

1 **Supplementary material for “soibean:**  
2 **High-Resolution Taxonomic Identification of**  
3 **Ancient Environmental DNA Using Mitochondrial**  
4 **Pangenome Graphs”**

5

6 Nicola Alexandra Vogel<sup>1\*</sup>, Joshua Daniel Rubin<sup>1</sup>, Anders Gorm Pedersen<sup>1</sup>, Peter  
7 Wad Sackett<sup>1</sup>, Mikkel Winther Pedersen<sup>2</sup>, Gabriel Renaud<sup>1\*</sup>

8 **1** Department of Health Technology, Section for Bioinformatics, Technical  
9 University of Denmark, Kongens Lyngby, Denmark

10 **2** Centre For Ancient Environmental Genomics, Globe Institute, University of  
11 Copenhagen, Copenhagen K, Denmark.

12 \* navo@dtu.dk, gabriel.reno@gmail.com

13 **Contents**

14 **1 Supplementary Results** **4**

15 1.1 Computational Performance . . . . . 4

16 1.2 Simulation Creation . . . . . 5

17 1.3 Benchmark . . . . . 7

18 1.4 Benchmarking Commands . . . . . 8

19 1.4.1 `soibean` . . . . . 8

20 1.4.2 `HAYSTAC` . . . . . 8

21 1.4.3 `pathPhynder - best path method` . . . . . 8

22 1.4.4 `pathPhynder - maximum likelihood method` . . . . . 9

23 1.5 Two Source Simulation . . . . . 10

24 1.6 Empirical Data: Commands . . . . . 11

25 1.6.1 `soibean` . . . . . 11

26 1.6.2 Phylogenetic Placement . . . . . 11

27 **2 Supplementary Methods** **12**

28 2.1 Database and Tree Construction . . . . . 12

29 2.2 HKY Model . . . . . 14

30 2.3 MCMC Sampling . . . . . 16

31 2.3.1 MCMC Diagnostics . . . . . 19

32 **3 Supplementary Table** **24**



# 34 1 Supplementary Results

## 35 1.1 Computational Performance

36 `soibean`'s computational performance highly depends on the used reference graph,  
37 the size of the input FASTQ and the number of iterations specified for the MCMC.  
38 Supplementary Figure 26 shows a comparison of user time in hours for three dif-  
39 ferent sizes of FASTQ input and three different numbers of iterations set for the  
40 MCMC. Each run multithreaded to 20 cores. `soibean` assesses the log-likelihood for  
41 the FASTQ-input in each iteration of the MCMC to take a new placement on the  
42 branch into account. Given our HKY model (see Method Section 2.2, our model  
43 needs to recompute the log-likelihood for every base of every fragment associated  
44 with the source. This is time-intensive; however, it provides us with extremely ac-  
45 curate estimations. Additionally, `soibean` can be multi-threaded, improving time  
46 performance.

47  
48 We have observed that `soibean`'s memory consumption depends primarily on the  
49 reference graph used. On average, `soibean` uses about 1.5Gb of memory. This varies  
50 depending on the size of the input graph.

51



## 52 1.2 Simulation Creation

53 We created six simulated samples to test `soibean`. The first sample was simu-  
54 lated to be a single source from an ancestral state sequence (N4) from the family  
55 of bears (Ursidae). Our ancestral states were simulated using `FastML` (Ashkenazy  
56 et al. 2012) (please find details about our database, including the ancestral state  
57 reconstruction, in the Supplementary Method Section ??). The ancestral state we  
58 sampled from for our first simulated data set is the most recent common ancestor  
59 (MRCA) to all Black bears (*Ursus thibetanus formosanus* - NCBI accession number  
60 NC\_009331.1, *Ursus thibetanus thibetanus* - NCBI accession number NC\_011118.1,  
61 *Ursus thibetanus* - NCBI accession number NC\_009971.1, *Ursus thibetanus ussuri-*  
62 *cus* - NCBI accession number NC\_011117.1, *Ursus thibetanus mupinensis* - NCBI  
63 accession number NC\_008753.1 and *Ursus americanus* - NCBI accession number  
64 NC\_003426.1), including the Sun bear (*Helarctos malayanus* - NCBI accession num-  
65 ber NC\_009968.1). We sampled the N4 sequence for five different levels of coverage  
66 (500, 250, 75, 50 and 10 reads). For our two-source mixture samples, we sampled  
67 the closely related bear species ( $\sim 98.8\%$  genome similarity) from two Asian Black  
68 bears, namely *Ursus thibetanus formosanus* - NCBI accession number NC\_009331.1  
69 and *Ursus thibetanus thibetanus* - NCBI accession number NC\_011118.1. The less  
70 closely related mixture samples ( $\sim 93.1\%$  genome similarity) were simulated from  
71 the Cave bear and the Brown bear mitogenome, *Ursus spelaeus* - NCBI accession  
72 number NC\_011112.1 and *Ursus arctos* - NCBI accession number NC\_003427.1. The  
73 distantly related two-source mixture sample ( $\sim 83.4\%$  genome similarity) was sim-

74 ulated from the American Black bear (*Ursus americanus* - NCBI accession number  
75 NC\_003426.1) and the Giant Panda bear (*Ailuropoda melanoleuca* - NCBI acces-  
76 sion number NC\_009492.1). Our three-source sample was simulated from the family  
77 Saturniidae (winged insects), where we sampled from two emperor moths *Gonimbra-*  
78 *sia belina* - NCBI accession number NC\_046032.1 and *Gonimbrasia tyrrhea* - NCBI  
79 accession number NC\_061326.1. The third source was sampled from the ancestral  
80 state N7 of said family. This ancestral state is the MRCA of five African moth  
81 species, namely *Nudaurelia wahlbergi* - NCBI accession number NC\_061330.1, *Go-*  
82 *nimbrasia cytherea* - NCBI accession number NC\_061331.1, *Bunaea alcinoe* - NCBI  
83 accession number NC\_061295.1 and the two previously named emperor moths. The  
84 last sample was simulated to contain four different sources from the family Phoci-  
85 dae. We simulated reads from the species *Phoca vitulia* - NCBI accession number  
86 NC\_001325.1, *Phoca fasciata* - NCBI accession number NC\_008428.1, *Phoca groen-*  
87 *landica* - NCBI accession number NC\_008429.1 and *Phoca largha* - NCBI accession  
88 number NC\_008430.1.

89

90 All simulated datasets were created using `gargammel`. We used the `fragSim`  
91 option to generate the wanted number of reads with a fragment length distribu-  
92 tion following a log-normal distribution with  $\mu = 3.7344$  and a  $\sigma = 0.35$  as com-  
93 monly seen in aDNA studies. We used `deamSim` with deamination rates taken from  
94 Günther et al. (2015) to simulate ancient damage and finally used `adptSim` to sim-  
95 ulate adapters. We used `ART` (Huang et al. 2012) to simulate Illumina sequenc-  
96 ing and finally used `leeHom` (Renaud et al. 2014) with ancient parameters to trim

97 adapters. Our mixture samples were downsampled using the `seqtk` (Shen et al.  
98 2016) subsampling function. All simulated datasets are in our provided test data  
99 <https://github.com/nicolaavogel/soibeanDatabase>.

100

### 101 **1.3 Benchmark**

102 We benchmarked against `pathPhynder` with standard parameters (Martiniano et al.  
103 2022) for our single-source simulations. We followed the tool’s recommended guide-  
104 lines and the approach taken in Kjær et al. (2022). We used `BWA ALN` (Li and Durbin  
105 2009) for mapping to the consensus sequence. The `pathPhynder` best path results  
106 for our single-source samples are visualised in Supplementary Figures 2, 3, 4, 5, and  
107 6. `pathPhynder`’s maximum likelihood results for the single-source samples are vi-  
108 sualised in Supplementary Figures 7, 8, 9, 10, and 11.

109

110 To the best of our knowledge, no existing tool estimates more than one source  
111 on a phylogenetic tree. Therefore, we used `HAYSTAC` as a baseline model. The  
112 tool specialises in classifying highly similar samples (Dimopoulos et al. 2022). We  
113 used `HAYSTAC`’s standard parameters and recommended thresholds for their coverage  
114 evenness filter; additionally, we specified that a source needs at least five uniquely  
115 mapped reads to be considered present. All results are summarised in Supplementary  
116 Table 1.

## 117 1.4 Benchmarking Commands

### 118 1.4.1 soibean

119 `./make_graph_files.sh [database name]`

120

121 `vgan soibean -fq1 [simulation.fq.gz] -t 20 --dbprefix [database name]`

122 `-o [output name]`

### 123 1.4.2 HAYSTAC

124 `haystac sample --output [sample_dir] --fastq [simulation.fq.gz]`

125 `--core 20 --trim-adapter False`

126

127 `haystac analyse --mode abundances --database [database]`

128 `--sample [sample_dir] --output [ouput_dir] --bowtie2-threads 20`

129 `--aDNA --cores 20`

### 130 1.4.3 pathPhynder - best path method

131 `bwa aln -l 1024 -n 0.001 -t 10 [consensus.fasta] [simulation.fq.gz]`

132 `| bwa samse [consensus.fasta] - [simulation.fq.gz] | samtools view`

133 `-F 4 -q 25 -@ 10 -uS - | samtools sort -@ 10 -o [output.sort.bam]`

134

135 `phynder -B -o[output.snp] [tree.newick] [input.vcf]`

136

```
137 pathPhynder -s prepare -i [tree.newick] -p [taxa_pathphynder_tree]
138 -f [output.snp] -r [consensus.fasta]
139
140 pathPhynder -s all -t 10 -i [tree.newick]
141 -p [taxa_pathphynder_tree] -b [output.sort.bam] -r [consensus.fasta]

142 1.4.4 pathPhynder - maximum likelihood method

143 bwa aln -l 1024 -n 0.001 -t 10 [consensus.fasta] [simulation.fq.gz]
144 | bwa samse [consensus.fasta] - [simulation.fq.gz] | samtools view
145 -F 4 -q 25 -@ 10 -uS - | samtools sort -@ 10 -o [output.sort.bam]
146
147 phynder -B -o[output.snp] [tree.newick] [input.vcf]
148
149 pathPhynder -s prepare -i [tree.newick] -p [taxa_pathphynder_tree]
150 -f [output.snp]
151
152 Rscript R/make_vcf.R intree_folder2/ 1 [sample.vcf]
153
154 ./phynder/phynder -q [output.vcf] -p 0.01 -o [query.output]
155 [tree.newick] [input.vcf]
```

## 156 1.5 Two Source Simulation

157 Additionally, we repeated this experiment with a distantly related mix of a Giant  
158 Panda bear and an American Black bear (83.4% similarity), as well as a very closely  
159 related mix of the Tibetan and Taiwan Black bear (98.8% similarity). Both mixtures  
160 were simulated in the same four proportions: 95% – 5%, 85% – 15%, 75% – 25% and  
161 55% – 45% with an average coverage of  $\sim 2.5X$  (1000 aDNA fragments) as in Pedersen  
162 et al. (2021). At this coverage, we obtain correct identification for the Giant Panda  
163 and American Black bear mixture at any proportion (Supplementary Figure 12),  
164 while we misidentify the closely related sample for the 95% – 5% and the 85% – 15%  
165 mixtures of the Asian Black bears (Supplementary Figure 13).

166

167 We repeat **soibean**'s robustness test by downsampling both mixtures to  $\sim 1.3X$ ,  $\sim$   
168  $0.7X$  and  $\sim 0.25X$  coverage. **soibean** can identify the mixture of the Giant Panda  
169 and the American Black bear for every mixture percentage until the lowest coverage  
170 (see Supplementary Figures 14, 15,16). For the most similar sample of Asian Black  
171 bears, we can observe that the 95% – 5% can never be identified, while the rest can  
172 be placed on the correct tree branch (see Supplementary Figures 20, 21 and 22),  
173 when downsampling.

174

## 175 1.6 Empirical Data: Commands

### 176 1.6.1 soibean

177 ./make\_graph\_files.sh [database name]

178

179 vgan soibean -fq1 [input.fastq] --dbprefix [database name] -t 20

180 -o [output name] -k 3

### 181 1.6.2 Phylogenetic Placement

182 SHRiMP\_2\_2\_2/bin/gmapper -N 4 -o 1 --single-best-mapping --sam-unaligned

183 --fastq --sam --no-qv-check --qv-offset 33 [input.fastq] [reference.fa]

184 | samtools view -bS -F4 /dev/stdin | samtools fillmd -b /dev/stdin

185 [reference.fa] | samtools sort /dev/stdin > [output.bam]

186 samtools index [output.bam]

187

188 bam2prof -double -5p [output\_5p.prof] -3p [output\_3p.prof] [output.bam]

189

190 endoCaller -seq [output.fa] -log [output.log] -name [name] -deam5p [output\_5p.prof]

191 -deam3p [output\_3p.prof] [reference.fa] [output.bam]

192

193 cat [output.fa] [phylogeneticSeqs.fa] > [output\_all.fa]

194

195 prank -d=[output\_all.fa] -o=[output\_all.prank] -showall -DNA

196

```
197 RAxML-8.2.12/raxmlHPC-PTHREADS -s [output_all.prank].best.fas --HKY85  
198 -m GTRGAMMA -n [output_prefix] -w [output_dir] -p 76 -T 40 -d
```

199

## 200 2 Supplementary Methods

### 201 2.1 Database and Tree Construction

202 The database for `soibean` uses the same taxa introduced in the database for `euka`  
203 (Vogel et al. 2023), which comprise 335 different taxa of tetrapodic and arthropodic  
204 eukaryotes (the taxon `Galloanserae` has been resorted since `euka`'s latest release).

205 We used the same reference genomes as input to construct a multiple sequence align-  
206 ment using `PRANK` version v.170427 (Löytynoja 2014), from which we constructed  
207 a phylogenetic tree using `RAxML` version 8.2.12 with the Hasegawa Kishino Yano  
208 (HKY) substitution model (Hasegawa et al. 1985) (`--HKY85`) and estimated base  
209 frequencies (Stamatakis 2014). The phylogenetic tree is used to infer the ancestral  
210 states for each taxon using `FastML` version 3.11 with the `-mr` flag for mitochon-  
211 drial genomes and the `-mh` flag to use the HKY model (Ashkenazy et al. 2012).

212 We use the HKY substitution model for our guide tree as it allows for both nu-  
213 cleotide frequencies and the transition/transversion ratio to vary (Hasegawa et al.  
214 1985). These features are crucial for the estimation of our guide tree and fitting our  
215 model because the transition/transversion ratio is much higher in mtDNA than in  
216 the autosomes (Tamura 1992) and because mtDNA is guanine-enriched on the heavy



217 strand to aid in stabilising RNA transcription (Wallace 1994; Levinstein Hallak et al.  
218 2018). The HKY model allowed us to represent the variety of taxa used accurately.  
219 **FastML** outputs a multiple sequence alignment (MSA), including reconstructed an-  
220 cestral states and the maximum-likelihood phylogenetic tree. We build a pangenome  
221 graph for each taxon from the new MSA with the added ancestral state sequences  
222 via the `vg construct -M` subcommand from the `vg` toolkit version 1.44.0 (Garrison  
223 et al. 2018) (Supplementary Figure 1), resulting in 326 subgraphs (i.e. connected  
224 components), each subgraph corresponding to a different taxon for **euka/soibean**.  
225 A pangenome graph is defined as  $G = (N, E, H)$ , where  $N$  are the nodes,  $E$  the  
226 edges, and  $H$  the graph’s embedded reference paths. Each node has a unique nu-  
227 merical identifier called node ID, which is internally assigned by `vg`. For a detailed  
228 overview of pangenome graphs applied to animal mitogenomes, see the Supplemen-  
229 tary Material from Vogel et al. (2023). Each path in the graph corresponds to a  
230 mitogenome in the MSA and, therefore, a node in the phylogenetic tree. Inferred  
231 ancestral genomes correspond to internal nodes, whereas the mitogenomes used as  
232 input to the MSA are leaves in the phylogenetic tree. We merged connected com-  
233 ponents for each pangenomic graph into a single set of multiple independent graphs  
234 using the `vg ids` subcommand (Garrison et al. 2018). This resulted in a combined  
235 graph with  $|N| = 6889846$ ,  $|E| = 10880449$  and  $|H| = 10992$ . We updated each  
236 taxon’s start and end node ID for our graph index and additionally empirically com-  
237 puted the base frequencies for each taxon to be used for our HKY-based likelihood  
238 model. A complete overview of the database construction and all scripts used can  
239 be found <https://github.com/nicolaavogel/soibeanDatabase>.

## 241 2.2 HKY Model

242 We are interested in finding the most likely placement of ancient species on a pre-  
 243 computed mitochondrial phylogenetic tree. To achieve this we use the HKY sub-  
 244 stitution model to compute the probability of an observed, ancient DNA sequence,  
 245 arising by substitutions occurring along a branch of the tree, some distance from one  
 246 of the existing nodes.

247 The instantaneous rate matrix for the HKY model can be written as (Hillis et al.  
 248 1996; Xia 2017):

$$249 \quad \mathbf{Q} = \mu \begin{bmatrix} - & \pi_C & \kappa\pi_G & \pi_T \\ \pi_A & - & \pi_G & \kappa\pi_T \\ \kappa\pi_A & \pi_C & - & \pi_T \\ \pi_A & \kappa\pi_C & \pi_G & - \end{bmatrix}$$

250 The rows and columns of  $\mathbf{Q}$  correspond to the nucleotides A, C, G, and T in that  
 251 order. An element,  $Q_{ij}$ , gives the rate at which nucleotide  $i$  changes to nucleotide  
 252  $j$ , and values are in units of substitutions per site per time unit. Elements on the  
 253 diagonal (not shown here) are set such that all rows sum to zero:  $Q_{ii} = -\sum_{j \neq i} Q_{ij}$ .  
 254 In the expressions above  $\kappa$  is the ratio between the rates of transition-type and  
 255 transversion-type substitutions, while  $\pi_A, \pi_C, \pi_G$ , and  $\pi_T$  are the equilibrium nu-  
 256 cleotide frequencies. The multiplicative factor  $\mu$  will, in part, depend on the units  
 257 of time used in the matrix (e.g., rates measured per century will be 100 times larger

258 than rates measured per year).

259 Given the rate matrix  $\mathbf{Q}$  it is possible to compute the transition probability  
260 matrix  $\mathbf{P}(t)$  for a given branch length  $t$ :  $\mathbf{P}(t) = e^{\mathbf{Q}t}$ . In this expression,  $t$  has to be  
261 expressed in the same time units as  $\mathbf{Q}$ . Since branch lengths in the pre-computed  
262 phylogeny are expressed as the expected change per site (and not in some unit of  
263 time), we first normalise  $\mathbf{Q}$  correspondingly. Specifically, this is done by setting  $\mu$   
264 such that the expected rate of change,  $\rho$ , is 1: If  $\rho = 1$  then the expected amount  
265 of change on a branch of length  $t$  will be:  $\nu = \rho t = t$ . This is equivalent to having  
266 chosen a unit of time such that the expected change on a branch ( $\nu$ , measured in  
267 substitutions per site) is numerically identical to the branch length ( $t$ , measured in  
268 those time units).

269 The expected rate of change for the HKY model is (Hillis et al. 1996):

$$\begin{aligned} 270 \quad \rho &= \sum_i \pi_i \sum_{j \neq i} Q_{ij} \\ 271 \quad &= 2(\kappa\pi_A\pi_G + \kappa\pi_C\pi_T + \pi_A\pi_C + \pi_A\pi_T + \pi_C\pi_G + \pi_G\pi_T)\mu \end{aligned}$$

272 Setting this equal to 1, and isolating  $\mu$ , we find:

$$273 \quad \mu = \frac{1}{2(\kappa\pi_A\pi_G + \kappa\pi_C\pi_T + \pi_A\pi_C + \pi_A\pi_T + \pi_G\pi_C + \pi_G\pi_T)}$$

274 Using  $\mathbf{P}(t) = e^{\mathbf{Q}t}$  one can derive the following expressions for the probability of  
275 observing nucleotide  $b$ , in an ancient taxon, given the presence of the nucleotide  $r$ ,  
276 at a distance of  $t$ , in a node on the tree:

$$\begin{cases}
\pi_r + \pi_r(\frac{1}{\Pi_r} - 1)e^{-\mu t} + (\frac{\Pi_r - \pi_r}{\Pi_r})e^{\mu t A} & (b = r), \\
\pi_r + \pi_r(\frac{1}{\Pi_r} - 1)e^{-\mu t} - (\frac{\pi_r}{\Pi_r})e^{-\mu t A} & (b \neq r, \text{ transition}), \\
\pi_r(1 - e^{-\mu t}) & (b \neq r, \text{ transversion})
\end{cases}$$

Here  $A = 1 + \Pi_r(\kappa - 1)$ , with  $\Pi_r = \pi_A + \pi_G$  (if  $r$  is A or G) and  $\Pi_r = \pi_C + \pi_T$  (if  $r$  is C or T). The branch length  $t$  is the distance between the ancient taxon and a nearby node on the tree and is expressed in units of expected change as explained above. Based on previous analyses of human mitochondrial evolution, we set  $\kappa$  to 22 (Levinstein Hallak et al. 2018). The equilibrium nucleotide frequencies are set to their empirical values in the data set at hand. The value of  $\mu$  is computed from  $\kappa$  and  $(\pi_A, \pi_C, \pi_G, \pi_T)$  according to the expression above.

### 2.3 MCMC Sampling

The MCMC process is run  $k$  times, corresponding to the number of significant signature node paths identified (or alternatively, the number specified by the user if one is provided). For each MCMC run, four chains are deployed: the first chain is initialized at the identified signature nodes, whereas the subsequent three incorporate random nodes from the tree as starting points. The initial run of the MCMC ( $k = 1$ ) initializes placement from the most prevalent signature node path. Subsequent runs incrementally include the next most frequent paths as starting points, continuing this

294 pattern until  $k$  equals the total count of significant signature node paths. The initial  
295 proportion of each source is set at  $1/k$ , and the initial branch position is sampled  
296 uniformly from  $[0, 1]$ , where 0 signifies the ancestral state, and 1 is the derived state  
297 of the branch.

298

299 We use a Metropolis-Hastings algorithm (Metropolis et al. 1953; Hastings 1970)  
300 for our MCMC sampling. Throughout the MCMC runs, new branch placements and  
301 source proportions are sampled. The branches are normalised to unit length  $t$ , with  
302 the  $N_A$  marked as 0 and the  $N_D$  as 1. Proposed distances  $d$  to new placements  $\beta^*$   
303 are drawn from a normal distribution, with a mean of 0 and the standard deviation  
304 linearly decreasing from 3 to 0.1 throughout the burn-in phase, then from 0.1 to  $1e^{-10}$   
305 in subsequent iterations. Should a proposed distance be negative, the transition is  
306 towards the root; conversely, a positive distance indicates a movement towards the  
307 leaves.

308

309 If the sum of our current position on the tree branch and the newly sampled  
310 distance is smaller than one ( $d + \beta_i < 1$ ) or the difference between the current posi-  
311 tion and the sampled distance is higher than 0 ( $d - \beta_i > 0$ ), we update our current  
312 position on the branch to the updated position on the branch ( $\beta_{i \pm d}$ ). If the sum  
313 of our current position on the tree branch and the newly sampled distance is larger  
314 than 1 ( $d + \beta_i > 1$ ), we move to either of the two child branches of  $N_D$  equiprobably.  
315 In case  $N_D$  is a leaf node, we sample a new move. When the difference between the  
316 current position on the tree branch and the sampled distance is negative ( $d - \beta_i > 0$ )

317 we move to either the parent branch or the sibling branch of  $N_A$  equiprobably. If  
318  $N_A$  is the root node, we sample a new move (see Figure ??).

319

320 If  $k \neq 1$ , we simultaneously sample a new proportion vector  $\theta^*$ . Every proportion  
321  $\theta_i^*$  is sampled from  $N(\theta_i, 0.1)$  and  $\theta_i^* \in [0, 1]$ . We use the sum of  $\theta^*$  to normalise each  
322  $\theta_i^*$  so that  $\sum \theta^* = 1$ .

323

324 We calculate the new likelihood for  $P(D|\beta^*, \theta^*)$ . If  $P(D|\beta^*, \theta^*) > P(D|\beta, \theta)$  we  
325 accept the move, append the new parameter values to the sample file, and continue  
326 to the next iteration. If  $P(D|\beta^*, \theta^*) < P(D|\beta, \theta)$  we accept the move with probabil-  
327 ity  $\frac{P(D|\beta^*, \theta^*)}{P(D|\beta, \theta)} \frac{q(\beta, \theta|\beta^*, \theta^*)}{q(\beta^*, \theta^*|\beta, \theta)}$ . Here  $q$  is the proposal distribution for moves in parameter  
328 space, and  $q(X|Y)$  is the probability of proposing  $X$  when the current parameter  
329 values are  $Y$ . If this distribution is symmetric, i.e.  $q(X|Y) = q(Y|X)$ , the q-terms  
330 cancel.

331

332 We present a quick argument that the proposal sampling distribution is symmet-  
333 ric. In other words, we aim to argue that for Markov states  $X$  and  $Y$  in our parameter  
334 space, sampling  $Y$  from state  $X$  has a probability equal to that of sampling  $X$  from  
335 state  $Y$ . To see this, recall that in our sampling scheme, we first fix a distance to  
336 move and a direction (ancestral or derived) in which to move that distance. Note  
337 also that for a given proposal move, we do not allow any self-intersecting paths on  
338 the tree. Therefore, given the bifurcating nature of tree topologies, there will always  
339 only be one possible way to sample  $Y$  from  $X$  for any Markov states  $X, Y$ . Along

340 this path, there will be some number of bifurcations in which one of the two possible  
341 forks will be followed. The probability of taking any given fork is  $\frac{1}{2}$ . Since there are  
342 the same number of forks on the unique path from  $Y$  to  $X$  as there are from  $X$  to  
343  $Y$ , it clearly follows that the probability of sampling  $X$  from  $Y$  is identical to that  
344 of sampling  $Y$  from  $X$ .

345

346 By default, we perform 500,000 MCMC iterations and discard the first 75,000  
347 iterations as burn-in. These values can be overridden by the user via command-line  
348 options. Notably, lower values of  $k$  should require a lower number of iterations for  
349 convergence.

350

### 351 **2.3.1 MCMC Diagnostics**

352 After inference is complete, it is important to assess whether the chains managed  
353 to converge to the true underlying posterior distribution. This can be assessed indi-  
354 vidually for each parameter of the model. To assess the convergence of phylogenetic  
355 placement, we keep track of the patristic distance from the accepted position to all  
356 the leaves on the tree (the patristic distance is the distance along the tree between  
357 two nodes, i.e., the sum of lengths of the branches on the path between those nodes  
358 (Fourment and Gibbs 2006)). Specifically, we compute a vector where the elements  
359 are the patristic distances from the accepted position to each of the leaf nodes. We  
360 also compute a vector giving the patristic distances from the *root node* to each of the  
361 leaves. This is done for every one hundred iterations after the burn-in period. We

362 then compute the Euclidean distance between these two vectors of patristic distances.  
363 If a trace plot of this Euclidean distance against MCMC iteration shows any upward  
364 or downward trend, then that would indicate a lack of convergence. If phylogenetic  
365 placement *has* converged then the trace plot should move around a constant value.  
366 We also provide the user with the chain statistics, including the mean, median, 95%  
367 credible intervals, variance, autocorrelation, and effective sample size (ESS) for each  
368 chain and each estimated parameter. We warn the user if the ESS for a chain is be-  
369 low 200 (Martino et al. 2017), as it could mean the estimates are insufficiently precise.

370

371 Another important and widely-used metric for assessing convergence is called  $\hat{R}$ .  
372 We compute  $\hat{R}$  for both  $\boldsymbol{\theta}$  and  $\boldsymbol{\beta}$  across chains. The measure essentially compares the  
373 within-chain variance with the between-chain variance; a value close to 1 indicates  
374 convergence (Gelman and Rubin 1992; Brooks and Gelman 1998). Values larger than  
375 1.05 (Vehtari et al. 2021) may indicate non-converging chains.

376



## 377 **References**

- 378 Ashkenazy, H. et al. (2012). “FastML: a web server for probabilistic reconstruction  
379 of ancestral sequences”. In: *Nucleic acids research* 40.W1, W580–W584.
- 380 Brooks, S. P. and Gelman, A. (1998). “General methods for monitoring convergence  
381 of iterative simulations”. In: *Journal of computational and graphical statistics* 7.4,  
382 pp. 434–455.
- 383 Dimopoulos, E. A. et al. (2022). “HAYSTAC: A Bayesian framework for robust  
384 and rapid species identification in high-throughput sequencing data”. In: *PLOS*  
385 *Computational Biology* 18.9, e1010493.
- 386 Fourment, M. and Gibbs, M. J. (2006). “PATRISTIC: a program for calculating  
387 patristic distances and graphically comparing the components of genetic change”.  
388 In: *BMC evolutionary biology* 6, pp. 1–5.
- 389 Garrison, E. et al. (2018). “Variation graph toolkit improves read mapping by repre-  
390 senting genetic variation in the reference”. In: *Nature biotechnology* 36.9, pp. 875–  
391 879.
- 392 Gelman, A. and Rubin, D. B. (1992). “Inference from iterative simulation using  
393 multiple sequences”. In: *Statistical science* 7.4, pp. 457–472.
- 394 Günther, T. et al. (2015). “Ancient genomes link early farmers from Atapuerca in  
395 Spain to modern-day Basques”. In: *Proceedings of the National Academy of Sci-*  
396 *ences* 112.38, pp. 11917–11922.

- 397 Hasegawa, M., Kishino, H., and Yano, T.-a. (1985). “Dating of the human-ape split-  
398 ting by a molecular clock of mitochondrial DNA”. In: *Journal of molecular evo-*  
399 *lution* 22, pp. 160–174.
- 400 Hastings, W. K. (1970). “Monte Carlo sampling methods using Markov chains and  
401 their applications”. In: *Biometrika*.
- 402 Hillis, D. M., Moritz, C., and Mable, B. K. (1996). *Molecular systematics*. Vol. 23.  
403 Sinauer.
- 404 Huang, W. et al. (2012). “ART: a next-generation sequencing read simulator”. In:  
405 *Bioinformatics* 28.4, pp. 593–594.
- 406 Kjær, K. H. et al. (2022). “A 2-million-year-old ecosystem in Greenland uncovered  
407 by environmental DNA”. In: *Nature* 612.7939, pp. 283–291.
- 408 Levinstein Hallak, K., Tzur, S., and Rosset, S. (2018). “Big data analysis of hu-  
409 man mitochondrial DNA substitution models: a regression approach”. In: *BMC*  
410 *genomics* 19.1, pp. 1–13.
- 411 Li, H. and Durbin, R. (2009). “Fast and accurate short read alignment with Burrows-  
412 Wheeler transform”. In: *bioinformatics* 25.14, pp. 1754–1760.
- 413 Löytynoja, A. (2014). “Phylogeny-aware alignment with PRANK”. In: *Multiple se-*  
414 *quence alignment methods*, pp. 155–170.
- 415 Martiniano, R. et al. (2022). “Placing ancient DNA sequences into reference phylo-  
416 genies”. In: *Molecular biology and evolution* 39.2, msac017.
- 417 Martino, L., Elvira, V., and Louzada, F. (2017). “Effective sample size for importance  
418 sampling based on discrepancy measures”. In: *Signal Processing* 131, pp. 386–401.

- 419 Metropolis, N. et al. (1953). “Equation of state calculations by fast computing ma-  
420 chines”. In: *The journal of chemical physics* 21.6, pp. 1087–1092.
- 421 Pedersen, M. W. et al. (2021). “Environmental genomics of Late Pleistocene black  
422 bears and giant short-faced bears”. In: *Current Biology* 31.12, pp. 2728–2736.
- 423 Renaud, G., Stenzel, U., and Kelso, J. (2014). “leeHom: adaptor trimming and merg-  
424 ing for Illumina sequencing reads”. In: *Nucleic acids research* 42.18, e141–e141.
- 425 Shen, W. et al. (2016). “SeqKit: a cross-platform and ultrafast toolkit for FASTA/Q  
426 file manipulation”. In: *PloS one* 11.10, e0163962.
- 427 Stamatakis, A. (2014). “RAxML version 8: a tool for phylogenetic analysis and post-  
428 analysis of large phylogenies”. In: *Bioinformatics* 30.9, pp. 1312–1313.
- 429 Tamura, K. (1992). “The rate and pattern of nucleotide substitution in *Drosophila*  
430 mitochondrial DNA.” In: *Molecular biology and evolution* 9.5, pp. 814–825.
- 431 Vehtari, A. et al. (2021). “Rank-normalization, folding, and localization: An improved  
432  $\hat{R}$  for assessing convergence of MCMC (with discussion)”. In: *Bayesian analysis*  
433 16.2, pp. 667–718.
- 434 Vogel, N. A. et al. (2023). “euka: Robust tetrapodic and arthropodic taxa detec-  
435 tion from modern and ancient environmental DNA using pangenomic reference  
436 graphs”. In: *Methods in Ecology and Evolution*.
- 437 Wallace, D. C. (1994). “Mitochondrial DNA mutations in diseases of energy metabolism”.  
438 In: *Journal of bioenergetics and biomembranes* 26, pp. 241–250.
- 439 Xia, X. (2017). “Deriving transition probabilities and evolutionary distances from  
440 substitution rate matrix by probability reasoning”. In: *J Genet Genome Res*  
441 3.031.

442 **3** **Supplementary Table**

Reads	Average coverage	Proportion	Simulated Source(s)	Predicted Source(s)
500	1.3X	1	Ancestral state $N_4$	<i>U. ameri-</i> <i>canus</i> <i>M. ursinus</i> <i>N13</i> <i>U. spelaeus</i> <i>U. arctos</i>
250	0.66X	1	Ancestral state $N_4$	<i>H.</i> <i>malayanus</i> <i>U. ameri-</i> <i>canus</i> <i>M. ursinus</i>
75	0.2X	1	Ancestral state $N_4$	<i>U. americanus</i>
50	0.13X	1	Ancestral state $N_4$	None
10	0.026X	1	Ancestral state $N_4$	None
1000	2.6X	55-45	<i>U. thibetanus</i> <i>formosanus</i> <i>U. thibetanus</i> <i>thibetanus</i>	<i>U. thibetanus</i> <i>formosanus</i> <i>U. thibetanus</i> <i>thibetanus</i>
500	1.3X	55-45	<i>U. thibetanus</i> <i>formosanus</i> <i>U. thibetanus</i> <i>thibetanus</i>	<i>U. thibetanus</i> <i>formosanus</i> <i>U. thibetanus</i> <i>thibetanus</i>
250	0.66X	55-45	<i>U. thibetanus</i> <i>formosanus</i> <i>U. thibetanus</i> <i>thibetanus</i>	<i>U. thibetanus</i> <i>formosanus</i> <i>U. thibetanus</i> <i>thibetanus</i>

100	0.26X	55-45	<i>U. thibetanus</i> <i>formosanus</i> <i>U. thibetanus</i> <i>thibetanus</i>	<i>U. thibetanus thibetanus</i>
1000	2.6X	55-45	<i>U. spelaeus</i> <i>U. arctos</i>	<i>U. spelaeus</i> <i>U. arctos</i>
500	1.3X	55-45	<i>U. spelaeus</i> <i>U. arctos</i>	<i>U. spelaeus</i> <i>U. arctos</i>
250	0.66X	55-45	<i>U. spelaeus</i> <i>U. arctos</i>	<i>U. spelaeus</i> <i>U. arctos</i>
100	0.26X	55-45	<i>U. spelaeus</i> <i>U. arctos</i>	<i>U. spelaeus</i> <i>U. arctos</i>
1000	2.6X	55-45	<i>A.</i> <i>melanoleuca</i> <i>U. ameri-</i> <i>canus</i>	<i>A.</i> <i>melanoleuca</i> <i>U. ameri-</i> <i>canus</i>
500	1.3X	55-45	<i>A.</i> <i>melanoleuca</i> <i>U. ameri-</i> <i>canus</i>	<i>A.</i> <i>melanoleuca</i> <i>U. ameri-</i> <i>canus</i>
250	0.66X	55-45	<i>A.</i> <i>melanoleuca</i> <i>U. ameri-</i> <i>canus</i>	<i>A.</i> <i>melanoleuca</i> <i>U. ameri-</i> <i>canus</i>
100	0.26X	55-45	<i>A.</i> <i>melanoleuca</i> <i>U. ameri-</i> <i>canus</i>	<i>A.</i> <i>melanoleuca</i> <i>U. ameri-</i> <i>canus</i>

1000	2.6X	75-25	<i>U. thibetanus</i> <i>formosanus</i> <i>U. thibetanus</i> <i>thibetanus</i>	<i>U. thibetanus</i> <i>formosanus</i> <i>U. thibetanus</i> <i>thibetanus</i>
500	1.3X	75-25	<i>U. thibetanus</i> <i>formosanus</i> <i>U. thibetanus</i> <i>thibetanus</i>	<i>U. thibetanus</i> <i>formosanus</i> <i>U. thibetanus</i> <i>thibetanus</i>
250	0.66X	75-25	<i>U. thibetanus</i> <i>formosanus</i> <i>U. thibetanus</i> <i>thibetanus</i>	<i>U. thibetanus</i> <i>formosanus</i> <i>U. thibetanus</i> <i>thibetanus</i>
100	0.26X	75-25	<i>U. thibetanus</i> <i>formosanus</i> <i>U. thibetanus</i> <i>thibetanus</i>	<i>U. thibetanus formosanus</i>
1000	2.6X	75-25	<i>U. spelaeus</i> <i>U. arctos</i>	<i>U. spelaeus</i> <i>U. arctos</i>
500	1.3X	75-25	<i>U. spelaeus</i> <i>U. arctos</i>	<i>U. spelaeus</i> <i>U. arctos</i>
250	0.66X	75-25	<i>U. spelaeus</i> <i>U. arctos</i>	<i>U. spelaeus</i> <i>U. arctos</i>
100	0.26X	75-25	<i>U. spelaeus</i> <i>U. arctos</i>	<i>U. spelaeus</i> <i>U. arctos</i>
1000	2.6X	75-25	<i>A.</i> <i>melanoleuca</i> <i>U. ameri-</i> <i>canus</i>	<i>A.</i> <i>melanoleuca</i> <i>U. ameri-</i> <i>canus</i>

500	1.3X	75-25	A. <i>melanoleuca</i> U. <i>ameri-</i> <i>canus</i>	A. <i>melanoleuca</i> U. <i>ameri-</i> <i>canus</i>
250	0.66X	75-25	A. <i>melanoleuca</i> U. <i>ameri-</i> <i>canus</i>	A. <i>melanoleuca</i> U. <i>ameri-</i> <i>canus</i>
100	0.26X	75-25	A. <i>melanoleuca</i> U. <i>ameri-</i> <i>canus</i>	A. <i>melanoleuca</i> U. <i>ameri-</i> <i>canus</i>
1000	2.6X	85-15	U. <i>thibetanus</i> <i>formosanus</i> U. <i>thibetanus</i> <i>thibetanus</i>	U. <i>thibetanus</i> <i>formosanus</i> U. <i>thibetanus</i> <i>thibetanus</i>
500	1.3X	85-15	U. <i>thibetanus</i> <i>formosanus</i> U. <i>thibetanus</i> <i>thibetanus</i>	U. <i>thibetanus</i> <i>formosanus</i> U. <i>thibetanus</i> <i>thibetanus</i>
250	0.66X	85-15	U. <i>thibetanus</i> <i>formosanus</i> U. <i>thibetanus</i> <i>thibetanus</i>	U. <i>thibetanus formosanus</i>
100	0.26X	85-15	U. <i>thibetanus</i> <i>formosanus</i> U. <i>thibetanus</i> <i>thibetanus</i>	U. <i>thibetanus formosanus</i>



1000	2.6X	85-15	<i>U. spelaeus</i> <i>U. arctos</i>	<i>U. spelaeus</i> <i>U. arctos</i>
500	1.3X	85-15	<i>U. spelaeus</i> <i>U. arctos</i>	<i>U. spelaeus</i> <i>U. arctos</i>
250	0.66X	85-15	<i>U. spelaeus</i> <i>U. arctos</i>	<i>U. spelaeus</i>
100	0.26X	85-15	<i>U. spelaeus</i> <i>U. arctos</i>	<i>U. spelaeus</i>
1000	2.6X	85-15	<i>A.</i> <i>melanoleuca</i> <i>U. ameri-</i> <i>canus</i>	<i>A.</i> <i>melanoleuca</i> <i>U. ameri-</i> <i>canus</i>
500	1.3X	85-15	<i>A.</i> <i>melanoleuca</i> <i>U. ameri-</i> <i>canus</i>	<i>A.</i> <i>melanoleuca</i> <i>U. ameri-</i> <i>canus</i>
250	0.66X	85-15	<i>A.</i> <i>melanoleuca</i> <i>U. ameri-</i> <i>canus</i>	<i>A. melanoleuca</i>
100	0.26X	85-15	<i>A.</i> <i>melanoleuca</i> <i>U. ameri-</i> <i>canus</i>	<i>A.</i> <i>melanoleuca</i> <i>U. ameri-</i> <i>canus</i>

1000	2.6X	95-5	<i>U. thibetanus</i> <i>formosanus</i> <i>U. thibetanus</i> <i>thibetanus</i> <i>M. ursinus</i> <i>A.</i> <i>melanoleuca</i>	<i>U. thibetanus</i> <i>formosanus</i> <i>U. thibetanus</i> <i>thibetanus</i> <i>M. ursinus</i> <i>A.</i> <i>melanoleuca</i>
500	1.3X	95-5	<i>U. thibetanus</i> <i>formosanus</i> <i>U. thibetanus</i> <i>thibetanus</i>	<i>U. thibetanus formosanus</i>
250	0.66X	95-5	<i>U. thibetanus</i> <i>formosanus</i> <i>U. thibetanus</i> <i>thibetanus</i>	<i>U. thibetanus formosanus</i>
100	0.26X	95-5	<i>U. thibetanus</i> <i>formosanus</i> <i>U. thibetanus</i> <i>thibetanus</i>	<i>U. thibetanus formosanus</i>
1000	2.6X	95-5	<i>U. spelaeus</i> <i>U. arctos</i>	<i>U. spelaeus</i> <i>U. arctos</i>
500	1.3X	95-5	<i>U. spelaeus</i> <i>U. arctos</i>	<i>U. spelaeus</i> <i>U. arctos</i>
250	0.66X	95-5	<i>U. spelaeus</i> <i>U. arctos</i>	<i>U. spelaeus</i>
100	0.26X	95-5	<i>U. spelaeus</i> <i>U. arctos</i>	<i>U. spelaeus</i>

1000	2.6X	95-5	A. <i>melanoleuca</i> U. <i>ameri-</i> <i>canus</i>	A. <i>melanoleuca</i> U. <i>ameri-</i> <i>canus</i>
500	1.3X	95-5	A. <i>melanoleuca</i> U. <i>ameri-</i> <i>canus</i>	A. <i>melanoleuca</i> U. <i>ameri-</i> <i>canus</i>
250	0.66X	95-5	A. <i>melanoleuca</i> U. <i>ameri-</i> <i>canus</i>	A. <i>melanoleuca</i> U. <i>ameri-</i> <i>canus</i>
100	0.26X	95-5	A. <i>melanoleuca</i> U. <i>ameri-</i> <i>canus</i>	U. <i>americanus</i>
1500	4X	47-33-20	<i>G. tyrrhea</i> <i>G. belina</i> N7	<i>G. tyrrhea</i> <i>G. belina</i> <i>G. maja</i> <i>B. alcinoe</i> <i>G. cytherea</i> <i>N. wahlbergi</i>
750	2X	47-33-20	<i>G. tyrrhea</i> <i>G. belina</i> N7	<i>G. tyrrhea</i> <i>G. belina</i> <i>B. alcinoe</i> <i>G. maja</i> <i>G. cytherea</i>

375	1X	47-33-20	<i>G. tyrrhea</i> <i>G. belina</i> <i>N7</i>	<i>G. tyrrhea</i> <i>G. belina</i> <i>B. alcinoe</i>
150	0.4X	47-33-20	<i>G. tyrrhea</i> <i>G. belina</i> <i>N7</i>	<i>G. tyrrhea</i> <i>G. belina</i>
2000	5.2X	25-25-25-25	<i>P. fasciata</i> <i>P. groen-</i> <i>landica</i> <i>P. vitulina</i> <i>P. largha</i>	<i>P. fasciata</i> <i>P. groen-</i> <i>landica</i> <i>P. vitulina</i> <i>P. largha</i>
1000	2.6X	25-25-25-25	<i>P. fasciata</i> <i>P. groen-</i> <i>landica</i> <i>P. vitulina</i> <i>P. largha</i>	<i>P. fasciata</i> <i>P. groen-</i> <i>landica</i> <i>P. vitulina</i> <i>P. largha</i>
500	1.3X	25-25-25-25	<i>P. fasciata</i> <i>P. groen-</i> <i>landica</i> <i>P. vitulina</i> <i>P. largha</i>	<i>P. fasciata</i> <i>P. groen-</i> <i>landica</i> <i>P. vitulina</i> <i>P. largha</i>
200	0.55X	25-25-25-25	<i>P. fasciata</i> <i>P. groen-</i> <i>landica</i> <i>P. vitulina</i> <i>P. largha</i>	<i>P. fasciata</i> <i>P. groen-</i> <i>landica</i> <i>P. vitulina</i>

Table 1: Results table for the baseline model HAYSTAC for every simulated test.

## 4 Supplementary Figures

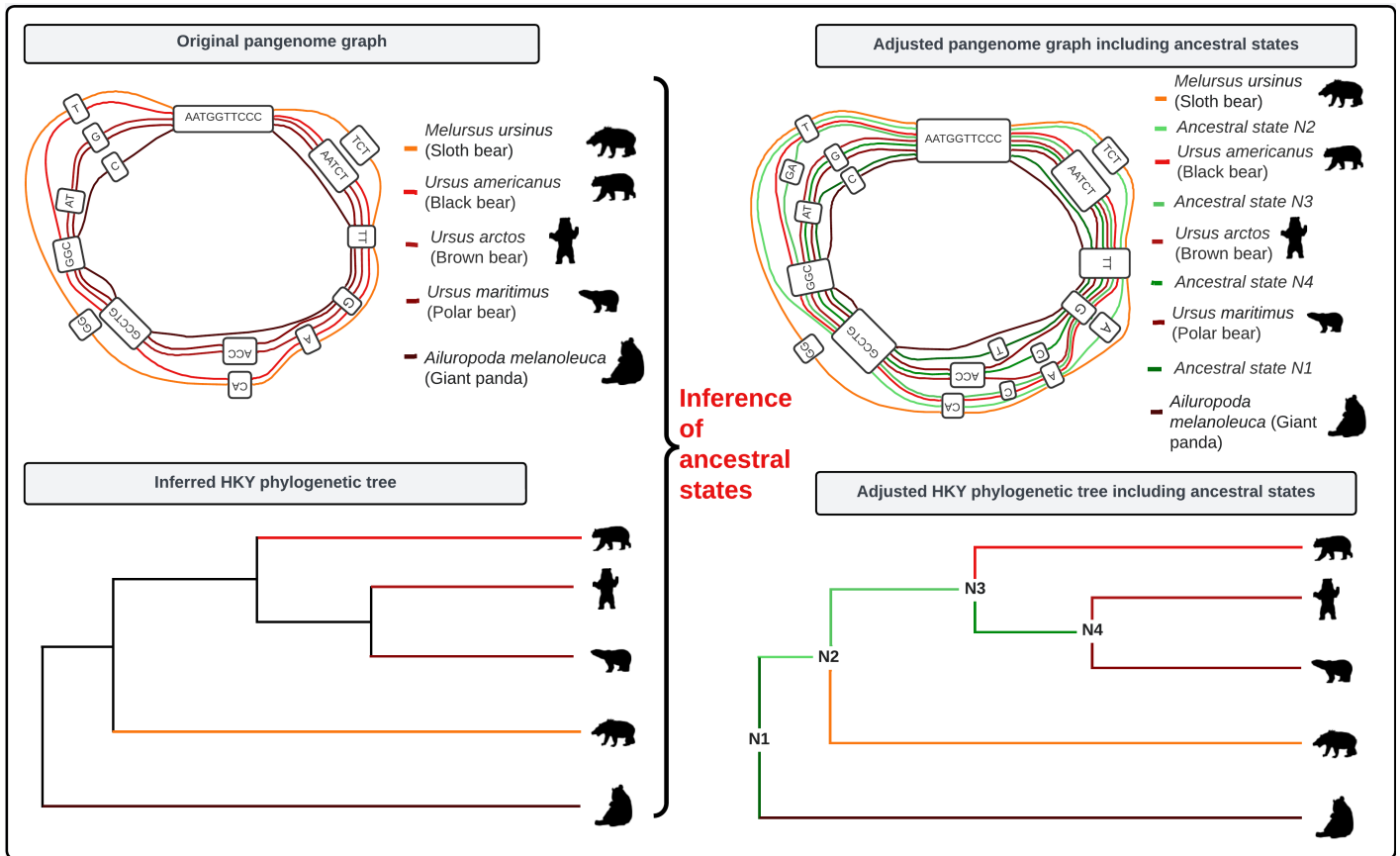


Figure 1: Workflow of a pangenome graph including ancestral states. The figure shows the original pangenome graph, including their phylogenetic tree, and follows with an adjusted pangenome graph, where ancestral state sequences are included and the corresponding phylogenetic tree.

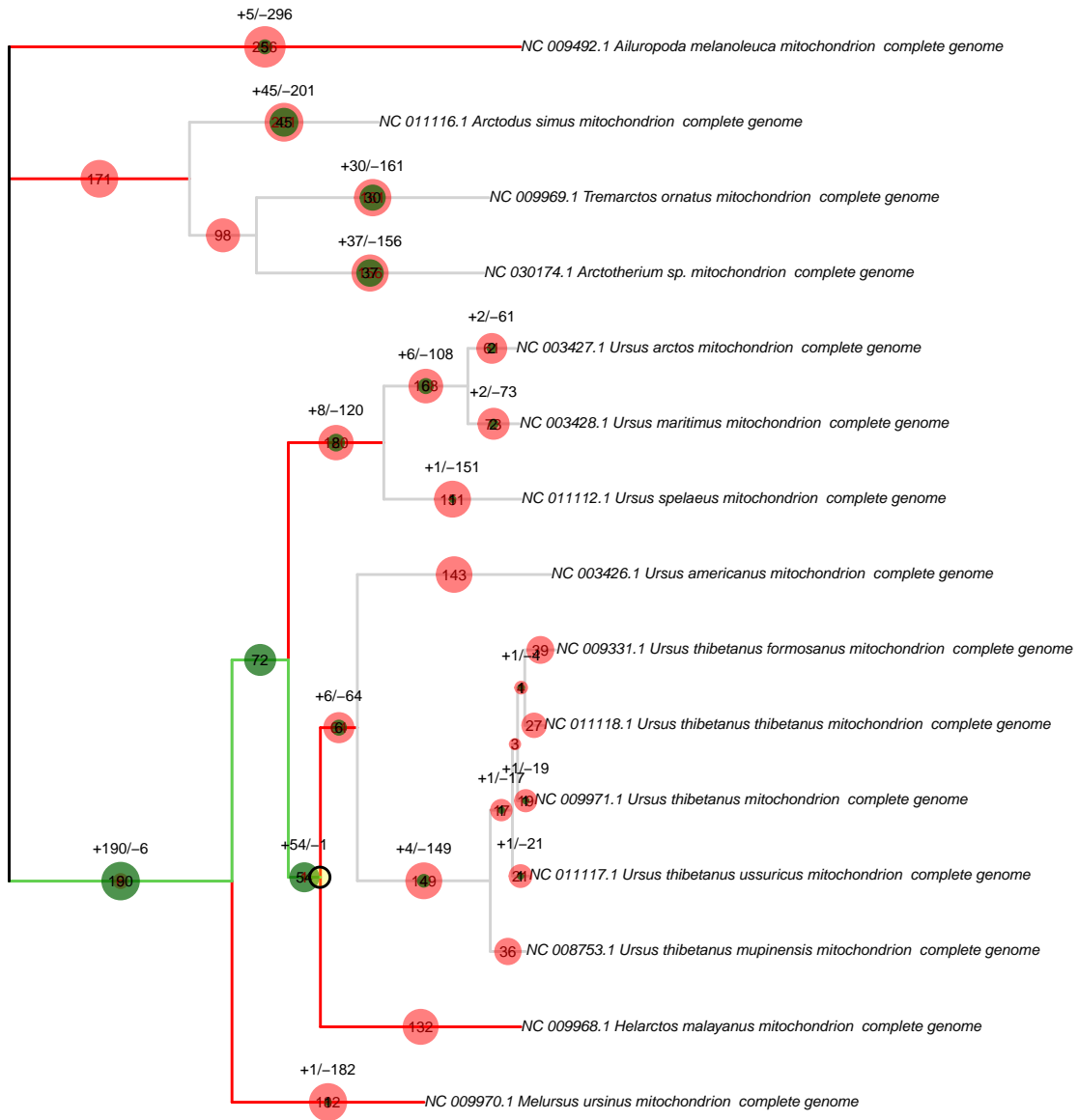


Figure 2: pathPhynder best path method results for the simulated ancient single-source sample N4 at  $\sim 1.3x$  coverage. pathPhynder shows the best path for the given sample, which ends at the correct ancestral node N4.

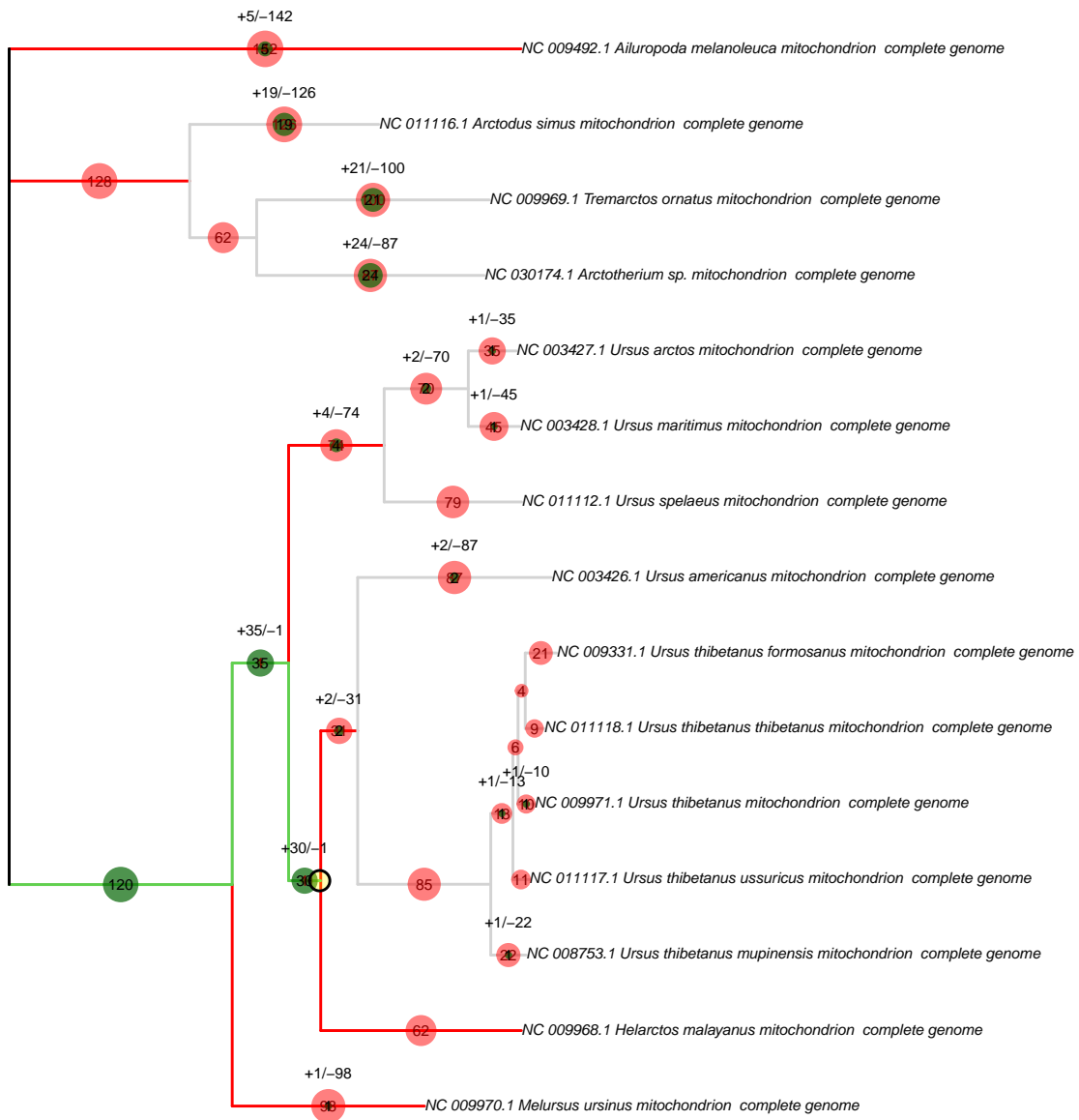


Figure 3: pathPhynder best path method results for the simulated ancient single-source sample N4 at  $\sim 0.66X$  coverage. pathPhynder shows the best path for the given sample, which ends at the correct ancestral node N4.



Figure 4: pathPhynder best path method results for the simulated ancient single-source sample N4 at  $\sim 0.2X$  coverage. pathPhynder shows the best path for the given sample, which ends at the correct ancestral node N4.



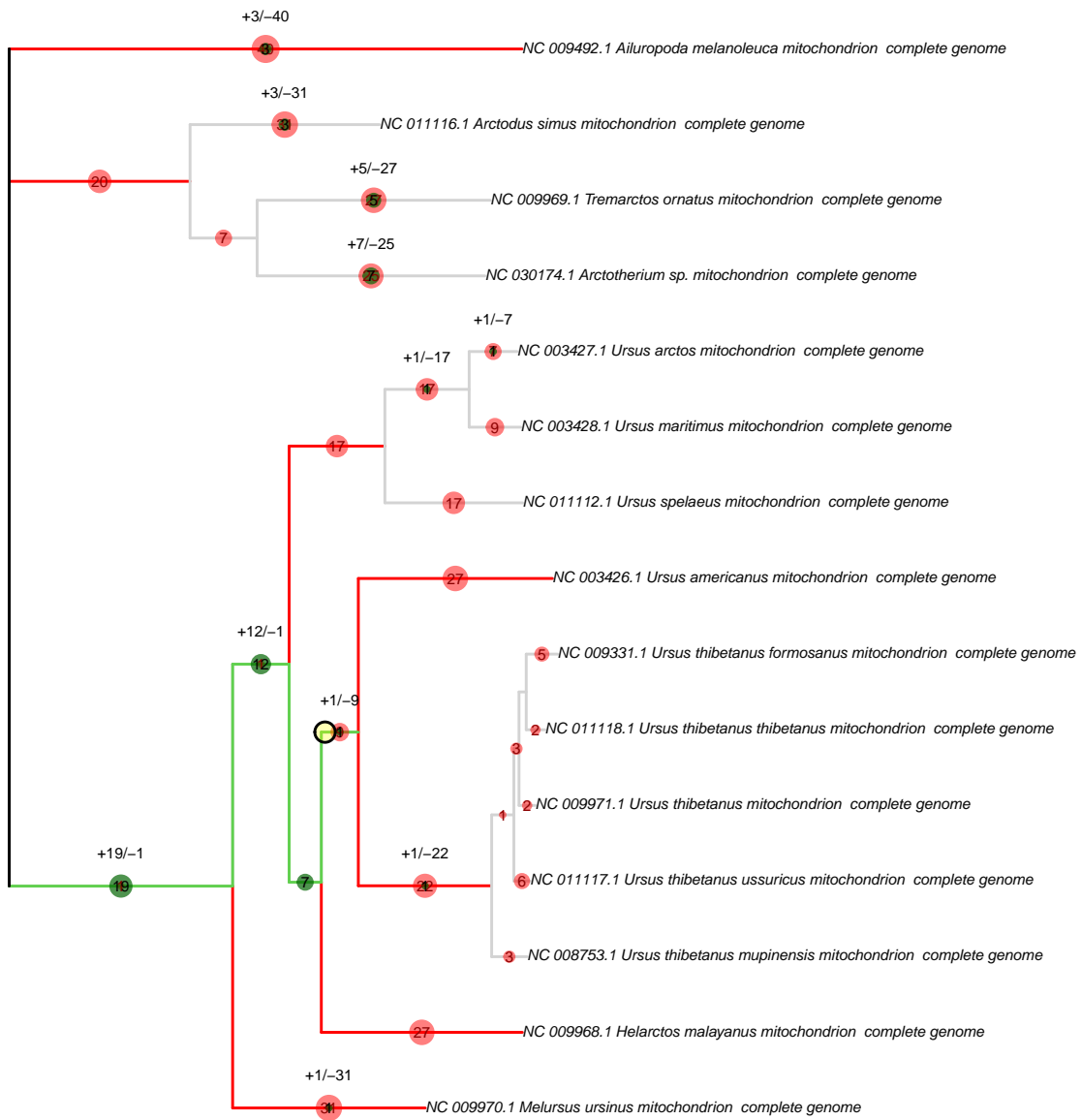


Figure 5: pathPhynder best path method results for the simulated ancient single-source sample N4 at  $\sim 0.13X$  coverage. pathPhynder shows the best path for the given sample, which ends at one ancestral node after the target one N4 at N5.

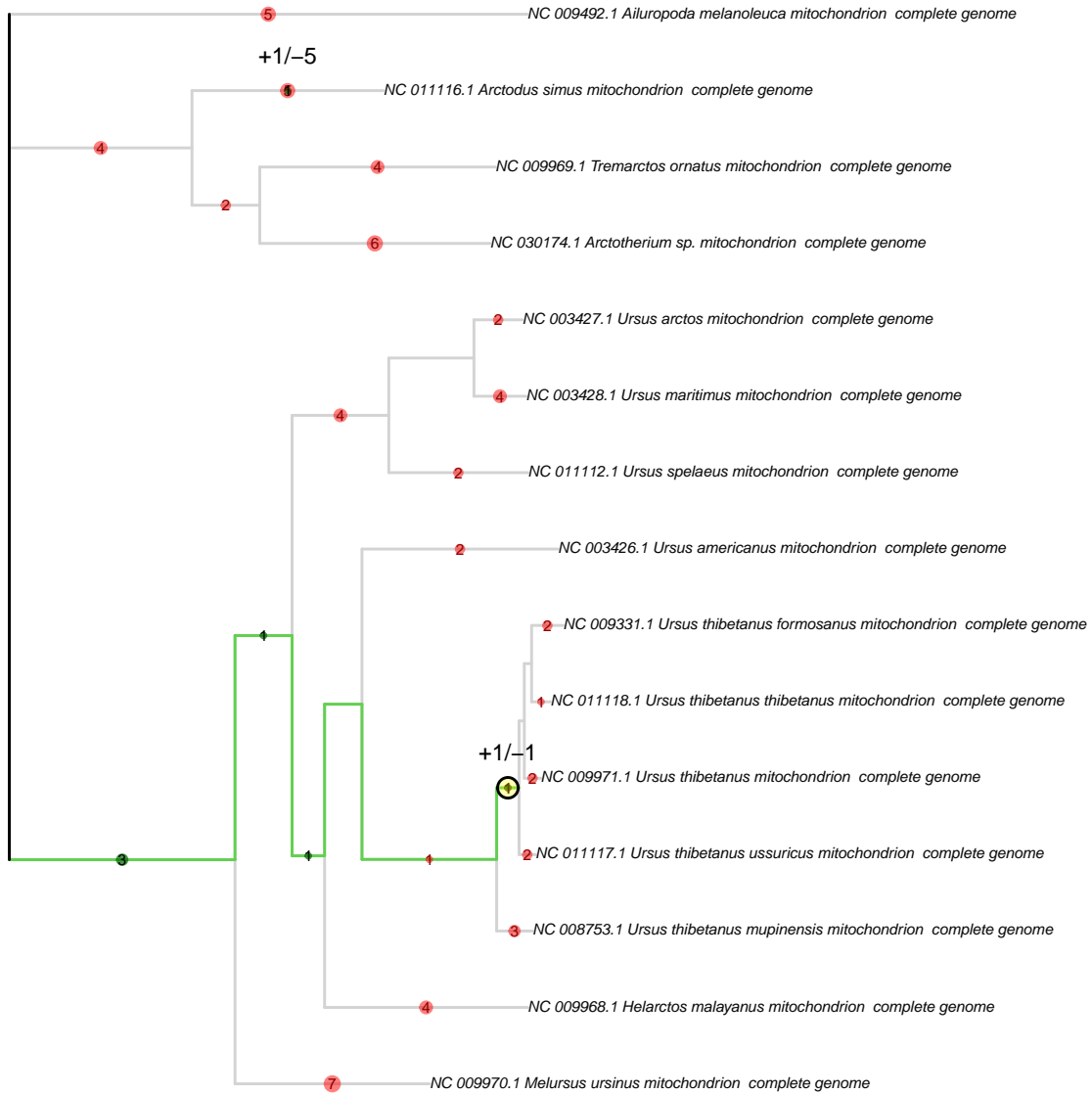


Figure 6: pathPhynder best path method results for the simulated ancient single-source sample N4 at  $\sim 0.026X$  coverage. pathPhynder shows the best path for the given sample, which ends at the ancestral node N6.

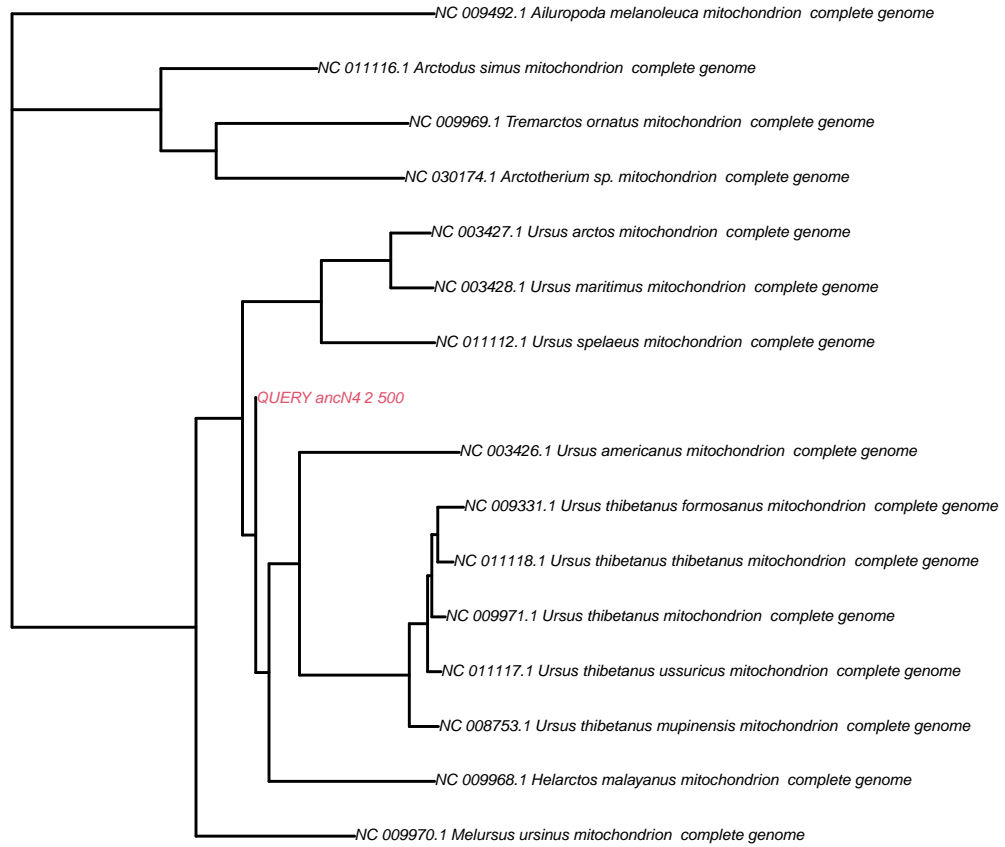


Figure 7: pathPhynder maximum likelihood method results for the simulated ancient single-source sample N4 at  $\sim 1.3x$  coverage. pathPhynder shows the best path for the given sample, which ends at the correct ancestral node N4.

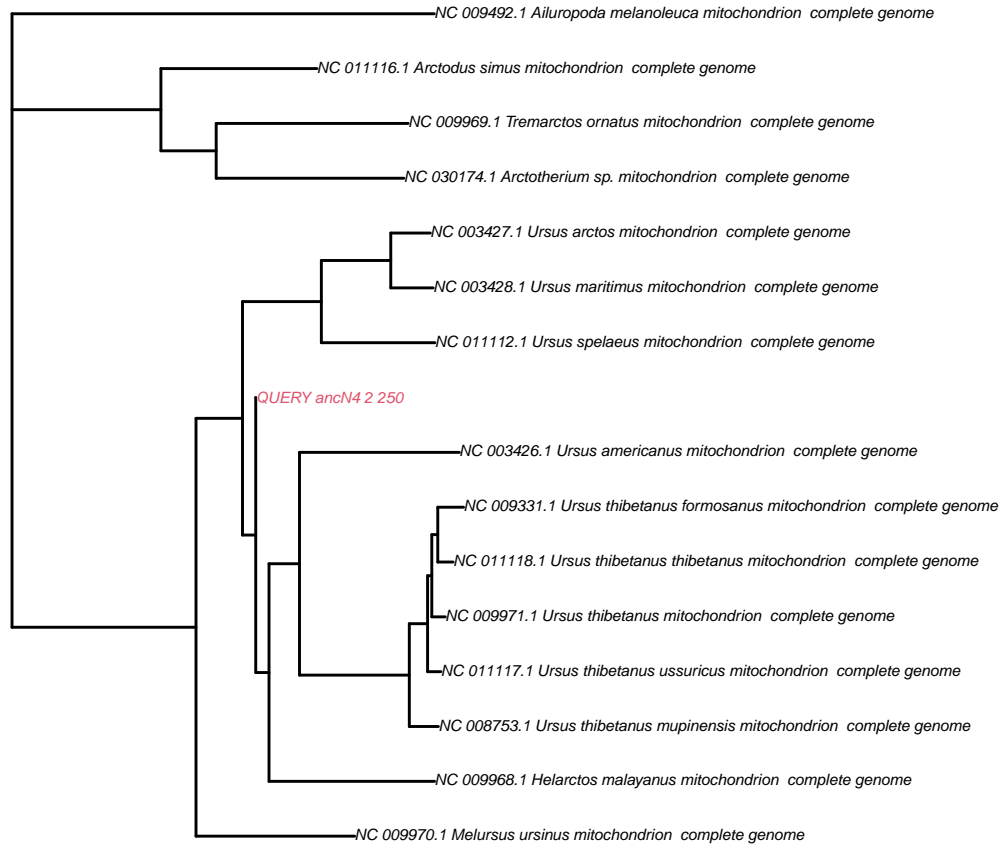


Figure 8: pathPhynder maximum likelihood method results for the simulated ancient single-source sample N4 at  $\sim 0.66X$  coverage. pathPhynder shows the best path for the given sample, which ends at the correct ancestral node N4.

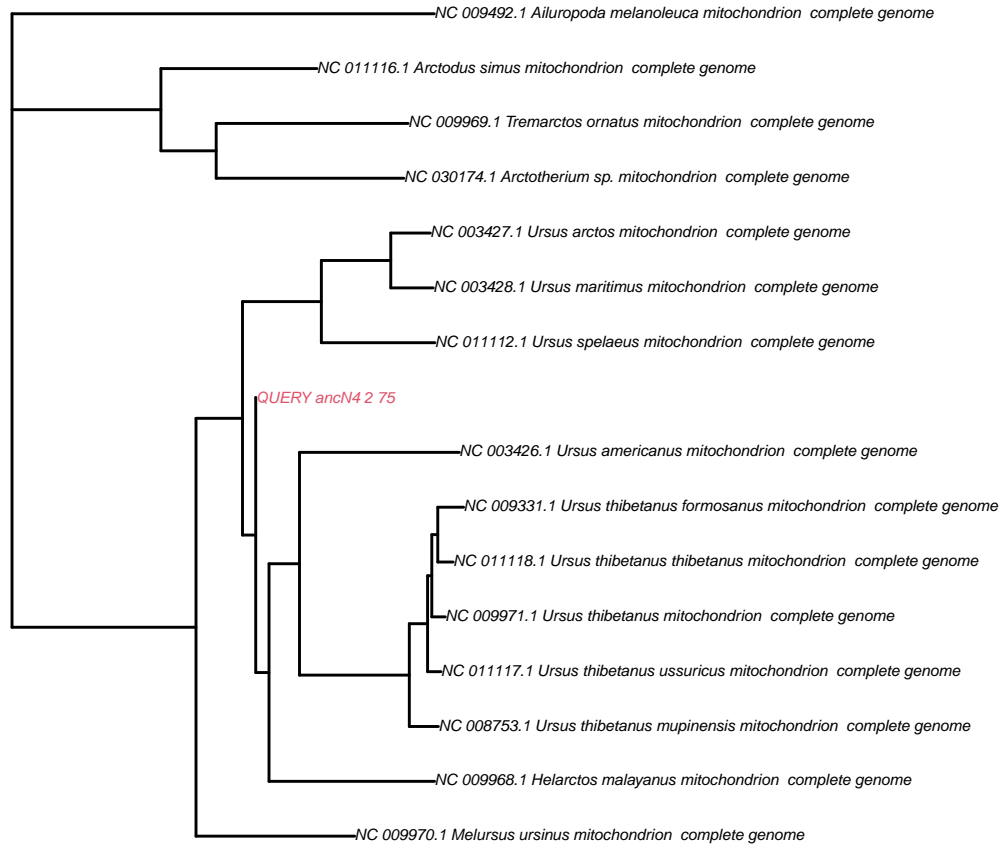


Figure 9: pathPhynder maximum likelihood method results for the simulated ancient single-source sample N4 at  $\sim 0.2X$  coverage. pathPhynder shows the best path for the given sample, which ends at the correct ancestral node N4.

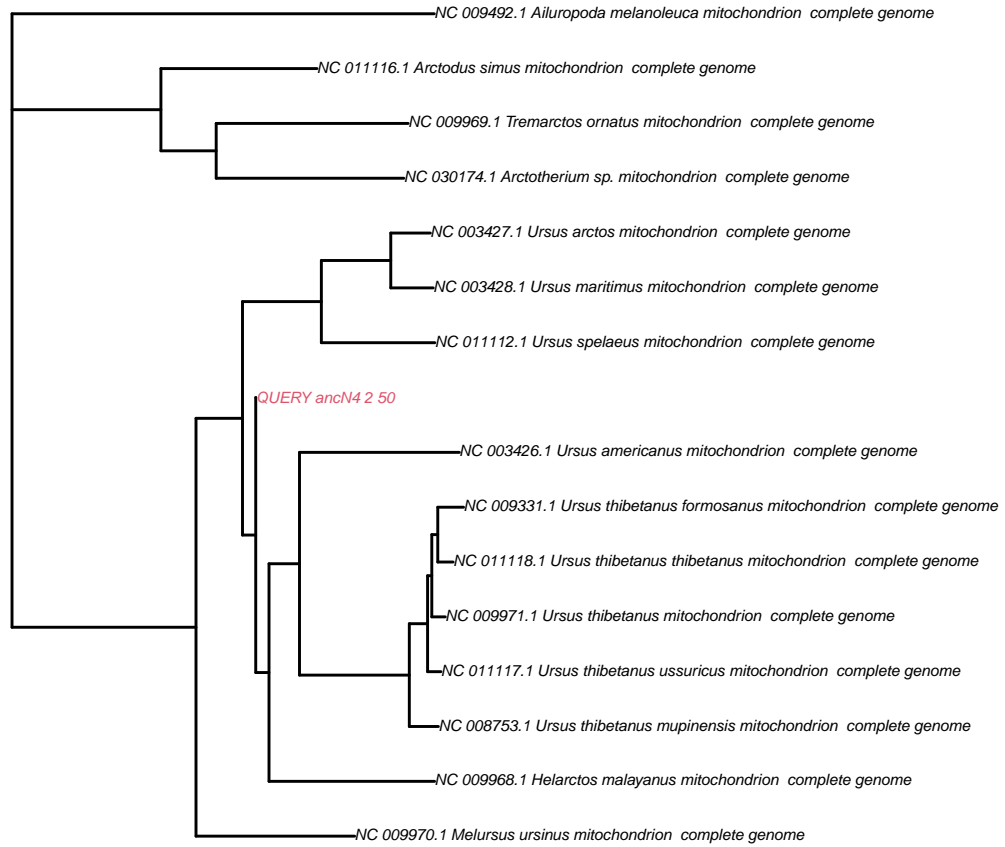


Figure 10: pathPhynder maximum likelihood method results for the simulated ancient single-source sample N4 at  $\sim 0.13X$  coverage. pathPhynder shows the best path for the given sample, which ends at the correct ancestral node N4.

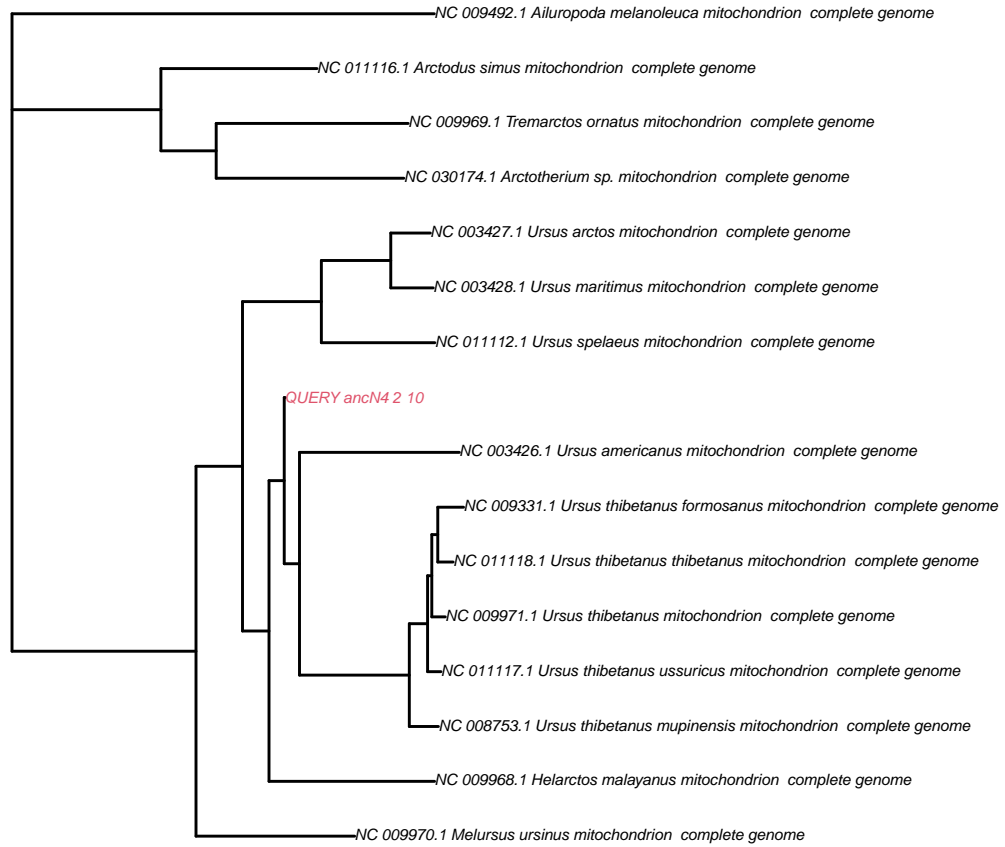


Figure 11: **pathPhynder** maximum likelihood method results for the simulated ancient single-source sample N4 at  $\sim 0.026X$  coverage. **pathPhynder** shows the best path for the given sample, which ends at the ancestral node N6.





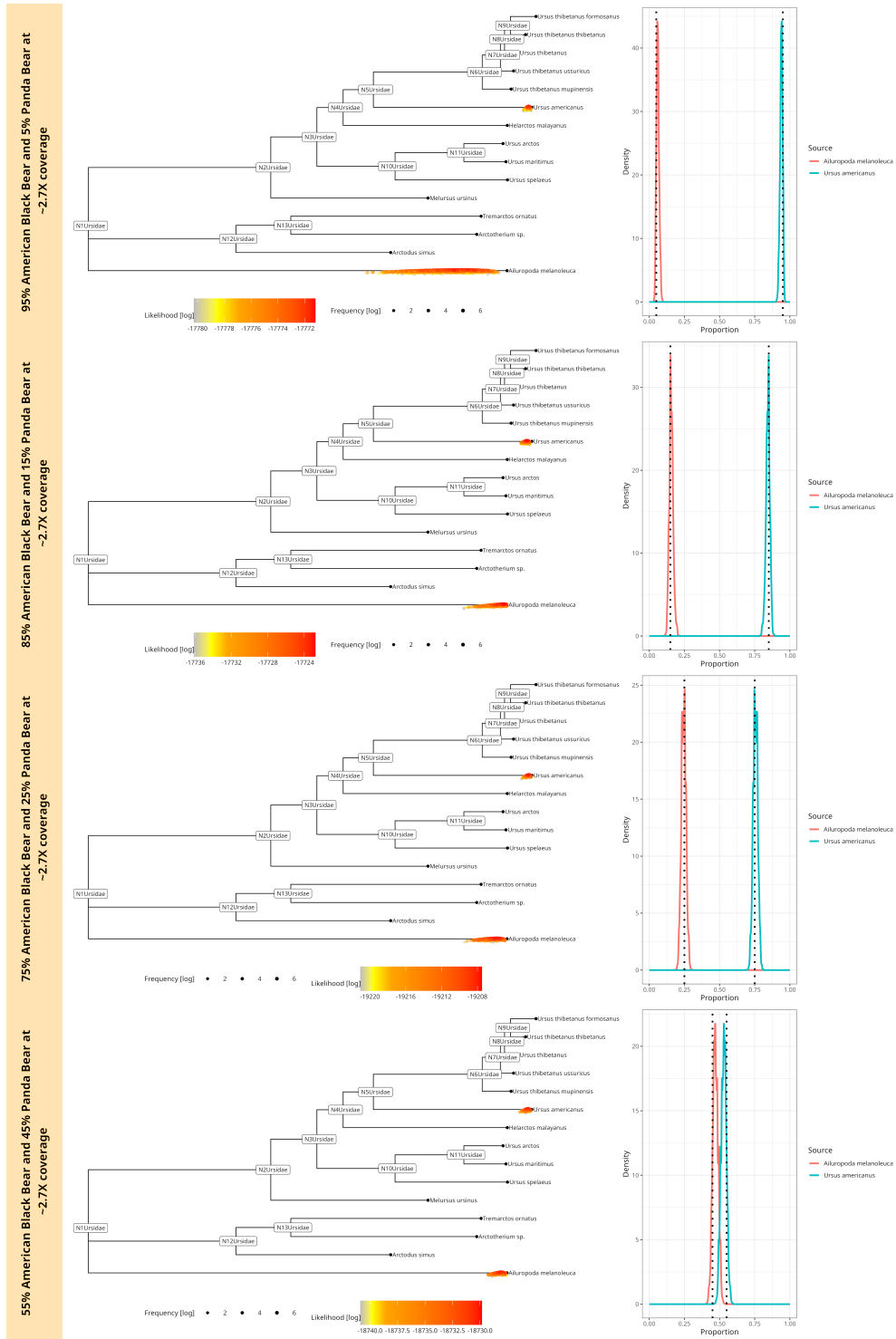


Figure 12: soibean results for simulated ancient fragments of two diverged species (83.4% similarity) from the family bears (Ursidae) at  $\sim 2.7X$  coverage. The plot shows four different mixtures at 55% – 45%, 75% – 25%, 85% – 15% and 95% – 5% of the American Black bear and Giant Panda bear. The corresponding phylogenetic trees are displayed on the left: we plotted every accepted MCMC move coloured by likelihood value on the tree. The accepted moves are positioned above or below the tree, corresponding to a higher or lower likelihood value than the median, respectively. Each neighbouring plot shows the posterior proportion distribution, including the simulated proportion with a black dotted line.



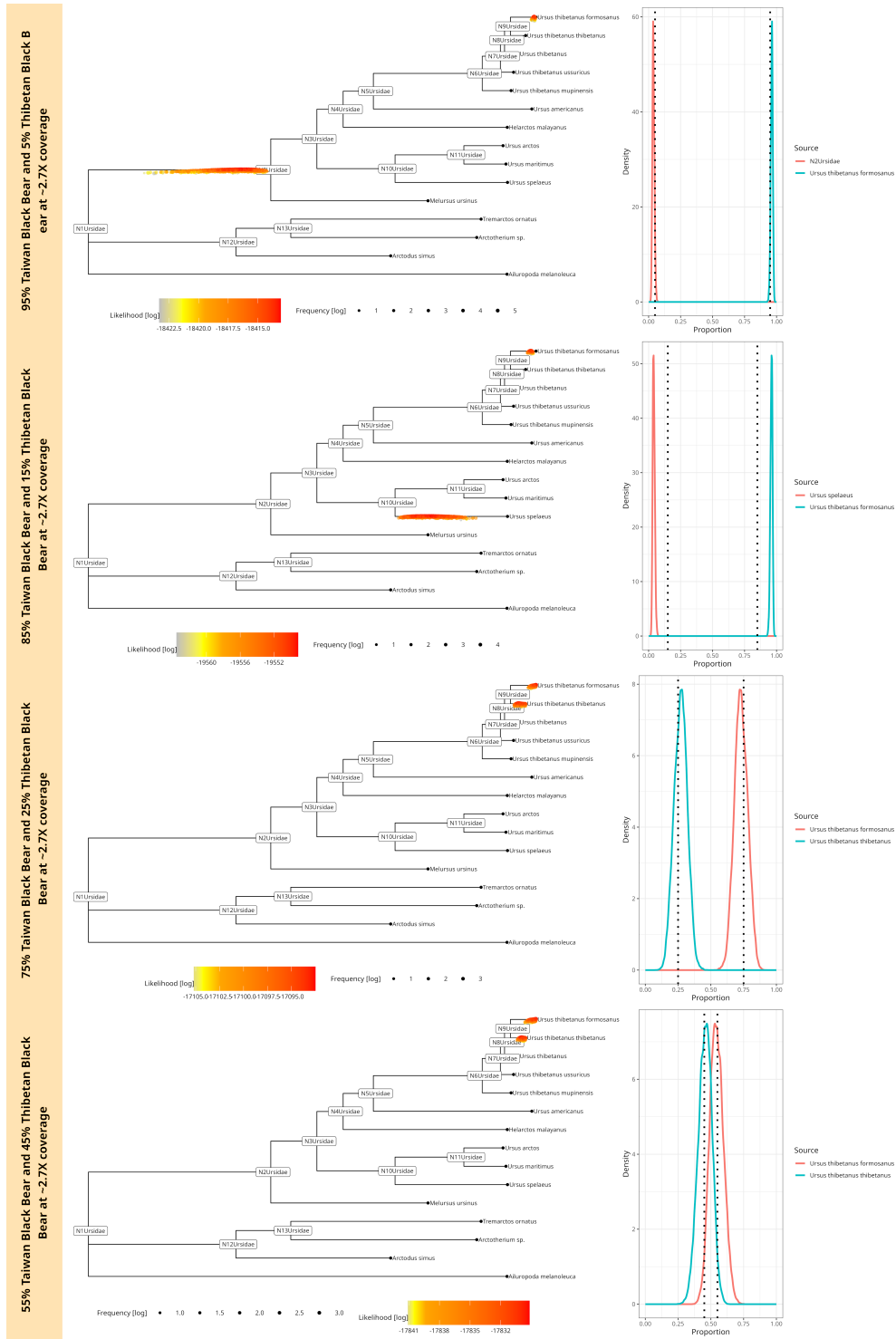


Figure 13: soibean results for simulated ancient fragments of two closely related species (98.8% similarity) from the family bears (Ursidae) at  $\sim 2.7X$  coverage. The plot shows four different mixtures at 55% – 45%, 75% – 25%, 85% – 15% and 95% – 5% of the Taiwan and Tibetan Black bear. The corresponding phylogenetic trees are displayed on the left: we plotted every accepted MCMC move coloured by likelihood value on the tree. The accepted moves are positioned above or below the tree, corresponding to a higher or lower likelihood value than the median, respectively. Each neighbouring plot shows the posterior proportion distribution, including the simulated proportion with a black dotted line



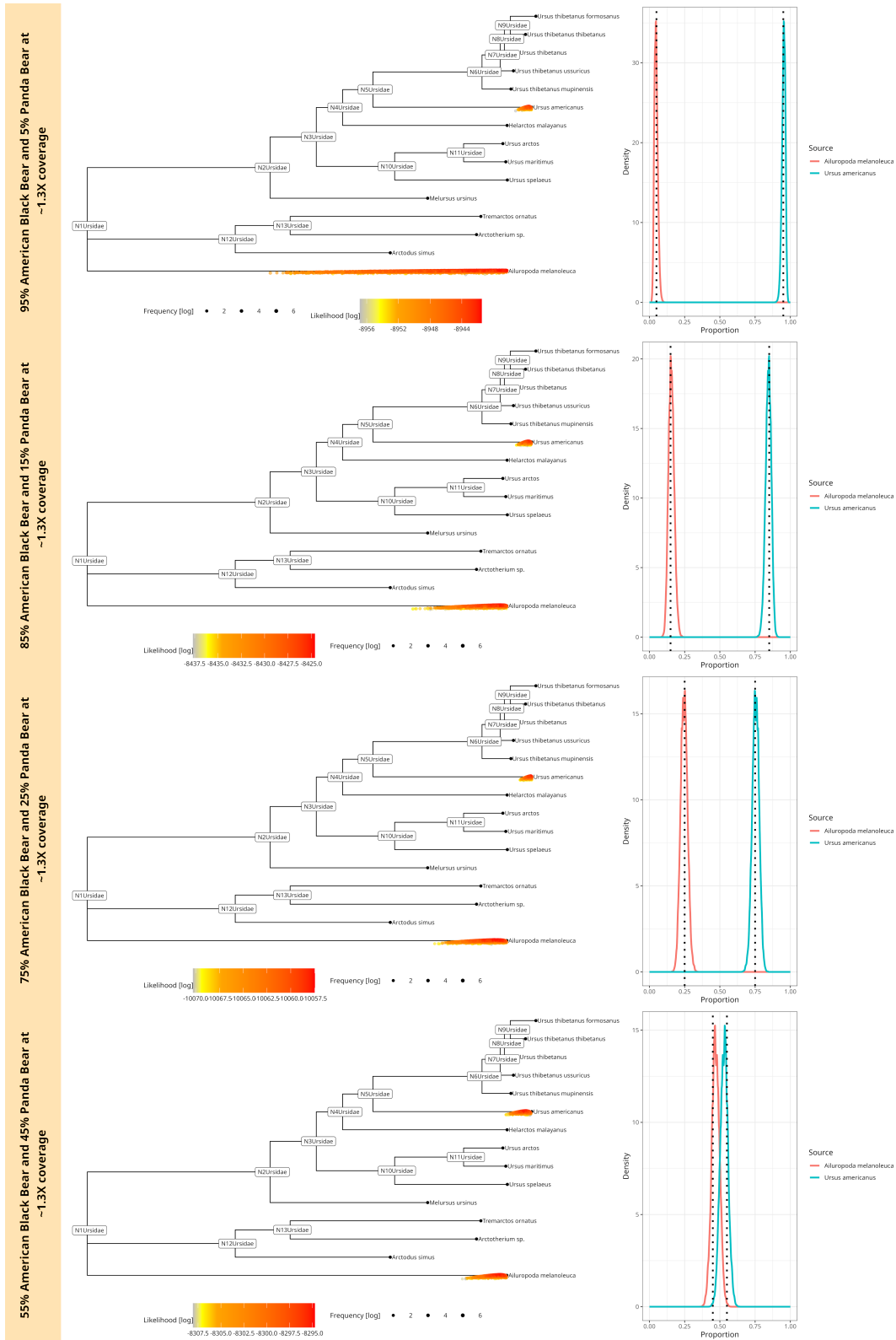


Figure 14: Mixtures of an American Black bear and a Giant Panda bear downsampled to  $\sim 1.3X$  coverage (500 aDNA fragments). The plot shows four different mixtures at 55% – 45%, 75% – 25%, 85% – 15% and 95% – 5%. The phylogenetic tree has a coloured point for each accepted MCMC move. The colour corresponds to the log-likelihood value. The posterior proportion distribution, including the simulated true proportion (black dotted line), is plotted on the right.



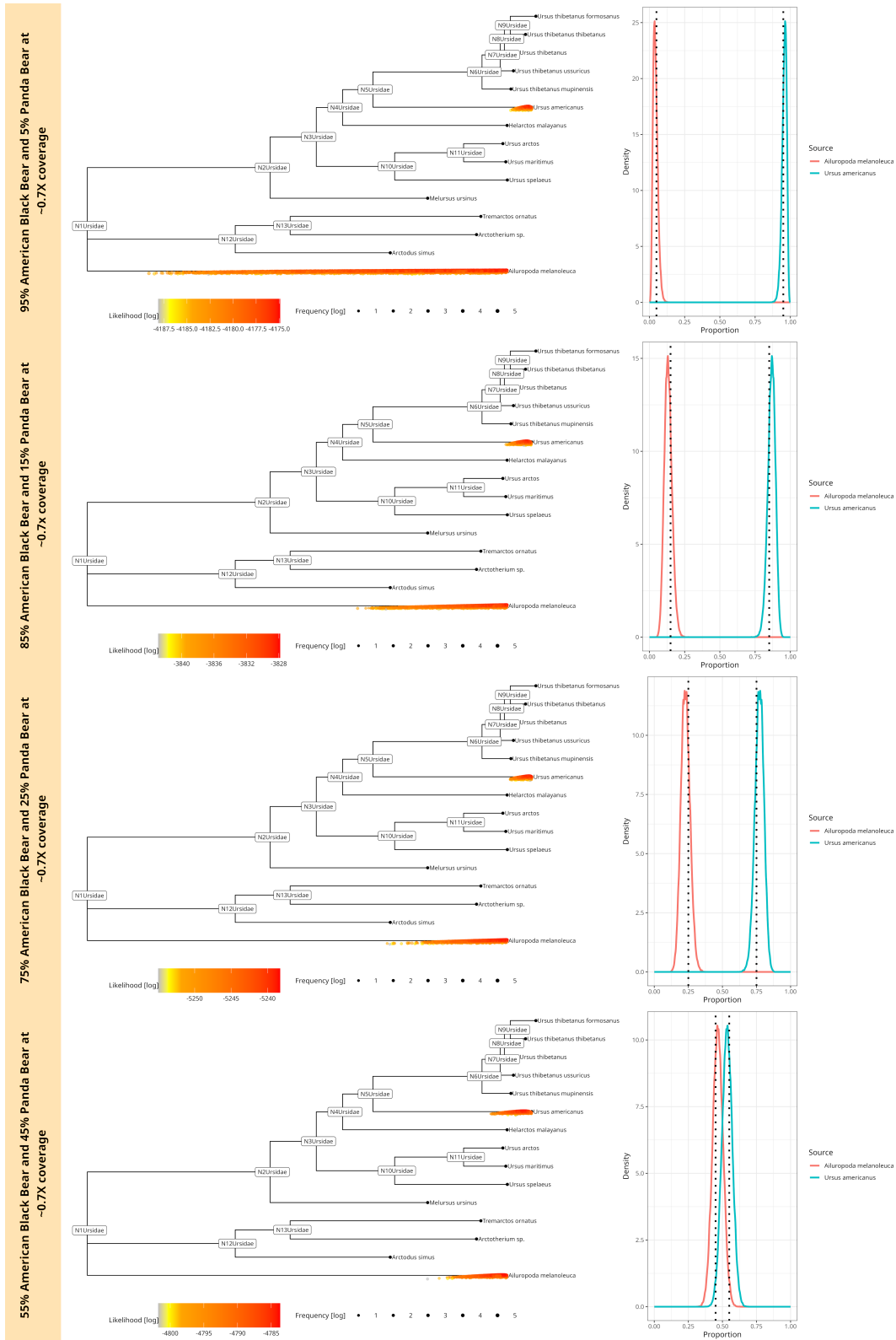


Figure 15: Mixtures of an American Black bear and a Giant Panda bear downsampled to  $\sim 0.7x$  coverage (250 aDNA fragments). The plot shows four different mixtures at 55% – 45%, 75% – 25%, 85% – 15% and 95% – 5%. The phylogenetic tree has a coloured point for each accepted MCMC move. The colour corresponds to the log-likelihood value. The posterior proportion distribution, including the simulated true proportion (black dotted line), is plotted on the right.





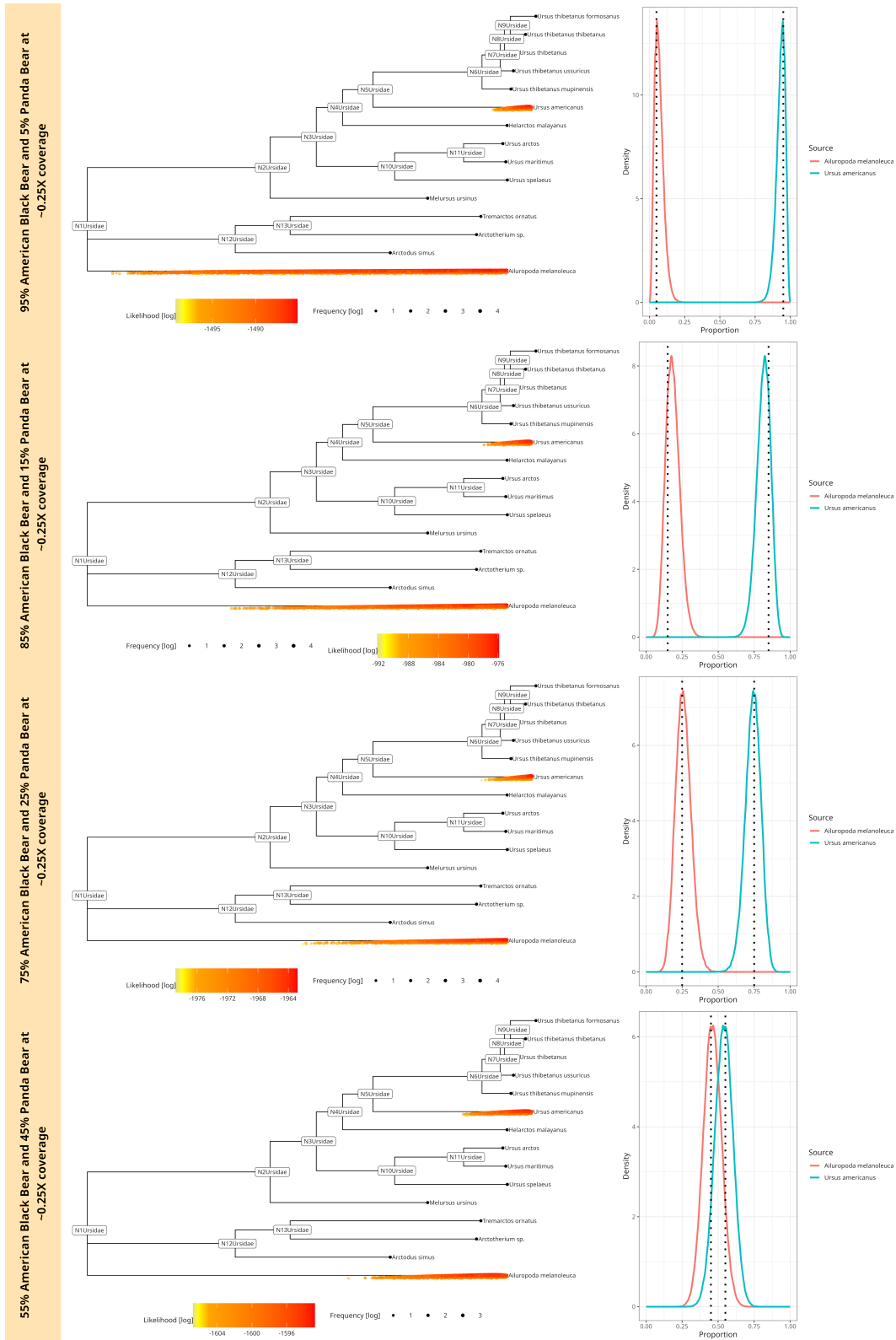


Figure 16: Mixtures of an American Black bear and a Giant Panda bear downsampled to  $\sim 0.25x$  coverage (100 aDNA fragments). The plot shows four different mixtures at 55% – 45%, 75% – 25%, 85% – 15% and 95% – 5%. The phylogenetic tree has a coloured point for each accepted MCMC move. The colour corresponds to the log-likelihood value. The posterior proportion distribution, including the simulated true proportion (black dotted line), is plotted on the right.



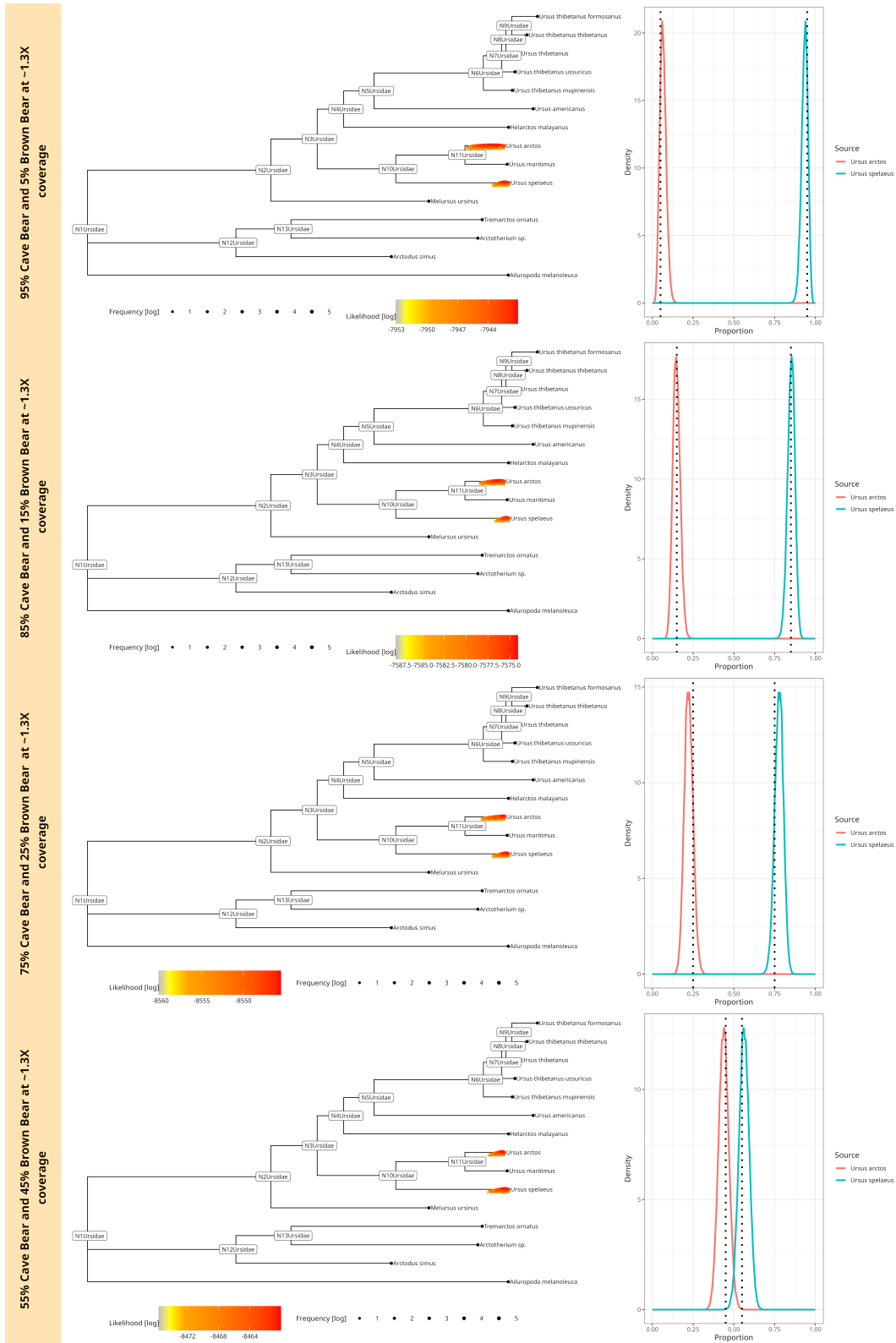


Figure 17: Cave bear and Brown bear mixture downsampled to  $\sim 1.3X$  coverage (500 aDNA fragments). The plot shows four different mixtures at 55% – 45%, 75% – 25%, 85% – 15% and 95% – 5%. The phylogenetic tree has a coloured point for each accepted MCMC move. The colour corresponds to the log-likelihood value. The posterior proportion distribution, including the simulated true proportion (black dotted line), is plotted on the right.



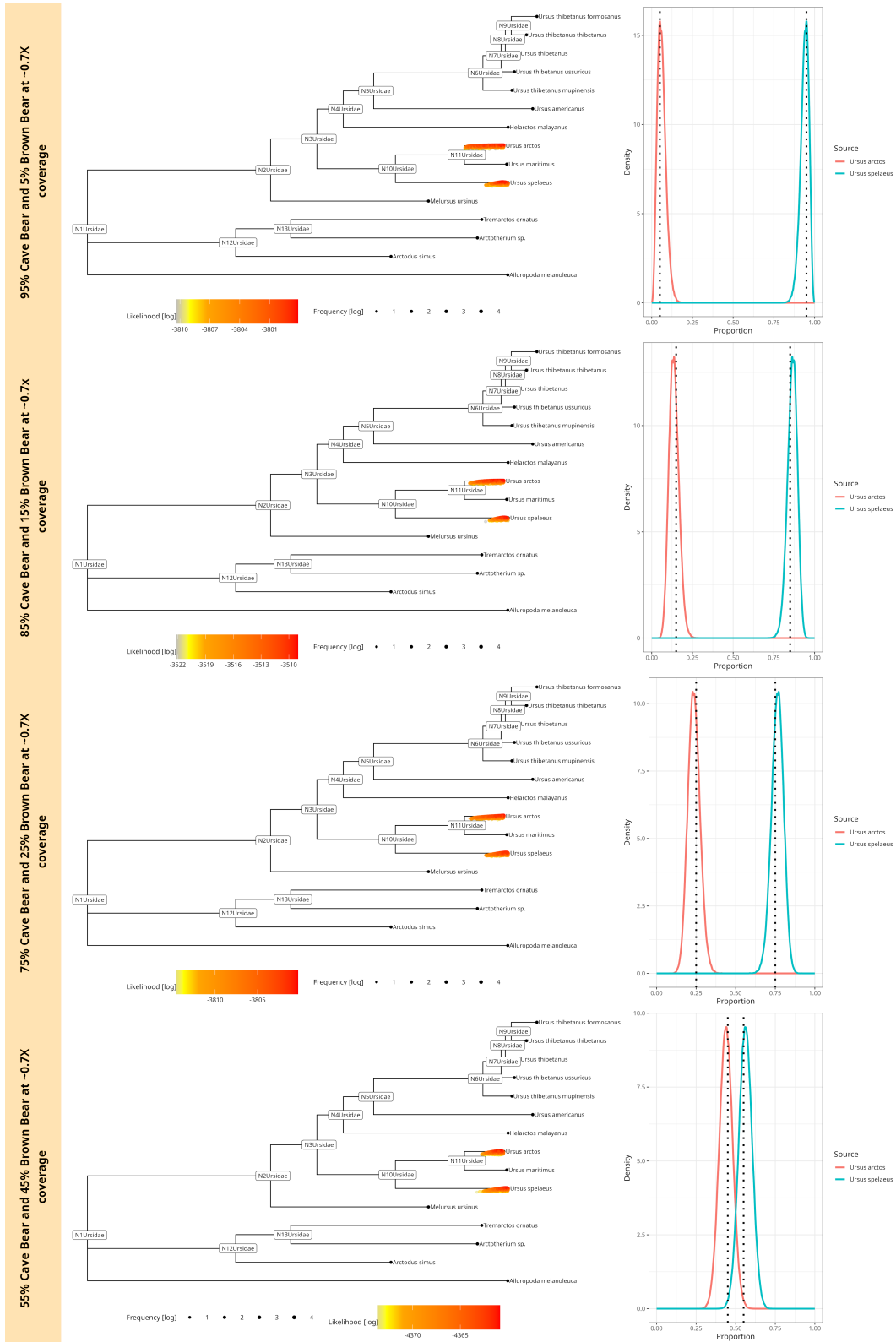


Figure 18: Cave bear and Brown bear mixture downsampled to  $\sim 0.7X$  coverage (250 aDNA fragments). The plot shows four different mixtures at 55% – 45%, 75% – 25%, 85% – 15% and 95% – 5%. The phylogenetic tree has a coloured point for each accepted MCMC move. The colour corresponds to the log-likelihood value. The posterior proportion distribution, including the simulated true proportion (black dotted line), is plotted on the right.



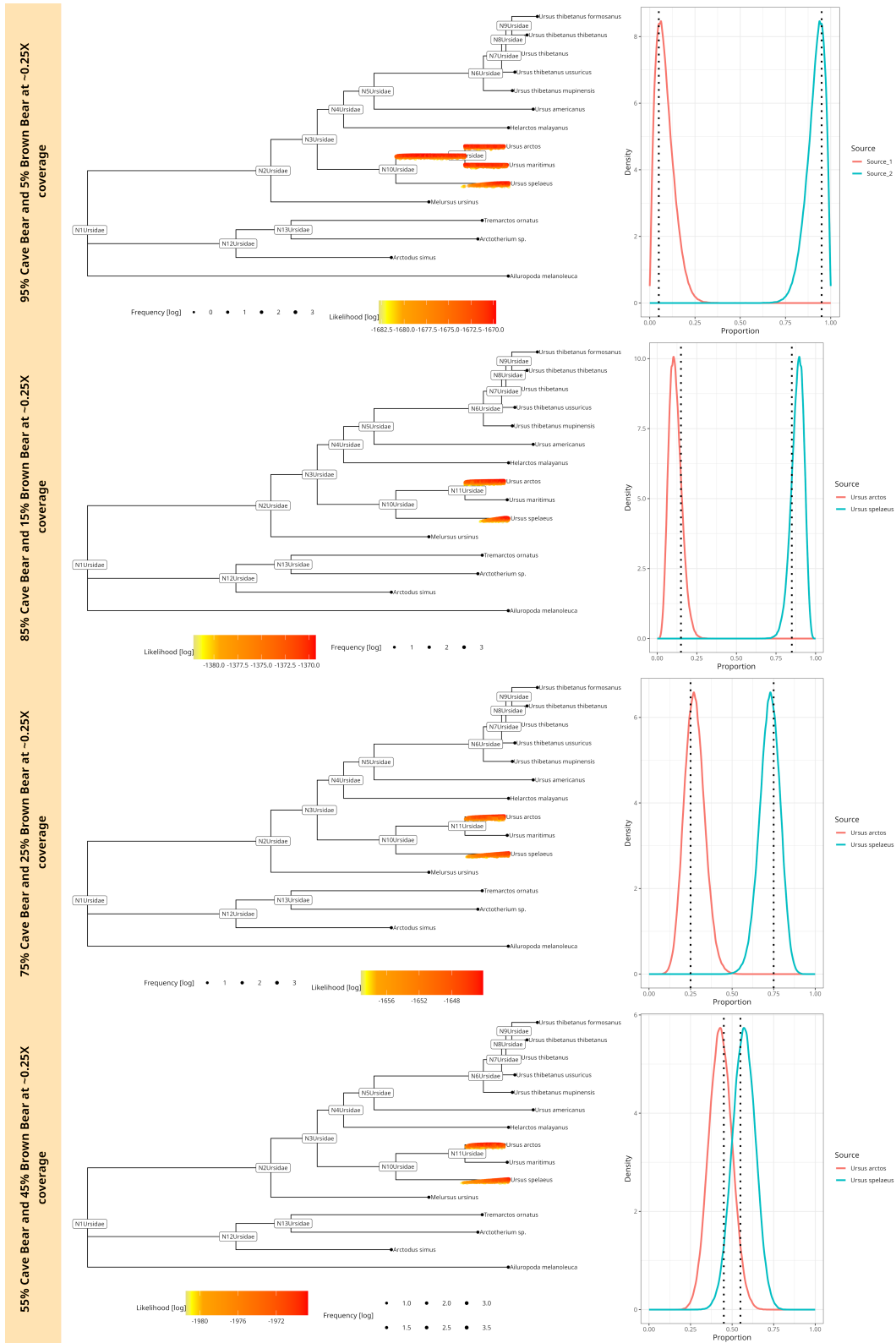


Figure 19: Mixtures of a Cave bear and a Brown bear downsampled to  $\sim 0.25X$  coverage (100 aDNA fragments). The plot shows four different mixtures at 55% – 45%, 75% – 25%, 85% – 15% and 95% – 5%. The phylogenetic tree has a coloured point for each accepted MCMC move. The colour corresponds to the log-likelihood value. The posterior proportion distribution, including the simulated true proportion (black dotted line), is plotted on the right.





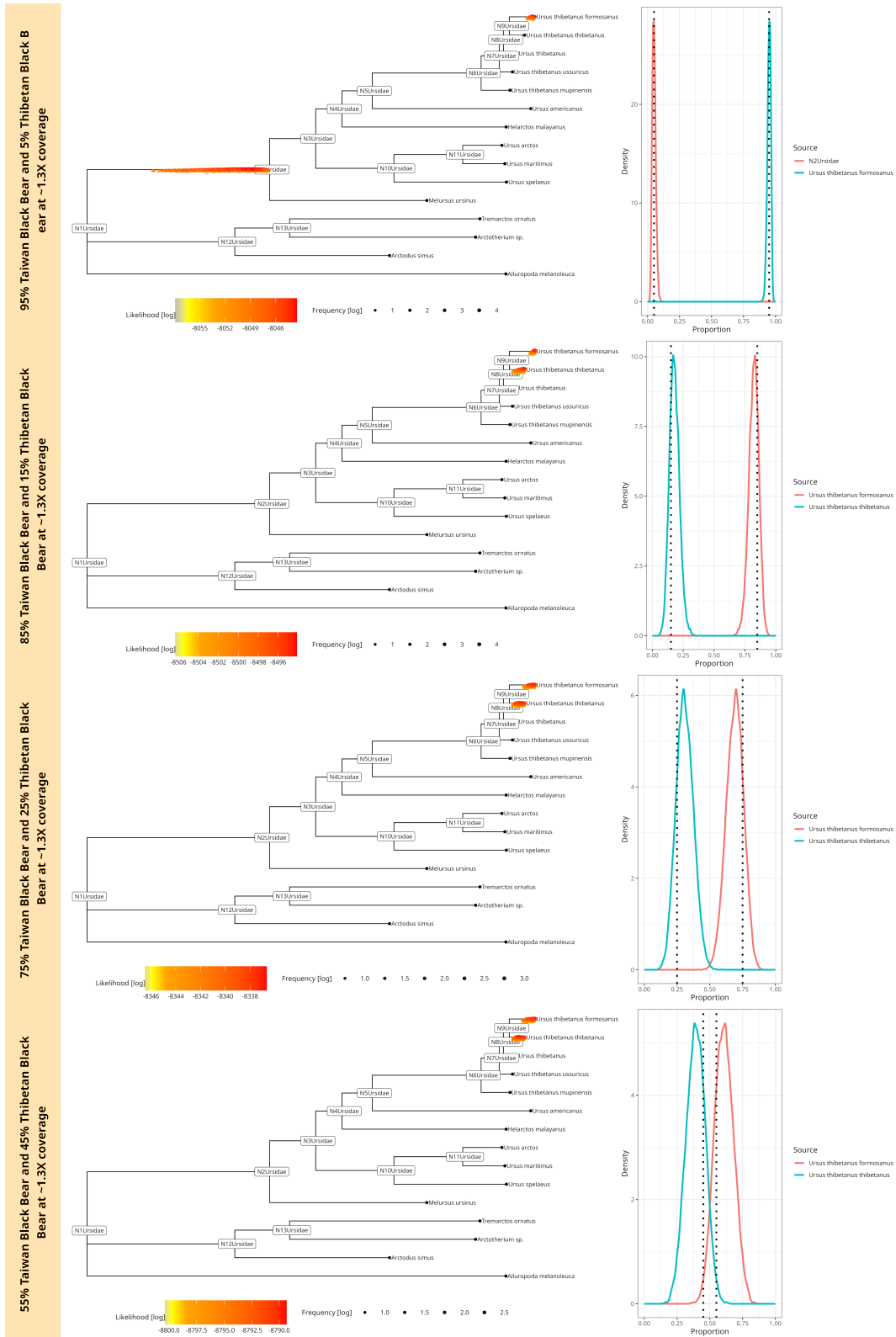


Figure 20: Mixtures of a Taiwan and Tibetan Black bear downsampled to  $\sim 1.3X$  coverage (500 aDNA fragments). The plot shows four different mixtures at 55% – 45%, 75% – 25%, 85% – 15% and 95% – 5%. The phylogenetic tree has a coloured point for each accepted MCMC move. The colour corresponds to the log-likelihood value. The posterior proportion distribution, including the simulated true proportion (black dotted line), is plotted on the right.



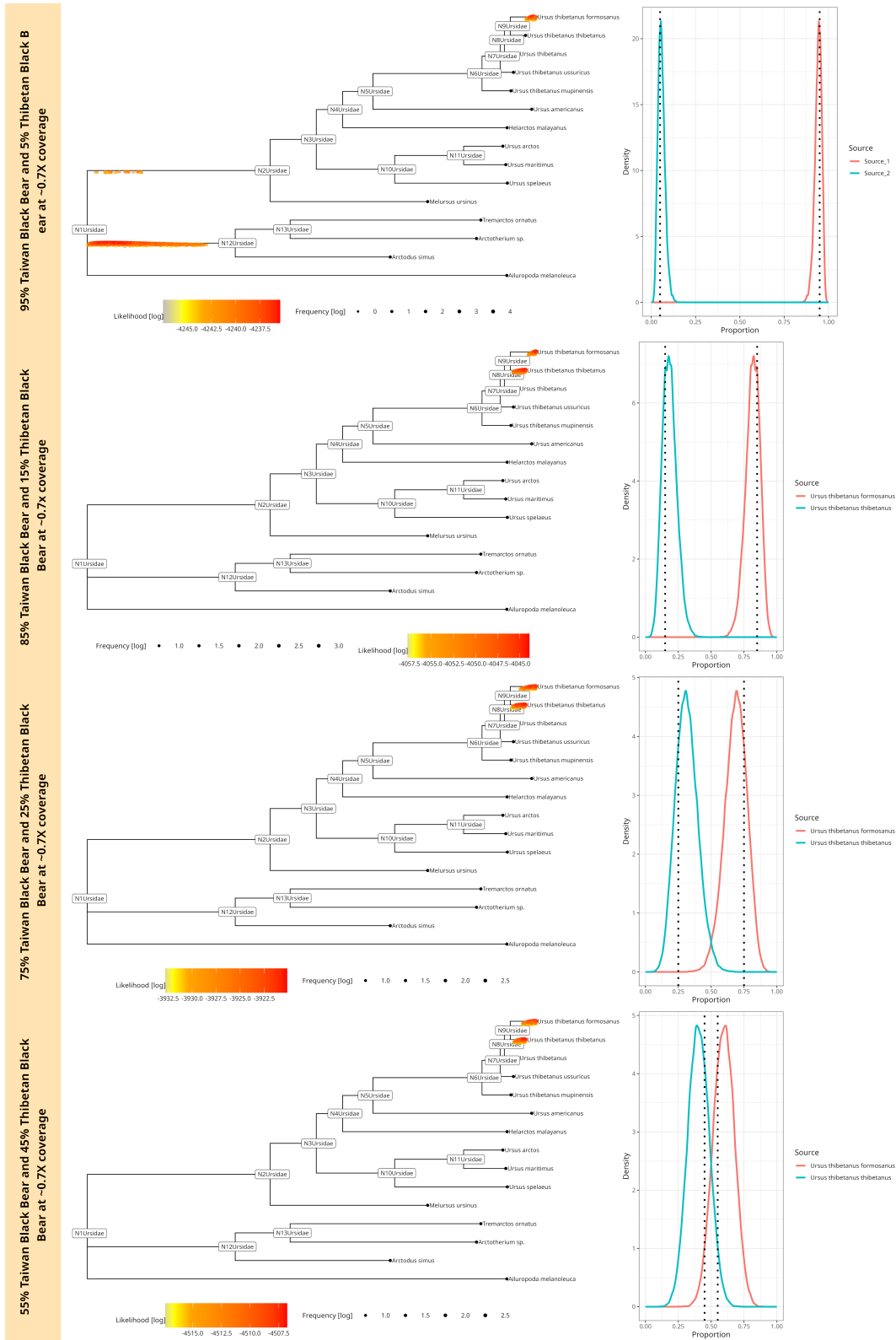


Figure 21: Mixtures of a Taiwan and Tibetan Black bear downsampled to  $\sim 0.7X$  coverage (250 aDNA fragments). The plot shows four different mixtures at 55% – 45%, 75% – 25%, 85% – 15% and 95% – 5%. The phylogenetic tree has a coloured point for each accepted MCMC move. The colour corresponds to the log-likelihood value. The posterior proportion distribution, including the simulated true proportion (black dotted line), is plotted on the right.



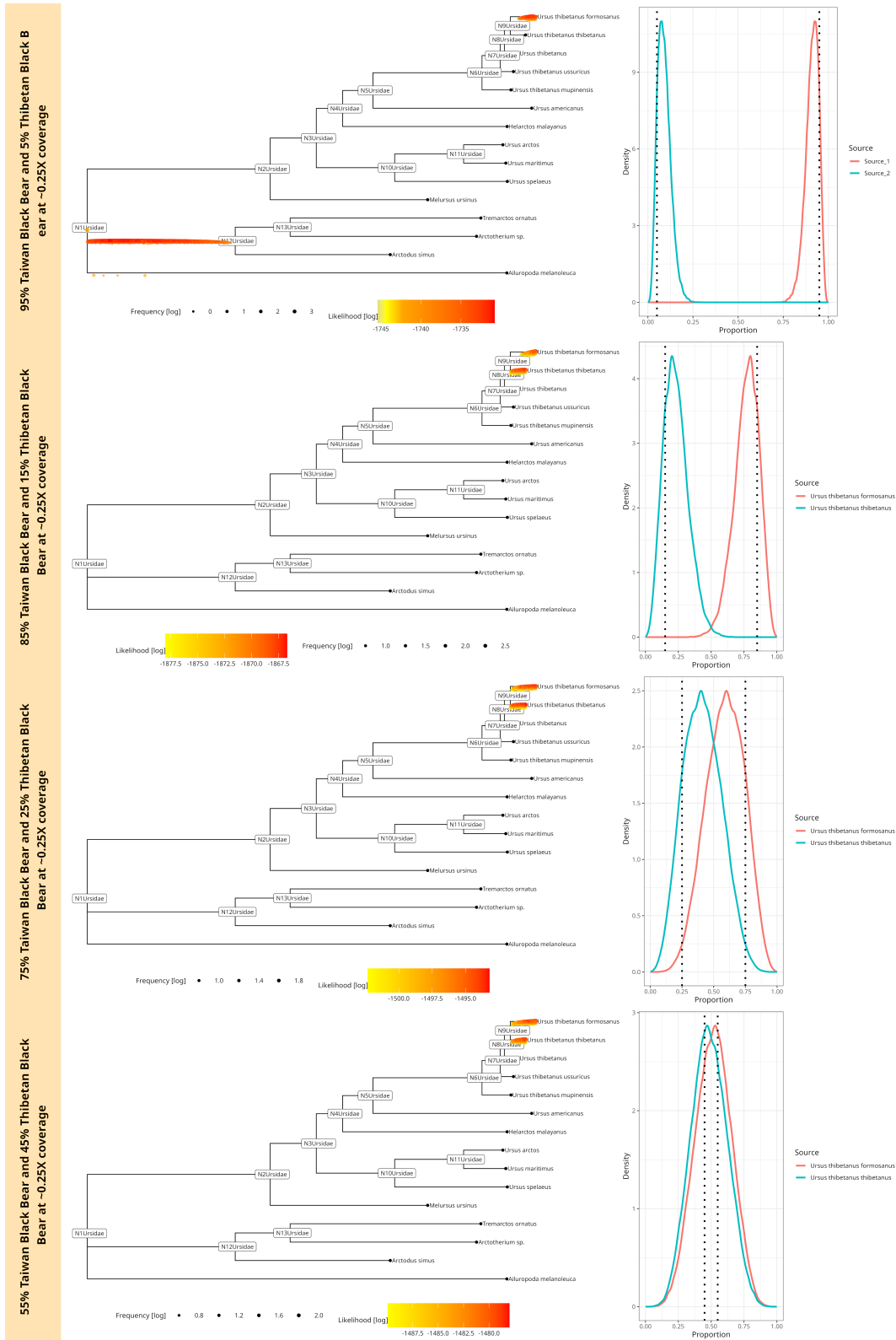


Figure 22: Mixtures of a Taiwan and Tibetan Black bear downsampled to  $\sim 0.25X$  coverage (100 aDNA fragments). The plot shows four different mixtures at 55% – 45%, 75% – 25%, 85% – 15% and 95% – 5%. The phylogenetic tree has a coloured point for each accepted MCMC move. The colour corresponds to the log-likelihood value. The posterior proportion distribution, including the simulated true proportion (black dotted line), is plotted on the right.

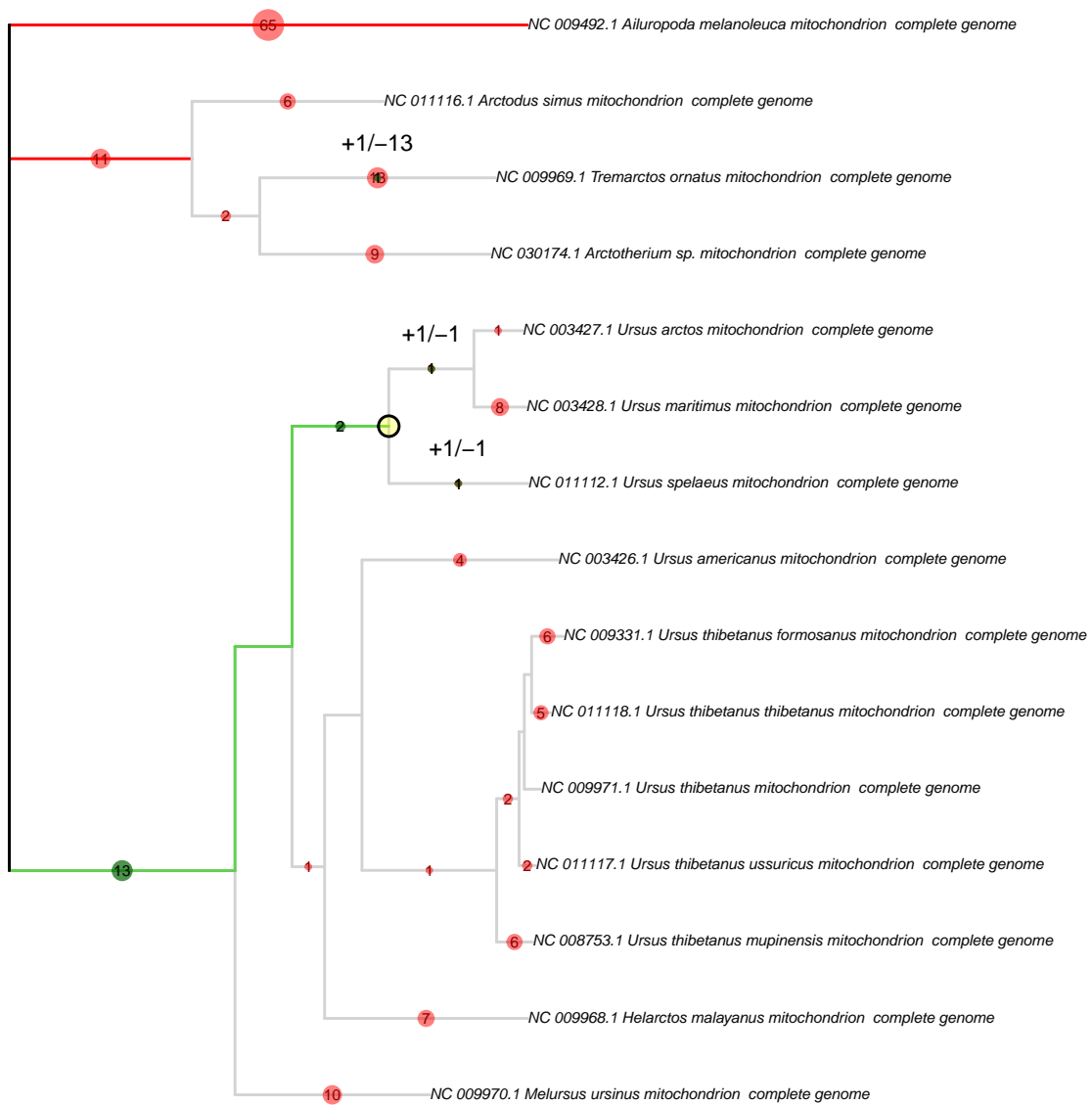


Figure 23: pathPhynder results for the simulated ancient mixture of a cave bear (*Ursus spelaeus*) 55% and the brown bear (*Ursus arctos*) 45% at  $\sim 2.7X$  coverage. pathPhynder shows the best path for the given sample, which ends at the lowest common ancestor for the given mixture.

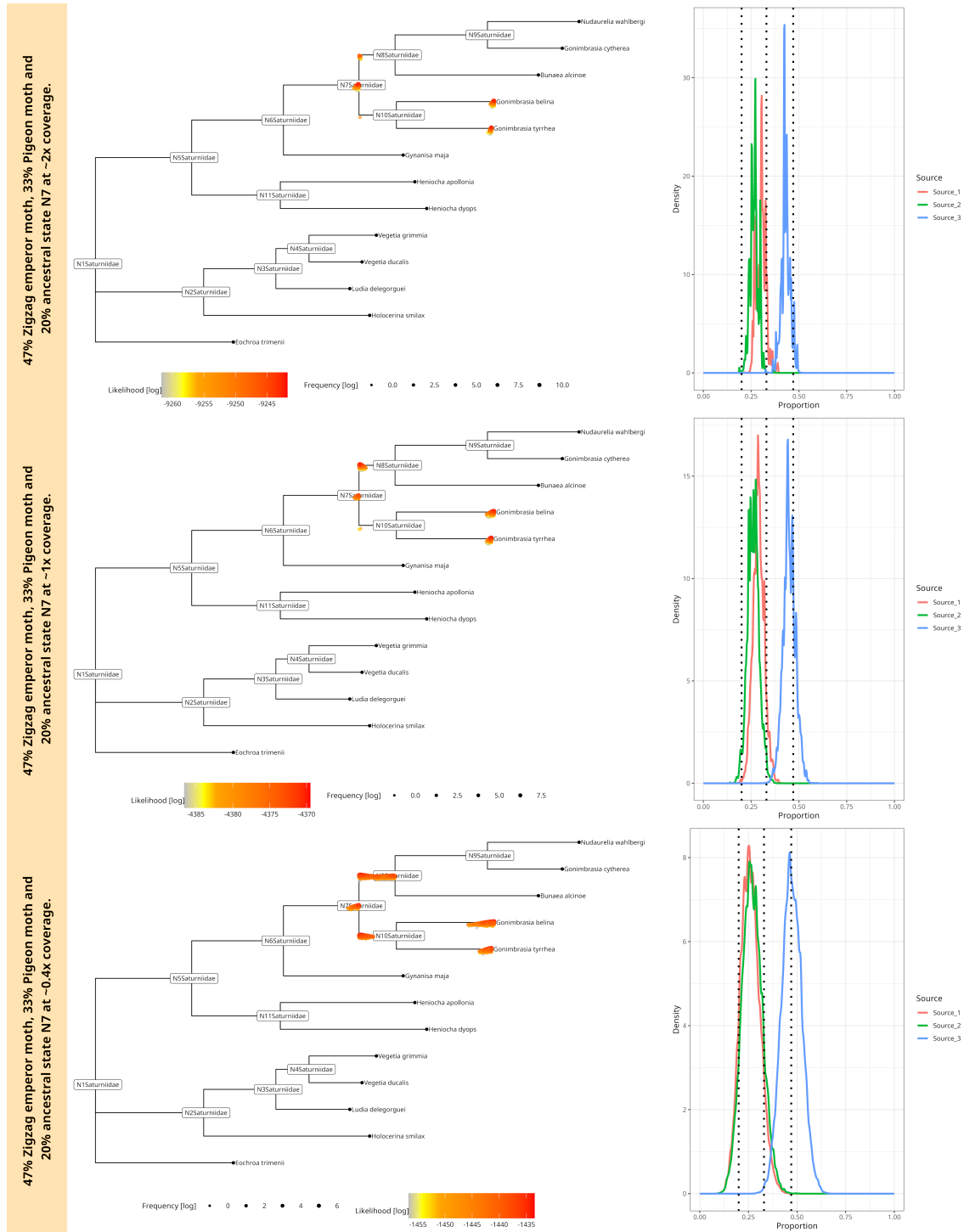


Figure 24: 44% – 30% – 23% mixture of a three-source simulated ancient sample from a family of winged insects (Saturniidae). The mixture contains two emperor moths (*Gonimbrasia tyrreha* and *Gonimbrasia belina*) as well as the ancestral state N7 and is downsampled to  $\sim 2X$  coverage (750 aDNA fragments),  $\sim 1X$  coverage (375 aDNA fragments) and  $\sim 0.4x$  coverage (150 aDNA fragments). The trees show every accepted MCMC move, coloured and placed by log-likelihood value (left side) and `soibean`'s proportion estimations on the left side, where the black dotted line shows the simulated proportion.

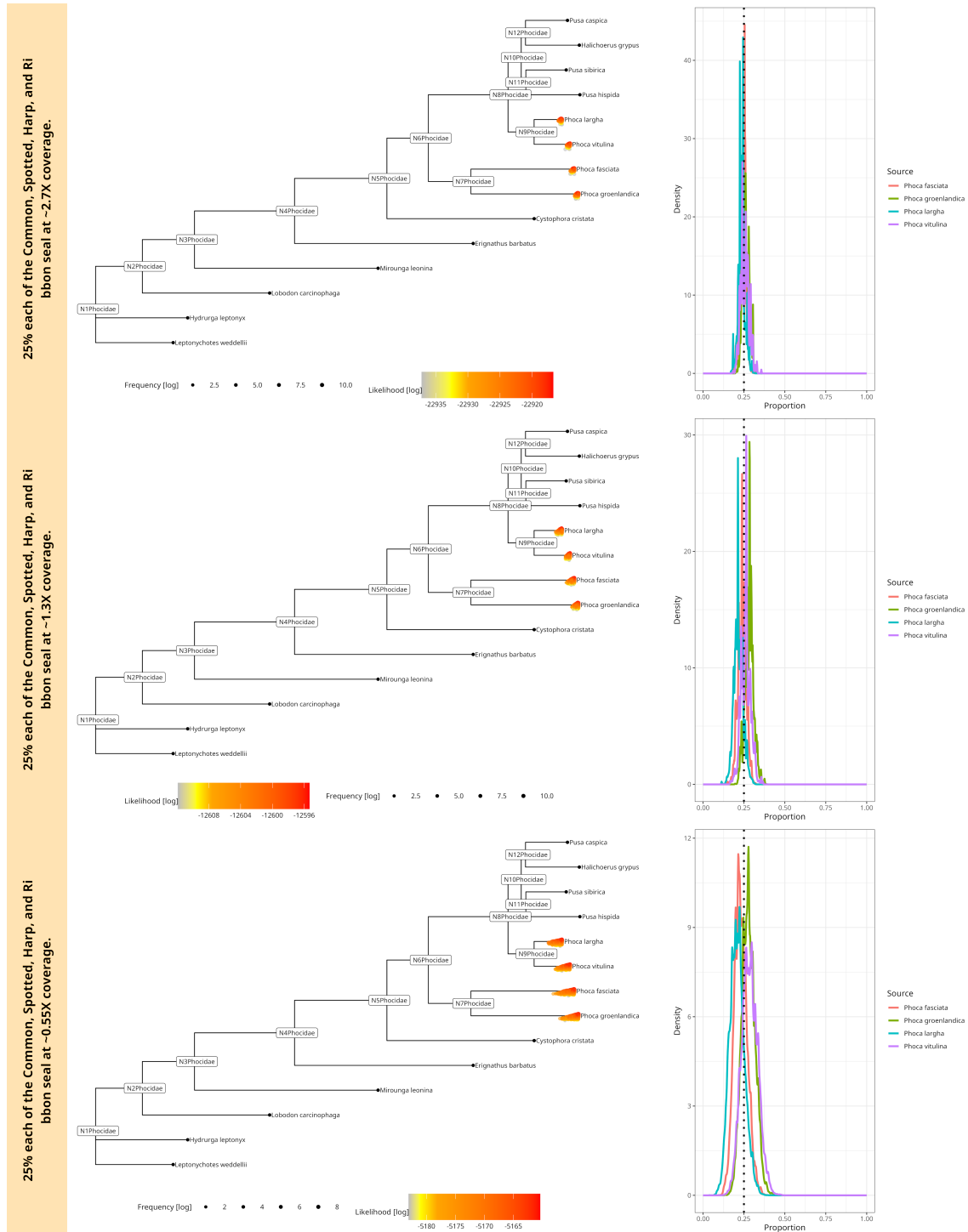


Figure 25: 25% – 25% – 25% – 25% mixture of a four-source simulated ancient sample from the family of seals (Phocidae). The mixture contains the earless seal species, namely *Phoca largha*, *Phoca vitulina*, *Phoca groenlandica* and *Pusa hispida*; it is downsampled to  $\sim 2.7X$  coverage (1000 aDNA fragments),  $\sim 1.3X$  coverage (500 aDNA fragments) and  $\sim 0.55X$  coverage (200 aDNA fragments). The trees show every accepted MCMC move, coloured and placed by log-likelihood value (left side) and `soibean`'s proportion estimations on the left side, where the black dotted line shows the simulated proportion.



Average User Time by Iterations and Number of Reads for three sources or less.

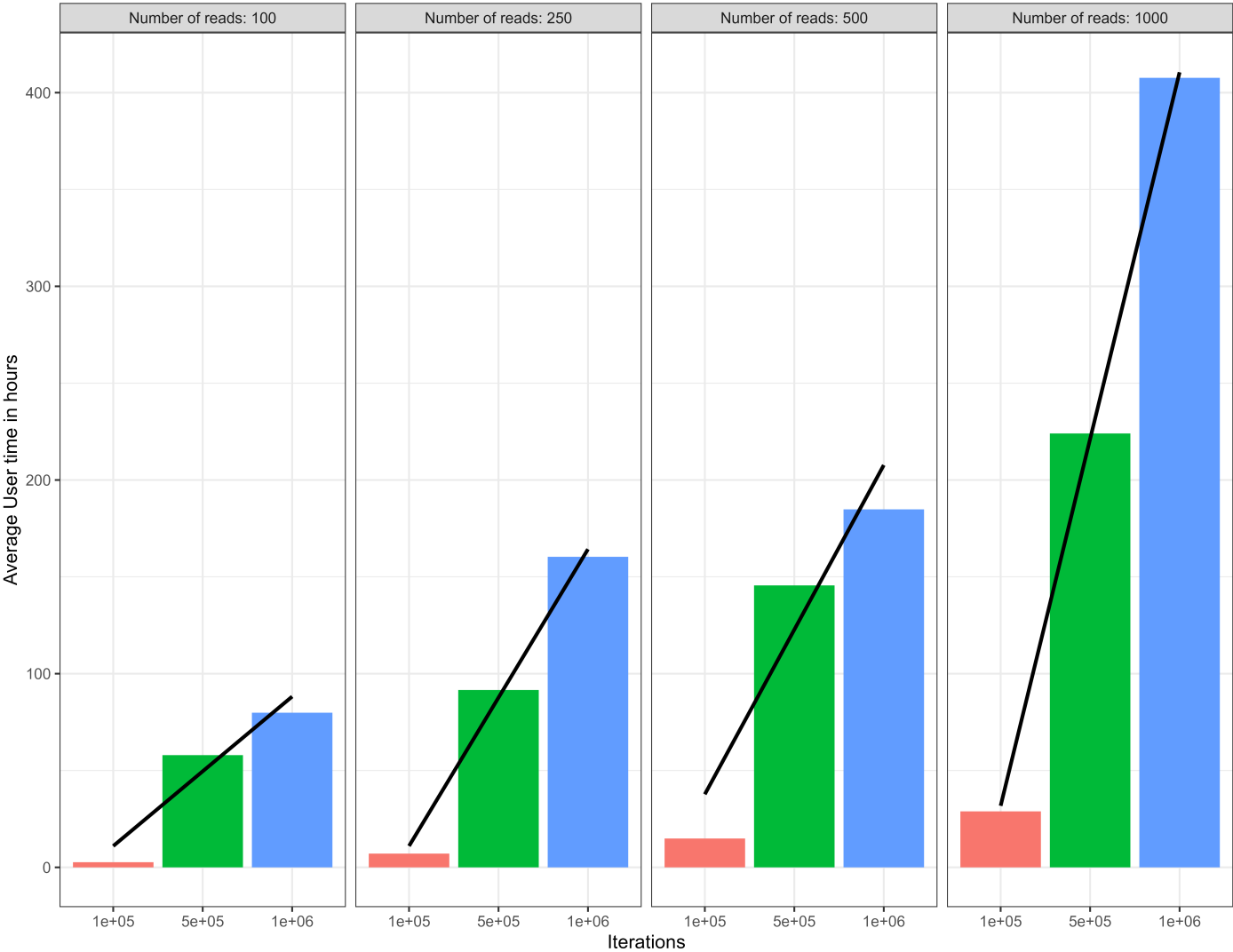


Figure 26: `soibean` user-time (hours) comparison for three different numbers of reads and three different numbers of iterations.

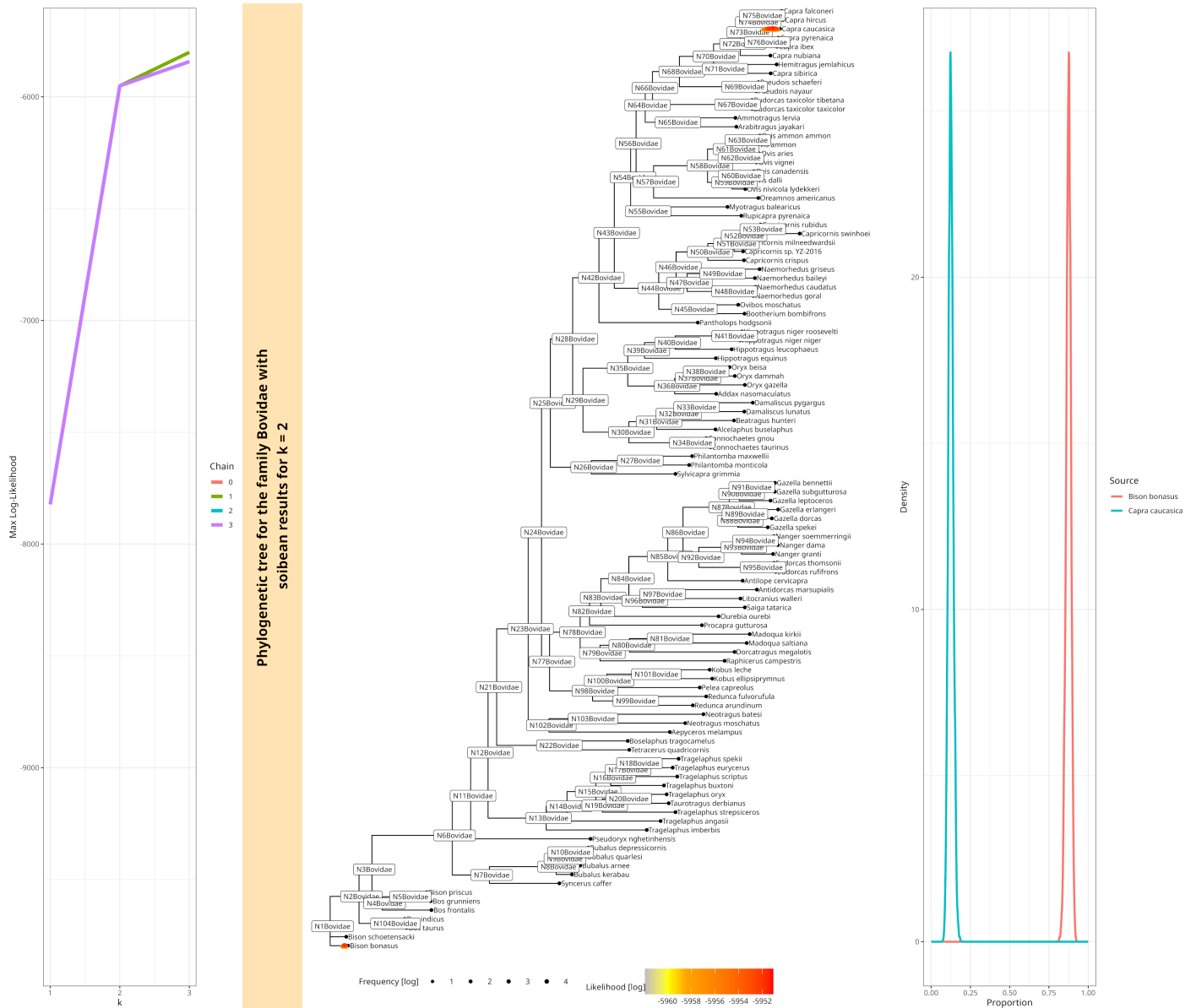


Figure 27: The results from soibean for the 25kya sediment sample from a Georgian cave. The k-curve highlights the presence of two contributing sources for the filtered Bovidae fragments. Analyzing the results, we can find the previously detected European Bison (*Bison bonasus*). Additionally, we can identify a low-coverage sample from the West Caucasian tur (*Capra caucasica*, with a proportion of about 10%, which was previously not reported).

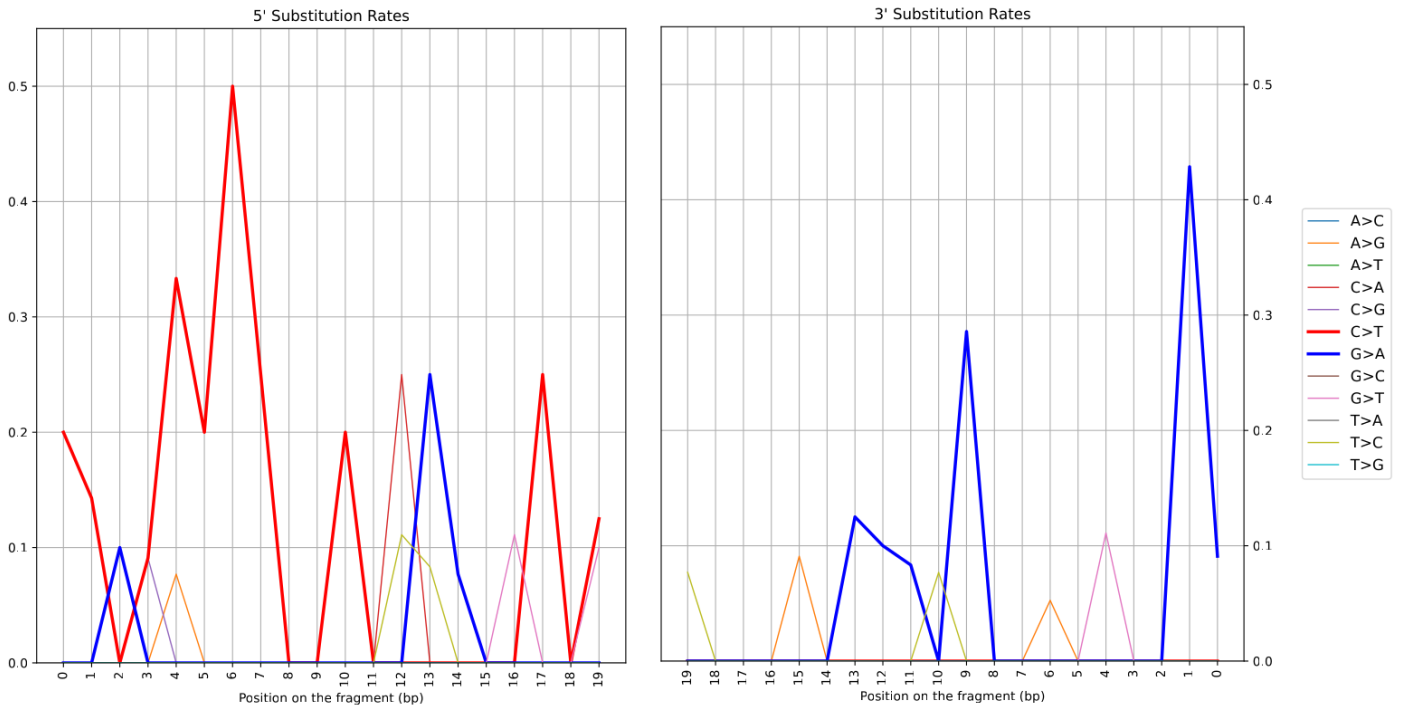


Figure 28: Deamination patterns for the alignments to the mitochondrial genome of *Capra caucasica*.

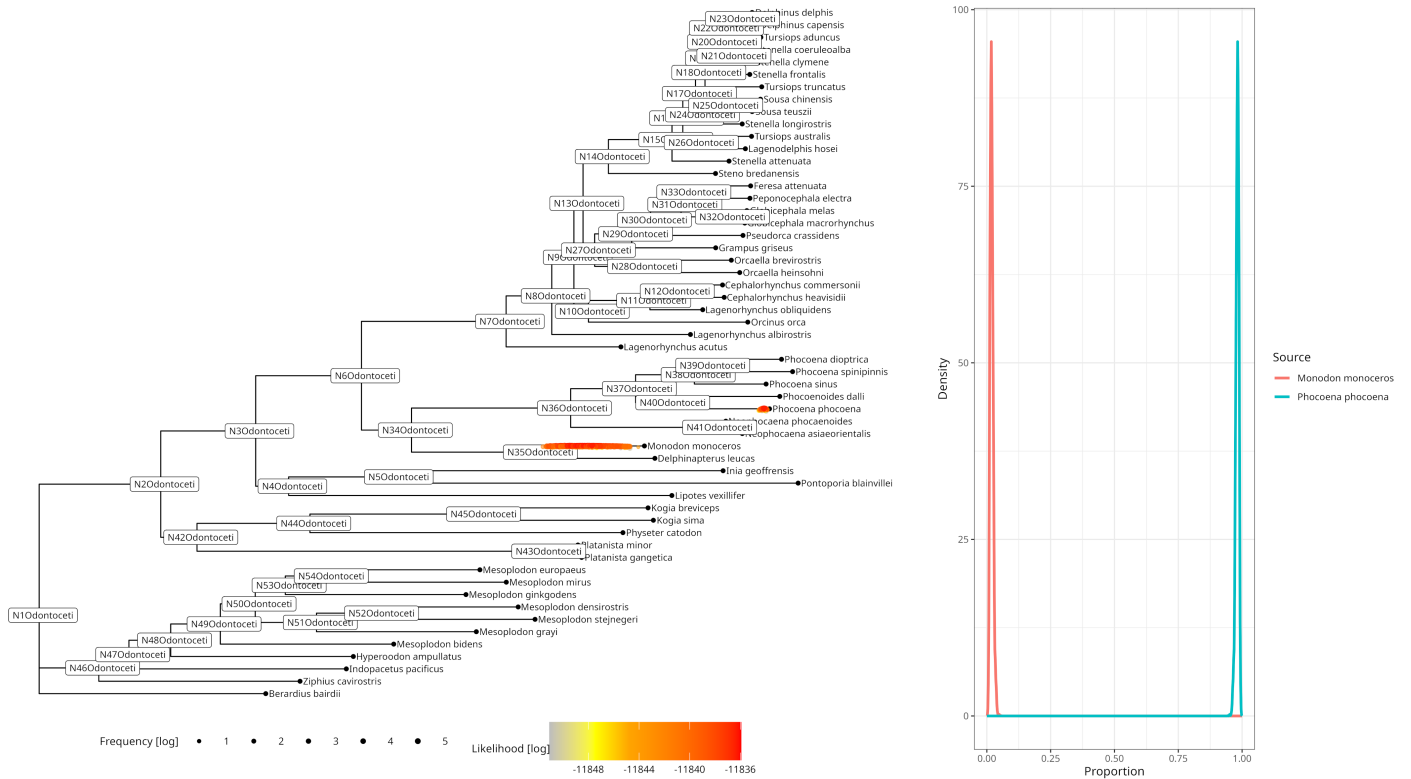


Figure 29: soibeant results for the 9500-year-old empirical metagenomic sample from pitch pieces found in Huseby Klev on the northwestern coast of Sweden for  $k = 3$ . The harbour porpoise (*Phocoena phocoena*) is still the most likely source. soibeant distributes one source at about 5% to an ancestral state of two Asiatic river dolphins. soibeant's branch and source proportion diagnostics showed an effective sample size below 200, indicating uncertainty for the second source.

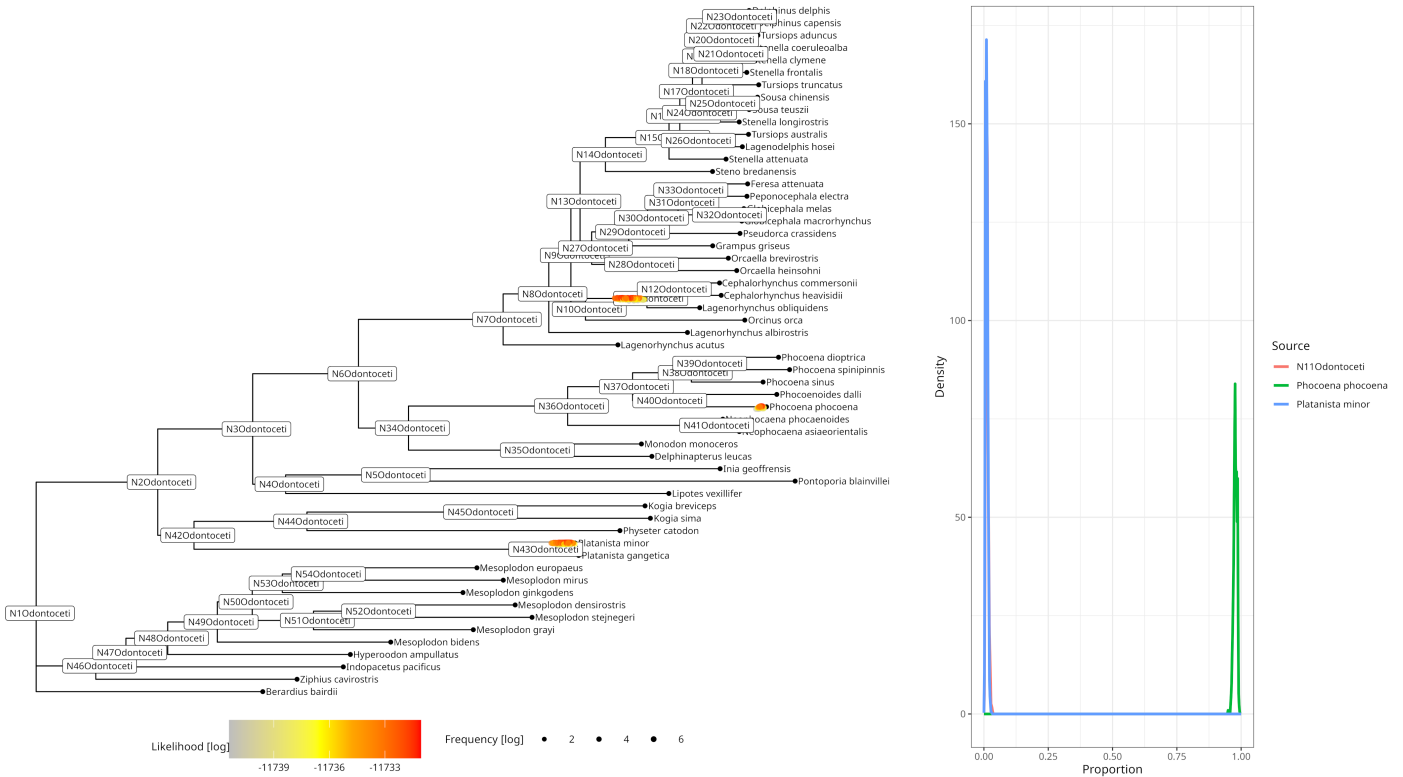


Figure 30: **soibean** results for the 9500-year-old empirical metagenomic sample from pitch pieces found in Huseby Klev on the northwestern coast of Sweden for  $k = 3$ . The harbour porpoise (*Phocoena phocoena*) is still the most likely source. **soibean** distributes two sources at about 2.5% to an ancestral state of two Asiatic river dolphins and a third dolphin species. Both additional sources have an effective sample size below 200 in **soibean**'s branch and source proportion diagnostics, making them unlikely contributing sources.

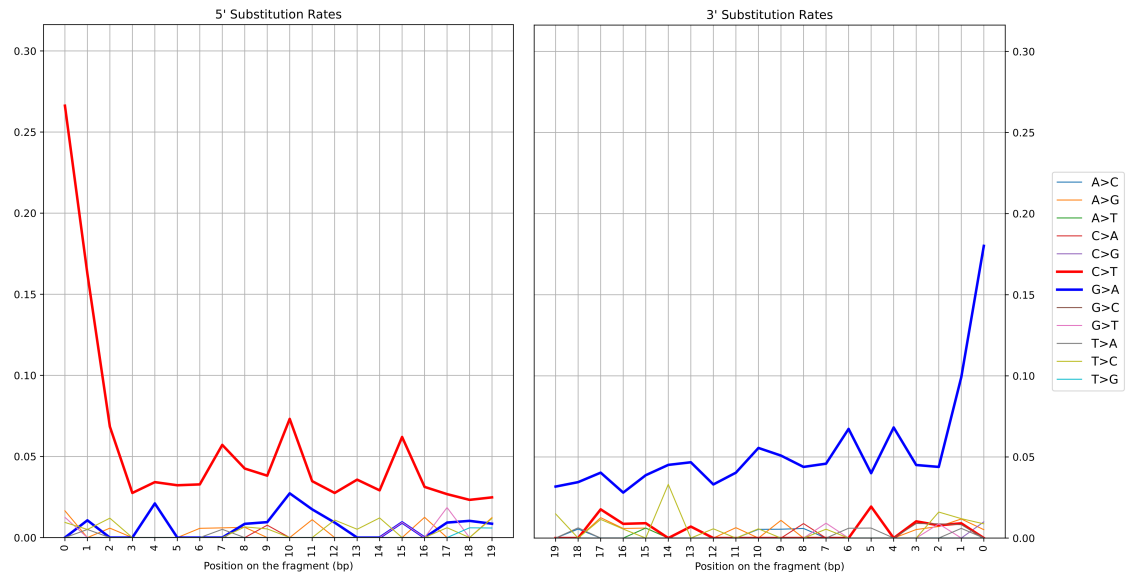


Figure 31: Complete deamination patterns for the alignments to the mitochondrial genome of *Phocoena phocoena*.