

1 **Supplemental Information for:**
2 **Biogeography in the deep: Hierarchical population genomic structure**
3 **of two beaked whale species**
4

5 Aubrie B. Onoufriou, Oscar E. Gaggiotti, Natacha Aguilar de Soto, Morgan L. McCarthy, Phillip A.
6 Morin, Massimiliano Rosso, Merel Dalebout, Nicholas Davison, Robin W. Baird, C. Scott Baker,
7 Simon Berrow, Andrew Brownlow, Daniel Burns, Florence Caurant, Diane Claridge, Rochelle
8 Constantine, Fabien Demaret, Sascha Dreyer, Martina Đuras, John Durban, Alexandros Frantzis,
9 Luis Freitas, Gabrielle Genty, Ana Galov, Sabine S. Hansen, Andrew C. Kitchener, Vidal Martin,
10 Antonio A. Mignucci-Giannoni, Valeria Montano, Aurelie Moulins, Carlos Olavarría, M. Michael
11 Poole, Cristel Reyes Suárez, Emer Rogan, Conor Ryan, Agustina Schiavi, Paola Tepsich, Jorge
12 Urban, Kristi West, Morten T. Olsen & Emma L. Carroll
13

14 **Electronic Supplementary Materials (ESM) Table of Contents:**

15 **Table of Contents**

16 ***Supplementary Spreadsheet Tables (SST) List:..... 1***
17 ***Electronic Supplementary Material (ESM) 1: Tissue Archive, DNA Extraction and Sample***
18 ***Selection 3***
19 ***ESM2: ddRAD and Mitogenome Library Preparation and Sequencing 6***
20 ***ESM3: ddRAD STACKS Parameter Optimisation 8***
21 ***ESM4: ddRAD STACKS Protocol, Loci Quality Control and Filtering Steps 9***
22 ***ESM5: Phylogenetic Trees with Southern Right Whale Outgroups 11***
23 ***ESM6: Isolation-by-Distance 14***
24 ***ESM7: ‘tess3r’ Cross-Entropy Scores 17***
25 ***ESM8: Discriminant Analysis of Principle Components..... 20***
26 ***ESM9: ddRAD Sequencing Results 24***
27 ***Supplementary References 25***
28

29 **Supplementary Spreadsheet Tables (SST) List:**
30 **Location: “Onoufriou_BW_Global_DNA_SupplementaryTables_Final.xlsx”**
31

32 **Worksheet 1 “SST1. Full Sample List”:** List of all samples included in both the ddRAD and
33 mitogenome analyses. Sample names are in the column headed “ITABW-ID”, indicating the
34 identifier used in the newly established ITAWB database. The following are also provided per

35 sample when available: species, sampling date, sampling location, sampling region, sampling
36 ocean, whether or not they were included in the ddRAD or mitogenome analysis, the
37 population defined using either method, and a detailed description of the sample's origin (who
38 contributed the sample, the origin of the sample, sample type, any permits or considerations
39 for collection of the sample and any impact minimisation or assessment steps taken during
40 sample collection).

41

42

43 **Worksheet 2 "SST2. BEAST Partitions"**: List of partitions inferred with PARTITIONFINDER and
44 analysed in a Bayesian phylogenetic framework implemented in BEAST.

45

46 **Worksheet 3 "SST3. NCBI mtDNA List"**: List of the already published mitogenome sequences
47 used in the mitogenome analysis by NCBI accession number, sample ID, species and sampling
48 locality.

49

50 **Worksheet 4 "SST4. mtDNA Seq. Stats"**: Summary of shotgun sequencing statistics for the
51 samples used in the mitogenome analyses. For each species, Copenhagen sample ID,
52 corresponding ITABW ID, sampling locality, number of total reads, reads mapped to the
53 reference mitogenome and coverage are given. Summaries of each read statistic are given
54 overall for each species.

55

56 **Worksheet 5 "SST5. mtDNA Div-Diff"**: Diversity and differentiation statistics for the mtDNA
57 data including whole mitogenomes and extracted control regions. Ocean-basin-level statistics
58 are provided including sample size (N), segregating sites (S), haplotypes (h), haplotype diversity
59 (Hd), nucleotide diversity (Nu) and fixed differences (FixDiff). Total and pair-wise differentiation
60 statistics are also provided: Fst and dA.

61

62 **Worksheet 6 "ST6. ddRAD Fst"**: Differentiation statistics (Fst) for the ddRAD data. Ocean-basin
63 and population-level statistics (Fst point estimate and 95% confidence interval) are provided
64 including sample size (n).

65

66 **Worksheet 7 "SST7. Haplotypes"**: List of samples and their mitogenome haplotypes.

67

68 Electronic Supplementary Material (ESM) 1: Tissue Archive, DNA 69 Extraction and Sample Selection

70 Tissue samples were either collected specifically for this project or provided on loan from
71 archives maintained by contributors. A large set of samples was provided by Dr. Merel
72 Dalebout, who collated samples for investigations into beaked whale systematics and genetic
73 diversity (Dalebout, 2002; Dalebout, Baker, Mead, Cockcroft, & Yamada, 2004; Dalebout et al.,
74 2007, 2014; Dalebout, Mead, Baker, Baker, & van Helden, 2002; Dalebout et al., 2005, 2003;
75 Dalebout, Ruzzante, Whitehead, & Øien, 2006; Dalebout, Steel, & Baker, 2008; Gomerčić et al.,
76 2006; Van Helden et al., 2002). All samples were either skin biopsies sampled directly from free-
77 swimming animals (Krützen et al., 2002; Lambertsen, 1987), tissue collected from dead beach-
78 cast or ship-strike individuals, or already extracted DNA provided by NOAA's Southwest
79 Fisheries Science Centre Marine Mammal and Turtle Molecular Research Sample Collection as
80 already extracted DNA ([https://www.fisheries.noaa.gov/west-coast/science-data/marine-
81 mammal-and-sea-turtle-research-tissue-collection](https://www.fisheries.noaa.gov/west-coast/science-data/marine-mammal-and-sea-turtle-research-tissue-collection)). Freshly collected tissue samples were
82 typically stored in either DMSO or 70-99% ethanol and stored at -20°C.

83
84 The purpose of the current study was to develop an understanding of the global baseline
85 genetic diversity and structure of Blainville's and Cuvier's beaked whales (*Mesoplodon*
86 *densirostris* and *Ziphius cavirostris*, respectively and henceforth 'Blainville's' and 'Cuvier's'),
87 providing context for investigating the impacts of anthropogenic activities on resident
88 populations in future studies. Balancing the coverage of samples across their respective
89 distributions and budgetary constraints, it was decided that five lanes of sequencing on a HiSeq
90 2500 (Illumina) would be used for the current study. The pilot study by Carroll et al. (2016)
91 concluded that sequencing up to 50 individuals per HiSeq2500 lane would generate ~10k
92 variable SNPs per individual, a number likely to detect structure in these two beaked whale
93 species.

94
95 Of the 89 Blainville's individuals in the sample archive, n=67 were available as tissue, n=21 were
96 already extracted DNA, and one individual had both a tissue and DNA sample. Of the 340
97 Cuvier's individuals, n=289 were available as tissue, n=29 were DNA, n=22 had both tissue and
98 already extracted DNA. DNA for the ddRAD and mitogenome analyses was extracted from
99 approximately 30-50mg of tissue, using the Phenol:Chloroform:Isoamyl Alcohol method
100 described by Sambrook, Fritsch, & Maniatis (1989) and modified for use in small tissue samples
101 by Baker et al. (1994). In addition, a subset of the Cuvier's samples used in the mitogenome
102 analysis were extracted using a KingFisher Duo™ (Thermo Scientific™) automated extraction
103 and purification instrument, following the manufacturer's instructions. Extracted DNA was
104 checked for quality using a NanoDrop™ (Thermo Scientific™) spectrophotometer and gel
105 electrophoresis, and quantified using a Qubit (Invitrogen) fluorometer.

106
107 As the protocol for double-digest restriction site-associated DNA sequencing (ddRAD) requires
108 both high molecular weight and high concentration DNA (>20ng/ul), a scoring system was
109 developed to rank samples (Table S1.1) prior to preparing libraries for pooling and sequencing.
110 DNA was run on 1.2% agarose gels to assess the overall quality of the sample and the

111 concentration was measured using spectrophotometry (NanoDrop) to obtain an approximate
 112 value. Some samples that yielded poor quality DNA were extracted up to two more times (n=4
 113 Blainville's, n=68 Cuvier's). More precise measurements of DNA quantity were made using
 114 fluorometry (Qubit) for n=88 Blainville's and n=302 Cuvier's that either had visible amounts of
 115 DNA on the agarose gel or quantifiable amounts of DNA on the spectrophotometer. Samples
 116 selected based on their DNA score were pooled into libraries with individuals that shared the
 117 same score whenever possible.

118
 119 Samples were also selected to ensure every geographical location possible was covered, and to
 120 fill in the rest of the sequencing lanes, samples from well-studied resident populations were
 121 prioritised (Hawai'i, Canary Islands, Bahamas, and Ligurian Sea). Although of poorer quality,
 122 DNA samples with scores lower than 'Good' were included in the libraries as many of them
 123 came from poorly sampled areas. Table S1.2 shows the geographical origin for the 170 Cuvier's
 124 and 55 Blainville's samples that were sequenced. Following sequencing and the bioinformatic
 125 steps outlined in Supplementary 3 and 4, a number of individuals was removed from the
 126 analysis that failed to pass quality control (QC). Table S1.3 shows the number of individuals
 127 from both Cuvier's and Blainville's that either passed or failed QC according to the score
 128 assigned to them based on DNA quality/quantity before sequencing.

129
 130 The final list of individuals used in this study is found in Supplementary Table 1 (SST1). Once
 131 duplicate individuals were removed, the number of unique individuals for ddRAD sequencing
 132 remaining was n=161 Cuvier's and n=55 Blainville's.

133
 134 Table S1.1. The scoring system developed to rank DNA samples before ddRAD library
 135 preparation based on the amount of DNA in the sample (measured using a Qubit fluorometer)
 136 and the molecular weight according to 1.2% agarose gels (HMW =High Molecular Weight,
 137 Smear=degraded DNA of varying sizes, LMW= Low Molecular Weight).

Score	DNA Concentration (Qubit)	DNA Gel Result
Great	>20ng/ul	HMW
Good	>20ng/ul	HMW + smear
Good	>20ng/ul	Smear
Good	>20ng/ul	No visible DNA
Good	15-20ng/ul	HMW
OK	>20ng/ul	LMW
OK	15-20ng/ul	Smear
OK	<15ng/ul	HMW
Some	<15ng/ul	Faint HMW
Some	<15ng/ul	Smear
Some	<15ng/ul	No visible DNA

148
 149
 150

151 Table S1.2. Broad geographic origin of the samples selected for ddRAD library preparation and
 152 sequencing for $n=170$ Cuvier's (*Ziphius cavirostris*) and $n=55$ Blainville's beaked whales
 153 (*Mesoplodon densirostris*).

	Cuvier's ($n=170$)	Blainville's ($n=55$)
North Atlantic	89	34
North Pacific	28	8
South Pacific	15	8
South Africa	2	5
Mediterranean	36	Not present

154
 155

156 Table S1.3. The number of Cuvier's (*Ziphius cavirostris*) and Blainville's beaked whales
 157 (*Mesoplodon densirostris*) that passed or failed the ddRAD quality control steps based on the
 158 DNA quality/quantity score described in Table S1.1.

Quality Score	Cuvier's ($n=170$)		Blainville's ($n=55$)	
	No. Passed QC	No. Failed QC	No. Passed QC	No. Failed QC
Great	102	29	39	9
Good	19	15	2	0
OK	2	3	1	3
Some	0	0	1	0

159
 160
 161
 162

163 ESM2: ddRAD and Mitogenome Library Preparation and Sequencing

164 ddRAD builds upon the earlier RADseq method (Baird et al., 2008) by adding a second
165 restriction enzyme (RE) to the digest and an explicit size-selection step, allowing researchers to
166 have more control over the fraction of the genome that is sequenced (Peterson, Weber, Kay,
167 Fisher, & Hoekstra, 2012). In this protocol, samples were digested with one RE that targeted a
168 commonly occurring motif (MspI, 4bp) and a rarely occurring motif (HindIII, 6bp). Unique P1
169 adaptors containing individual barcodes and PCR primers, and universal P2 adaptors with PCR
170 primers, were ligated to both ends of the digested DNA. The samples were then cleaned,
171 pooled and size selected using a Pippin Prep (Sage Science). A PCR step added a secondary
172 identifier (reverse index) to the P2 end and Illumina flow cell annealing sequences on both
173 ends. After this step, samples were cleaned, pooled and sequenced. By adding a second RE
174 digest, and eliminating random mechanical shearing and broad size selection, individual studies
175 are much more reproducible and precise. Limiting the DNA window that is sequenced and
176 precisely selecting for size, means that the sequenced fragments from different individuals are
177 more likely to be recovered from the same region of the genome (Peterson et al., 2012).

178
179 The following ddRAD protocol was optimised for beaked whale tissue by Carroll et al., (2021,
180 2016). Samples selected for sequencing were grouped based on their score, normalised to
181 20ng/μl and a total of 250ng of DNA per individual underwent an overnight digestion at 37°C
182 with MspI and HindIII. After a 20-minute heat kill step at 65°C, adaptors were ligated with one
183 of 10 forward barcodes per sample, using the temperature profile of 22°C for 2 hours followed
184 by 65°C for 20 minutes. With unique barcodes now ligated, up to 10 samples, grouped
185 according to quality classification, were pooled to form a library and cleaned using three
186 PureLink PCR Micro Kit columns (Invitrogen) per library. Following the final elution step, 30μl of
187 cleaned ligate underwent size selection to a 300-400 bp range using a Pippin Prep (Sage
188 Science). The resulting size-selected ligate was divided into 8 wells and library-specific reverse
189 indices were annealed during low-cycle number PCR using a Phusion High-Fidelity PCR kit
190 (ThermoFisher). PCR products were pooled and cleaned using AMPURE-XP (Beckman-Coulter)
191 beads and eluted to a final volume of 15μl in EB buffer (Qiagen). The final libraries were sent to
192 the National High-Throughput DNA Sequencing Centre at the University of Copenhagen where
193 the quantities and quality were determined with qPCR and Bioanalyzer (Agilent Genomics).
194 Finally, the libraries were normalised and up to five libraries (~50 samples) were pooled into
195 sequencing lanes, and sequenced on a HiSeq2500 V4 chemistry (Illumina®).

196
197 For the mitogenome sequencing, we used the Carøe et al. (2018) protocol. Briefly, DNA was
198 fragmented to approximately 350 bp, using the M220 Focused-Ultrasonicator™ (Covaris),
199 according to the manufacturer's protocol. After fragmentation, samples were quantified on the
200 Agilent 2200 TapeStation according to the protocol for genomic DNA, to verify fragmentation
201 success. DNA libraries were built, using the blunt-end single-tube protocol described by Carøe
202 et al. (2018) with a few modifications. To each library, 2 μl of 10 μM Illumina® adaptors were
203 added to the fragmented DNA, followed by a MiniElute (Qiagen) clean-up step before indexing
204 with P5 and P7 indices. Libraries were sequenced using paired-end (PE) 150 bp chemistry on
205 two lanes of Illumina® HiSeq 4000 at the National High-Throughput Sequencing Centre at

206 University of Copenhagen, Denmark. In addition, 16 libraries characterized by average fragment
207 lengths <300bp were sequenced on a single lane of HiSeq4000 using single-end 80 bp
208 chemistry.
209

210 **ESM3: ddRAD STACKS Parameter Optimisation**

211 The STACKS SNP discovery pipeline (Rochette, Rivera-Colón, & Catchen, 2019) implements the
212 Bayesian genotype caller (BGC) algorithm of Maruki & Lynch (2015, 2017) in the command
213 'gstacks'. This algorithm uses a Bayesian genotype-frequency prior that takes into account
214 population-level allele frequencies, does not assume Hardy-Weinberg equilibrium and estimates
215 error rates directly from the sequence data (not from read quality scores) when calling
216 genotypes (Maruki & Lynch, 2015, 2017). The program first estimates significant polymorphic
217 loci from the read alignments with confidence set by the parameter "var-alpha" in 'gstacks'
218 (Maruki & Lynch, 2015; Rochette et al., 2019). Genotypes for each individual are called at these
219 loci, using a method that takes into account allele balance and read depth (Maruki & Lynch,
220 2017). Confidence in the genotype calling is done using a likelihood ratio test which compares
221 the likelihood of the two most likely genotypes (Maruki & Lynch, 2015).

222
223 A subset of demultiplexed and quality-controlled (QC) reads from Cuvier's ($n=40$) and
224 Blainville's ($n=55$) individuals was selected to optimize the parameters used in the STACKS SNP
225 discovery pipeline. The subset of high-quality samples (>98% retained reads and >1,000,000
226 retained reads) were selected to cover the widest geographical range of Cuvier's and reduce
227 computing time, while all Blainville's individuals were selected. In summary, the following
228 parameters were modified to optimise the 'gstacks' command for each dataset: "min-mapq"
229 (minimum mapping quality score to consider a read; 10,20), "max-clipped" (maximum soft-
230 clipping level as a fraction of the read length; 0.1, 0.2), "var-alpha" (SNP discovery threshold;
231 0.05, 0.01) and "gt-alpha" (genotype calling threshold; 0.05, 0.01). Samples were selected for
232 optimisation based on the number and proportion of retained reads and to ensure an even
233 distribution from all geographical regions. The best combination of parameters was selected
234 based on resulting datasets with the highest number of SNP loci and the lowest amount of
235 missing data.

236
237 At the end of all STACKS and filtering steps, the optimal parameters were selected based on
238 maximizing the total number of final SNP loci and reducing the amount of missing data per
239 species. The final optimised parameters were the same for both Cuvier's and Blainville's
240 samples: mapq=10, sclip=0.2, var_alpha=0.05, gt_alpha=0.05.

241

242 **ESM4: ddRAD STACKS Protocol, Loci Quality Control and Filtering Steps**

243 Following the discovery and genotyping of SNPs in ‘gstacks’, individuals and their genotyped
 244 loci can be analysed in a framework incorporating some sort of group assignment (such as
 245 geographical origin or sex) and then filtered according to minor allele frequency or locus
 246 frequency within the entire population in the STACKS ‘populations’ module. To reduce bias from
 247 potentially arbitrary population designations, no such population framework was provided in
 248 the current study. Massively parallel sequencing can lead to high error rates and genotypic
 249 uncertainties that can be introduced at any step throughout the analysis (O’Leary, Puritz, Willis,
 250 Hollenbeck, & Portnoy, 2018). Fortunately, many errors can be overcome by employing
 251 rigorous filtering to identify and reduce errors before analyzing the final dataset (O’Leary et al.,
 252 2018). In this study, we took a tiered approach to filtering, starting with low cut-off values for
 253 missing data (applied separately per locus and individual) and finalizing the dataset with higher
 254 thresholds. This alternative and iterative filtering method, whereby you increase the cut-off
 255 threshold, has been shown to retain more loci and individuals as poor-quality individuals can
 256 deflate genotype call rates in otherwise acceptable loci, while poor-quality loci can increase the
 257 amount of missing data in otherwise acceptable individuals (O’Leary et al., 2018). Below we
 258 describe each of the steps that were implemented using R v. 3.6.0 (R Core Team, 2019) and
 259 VCFTOOLS v. 0.1.12a (Danecek et al., 2011) to filter individuals and loci based on the amount of
 260 missing data, read depth, and quality score (Table S4.1). In Table S4.2, each bioinformatic step
 261 is listed, with the resulting number of loci and individuals remaining throughout the process for
 262 both Cuvier’s and Blainville’s datasets.

263
 264 Table S4.1. List of filtering commands and steps used in the program VCFTOOLS to filter loci and
 265 individuals based on locus depth, genotype quality, minor allele frequency (MAF) and
 266 missingness. Low stringency indicates that lower cut-off values are used to filter out missing
 267 data before iteratively increasing cut-off values, a strategy shown by O’Leary et al. (2018) to
 268 increase the proportion of retained loci and individuals.

Stringency	VCFTOOLS Command	Description
Low	--minDP 5 --minGQ20	Recode genotypes with quality <20 and depth <5 to zero
	--maf 0.001	Remove the sites made monomorphic by previous step.
	--max-missing 0.5	Remove sites with >50% missing data
	--missing-indv	Calculate missingness per individual, write a list of individuals with >50% missing data
	--remove	Remove individuals on list with >50% missing data
High	--site-depth	Calculate site depth, write a list of loci with mean site depth >3x the overall mean
	--exclude-positions	Remove sites with site depth >3x the overall mean
	--max-missing 0.75	Remove sites with more than 75% missing data
	--missing-indv	Calculate missingness per individual, write a list of individuals with >25% missing data
	--remove	Remove individuals on list with >25% missing data

269

270 Table S4.2. Summary of each bioinformatic step to discover, genotype and filter loci based on
 271 the steps described in Table 1. Data are presented for the Cuvier's and Blainville's beaked whale
 272 (*Ziphius cavirostris* and *Mesoplodon densirostris*, respectively) datasets with a summary of each
 273 step and the program that was used.

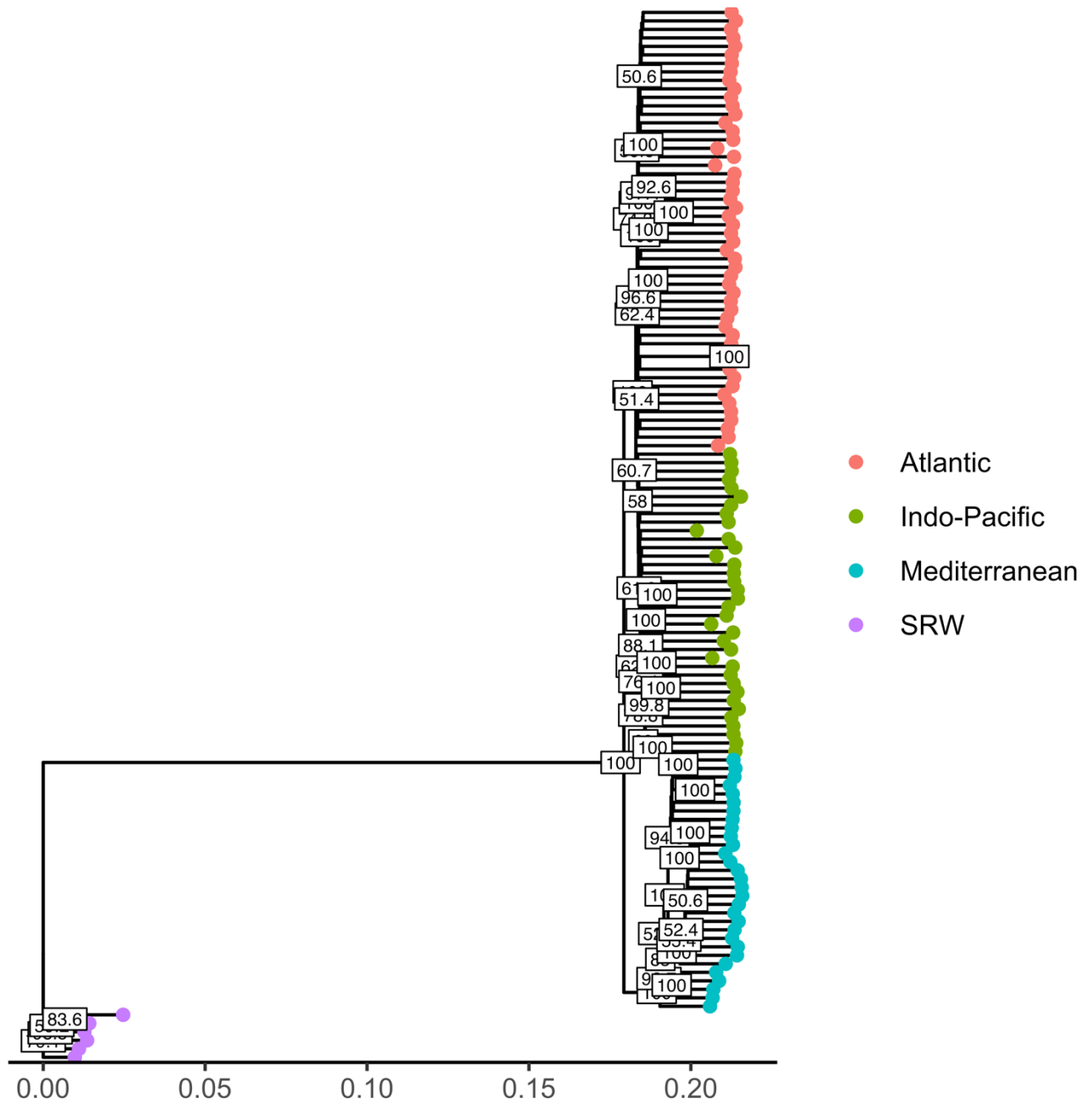
	Species starting sample size	Cuvier's 170	Blainville's 56
optimised gstacks parameters	mapq sclip var_alpha gt_alpha	10 0.2 0.05 0.05	10 0.2 0.05 0.05
Process RadTags (STACKS)	Total PE reads Retained PE reads Mean PE retained reads/sample % Retained (across all samples)	1112692388 1063993201 6258784 95.6%	339961638 322949310 5766952 95.0%
gstacks (STACKS)	Individuals Remaining Alignments Read Alignments Kept % Alignments Kept Loci Built and Genotyped	154 1085038314 778173558 71.7% 1801998	54 334888850 249685545 74.6% 979828
Populations (STACKS)	Loci Kept No. Sites No. Polymorphic Sites	1795750 477045540 2087305	977060 184528823 1054577
Filtering (VCFTOOLS)	Sites with >5x depth, >20 genotype quality, >0.001 MAF Sites with <50% missing data Sites with depth <3x overall mean depth Sites with <75% missing data	1326391 327270 326459 262482	735883 340095 339441 296250
Populations (STACKS)	Loci passed filtering (whitelist) Loci present in >80% individuals with >0.01 MAF No. Sites No. Polymorphic Sites No. Genomic Sites No. Individuals remaining	34264 31734 9994609 30479 9610872 123	37617 32610 9527357 271983 9504054 49
gPlot/Removing Duplicates ('Adegenet')	Final no. Individuals Final no. Loci	123 30479	43 13988

274
 275
 276
 277
 278

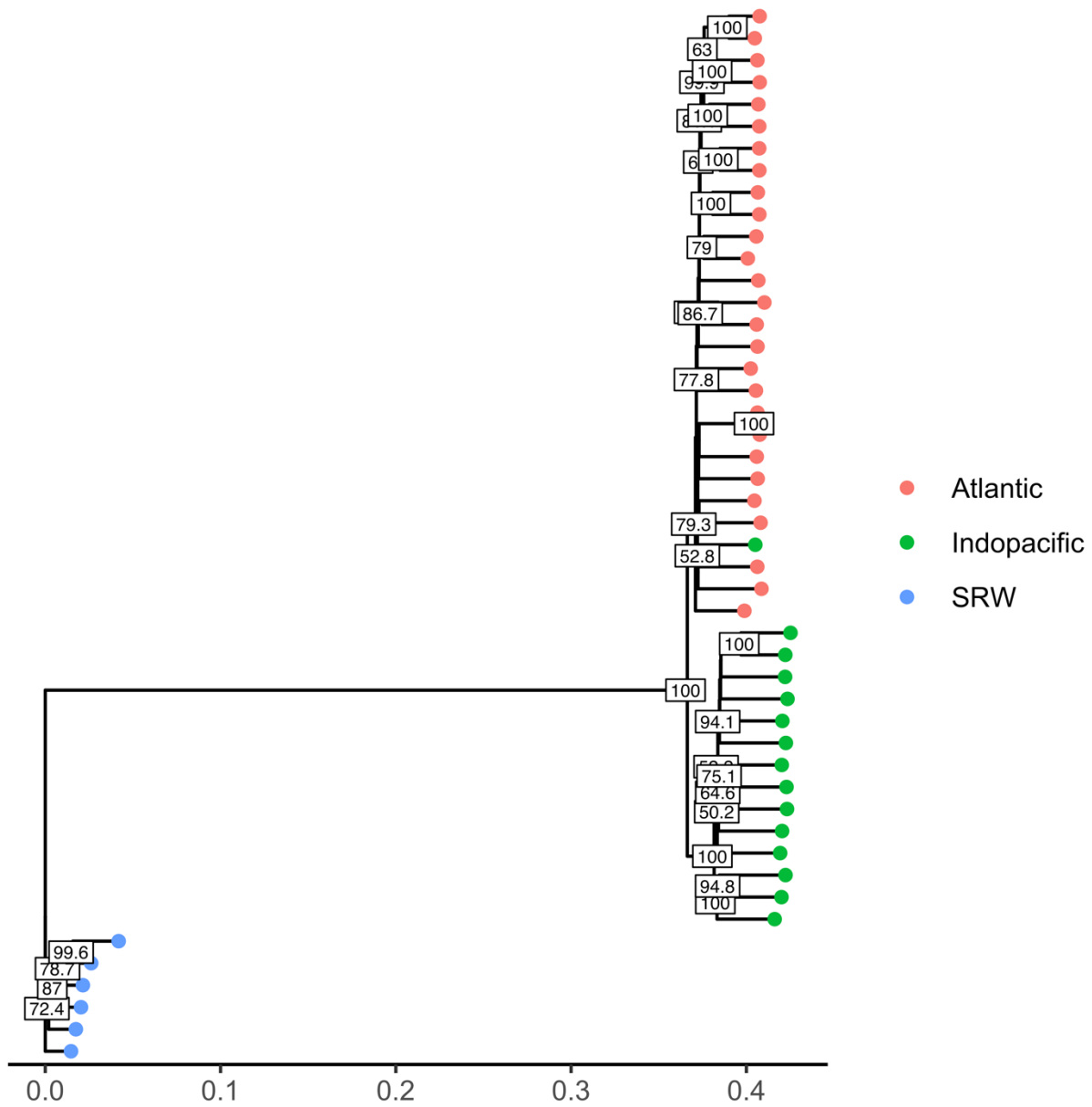
279 ESM5: Phylogenetic Trees with Southern Right Whale Outgroups

280 Phylogenetic trees using the ddRAD SNP data were generated for both Cuvier's and Blainville's
281 incorporating data from Southern right whales (SRW, *Eubalaena australis*) as the outgroup. The
282 full SNP discovery pipeline as described in SM4 was repeated for the final n=123 Cuvier's and
283 n=43 Blainville's, each time including sequence data from six SRWs. The Cuvier's + SRW dataset
284 was aligned to the same Cuvier's genome (NCBI Genbank database accession: PRJNA399469)
285 and the Blainville's + SRW sequences were aligned to the same Sowerby's genome
286 (*Mesoplodon bidens*: PRJNA399476). All bioinformatic steps were followed as before, with one
287 exception.

288
289 The final VCF file was uploaded into R and converted to a "genlight" file as before and BIONJ
290 trees with bootstrap support (in %, based on 100 bootstraps) were produced in R ('poppr'
291 v2.8.5; Kamvar, Brooks, & Grünwald, 2015 and 'ggtree'v.2.0.2; Yu, Smith, Zhu, Guan, & Lam,
292 2017). The trees were rooted using one of the SRW sequences and re-plotted. The SRW
293 individuals were then dropped from the tree, to better visualize the ocean-level phylogenetic
294 patterns of Cuvier's and Blainville's. The original BIONJ trees (without dropping the SRW
295 outgroup are found in Figures S5.1 and S5.2. The final datasets (including 6 SRW samples) were
296 n=118 Cuvier's (n=33137 SNPs) and n=42 Blainville's (n=29904 SNPs).



297 Figure S5.1. BIONJ phylogenetic tree of 118 Cuvier's beaked whales (*Ziphius cavirostris*) and 6
 298 Southern right whales (*Eubalaena australis*) as outgroups, generated using n=33137 ddRAD
 299 SNPs.



300 Figure S5.2. BIONJ phylogenetic tree of 42 Blainville's beaked whales (*Mesoplodon densirostris*)
 301 and 6 Southern right whales (*Eubalaena australis*) as outgroups, generated using n=329904
 302 ddRAD SNPs.

303
 304
 305
 306
 307
 308
 309
 310

311 **ESM6: Isolation-by-Distance**

312 Isolation by distance (IBD) was calculated per species, and within ocean basin per species, using
 313 a Mantel test in ‘ade4’ v1.7-16 in R (Dray & Dufour, 2007) and geographical distances calculated
 314 via the least cost (LC) path distance over seawater in ‘marmap’ v1.0.5 (Pante & Simon-Bouhet,
 315 2013). To calculate the LC distance over seawater, sampling locations were plotted on a global
 316 bathymetry map (with 4 minute resolution) using the “getNOAA.bathy” function in ‘marmap’.
 317 The resolution of the world bathymetry map resulted in some sampling locations of stranded
 318 individuals to be on land, and therefore incur a great coast in the LC path. Sample coordinates
 319 were therefore adjusted to the nearest -200m isobath using the “dist2isobath” function in
 320 ‘marmap’. The updated sample coordinates and bathymetry map were used to calculate a
 321 transition matrix using “trans.mat” in ‘marmap’, requiring the LC path to have a minimum
 322 depth of 200m. Finally, the LC path distance between each individual was calculated using
 323 “lc.dist” in ‘marmap’. The resulting pairwise matrix of geographic distances was used in
 324 combination with pairwise genetic distance (Euclidean) to run a Mantel test using
 325 “mantel.randtest” with 999 permutations. Mantel tests were conducted based on 999
 326 replicates for all Cuvier’s or Blainville’s combined, and for individual ocean basins (Atlantic,
 327 Indo-Pacific, and Mediterranean- Cuvier’s only) (Table S6.1). All correlation values were positive
 328 and all but the Indo-Pacific Blainville’s were significant (p<0.05). The genetic and geographic
 329 distance matrices were plotted with a 2 dimensional kernel density estimation to visualize
 330 whether the apparent IBD was the result of a continuous cline or population clustering (Figures
 331 S6.1 and S6.2).

332

333 Table S6.1. The observation correlation and associated p-value of Mantel tests for Isolation by
 334 distance.

		<i>n</i>	Observation correlation (<i>r</i>)	P-value
Cuvier’s	All	123	0.308	0.001
	Atlantic	54	0.154	0.001
	Indo-Pacific	36	0.162	0.028
	Mediterranean	33	0.218	0.002
Blainville’s	All	43	0.665	0.001
	Atlantic	28	0.110	0.03
	Indo-Pacific	15	0.014	0.427

335

336

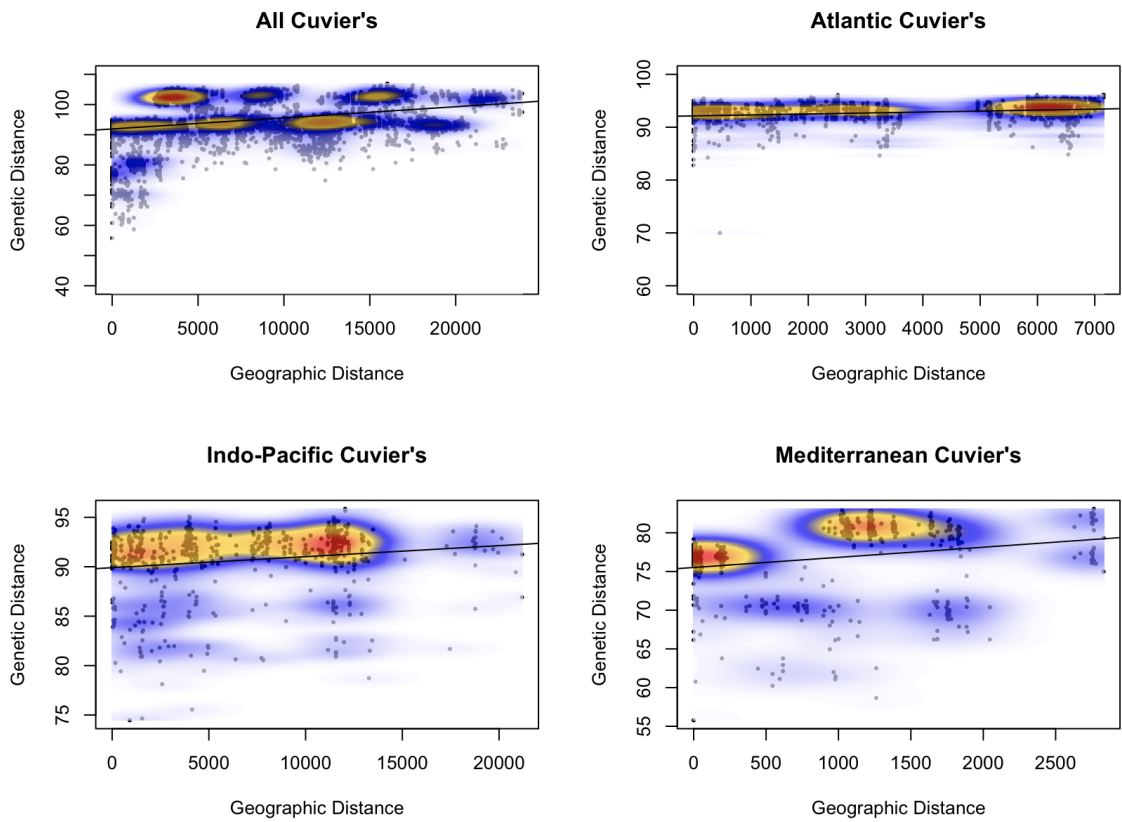
337

338

339

340

341

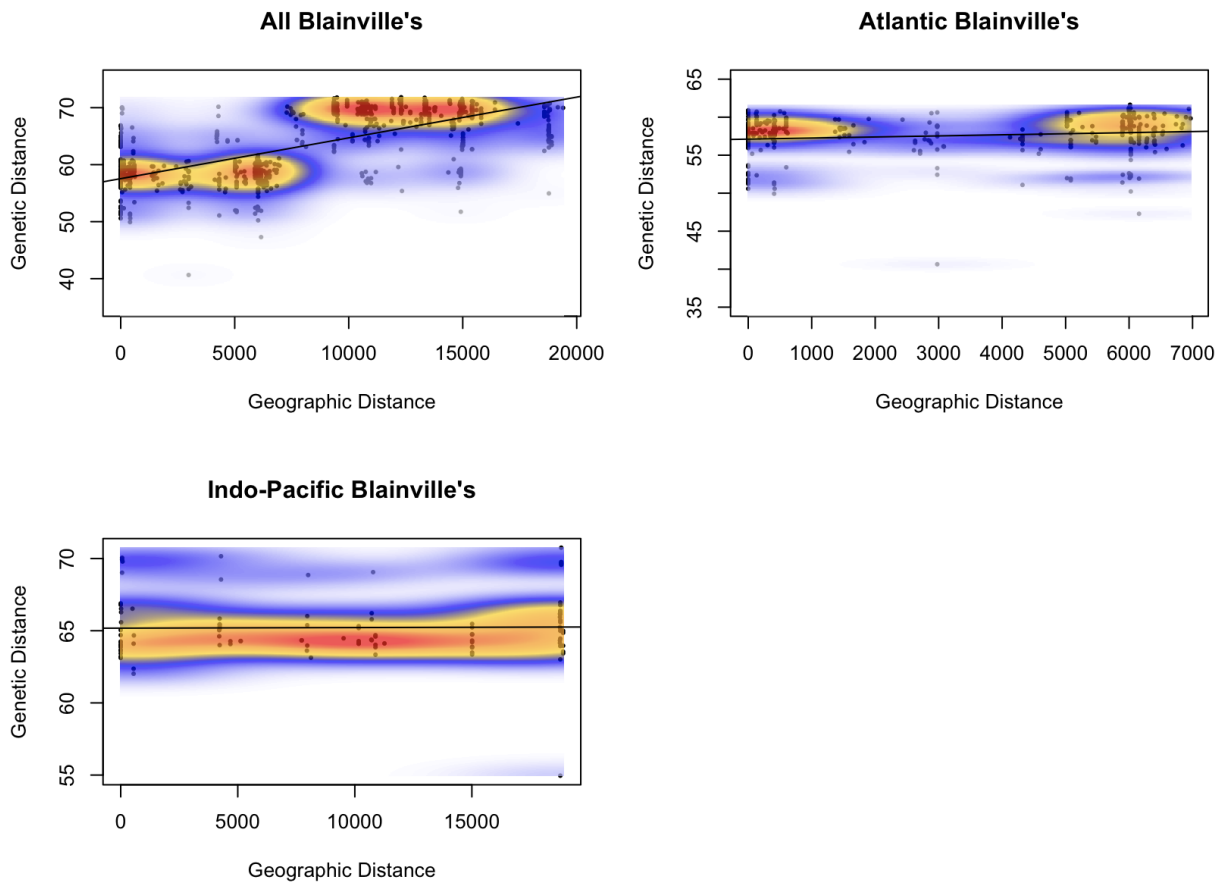


342
 343 Figure S6.1. Scatterplots of genetic distance (Euclidian) and geographic distance (least cost (LC)
 344 path distance over seawater) overlaid with 2-dimensional kernel density estimation from n=123
 345 Cuvier's beaked whales (*Ziphius cavirostris*, n=30479 SNPs).
 346

347
 348
 349
 350
 351
 352
 353
 354
 355
 356
 357
 358

359

360



361 Figure S6.2 Scatterplots of genetic distance (Euclidian) and geographic distance (least cost (LC)
362 path distance over seawater) overlaid with 2-dimensional kernel density estimation from n=43
363 Blainville's beaked whales (*Mesoplodon densirostris*, n=13988 SNPs).

364

365

366 ESM7: 'tess3r' Cross-Entropy Scores

367 The R package 'tess3r' (Caye, Jay, Michel, & Francois, 2018) incorporates genotypic and
368 geographical information (latitude and longitude coordinates for each sample) in a spatially
369 explicit, least-squares optimization approach to estimate ancestry. The user defines and
370 compares results from a range of ancestral population clusters (K) with the resulting bar plots
371 displaying ancestry coefficients reflecting the probability of population membership and
372 estimates of admixture. As opposed to the initial version of 'tess' (Chen, Durand, Forbes, &
373 François, 2007; Durand, Jay, Gaggiotti, & François, 2009), there is no biological model
374 underlying this version however, the model does expect that individuals sampled in close
375 geographical proximity are more likely to share ancestry than those sampled from further
376 away.

377
378 For both Cuvier's and Blainville's, tess3r was run for $K=2-10$ and cross-entropy scores were
379 plotted against K values to infer the most likely number of genetic clusters. In cross-entropy
380 plots, smaller values indicate better fit with the best estimate of K corresponding to the value at
381 which the curve reaches a plateau or starts to increase. In cases where a clear minimum or
382 plateau is not observed, the K value that leads to the most parsimonious assignment of
383 individuals (least amount of admixture) to populations can be considered as selection criteria.
384 The figures below display the cross-entropy scores for: Cuvier's (Figure S7.1) and Blainville's
385 (Figure S7.2).

386

387

388

389

390

391

392

393

394

395

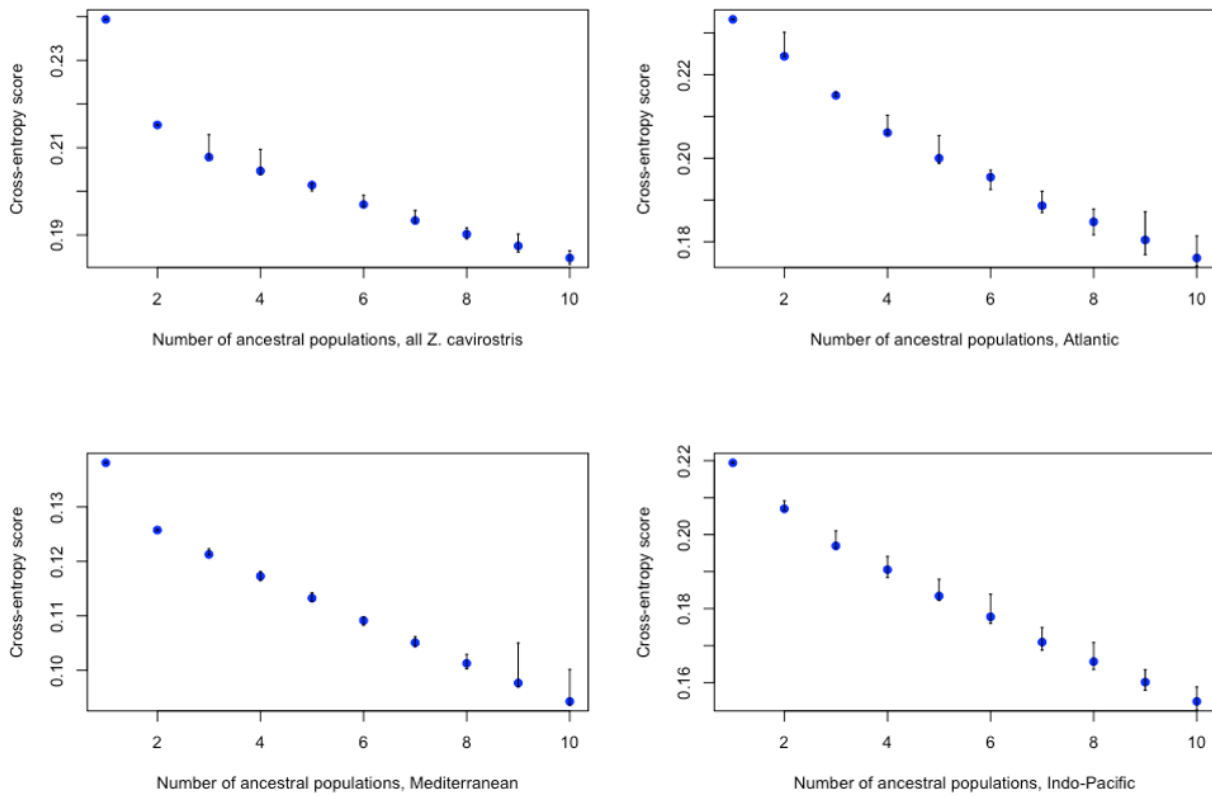
396

397

398

399

400



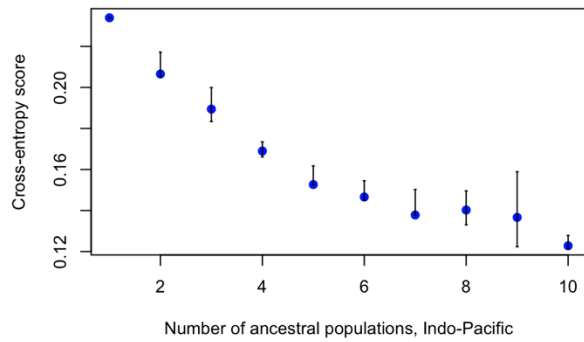
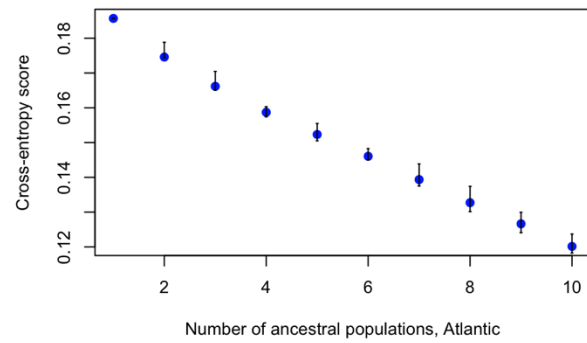
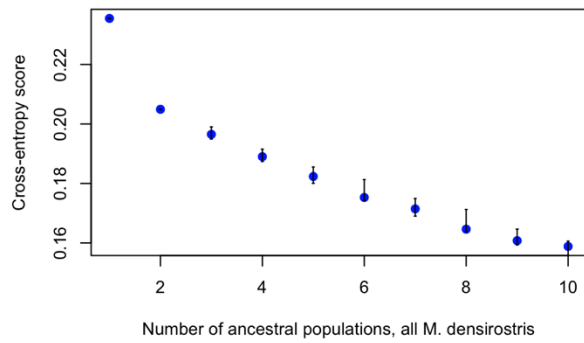
401

402

403

404

Figure S7.1. Cross-entropy scores of $K=2-10$ genetic clusters generated using 'tess3r' for $n=123$ Cuvier's beaked whales (*Ziphius cavirostris*) sampled from across their global range (top left), the Atlantic (top right), Mediterranean (bottom left) and Indo-Pacific (bottom right).



405
 406
 407
 408
 409

Figure S7.2. Cross-entropy scores of $K=2-10$ genetic clusters generated using 'tess3r' for $n=43$ Blainville's beaked whales (*Mesoplodon densirostris*) sampled from across their global range (top left), the Atlantic (top right), and Indo-Pacific (bottom left).

410 **ESM8: Discriminant Analysis of Principle Components**

411 Discriminant Analysis of Principle Components (DAPC) was conducted in the R package,
412 'adegenet' (Jombart, Devillard, & Balloux, 2010). DAPC is a useful tool to summarise the
413 amount of genetic differentiation between groups (either determined *a priori* or *de novo* using
414 K-means clustering) while ignoring the amount of variation within groups. To optimise the
415 number of principle components (PCs) to retain from the analysis, 'adegenet' offers a cross-
416 validation tool with "xvalDAPC". This command subsets the data to use as a training set, runs
417 the analysis over a pre-determined number of repeats (n=30), and determines the best number
418 of PCs to retain based on whichever yields the highest predictive success of the training data
419 with the lowest root mean squared error (RMSE). The resulting DAPC can be plotted to observe
420 the spatial structure of SNP genotypes across clusters and the function "assignplot" can be used
421 to visualize the proportion of successful reassignment to the prior groups.

422

423 DAPC with cross-validation was conducted for n=118 Cuvier's (excluding Atl_Spain and
424 Indo_Mix). The highest mean success and lowest MSE was achieved when 20 PCs were
425 retained. The resulting scatterplot and assignment plot are found in figure S8.1. DAPC with
426 cross-validation was conducted for n=43 Blainville's and the optimal number of PCs to retain
427 was 10. The resulting scatter and assignment plots are in figure S8.2.

428

429 The presence of hierarchical structure requires investigation of more than the first and second
430 axes to resolve finer scales. Scatter plots of the 2nd vs 3rd (Figure S8.3) and 3rd vs 4th (Figure
431 S8.4) axes help to discriminate between the genetic clusters found in Cuvier's within the
432 Atlantic and Indo-Pacific.

433
 434
 435
 436
 437
 438
 439
 440
 441
 442
 443
 444
 445
 446
 447
 448
 449
 450
 451
 452
 453
 454
 455
 456
 457
 458
 459
 460
 461
 462
 463
 464
 465
 466
 467
 468
 469
 470
 471
 472
 473
 474
 475
 476

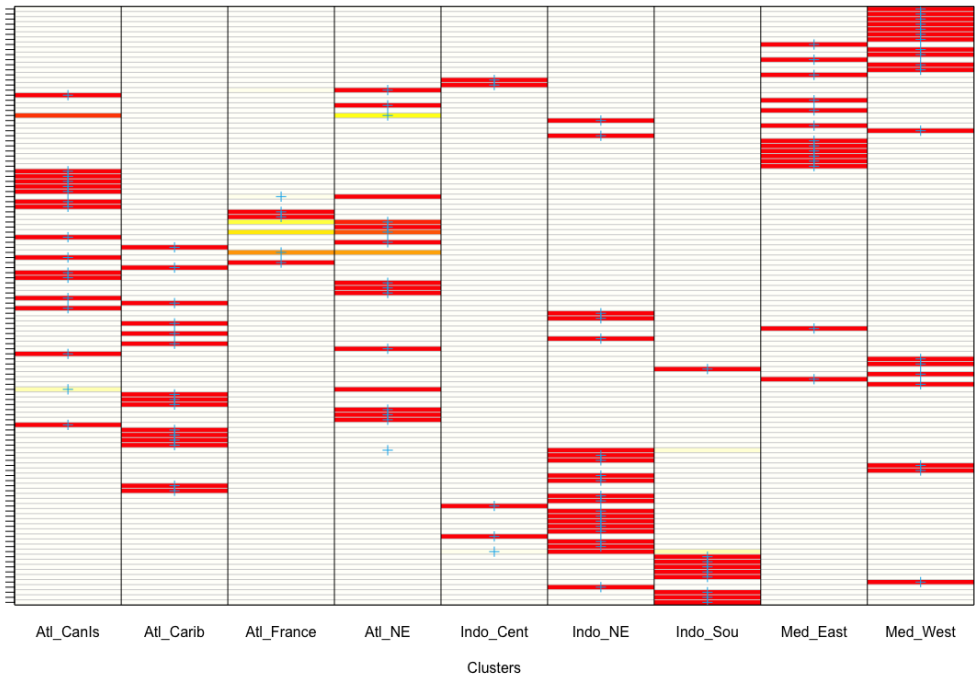
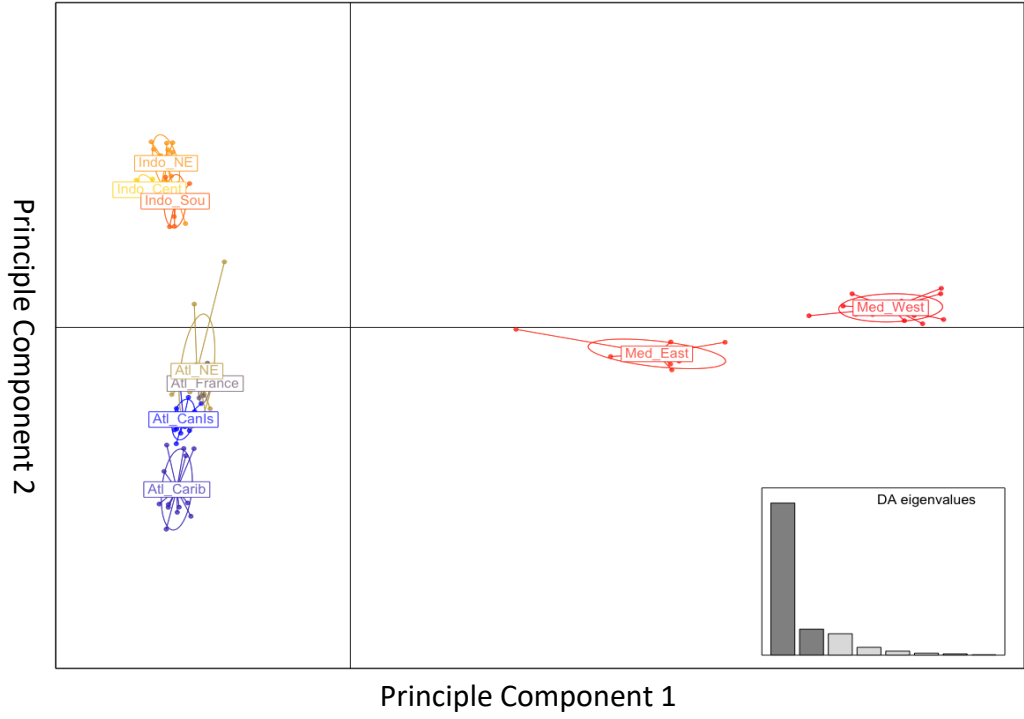


Figure S8.1. DAPC scatter (top) and assignment (bottom) plot for n=118 Cuvier’s beaked whales (*Ziphius cavirostris*) generated using cross-validation and retaining 20 PCs. In the assignment plot, each row represents an individual, the blue cross indicates the prior cluster assignment and the colours represent membership probability (red=1, white=0).

477
 478
 479
 480
 481
 482
 483
 484
 485
 486
 487
 488
 489
 490
 491
 492
 493
 494
 495
 496
 497
 498
 499
 500
 501
 502
 503
 504
 505
 506
 507
 508
 509
 510
 511
 512
 513
 514
 515
 516
 517
 518
 519
 520
 521
 522

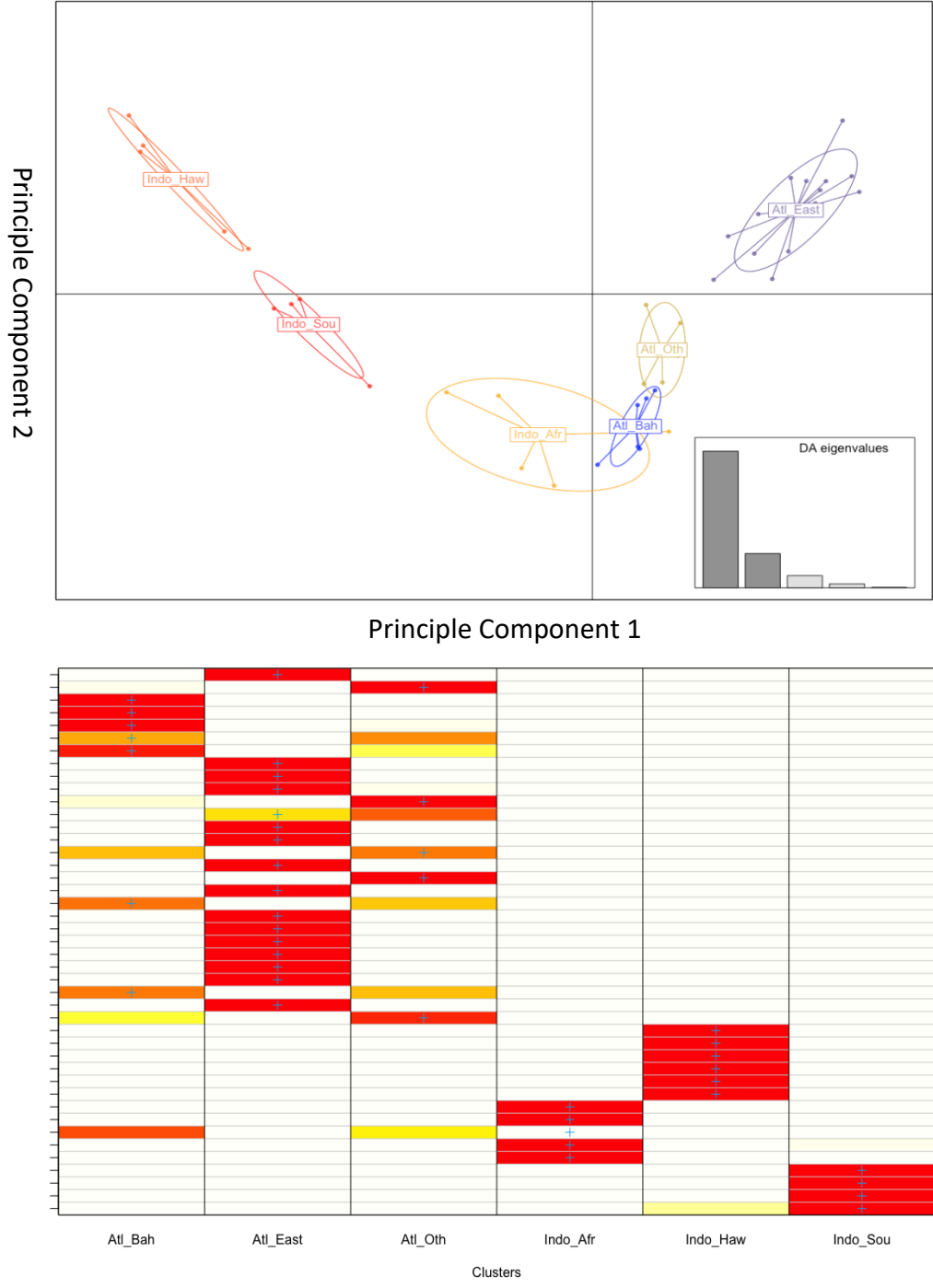
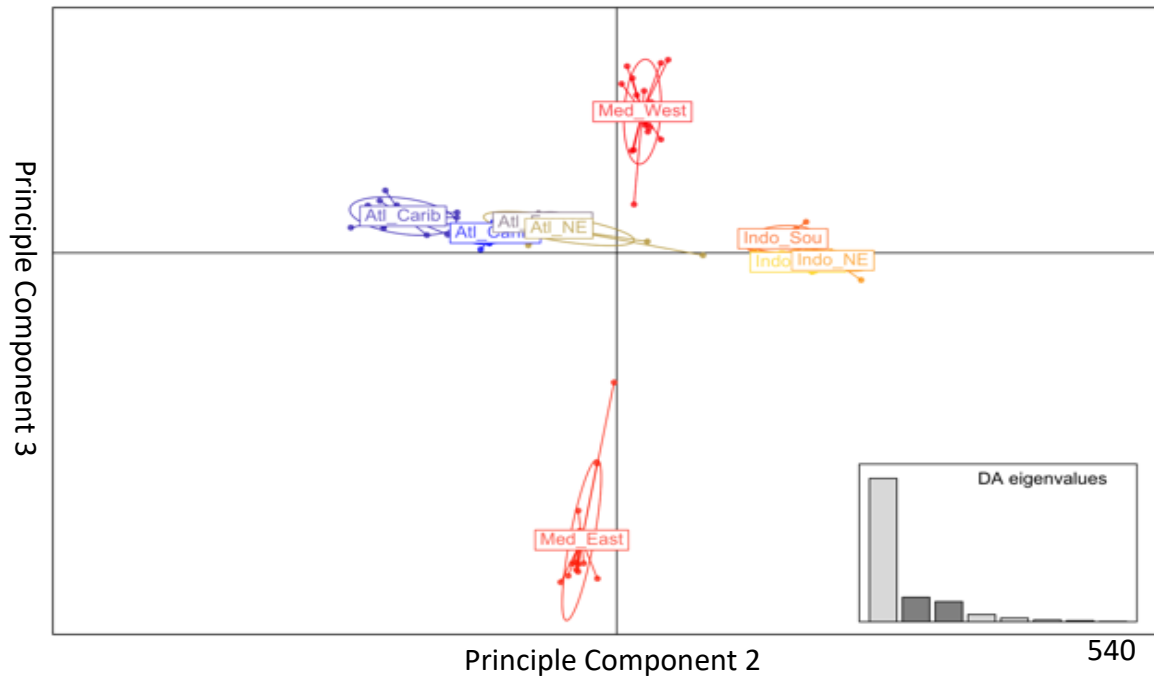
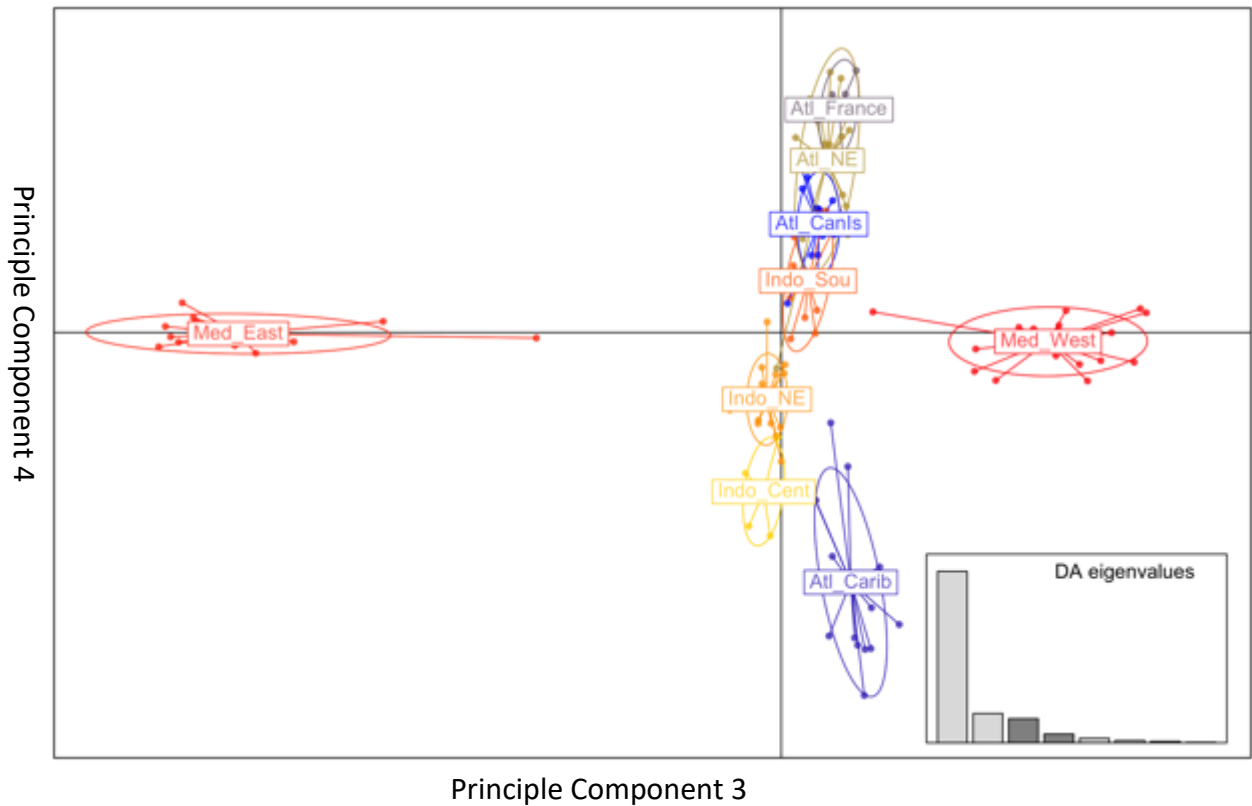


Figure S8.2. DAPC scatter (top) and assignment (bottom) plot for n=43 Blainville's beaked whales (*Mesoplodon densirostris*) generated using cross-validation and retaining 10 PCs. In the assignment plot, each row represents an individual, the blue cross indicates the prior cluster assignment, and the colours represent membership probability (red=1, white=0).



541
 542 Figure S8.3. DAPC scatter plot (2nd and 3rd axes) for n=118 Cuvier's beaked whales (*Ziphius*
 543 *cavirostris*) generated using cross-validation and retaining 20 PCs.



544
 545
 546 Figure S8.4. DAPC scatter plot (3rd and 4th axes) for n=118 Cuvier's beaked whales (*Ziphius*
 547 *cavirostris*) generated using cross-validation and retaining 20 PCs.

548 ESM9: ddRAD Sequencing Results

549 From the ITABW collection, 225 individuals were selected for ddRAD analysis (Cuvier's $n=170$,
550 Blainville's $n=55$), balancing DNA quality and quantity, and covering as much of each species'
551 broad geographical ranges as possible (See ESM1 for sample selection process). The samples
552 were run across five HiSeq 2500 lanes, generating a total of 340 million and 1.113 billion PE
553 reads across libraries of Blainville's and Cuvier's samples, respectively. Following demultiplexing
554 and initial QC, 323 million and 1.064 billion PE reads were retained for Blainville's and Cuvier's
555 samples, respectively. The number of reads and proportion of those retained were consistent
556 across libraries. In the Blainville's libraries ($n=7$), the mean number of retained reads was 5.77
557 million per sample (97% retained, standard error of the mean (SE)=1.4%) and in the Cuvier's
558 libraries ($n=18$), the mean number of retained reads was 6.26 million per sample (94.6%
559 retained, SE=1.0%).

560

561 The final ddRAD QC dataset included 123 Cuvier's individuals and 30,479 (72.4% individuals
562 retained) and 43 Blainville's individuals and 13,988 SNPs (76.8% individuals retained) (Figure 1).
563 Each Cuvier's individual was genotyped at an average of 29,697 (SE=90.5) SNPs with a mean per
564 locus read depth of 59x (SE=3.6). Each Blainville's individual was genotyped at an average of
565 13,760 (SE=63.59) SNPs with a mean per locus read depth of 53x (SE=4.4). Overall, both
566 datasets had low levels of missing data (2.6% in Cuvier's and 1.6% in Blainville's). Only SNPs
567 with a genotype quality greater than 20 were kept (99% base call accuracy).

568

569

570

571

572

573 Supplementary References

- 574 Baird, N. A., Etter, P. D., Atwood, T. S., Currey, M. C., Shiver, A. L., Lewis, Z. A., ... Johnson, E. A.
575 (2008). Rapid SNP discovery and genetic mapping using sequenced RAD markers. *PLoS*
576 *ONE*, 3(10), 1–7. <https://doi.org/10.1371/journal.pone.0003376>
- 577 Baker, C. S., Slade, R. W., Bannister, J. L., Abernethy, R. B., Weinrich, M. T., Lien, J., ... Palumbi, S.
578 R. (1994). Hierarchical structure of mitochondrial DNA gene flow among humpback whales
579 *Megaptera novaeangliae*, world-wide. *Molecular Ecology*, 3(4), 313–327.
580 <https://doi.org/10.1111/j.1365-294X.1994.tb00071.x>
- 581 Carøe, C., Gopalakrishnan, S., Vinner, L., Mak, S. S. T., Sinding, M.-H. S., Samaniego, J. A., ...
582 Gilbert, M. T. P. (2018). Single-tube library preparation for degraded DNA. *Methods in*
583 *Ecology and Evolution*, 9(2), 410–419. <https://doi.org/10.1111/2041-210X.12871>
- 584 Carroll, E. L., McGowen, M. R., McCarthy, M. L., Marx, F. G., Aguilar de Soto, N., Dalebout, M. L.,
585 ... Olsen, M. T. (2021). Speciation in the deep: Genomics and morphology reveal a new
586 species of beaked whale *Mesoplodon eueu*. *Proceedings of the Royal Society B: Biological*
587 *Sciences*, 288(1961). <https://doi.org/10.1098/rspb.2021.1213>
- 588 Carroll, E. L., Reyes, C., Gaggiotti, O. E., Olsen, M. T., Maaholm, D. J., Rosso, M., ... Aguilar de
589 Soto, N. (2016). Pilot study to assess the utility of ddRAD sequencing in identifying species-
590 specific and shared SNPs among Blainville's (*Mesoplodon densirostris*) and Cuvier's (*Ziphius*
591 *cavirostris*) beaked whales. *International Whaling Commission*, (July).
592 <https://doi.org/10.13140/RG.2.1.2286.5527>
- 593 Caye, K., Jay, F., Michel, O., & Francois, O. (2018). Fast inference of individual admixture
594 coefficients using geographic data. *Annals of Applied Statistics*, 12(1), 586–608.
595 <https://doi.org/https://doi.org/10.1214/17-AOAS1106>
- 596 Chen, C., Durand, E., Forbes, F., & François, O. (2007). Bayesian clustering algorithms
597 ascertaining spatial population structure: A new computer program and a comparison
598 study. *Molecular Ecology Notes*, 7(5), 747–756. [https://doi.org/10.1111/j.1471-](https://doi.org/10.1111/j.1471-8286.2007.01769.x)
599 [8286.2007.01769.x](https://doi.org/10.1111/j.1471-8286.2007.01769.x)
- 600 Dalebout, M. L. (2002). *Species identity, genetic diversity and molecular systematic relationships*
601 *among the Ziphiidae (beaked whales)*. University of Auckland.
- 602 Dalebout, M. L., Baker, C. S., Mead, J. G., Cockcroft, V. G., & Yamada, T. K. (2004). A
603 comprehensive and validated molecular taxonomy of beaked whales, family Ziphiidae.
604 *Journal of Heredity*, 95(6), 459–473. <https://doi.org/10.1093/jhered/esh054>
- 605 Dalebout, M. L., Baker, C. S., Steel, D. J., Robertson, K. M., Chivers, S. J., Perrin, W. F., ...
606 Schofield Jr., T. D. (2007). A divergent mtDNA lineage among *Mesoplodon* beaked whales:
607 Molecular evidence for a new species in the tropical pacific? *Marine Mammal Science*,
608 23(4), 954–966. <https://doi.org/10.1111/j.1748-7692.2007.00143.x>
- 609 Dalebout, M. L., Baker, C. S., Steel, D. J., Thompson, K. F., Robertson, K. M., Chivers, S. J., ...
610 Yamada, T. K. (2014). Resurrection of *Mesoplodon hotaula* Deraniyagala 1963: A new
611 species of beaked whale in the tropical Indo-Pacific. *Marine Mammal Science*, 30(3), 1081–
612 1108. <https://doi.org/10.1111/mms.12113>
- 613 Dalebout, M. L., Mead, J. G., Baker, C. S., Baker, A. N., & van Helden, A. L. (2002). A new species
614 of beaked whale *Mesoplodon perrini* sp. n. (Cetacea: Ziphiidae) discovered through
615 phylogenetic analyses of mitochondrial DNA sequences. *Marine Mammal Science*, 18(3),

577–608. <https://doi.org/10.1111/j.1748-7692.2002.tb01061.x>

617 Dalebout, M. L., Robertson, K. M., Frantzis, A., Engelhaupt, D., Mignucci-Giannoni, A. A.,
618 Rosario-Delestre, R. J., & Baker, C. S. (2005). Worldwide structure of mtDNA diversity
619 among Cuvier's beaked whales (*Ziphius cavirostris*): Implications for threatened
620 populations. *Molecular Ecology*, *14*(11), 3353–3371. <https://doi.org/10.1111/j.1365-294X.2005.02676.x>

622 Dalebout, M. L., Ross, G. J. B., Baker, C. S., Anderson, R. C., Best, P. B., Cockcroft, V. G., ...
623 Pitman, R. L. (2003). Appearance, Distribution, and Genetic Distinctiveness of Longman's
624 Beaked Whale, *Indopacetus pacificus*. *Marine Mammal Science*, *19*(July), 421–461.

625 Dalebout, M. L., Ruzzante, D. E., Whitehead, H., & Øien, N. I. (2006). Nuclear and mitochondrial
626 markers reveal distinctiveness of a small population of bottlenose whales (*Hyperoodon*
627 *ampullatus*) in the western North Atlantic. *Molecular Ecology*, *15*(11), 3115–3129.
628 <https://doi.org/10.1111/j.1365-294X.2006.03004.x>

629 Dalebout, M. L., Steel, D. J., & Baker, C. S. (2008). Phylogeny of the beaked whale genus
630 *Mesoplodon* (Ziphiidae: Cetacea) revealed by nuclear introns: Implications for the
631 evolution of male tusks. *Systematic Biology*, *57*(6), 857–875.
632 <https://doi.org/10.1080/10635150802559257>

633 Danecek, P., Auton, A., Abecasis, G., Albers, C. A., Banks, E., DePristo, M. A., ... Durbin, R. (2011).
634 The variant call format and VCFtools. *Bioinformatics*, *27*(15), 2156–2158.
635 <https://doi.org/10.1093/bioinformatics/btr330>

636 Dray, S., & Dufour, A. B. (2007). The ade4 package: Implementing the duality diagram for
637 ecologists. *Journal of Statistical Software*, *22*(4), 1–20.
638 <https://doi.org/10.18637/jss.v022.i04>

639 Durand, E., Jay, F., Gaggiotti, O. E., & François, O. (2009). Spatial inference of admixture
640 proportions and secondary contact zones. *Molecular Biology and Evolution*, *26*(9), 1963–
641 1973. <https://doi.org/10.1093/molbev/msp106>

642 Gomerčić, H., Duras Gomerčić, M., Gomerčić, T., Lucić, H., Dalebout, M. L., Galov, A., ... Huber,
643 D. (2006). Biological aspects of Cuvier's beaked whale (*Ziphius cavirostris*) recorded in the
644 Croatian part of the Adriatic Sea. *European Journal of Wildlife Research*, *52*(3), 182–187.
645 <https://doi.org/10.1007/s10344-006-0032-8>

646 Jombart, T., Devillard, S., & Balloux, F. (2010). Discriminant analysis of principal components: a
647 new method for the analysis of genetically structured populations. *BMC Genetics*, *11*(94),
648 1–15. <https://doi.org/10.1186/1471-2156-11-94>

649 Kamvar, Z. N., Brooks, J. C., & Grünwald, N. J. (2015). Novel R tools for analysis of genome-wide
650 population genetic data with emphasis on clonality. *Frontiers in Genetics*, *6*(JUN), 1–10.
651 <https://doi.org/10.3389/fgene.2015.00208>

652 Krützen, M., Barré, L. M., Möller, L. M., Heithaus, M. R., Simms, C., & Sherwin, W. B. (2002). A
653 biopsy system for small cetaceans: Darting success and wound healing in *Tursiops* spp.
654 *Marine Mammal Science*, *18*(4), 863–878. <https://doi.org/10.1111/j.1748-7692.2002.tb01078.x>

656 Lambertsen, R. H. (1987). A biopsy system for large whales and its use for cytogenetics. *Journal*
657 *of Mammalogy*, *68*(2), 443–445.

658 Maruki, T., & Lynch, M. (2015). Genotype-frequency estimation from high-throughput

659 sequencing data. *Genetics*, 201(2), 473–486. <https://doi.org/10.1534/genetics.115.179077>
660 Maruki, T., & Lynch, M. (2017). Genotype calling from population-genomic sequencing data. *G3: Genes, Genomes, Genetics*, 7(5), 1393–1404. <https://doi.org/10.1534/g3.117.039008>
661 O’Leary, S. J., Puritz, J. B., Willis, S. C., Hollenbeck, C. M., & Portnoy, D. S. (2018). These aren’t
662 the loci you’e looking for: Principles of effective SNP filtering for molecular ecologists.
663 *Molecular Ecology*, (June), 3193–3206. <https://doi.org/10.1111/mec.14792>
664 Pante, E., & Simon-Bouhet, B. (2013). marmap: A Package for Importing, Plotting and Analyzing
665 Bathymetric and Topographic Data in R. *PLoS ONE*, 8(9), 6–9.
666 <https://doi.org/10.1371/journal.pone.0073051>
667 Peterson, B. K., Weber, J. N., Kay, E. H., Fisher, H. S., & Hoekstra, H. E. (2012). Double digest
668 RADseq: An inexpensive method for de novo SNP discovery and genotyping in model and
669 non-model species. *PLoS ONE*, 7(5). <https://doi.org/10.1371/journal.pone.0037135>
670 R Core Team. (2019). R: A Language and Environment for Statistical Computing. Vienna, Austria:
671 R Foundation for Statistical Computing.
672 Rochette, N. C., Rivera-Colón, A. G., & Catchen, J. M. (2019). Stacks 2: Analytical methods for
673 paired-end sequencing improve RADseq-based population genomics. *Molecular Ecology*,
674 28(21), 4737–4754. <https://doi.org/10.1111/mec.15253>
675 Sambrook, J., Fritsch, E. F., & Maniatis, T. (1989). *Molecular Cloning, A Laboratory Manual. Second Edition*. Cold Spring Harbor Laboratory Press, Cold Spring Harbor, New York.
676 Van Helden, A. L., Baker, A. N., Dalebout, M. L., Reyes, J. C., Van Waerebeek, K., & Baker, C. S.
677 (2002). Resurrection of *Mesoplodon traversii* (Gray, 1874), senior synonym of *M.*
678 *bahamondi* Reyes, Van Waerebeek, Cardenas and Yanez, 1995 (Cetacea: Ziphiidae). *Marine*
679 *Mammal Science*, 18(July), 609–621.
680 Yu, G., Smith, D. K., Zhu, H., Guan, Y., & Lam, T. T. Y. (2017). GGTREE: an R package for
681 visualization and annotation of phylogenetic trees with their covariates and other
682 associated data. *Methods in Ecology and Evolution*, 8, 28–36.
683 <https://doi.org/10.1111/2041-210X.12628>
684
685
686
687