

SC/66a/BRG/12

Continued development of MTDNA and SNP databases for bowhead whales

Amy B. Baird, Alesha R. Rimmelin, Christina Q. Tran,
John C. George, Robert S. Suydam, and John W.
Bickham



INTERNATIONAL
WHALING COMMISSION

1 **CONTINUED DEVELOPMENT OF MTDNA AND SNP DATABASES FOR BOWHEAD WHALES**

2 Amy B. Baird¹, Alesha R. Rimmelin¹, Christina Q. Tran¹, John C. George², Robert S. Suydam²,
3 and John W. Bickham³

4 ¹Department of Natural Sciences, University of Houston – Downtown, 1 Main St., Houston, TX
5 77002, USA; ²North Slope Borough, Department of Wildlife Management, Barrow, AK 99723,
6 USA; ³Battelle Memorial Institute, Houston, TX, USA

7 **ABSTRACT**

8 This paper details the current status of ongoing genetic studies of bowhead whales. Our research
9 has two main objectives: building a mitochondrial DNA (mtDNA) database and the development
10 of a single nucleotide polymorphism (SNP) panel for bowheads. Here we report the current
11 status of the mtDNA database, as well as the pros and cons of several methods of SNP
12 genotyping we tested on bowheads. Our mtDNA database currently includes sequences from
13 711 whales (447 complete cytochrome-b sequences, 427 complete ND1 sequences, and 638
14 HVR1 sequences). Of these, 345 whales have been sequenced for all three loci. We report on
15 the methods and criteria used to select a panel of 96 bowhead SNPs, which currently are being
16 analyzed.

17 **INTRODUCTION**

18 Genetic data have played an important role in bowhead whale conservation, management and in
19 establishing an aboriginal quota through the IWC to meet nutritional and cultural needs of Inuit
20 people in Alaska, USA, and Chukotka, Russia. The long-term acquisition of mitochondrial DNA
21 (mtDNA) data has been integral in studies examining stock structure and also provides a long
22 term perspective on population size estimates as a means of estimating effective population size,
23 N_e . Initially, studies on bowhead mtDNA focused on a single locus, the control region (e.g.
24 Rooney et al. 2001, LeDuc et al. 2005). Recent studies have expanded the mtDNA sequencing
25 to two additional loci: ND1 and cyt-b (Phillips et al. 2012). The expansion of the mtDNA
26 database has increased the number of variable, informative sites, which has allowed more
27 sophisticated analyses to be conducted, with more accurate results obtained. This, in turn, has
28 shed more light on the evolutionary history and population structure of bowhead whales.
29 Specifically, Phillips et al. (2012) used the 3 mtDNA loci to examine historical demography of
30 bowheads to investigate the impact of commercial whaling and concluded that no genetic
31 signature of recent population depletion was observed. The same data were able to show
32 evidence, however, for historical population expansion and subsequent decline due most likely to
33 Pleistocene climate change. The analyses using data from the 3 mtDNA loci provided a much
34 more accurate picture of historical relative population size estimates than similar studies using
35 only the mtDNA control region (Rooney et al. 1999 and Rooney et al. 2001).

36 Moreover, nuclear data are an important source of genetic information to answer questions that
37 mtDNA alone cannot address, such as genetic interchange between populations, and more
38 accurate estimates of relatedness between individuals, and population subdivision. Previous
39 studies of bowheads have utilized different nuclear markers to address these questions. In the
40 past, microsatellites were used with good success (LeDuc et al. 2005, Givens et al. 2010);
41 however, these markers can be problematic in several ways. Microsatellites are considered to
42 have a high error rate and are not always replicable between labs. More recently, single
43 nucleotide polymorphisms (SNPs) have been used on bowheads with good success. When
44 directly compared to microsatellite data in bowheads, SNPs were shown to be successful at
45 answering questions about population size and structure (Morin et al. 2012). Morin et al. (2012)
46 concluded that SNPs are the preferred method for genotyping bowheads because of their ability
47 to be replicated across labs, their ability to be used on old or degraded samples, and their cost
48 effectiveness.

49 This project has two objectives. The first is to continue to build a mtDNA database that includes
50 sequences from three genes. We continue to sequence the commonly used mtDNA control
51 region, as well as two protein coding genes (cytochrome-b and ND1). The data presented here
52 build upon mtDNA sequence data presented in SC/64/AWMP9 (Bickham et al. 2012) and
53 SC/65b/BRG13 (Baird et al. 2014).

54 The second focus of this project is on the continued development of a SNP panel. Previous
55 studies have reported on the identification of these variable loci for bowheads, including
56 SC/65a/BRG22 and SC/65b/BRG13. The most recent of these papers described the 155 SNP
57 loci that have been identified for bowheads. Here, the goal was to select a panel of 96 of those
58 previously identified loci and determine the most cost-effective method to analyze them for large
59 sample sizes of whales.

60 The combination of mtDNA and SNP data, along with the genetic resources database provided
61 by the bowhead genome project (Seim et al., 2014; Keane et al., 2015), will allow us to continue
62 addressing issues related to bowhead stock structure (such as estimating potential genetic
63 interchange between BCB and Canadian or Sea of Okhotsk populations), revisit issues of
64 historical demography and molecular ecology, and continue to monitor the genetic health of the
65 BCB bowhead population.

66 **METHODS**

67 *New mtDNA sequence data acquisition.* – Bowhead whale DNA was extracted, amplified, and
68 sequenced following the methods of LeDuc et al. (2008) and Phillips et al. (2011). The hyper
69 variable region-1 (HVR1) region of the mtDNA control region was sequenced (397bp) along
70 with the complete cytochrome-b (cytb) gene (1140bp) and the complete ND1 gene (957bp).
71 Individual sequences were compared to existing haplotypes from previous studies deposited on
72 GenBank and haplotype codes were assigned.

73 *Single nucleotide polymorphism (SNP) selection.* – We examined all 155 SNP loci that were
74 reported in previous papers to determine the best loci to use in a SNP panel. Our evaluation of
75 each locus included noting position of the SNP in the genome; proximity of a given SNP to any
76 other known SNPs; proximity of the SNP to repeated genetic elements that might make
77 amplification difficult; and whether the SNP was a base substitution or an insertion/deletion.

78 Because some of our overall goals with the bowhead genetics program are to estimate current
79 and historical effective population size and examine whether there is genetic interchange
80 between stocks, a priority was to find SNP markers that would inform these questions. In
81 combination with our mtDNA data, some of the most informative markers for these questions
82 will be SNPs from the sex chromosomes. Therefore, we chose all usable variable loci on the X
83 chromosome (12 loci) and Y chromosome (1 locus).

84 Next, we wanted to minimize the linkage between any of the autosomal SNPs that we chose, and
85 we therefore prioritized SNPs that were not reported from the same amplicons in the papers that
86 originally described the SNP loci.

87 *Selection of SNP method of analysis.*—We tested three different methods of SNP genotyping and
88 considered a fourth method without testing it. The methods tried included TaqMan genotyping,
89 High Resolution Melting analysis and the Fluidigm SNPtype assay. We considered the Illumina
90 MiSeq method, but it was ruled out before we attempted to use it (see below). To select which
91 method to use, we considered several factors: cost per sample; ease of data interpretation;
92 reliability upon replication; and labor intensity.

93 TaqMan genotyping was performed using an Eppendorf real time thermal cycler using 6 loci
94 described by Morin et al. (2010).

95 The High Resolution Melting analysis was performed using the same thermal cycler on 4 loci.

96 Fluidigm SNPtype assay primers were designed for 96 SNPs using a combination of loci
97 described in Morin et al. (2010), SC/65b/BRG13 (Baird et al. 2014), and SC/65a/BRG22
98 (Bickham et al. 2013).

99 **RESULTS**

100 *Mitochondrial DNA database.* – To date, we have samples representing 711 whales. Sequences
101 in the database include data from 447 whales with complete cyt-b sequences, 427 whales with
102 complete ND1 sequences, and 638 whales with complete HVR1 sequences. Of these, 345
103 whales have been sequenced for all three loci. The majority of our samples come from the BCB
104 stock, but we also have representatives from the Eastern Canadian Arctic and Sea of Okhotsk.

105 *Selection of SNP method of analysis.*—We initially considered the Illumina MiSeq method of
106 SNP genotyping (<http://www.illumina.com/>). This method is a type of next-generation
107 sequencing that can produce 2x300 bp paired-end reads in one run, allowing for the collection of

108 a large amount of data in a single run. For our purposes, we would initially need to design
109 primers for each SNP locus that would be compatible with the MiSeq platform. Depending on
110 the service company, it would be possible to design primers for up to 90 SNPs that could all be
111 run on a single run for 96 samples. The benefits of this method include the ability to run a large
112 number of samples/SNPs at a time and that many service companies have the ability to use this
113 platform. The major downside of this method was the cost. For the development of ~90 primers
114 there would be a one-time cost of \$15,000. Then one would need to amplify each sample for
115 each primer pair, resulting in significant costs for consumables. Finally, each miSeq run of 90
116 primers x 96 samples costs approximately ~\$3500 (depending on the service company).

117 We then considered the Life Technologies TaqMan genotyping method
118 [http://www.lifetechnologies.com/us/en/home/life-science/pcr/real-time-pcr/real-time-pcr-
120 assays/snp-genotyping-taqman-assays.html](http://www.lifetechnologies.com/us/en/home/life-science/pcr/real-time-pcr/real-time-pcr-
119 assays/snp-genotyping-taqman-assays.html)). This method is a real-time PCR-based method that
121 works by using allele-specific primers each tagged with a unique fluorescence. Upon annealing
122 and extension of an allele-specific primer, it (or they, in the case of a heterozygote) will
123 fluoresce. We initially tested the method on 6 loci described in Morin et al. 2010. The
124 advantages of this method for our purposes included the ability to run each locus individually, so
125 it was not a large cost burden up front. We could also do all of the work in-house because we
126 own a real-time PCR machine. However, we had several reservations about the feasibility of this
127 method for our purposes. Although the up-front costs per locus were not burdensome (each
128 custom primer pair cost \$242.50 which was enough product to run 1500 samples), when summed
129 for the number of loci needed to produce a usable panel of SNPs for a large number of samples,
130 as well as costs associated with real-time PCR enzyme master mixes, the costs for this method
131 would quickly accumulate. More worrisome for us was that we could not reliably replicate our
132 data with this method. We are not completely sure of the cause of this, but we suspect it had
133 more to do with the real-time PCR machine we were using (Eppendorf) rather than the TaqMan
134 technology itself. Because we were trying to detect single base-pair changes, the machine needs
135 to be very sensitive to those changes in fluorescence. We believe that our machine was unable to
136 accurately distinguish between different bases at a given SNP in many instances.

137 Next, we tested another real-time PCR-based method: high resolution melting analysis. We used
138 the SensiFAST HRM Kit from Bioline for our test runs ([http://www.bioline.com/us/sensifast-
140 hrm-kit.html](http://www.bioline.com/us/sensifast-
139 hrm-kit.html)). This method works by amplifying the SNP target (and surrounding sequence),
141 then slowly melting the amplicon. Amplicons with different sequences will melt at slightly
142 different temperatures. The kit cost is quite modest (~\$100), and has enough to run 200
143 reactions, or a larger kit for \$1050 (for 2000 reactions). The kit can then be used with normal
144 primers, which are significantly cheaper than fluorescent primers. The major benefits of this
145 method included a much lower cost than the methods we tested, and it could be done in-house
146 with our real-time PCR machine. The downside was that, again, we could not reliably replicate
147 the results. We again suspect that the problem may have had more to do with our PCR machine
148 than with the chemistry/technology of the method.

147 Finally, after testing the other methods and consulting with colleagues who have also searched
148 for a cost-effective and reliable method for SNP genotyping, we settled on the Fluidigm SNPtype
149 assay. This method is based on allele-specific PCR SNP detection. It requires a Fluidigm
150 platform to run/read the assay, requiring samples to be run at a service facility. The method is
151 quite cost effective, fast, and results are reproducible (according to colleagues who have tested it
152 extensively). The method can handle multiplexing SNP loci in groups of 96 or 192 for 96
153 samples at a time. Each set of 96 loci costs just under \$5000 and includes enough product to run
154 about 14,000 samples, so this cost is a one-time expense. Then, each run of 96 samples (for all
155 96 loci) costs approximately \$1300 (depending on the service company). There is essentially no
156 labor involved if they are run by a service company, other than extracting and quantifying the
157 DNA.

158 *Selection of SNP loci.*—After having settled on the Fluidigm method, the most cost effective
159 approach was to design primers for 96 loci of the 155 SNPs identified for bowheads. The first
160 step was primer design (described below) and quality control.

161 Fluidigm’s primer design software required having a minimum of 60bp sequence both upstream
162 and downstream of the desired SNP. Indels could not be greater than 10bp. A desired SNP
163 could also not be within 30bp of another variable locus. Using these criteria, along with our own
164 criteria for our desired SNP loci (see Methods), we designed primers using the Fluidigm
165 software.

166 First, we acquired all available sequence data surrounding the SNPs of interest. These data were
167 acquired from GenBank (for loci described in Morin et al. 2010) or from positional information
168 reported in previous papers and aligned to the bowhead whale genome (Keane et al. 2015;
169 <http://www.bowhead-whale.org/>). Several of the 155 available SNPs were described from single
170 amplicons, which placed some SNPs close to one another (sometimes less than 30bp apart). To
171 avoid linkage between autosomal SNPs, we first only selected one SNP per amplicon. The
172 flanking sequence was then submitted to Fluidigm for quality control.

173 A small number of our initially submitted loci were rejected by Fluidigm based on quality
174 control concerns. These included a target close to an area of “repeats and/or difficult genomic
175 diversity” and indels being too complicated. To then get to the target of 96 loci, this required us
176 to select some SNPs that would have otherwise not met our criteria, particularly regarding the
177 concern over linkage. We ended up selecting 9 autosomal loci that were located in proximity to
178 other SNPs in the panel. In total, we selected 13 X- and Y-chromosome SNPs and 83 autosomal
179 SNPs.

180 **DISCUSSION**

181 The continued expansion of a mitochondrial DNA database for bowhead whales has been, and
182 will continue to be, a valuable tool in analyzing population size and structure for this species.
183 The inclusion of two protein-coding genes has proven to add considerable value to the mtDNA

184 database, rather than including only the HVR1 region of the control region (Bickham et al.
185 2012). We are continuing to expand this database, with the eventual goal of having all samples
186 sequenced for all three loci.

187 Having selected a cost-effective and reliable method for SNP analysis will also result in valuable
188 data to inform our analyses of population and evolutionary genetics of BCB bowheads, as well as
189 examining the potential for genetic interchange between BCB and other populations. These data
190 will be important for having a firm understanding for what stock boundaries are and whether
191 there is any gene flow between them. This is clearly an important issue for management and
192 conservation in this species.

193 The 96-locus SNP panel we designed has utilized the best data and ascertainment schemes we
194 have available to date. The 42 SNPs identified by Morin et al. (2010) were obtained through a
195 random locus approach and a targeted gene approach using bowhead samples from eastern
196 Russia, St. Lawrence Island, Barrow, Alaska and eastern Canada. The sex chromosome SNPs
197 from SC/65a/BRG22 (Bickham et al. 2013) were obtained using a targeted gene approach (exon-
198 primed intron crossing) from Alaskan whales. The remaining autosomal SNPs described in
199 SC/65b/BRG13 (Baird et al. 2014) came from examining variable sites in the bowhead
200 transcriptome, which was based on multiple individuals of Alaskan whales. Therefore, our
201 subset of 96 SNPs contains loci identified by a variety of methods and from a comparison of
202 whales from multiple populations.

203 SNP ascertainment schemes can significantly affect conclusions regarding population history
204 (McTavish and Hillis 2015). This has long been known to be true of other sources of population
205 genetics data, including microsatellites (Barbara et al., 2007; Putman and Carbone, 2014).
206 Biased selection of SNP loci can affect estimates of F_{ST} , population admixture, and population
207 histories. Ascertainment schemes can be potentially biased toward samples from a geographic
208 region (which would overestimate polymorphisms in that region) or toward loci polymorphic in
209 multiple groups. Simulated data presented in McTavish and Hillis (2015) demonstrated that
210 ascertainment schemes inflating geographic polymorphisms led to decreased differentiation
211 among populations using F_{ST} . When the ascertainment bias favors shared polymorphisms among
212 multiple groups, population differentiation was obscured.

213 Keeping in mind the above caveats regarding ascertainment bias, we sought to include SNPs
214 identified from multiple populations whenever possible. A total of 33 of the 96 SNPs in our
215 panel come from the Morin et al. (2010) study which used whales from multiple populations for
216 SNP identification. The remaining SNPs come from studies which used only individuals from
217 the BCB population. To improve this potential bias toward BCB polymorphisms, we plan to
218 develop a second 96-SNP panel that will be derived from bowhead genome/transcriptome data
219 containing data from both Alaskan and Greenland whales (Keane et al. 2015;
220 <http://www.bowhead-whale.org/>).

221 **ACKNOWLEDGMENTS**

222 We thank the Alaska Eskimo Whaling Commission (AEWC) and the 11 villages of the Whaling
223 Captains' Associations for their confidence, guidance, and support of our research. We
224 gratefully acknowledge funding provided by the North Slope Borough Department of Wildlife
225 Management and the many people who helped collect the tissue samples for this work. We
226 thank Donald Stoeckel and Andrew DeWoody for sharing their experience with different
227 methods of SNP genotyping.

228 **LITERATURE CITED**

- 229 Baird, A. B., R. S. Suydam, J. C. George, and J. W. Bickham. 2014. Update on mtDNA and
230 SNP database for bowhead whales. Paper SC/65b/BRG13 submitted to the International
231 Whaling Commission Scientific Committee.
- 232 Bickham, J. W., R. M. Huebinger, C. D. Philips, J. C. Patton, L. D. Postma, J. C. George, and R.
233 S. Suydam. 2012. Assessing molecular substitution patterns in mitochondrial control
234 region compared to protein coding genes in bowhead whales: update of SC/63/BRG13.
235 Paper SC/64/AWMP9 submitted to the International Whaling Commission Scientific
236 Committee.
- 237 Bickham, J. W., G. W. Stuart, H. K. Downing, J. C. Patton, J. C. George, and R. S. Suydam.
238 2013. Comparison of methods for molecular assessment of sex chromosome
239 polymorphisms and levels of genetic divergence in the bowhead whale. Paper
240 SC/65a/BRG22 submitted to the International Whaling Commission Scientific
241 Committee.
- 242 Barbara T, Palmer-Silva C, Paggi GM, Bered F, Fay MF, et al. (2007) Cross-species transfer of
243 nuclear microsatellite markers: potential and limitations. *Mol Ecol* 16: 3759–3767.
- 244 Givens, G. H., R. M. Huebinger, J. C. Patton, L. D. Postma, M. Lindsay, R. S. Suydam, J. C.
245 George, C. W. Matson, and J. W. Bickham. 2010. Population genetics of bowhead
246 whales (*Balaena mysticetus*) in the western Arctic. *Arctic* 63: 1-12.
- 247 Keane M, Semeiks J, Webb AE, Li YI, Quesada V, Craig T, Madsen LB, van Dam S, Brawand
248 D, Marques PI, Michalak P, Kang L, Bhak J, Yim HS, Grishin NV, Nielsen NH, Heide-
249 Jorgensen MP, Oziolor EM, Matson CW, Church GM, Stuart GW, Patton JC, George JC,
250 Suydam R, Larsen K, Lopez-Otin C, O'Connell MJ, Bickham JW, Thomsen B, de
251 Magalhaes JP (2015) Insights into the evolution of longevity from the bowhead whale
252 genome. *Cell Reports*. 10, 112-120.
- 253 LeDuc RG, Martien KK, Morin PA, Hedrick N, Robertson KM, Taylor BL, Mugue NS, Borodin
254 RG, Zelenina DA, Litovka D, George JC (2008) Mitochondrial genetic variation in

255 bowhead whales in the western Arctic. *Journal of Cetacean Research and Management*
256 10, 93–97

257 LeDuc, R. G., A. E. Dizon, A. M. Burdin, S. A. Blokhin, J. C. George, and R. L. Brownell, Jr.
258 2005. Genetic analysis (mtDNA and microsatellites) of Okhotsk and
259 Bering/Chukchi/Beaufort Seas populations of bowhead whales. *Journal of Cetacean*
260 *Resource Management* 7: 107-111.

261 McTavish, E. J., and D. M. Hillis. 2015. How do SNP ascertainment schemes and population
262 demographics affect inferences about population history? *BMC Genomics* 16:266.

263 Morin, P. A., V. L. Pease, B. L. Hancock, K. M. Robertson, C. W. Antolik. 2010.
264 Characterization of 42 single nucleotide polymorphism (SNP) markers for the bowhead
265 whale (*Balaena mysticetus*) for use in discriminating populations. *Marine Mammal*
266 *Science* 26: 716-732.

267 Morin, P. A., F. I. Archer, V. L. Pease, B. L. Hancock-Hanser, K. M. Robertson, R. M.
268 Huebinger, K. K. Martien, J. W. Bickham, J. C. George, L. D. Postma, B. L. Taylor.
269 2012. Empirical comparison of single nucleotide polymorphisms and microsatellites for
270 population and demographic analyses of bowhead whales. *Endangered Species Research*
271 19: 129-147.

272 Phillips, C. D., T. S. Gelatt, J. C. Patton, and J. W. Bickham. 2011. Phylogeography of Steller
273 sea lions: relationships among climate change, effective population size, and genetic
274 diversity. *Journal of Mammalogy* 92:1091-1104.

275 Phillips, C. D., J. I. Hoffman, J. C. George, R. S. Suydam, R. M. Huebinger, J. C. Patton, and J.
276 W. Bickham. 2012. Molecular insights into the historic demography of bowhead
277 whales: understanding the evolutionary basis of contemporary management practices.
278 *Ecology and Evolution* 1-20.

279 Putman, A. I. and I. Carbone. 2014. Challenges in analysis and interpretation of microsatellite
280 data for population genetic studies. *Ecology and Evolution* 4: 4399-4428.

281 Rooney, A. P., R. L. Honeycutt, S. K. Davis, J. N. Derr. 1999. Evaluating a putative bottleneck
282 in a population of bowhead whales from patterns of microsatellite diversity and genetic
283 disequilibria. *Journal of Molecular Evolution* 49: 682-690.

284 Rooney, A. P., R. L. Honeycutt, and J. N. Derr. 2001. Historical population size change of
285 bowhead whales inferred from DNA sequence polymorphism data. *Evolution* 55: 1678-
286 1685.

287 Seim, I., S. Ma, X. Zhou, M. V. Gerashchenko, S.-G. Lee, R. Suydam, J. C. George, John W.
288 Bickham, and V. N. Gladyshev. 2014. The transcriptome of the bowhead whale *Balaena*
289 *mysticetus* reveals adaptations of the longest-lived mammal. *Aging* 6:879-899.