.0

1.5

1.0

0.5

1000  10000  20000  30000  40000

Population size N

If an oracle gives us the true relationship tree $G$ then we can calculate the population size $N$.

$$p(G|N, n) = \prod_{k=2}^{n} \exp\left(-u_k \frac{k(k-1)}{4N}\right) \frac{2}{4N}$$

Prob( G | N) [$\cdot 10^{-43}$]

If an oracle gives us the true relationship tree $G$ then we can calculate the population size $N$.

$$p(G|N, n) = \prod_{k=2}^{n} \exp\left(-u_k \frac{k(k-1)}{4N}\right) \frac{2}{4N}$$

Prob( G | N)
$[\cdot 10^{-43}]$

2.0

1.5

1.0

0.5

1000   10000   20000   30000   40000

Population size N

If an oracle gives us the true relationship tree $G$ then we can calculate the population size $N$.

$$p(G|N, n) = \prod_{k=2}^{n} \exp\left(-u_k \frac{k(k-1)}{4N}\right) \frac{2}{4N}$$

There are at least two problems with the oracle-approach:

- There is no oracle to gives us clear information!

- We do not record genealogies, our data are sequences, microsatellite loci!

- What about the variability of the coalescence process?

All genealogies were simulated with the same population size $N_e = 10,000$

- freq. $[10^{-6}]$
- Time to MRCA
- $[10^3$ generations$]$

MRCA = most recent common ancestor (last node in the genealogy)

- All individuals have the same fitness (no selection).

- All individuals have the same chance to be in the sample (random sampling).

- The coalescent allows only merging two lineages per generation. This restricts us to to have a much smaller sample size than the population size.
$$n << N$$

- Yun-Xin Fu (2005) described the exact coalescent for the Wright-Fisher model and derived a maximal sample size $n < \sqrt{4N}$ for a diploid population. Although this may look like a severe restriction for the use of the coalescence in small populations, it turned out that the coalescence is rather robust and that even sample sizes close to the effective population size are not biasing immensely.

- Large samples coalesce on average in $4N$ generations.

- The time to the most recent common ancestor (TMRCA) has a large variance

- Even a sample with few individuals can most often recover the same TMRCA as a large sample.

- The sample size should be much smaller than the population size, although severe problems appear only with sample sizes of the same magnitude as the population size, or with non-random samples because Kingman's coalescence process assumes that maximally two sample lineages coalesce in any generation.

- With a known genealogy we can estimate the population size. Unfortunately, the true genealogy of a sample is rarely known.

Finding the best genealogy from such data is difficult

# Genetic data and the coalescent

- Finite populations loose alleles due to genetic drift

- Mutation introduces new alleles into a population at rate $\mu$

- With $2N$ chromosomes we can expect to see every generation $2N\mu$ new mutations. The population size $N$ is positively correlated with the mutation rate $\mu$.

- With genetic data sampled from several individuals we can use the mutational variability to estimate the population size.

The observed genetic variability

$$\mathcal{S} = f(N, \mu, n).$$

Different $N$ and appropriate $\mu$ can give the same number of mutations. For example, for 100 loci sampled from 20 individuals with 1000bp each, we get :

| $N$ | $\mu$ | $4N\mu$ | $\hat{S}$ | $\sigma_S^2$ |
|---|---|---|---|---|
| 1250 | $10^{-5}$ | 0.05 | 153.95 | 16.25 |
| 12500 | $10^{-6}$ | 0.05 | 152.89 | 16.05 |

Using genetic variability alone therefore does not allow to disentangle $N$ and $\mu$.

With multiple dated samples and known generation time we can estimate $N$ and $\mu$ independently.

# Mutation-scaled population size

By convention we express most results as the compound $N\mu$ and an inheritance scalar $x$, for simplicity we call this the mutation-scaled population size

$$\Theta = xN\mu,$$

where $\mu$ is the mutation rate per generation and per site. With a mutation rate per locus we use $\theta$.

- for diploids: $\Theta = 4N\mu$.

- for haploids: $\Theta = 2N\mu$.

- For mtDNA in diploids with strictly maternal inheritance this leads to $\Theta = 2N_f\mu$, and if the sex ratio is $1:1$ then $\Theta = N\mu$

Most real populations do not behave exactly like Wright-Fisher populations, therefore we subscript $N$ and call it the effective population size $N_e$, and consider $\Theta$ the mutation-scaled EFFECTIVE population size.

# Mutation-scaled population size

By convention we express most results as the compound $N\mu$ and an inheritance scalar $x$, for simplicity we call this the mutation-scaled population size
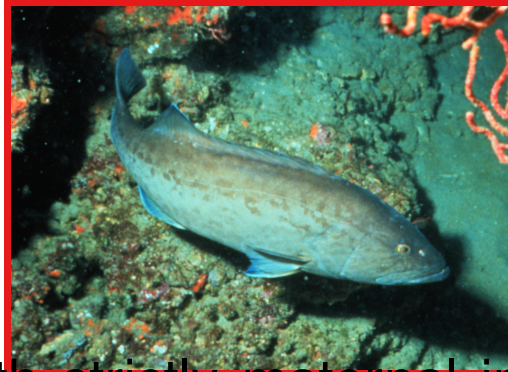
$$\Theta = xN\mu,$$

where $\mu$ is the mutation rate per generation and per site. With a mutation rate per locus we use $\theta$.

- for diploids: $\Theta = 4N\mu$.

- for haploids: $\Theta = 2N\mu$.



Gag Grouper starts out as a female and later in live becomes male.

- For mtDNA in diploids with strictly maternal inheritance this leads to $\Theta = 2N_f\mu$, and if the sex ratio is $1:1$ then $\Theta = N\mu$

Most real populations do not behave exactly like Wright-Fisher populations, therefore we subscript $N$ and call it the effective population size $N_e$, and consider $\Theta$ the mutation-scaled EFFECTIVE population size.

Humpback whales in the North Atlantic: Census population size around 12,000.

# Historical humpback whale population size

using the data by Joe Roman and Stephen R. Palumbi (Science 2003 301: 508-510)

| | | |
|---|---|---|
| $\Theta = 2N_{\female}\mu$ | 0.01529 | Population size of the North Atlantic population, estimated using migrate |
| $N_{\female} = \frac{\Theta}{2\mu}$ | 31,854 | with $\mu = 2.0 \times 10^{-8} \mathrm{bp}^{-1}\mathrm{year}^{-1}$ and a generation time of 12 years |
| $N_e = N_{\female} + N_{\male}$ | 63,708 | Sex ratio is 1:1 |
| $N_B = 2N_e$ | 127,417 | ratio $N_B/N_e$ assumed, using other data |
| $N_T = N_B \frac{N_{\text{juveniles}} + N_{\text{adults}}}{N_{\text{adults}}}$ | 203,867 | from catch and survey data (used a ratio of 1.6) |

More modern estimates for mtDNA: $150,000$ [improved estimation of mutation rate]; for nucDNA: $112,000 (45,000 - 235,000)$ [Conservation Genetics (2013) 14:103114]

Using the infinite sites model we use the number of variable sites $S$ per locus to calculate the mutation-scaled population size:

$$\theta_W = \frac{S}{\sum\limits_{k=1}^{n-1} \frac{1}{k}}$$

from a sample of $n$ individuals. For a single population the Watterson's estimator works marvelously well, but it is vulnerable to population structure.

Watterson's $\theta_W$ uses a mutation rate per locus! To compare with other work use mutation rate per site.

For Bayesian inference we want to calculate the probability of the model parameters given the data $\mathrm{p}(\mathrm{model}|\mathrm{D})$.
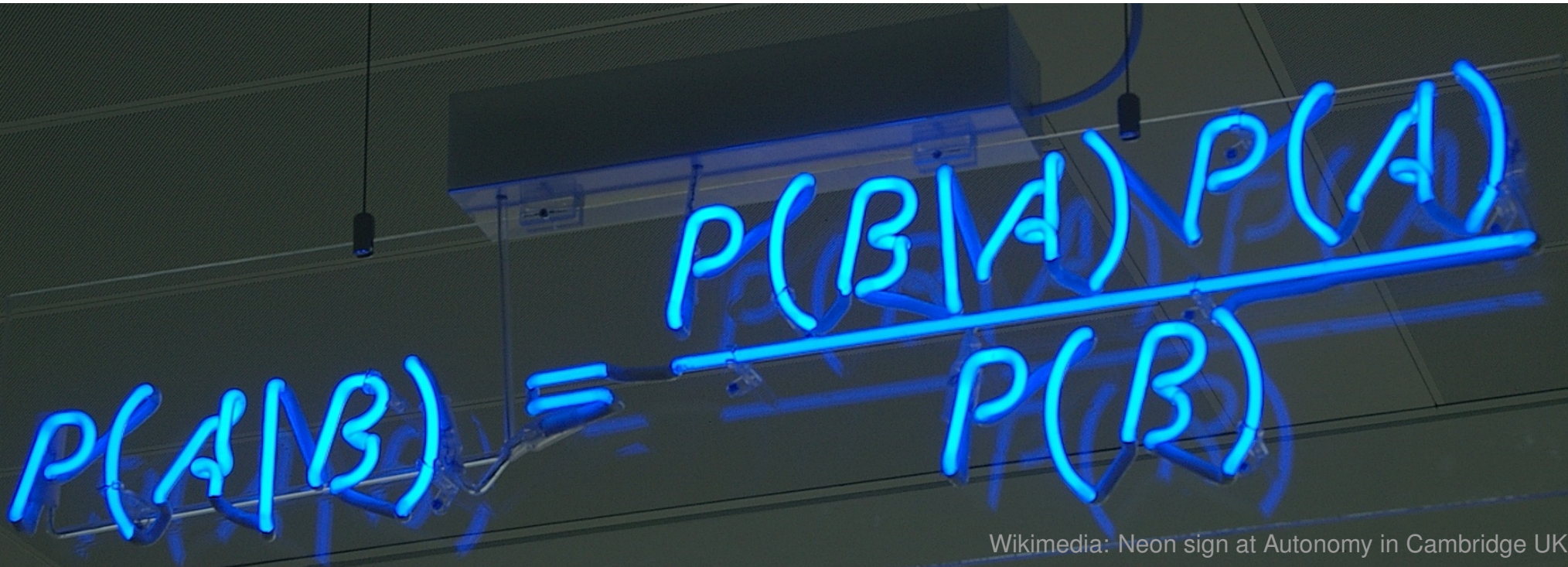
Coalescent         to describe the population genetic processes.

Mutation model     to describe the change of genetic material over time.

We calculate the Posterior distribution $p(\Theta|D)$ using Bayes' rule

$$p(\Theta|D) = \frac{p(\Theta)p(D|\Theta)}{p(D)}$$

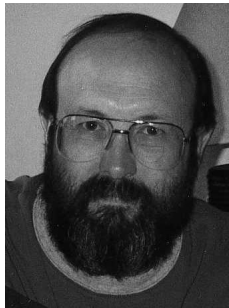where $p(D|\Theta)$ is the likelihood of the parameters.

# (almost) Felsenstein equation

$$p(D|\boldsymbol{\Theta}, G) = \mathrm{p}(\mathrm{G}|\boldsymbol{\Theta})\mathrm{p}(\mathrm{D}|\mathrm{G})$$

$p(G|\boldsymbol{\Theta})$    The probability density of a genealogy given parameters.

$\mathrm{p}(\mathrm{D}|\mathrm{G})$    The probability density of the data for a given genealogy. Phylogeneticists know this as the tree-likelihood.
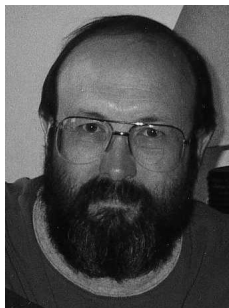
$$p(D|\boldsymbol{\Theta}) = \int_G p(G|\boldsymbol{\Theta})p(D|G)dG$$

$p(G|\boldsymbol{\Theta})$    The probability density of a genealogy given parameters.

$p(D|G)$    The probability density of the data for a given genealogy. Phylogeneticists know this as the tree-likelihood.
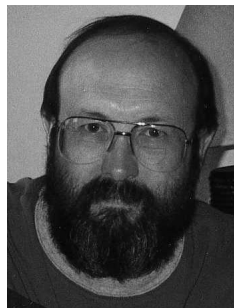
$$p(D|\mathbf{\Theta}) = \sum_{G} \mathrm{p}(\mathrm{G}|\mathbf{\Theta})\mathrm{p}(\mathrm{D}|\mathrm{G})$$

$p(G|\mathbf{\Theta})$   The probability of a genealogy given parameters.

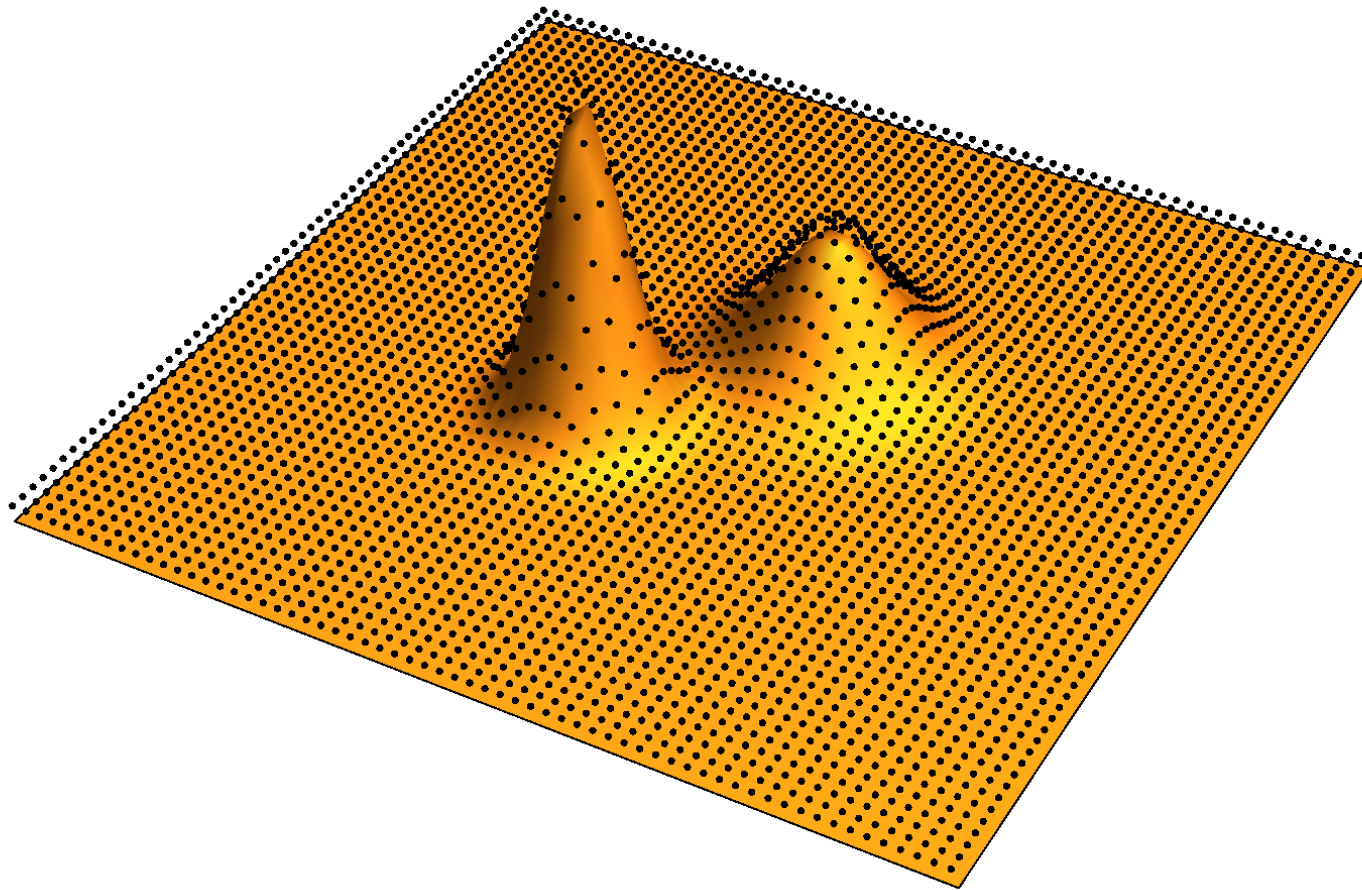$\mathrm{p}(\mathrm{D}|\mathrm{G})$   The probability of the data for a given genealogy. Phylogeneticists know this as the tree-likelihood.
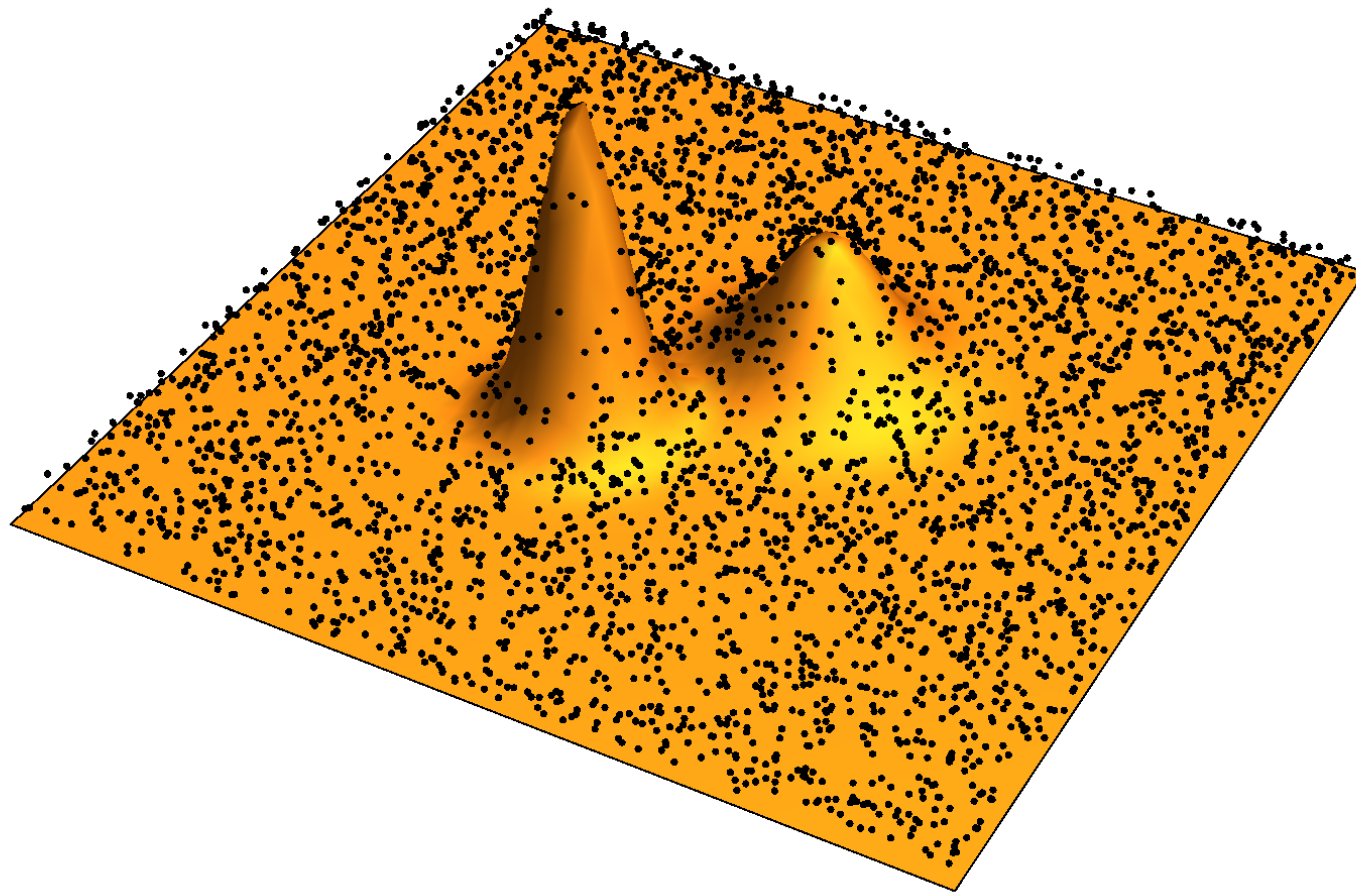
# Problem with integration formula
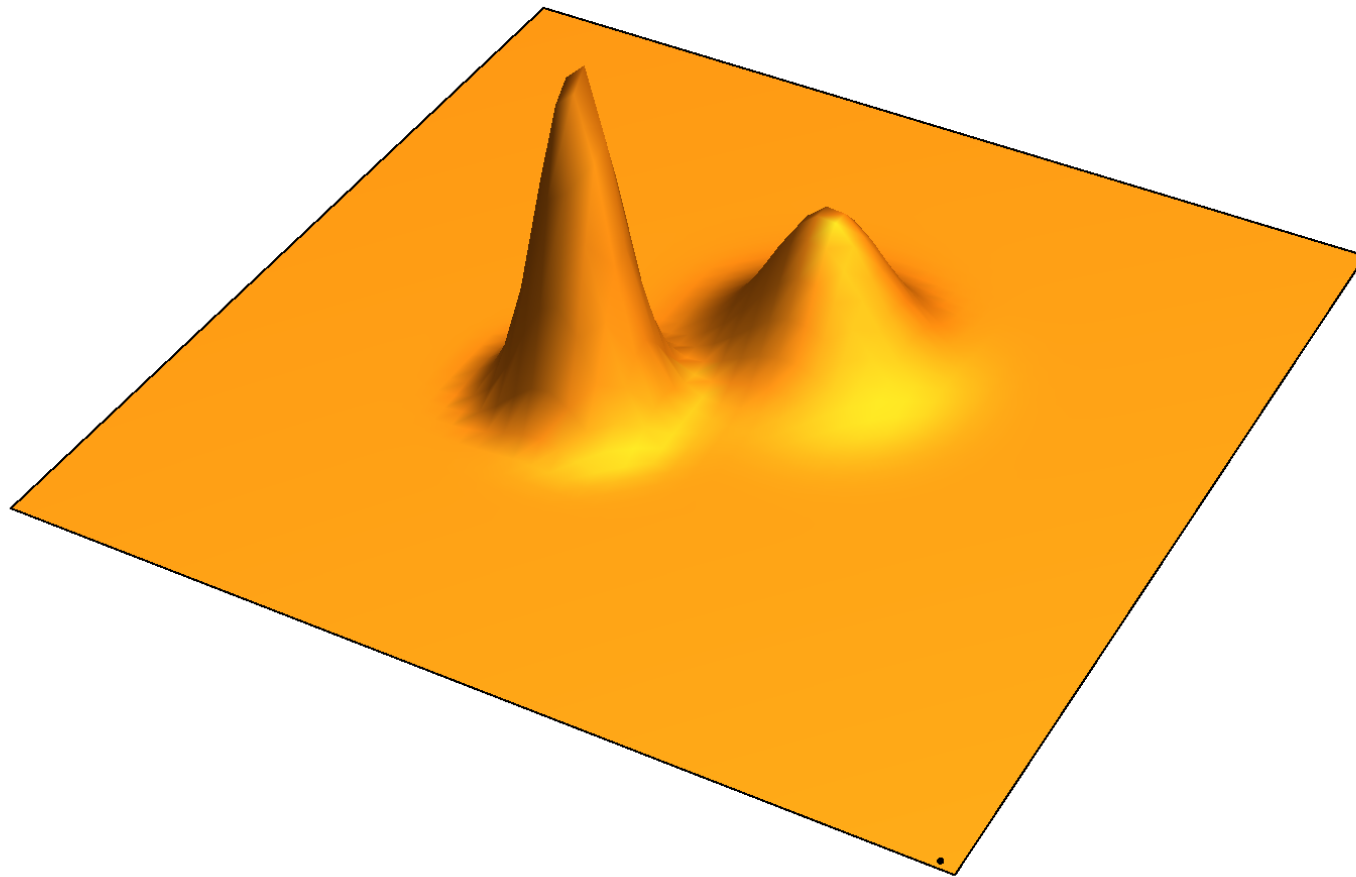
$$p(D|\Theta) = \int_G p(G|\Theta)p(D|G)dG$$

The number of possible genealogies is very large and for realistic data sets, programs need to use Markov chain Monte Carlo methods.
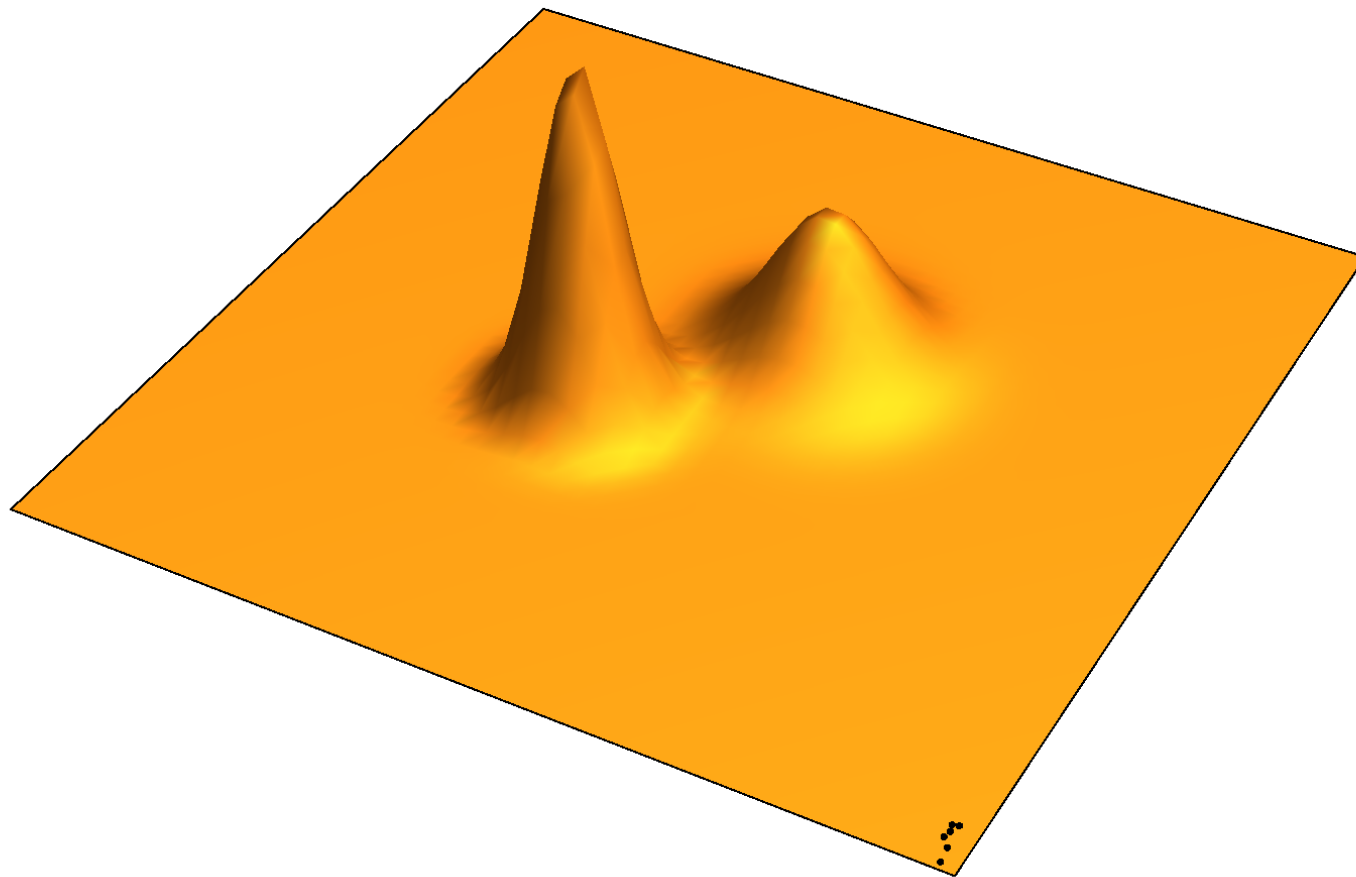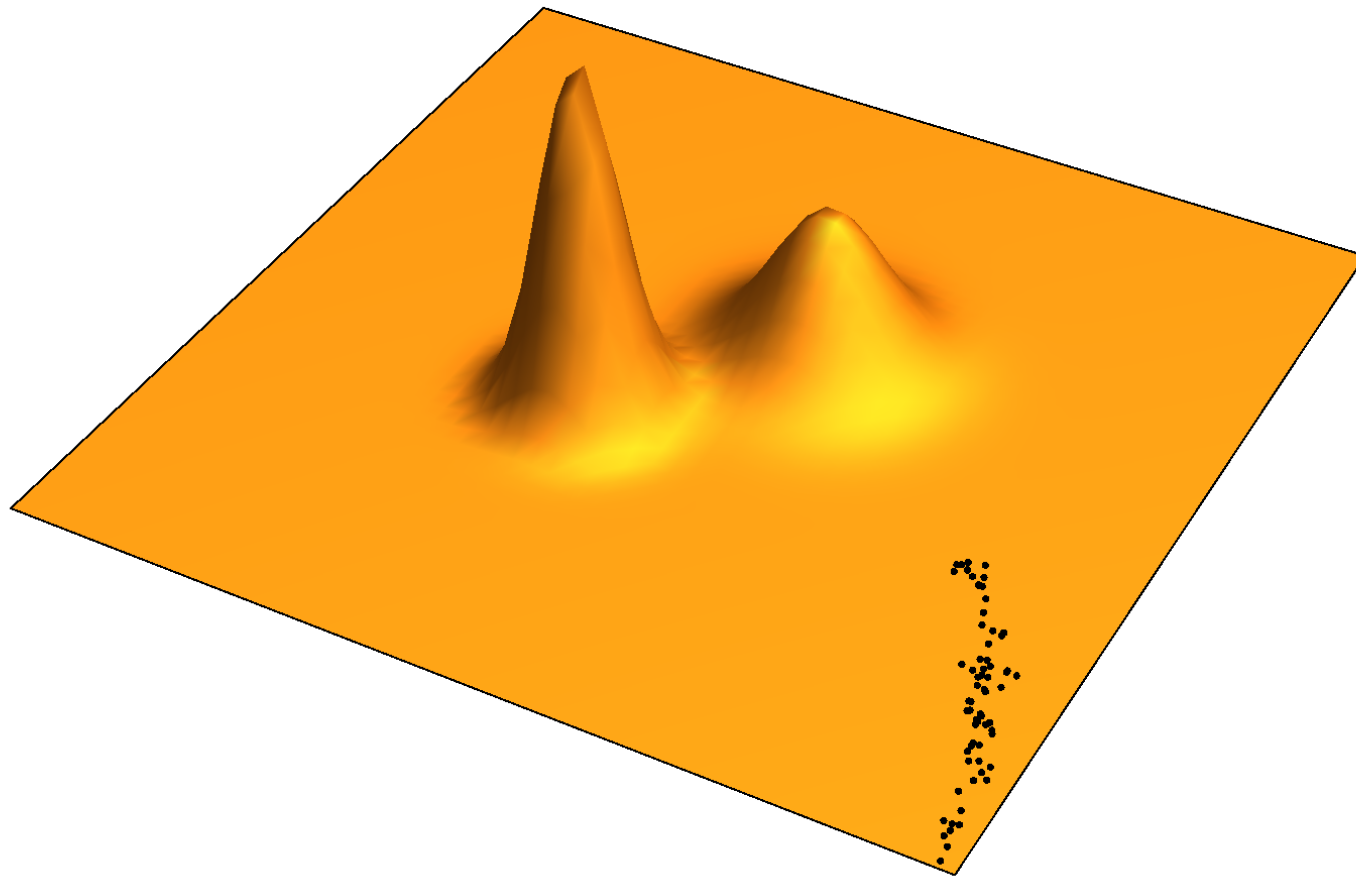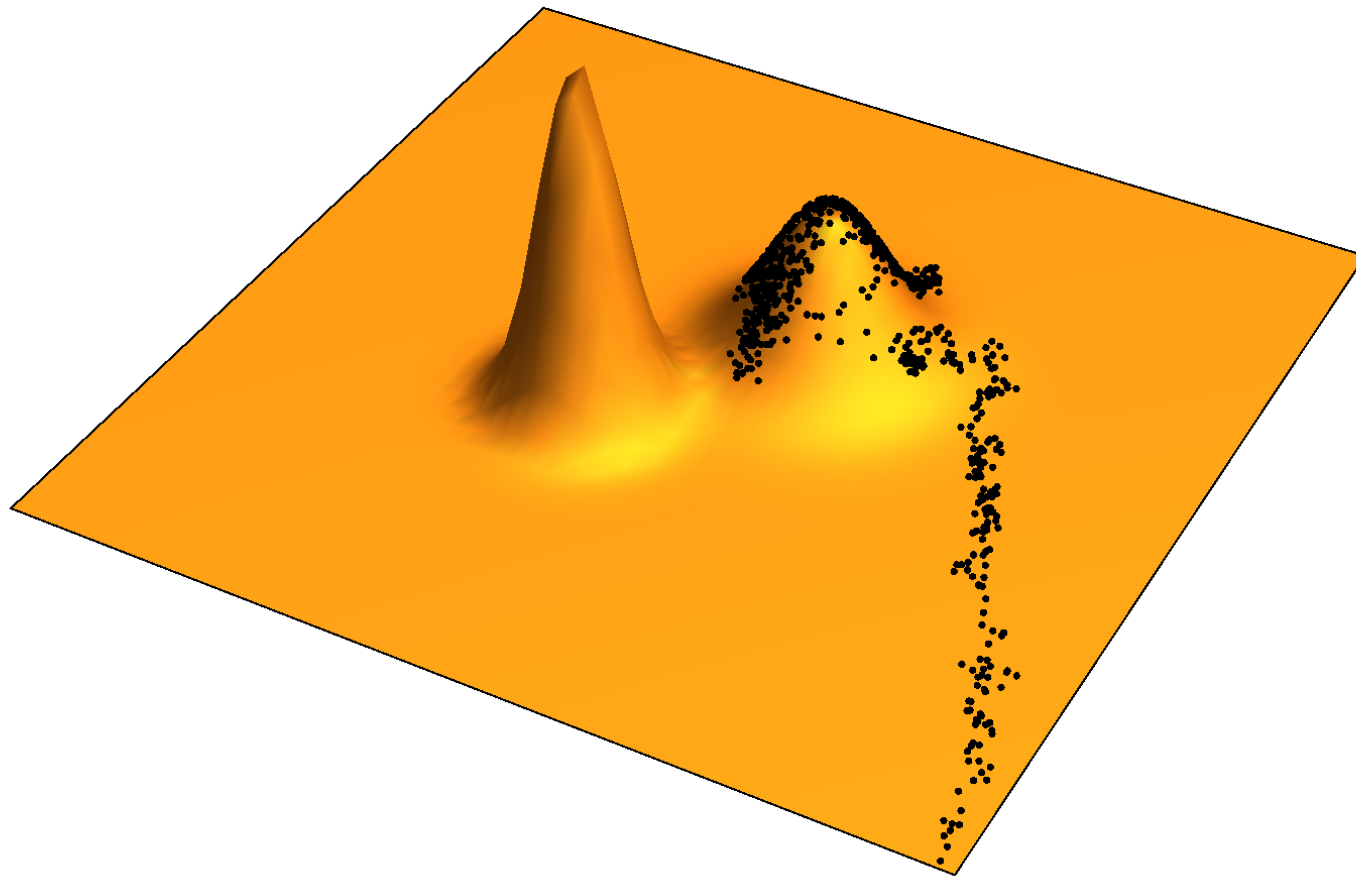
# Metropolis-Hastings algorithm

2.5% percentile=0.007 — Mode=0.00903

Median=0.00934
Mean=0.00934

97.5% percentile=0.0118



Bayesian inference: $\Theta = 0.00903$

Watterson Estimator $\Theta_W = 0.01003$

# References

Coalescent:

Nuu-Cha-Nulth population size: J. Felsenstein. 1971. Inbreeding and variance effective numbers in populations with overlapping generations. Genetics 68:581-597;

R. H. Ward, B. L. Frazier, Kerry Dew-Jager, and S. Pääbo. 1991. Extensive mitochondrial diversity within a single Amerindian tribe. PNAS 88:8780-8724;

Sigurğardóttir S, Helgason A, Gulcher JR, Stefansson K, Donnelly P. 2000. The mutation rate in the human mtDNA control region. Am J Hum Genet. 66:1599-609;

S. Matsumura and P. Forster. 2008. Generation time and effective population size in Polar Eskimos. Proc. R. Soc. B 275:1501-1508.

Sample size:

Felsenstein, J.2005. Accuracy of coalescent likelihood estimates: Do we need more sites, more sequences, or more loci? MBE 23: 691-700.

Pluzhnikov A, Donnelly P. 1996. Optimal sequencing strategies for surveying molecular genetic diversity. Genetics 144: 1247-1262.