# Edinburgh Research Explorer

## Genomic methods take the plunge: recent advances in high-throughput sequencing of marine mammals

**Published In:**
Journal of Heredity

OXFORD
UNIVERSITY PRESS | Journal of Heredity

# Genomic methods take the plunge: recent advances in high-throughput sequencing of marine mammals

SCHOLARONE™
Manuscripts

1
2
3
4  **Genomic methods take the plunge: recent advances in high-throughput sequencing of**
5  **marine mammals**
6
7

4  KRISTINA M. CAMMEN[1]*, KIMBERLY R. ANDREWS[2], EMMA L. CARROLL[3], ANDREW D.
5  FOOTE[4], EMILY HUMBLE[5,6], JANE I. KHUDYAKOV[7], MARIE LOUIS[3], MICHAEL R.
6  MCGOWEN[8], MORTEN TANGE OLSEN[9], AND AMY M. VAN CISE[10]
7

8  [1]School of Marine Sciences, University of Maine, Orono, Maine 04469, USA
9  [2]Department of Fish and Wildlife Sciences, University of Idaho, 875 Perimeter Drive MS 1136,
10 Moscow, Idaho 83844-1136, USA
11 [3]Scottish Oceans Institute, University of St Andrews, East Sands, St Andrews, Fife KY16 8LB,
12 UK
13 [4]Computational and Molecular Population Genetics Lab, Institute of Ecology and Evolution,
14 University of Bern, Bern CH-3012, Switzerland
15 [5]Department of Animal Behaviour, University of Bielefeld, Postfach 100131, 33501 Bielefeld,
16 Germany
17 [6]British Antarctic Survey, High Cross, Madingley Road, Cambridge CB3 OET, UK
18 [7]Department of Biology, Sonoma State University, Rohnert Park, California 94928, USA
19 [8]School of Biological and Chemical Sciences, Queen Mary University of London,
20 Mile End Road, London E1 4NS, UK
21 [9]Evolutionary Genomics Section, Natural History Museum of Denmark, University of
22 Copenhagen, DK-1353 Copenhagen K, Denmark
23 [10]Scripps Institution of Oceanography, 8622 Kennel Way, La Jolla, California 92037, USA
24
25 *Corresponding author: kristina.cammen@maine.edu
26
27 Running title: Marine mammal genomics

28 **Abstract**

29 The dramatic increase in the application of genomic techniques to non-model organisms over the

30 past decade has yielded numerous valuable contributions to evolutionary biology and ecology,

31 many of which would not have been possible with traditional genetic markers. We review this

32 recent progression with a particular focus on genomic studies of marine mammals, a group of

33 taxa that represent key macroevolutionary transitions from terrestrial to marine environments and

34 for which available genomic resources have recently undergone notable rapid growth. Genomic

35 studies of non-model organisms utilize an expanding range of approaches, including whole

36 genome sequencing, restriction site-associated DNA sequencing, array-based sequencing of

37 single nucleotide polymorphisms and target sequence probes (e.g., exomes), and transcriptome

38 sequencing. These approaches generate different types and quantities of data, and many can be

39 applied with limited or no prior genomic resources, thus overcoming one traditional limitation of

40 research on non-model organisms. Within marine mammals, such studies have thus far yielded

41 significant contributions to the fields of phylogenomics and comparative genomics, as well as

42 enabled investigations of fitness, demography, and population structure. Here we review the

43 primary options for generating genomic data, introduce several emerging techniques, and discuss

44 the suitability of each approach for different applications in the study of non-model organisms.

45

46 **Keywords**: RADseq, SNP array, target sequence capture, whole genome sequencing, RNAseq,

47 non-model organisms

48 **Introduction**

49 Recent advances in sequencing technologies, coincident with dramatic declines in cost, have

50 increasingly enabled the application of genomic sequencing in non-model systems (Ekblom and

51 Galindo 2011; Ellegren 2014). These advances in molecular technologies have in many ways

52 begun to blur the distinction between model and non-model organisms (Armengaud et al. 2014).

53 Non-model organisms (NMOs) have traditionally been defined as those for which whole-

54 organism experimental manipulation is rarely, if ever, possible due to logistical and/or ethical

55 constraints (Ankeny and Leonelli 2011). Further, NMOs have typically been characterized by

56 limited genomic resources, but this is becoming increasingly less so as the number of NMO

57 reference genomes grows rapidly, for example through efforts like the Genome 10K Project

58 (Koepfli et al. 2015). In fact, in some taxonomic orders, we are approaching the point at which

59 all species have at least one representative reference genome available for a closely related

60 species (Fig 1).

61

62 Despite the limitations of working with NMOs, including potentially small sample sizes, low

63 DNA quantity, and limited information on gene function, genetic and genomic investigations of

64 NMOs have yielded numerous valuable contributions to understanding their evolutionary

65 biology and ecology. For the past several decades, traditional genetic markers such as

66 microsatellites and short fragments of mitochondrial DNA (e.g., the control region) have been

67 extensively used in molecular ecology. These markers, which typically evolve under neutral

68 expectations, have proven useful for identifying population structure and reconstructing

69 population demographic history (Hedrick 2000). However, the power of such studies is limited

70 by the number of markers that can feasibly be evaluated using traditional approaches. The advent

71 of low-cost high-throughput sequencing has led to dramatic increases in the number of neutral

72 markers that can be evaluated, in many cases improving our power to resolve fine-scale or

73 cryptic population structure in species with high dispersal capability (e.g., Corander et al. 2013)

74 and improving the accuracy of estimating some (though not all) demographic parameters (Li and

75 Jakobsson 2012; Shafer et al. 2015). Importantly, high-throughput sequencing has also further

76 enabled genomic studies of non-neutral processes in NMOs, for example, characterizing both

77 deleterious and adaptive variation within and across species (Stinchcombe and Hoekstra 2008;

78   Künstner et al. 2010). It is increasingly evident that genomic analyses of NMOs can and have

79   provided important insights that could not be identified with traditional genetic markers.

80

81   Many molecular ecologists now face the challenge of deciding which of the broad range of

82   genomic approaches to apply to their study systems. Here we review the primary options for

83   generating genomic data and their relative suitability for different applications in the study of

84   NMOs. We focus on marine mammals, which represent several mammalian clades with notably

85   rapid growth in available genomic resources in recent years. This growth is clearly evident in

86   both publication rate (Fig 2) and the rise in number and size of genomic sequences deposited in

87   public resources (Fig 3). We comprehensively review the literature on marine mammal

88   genomics, highlighting recent trends in methodology and applications, and then describe in detail

89   the molecular approaches that are most commonly applied to studies of NMO genomics. Our

90   hope is that this review will highlight the promise of genomics for NMOs and offer guidance to

91   researchers considering the application of genomic techniques in their non-model study system

92   of choice.

93

94   **Why study marine mammal genomics?**

95   Marine mammals represent key macroevolutionary transitions from terrestrial to marine

96   environments (McGowen et al. 2014) and accordingly are an exemplary system for investigating

97   the evolution of several morphological and physiological adaptations (Foote et al. 2015)

98   associated with locomotion (Shen et al. 2012), sight (Meredith et al. 2013), echolocation (Parker

99   et al. 2013; Zou and Zhang 2015), deep diving (Mirceta et al. 2013), osmoregulation (Ruan et al.

100  2015), and cognition (McGowen et al. 2012). Furthermore, studies of marine mammal evolution

101  to date have characterized several unique aspects of their genome evolution that merit further

102  investigation, including low genomic diversity and a relatively slow molecular clock, especially

103  in cetaceans (Jackson et al. 2009; McGowen et al. 2012; Zhou et al. 2013). As many cetacean

104  species are highly mobile with no obvious physical geographic barriers to dispersal, they provide

105  a unique opportunity to study the role of behavior and culture in shaping population structure and

106  genetic diversity (Riesch et al. 2012; Carroll et al. 2015; Alexander et al. 2016). Though highly

107  mobile, many marine mammals exhibit evidence of local adaptation; for example, several species

108  show parallel divergent morphological and behavioral adaptations to coastal and pelagic

109    environments (Moura et al. 2013; Louis et al. 2014; Viricel and Rosel 2014). These species may

110    be studied across ocean basins as emerging examples of ecological adaptation and speciation

111    (Morin et al. 2010a).

112

113    Beyond their value as systems of evolutionary study, many marine mammals are also of broader

114    interest relating to their historical and present conservation status. Many marine mammal

115    populations share histories of dramatic decline due to hunting and other human impacts.

116    Genomics provides a promising tool with which to expand our insights into these historical

117    population changes, which so far primarily have been derived from archival review and

118    traditional genetic approaches (Ruegg et al. 2013; Sremba et al. 2015). More recently, since the

119    implementation of national and international protections, many marine mammal populations

120    have partially or fully recovered (Magera et al. 2013), yet the conservation status of certain

121    marine mammal populations remains of concern. Such vulnerable populations could benefit

122    greatly from an improved understanding of their genetic diversity and evolution, especially in

123    ways that can inform predictions of adaptive capacity to anthropogenic pressures and expand the

124    toolkit for conservation policy (Garner et al. 2016; Taylor and Gemmell 2016).

125

126    **Recent trends in marine mammal genomics**

127    We conducted a meta-analysis of the peer-reviewed marine mammal genomics literature to

128    evaluate trends in publication rates across research methodologies and aims. A search of the Web

129    of Science database using the term "genom*" and one of the following terms indicating study

130    species - "marine mammal", "pinniped", "seal", "sea lion", "sea otter", "whale", "dolphin",

131    "polar bear", "manatee" - identified 825 records on December 11, 2015. We excluded 77% of the

132    search results that were not directly related to genomic studies in marine mammal systems. The

133    remaining 101 articles that were relevant to marine mammal genomics were further categorized

134    by primary research methodology and general research aim. A subset of these articles is

135    described briefly in Supplemental Table 1.

136

137    From the early 1990s through 2015, published literature in the field shifted from an early focus

138    on mitogenome sequencing to more sequence-intensive approaches, such as transcriptome and

139    whole genome sequencing (Figs 2 and 4). This trajectory closely follows trends in sequencing

140 technologies, from Sanger sequencing of short- and long-range PCR products for mitogenome

141 sequencing (Arnason et al. 1991) and SNP discovery (Olsen et al. 2011), to high-throughput

142 sequencing of reduced-representation genomic libraries (RRLs) that consist of selected subsets

143 of the genome (e.g., Viricel et al. 2014), to high-throughput sequencing of whole genomes with

144 varying levels of depth, coverage, and contiguity. Today, high-throughput sequencing can be

145 used both to generate high-quality reference genome assemblies (Yim et al. 2014; Foote et al.

146 2015; Humble et al. 2016) and to re-sequence whole genomes at a population scale (Liu et al.

147 2014a; Foote et al. 2016). Similarly, the scale of gene expression studies has increased from

148 quantitative real-time PCR of candidate genes (Tabuchi et al. 2006) to microarrays containing

149 hundreds to thousands of genes (Mancia et al. 2007) and high-throughput RNAseq that evaluates

150 hundreds of thousands of contigs across the genome (Khudyakov et al. 2015b). As the cost of

151 high-throughput sequencing continues to decline, we anticipate an increase in studies that

152 sequence RRLs, whole genomes, and transcriptomes in NMOs at a population scale.

153

154 Marine mammal genomic studies thus far have primarily contributed to the fields of

155 phylogenomics and comparative genomics (Fig 2, Table S1). Several of these comparative

156 genomics studies have aimed to improve our understanding of the mammalian transition to an

157 aquatic lifestyle and describe the evolutionary relationships within and among marine mammals

158 and their terrestrial relatives (McGowen et al. 2014; Foote et al. 2015). Whereas such studies

159 require only a single representative genome per species, an emerging class of studies applying

160 genomic techniques at a population scale enables further investigations of fitness, demography,

161 and population structure within species (Table S1). However, expanding the scale of genomic

162 studies requires careful selection of an appropriate method for data generation and analysis from

163 a growing number of approaches that are becoming available to non-model systems.

164

## Data generation

166 Our review of marine mammal genomics highlights an increasing number of options for the

167 generation and analysis of genomic data. Choosing which of these sequencing strategies to apply

168 is a key step in any genomics study. Here we describe approaches that have been used

169 successfully in order to help guide future studies of ecological, physiological, and evolutionary

170 genomics in NMOs. Across data generation methods, we highlight approaches that can be used

171   with limited or no prior genomic resources, overcoming one traditional challenge of genomic

172   studies of NMOs (the need for a reference genome to which sequencing reads can be mapped).

173   These methods produce a range in quantity and type of data output, from hundreds of SNPs to

174   whole genome sequences, and from single individuals to population samples, reflecting the

175   trade-off between number of samples and amount of data generated per sample.

176

177   <u>Sample collection, storage and extraction</u>

178   Prior to starting a genomic study, researchers must recognize that many recent methods for high-

179   throughput sequencing require genetic material of much higher quality and quantity than

180   techniques used to characterize traditional genetic markers. These more stringent sample

181   requirements necessitate new standards for tissue sampling, storage, and DNA/RNA extraction.

182   Ideally, samples should be collected from live or newly deceased individuals and stored at -80°C,

183   or when this is not possible at -20°C in RNAlater, Trizol, ethanol, salt-saturated DMSO, or dry,

184   depending on the intended application. Given the sensitivity of new sequencing methods, great

185   care should be taken to minimize cross-contamination during sampling, as even minute amounts

186   of genetic material from another individual can bias downstream analyses, for example variant

187   genotyping and gene expression profiles. Choice of extraction method varies with sample type

188   and study aim, but typically genomic methods require cleanup and treatment with RNase to yield

189   pure extracts, whereas RNAseq methods require rigorous DNase treatment to remove genomic

190   contamination that can bias expression results. Depending on the genomic methodology, target

191   quantities for a final sample may range from as low as 50 ng of DNA for some RRL sequencing

192   methods (Andrews et al. 2016) up to ~1 mg for sequencing the full set of libraries (of different

193   insert sizes) necessary for high-quality genome assemblies (Ekblom and Wolf 2014). Most

194   commercial RNAseq library preparation services require at least 500-1,000 ng of pure total RNA

195   that shows minimal degradation as measured by capillary gel electrophoresis (RNA Integrity

196   Number (RIN) $\geq$ 8). Samples should ideally consist of high molecular weight genetic material

197   (with little shearing), though continuing molecular advances enable genomic sequencing even of

198   low quantity or poor quality starting material. Extreme examples of the latter include

199   successfully sequenced whole genomes from ancient material (e.g., Rasmussen et al. 2010;

200   Meyer et al. 2012; Allentoft et al. 2015), including a more than 500,000-year-old horse (Orlando

201   et al. 2013).

202

203    Reduced-representation genome sequencing

204    *i. RADseq*

205    Reduced-representation sequencing methods evaluate only a small portion of the genome,

206    allowing researchers to sequence samples from a larger number of individuals within a given

207    budget in comparison to sequencing whole genomes. Restriction site-associated DNA

208    sequencing (RADseq) is currently the most widely used RRL sequencing method for NMOs

209    (Davey et al. 2011; Narum et al. 2013; Andrews et al. 2016). RADseq generates sequence data

210    from short regions adjacent to restriction cut sites and therefore targets markers that are

211    distributed relatively randomly across the genome and occur primarily in non-coding regions.

212    This method allows simultaneous discovery and genotyping of thousands of genetic markers for

213    virtually any species, regardless of availability of prior genomic resources. Of greatest interest

214    are variable markers, characterized either as single SNPs or phased alleles that can be resolved

215    from the identification of several variants within a single locus.

216

217    The large number of markers generated by RADseq dramatically increases genomic resolution

218    and statistical power for addressing many ecological and evolutionary questions when compared

219    to studies using traditional markers (Table S1). For example, heterozygosity fitness correlations

220    in harbor seals (*Phoca vitulina*) were nearly fivefold higher when using 14,585 RADseq SNPs

221    than when using 27 microsatellite loci (Hoffman et al. 2014). A recent study on the Atlantic

222    walrus (*Odobenus rosmarus rosmarus*) using 4,854 RADseq SNPs to model demographic

223    changes in connectivity and effective population size associated with the Last Glacial Maximum

224    (Shafer et al. 2015) both supported and extended inferences from previous studies using

225    traditional markers (Shafer et al. 2010; Shafer et al. 2014).

226

227    Furthermore, RADseq can provide sufficient numbers of markers across the genome to identify

228    genomic regions influenced by natural selection. These analyses require large numbers

229    (thousands to tens of thousands) of markers to ensure that some markers will be in linkage

230    disequilibrium with genomic regions under selection and to minimize false positives, particularly

231    under non-equilibrium demographic scenarios (Narum and Hess 2011; De Mita et al. 2013;

232    Lotterhos and Whitlock 2014). Extreme demographic shifts, as experienced by many marine

233    mammal populations (e.g., killer whales, Foote et al. 2016), can drive shifts in allele frequencies

234    that confound the distinction of drift and selection and make it difficult to detect genomic

235    signatures of selection (Poh et al. 2014). Proof of concept of the application of RADseq for

236    identifying genomic signatures of selection in wild populations was demonstrated in three-spined

237    sticklebacks (*Gasterosteus aculeatus*), for which analyses of over 45,000 SNPs (Hohenlohe et al.

238    2010) identified genomic regions of known evolutionary importance associated with differences

239    between marine and freshwater forms (Colosimo et al. 2005; Barrett et al. 2008). RADseq

240    studies with similar aims in marine mammals have resulted in comparatively sparser sampling of

241    SNPs (<10,000), likely due to both methodological differences and generally low genetic

242    diversity particularly among cetaceans. Nonetheless, genomic regions associated with resistance

243    to harmful algal blooms in common bottlenose dolphins (*Tursiops truncatus*) were identified

244    across multiple pairwise comparisons using 7,431 RADseq SNPs (Cammen et al. 2015), and

245    genomic regions associated with habitat use and resource specialization in killer whales (*Orcinus*

246    *orca*) were identified using 3,281 RADseq SNPs (Moura et al. 2014a). Some of these RADseq

247    SNPs associated with diet in killer whales were later also confirmed as occurring in genomic

248    regions of high differentiation and reduced diversity consistent with a signature of selection

249    identified in a study utilizing whole genome re-sequencing (Foote et al. 2016). It will remain

250    important for further studies of genomic signatures of selection in NMOs to carefully consider

251    which approach will generate a sufficiently large number of SNPs to accurately identify the

252    range of putatively neutral $F_{ST}$ values (and thus outliers) given the demographic history of the

253    population (Lotterhos and Whitlock 2014).

254

255    Numerous laboratory methods have been developed for generating RADseq data (reviewed in

256    Andrews et al. 2016), with the most popular library preparation methods currently being the

257    original RAD (Miller et al. 2007; Baird et al. 2008), Genotyping by Sequencing (GBS, Elshire et

258    al. 2011; Poland et al. 2012), and double digest RAD (ddRAD, Peterson et al. 2012). All

259    RADseq methods share the common goal of sequencing regions adjacent to restriction cut sites

260    across the genome, but differ in technical details, such as the number and type of restriction

261    enzymes used, the mechanisms for reducing genomic DNA fragment sizes, and the strategies for

262    attaching sequencing adapters to the target DNA fragments. For example, both the original RAD

263    method and GBS use a single enzyme digest, but the original RAD method uses a rare-cutting

264   enzyme and mechanical shearing to reduce DNA fragment size (Baird et al. 2008), whereas GBS

265   uses a more frequent-cutting enzyme and relies on preferential PCR amplification of shorter

266   fragments for indirect size selection (Elshire et al. 2011). These modifications lead to differences

267   across methods in the time and cost of library preparation, the number and lengths of loci

268   produced, and the types of error and bias present in the resulting data. Different RADseq

269   methods will be better suited to different research questions, study species, and research budgets,

270   and therefore researchers embarking on a RADseq study should carefully consider the suitability

271   of each method for their individual projects. Further details on the advantages and disadvantages

272   of each method are described in Andrews et al. (2016).

273

274   *ii. SNP arrays*

275   An alternative high-throughput reduced-representation genotyping approach involves the use of

276   custom arrays designed to capture and sequence targeted regions of the genome. Such array-

277   based approaches may provide certain advantages over RADseq, including the ability to easily

278   estimate genotyping error rates, scalability to thousands of samples, lower requirements for DNA

279   quantity/quality and technical effort, greater comparability of markers across studies, and the

280   ability to genotype SNPs within candidate genomic regions. However, unlike RADseq, array-

281   based techniques require prior knowledge of the study system's genome or the genome of a

282   closely related species, which remains unavailable for some NMOs. Furthermore, SNP arrays

283   must take into account the potential for ascertainment bias (e.g., Malenfant et al. 2015), whereas

284   RADseq avoids ascertainment bias by simultaneously discovering and genotyping markers.

285

286   To identify SNPs for NMO array development, researchers must rely on existing genomic

287   resources or generate new reference sequences, in the form of whole or reduced-representation

288   genomes or transcriptomes (Hoffman et al. 2012; Malenfant et al. 2015). When a whole genome

289   reference assembly is available for the target species or a related species, multiplex shotgun

290   sequencing can facilitate the rapid discovery of hundreds of thousands of SNPs for array

291   development. This SNP discovery approach involves high-throughput sequencing of sheared

292   genomic DNA that can be sequenced at a low depth of coverage (i.e., low mean read depth

293   across the genome) if suitable genotype likelihood-based methods (O'Rawe et al. 2015) are used

294   to identify polymorphic sites. Thus, this approach is less restrictive in terms of DNA quality. For

295    example, shotgun sequencing of 33 Northeast Atlantic common bottlenose dolphins, which

296    included degraded DNA collected from stranded specimens, on one Illumina HiSeq2000 lane of

297    100 bp single-end sequencing identified 440,718 high-quality SNPs (M. Louis unpublished data).

298    Such dense sampling of SNPs is essential for studies of population genomics that require a large

299    number of markers, such as for inferences of demographic history (Gutenkunst et al. 2009;

300    Excoffier et al. 2013; Liu and Fun 2015) and selective sweeps (Chen et al. 2010). Once a set of

301    putative markers has been identified, hybridization probes can be designed from their flanking

302    sequences and printed onto a SNP array. The two principal SNP genotyping platforms supporting

303    thousands to millions of SNPs are the Illumina Infinium iSelect® and Affymetrix Axiom®

304    arrays.

305

306    The use of SNP arrays in NMOs has thus far been somewhat limited, potentially due to low SNP

307    validation rates (Chancerel et al. 2011; Helyar et al. 2011), issues of ascertainment bias

308    (Albrechtsen et al. 2010; McTavish and Hillis 2015), and cost of SNP discovery. However, using

309    both SNP data and whole genome sequence from the Antarctic fur seal (*Arctocephalus gazella*),

310    Humble et al. (2016) recently demonstrated that careful filtering based on SNP genomic context

311    prior to array development has the potential to substantially increase assay success rates. Further,

312    ascertainment bias can be reduced by selecting samples for SNP discovery that span the

313    geographic range of populations that will be target sequenced (Morin et al. 2004). By accounting

314    for ascertainment bias, Malenfant et al. (2015) were able to demonstrate population structure in

315    Canadian polar bears (*Ursus maritimus*) more clearly using a 9K SNP array than 24

316    microsatellite markers.

317

318    *iii. Target sequence capture*

319    Target sequence capture (TSC, also called target enrichment, direct selection, or Hyb-seq) has

320    many of the same advantages and disadvantages as the array-based SNP approaches described

321    above, but differs in library preparation, sequencing platform, and resulting sequence data. While

322    SNP arrays genotype single variable positions, TSC can be used to sequence selected short

323    fragments. With TSC, researchers can amplify and sequence up to a million target probes on

324    solid-state arrays, and even more if in-solution arrays are used. This gives the user the ability to

325    choose to sequence many samples in parallel (Cummings et al. 2010), as many as 100-150 per

326    Illumina HiSeq lane, or to sequence many regions per individual. Recent advances in target

327    enrichment, such as genotyping in thousands (Campbell et al. 2015), anchored hybrid enrichment

328    (Lemmon et al. 2012), and target capture of ultraconserved elements (UCEs, Faircloth et al.

329    2012; McCormack et al. 2012), have further increased the number of regions and individuals that

330    can be sampled in a single lane. In addition, UCEs overcome the need for a reference genome,

331    enabling their wide application across many NMOs (though designing custom probe sets from

332    closely related species will remain preferable in many cases (Hancock-Hanser et al. 2013)).

333    Although a number of methodological variants have been developed and optimized (Bashiardes

334    et al. 2005; Noonan et al. 2006; Hodges et al. 2009; Cummings et al. 2010; Mamanova et al.

335    2010; Hancock-Hanser et al. 2013), TSC generally relies on hybridization and amplification of

336    specially prepared libraries consisting of fragmented genomic DNA. Many companies offer kits

337    for TSC, such as Agilent (SureSelect) and MYcroarray (MYbaits), with MYcroarray specifically

338    marketing their kits for use with NMOs.

339

340    The most common use of TSC has been the capture of whole exomes in model organisms,

341    including humans (Ng et al. 2009). However, as costs have plummeted, TSC is increasingly

342    being used in investigations of NMOs. TSC is particularly useful in sequencing ancient DNA,

343    where it can enrich the sample for endogenous DNA content relative to exogenous DNA (i.e.,

344    contamination) and thereby increase the relative DNA yield (Ávila-Arcos et al. 2011; Enk et al.

345    2014). For example, TSC has been used to generate mitogenome sequences from subfossil killer

346    whale specimens originating from the mid-Holocene for comparison with modern lineages

347    (Foote et al. 2013). TSC was also recently utilized to compare >30 kb of exonic sequence from

348    museum specimens of the extinct Steller's sea cow (*Hydrodamalis gigas*) and a modern dugong

349    (*Dugong dugon*) specimen to investigate evolution within Sirenia (Springer et al. 2015). Springer

350    et al. (2016) further used TSC to examine gene evolution related to dentition across edentulous

351    mammals, including mysticetes. Finally, TSC of both exonic and intronic regions has been used

352    to assess genetic divergence across cetacean species (Hancock-Hanser et al. 2013; Morin et al.

353    2015). These studies show the potential use of TSC across evolutionary timescales for population

354    genomics, phylogenomics, and studies of selection and gene loss across divergent lineages

355    (Table S1).

356

357 <u>Whole genome sequencing</u>

358 Beyond advances enabled by the reduced-representation methods presented above, our power

359 and resolution to elucidate evolutionary processes, including selection and demographic shifts,

360 can be further increased by sequencing whole genomes.

361

362 *i. Reference genome sequencing*

363 At the time of publication, there are 12 publicly available[1] whole (or near-whole) marine

364 mammal genomes of varying quality representing 10 families, including 7 cetaceans (Fig 1A), 3

365 pinnipeds (Fig 1B), the West Indian manatee (*Trichechus manatus*), and the polar bear. The first

366 sequenced marine mammal genome was that of the common bottlenose dolphin, which was

367 originally sequenced to ~2.5x depth of coverage using Sanger sequencing (Lindblad-Toh et al.

368 2011). This genome was later improved upon by adding both 454 and Illumina HiSeq data

369 (Foote et al. 2015). Other subsequent marine mammal genomes were produced solely using

370 Illumina sequencing and mate-paired or paired-end libraries with varied insert sizes (Miller et al.

371 2012; Zhou et al. 2013; Yim et al. 2014; Foote et al. 2015; Keane et al. 2015; Kishida et al. 2015;

372 Humble et al. 2016).

373

374 Whole genome sequencing has been used to address many issues in marine mammal genome

375 evolution, usually by comparison with other existing mammalian genomes. Biological insights

376 discussed in the genome papers listed above include the evolution of transposons and repeat

377 elements, gene evolution and positive selection, predicted population structure through time,

378 SNP validation, molecular clock rates, and convergent molecular evolution (Table S1). For

379 example, analyses of the Yangtze river dolphin (*Lipotes vexillifer*) genome confirmed that a

380 bottleneck occurred in this species during the last period of deglaciation (Zhou et al. 2013). In

381 addition, following upon earlier smaller-scale studies (e.g., Deméré et al. 2008; McGowen et al.

382 2008; Hayden et al. 2010), genomic analyses have confirmed the widespread decay of gene

383 families involved in olfaction, gustation, enamelogenesis, and hair growth in some cetaceans

384 (Yim et al. 2014; Kishida et al. 2015). Perhaps the most widespread use of whole genome studies

[1] These genomes are available on NCBI's online genome database or Dryad, but they have not all been published. As agreed upon in the Fort Lauderdale Convention, the community standard regarding such unpublished genomic resources is to respect the data generators' right to publish with these data first.

385    has been the use of models of selection to detect protein-coding genes that show evidence of

386    natural selection in specific lineages. A recent study by Foote et al. (2015) extended this

387    approach to investigate convergent positive selection among cetaceans, pinnipeds, and sirenians.

388    This study exemplifies a trend in recent genomic studies that sequence multiple genomes to

389    address a predetermined evolutionary question, in this case, the molecular signature of aquatic

390    adaptation.

391

392    In addition to these evolutionary insights that typically stem from a comparative genomics

393    approach, the development of high-quality reference genome assemblies provide an important

394    resource that facilitates mapping of reduced-representation genomic data (see previous section)

395    as well as short-read sequencing data with relatively low depth of coverage (see following

396    section). These data types can be generated at relatively low cost on larger sample sizes enabling

397    population-scale genomic studies. In many cases, genome assemblies from closely related

398    species are sufficient for use as a reference. Particularly among marine mammals, given their

399    generally slow rate of nucleotide divergence, it is therefore likely unnecessary to sequence a

400    high-quality reference genome assembly for every species. Instead, resources could be allocated

401    toward population-scale studies, including genome re-sequencing efforts.

402

403    *ii. Population-level genome re-sequencing*

404    In contrast to reference genome sequencing that today often exceeds 100x mean read depth and

405    typically combines long- and short-insert libraries to generate high-quality assemblies for one to

406    a few individuals, genome re-sequencing studies aim to achieve only ≥2x mean read depth on

407    tens to hundreds of individuals from short-insert libraries whose reads are anchored to existing

408    reference assemblies. Despite the inherent trade-offs between cost, read depth, coverage, and

409    sample size, genome re-sequencing of large numbers of individuals for population-level

410    inference can be conducted at a relatively low cost. In the past five years, several influential

411    studies have used genome re-sequencing to advance our understanding of the genomic

412    underpinnings of different biological questions in model systems. For example, population

413    genomics of *Heliconius* butterflies highlighted the exchange of genes between species that

414    exhibit convergent wing patterns (The *Heliconius* Genome Consortium 2012); whole genome re-

415    sequencing of three-spined sticklebacks highlighted the re-use of alleles in replicated

416  divergences associated with ecological speciation and local adaptation (Jones et al. 2012); and

417  combined population genomics and phylogenomics have identified regions of the genome

418  associated with variation in beak shape and size in Darwin's finches (Lamichhaney et al. 2015).

419

420  To date only two marine mammal population genomics studies using whole genome re-

421  sequencing have been published. These studies involved re-sequencing the genomes of 79

422  individuals from three populations of polar bears (Liu et al. 2014a) and 48 individuals from five

423  evolutionarily divergent ecotypes of killer whale (Foote et al. 2016). The findings of Foote et al.

424  (2016) confirmed results of population differentiation that had previously been established using

425  traditional genetic markers (Morin et al. 2010a). However, the study also provided new insights

426  into the demographic history, patterns of selection associated with ecological niche, and evidence

427  of episodic ancestral admixture that could not have been obtained using traditional markers.

428

429  Several new resources have made such population genomic studies economically possible for a

430  greater number of NMOs, including the availability of reference genome assemblies (see section

431  above), relatively low-cost high-throughput sequencing (further increases in throughput expected

432  with the new Illumina HiSeq X Ten (van Dijk et al. 2014)), and crucially, the development of

433  likelihood-based methods that allow estimation of population genetic metrics from re-sequencing

434  data (Fumagalli et al. 2013; O'Rawe et al. 2015). One last consideration is the ease of laboratory

435  methods necessary to generate whole genome re-sequencing data when compared to other

436  methods such as RADseq or TSC. DNA simply needs to be extracted from the samples and,

437  using proprietary kits, built into individually index-amplified libraries that are equimolarly

438  pooled and submitted for sequencing.

439

440  Many population genomic analyses are based on the coalescent model that gains most

441  information from the number of independent genetic markers, not the number of individuals

442  sampled. Sample sizes of ~10 individuals are usually considered sufficient (Robinson et al.

443  2014) and have been standard in many genome-wide studies in the eco-evolutionary sciences

444  (Ellegren et al. 2012; Jones et al. 2012). Thus, sampling fewer individuals by whole genome re-

445  sequencing is a salient approach that allows us to consider many more gene trees, whilst

446  continuing to provide robust estimates of per-site genetic metrics (e.g., $F_{ST}$). The robustness of

447    inference from data with low mean read depth across the genome was recently confirmed using a

448    comparison of per-site $F_{ST}$ estimates for the same sites from high-depth (≥20x) RADseq data and

449    low-depth (≈2x) whole genome re-sequencing data in pairwise comparisons between the same

450    two killer whale ecotypes (Foote et al. 2016).

451

452    Beyond the increased power afforded by sequencing more polymorphic sites, whole genome re-

453    sequencing also allows inference of demographic history from the genome of even just a single

454    individual by identifying Identical By Descent (IBD) segments and runs of homozygosity (Li

455    and Durbin 2011; Harris and Nielsen 2013). For example, Liu et al. (2014a) found evidence for

456    ongoing gene flow from polar bears into brown bears after the two species initially diverged.

457    Genome re-sequencing of sufficient numbers of individuals also facilitates haplotype phasing,

458    which has many applications, including the detection of ongoing selective sweeps (Ferrer-

459    Admetlla et al. 2014) and the inference of demographic history of multiple populations based on

460    coalescence of pairs of haplotypes in different individuals (Schiffels and Durbin 2014).

461    However, haplotype phasing typically requires genomic data with higher mean read depth (~20x)

462    from tens of individuals (though recent advances in genotype imputation suggest success with

463    data of lower mean read depth (VanRaden et al. 2015)). Thus far, phasing has been restricted to

464    relatively few NMO studies, and no marine mammal studies to the best of our knowledge.

465

466    Transcriptome sequencing

467    In comparison with the DNA-based genomic approaches described above, RNA-based genomic

468    approaches are a relatively new and emerging application in NMOs such as marine mammals.

469    Transcriptomics by RNA sequencing (RNAseq) can rapidly generate vast amounts of

470    information regarding genes and gene expression without any prior genomic resources. This

471    approach can resolve differences in global gene expression patterns between populations,

472    individuals, tissues, cells, and physiological or environmental conditions, and can yield insights

473    into the molecular basis of environmental adaptation and speciation in wild animals (Wolf 2013;

474    Alvarez et al. 2015). Furthermore, RNAseq is a valuable tool for resource development, for

475    example as a precursor to designing SNP and TSC arrays (e.g., Hoffman et al. 2012). However,

476    applying RNAseq to NMOs requires several unique considerations in comparison to the DNA-

477    based methods described above. Most importantly, the labile nature of gene transcription and

478   high detection sensitivity of RNAseq have the potential to amplify transcriptional "noise" and

479   are thus extremely sensitive to experimental design.

480

481   If the experimental goal is to capture a comprehensive transcriptome profile for a study

482   organism, multiple tissues from individuals of varied life history stages should be sampled.

483   However, if the aim is to characterize transcriptional responses to physiological or environmental

484   stimuli, efforts should focus on minimizing variability in individuals and sampling conditions

485   (Wolf 2013). For differential expression analyses, pairwise comparisons should be made within

486   the same individual if at all possible (e.g., before and after treatment, between two

487   developmental stages). As RNAseq only captures a 'snapshot' of gene expression in time,

488   repeated sampling or time-course studies are necessary to obtain a more complete picture of

489   cellular responses to the condition(s) in question (Spies and Ciaudo 2015). Sampling and

490   sequencing depth requirements will depend on the study design. Simulation studies have shown

491   that a minimum of 5-6 biological replicates sequenced at a depth of 10-20 million reads per

492   sample is necessary for differential expression analysis (Liu et al. 2014b; Schurch et al. 2015).

493   RNAseq can also be used for biomarker development to expand molecular toolkits for NMOs

494   without sequenced genomes (Hoffman et al. 2013). In this case, higher sequencing depths of 30-

495   60 million reads per sample are recommended for SNP discovery and genotyping (De Wit et al.

496   2015).

497

498   Following sequence generation, transcript annotation remains a challenge for NMOs without

499   reference transcriptomes or genomes. *De novo* transcriptomes can be annotated through detection

500   of assembled orthologs of highly conserved proteins, but these analyses remain limited by the

501   quality of reference databases. As a result, NMO transcriptomes are biased in favor of highly

502   conserved terrestrial mammal genes and therefore provide an incomplete understanding of

503   animal adaptations to natural environments (Evans 2015). For example, while 70.0% of northern

504   elephant seal (*Mirounga angustirostris*) skeletal muscle transcripts had BLASTx hits to mouse

505   genes, only 54.1% of blubber transcripts could be annotated due to poor representation of this

506   tissue in terrestrial mammal reference proteomes (Khudyakov et al. 2015b).

507

508  To date, RNAseq has been used for gene discovery and phylogenomics analyses in Antarctic fur

509  seal (Hoffman 2011; Hoffman et al. 2013), polar bear (Miller et al. 2012), Indo-Pacific

510  humpback dolphin (*Sousa chinensis* (Gui et al. 2013)), spotted seal (*Phoca largha* (Gao et al.

511  2013)), bowhead whale (*Balaena mysticetus* (Seim et al. 2014)), narrow-ridged finless porpoise

512  (*Neophocaena asiaeorientalis* (Ruan et al. 2015)), and humpback whale (*Megaptera*

513  *novaeangliae* (Tsagkogeorga et al. 2015)) (Table S1). Due to the challenges of repeated

514  sampling of wild marine mammals, few studies have examined cetacean or pinniped

515  transcriptome responses to environmental or experimental stimuli. The majority of such

516  functional gene expression studies have used microarrays (Mancia et al. 2008; Mancia et al.

517  2012; Mancia et al. 2015); however, RNAseq has been employed to profile sperm whale

518  (*Physeter macrocephalus*) skin cell response to hexavalent chromium (Pabuwal et al. 2013) and

519  free-ranging northern elephant seal skeletal muscle response to an acute stress challenge

520  (Khudyakov et al. 2015a; Khudyakov et al. 2015b). With decreasing sequencing costs and

521  improvements in bioinformatics tools, RNAseq has the potential to accelerate molecular

522  discoveries in marine mammal study systems and supplement existing functional genomics

523  approaches.

524

525  <u>Emerging techniques</u>

526  In addition to the relatively proven NMO genomic data generation techniques described above, a

527  suite of emerging techniques is entering the field, with exciting promise for exploration of

528  existing and new research areas. For example, high-throughput shotgun sequencing is

529  increasingly being used to identify genetic material from multiple species in a single sample

530  (metagenomics and metatranscriptomics), rather than focus on characterizing variation in a

531  single target individual. These multi-species approaches can be used, for example, to

532  characterize diet from fecal samples (Deagle et al. 2009) and to investigate microbiomes (Nelson

533  et al. 2015), objectives with implications for improving our understanding of both basic ecology

534  and health in natural populations of NMOs. Furthermore, high-throughput sequencing of

535  environmental DNA dramatically increases the throughput of NMO detection in environmental

536  (e.g., seawater) samples (Thomsen et al. 2012), using degenerate primers for multi-species

537  detection rather than requiring the design and implementation of numerous single-species

538  protocols (Foote et al. 2012).

539

540 A second broad area of emerging interest moves beyond the study of variation at the DNA and

541 RNA levels to examine epigenetic effects of histone modification on gene regulation and

542 evolution. Epigenomic studies often examine changes in DNA methylation in association with

543 processes such as cancer and ageing. Such approaches, from targeted gene to genome-wide, have

544 only very recently and not yet frequently been applied in NMOs. Polanowski et al. (2014) used a

545 targeted gene approach to examine changes in DNA methylation in age-associated genes,

546 previously identified in humans and mice, in humpback whales of known age. The most

547 informative markers were able to estimate humpback whale ages with standard deviations of

548 approximately 3-5 years, demonstrating the potential transferability of these approaches from

549 model to non-model organism. Villar et al. (2015) utilized a genome-wide approach – chromatin

550 immunoprecipitation followed by high-throughput sequencing (ChIPseq) – to examine gene

551 regulatory element evolution across mammals, including four species of cetaceans. This study

552 identified highly conserved gene regulatory elements based on their histone modifications

553 (H3K27ac and H3K4me3), showed that recently evolved enhancers were associated with genes

554 under positive selection in marine mammals, and identified unique *Delphinus*-specific enhancers.

555 Finally, reduced-representation epigenomic approaches have also been developed (Gu et al.

556 2011), and although they have not yet been used in marine mammals to our knowledge, these

557 techniques could facilitate future studies of how changes in DNA methylation patterns affect

558 other biological processes, such as stress levels or pregnancy.

559

560 **Data analysis**

561 Following the generation of genomic data, researchers must select the most appropriate genomic

562 analysis (i.e., bioinformatics) pipelines, which often differ significantly from those used in

563 traditional genetic studies of NMOs. The choice of analysis pipeline will depend on multiple

564 factors including the availability of a reference genome, the level of diversity within the dataset

565 (e.g., single- or multi-species), the type of data generated (e.g., single- or paired-end), and the

566 computing resources available. The computational needs, both in terms of hardware and

567 competency in computer science, for analysis of genomic data typically far exceed those

568 necessary for traditional genetic markers. On the smaller end of the spectrum, one lane of 50 bp

569 single-end sequencing on an Illumina HiSeq 2500 can produce tens of gigabytes of data, while

570  data files associated with a single high-quality vertebrate genome may reach hundreds of

571  gigabytes in size (Ekblom and Wolf 2014). Computing resources necessary for the analysis of

572  these genomic datasets can range from ~10 gigabytes for a pilot study using a reduced-

573  representation sequencing approach to over a terabyte for whole genome sequence assembly

574  (Ekblom and Wolf 2014). Fortunately, university computing clusters, cloud-based (Stein 2010)

575  and high-performance computing clusters (e.g., XSEDE; Towns et al. 2014), and open web-

576  based platforms for genomic research (e.g., Galaxy; Goecks et al. 2010) are becoming

577  increasingly accessible. Furthermore, new pipelines are continuously being developed and

578  improved, and there are a growing number of resources aimed at training molecular ecologists

579  and evolutionary biologists in computational large-scale data analysis (Andrews and Luikart

580  2014; Belcaid and Toonen 2015; Benestan et al. 2016). We provide an indicative list of the

581  current, most commonly used analysis pipelines that are specific to each data generation method

582  in Table 1. Here we briefly summarize current genomic data analysis pipelines and discuss

583  considerations that are likely to be similar across multiple data generation methods.

584

585  Genomic data analysis often involves multiple steps, and the choice of analysis tool for each step

586  can greatly affect the outcome, with different tools producing different (though usually

587  overlapping) sets of results (e.g., Schurch et al. 2015). All analyses begin by evaluating data

588  quality, trimming sequences if necessary to remove erroneous nucleotides (MacManes 2014),

589  and implementing appropriate data quality filters (e.g., phred scores, read length, and/or read

590  depth). Raw reads also need to be demultiplexed based on unique barcodes if pools of

591  individuals were sequenced in a single lane. Analyses then usually proceed in a *de novo* or

592  genome-enabled manner, depending on available resources. Briefly, sequences can be compared

593  (e.g., to identify variants) by mapping all reads to a reference genome or *de novo* assembling

594  stacks of sequences putatively derived from the same locus based on sequence similarity. *De*

595  *novo* methods are sensitive to sequencing error, as well as true genetic variation, and therefore

596  can erroneously assemble polymorphic sequences as separate loci or transcripts, requiring further

597  filtering to remove redundancy. The opposite problem can also occur in both *de novo* and

598  reference mapping approaches, where two distinct loci (e.g., paralogous loci) may assemble as a

599  single locus or map to the same reference location. Researchers should therefore recognize the

600    inherent trade-offs when carefully selecting their thresholds for acceptable levels of variation

601    within and among loci.

602

603    Considerations relevant to the selection of subsequent downstream analyses are specific to the

604    type of data generated and the research objective. For example, RADseq analysis pipelines differ

605    in the algorithms used to genotype variants (Table 1). Similarly, there are several gene

606    expression analysis pipelines for RNAseq data that compare transcript abundance between

607    samples (Table 1). Analysis of TSC data usually uses standard *de novo* assemblers (e.g., Trinity,

608    Velvet); these assemblers can be run using packages such as PHYLUCE (Faircloth 2015), which

609    is designed specifically for use with ultraconserved elements. Unfortunately, for most analyses,

610    there are no unifying recommendations currently available and researchers must evaluate several

611    approaches, each with their own advantages and disadvantages, in order to select the most

612    appropriate tool for their particular experiment and system. Furthermore, we can expect that the

613    recommendations for analysis tools will continue to evolve as new programs become available in

614    the future.

615

616    <u>Guidelines for data quality control and sharing</u>

617    With rapid growth in sequencing platforms and bioinformatics analysis pipelines comes the need

618    to extend existing principles (e.g., Bonin et al. 2004) on quality control, analysis, and

619    transparency. General recommendations for sample and data handling, library preparation, and

620    sequencing have been discussed elsewhere (Paszkiewicz et al. 2014). We therefore focus on the

621    need to produce guidelines on data quality evaluation and reporting for genomic data (e.g.,

622    Morin et al. 2010b). A primary challenge in this area is that quality metrics vary widely across

623    sequencing technologies. Yet, regardless of sequencing platform, the quality of sequencing reads

624    must be evaluated (e.g., using FastQC; Andrews 2010) and reported.

625

626    Best practices guidelines for reference genome sequencing and RNAseq data generation,

627    analysis, and reporting are available from the human-centric ENCODE consortium

628    (www.encodeproject.org). These include minimum depth of sequencing and number and

629    reproducibility of biological replicates. For RNAseq experiments, evaluation of *de novo*

630    assembly quality remains a challenge. Suggested quality metrics include percentage of raw reads

631   mapping back to the assembly and number of assembled transcripts with homology to known

632   proteins (MacManes 2016). Emerging tools such as Transrate (Smith-Unna et al. 2015) attempt

633   to integrate these and other metrics into a comprehensive assembly quality score.

634

635   In contrast, there is not yet any standard way to estimate or report error rates with RADseq or

636   genome re-sequencing methods (but see Mastretta-Yanes et al. 2015; Fountain et al. 2016).

637   Recommendations to improve confidence in genotyping include using methods that account for

638   population-level allele frequencies when calling individual genotypes, mapping reads to

639   reference genomes rather than *de novo* assembly (Nadeau et al. 2014; Fountain et al. 2016),

640   filtering out PCR duplicates (Andrews et al. 2014), identifying and removing markers in possible

641   repeat regions, and filtering data to include only those with high read depth (>10-20x per locus

642   per individual) (Nielsen et al. 2011). Other analysis methods, such as robust Bayesian methods

643   and likelihood-based approaches that account for read quality in calculations of posterior

644   probabilities of genotypes and per-site allele frequencies utilizing the sample mean site

645   frequency spectrum as a prior (Fumagalli et al. 2013), can account for uncertainty and/or error in

646   the data, and are therefore suitable for use with low to moderate read depths (2-20x per locus;

647   e.g., Han et al. 2015; O'Rawe et al. 2015).

648

649   Due to the large number of analysis tools that are available, data quality and reproducibility

650   ultimately depend on methods and data transparency. All raw sequencing reads should be

651   publicly archived, for example deposited in the NCBI Sequence Read Archive. Many journals,

652   including the *Journal of Heredity* (Baker 2013), now also require that primary data supporting

653   the published results and conclusions (e.g., SNP genotypes, assemblies) be publicly archived in

654   online data repositories (e.g., Dryad). We further recommend making public the analysis

655   pipelines, scripts (e.g., using GitHub), and additional outputs, as appropriate, in order for

656   analyses to be fully reproducible and transparent, which is the cornerstone of the scientific

657   method (Nosek et al. 2015).

658

**Future directions**

660   As demonstrated here for one group of mammalian taxa, the rapid growth of the field of non-

661   model genomics has been both impressive and empowering. As we approach a point of relative

662 saturation in reference genomes, we anticipate an increase in population-scale genomic studies

663 that produce lower depth or coverage datasets per individual but across larger sample sizes. In

664 addition (or alternatively), we hope to see increasing efforts to sequence reference transcriptomes

665 and improve NMO genome annotation in ways beyond the inherently limited approach of

666 comparison to gene lists from a few model organisms. Population-scale genomic studies will

667 facilitate greater ecological understanding of natural populations, while efforts to improve

668 annotation will address persistent limitations in our understanding of gene function for NMOs.

669 Ultimately, improving our understanding of local adaptation, adaptive potential, and

670 demographic history through the use of genomic toolkits such as those described here is likely to

671 have important implications for the future conservation of these populations.

672

673 Advances in sequencing technologies and analytical tools will no doubt continue, in some cases

674 drawing on established techniques in model organisms, posing both new opportunities and new

675 challenges for researchers in NMO genomics. Likely the most persistent challenge will remain

676 selecting the data generation and experimental design that is most appropriate for the respective

677 research objective. Our review identified few cases that exhibit relative dominance of a single

678 methodology and analytical pipeline (e.g., RADseq and STACKS, RNAseq and Trinity); rather,

679 more often we found a diversity of approaches even within each category of data generation. In

680 fact, such diversity of approaches has its benefits, with each approach promoting its own

681 advantages (and limitations). Overall, our reflections on lessons learned from the past decade of

682 NMO genomics in one well-studied group of mammalian taxa clearly demonstrate the value,

683 increased ease, and future promise of applying genomic techniques across a wide range of non-

684 model species to gain previously unavailable insights into evolution, population biology, and

685 physiology on a genome-wide scale.

686

687 **Acknowledgements**

709 **References**

710 Albrechtsen A, Nielsen FC, Nielsen R. 2010. Ascertainment biases in SNP chips affect measures
711     of population divergence. *Mol Biol Evol*. 27:2534-2547.
712 Alexander A, Steel D, Hoekzema K, Mesnick S, Engelhaupt D, Kerr I, Payne R, Baker CS.
713     2016. What influences the worldwide genetic structure of sperm whales (*Physeter*
714     *macrocephalus*)? *Mol Ecol*.
715 Allentoft ME, Sikora M, Sjögren K-G, Rasmussen S, Rasmussen M, Stenderup J, Damgaard PB,
716     Schroeder H, Ahlstrom T, Vinner L*, et al.* 2015. Population genomics of Bronze Age
717     Eurasia. *Nature*. 522:167-172.
718 Alvarez M, Schrey AW, Richards CL. 2015. Ten years of transcriptomics in wild populations:
719     what have we learned about their ecology and evolution? *Mol Ecol*. 24:710-725.
720 Anders S, Huber W. 2010. Differential expression analysis for sequence count data. *Genome*
721     *Biol*. 11:R106.
722 Andrews K, Good JM, Miller MR, Luikart G, Hohenlohe PA. 2016. Harnessing the power of
723     RADseq for ecological and evolutionary genomics. *Nat Rev Genet*. 17:81-92.
724 Andrews KR, Hohenlohe PA, Miller MR, Hand BK, Seeb JE, Luikart G. 2014. Trade-offs and
725     utility of alternative RADseq methods: Reply to Puritz *et al*. 2014. *Mol Ecol*. 23:5943-
726     5946.
727 Andrews KR, Luikart G. 2014. Recent novel approaches for population genomics data analysis.
728     *Mol Ecol*. 23:1661-1667.
729 Andrews S. 2010. FastQC: a quality control tool for high throughput sequence data. Available
730     online at: http://www.bioinformatics.babraham.ac.uk/projects/fastqc

731    Ankeny RA, Leonelli S. 2011. What's so special about model organisms? *Studies in History and*
732         *Philosophy of Science*. 42:313-323.
733    Armengaud J, Trapp J, Pible O, Geffard O, Chaumot A, Hartmann EM. 2014. Non-model
734         organisms, a species endangered by proteogenomics. *J Proteomics*. 105:5-18.
735    Arnason U, Adegoke JA, Bodin K, Born EW, Esa YB, Gullberg A, Nilsson M, Short RV, Xu X,
736         Janke A. 2002. Mammalian mitogenomic relationships and the root of the eutherian tree.
737         *Proc Natl Acad Sci USA*. 99:8151-8156.
738    Arnason U, Gullberg A, Widegren B. 1991. The complete nucleotide sequence of the
739         mitochondrial DNA of the fin whale, *Balaenoptera physalus*. *J Mol Evol*. 33:556-568.
740    Ávila-Arcos M, Cappellini E, Romero-Navarro JA, Wales N, Moreno-Mayar JV, Rasmussen M,
741         Fordyce SL, Montiel R, Vielle-Calzada J-P, Willerslev E, *et al.* 2011. Application and
742         comparison of large-scale solution-based DNA capture-enrichment methods on ancient
743         DNA. *Sci Rep*. 1:74.
744    Baird NA, Etter PD, Atwood TS, Currey MC, Shiver AL, Lewis ZA, Selker EU, Cresko WA,
745         Johnson EA. 2008. Rapid SNP discovery and genetic mapping using sequenced RAD
746         markers. *PLoS One*. 3:e3376.
747    Baker CS. 2013. *Journal of Heredity* adopts Joint Data Archiving Policy. *J Hered*. 104:1.
748    Barrett RDH, Rogers SM, Schluter D. 2008. Natural selection on a major armor gene in
749         threespine stickleback. *Science*. 322:255-257.
750    Bashiardes S, Veile R, Helms C, Mardis ER, Bowcock AM, Lovett M. 2005. Direct genomic
751         selection. *Nat Methods*. 2:63-69.
752    Belcaid M, Toonen RJ. 2015. Demystifying computer science for molecular ecologists. *Mol*
753         *Ecol*. 24:2619-2640.
754    Benestan LM, Ferchaud A-L, Hohenlohe PA, Garner BA, Naylor GJP, Baums IB, Schwartz MK,
755         Kelley JL, Luikart G. 2016. Conservation genomics of natural and managed populations:
756         building a conceptual and practical framework. *Mol Ecol*.
757    Bolger AM, Lohse M, Usadel B. 2014. Trimmomatic: a flexible trimmer for Illumina Sequence
758         Data. *Bioinformatics*. 30:2114-2120.
759    Bonin A, Bellemain E, Bronken Eidesen P, Pompanon F, Brochmann C, Taberlet P. 2004. How
760         to track and assess genotyping errors in population genetics studies. *Mol Ecol*. 13:3261-
761         3273.
762    Brown CT, Howe A, Zhang Q, Pyrkosz AB, Brom TH. 2012. A reference-free algorithm for
763         computational normalization of shotgun sequencing data. *arXive*. 1203:4802.
764    Cammen KM, Schultz TF, Rosel PE, Wells RS, Read AJ. 2015. Genomewide investigation of
765         adaptation to harmful algal blooms in common bottlenose dolphins (*Tursiops truncatus*).
766         *Mol Ecol*. 24:4697-4710.
767    Campbell NR, Harmon SA, Narum SR. 2015. Genotyping-in-Thousands by sequencing (GT-
768         seq): a cost effective SNP genotyping method based on custom amplicon sequencing.
769         *Mol Ecol Resour*. 15:855-867.
770    Carroll EL, Baker CS, Watson M, Alderman R, Bannister J, Gaggiotti OE, Gröcke DR,
771         Patenaude N, Harcourt R. 2015. Cultural traditions across a migratory network shape the
772         genetic structure of southern right whales around Australia and New Zealand. *Sci Rep*.
773         5:16182.
774    Catchen JM, Amores A, Hohenlohe PA, Cresko WA, Postlethwait JH. 2011. *Stacks*: building
775         and genotyping loci *de novo* from short-read sequences. *G3*. 1:171-182.

776    Catchen JM, Hohenlohe PA, Bassham S, Amores A, Cresko WA. 2013. Stacks: an analysis tool
777        set for population genomics. *Mol Ecol*. 22:3124-2140.
778    Chancerel E, Lepoittevin C, Le Provost G, Lin Y-C, Jaramillo-Correa JP, Eckert AJ, Wegrzyn
779        JL, Zelenika D, Boland A, Frigerio J-M*, et al.* 2011. Development and implementation of
780        a highly-multiplexed SNP array for genetic mapping in maritime pine and comparative
781        mapping with loblolly pine. *BMC Genomics*. 12:368.
782    Chen H, Patterson N, Reich D. 2010. Population differentiation as a test for selective sweeps.
783        *Genome Res*. 20:393-402.
784    Colosimo PF, Hosemann KE, Balabhadra S, Villarreal G, Grimwood J, Schmutz J, Myers RM,
785        Schluter D, Kingsley DM. 2005. Widespread parallel evolution in sticklebacks by
786        repeated fixation of ectodysplasin alleles. *Science*. 307:1928-1933.
787    Conesa A, Gotz S, Garcia-Gomez JM, Terol J, Talon M, Robles M. 2005. Blast2GO: a universal
788        tool for annotation, visualization and analysis in functional genomics research.
789        *Bioinformatics*. 21:3674-3676.
790    Corander J, Majander KK, Cheng L, Merilä J. 2013. High degree of cryptic population
791        differentiation in the Baltic Sea herring *Clupea harengus*. *Mol Ecol*. 22:2931-2940.
792    Cummings N, King R, Rickers A, Kaspi A, Lunke S, Haviv I, Jowett JBM. 2010. Combining
793        target enrichment with barcode multiplexing for high throughput SNP discovery. *BMC
794        Genomics*. 11:641.
795    Davey JW, Hohenlohe PA, Etter PD, Boone JQ, Catchen JM, Blaxter ML. 2011. Genome-wide
796        genetic marker discovery and genotyping using next-generation sequencing. *Nat Rev
797        Genet*. 12:499-510.
798    De Mita S, Thuillet A-C, Gay L, Ahmadi N, Manel S, Ronfort J, Vigouroux Y. 2013. Detecting
799        selection along environmental gradients: analysis of eight methods and their effectiveness
800        for outbreeding and selfing populations. *Mol Ecol*. 22:1383-1399.
801    De Wit P, Pespeni MH, Palumbi SR. 2015. SNP genotyping and population genomics from
802        expressed sequences - current advances and future possibilities. *Mol Ecol*. 24:2310-2323.
803    Deagle BE, Kirkwood R, Jarman SN. 2009. Analysis of Australian fur seal diet by
804        pyrosequencing prey DNA in faeces. *Mol Ecol*. 18:2022-2038.
805    Deméré TA, McGowen MR, Berta A, Gatesy J. 2008. Morphological and molecular evidence for
806        a stepwise evolutionary transition from teeth to baleen in mysticete whales. *Syst Biol*.
807        57:15-37.
808    DePristo MA, Banks E, Poplin R, Garimella KV, Maguire JR, Hartl C, Philippakis AA, del
809        Angel G, Rivas MA, Hanna M*, et al.* 2011. A framework for variation discovery and
810        genotyping using next-generation DNA sequencing data. *Nat Genet*. 43:491-498.
811    Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, Batut P, Chaisson M, Gingeras
812        TR. 2013. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics*. 29:15-21.
813    Eaton DAR. 2014. PyRAD: assembly of *de novo* RADseq loci for phylogenetic analysis.
814        *Bioinformatics*. 30:1844-1849.
815    Ekblom R, Galindo J. 2011. Applications of next generation sequencing in molecular ecology of
816        non-model organisms. *Heredity*. 107:1-15.
817    Ekblom R, Wolf JBW. 2014. A field guide to whole-genome sequencing, assembly and
818        annotation. *Evolutionary Applications*. 7:1026-1042.
819    Ellegren H. 2014. Genome sequencing and population genomics in non-model organisms.
820        *Trends Ecol Evol*. 29:51-63.

821 Ellegren H, Smeds L, Burri R, Olason PI, Backström N, Kawakami T, Künstner A, Mäkinen H,
822     Nadachowska-Brzyska K, Qvarnström A, *et al.* 2012. The genomic landscape of species
823     divergence in *Ficedula* flycatchers. *Nature*. 491:756-760.
824 Elshire RJ, Glaubitz JC, Sun Q, Poland JA, Kawamoto K, Buckler ES, Mitchell SE. 2011. A
825     robust, simple genotyping-by-sequencing (GBS) approach for high diversity species.
826     *PLoS One*. 6:e19379.
827 Enk J, Devault A, Kuch M, Murgha Y, Rouillard J-M, Poinar H. 2014. Ancient whole genome
828     enrichment using baits built from modern DNA. *Mol Biol Evol*. 31:1292-1294.
829 Evans TG. 2015. Considerations for the use of transcriptomics in identifying the 'genes that
830     matter' for environmental adaptation. *J Exp Biol*. 218:1925-1935.
831 Excoffier L, Dupanloup I, Huerta-Sánchez E, Sousa VC, Foll M. 2013. Robust demographic
832     inference from genomic and SNP data. *PLoS Genetics*. 9:e1003905.
833 Faircloth BC. 2015. PHYLUCE is a software package for the analysis of conserved genomic
834     loci. *Bioinformatics*. 32:786-788.
835 Faircloth BC, McCormack JE, Crawford NG, Harvey MG, Brumfield RT, Glenn TC. 2012.
836     Ultraconserved elements anchor thousands of genetic markers spanning multiple
837     evolutionary timescales. *Syst Biol*. 61:717-726.
838 Ferrer-Admetlla A, Liang M, Korneliussen T, Nielsen R. 2014. On detecting incomplete soft or
839     hard selective sweeps using haplotype structure. *Mol Biol Evol*. 31:1275-1291.
840 Flicek P, Birney E. 2009. Sense from sequence reads: methods for alignment and assembly. *Nat
841     Methods*. 6:S6-S12.
842 Foote AD, Liu Y, Thomas GWC, Vinař Ts, Alföldi J, Deng J, Dugan S, van Elk CE, Hunter ME,
843     Joshi V, *et al.* 2015. Convergent evolution of the genomes of marine mammals. *Nat
844     Genet*. 47:272-275.
845 Foote AD, Newton J, Ávila-Arcos MC, Kampmann M-L, Samaniego JA, Post K, Rosing-Asvid
846     A, Sinding M-HS, Gilbert MTP. 2013. Tracking niche variation over millennial
847     timescales in sympatric killer whale lineages. *Proc R Soc Lond B Biol Sci*. 280:20131481.
848 Foote AD, Thomsen PF, Sveegaard S, Wahlberg M, Kielgast J, Kyhn LA, Salling AB, Galatius
849     A, Orlando L, Gilbert MTP. 2012. Investigating the potential use of environmental DNA
850     (eDNA) for genetic monitoring of marine mammals. *PLoS One*. 7:e41781.
851 Foote AD, Vijay N, Ávila-Arcos M, Baird RW, Durban JW, Fumagalli M, Gibbs RA, Hanson
852     MB, Korneliussen TS, Martin MD, *et al.* 2016. Genome-culture coevolution promotes
853     rapid divergence of killer whale ecotypes. *Nat Commun*. 7:11693.
854 Fountain ED, Pauli JN, Reid BN, Palsbøll PJ, Peery MZ. 2016. Finding the right coverage: the
855     impact of coverage and sequence quality on single nucleotide polymorphism genotyping
856     error rates. *Mol Ecol Resour*.
857 Fumagalli M, Vieira FG, Korneliussen TS, Linderoth T, Huerta-Sánchez E, Albrechtsen A,
858     Nielsen R. 2013. Quantifying population genetic differentiation from next-generation
859     sequencing data. *Genetics*. 195:979-992.
860 Fumagalli M, Vieira FG, Linderoth T, Nielsen R. 2014. *ngsTools*: methods for population
861     genetics analyses from Next-Generation Sequencing data. *Bioinformatics*. 30:1486-1487.
862 Gao X, Han J, Lu Z, Li Y, He C. 2013. *De novo* assembly and characterization of spotted seal
863     *Phoca largha* transcriptome using Illumina paired-end sequencing. *Comp Biochem
864     Physiol D Genom Proteom*. 8:103-110.

865 Garner BA, Hand BK, Amish SJ, Bernatchez L, Foster JT, Miller KM, Morin PA, Narum SR,
866     O'Brien SJ, Roffler G*, et al.* 2016. Genomics in conservation: case studies and bridging
867     the gap between data and application. *Trends Ecol Evol*. 31:81-83.
868 Glaubitz JC, Casstevens TM, Lu F, Harriman J, Elshire RJ, Sun Q, Buckler ES. 2014. TASSEL-
869     GBS: a high capacity genotyping by sequencing analysis pipeline. *PLoS One*. 9:e90346.
870 Gnerre S, MacCallum I, Przybylski D, Ribeiro FJ, Burton JN, Walker BJ, Sharpe T, Hall G, Shea
871     TP, Sykes S*, et al.* 2011. High-quality draft assemblies of mammalian genomes from
872     massively parallel sequence data. *Proc Natl Acad Sci USA*. 108:1513-1518.
873 Goecks J, Nekrutenko A, Taylor J, The Galaxy Team. 2010. Galaxy: a comprehensive approach
874     for supporting accessible, reproducible, and transparent computational research in the life
875     sciences. *Genome Biol*. 11:R86.
876 Gu H, Smith ZD, Bock C, Boyle P, Gnirke A, Meissner A. 2011. Preparation of reduced
877     representation bisulfite sequencing libraries for genome-scale DNA methylation
878     profiling. *Nat Protoc*. 6:468-481.
879 Gui D, Jia K, Xia J, Yang L, Chen J, Wu Y, Yi M. 2013. *De novo* assembly of the Indo-Pacific
880     humpback dolphin leucocyte transcriptome to identify putative genes involved in the
881     aquatic adaptation and immune response. *PLoS One*. 8:e72417.
882 Gutenkunst RN, Hernandez RD, Williamson SH, Bustamante CD. 2009. Inferring the joint
883     demographic history of multiple populations from multidimensional SNP frequency data.
884     *PLoS Genetics*. 5:e1000695.
885 Haas BJ, Papanicolaou A, Yassour M, Grabherr M, Blood PD, Bowden J, Couger MB, Eccles D,
886     Li B, Lieber M*, et al.* 2013. *De novo* transcript sequence reconstruction from RNA-seq
887     using the Trinity platform for reference generation and analysis. *Nat Protoc*. 8:1494-
888     1512.
889 Han E, Sinsheimer JS, Novembre J. 2015. Fast and accurate site frequency spectrum estimation
890     from low coverage sequence data. *Bioinformatics*. 31:720-727.
891 Hancock-Hanser BL, Frey A, Leslie MS, Dutton PH, Archer FI, Morin PA. 2013. Targeted
892     multiplex next-generation sequencing: advances in techniques of mitochondrial and
893     nuclear DNA sequencing for population genomics. *Mol Ecol Resour*. 13:254-268.
894 Harris K, Nielsen R. 2013. Inferring demographic history from a spectrum of shared haplotype
895     lengths. *PLoS Genetics*. 9:e1003521.
896 Hayden S, Bekaert M, Crider TA, Mariani S, Murphy WJ, Teeling EC. 2010. Ecological
897     adaptation determines functional mammalian olfactory subgenomes. *Genome Res*. 20:1-
898     9.
899 Hedrick PW. 2000 *Genetics of Populations*. Jones and Bartlett Publishers, Sudbury, MA.
900 Helyar SJ, Hemmer-Hansen J, Bekkevold D, Taylor MI, Ogden R, Limborg MT, Cariani A,
901     Maes GE, Diopere E, Carvalho GR*, et al.* 2011. Application of SNPs for population
902     genetics of nonmodel organisms: new opportunities and challenges. *Mol Ecol Resour*.
903     11:123-136.
904 Higdon JW, Bininda-Emonds ORP, Beck RMD, Ferguson SH. 2007. Phylogeny and divergence
905     of the pinnipeds (Carnivora: Mammalia) assessed using a multigene dataset. *BMC Evol
906     Biol*. 7:216.
907 Hodges E, Rooks M, Xuan Z, Bhattacharjee A, Gordon DB, Brizuela L, McCombie WR,
908     Hannon GJ. 2009. Hybrid selection of discrete genomic intervals on custom-designed
909     microarrays for massively parallel sequencing. *Nat Protoc*. 4:960-974.

910 Hoffman JI. 2011. Gene discovery in the Antarctic fur seal (*Arctocephalus gazella*) skin
911     transcriptome. *Mol Ecol Resour*. 11:703-710.
912 Hoffman JI, Nicholas HJ. 2011. A novel approach for mining polymorphic microsatellite
913     markers *in silico*. *PLoS One*. 6:e23283.
914 Hoffman JI, Simpson F, David P, Rijks JM, Kuiken T, Thorne MAS, Lacy RC, Dasmahapatra
915     KK. 2014. High-throughput sequencing reveals inbreeding depression in a natural
916     population. *Proc Natl Acad Sci USA*. 111:3775-3780.
917 Hoffman JI, Thorne MAS, Trathan PN, Forcada J. 2013. Transcriptome of the dead:
918     characterisation of immune genes and marker development from necropsy samples in a
919     free-ranging marine mammal. *BMC Genomics*. 14:52.
920 Hoffman JI, Tucker R, Bridgett SJ, Clark MS, Forcada J, Slate J. 2012. Rates of assay success
921     and genotyping error when single nucleotide polymorphism genotyping in non-model
922     organisms: a case study in the Antarctic fur seal. *Mol Ecol Resour*. 12:861-872.
923 Hohenlohe PA, Bassham S, Etter PD, Stiffler N, Johnson EA, Cresko WA. 2010. Population
924     genomics of parallel adaptation in threespine stickleback using sequenced RAD tags.
925     *PLoS Genet*. 6:e1000862.
926 Holt C, Yandell M. 2011. MAKER2: an annotation pipeline and genome-database management
927     tool for second-generation genome projects. *BMC Bioinformatics*. 12:491.
928 Humble E, Martinez-Barrio A, Forcada J, Trathan PN, Thorne MAS, Hoffmann M, Wolf JBW,
929     Hoffman JI. 2016. A draft fur seal genome provides insights into factors affecting SNP
930     validation and how to mitigate them. *Mol Ecol Resour*.
931 Jackson JA, Baker CS, Vant M, Steel DJ, Medrano-González L, Palumbi SR. 2009. Big and
932     slow: phylogenetic estimates of molecular evolution in baleen whales (suborder
933     Mysticeti). *Mol Biol Evol*. 26:2427-2440.
934 Jones FC, Grabherr MG, Chan YF, Russell P, Mauceli E, Johnson J, Swofford R, Pirun M, Zody
935     MC, White S*, et al.* 2012. The genomic basis of adaptive evolution in threespine
936     sticklebacks. *Nature*. 484:55-61.
937 Kajitani R, Toshimoto K, Noguchi H, Toyoda A, Ogura Y, Okuno M, Yabana M, Harada M,
938     Nagayasu E, Maruyama H*, et al.* 2014. Efficient *de novo* assembly of highly
939     heterozygous genomes from whole-genome shotgun short reads. *Genome Res*. 24:1384-
940     1395.
941 Keane M, Semeiks J, Webb AE, Li YI, Quesada V, Craig T, Madsen LB, van Dam S, Brawand
942     D, Marques PI*, et al.* 2015. Insights into the evolution of longevity from the bowhead
943     whale genome. *Cell Reports*. 10:112-122.
944 Khudyakov JI, Champagne CD, Preeyanon L, Ortiz RM, Crocker DE. 2015a. Muscle
945     transcriptome response to ACTH administration in a free-ranging marine mammal.
946     *Physiol Genomics*. 47:318-330.
947 Khudyakov JI, Preeyanon L, Champagne CD, Ortiz RM, Crocker DE. 2015b. Transcriptome
948     analysis of northern elephant seal (*Mirounga angustirostris*) muscle tissue provides a
949     novel molecular resource and physiological insights. *BMC Genomics*. 16:64.
950 Kishida T, Thewissen JGM, Hayakawa T, Imai H, Agata K. 2015. Aquatic adaptation and the
951     evolution of smell and taste in whales. *Zoolog Lett*. 1:9.
952 Koepfli K-P, Paten B, Genome 10K Community of Scientists, O'Brien SJ. 2015. The Genome
953     10K Project: a way forward. *Annu Rev Anim Biosci*. 3:57-111.
954 Korneliussen TS, Albrechtsen A, Nielsen R. 2014. ANGSD: Analysis of Next Generation
955     Sequencing Data. *BMC Bioinformatics*. 15:356.

956 Künstner A, Wolf JBW, Backström N, Whitney O, Balakrishnan CN, Day L, Edwards SV, Janes
957     DE, Schlinger BA, Wilson RK, *et al.* 2010. Comparative genomics based on massive
958     parallel transcriptome sequencing reveals patterns of substitution and selection across 10
959     bird species. *Mol Ecol*. 19:266-276.
960 Lamichhaney S, Berglund J, Almén MS, Maqbool K, Grabherr M, Martinez-Barrio A,
961     Promerová M, Rubin C-J, Wang C, Zamani N, *et al.* 2015. Evolution of Darwin's finches
962     and their beaks revealed by genome sequencing. *Nature*. 518:371-375.
963 Langmead B, Trapnell C, Pop M, Salzberg SL. 2009. Ultrafast and memory-efficient alignment
964     of short DNA sequences to the human genome. *Genome Biol*. 10:R25.
965 Lemmon AR, Emme SA, Lemmon EM. 2012. Anchored hybrid enrichment for massively high-
966     throughput phylogenomics. *Syst Biol*. 61:727-744.
967 Li B, Dewey CN. 2011. RSEM: accurate transcript quantification from RNA-Seq data with or
968     without a reference genome. *BMC Bioinformatics*. 12:323.
969 Li H, Durbin R. 2009. Fast and accurate short read alignment with Burrows-Wheeler transform.
970     *Bioinformatics*. 25:1754-1760.
971 Li H, Durbin R. 2011. Inference of human population history from individual whole-genome
972     sequences. *Nature*. 475:493-496.
973 Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R,
974     1000 Genome Project Data Processing Subgroup. 2009. The Sequence Alignment/Map
975     form and SAMtools. *Bioinformatics*. 25:2078-2079.
976 Li S, Jakobsson M. 2012. Estimating demographic paramaters from large-scale population
977     genomic data using Approximate Bayesian Computation. *BMC Genet*. 13:22.
978 Li Y, Hu Y, Bolund L, Wang J. 2010. State of the art *de novo* assembly of human genomes from
979     massively parallel sequencing data. *Human Genomics* 4:271-277.
980 Lindblad-Toh K, Garber M, Zuk O, Lin MF, Parker BJ, Washietl S, Kheradpour P, Ernst J,
981     Jordan G, Mauceli E, *et al.* 2011. A high-resolution map of human evolutionary
982     constraint using 29 mammals. *Nature*. 478:476-482.
983 Lindqvist C, Schuster SC, San Y, Talbot SL, Qi J, Ratan A, Tomsho LP, Kasson L, Zeyl E, Aars
984     J, *et al.* 2010. Complete mitochondrial genome of a Pleistocene jawbown unveils the
985     origin of polar bear. *Proc Natl Acad Sci USA*. 107:5053-5057.
986 Liu S, Lorenzen ED, Fumagalli M, Li B, Harris K, Xiong Z, Zhou L, Korneliussen TS, Somel M,
987     Babbitt C, *et al.* 2014a. Population genomics reveal recent speciation and rapid
988     evolutionary adaptation in polar bears. *Cell*. 157:785-794.
989 Liu X, Fun Y-X. 2015. Exploring population size changes using SNP frequency spectra. *Nat
990     Genet*. 47:555-559.
991 Liu Y, Zhou J, White KP. 2014b. RNA-seq differential expression studies: more sequence or
992     more replication? *Bioinformatics*. 30:301-304.
993 Lotterhos KE, Whitlock MC. 2014. Evaluation of demographic history and neutral
994     parameterization on the performance of $F_{ST}$ outlier tests. *Mol Ecol*. 23:2178-2192.
995 Louis M, Viricel A, Lucas T, Peltier H, Alfonsi E, Berrow S, Brownlow A, Covelo P, Dabin W,
996     Deaville R, *et al.* 2014. Habitat-driven population structure of bottlenose dolphins,
997     *Tursiops truncatus*, in the North-east Atlantic. *Mol Ecol*. 23:857-874.
998 Love MI, Huber W, Anders S. 2014. Moderated estimation of fold change and dispersion for
999     RNA-seq data with DESeq2. *Genome Biol*. 15:550.
1000 MacManes MD. 2014. On the optimal trimming of high-throughput mRNA sequence data. *Front
1001     Genet*. 5:13.

1002   MacManes MD. 2016. Establishing evidence-based best practice for the *de novo* assembly and
1003          evaluation of transcriptomes from non-model organisms. *bioRxiv*.  doi:
1004          http://dx.doi.org/10.1101/035642.
1005   Magera AM, Mills Flemming JE, Kaschner K, Christensen LB, Lotze HK. 2013. Recovery
1006          trends in marine mammal populations. *PLoS One*. 8:e77908.
1007   Malenfant RM, Coltman DW, Davis CS. 2015. Design of a 9K Illumina BeadChip for polar
1008          bears (*Ursus maritimus*) from RAD and transcriptome sequencing. *Mol Ecol Resour*.
1009          15:587-600.
1010   Mamanova L, Coffey AJ, Scott CE, Kozarewa I, Turner EH, Kumar A, Howard E, Shendure J,
1011          Turner DJ. 2010. Target-enrichment strategies for next-generation sequencing. *Nat*
1012          *Methods*. 7:111-118.
1013   Mancia A, Abelli L, Kucklick JR, Rowles TK, Wells RS, Balmer BC, Hohn AA, Baatz JE, Ryan
1014          JC. 2015. Microarray applications to understand the impact of exposure to environmental
1015          contaminants in wild dolphins (*Tursiops truncatus*). *Mar Genomics*. 19:47-57.
1016   Mancia A, Lundqvist ML, Romano TA, Peden-Adams MM, Fair PA, Kindy MS, Ellis BC,
1017          Gattoni-Celli S, McKillen DJ, Trent HF, *et al.* 2007. A dolphin peripheral blood
1018          leukocyte cDNA microarray for studies of immune function and stress reactions. *Dev*
1019          *Comp Immunol*. 31:520-529.
1020   Mancia A, Ryan JC, Chapman RW, Wu Q, Warr GW, Gulland FMD, Van Dolah FM. 2012.
1021          Health status, infection and disease in California sea lions (*Zalophus californianus*)
1022          studied using a canine microarray platform and machine-learning approaches. *Dev Comp*
1023          *Immunol*. 36:629-637.
1024   Mancia A, Warr GW, Chapman RW. 2008. A transcriptomic analysis of the stress induced by
1025          capture-release health assessment studies in wild dolphins (*Tursiops truncatus*). *Mol*
1026          *Ecol*. 17:2581-2589.
1027   Mastretta-Yanes A, Arrigo N, Alvarez N, Jorgensen TH, Piñero D, Emerson BC. 2015.
1028          Restriction site-associated DNA sequencing, genotyping error estimation and *de novo*
1029          assembly optimization for population genetic inference. *Mol Ecol Resour*. 15:28-41.
1030   McCormack JE, Faircloth BC, Crawford NG, Gowaty PA, Brumfield RT, Glenn TC. 2012.
1031          Ultraconserved elements are novel phylogenomic markers that resolve placental mammal
1032          phylogeny when combined with species-tree analysis. *Genome Res*. 22:746-754.
1033   McGowen MR. 2011. Toward the resolution of an explosive radiation - a multilocus phylogeny
1034          of oceanic dolphins (Delphinidae). *Mol Phylogenet Evol*. 60:345-357.
1035   McGowen MR, Clark C, Gatesy J. 2008. The vestigial olfactory receptor subgenome of
1036          odontocete whales: phylogenetic congruence between gene-tree reconciliation and
1037          supermatrix methods. *Syst Biol*. 57:574-590.
1038   McGowen MR, Gatesy J, Wildman DE. 2014. Molecular evolution tracks macroevolutionary
1039          transitions in Cetacea. *Trends Ecol Evol*. 29:336-346.
1040   McGowen MR, Grossman LI, Wildman DE. 2012. Dolphin genome provides evidence for
1041          adaptive evolution of nervous system genes and a molecular rate slowdown. *Proc R Soc*
1042          *Lond B Biol Sci*. 279:3643-3651.
1043   McGowen MR, Spaulding M, Gatesy J. 2009. Divergence date estimation and a comprehensive
1044          molecular tree of extant cetaceans. *Mol Phylogenet Evol*. 53:891-906.
1045   McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, Garimella K,
1046          Altshuler D, Gabriel S, Daly M*, et al.* 2010. The Genome Analysis Toolkit: A

MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res*. 20:1297-1303.

McTavish EJ, Hillis DM. 2015. How do SNP ascertainment schemes and population demographics affect inferences about population history? *BMC Genomics*. 16:266.

Meredith RW, Gatesy J, Emerling CA, York VM, Springer MS. 2013. Rod monochromacy and the coevolution of cetacean retinal opsins. *PLoS Genetics*. 9:e1003432.

Meyer M, Kircher M, Gansauge M-T, Li H, Racimo F, Mallick S, Schraiber JG, Jay F, Prüfer K, de Filippo C*, et al.* 2012. A high-coverage genome sequence from an archaic Denisovan individual. *Science*. 338:222-226.

Miller MR, Dunham JP, Amores A, Cresko WA, Johnson EA. 2007. Rapid and cost-effective polymorphism identification and genotyping using restriction site associated DNA (RAD) markers. *Genome Res*. 17:240-248.

Miller W, Schuster SC, Welch AJ, Ratan A, Bedoya-Reina OC, Zhao F, Kim HL, Burhans RC, Drautz DI, Wittekindt NE*, et al.* 2012. Polar and brown bear genomes reveal ancient admixture and demographic footprints of past climate change. *Proc Natl Acad Sci USA*. 109:E2382-E2390.

Mirceta S, Signore AV, Burns JM, Cossins AR, Campbell KL, Berenbrink M. 2013. Evolution of mammalian diving capacity traced by myoglobin net surface charge. *Science*. 340:1234192.

Morin PA, Archer FI, Foote AD, Vilstrup J, Allen EE, Wade P, Durban JW, Parsons K, Pitman R, Li L*, et al.* 2010a. Complete mitochondrial genome phylogeographic analysis of killer whales (*Orcinus orca*) indicates multiple species. *Genome Res*. 20:908-916.

Morin PA, Luikart G, Wayne RK, SNP workshop group. 2004. SNPs in ecology, evolution and conservation. *Trends Ecol Evol*. 19:208-216.

Morin PA, Martien KK, Archer FI, Cipriano F, Steel D, Jackson J, Taylor BL. 2010b. Applied conservation genetics and the need for quality control and reporting of genetic data used in fisheries and wildlife management. *J Hered*. 101:1-10.

Morin PA, Parsons KM, Archer FI, Ávila-Arcos M, Barrett-Lennard LG, Dalla Rosa L, Duchêne S, Durban JW, Ellis GM, Ferguson SH*, et al.* 2015. Geographic and temporal dynamics of a global radiation and diversification in the killer whale. *Mol Ecol*. 24:3964-3979.

Moura AE, Kenny JG, Chaudhuri R, Hughes MA, Welch AJ, Reisinger RR, de Bruyn PJN, Dahlheim ME, Hall N, Hoelzel AR. 2014a. Population genomics of the killer whale indicates ecotype evolution in sympatry involving both selection and drift. *Mol Ecol*. 23:5179-5192.

Moura AE, Nielsen SCA, Vilstrup JT, Moreno-Mayar JV, Gilbert MTP, Gray HWI, Natoli A, Möller L, Hoelzel AR. 2013. Recent diversification of a marine genus (*Tursiops* spp.) tracks habitat preference and environmental change. *Syst Biol*. 62:865-877.

Moura AE, van Rensburg CJ, Pilot M, Tehrani A, Best PB, Thornton M, Plön S, de Bruyn PJN, Worley KC, Gibbs RA*, et al.* 2014b. Killer whale nuclear genome and mtDNA reveal widespread population bottleneck during the last glacial maximum. *Mol Biol Evol*. 31:1121-1131.

Nadeau NJ, Ruiz M, Salazar P, Counterman B, Alejandro Medina J, Ortiz-Zuazaga H, Morrison A, McMillan WO, Jiggins CD, Papa R. 2014. Population genomics of parallel hybrid zones in the mimetic butterflies, *H. melpomene* and *H. erato*. *Genome Res*. 24:1316-1333.

1092 Narum SR, Buerkle CA, Davey JW, Miller MR, Hohenlohe PA. 2013. Genotyping-by-
1093     sequencing in ecological and conservation genomics. *Mol Ecol*. 22:2841-2847.
1094 Narum SR, Hess JE. 2011. Comparison of $F_{ST}$ outlier tests for SNP loci under selection. *Mol*
1095     *Ecol Resour*. 11:184-194.
1096 Nelson TM, Apprill A, Mann J, Rogers TL, Brown MV. 2015. The marine mammal microbiome:
1097     current knowledge and future directions. *Microbiology Australia*. 36:8-13.
1098 Ng SB, Turner EH, Robertson PD, Flygare SD, Bigham AW, Lee C, Shaffer T, Wong M,
1099     Bhattacharjee A, Eichler EE, *et al.* 2009. Targeted capture and massively parallel
1100     sequencing of twelve human exomes. *Nature*. 461:272-276.
1101 Nielsen R, Paul JS, Anders A, Song YS. 2011. Genotype and SNP calling from next-generation
1102     sequencing data. *Nat Rev Genet*. 12:433-451.
1103 Noonan JP, Coop G, Kudaravalli S, Smith D, Krause J, Alessi J, Chen F, Platt D, Pääbo S,
1104     Pritchard JK, *et al.* 2006. Sequencing and analysis of Neanderthal genomic DNA.
1105     *Science*. 314:1113-1118.
1106 Nosek BA, Alter G, Banks GC, Borsboom D, Bowman SD, Breckler SJ, Buck S, Chambers CD,
1107     Chin G, Christensen G, *et al.* 2015. Promoting an open research culture: Author
1108     guidelines for journals could help to promote transparency, openness, and reproducibility.
1109     *Science*. 348:1422-1425.
1110 O'Rawe JA, Ferson S, Lyon GJ. 2015. Accounting for uncertainty in DNA sequencing data.
1111     *Trends Genet*. 31:61-66.
1112 Olsen MT, Volny VH, Bérubé M, Dietz R, Lydersen C, Kovacs KM, Dodd RS, Palsbøll PJ.
1113     2011. A simple route to single-nucleotide polymorphisms in a nonmodel species:
1114     identification and characterization of SNPs in the Arctic ringed seal (*Pusa hispida*
1115     *hispida*). *Mol Ecol Resour*. 11:9-19.
1116 Orlando L, Ginolhac A, Zhang G, Froese D, Albrechtsen A, Stiller M, Schubert M, Cappellini E,
1117     Petersen B, Moltke I, *et al.* 2013. Recalibrating *Equus* evolution using the genome
1118     sequence of an early Middle Pleistocene horse. *Nature*. 499:74-78.
1119 Pabuwal V, Boswell M, Pasquali A, Wise SS, Kumar S, Shen Y, Garcia T, Lacerte C, Wise JP,
1120     Jr., Wise JP, Sr., *et al.* 2013. Transcriptomic analysis of cultured whale skin cells exposed
1121     to hexavalent chromium [Cr(VI)]. *Aquat Toxicol*. 134-135:74-81.
1122 Parker J, Tsagkogeorga G, Cotton JA, Liu Y, Provero P, Stupka E, Rossiter SJ. 2013. Genome-
1123     wide signatures of convergent evolution in echolocating mammals. *Nature*. 502:228-231.
1124 Paszkiewicz KH, Farbox A, O'Neill P, Moore K. 2014. Quality control on the frontier. *Front*
1125     *Genet*. 5:157.
1126 Patro R, Duggal G, Kingsford C. 2015. Accurate, fast, and model-aware transcript expression
1127     quantification with Salmon. *bioRxiv*. doi: http://dx.doi.org/10.1101/021592.
1128 Peterson BK, Weber JN, Kay EH, Fisher HS, Hoekstra HE. 2012. Double digest RADseq: an
1129     inexpensive method for *de novo* SNP discovery and genotyping in model and non-model
1130     species. *PLoS One*. 7:e37135.
1131 Poh Y-P, Domingues VS, Hoekstra HE, Jensen JD. 2014. On the prospect of identifying adaptive
1132     loci in recently bottlenecked populations. *PLoS One*. 9:e110579.
1133 Poland JA, Brown PJ, Sorrells ME, Jannink J-L. 2012. Development of high-density genetic
1134     maps for barley and wheat using a novel two-enzyme genotyping-by-sequencing
1135     approach. *PLoS One*. 7:e32253.
1136 Polanowski AM, Robbins J, Chandler D, Jarman SN. 2014. Epigenetic estimation of age in
1137     humpback whales. *Mol Ecol Resour*. 14:976-987.

1138  Puritz JB, Hollenbeck CM, Gold JR. 2014. *dDocent*: a RADseq, variant-calling pipeline
1139      designed for population genomics of non-model organisms. *PeerJ*. 2:e431.
1140  Rasmussen M, Li Y, Lindgreen S, Pedersen JS, Albrechtsen A, Moltke I, Metspalu M, Metspalu
1141      E, Kivisild T, Gupta R, *et al.* 2010. Ancient human genome sequence of an extinct
1142      Palaeo-Eskimo. *Nature*. 463:757-762.
1143  Riesch R, Barrett-Lennard LG, Ellis GM, Ford JKB, Deecke VB. 2012. Cultural traditions and
1144      the evolution of reproductive isolation: ecological speciation in killer whales? *Biol J Linn*
1145      *Soc Lond*. 2012:1-17.
1146  Robinson JD, Coffman AJ, Hickerson MJ, Gutenkunst RN. 2014. Sampling strategies for
1147      frequency spectrum-based population genomic inference. *BMC Evol Biol*. 14:254.
1148  Robinson MD, McCarthy DJ, Smyth GK. 2010. edgeR: a Bioconductor package for differential
1149      expression analysis of digital gene expression data. *Bioinformatics*. 26:139-140.
1150  Ruan R, Guo A-H, Hao Y-J, Zheng J-S, Wang D. 2015. *De novo* assembly and characterization
1151      of narrow-ridged finless porpoise renal transcriptome and identification of candidate
1152      genes involved in osmoregulation. *Int J Mol Sci*. 16:2220-2238.
1153  Ruegg K, Rosenbaum HC, Anderson EC, Engel M, Rothschild A, Baker CS, Palumbi SR. 2013.
1154      Long-term population size of the North Atlantic humpback whale within the context of
1155      worldwide population structure. *Cons Gen*. 14:103-114.
1156  Schiffels S, Durbin R. 2014. Inferring human population size and separation history from
1157      multiple genome sequences. *Nat Genet*. 46:919-925.
1158  Schubert M, Lindgreen S, Orlando L. 2016. AdapterRemoval v2: rapid adapter trimming,
1159      identification, and read merging. *BMC Res Notes*. 9:88.
1160  Schurch NJ, Schofield P, Gierlinski M, Cole C, Sherstnev A, Singh V, Wrobel N, Gharbi K,
1161      Simpson GG, Owen-Hughes T, *et al.* 2015. Evaluation of tools for differential gene
1162      expression analysis by RNA-seq on a 48 biological replicate experiment. *arXive*.
1163      1505:02017.
1164  Seim I, Ma S, Zhou X, Gerashchenko MV, Lee SG, Suydam R, George JC, Bickham JW,
1165      Gladyshev VN. 2014. The transcriptome of the bowhead whale *Balaena mysticetus*
1166      reveals adaptations of the longest-lived mammal. *Aging*. 6:879-899.
1167  Shafer ABA, Cullingham CI, Côté SD, Coltman DW. 2010. Of glaciers and refugia: a decade of
1168      study sheds new light on the phylogeographic patterns of northwestern North America.
1169      *Mol Ecol*. 19:4589-4621.
1170  Shafer ABA, Davis CS, Coltman DW, Stewart REA. 2014. Microsatellite assessment of walrus
1171      (*Odobenus rosmarus rosmarus*) stocks in Canada. *NAMMCO Scientific Publications*. 9.
1172  Shafer ABA, Gattepaille LM, Stewart REA, Wolf JBW. 2015. Demographic inferences using
1173      short-read genomic data in an approximate Bayesian computation framework: *in silico*
1174      evaluation of power, biases and proof of concept in Atlantic walrus. *Mol Ecol*. 24:328-
1175      345.
1176  Shen Y-Y, Zhou W-P, Zhou T-C, Zeng Y-N, Li G-M, Irwin DM, Zhang Y-P. 2012. Genome-
1177      wide scan for bats and dolphin to detect their genetic basis for new locomotive styles.
1178      *PLoS One*. 7:e46455.
1179  Smith-Unna RD, Boursnell C, Patro R, Hibberd JM, Kelly S. 2015. TransRate: reference free
1180      quality assessment of *de-novo* transcriptome assemblies. *bioRxiv*.
1181  Spies D, Ciaudo C. 2015. Dynamics in transcriptomics: advancements in RNA-seq time course
1182      and downstream analysis. *Comput Struct Biotechnol J*. 13:469-477.

Springer MS, Signore AV, Paijmans JLA, Vélez-Juarbe J, Domning DP, Bauer CE, He K, Crerar L, Campos PF, Murphy WJ, *et al.* 2015. Interordinal gene capture, the phylogenetic position of Steller's sea cow based on molecular and morphological data, and the macroevolutionary history of Sirenia. *Mol Phylogenet Evol*. 91:178-193.

Springer MS, Starrett J, Morin PA, Lanzetti A, Hayashi C, Gatesy J. 2016. Inactivation of *C4orf26* in toothless placental mammals. *Mol Phylogenet Evol*. 95:34-45.

Sremba AL, Martin AR, Baker CS. 2015. Species identification and likely catch time preiod of whale bones from South Georgia. *Mar Mamm Sci*. 31:122-132.

Stanke M, Keller O, Gunduz I, Hayes A, Waack S, Morgenstern B. 2006. AUGUSTUS: *ab initio* prediction of alternative transcripts. *Nucleic Acids Res*. 34:W435-W439.

Stein LD. 2010. The case for cloud computing in genome informatics. *Genome Biol*. 11:207.

Stinchcombe JR, Hoekstra HE. 2008. Combining population genomics and quantitative genetics: finding genes underlying ecologically important traits. *Heredity*. 100:158-170.

Tabuchi M, Veldhoen N, Dangerfield N, Jeffries S, Helbing CC, Ross PS. 2006. PCB-related alteration of thyroid hormones and thyroid hormone receptor gene expression in free-ranging harbor seals (*Phoca vitulina*). *Environ Health Perspect*. 114:1024-1031.

Taylor BL, Gemmell NJ. 2016. Emerging technologies to conserve biodiversity: further opportunities via genomics. Response to Pimm *et al. Trends Ecol Evol*. 31:171-172.

The *Heliconius* Genome Consortium. 2012. Butterfly genome reveals promiscuous exchange of mimicry adaptations among species. *Nature*. 487:94-98.

Thomsen PF, Kielgast J, Iversen LL, Møller PR, Rasmussen M, Willerslev E. 2012. Detection of a diverse marine fish fauna using environmental DNA from seawater samples. *PLoS One*. 7:e41732.

Towns J, Cockerill T, Dahan M, Foster I, Gaither K, Grimshaw A, Hazlewood V, Lathrop S, Lifka D, Peterson GD, *et al.* 2014. XSEDE: accelerating scientific discovery. *Computing in Science and Engineering*. 16:62-74.

Tsagkogeorga G, McGowen MR, Davies KT, Jarman S, Polanowski A, Bertelsen MF, Rossiter SJ. 2015. A phylogenomic analysis of the role and timing of molecular adaptation in the aquatic transition of cetartiodactyl mammals. *R Soc Open Sci*. 2:150156.

van Dijk EL, Auger H, Jaszczyzyn Y, Thermes C. 2014. Ten years of next-generation sequencing technology. *Trends Genet*. 30:418-426.

VanRaden PM, Sun C, O'Connell JR. 2015. Fast imputation using medium or low-coverage sequence data. *BMC Genet*. 16:82.

Villar D, Berthelot C, Aldridge S, Rayner TF, Lukk M, Pignatelli M, Park TJ, Deaville R, Erichsen JT, Jasinska AJ, *et al.* 2015. Enhancer evolution across 20 mammalian species. *Cell*. 160:554-566.

Viricel A, Pante E, Dabin W, Simon-Bouhet B. 2014. Applicability of RAD-tag genotyping for interfamilial comparisons: empirical data from two cetaceans. *Mol Ecol Resour*. 14:597-605.

Viricel A, Rosel PE. 2014. Hierarchical population structure and habitat differences in a highly mobile marine species: the Atlantic spotted dolphin. *Mol Ecol*. 23:5018-5035.

Wolf JB. 2013. Principles of transcriptome analysis and gene expression quantification: an RNA-seq tutorial. *Mol Ecol Resour*. 13:559-572.

Xiong Y, Brandley MC, Xu S, Zhou K, Yang G. 2009. Seven new dolphin mitochondrial genomes and a time-calibrated phylogeny of whales. *BMC Evol Biol*. 9:20.

1228    Yandell M, Ence D. 2012. A beginner's guide to eukaryotic genome annotation. *Nat Rev Genet*.
1229            13:329-342.
1230    Yeh R-F, Lim LP, Burge CB. 2001. Computational inference of homologous gene structures in
1231            the human genome. *Genome Res*. 11:803-816.
1232    Yim H-S, Cho YS, Guang X, Kang SG, Jeong J-Y, Cha S-S, Oh H-M, Lee J-H, Yang EC, Kwon
1233            KK*, et al.* 2014. Minke whale genome and aquatic adaptation in cetaceans. *Nat Genet*.
1234            46:88-92.
1235    Zhao Q-Y, Wang Y, Kong Y-M, Luo D, Li X, Hao P. 2011. Optimizing *de novo* transcriptome
1236            assembly from short-read RNA-Seq data: a comparative study. *BMC Bioinformatics*.
1237            12:S2.
1238    Zhou X, Sun F, Xu S, Fan G, Zhu K, Liu X, Chen Y, Shi C, Yang Y, Huang Z*, et al.* 2013. Baiji
1239            genomes reveal low genetic variability and new insights into secondary aquatic
1240            adaptations. *Nat Commun*. 4:2708.
1241    Zou Z, Zhang J. 2015. No genome-wide protein sequence convergence for echolocation. *Mol
1242            Biol Evol*. 32:1237-1241.
1243

1244 Table 1. Current and commonly used tools for analysis of genomic data generated in non-model organisms. Please note that this list is
1245 not exhaustive and new computational tools are continuously being developed.
1246

| Computational Tool | Purpose | Strengths/Weaknesses | Reference |
|---|---|---|---|
| *RADseq** | | | |
| STACKS | quality filtering, *de novo* assembly or reference-aligned read mapping, variant genotyping | scalable (new data can be compared against existing locus catalog); flexible filtering and export options; recently implemented a gapped alignment algorithm to process insertion-deletion (indel) mutations; secondary algorithm adjusts SNP calls using population-level allele frequencies; compatible with input data from multiple RADseq methods | Catchen et al. (2011; 2013), http://catchenlab.life.illinois.edu/stacks/ |
| PyRAD | quality filtering, *de novo* assembly, read mapping, variant genotyping | efficiently processes indel mutations, thus optimal for analysis of highly divergent species; high speed and quality of paired-end library assemblies; compatible with input data from multiple RADseq methods | Eaton (2014) |
| TASSEL-GBS | quality filtering, reference-aligned read mapping, variant genotyping | optimized for single-end data from large sample sizes (tens of thousands of individuals) with a reference genome; performs genome-wide association studies | Glaubitz et al. (2014) |
| dDocent | quality trimming, *de novo* assembly, read mapping, variant genotyping | beneficial in analysis of paired-end data; identifies both SNP and indel variants; most appropriate for ezRAD and ddRAD data | Puritz et al. (2014) |
| AftrRAD | quality filtering, *de novo* assembly, read mapping, variant genotyping | identifies both SNP and indel variants; computationally faster than STACKS and PyRAD | Sovic et al. (2015) |
| *Array-based high-throughput sequencing* | | | |
| Affymetrix Axiom™ Analysis Suite, Illumina® GenomeStudio | genotype scoring | visualization of genotype clusters; quality scores assigned to genotype calls allow user-specific filtering; manual editing possible | |
| *Whole genome sequencing* | | | |
| AdapterRemoval v2, Trimmomatic | trim raw sequences | remove adapter sequences and low-quality bases prior to assembly | Bolger et al. (2014), Schubert et al. (2016) |
| ALLPATHS-LG, PLATANUS, SOAPdenovo | *de novo* genome assembly | designed for short-read sequences of large heterozygous genomes | Li et al. (2010), Gnerre et al. (2011), Kajitani et al. (2014) |
| AUGUSTUS, GenomeScan, MAKER2 | gene annotation | highly accurate evidence-driven or BLASTX-guided gene prediction (Yandell and Ence 2012) | Yeh et al. (2001), Stanke et al. (2006), Holt and Yandell (2011) |

| Software | Purpose | Description | Reference |
|---|---|---|---|
| Bowtie, bwa | read mapping | rapid short-read alignment with compressed reference genome index, but limited number of acceptable mismatches per alignment (Flicek and Birney 2009) | Langmead et al. (2009), Li and Durbin (2009) |
| SAMtools | data processing, variant calling | multi-purpose tool that conducts file conversion, alignment sorting, PCR duplicate removal, and variant (SNP and indel) calling for SAM/BAM/CRAM files | Li et al. (2009) |
| GATK | data processing and quality control, variant calling | suitable for data with low to high mean read depth across the genome; initially optimized for large human datasets, then modified for use with non-model organisms | McKenna et al. (2010), DePristo et al. (2011) |
| ANGSD/NGStools | data processing, variant calling, estimation of diversity metrics, population genomic analyses | suitable for data with low mean read depth, including palaeogenomic data; allow downstream analyses such as D-statistics and SFS estimation | Fumagalli et al. (2014), Korneliussen et al. (2014) |
| *RNAseq* | | | |
| Fastx Toolkit, Trimmomatic | trim raw sequences | remove erroneous nucleotides from reads prior to assembly | MacManes (2014) |
| khmer diginorm, Trinity normalization | *in silico* read normalization | reduce memory requirements for assembly, but can result in fragmented assemblies and collapse heterozygosity | Brown et al. (2012); Haas et al. (2013) |
| Trinity | *de novo* and genome-guided transcriptome assembly | accurate assembly across conditions, but requires long runtime if normalization is not used (Zhao et al. 2011) | Haas et al. (2013) |
| bowtie, bowtie2, STAR | read alignment to genome or transcriptome assembly | required for many downstream analyses, but bowtie is computationally intensive and all produce very large output BAM files | Langmead et al. (2009), Dobin et al. (2013) |
| eXpress, kallisto, RSEM, Sailfish, Salmon | estimation of transcript abundance | RSEM requires computationally intensive read mapping back to the assembly; the others are faster streaming alignment, quasi-alignment, or alignment-free algorithms | Li and Dewey (2011), Patro et al. (2015) |
| DESeq, DESeq2, edgeR | differential expression analysis | exhibit highest true positive and lowest false positive rates in experiments with smaller sample sizes (Schurch et al. 2015) | Anders and Huber (2010), Robinson et al. (2010), Love et al. (2014) |
| blast2GO, Trinotate | functional annotation of assembled transcripts | complete annotation pipelines including gene ontology and pathway enrichment analyses | Conesa et al. (2005), Haas et al. (2013) |

1247 * This is a non-exhaustive list of software that focuses on *de novo* loci assembly and genotype calling for RADseq data, as many practitioners working on NMOs
1248 will not have access to a reference genome. Other programs (e.g., GATK and ANGSD) that undertake genotype calling using reference-aligned loci are described
1249 in the whole genome sequencing section.

**Figure 1**



Figure 1. Phylogenetic tree showing current genomic resources available for (A) cetaceans and (B) pinnipeds; relationships and branch lengths are based on molecular dating estimates from McGowen et al. (2009), McGowen (2011), and Higdon et al. (2007). Scale is in millions of years ago (MYA). Red circles indicate species with high-quality reference genomes; green stars indicate whole genome re-sequencing data; blue triangles indicate transcriptomes (generated by microarray or RNAseq); and black squares indicate RADseq data.

**Figure 2**



Figure 2. Number of marine mammal genomics publications from 1990 to 2015, categorized by primary methodology and research aim. Genomic methodologies include high-throughput single nucleotide polymorphism (SNP) genotyping and sequencing of mitogenomes, whole genomes (WGS), transcriptomes (generated by microarray or RNAseq), and reduced-representation genomic libraries (RRL). The "Other" category includes studies of microbiomes, BAC libraries, and large (~100) gene sets.

**Figure 3**



Figure 3. Number of BioProjects (gray bars) related to marine mammal genomics submitted from 2006 to 2015 to an online public database maintained by NCBI. Early BioProjects were largely microarray datasets. The number of projects created each year, as well as the yearly average (black dots ± SE) and maximum (×) size of data submitted in each BioProject, increased dramatically after 2011, reflecting advances in high-throughput sequencing technologies that facilitated their use in non-model systems.

**Figure 4**

| | 2002 --- 2009 | 2010 | 2011 | 2012 | 2013 | 2014 | 2015 | 2016 |
|---|---|---|---|---|---|---|---|---|
| Common bottlenose dolphin (*Tursiops truncatus*) | Sequence mitogenome[1] | | Sequence genome (2.8x)[2] | | Population mitogenomics[3] | | RADseq[4] — Improve genome (3.5x 454, 30x Illumina)[5] | Shotgun sequencing for SNP discovery[6] |
| Killer whale (*Orcinus orca*) | | Population mitogenomics[7] | | | | RADseq[8] — Sequence genome (20x)[9] | Sequence genome (200x)[5] | Population genome re-sequencing (avg 2x, N = 48)[10] |
| Antarctic fur seal (*Arctocephalus gazella*) | | | Sequence transcriptome[11] — Microsat. discovery[12] | SNP discovery[13] | | | | Sequence genome (200x) to validate SNPs[14] |
| Polar bear (*Ursus maritimus*) | Sequence mitogenome[15] | Ancient DNA mitogenome[16] | | Sequence genome (100x) & transcriptome[17] | | Population genome re-sequencing (avg 3.5x, N=61)[18] — RADseq & transcriptome sequencing for SNP discovery[19] | | |

1) Xiong et al. 2009; 2) Lindblad-Toh et al. 2011; 3) Moura et al. 2013; 4) Cammen et al. 2015; 5) Foote et al. 2015; 6) Louis et al. unpubl. data; 7) Morin et al. 2010; 8) Moura et al. 2014a; 9) Moura et al. 2014b; 10) Foote et al. 2016; 11) Hoffman 2011; 12) Hoffman and Nicholas 2011; 13) Hoffman et al. 2012; 14) Humble et al. 2016; 15) Arnason et al. 2002; 16) Lindqvist et al. 2010; 17) Miller et al. 2012; 18) Liu et al. 2014; 19) Malenfant et al. 2015

Figure 4. Timelines depicting the independent progression of genomic studies for four representative marine mammal species. Trajectories show the common progression for non-model species from mitogenome sequencing to whole genome sequencing, as well as from sequencing reference specimens to population-scale genomic sequencing. In addition, the timelines reveal the utility of genomic and transcriptomic sequencing for subsequent genetic marker development.

**Manuscripts submitted to Journal of Heredity**

Research aim:

- Phylogenomics
- Evolutionary adaptation & genome evolution
- Resource development
- Fitness
- Population genomics & demography
- Microbiome

- Other
- SNP
- RRL
- Transcriptome
- WGS
- Mitogenome

Number of Publications

Manuscripts submitted to Journal of Heredity

| Species | 2002 | 2009 | 2010 | 2011 | 2012 | 2013 | 2014 | 2015 | 2016 |
|---|---|---|---|---|---|---|---|---|---|
| Common bottlenose dolphin (*Tursiops truncatus*) | | Sequence mitogenome[1] | | Sequence genome (2.8x)[2] | | Population mitogenomics[3] | | RADseq[4]; Improve genome (3.5x 454, 30x Illumina)[5] | Shotgun sequencing for SNP discovery[6] |
| Killer whale (*Orcinus orca*) | | | Population mitogenomics[7] | | | | RADseq[8]; Sequence genome (20x)[9] | Sequence genome (200x)[5] | Population genome re-sequencing (avg 2x, N = 48)[10] |
| Antarctic fur seal (*Arctocephalus gazella*) | | | | Sequence transcriptome[11]; Microsat. discovery[12] | SNP discovery[13] | | | | Sequence genome (200x) to validate SNPs[14] |
| Polar bear (*Ursus maritimus*) | Sequence mitogenome[15] | | Ancient DNA mitogenome[16] | | Sequence genome (100x) & transcriptome[17] | | Population genome re-sequencing (avg 3.5x, N=61)[18] | RADseq & transcriptome sequencing for SNP discovery[19] | |

1) Xiong et al. 2009; 2) Lindblad-Toh et al. 2011; 3) Moura et al. 2013; 4) Cammen et al. 2015; 5) Foote et al. 2015; 6) Louis et al. unpubl. data; 7) Moura et al. 2010; 8) Moura et al. 2014a; 9) Moura et al. 2014b; 10) Foote et al. 2016; 11) Hoffman 2011; 12) Hoffman and Nicholas 2011; 13) Hoffman et al. 2012; 14) Humble et al. 2016; 15) Arnason et al. 2002; 16) Lindqvist et al. 2010; 17) Miller et al. 2012; 18) Liu et al. 2014; 19) Malenfant et al. 2015

Table 1. Current and commonly used tools for analysis of genomic data generated in non-model organisms. Please note that this list is not exhaustive and new computational tools are continuously being developed.

| Computational Tool | Purpose | Strengths/Weaknesses | Reference |
|---|---|---|---|
| *RADseq\** | | | |
| STACKS | quality filtering, *de novo* assembly or reference-aligned read mapping, variant genotyping | scalable (new data can be compared against existing locus catalog); flexible filtering and export options; recently implemented a gapped alignment algorithm to process insertion-deletion (indel) mutations; secondary algorithm adjusts SNP calls using population-level allele frequencies; compatible with input data from multiple RADseq methods | Catchen et al. (2011; 2013), http://catchenlab.life.illinois.edu/stacks/ |
| PyRAD | quality filtering, *de novo* assembly, read mapping, variant genotyping | efficiently processes indel mutations, thus optimal for analysis of highly divergent species; high speed and quality of paired-end library assemblies; compatible with input data from multiple RADseq methods | Eaton (2014) |
| TASSEL-GBS | quality filtering, reference-aligned read mapping, variant genotyping | optimized for single-end data from large sample sizes (tens of thousands of individuals) with a reference genome; performs genome-wide association studies | Glaubitz et al. (2014) |
| dDocent | quality trimming, *de novo* assembly, read mapping, variant genotyping | beneficial in analysis of paired-end data; identifies both SNP and indel variants; most appropriate for ezRAD and ddRAD data | Puritz et al. (2014) |
| AftrRAD | quality filtering, *de novo* assembly, read mapping, variant genotyping | identifies both SNP and indel variants; computationally faster than STACKS and PyRAD | Sovic et al. (2015) |
| *Array-based high-throughput sequencing* | | | |
| Affymetrix Axiom™ Analysis Suite, Illumina® GenomeStudio | genotype scoring | visualization of genotype clusters; quality scores assigned to genotype calls allow user-specific filtering; manual editing possible | |
| *Whole genome sequencing* | | | |
| AdapterRemoval v2, Trimmomatic | trim raw sequences | remove adapter sequences and low-quality bases prior to assembly | Bolger et al. (2014), Schubert et al. (2016) |
| ALLPATHS-LG, PLATANUS, SOAPdenovo | *de novo* genome assembly | designed for short-read sequences of large heterozygous genomes | Li et al. (2010), Gnerre et al. (2011), Kajitani et al. (2014) |
| AUGUSTUS, GenomeScan, MAKER2 | gene annotation | highly accurate evidence-driven or BLASTX-guided gene prediction (Yandell and Ence 2012) | Yeh et al. (2001), Stanke et al. (2006), Holt and Yandell (2011) |

| Software | Function | Description | Reference |
|---|---|---|---|
| Bowtie, bwa | read mapping | rapid short-read alignment with compressed reference genome index, but limited number of acceptable mismatches per alignment (Flicek and Birney 2009) | Langmead et al. (2009), Li and Durbin (2009) |
| SAMtools | data processing, variant calling | multi-purpose tool that conducts file conversion, alignment sorting, PCR duplicate removal, and variant (SNP and indel) calling for SAM/BAM/CRAM files | Li et al. (2009) |
| GATK | data processing and quality control, variant calling | suitable for data with low to high mean read depth across the genome; initially optimized for large human datasets, then modified for use with non-model organisms | McKenna et al. (2010), DePristo et al. (2011) |
| ANGSD/NGStools | data processing, variant calling, estimation of diversity metrics, population genomic analyses | suitable for data with low mean read depth, including palaeogenomic data; allow downstream analyses such as D-statistics and SFS estimation | Fumagalli et al. (2014), Korneliussen et al. (2014) |
| *RNAseq* | | | |
| Fastx Toolkit, Trimmomatic | trim raw sequences | remove erroneous nucleotides from reads prior to assembly | MacManes (2014) |
| khmer diginorm, Trinity normalization | *in silico* read normalization | reduce memory requirements for assembly, but can result in fragmented assemblies and collapse heterozygosity | Brown et al. (2012); Haas et al. (2013) |
| Trinity | *de novo* and genome-guided transcriptome assembly | accurate assembly across conditions, but requires long runtime if normalization is not used (Zhao et al. 2011) | Haas et al. (2013) |
| bowtie, bowtie2, STAR | read alignment to genome or transcriptome assembly | required for many downstream analyses, but bowtie is computationally intensive and all produce very large output BAM files | Langmead et al. (2009), Dobin et al. (2013) |
| eXpress, kallisto, RSEM, Sailfish, Salmon | estimation of transcript abundance | RSEM requires computationally intensive read mapping back to the assembly; the others are faster streaming alignment, quasi-alignment, or alignment-free algorithms | Li and Dewey (2011), Patro et al. (2015) |
| DESeq, DESeq2, edgeR | differential expression analysis | exhibit highest true positive and lowest false positive rates in experiments with smaller sample sizes (Schurch et al. 2015) | Anders and Huber (2010), Robinson et al. (2010), Love et al. (2014) |
| blast2GO, Trinotate | functional annotation of assembled transcripts | complete annotation pipelines including gene ontology and pathway enrichment analyses | Conesa et al. (2005), Haas et al. (2013) |

\* This is a non-exhaustive list of software that focuses on *de novo* loci assembly and genotype calling for RADseq data, as many practitioners working on NMOs will not have access to a reference genome. Other programs (e.g., GATK and ANGSD) that undertake genotype calling using reference-aligned loci are described in the whole genome sequencing section.

2

Cammen_SupMat_TableS1 - Marine mammal genomics - *JHered*

Table S1. Broad applications of genomic tools in studies of non-model organisms are provided with concrete examples of research areas drawn from the field of marine mammal genomics. The number of loci used in each study provides an estimate of the scope of the respective genomic tools and study, but represents the outcome of several filtering steps from raw sequence data that vary across studies. Further details of each method can be found in the listed references. Please note that this is not an exhaustive list. GBS: Genotyping by Sequencing; RADseq: restriction site-associated DNA sequencing; SNP: single nucleotide polymorphism; TSC: target sequence capture; WGS: whole genome sequencing.

| Method | # loci | Research area | Reference |
|---|---|---|---|
| *Evolutionary genomics: describe evolutionary history and adaptation* | | | |
| Mitogenome sequencing | Mitogenome | Cetacean phylogenomics | McGowen et al. (2009) |
| TSC | Mitogenome | Comparison of sub-fossil and modern killer whales | Foote et al. (2013) |
| TSC | >30kb coding sequence | Evolution of Sirenia | Springer et al. (2015) |
| WGS | Whole genome | Yangtze river dolphin genome analysis | Zhou et al. (2013) |
| WGS | Whole genome | Minke whale genome analysis | Yim et al. (2014) |
| WGS | Whole genome | Bowhead whale genome analysis | Keane et al. (2015) |
| WGS | Whole genome | Analysis of convergent evolution in marine mammal lineages | Foote et al. (2015) |
| WGS | 10,025 coding sequences | Positive selection in common bottlenose dolphin genome | McGowen et al. (2012) |
| WGS | Sensory genes | Analysis of gene loss in olfaction and taste in Antarctic minke whale | Kishida et al. (2015) |
| Genome re-seq | Whole genome | Speciation and adaptation in brown and polar bears | Liu et al. (2014) |
| Transcriptomics | 9,395 genes | Evolution of longevity in bowhead whales | Seim et al. (2014) |
| Transcriptomics | 103,077 unigenes | Osmoregulatory divergence in narrow-ridged finless porpoise | Ruan et al. (2015) |
| *Population genomics: characterize population structure and investigate demography* | | | |
| RADseq | 3,281 SNPs | Killer whale ecotype divergence | Moura et al. (2014) |
| RADseq (GBS) | 24,996 loci; 4,854 SNPs | Historical demography in Atlantic walrus | Shafer et al. (2015) |
| TSC | Mitogenome and 43-118 nuclear loci | Phylogeography and population genomics of cetaceans | Hancock-Hanser et al. (2013); Morin et al. (2015) |
| Genome re-seq | Whole genome | Demographic history, population differentiation, and ecotype divergence in killer whales | Foote et al. (2016) |

Cammen_SupMat_TableS1 - Marine mammal genomics - *JHered*

| | | | |
|---|---|---|---|
| *Adaptation genomics: describe relationships between genomic variation and fitness* | | | |
| RADseq | 83,148 loci; 14,585 SNPs | Effect of inbreeding depression on parasite infection in harbor seals | Hoffman et al. (2014) |
| RADseq | 129,494 loci; 7,431 SNPs | Common bottlenose dolphin adaptation to harmful algal blooms | Cammen et al. (2015) |
| Transcriptomics | 11,286 contigs | Sperm whale skin cell response to hexavalent chromium | Pabuwal et al. (2013) |
| Transcriptomics | 164,966 contigs | Physiological stress response in northern elephant seals | Khudyakov et al. (2015a; 2015b) |
| *Develop molecular resources* | | | |
| RADseq | 3,595 loci | Comparison of short-beaked common dolphin and harbor porpoise | Viricel et al. (2014) |
| Shotgun sequencing | 440,718 SNPs | SNP discovery in Northeast Atlantic common bottlenose dolphins | M. Louis (unpubl. data) |
| WGS | 144 SNPs | SNP validation in Antarctic fur seal | Humble et al. (2016) |
| Transcriptomics | 23,096 contigs; 144 SNPs | Gene and SNP discovery in Antarctic fur seal | Hoffman et al. (2011; 2012; 2013) |
| Transcriptomics & RADseq | 9,000 SNPs | Development of SNP array for polar bear and demonstration of utility in population genomics | Malenfant et al. (2015) |

Cammen_SupMat_TableS1 - Marine mammal genomics - *JHered*

References

Cammen KM, Schultz TF, Rosel PE, Wells RS, Read AJ. 2015. Genomewide investigation of adaptation to harmful algal blooms in common bottlenose dolphins (*Tursiops truncatus*). *Mol Ecol*. 24:4697-4710.

Foote AD, Liu Y, Thomas GWC, Vinař Ts, Alföldi J, Deng J, Dugan S, van Elk CE, Hunter ME, Joshi V, *et al.* 2015. Convergent evolution of the genomes of marine mammals. *Nat Genet*. 47:272-275.

Foote AD, Newton J, Ávila-Arcos MC, Kampmann M-L, Samaniego JA, Post K, Rosing-Asvid A, Sinding M-HS, Gilbert MTP. 2013. Tracking niche variation over millennial timescales in sympatric killer whale lineages. *Proc R Soc Lond B Biol Sci*. 280:20131481.

Foote AD, Vijay N, Ávila-Arcos M, Baird RW, Durban JW, Fumagalli M, Gibbs RA, Hanson MB, Korneliussen TS, Martin MD, *et al.* 2016. Genome-culture coevolution promotes rapid divergence of killer whale ecotypes. *Nat Commun*. 7:11693.

Hancock-Hanser BL, Frey A, Leslie MS, Dutton PH, Archer FI, Morin PA. 2013. Targeted multiplex next-generation sequencing: advances in techniques of mitochondrial and nuclear DNA sequencing for population genomics. *Mol Ecol Resour*. 13:254-268.

Hoffman JI. 2011. Gene discovery in the Antarctic fur seal (*Arctocephalus gazella*) skin transcriptome. *Mol Ecol Resour*. 11:703-710.

Hoffman JI, Simpson F, David P, Rijks JM, Kuiken T, Thorne MAS, Lacy RC, Dasmahapatra KK. 2014. High-throughput sequencing reveals inbreeding depression in a natural population. *Proc Natl Acad Sci USA*. 111:3775-3780.

Hoffman JI, Thorne MAS, Trathan PN, Forcada J. 2013. Transcriptome of the dead: characterisation of immune genes and marker development from necropsy samples in a free-ranging marine mammal. *BMC Genomics*. 14:52.

Hoffman JI, Tucker R, Bridgett SJ, Clark MS, Forcada J, Slate J. 2012. Rates of assay success and genotyping error when single nucleotide polymorphism genotyping in non-model organisms: a case study in the Antarctic fur seal. *Mol Ecol Resour*. 12:861-872.

Humble E, Martinez-Barrio A, Forcada J, Trathan PN, Thorne MAS, Hoffmann M, Wolf JBW, Hoffman JI. 2016. A draft fur seal genome provides insights into factors affecting SNP validation and how to mitigate them. *Mol Ecol Resour*.

Keane M, Semeiks J, Webb AE, Li YI, Quesada V, Craig T, Madsen LB, van Dam S, Brawand D, Marques PI, *et al.* 2015. Insights into the evolution of longevity from the bowhead whale genome. *Cell Reports*. 10:112-122.

Khudyakov JI, Champagne CD, Preeyanon L, Ortiz RM, Crocker DE. 2015a. Muscle transcriptome response to ACTH administration in a free-ranging marine mammal. *Physiol Genomics*. 47:318-330.

Khudyakov JI, Preeyanon L, Champagne CD, Ortiz RM, Crocker DE. 2015b. Transcriptome analysis of northern elephant seal (*Mirounga angustirostris*) muscle tissue provides a novel molecular resource and physiological insights. *BMC Genomics*. 16:64.

Kishida T, Thewissen JGM, Hayakawa T, Imai H, Agata K. 2015. Aquatic adaptation and the evolution of smell and taste in whales. *Zoolog Lett*. 1:9.

Liu S, Lorenzen ED, Fumagalli M, Li B, Harris K, Xiong Z, Zhou L, Korneliussen TS, Somel M, Babbitt C, *et al.* 2014. Population genomics reveal recent speciation and rapid evolutionary adaptation in polar bears. *Cell*. 157:785-794.

Cammen_SupMat_TableS1 - Marine mammal genomics - *JHered*

Malenfant RM, Coltman DW, Davis CS. 2015. Design of a 9K Illumina BeadChip for polar bears (*Ursus maritimus*) from RAD and transcriptome sequencing. *Mol Ecol Resour*. 15:587-600.

McGowen MR, Grossman LI, Wildman DE. 2012. Dolphin genome provides evidence for adaptive evolution of nervous system genes and a molecular rate slowdown. *Proc R Soc Lond B Biol Sci*. 279:3643-3651.

McGowen MR, Spaulding M, Gatesy J. 2009. Divergence date estimation and a comprehensive molecular tree of extant cetaceans. *Mol Phylogenet Evol*. 53:891-906.

Morin PA, Parsons KM, Archer FI, Ávila-Arcos M, Barrett-Lennard LG, Dalla Rosa L, Duchêne S, Durban JW, Ellis GM, Ferguson SH, *et al.* 2015. Geographic and temporal dynamics of a global radiation and diversification in the killer whale. *Mol Ecol*. 24:3964-3979.

Moura AE, Kenny JG, Chaudhuri R, Hughes MA, Welch AJ, Reisinger RR, de Bruyn PJN, Dahlheim ME, Hall N, Hoelzel AR. 2014. Population genomics of the killer whale indicates ecotype evolution in sympatry involving both selection and drift. *Mol Ecol*. 23:5179-5192.

Pabuwal V, Boswell M, Pasquali A, Wise SS, Kumar S, Shen Y, Garcia T, Lacerte C, Wise JP, Jr., Wise JP, Sr., *et al.* 2013. Transcriptomic analysis of cultured whale skin cells exposed to hexavalent chromium [Cr(VI)]. *Aquat Toxicol*. 134-135:74-81.

Ruan R, Guo A-H, Hao Y-J, Zheng J-S, Wang D. 2015. *De novo* assembly and characterization of narrow-ridged finless porpoise renal transcriptome and identification of candidate genes involved in osmoregulation. *Int J Mol Sci*. 16:2220-2238.

Seim I, Ma S, Zhou X, Gerashchenko MV, Lee S-G, Suydam R, George JC, Bickham JW, Gladyshev VN. 2014. The transcriptome of the bowhead whale *Balaena mysticetus* reveals adaptations of the longest-lived mammal. *Aging*. 6:879-899.

Shafer ABA, Gattepaille LM, Stewart REA, Wolf JBW. 2015. Demographic inferences using short-read genomic data in an approximate Bayesian computation framework: *in silico* evaluation of power, biases and proof of concept in Atlantic walrus. *Mol Ecol*. 24:328-345.

Springer MS, Signore AV, Paijmans JLA, Vélez-Juarbe J, Domning DP, Bauer CE, He K, Crerar L, Campos PF, Murphy WJ, *et al.* 2015. Interordinal gene capture, the phylogenetic position of Steller's sea cow based on molecular and morphological data, and the macroevolutionary history of Sirenia. *Mol Phylogenet Evol*. 91:178-193.

Viricel A, Pante E, Dabin W, Simon-Bouhet B. 2014. Applicability of RAD-tag genotyping for interfamilial comparisons: empirical data from two cetaceans. *Mol Ecol Resour*. 14:597-605.

Yim H-S, Cho YS, Guang X, Kang SG, Jeong J-Y, Cha S-S, Oh H-M, Lee J-H, Yang EC, Kwon KK, *et al.* 2014. Minke whale genome and aquatic adaptation in cetaceans. *Nat Genet*. 46:88-92.

Zhou X, Sun F, Xu S, Fan G, Zhu K, Liu X, Chen Y, Shi C, Yang Y, Huang Z, *et al.* 2013. Baiji genomes reveal low genetic variability and new insights into secondary aquatic adaptations. *Nat Commun*. 4:2708.

1 **Genomic methods take the plunge: recent advances in high-throughput sequencing of**
2 **marine mammals**
3

4 KRISTINA M. CAMMEN[1]*, KIMBERLY R. ANDREWS[2], EMMA L. CARROLL[3], ANDREW D.
5 FOOTE[4], EMILY HUMBLE[5,6], JANE I. KHUDYAKOV[7], MARIE LOUIS[3], MICHAEL R.
6 MCGOWEN[8], MORTEN TANGE OLSEN[9], AND AMY M. VAN CISE[10]
7

8 [1]School of Marine Sciences, University of Maine, Orono, Maine 04469, USA
9 [2]Department of Fish and Wildlife Sciences, University of Idaho, 875 Perimeter Drive MS 1136,
10 Moscow, Idaho 83844-1136, USA
11 [3]Scottish Oceans Institute, University of St Andrews, East Sands, St Andrews, Fife KY16 8LB,
12 UK
13 [4]Computational and Molecular Population Genetics ~~CMPG L~~lab, Institute of Ecology and
14 Evolution, University of Bern, Bern CH-3012, Switzerland
15 [5]Department of Animal Behaviour, University of Bielefeld, Postfach 100131, 33501 Bielefeld,
16 Germany
17 [6]British Antarctic Survey, High Cross, Madingley Road, Cambridge CB3 OET, UK
18 [7]Department of Biology, Sonoma State University, Rohnert Park, California 94928, USA
19 [8]School of Biological and Chemical Sciences, Queen Mary University of London,
20 Mile End Road, London E1 4NS, UK
21 [9]Evolutionary Genomics Section, Natural History Museum of Denmark, University of
22 Copenhagen, DK-1353 Copenhagen K, Denmark
23 [10]Scripps Institution of Oceanography, 8622 Kennel Way, La Jolla, California 92037, USA
24

25 *Corresponding author: kristina.cammen@maine.edu
26

27 Running title: Marine mammal genomics

**Abstract**

The dramatic increase in the application of genomic techniques to non-model organisms over the past decade has yielded numerous valuable contributions to evolutionary biology and ecology, many of which would not have been possible with traditional genetic markers. We review this recent progression with a particular focus on genomic studies of marine mammals, a group of taxa that represent key macroevolutionary transitions from terrestrial to marine environments and for which available genomic resources have recently undergone notable rapid growth. Genomic studies of non-model organisms utilize an expanding range of approaches, including ~~low- and high-coverage~~ whole genome sequencing, restriction site-associated DNA sequencing, array-based ~~high-throughput~~ sequencing of single nucleotide polymorphisms and target sequence probes (e.g., exomes), and transcriptome sequencing. These approaches generate different types and quantities of data, and many can be applied with limited or no prior genomic resources, thus overcoming <u>one traditional</u> limitation~~s~~ of research on non-model organisms. Within marine mammals, such studies have thus far yielded significant contributions to the fields of phylogenomics and comparative genomics, as well as enabled investigations of fitness, demography, and population structure ~~in natural populations~~. Here~~,~~ we review the primary options for generating genomic data, introduce several emerging techniques, and discuss the suitability of each approach for different applications in the study of non-model organisms.

**Keywords**: RADseq, SNP array, target sequence capture, whole genome sequencing, RNAseq, non-model organisms

49  **Introduction**

50  Recent advances in sequencing technologies, coincident with dramatic declines in cost, have

51  increasingly enabled the application of genomic sequencing in non-model systems (Ekblom and

52  Galindo 2011; Ellegren 2014). These advances in molecular technologies have in many ways

53  begun to blur the distinction between model and non-model organisms (Armengaud et al. 2014).

54  Non-model organisms (NMOs) have traditionally been defined as those for which whole-

55  organism experimental manipulation is rarely, if ever, possible due to logistical and/or ethical

56  constraints (Ankeny and Leonelli 2011). Further, NMOs have typically been characterized by

57  limited genomic resources, but this is becoming increasingly less so as the number of NMO

58  reference genomes grows rapidly, for example through efforts like the Genome 10K Project

59  (Koepfli et al. 2015). In fact, in some taxonomic orders, we are approaching the point at which

60  all species have at least one representative reference genome available for a closely related

61  species (Fig 1).

62

63  Despite the limitations of working with NMOs, including potentially small sample sizes, low

64  DNA quantity, and limited information on gene function, genetic and genomic investigations of

65  NMOs have yielded numerous valuable contributions to understanding their evolutionary

66  biology and ecology. For the past several decades, traditional genetic markers such as

67  microsatellites and short fragments of mitochondrial DNA (e.g., the control region) have been

68  extensively used in molecular ecology. These markers, which typically evolve under neutral

69  expectations, have proven useful for identifying population structure and reconstructing

70  population demographic history (Hedrick 2000). However, the power of such studies is limited

71  by the number of markers that can feasibly be evaluated using traditional approaches. The advent

72  of low-cost high-throughput sequencing has led to dramatic increases in the number of neutral

73  markers that can be evaluated, in many cases improving our power to resolve fine-scale or

74  cryptic population structure in species with high dispersal capability (e.g., Corander et al. 2013)

75  and improving the accuracy of estimating some (though not all) demographic parameters (Li and

76  Jakobsson 2012; Shafer et al. 2015). Importantly, high-throughput sequencing has also further

77  enabled genomic studies of non-neutral processes in NMOs, for example, characterizing both

78  deleterious and adaptive variation within and across species (Stinchcombe and Hoekstra 2008;

79    Künstner et al. 2010). It is increasingly evident that genomic analyses of NMOs can and have

80    provided important insights that could not be identified with traditional genetic markers.

81

82    Many molecular ecologists now face the challenge of deciding which of the broad range of

83    genomic approaches to apply to their study systems. Here we review the primary options for

84    generating genomic data and their relative suitability for different applications in the study of

85    NMOs. We focus on marine mammals, which represent several mammalian clades with notably

86    rapid growth in available genomic resources in recent years. This growth is clearly evident in

87    both publication rate (Fig 2) and the rise in number and size of genomic sequences deposited in

88    public resources (Fig 3). We comprehensively review the literature on marine mammal

89    genomics, highlighting recent trends in methodology and applications, and then describe in detail

90    the molecular approaches that are most commonly applied to studies of ~~non-model~~NMO

91    genomics. Our hope is that this review will highlight the promise of genomics for NMOs and

92    offer guidance to researchers considering the application of genomic techniques in their non-

93    model study system of choice.

94

95    **Why study marine mammal genomics?**

96    Marine mammals represent key macroevolutionary transitions from terrestrial to marine

97    environments (McGowen et al. 2014) and accordingly are an exemplary system for investigating

98    the evolution of several morphological and physiological adaptations (Foote et al. 2015)

99    associated with locomotion (Shen et al. 2012), sight (Meredith et al. 2013), echolocation (Parker

100   et al. 2013; Zou and Zhang 2015), deep diving (Mirceta et al. 2013), osmoregulation (Ruan et al.

101   2015), and cognition (McGowen et al. 2012). Furthermore, studies of marine mammal evolution

102   to date have characterized several unique aspects of their genome evolution that merit further

103   investigation, including low genomic diversity and a relatively slow molecular clock, especially

104   in cetaceans (Jackson et al. 2009; McGowen et al. 2012; Zhou et al. 2013). As many cetacean

105   species are highly mobile with no obvious physical geographic barriers to dispersal, they provide

106   a unique opportunity to study the role of behavior and culture in shaping population structure and

107   genetic diversity (Riesch et al. 2012; Carroll et al. 2015; Alexander et al. 2016). ~~Finally,~~ T~~t~~hough

108   highly mobile, many marine mammals exhibit evidence of local adaptation; for example, several

109   species show parallel divergent morphological and behavioral adaptations to coastal and pelagic

110   environments (Moura et al. 2013; Louis et al. 2014; Viricel and Rosel 2014). These species may

111   be studied across ocean basins as emerging examples of ecological adaptation and speciation

112   (Morin et al. 2010a).

113

114   Beyond their value as systems of evolutionary study, many marine mammals are also of broader

115   interest relating to their historical and present conservation status. Many marine mammal

116   populations share histories of dramatic decline due to hunting and other human impacts.

117   Genomics provides a promising tool with which to expand our insights into these historical

118   population changes, which so far primarily have been derived from archival review and

119   traditional genetic approaches (Ruegg et al. 2013; Sremba et al. 2015). More recently, since the

120   implementation of national and international protections, many marine mammal populations

121   have partially or fully recovered (Magera et al. 2013), yet the conservation status of certain

122   marine mammal populations remains of concern. Such vulnerable populations could benefit

123   greatly from an improved understanding of their genetic diversity and evolution, especially in

124   ways that can inform predictions of adaptive capacity to anthropogenic pressures and expand the

125   toolkit for conservation policy (Garner et al. 2016; Taylor and Gemmell 2016).

126

**Recent trends in marine mammal genomics**

128   We conducted a meta-analysis of the peer-reviewed marine mammal genomics literature to

129   evaluate trends in publication rates across research methodologies and aims. A search of the Web

130   of Science database using the term "genom*" and one of the following terms indicating study

131   species - "marine mammal", "pinniped", "seal", "sea lion", "sea otter", "whale", "dolphin",

132   "polar bear", "manatee" - identified 825 records on December 11, 2015. We excluded 77% of the

133   search results that were not directly related to genomic studies in marine mammal systems. The

134   remaining 101 articles that were relevant to marine mammal genomics were further categorized

135   by primary research methodology and general research aim. A subset of these articles is

136   described briefly in Supplemental Table 1.

137

138   From the early 1990s through 2015, published literature in the field ~~has~~ shifted from an early

139   focus on mitogenome sequencing to more sequence-intensive approaches, such as transcriptome

140   and whole genome sequencing (Figs 2 and 4). This trajectory closely follows trends in

141　sequencing technologies, from Sanger sequencing of short- and long-range PCR products for

142　mitogenome sequencing (Arnason et al. 1991) and SNP discovery (Olsen et al. 2011), to high-

143　throughput sequencing of reduced-representation genomic libraries (RRLs) that consist of

144　selected subsets of the genome (e.g., Viricel et al. 2014), to high-throughput sequencing of whole

145　genomes with varying levels of depth, ~~of~~ coverage, and contiguity. Today, high-throughput

146　sequencing can be used both to generate high-quality reference genome assemblies (Yim et al.

147　2014; Foote et al. 2015; Humble et al. 2016) and to re-sequence whole genomes at a population

148　scale (Liu et al. 2014a; Foote et al. 2016). Similarly, the scale of gene expression studies has

149　increased from quantitative real-time PCR of candidate genes (Tabuchi et al. 2006) to

150　microarrays containing hundreds to thousands of genes (Mancia et al. 2007) and high-throughput

151　RNAseq that evaluates hundreds of thousands of contigs across the genome (Khudyakov et al.

152　2015b). As the cost of high-throughput sequencing continues to decline, we anticipate an

153　increase in studies that sequence RRLs, whole genomes, and transcriptomes in NMOs at a

154　population scale.

155

156　Marine mammal genomic studies thus far have primarily contributed to the fields of

157　phylogenomics and comparative genomics (Fig 2, Table S1). Several of these comparative

158　genomics studies have aimed to improve our understanding of the mammalian transition to an

159　aquatic lifestyle and describe the evolutionary relationships within and among marine mammals

160　and their terrestrial relatives (McGowen et al. 2014; Foote et al. 2015). Whereas such studies

161　require only a single representative genome per species, an emerging class of studies applying

162　genomic techniques at a population scale enables further investigations of fitness, demography,

163　and population structure within ~~a~~ species (Table S1). However, expanding the scale of genomic

164　studies requires careful selection of an appropriate method for data generation and analysis, from

165　a growing number of approaches that are becoming available to non-model systems.

166

**Data generation**

167

168　Our review of marine mammal genomics highlights an increasing number of options for the

169　generation and analysis of genomic data. Choosing which of these sequencing strategies to apply

170　is a key step in any genomics study. Here, we describe approaches that have been used

171　successfully in order to help guide future studies of ecological, physiological, and evolutionary

172    genomics in NMOs. Across data generation methods, we highlight approaches that can be used

173    with limited or no prior genomic resources, overcoming one traditional challenge of genomic

174    studies of NMOs (the need for a reference genome to which sequencing reads can be mapped).

175    These methods produce a range in quantity and type of data output, from hundreds of SNPs to

176    whole genome sequences, and from single individuals to population samples, reflecting the

177    trade-off between number of samples and amount of data generated per sample.

178

179    <u>Sample collection, storage and extraction</u>

180    Prior to starting a genomic study, researchers must recognize that many recent methods for high-

181    throughput sequencing require genetic material of much higher quality and quantity than

182    techniques used to characterize traditional genetic markers. These more stringent sample

183    requirements necessitate new standards for tissue sampling, storage, and DNA/RNA extraction.

184    Ideally, samples should be collected from live or newly deceased individuals and stored at -80°C,

185    or when this is not possible at -20°C in RNAlater, Trizol, ethanol, salt-saturated DMSO, or dry,

186    depending on the intended application. Given the sensitivity of new sequencing methods, great

187    care should be taken to minimize cross-contamination during sampling, as even minute amounts

188    of genetic material from another individual can bias downstream analyses, for example variant

189    genotyping and gene expression profiles. Choice of extraction method varies with sample type

190    and study aim, but typically genomic methods require cleanup and treatment with RNase to yield

191    pure extracts, whereas RNAseq methods require rigorous DNase treatment to remove genomic

192    contamination that can bias expression results. Depending on the genomic methodology, target

193    quantities for a final sample may range from as low as 50 ng of DNA for some RRL sequencing

194    methods (Andrews et al. 2016) up to ~1 mg for sequencing the full set of libraries (of different

195    insert sizes) necessary for high-quality genome assemblies (Ekblom and Wolf 2014). Most

196    commercial RNAseq library preparation services require at least 500-1,000 ng of pure total RNA

197    that shows minimal degradation as measured by capillary gel electrophoresis (RNA Integrity

198    Number (RIN) $\geq$ 8). Samples should ideally consist of high molecular weight genetic material

199    (with little shearing), though continuing molecular advances enable genomic sequencing even of

200    low quantity or poor quality starting material. Extreme examples of the latter include

201    successfully sequenced whole genomes from ancient material (e.g., Rasmussen et al. 2010;

202  Meyer et al. 2012; Allentoft et al. 2015), including a more than 500,000-year-old horse (Orlando

203  et al. 2013).

204

205  <u>Reduced-representation genome sequencing</u>

206  *i. RADseq*

207  Reduced-representation sequencing methods evaluate only a small portion of the genome,

208  allowing researchers to sequence samples from a larger number of individuals within a given

209  budget in comparison to sequencing whole genomes. Restriction site-associated DNA

210  sequencing (RADseq) is currently the most widely -used RRL sequencing method for NMOs

211  (Davey et al. 2011; Narum et al. 2013; Andrews et al. 2016). RADseq generates sequence data

212  from short regions adjacent to restriction cut sites and therefore targets markers that are

213  distributed relatively randomly across the genome and occur primarily in non-coding regions.

214  This method allows simultaneous discovery and genotyping of thousands of genetic markers for

215  virtually any species, regardless of availability of prior genomic resources. Of greatest interest

216  are variable markers, characterized either as single SNPs or phased alleles that can be resolved

217  from the identification of several ~~SNPs~~ variants within a single locus.

218

219  The large number of markers generated by RADseq dramatically increases genomic resolution

220  and statistical power for addressing many ecological and evolutionary questions when compared

221  to studies using traditional markers (Table S1). For example, heterozygosity -fitness ~~associations~~

222  correlations in harbor seals (*Phoca vitulina*) were nearly fivefold higher when using 14,585

223  RADseq SNPs than when using 27 microsatellite loci (Hoffman et al. 2014). A recent study on

224  the Atlantic walrus (*Odobenus rosmarus rosmarus*) using 4,854 RADseq SNPs to model

225  demographic changes in connectivity and effective population size associated with the Last

226  Glacial Maximum (Shafer et al. 2015) both supported and extended inferences from previous

227  studies using traditional markers (Shafer et al. 2010; Shafer et al. 2014).

228

229  Furthermore, RADseq can provide sufficient numbers of markers across the genome to identify

230  genomic regions influenced by natural selection ~~in some cases~~. These analyses require large

231  numbers (thousands to tens of thousands) of markers to ensure that some markers will be in

232  linkage disequilibrium with genomic regions under selection and to minimize false positives,

233　particularly under non-equilibrium demographic scenarios (Narum and Hess 2011; De Mita et al.

234　2013; Lotterhos and Whitlock 2014). Extreme demographic shifts, as experienced by many

235　marine mammal populations (e.g., killer whales, Foote et al. 2016), can drive shifts in allele

236　frequencies that confound the distinction of drift and selection and make it difficult to detect

237　genomic signatures of selection (Poh et al. 2014). Proof of concept of the application of RADseq

238　for identifying genomic signatures of selection in wild populations was demonstrated in three-

239　spined sticklebacks (*Gasterosteus aculeatus*), for which analyses of over 45,000 SNPs

240　(Hohenlohe et al. 2010) identified genomic regions of known evolutionary importance associated

241　with differences between marine and freshwater forms (Colosimo et al. 2005; Barrett et al.

242　2008). RADseq studies with similar aims in marine mammals have resulted in comparatively

243　sparser sampling of SNPs (<10,000), likely due to both methodological differences and generally

244　low genetic diversity particularly among cetaceans. Nonetheless, genomic regions associated

245　with resistance to harmful algal blooms in common bottlenose dolphins (*Tursiops truncatus*)

246　were identified across multiple pairwise comparisons using 7,431 RADseq SNPs (Cammen et al.

247　2015), and genomic regions associated with habitat use and resource specialization in killer

248　whales (*Orcinus orca*) were identified using 3,281 RADseq SNPs (Moura et al. 2014a). Some of

249　these RADseq SNPs associated with diet in killer whales were later also confirmed as occurring

250　in genomic regions of high differentiation and reduced diversity consistent with a signature of

251　selection identified in a study utilizing ~~low-coverage~~ whole genome re-sequencing (Foote et al.

252　2016). It will remain important for further studies of genomic signatures of selection in NMOs to

253　carefully consider which approach~~es~~ will generate a sufficiently large number of SNPs to

254　accurately identify the range of putatively neutral $F_{ST}$ values (and thus outliers) given the

255　demographic history of the population (Lotterhos and Whitlock 2014).

256

257　Numerous laboratory methods have been developed for generating RADseq data (reviewed in

258　Andrews et al. 2016), with the most popular library preparation methods currently being the

259　original RAD (Miller et al. 2007; Baird et al. 2008), Genotyping by Sequencing (GBS, Elshire et

260　al. 2011; Poland et al. 2012), and double digest RAD (ddRAD, Peterson et al. 2012). All

261　RADseq methods share the common goal of sequencing regions adjacent to restriction cut sites

262　across the genome, but differ in technical details, such as the number and type of restriction

263　enzymes used, the mechanisms for reducing genomic DNA fragment sizes, and the strategies for

264    attaching sequencing adapters to the target DNA fragments. For example, both the original RAD

265    method and GBS use a single enzyme digest, but the original RAD ~~protocol~~ method uses a rare-

266    cutting enzyme and mechanical shearing to reduce DNA fragment size (Baird et al. 2008),

267    whereas  GBS uses a more frequent-cutting enzyme and relies on preferential PCR amplification

268    of shorter fragments for indirect size selection (Elshire et al. 2011). These ~~types of~~

269    ~~variation~~modifications lead to differences across methods in the time and cost of library

270    preparation, the number and lengths of loci produced, and the types of error and bias present in

271    the resulting data. Different RADseq methods will be better suited to different research

272    questions, study species, and research budgets, and therefore researchers embarking on a

273    RADseq study should carefully consider the suitability of each method for their individual

274    projects. Further details on the advantages and disadvantages of each method are described in

275    Andrews et al. (2016).

276

277    *ii. SNP arrays*

278    An alternative high-throughput reduced-representation genotyping approach involves the use of

279    custom arrays designed to capture and sequence targeted regions of the genome. Such array-

280    based approaches may provide certain advantages over RADseq, including the ability to easily

281    estimate genotyping error rates, scalability to thousands of samples, lower requirements for DNA

282    quantity/quality and technical effort, greater comparability of markers across studies, and the

283    ability to genotype SNPs within candidate genomic regions. However, unlike RADseq, array-

284    based techniques require prior knowledge of the study system's genome or the genome of a

285    closely related species, which remains unavailable for some NMOs. Furthermore, SNP arrays

286    must take into account the potential for ascertainment bias (e.g., Malenfant et al. 2015), whereas

287    RADseq avoids ascertainment bias by simultaneously discovering and genotyping markers.

288

289    To identify SNPs for NMO array development, researchers must rely on existing genomic

290    resources or generate new reference sequences, in the form of whole or reduced-representation

291    genomes or transcriptomes (Hoffman et al. 2012; Malenfant et al. 2015). When a whole genome

292    reference assembly is available for the target species or a related species, multiplex shotgun

293    sequencing can facilitate the rapid discovery of hundreds of thousands of SNPs for array

294    development. This SNP discovery approach involves high-throughput sequencing of sheared

295  genomic DNA, ~~which~~ that can be sequenced at a low depth of coverage (i.e., low mean read
296  depth across the genome) if suitable genotype likelihood-based methods (O'Rawe et al. 2015) are
297  used to identify polymorphic sites. Thus, this approach is less restrictive in terms of DNA
298  quality. For example, shotgun sequencing of 33 Northeast Atlantic common bottlenose dolphins,
299  which included degraded DNA collected from stranded specimens, on one Illumina HiSeq2000
300  lane of 100-bp single-end sequencing identified 440,718 high-quality SNPs (M. Louis
301  unpublished data). Such dense sampling of SNPs is essential for studies of population genomics
302  that require a large number of markers, such as for inferences of demographic history
303  (Gutenkunst et al. 2009; Excoffier et al. 2013; Liu and Fun 2015) and selective sweeps (Chen et
304  al. 2010). Once a set of putative markers has been identified, hybridization probes can be
305  designed from their flanking sequences and printed onto a SNP array. The two principal SNP
306  genotyping platforms supporting thousands to millions of SNPs are the Illumina Infinium
307  iSelect® and Affymetrix Axiom® arrays.

308

309  The use of SNP arrays in NMOs has thus far been somewhat limited, potentially due to low SNP
310  validation rates (Chancerel et al. 2011; Helyar et al. 2011), issues of ascertainment bias
311  (Albrechtsen et al. 2010; McTavish and Hillis 2015), and cost of SNP discovery. However, using
312  both SNP data and whole genome sequence from the Antarctic fur seal (*Arctocephalus gazella*),
313  Humble et al. (2016) recently demonstrated that careful filtering based on SNP genomic context
314  prior to array development has the potential to substantially increase assay success rates. Further,
315  ascertainment bias can be reduced by selecting samples for SNP discovery that span the
316  geographic range of populations that will be target-sequenced (Morin et al. 2004). By
317  accounting for ascertainment bias, Malenfant et al. (2015) were able to demonstrate population
318  structure in Canadian polar bears (*Ursus maritimus*) more clearly using a 9K SNP array than 24
319  microsatellite markers.

320

321  *iii. Target sequence capture*
322  Target sequence capture (TSC, also called target enrichment, direct selection, or Hyb-seq) has
323  many of the same advantages and disadvantages as the array-based SNP approaches described
324  above, but differs in library preparation, sequencing platform, and resulting sequence data. While
325  SNP arrays genotype single variable positions, TSC can be used to sequence selected short

326    fragments. With TSC, researchers can amplify and sequence up to a million target probes on

327    solid-state arrays, and even more if in-solution arrays are used. This gives the user the ability to

328    choose to sequence many samples in parallel (Cummings et al. 2010), as many as 100-150 per

329    Illumina HiSeq lane, or to sequence many regions per individual. Recent advances in target

330    enrichment, such as genotyping in thousands (Campbell et al. 2015), anchored hybrid enrichment

331    (Lemmon et al. 2012), and target capture of ultra-conserved elements (UCEs, Faircloth et al.

332    2012; McCormack et al. 2012), have further increased the number of regions and individuals that

333    can be sampled in a single lane. In addition, UCEs overcome the need for a reference genome,

334    enabling their wide application across many NMOs (though designing custom probe sets from

335    closely related species will remain preferable in many cases (Hancock-Hanser et al. 2013)).

336    Although a number of methodological variants have been developed and optimized (Bashiardes

337    et al. 2005; Noonan et al. 2006; Hodges et al. 2009; Cummings et al. 2010; Mamanova et al.

338    2010; Hancock-Hanser et al. 2013), TSC generally relies on hybridization and amplification of

339    specially prepared libraries consisting of fragmented genomic DNA. Many companies offer kits

340    for TSC, such as Agilent (SureSelect) and MYcroarray (MYbaits), with MYcroarray specifically

341    marketing their kits for use with NMOs.

342

343    The most common use of TSC has been the capture of whole exomes in model organisms,

344    including humans (Ng et al. 2009). However, as costs have plummeted, TSC is increasingly

345    being used in investigations of NMOs. TSC is particularly useful in sequencing ancient DNA,

346    where it can enrich the sample for endogenous DNA content relative to exogenous DNA (i.e.,

347    contamination) and thereby increase the relative DNA yield (Ávila-Arcos et al. 2011; Enk et al.

348    2014). For example, TSC has been used to generate mitogenome sequences from subfossil killer

349    whale specimens originating from the mid-Holocene, for comparison with modern lineages

350    (Foote et al. 2013). TSC was also recently utilized to compare >30 kb of exonic sequence from

351    museum specimens of the extinct Steller's sea cow (*Hydrodamalis gigas*) and a modern dugong

352    (*Dugong dugon*) specimen to investigate evolution within Sirenia (Springer et al. 2015). Springer

353    et al. (2016) further used TSC to examine gene evolution related to dentition across edentulous

354    mammals, including mysticetes. Finally, TSC of both exonic and intronic regions has been used

355    to assess genetic divergence across cetacean species (Hancock-Hanser et al. 2013; Morin et al.

356    2015). These studies show the potential use of TSC across evolutionary time-scales for

357 population genomics, phylogenomics, and studies of selection and gene loss across divergent

358 lineages (Table S1).

359

360 <u>Whole genome sequencing</u>

361 Beyond advances enabled by the reduced-representation methods presented above, our power

362 and resolution to elucidate evolutionary processes, including selection and demographic shifts,

363 can be further increased by sequencing whole genomes.

364

365 *i. ~~High-coverage~~ Rreference genome sequencing*

366 At the time of publication, there ~~exist~~ are 12 publicly available[1] whole (or near-whole) marine

367 mammal genomes of varying quality representing 10 families, including 7 cetaceans (Fig 1A), 3

368 pinnipeds (Fig 1B), the West Indian manatee (*Trichechus manatus*), and the polar bear. The first

369 sequenced marine mammal genome was that of the common bottlenose dolphin, which was

370 originally sequenced to ~2.5x depth of coverage using Sanger sequencing (Lindblad-Toh et al.

371 2011). This genome was later improved upon by adding both 454 and Illumina HiSeq data

372 (Foote et al. 2015). Other subsequent marine mammal genomes were produced solely using

373 Illumina sequencing and mate-paired or paired-end libraries with varied insert sizes (Miller et al.

374 2012; Zhou et al. 2013; Yim et al. 2014; Foote et al. 2015; Keane et al. 2015; Kishida et al. 2015;

375 Humble et al. 2016).

376

377 Whole genome sequencing has been used to address many issues in marine mammal genome

378 evolution, usually by comparison with other existing mammalian genomes. Biological insights

379 discussed in the genome papers listed above include the evolution of transposons and repeat

380 elements, gene evolution and positive selection, predicted population structure through time,

381 SNP validation, molecular clock rates, and convergent molecular evolution (Table S1). For

382 example, analyses of the Yangtze river dolphin (*Lipotes vexillifer*) genome confirmed that a

383 bottleneck occurred in this species during the last period of deglaciation (Zhou et al. 2013). In

384 addition, following upon earlier smaller-scale studies (e.g., Deméré et al. 2008; McGowen et al.

---

[1] These genomes are available on NCBI's online genome database or Dryad, but they have not all been published. As agreed upon in the Fort Lauderdale Convention, the community standard regarding such unpublished genomic resources is to respect the data generators' right to publish with these data first.

385  2008; Hayden et al. 2010), genomic analyses have confirmed the widespread decay of gene

386  families involved in olfaction, gustation, enamelogenesis, and hair growth in some cetaceans

387  (Yim et al. 2014; Kishida et al. 2015). Perhaps the most widespread use of whole genome studies

388  has been the use of models of selection to detect protein-coding genes that show evidence of

389  natural selection in specific lineages. A recent study by Foote et al. (2015) ~~has~~ extended this

390  approach to investigate convergent positive selection among cetaceans, pinnipeds, and sirenians.

391  This study exemplifies a trend in recent genomic studies, ~~which~~ that sequence multiple genomes

392  to address a predetermined evolutionary question, in this case, the molecular signature of aquatic

393  adaptation.

394

395  In addition to these evolutionary insights that typically stem from a comparative genomics

396  approach, the development of high--quality reference genome assemblies provide an important

397  resource that facilitates mapping of reduced-representation genomic data (see previous section)

398  as well as ~~relatively low-coverage,~~ short-read sequencing data with relatively low depth of

399  coverage (see following section). These data types can be generated at relatively low cost on

400  larger sample sizes enabling population--scale genomic studies. In many cases, genome

401  assemblies from closely related species are sufficient for use as a reference. Particularly among

402  marine mammals, given their generally slow rate of nucleotide divergence, it is therefore likely

403  unnecessary to sequence a high--quality reference genome assembly for every species. Instead,

404  resources could be allocated toward population--scale studies, including ~~low coverage~~ genome

405  re-sequencing efforts.

406

407  *ii. Population-level ~~low-coverage~~ genome re-sequencing*

408  In contrast to ~~high-coverage~~reference genome sequencing that today often exceeds 100x mean

409  read ~~coverage~~ depth and typically combines long- and short-insert libraries to generate high-

410  quality assemblies for one to a few individuals, ~~low-coverage~~ genome re-sequencing studies

411  ~~capitalize on existing reference assemblies and~~ aim to achieve only ≥2x ~~coverage~~ mean read

412  depth on tens to hundreds of individuals from short-insert libraries ~~which~~ ~~that are then~~whose

413  reads are anchored to ~~the~~ existing reference assembl~~ies~~y. ~~Given~~ Despite the inherent trade-offs

414  between cost, read depth, ~~coverage,~~ and sample size, ~~low-coverage~~ genome re-sequencing of

415  large numbers of individuals for population-level inference can be conducted at a relatively low

416    cost. In the past five years, several influential studies have used genome re-sequencing to

417    advance our understanding of the genomic underpinnings of different biological questions in

418    model systems. For example, population genomics of *Heliconius* butterflies highlighted the

419    exchange of genes between species that exhibit convergent wing patterns (The *Heliconius*

420    Genome Consortium 2012); whole genome re-sequencing of three-spined sticklebacks

421    highlighted the re-use of alleles in replicated divergences associated with ecological speciation

422    and local adaptation (Jones et al. 2012); and combined population genomics and phylogenomics

423    have identified regions of the genome associated with variation in beak shape and size in

424    Darwin's finches (Lamichhaney et al. 2015).

425

426    To date only two marine mammal population genomics studies using whole genome re-

427    sequencing have been published. These studies involved re-sequencing the genomes of 79

428    individuals from three populations of polar bears (Liu et al. 2014a) and 48 individuals from five

429    evolutionarily divergent ecotypes of killer whale (Foote et al. 2016). The findings of Foote et al.

430    (2016) confirmed results of population differentiation that had previously been established using

431    traditional genetic markers (Morin et al. 2010a). However, the study also provided new insights

432    into the demographic history, patterns of selection associated with ecological niche, and evidence

433    of episodic ancestral admixture that could not have been obtained using traditional markers.

434

435    Several new resources have made such population genomic studies economically possible for a

436    greater number of NMOs, including the availability of a reference genome assembliesy (see

437    section above), relatively low-cost high-throughput sequencing (further increases in throughput

438    expected with the new Illumina HiSeq X Ten (van Dijk et al. 2014)), and crucially, the

439    development of likelihood-based methods that allow estimation of population genetic metrics

440    from low-coverage re-sequencing data (Fumagalli et al. 2013; O'Rawe et al. 2015). One last

441    consideration is the ease of laboratory methods necessary to generate whole genome re-

442    sequencing data when compared to other methods such as RADseq or TSC. DNA simply needs

443    to be extracted from the samples and, using proprietary kits, built into individually index-

444    amplified libraries using proprietary kits, which that are then equimolarly pooled and submitted

445    for sequencing.

446

447 Many population genomic analyses are based on the coalescent model that gains most

448 information from the number of independent genetic markers, not the number of individuals

449 sampled. Sample sizes of ~10 individuals are usually considered sufficient (Robinson et al.

450 2014) and have been standard in many genome-wide studies in the eco-evolutionary sciences

451 (Ellegren et al. 2012; Jones et al. 2012). Thus, sampling fewer individuals ~~at lower coverage but~~

452 ~~for orders of magnitude more data~~by whole genome re-sequencing is a salient approach~~, which~~

453 that allows us to consider many more gene trees, whilst continuing to provide robust estimates of

454 per-site genetic metrics (e.g., $F_{ST}$). The robustness of inference from ~~low-coverage~~ data with low

455 mean read depth across the genome was recently confirmed using a comparison of per-site $F_{ST}$

456 estimates for the same sites from high~~-coverage~~depth (≥20x) RADseq data and low~~-~~

457 ~~coverage~~depth (≈2x) whole genome re-sequencing data in pairwise comparisons between the

458 same two killer whale ecotypes (Foote et al. 2016).

459

460 Beyond the increased power afforded by sequencing more polymorphic sites, whole genome re-

461 sequencing also allows inference of demographic history from the genome of even just a single

462 individual by identifying Identical By Descent (IBD) segments and runs of homozygosity (Li

463 and Durbin 2011; Harris and Nielsen 2013). For example, Liu et al. (2014a) found evidence for

464 ongoing gene flow from polar bears into brown bears after the two species initially diverged.

465 Genome re-sequencing of sufficient numbers of individuals also facilitates haplotype phasing,

466 which has many applications, including the detection of ongoing selective sweeps (Ferrer-

467 Admetlla et al. 2014) and the inference of demographic history of multiple populations based on

468 coalescence of pairs of haplotypes in different individuals (Schiffels and Durbin 2014).

469 However, haplotype phasing ~~has~~ typically require~~s~~d genomic ~~higher coverage~~ data with higher

470 mean read depth (~20x) from tens of individuals (though recent advances in genotype imputation

471 suggest success with ~~lower coverage~~ data of lower mean read depth (VanRaden et al. 2015)).

472 Thus far, phasing has been restricted to relatively few NMO studies, and no marine mammal

473 studies to the best of our knowledge.

474

475 Transcriptome sequencing

476 In comparison with the DNA-based genomic approaches described above, RNA-based genomic

477 approaches are a relatively new and emerging application in NMOs such as marine mammals.

478    Transcriptomics by RNA sequencing (RNAseq) can rapidly generate vast amounts of

479    information regarding genes and gene expression without any prior genomic resources. This

480    approach can resolve differences in global gene expression patterns between populations,

481    individuals, tissues, cells, and physiological or environmental conditions, and can yield insights

482    into the molecular basis of environmental adaptation and speciation in wild animals (Wolf 2013;

483    Alvarez et al. 2015). Furthermore, RNAseq is a valuable tool for resource development, for

484    example as a precursor to designing SNP and TSC arrays (e.g., Hoffman et al. 2012). However,

485    applying RNAseq to NMOs requires several unique considerations in comparison to the DNA-

486    based methods described above. Most importantly, the labile nature of gene transcription and

487    high detection sensitivity of RNAseq have the potential to amplify transcriptional "noise" and

488    are thus extremely sensitive to experimental design.

489

490    If the experimental goal is to capture a comprehensive transcriptome profile for a study

491    organism, multiple tissues from individuals of varied life history stages should be sampled.

492    However, if the aim is to characterize transcriptional responses to physiological or environmental

493    stimuli, efforts should focus on minimizing variability in individuals and sampling conditions

494    (Wolf 2013). For differential expression analyses, pairwise comparisons should be made within

495    the same individual if at all possible (e.g., before and after treatment, between two

496    developmental stages). As RNAseq only captures a 'snapshot' of gene expression in time,

497    repeated sampling or time-course studies are necessary to obtain a more complete picture of

498    cellular responses to the condition(s) in question (Spies and Ciaudo 2015). Sampling and

499    sequencing depth requirements will depend on the study design. Simulation studies have shown

500    that a minimum of 5-6 biological replicates sequenced at a depth of 10-20 million reads per

501    sample is necessary for differential expression analysis (Liu et al. 2014b; Schurch et al. 2015).

502    RNAseq can also be used for biomarker development to expand molecular toolkits for NMOs

503    without sequenced genomes (Hoffman et al. 2013). In this case, higher sequencing depths of 30-

504    60 million reads per sample are recommended for SNP discovery and genotyping (De Wit et al.

505    2015).

506

507    Following sequence generation, transcript annotation remains a challenge for NMOs without

508    reference transcriptomes or genomes. *De novo* transcriptomes can be annotated through detection

509    of assembled orthologs of highly conserved proteins, but these analyses remain limited by the

510    quality of reference databases. As a result, NMO transcriptomes are biased in favor of highly

511    conserved terrestrial mammal genes and therefore provide an incomplete understanding of

512    animal adaptations to natural environments (Evans 2015). For example, while 70.0% of northern

513    elephant seal (*Mirounga angustirostris*) skeletal muscle transcripts had BLASTx hits to mouse

514    genes, only 54.1% of blubber transcripts could be annotated due to poor representation of this

515    tissue in terrestrial mammal reference proteomes (Khudyakov et al. 2015b).

516

517    To date, RNAseq has been used for gene discovery and phylogenomics analyses in Antarctic fur

518    seal (Hoffman 2011; Hoffman et al. 2013), polar bear (Miller et al. 2012), Indo-Pacific

519    humpback dolphin (*Sousa chinensis* (Gui et al. 2013)), spotted seal (*Phoca largha* (Gao et al.

520    2013)), bowhead whale (*Balaena mysticetus* (Seim et al. 2014)), narrow-ridged finless porpoise

521    (*Neophocaena asiaeorientalis* (Ruan et al. 2015)), and humpback whale (*Megaptera*

522    *novaeangliae* (Tsagkogeorga et al. 2015)) (Table S1). Due to the challenges of repeated

523    sampling of wild marine mammals, few studies have examined cetacean or pinniped

524    transcriptome responses to environmental or experimental stimuli. The majority of such

525    functional gene expression studies have used microarrays (Mancia et al. 2008; Mancia et al.

526    2012; Mancia et al. 2015); however, RNAseq has been employed to profile sperm whale

527    (*Physeter macrocephalus*) skin cell response to hexavalent chromium (Pabuwal et al. 2013) and

528    free-ranging northern elephant seal skeletal muscle response to an acute stress challenge

529    (Khudyakov et al. 2015a; Khudyakov et al. 2015b). With decreasing sequencing costs and

530    improvements in bioinformatics tools, RNAseq has the potential to accelerate molecular

531    discoveries in marine mammal study systems and supplement existing functional genomics

532    approaches.

533

534    Emerging techniques

535    In addition to the relatively proven NMO genomic data generation techniques described above, a

536    suite of emerging techniques is entering the field, with exciting promise for exploration of

537    existing and new research areas. For example, high-throughput shotgun sequencing is

538    increasingly being used to identify genetic material from multiple species in a single sample

539    (metagenomics and metatranscriptomics), rather than focus on characterizing variation in a

540    single target individual. These multi-species approaches can be used, for example, to

541    characterize diet from fecal samples (Deagle et al. 2009) and to investigate microbiomes (Nelson

542    et al. 2015), objectives with implications for improving our understanding of both basic ecology

543    and health in natural populations of NMOs. Furthermore, high-throughput sequencing of

544    environmental DNA dramatically increases the throughput of NMO detection in environmental

545    (e.g., seawater) samples (Thomsen et al. 2012), using degenerate primers for multi-species

546    detection rather than requiring the design and implementation of numerous single-species

547    protocols (Foote et al. 2012).

548

549    A second broad area of emerging interest moves beyond the study of variation at the DNA and

550    RNA levels to examine epigenetic effects of histone modification on gene regulation and

551    evolution. Epigenomic studies often examine changes in DNA methylation in association with

552    processes such as cancer and ageing. Such approaches, from targeted gene to genome-wide, have

553    only very recently and not yet frequently been applied in NMOs. Polanowski et al. (2014) used a

554    targeted gene approach to examine changes in DNA methylation in age-associated genes,

555    previously identified in humans and mice, in humpback whales of known age. The most

556    informative markers were able to estimate humpback whale ages with standard deviations of

557    approximately 3-5 years, demonstrating the potential transferability of these approaches from

558    model to non-model organism. Villar et al. (2015) utilized a genome-wide approach – chromatin

559    immunoprecipitation followed by high-throughput sequencing (ChIPseq) – to examine gene -

560    regulatory element evolution across mammals, including four species of cetaceans. This study

561    identified highly conserved gene -regulatory elements based on their histone modifications

562    (H3K27ac and H3K4me3), showed that recently evolved enhancers were associated with genes

563    under positive selection in marine mammals, and identified unique *Delphinus*-specific enhancers.

564    Finally, reduced--representation epigenomic approaches have also been developed (Gu et al.

565    2011), and although they have not yet been used in marine mammals to our knowledge, these

566    techniques could facilitate future studies of how changes in DNA methylation patterns affect

567    other biological processes, such as stress levels or pregnancy.

568

569    **Data analysis**

570    Following the generation of genomic data, researchers must select the most appropriate genomic

571    analysis (i.e., bioinformatics) pipelines, which often differ significantly from those used in

572    traditional genetic studies of NMOs. The choice of analysis pipeline will depend on multiple

573    factors including the availability of a reference genome, the level of diversity within the dataset

574    (e.g., single- or multi-ple species), the type of data generated (e.g., single end vs. or paired-

575    end), and the computing resources available. The computational needs, both in terms of hardware

576    and competency in computer science, for analysis of genomic data typically far exceed those

577    necessary for traditional genetic markers. On the smaller end of the spectrum, one lane of 50 bp

578    single-end sequencing on an Illumina HiSeq 2500 can produce tens of gigabytes of data, while

579    data files associated with a single high-coverage quality vertebrate genome may reach hundreds

580    of gigabytes in size (Ekblom and Wolf 2014). Computing resources necessary for the analysis of

581    these genomic datasets can range from ~10 gigabytes for a pilot study using a reduced-

582    representation sequencing approach to over a terabyte for whole -genome sequence assembly

583    (Ekblom and Wolf 2014). Fortunately, university computing clusters, cloud-based (Stein 2010)

584    and high-performance computing clusters (e.g., XSEDE; Towns et al. 2014), and open web-

585    based platforms for genomic research (e.g., Galaxy; Goecks et al. 2010) are becoming

586    increasingly accessible. Furthermore, new pipelines are continuously being developed and

587    improved, and there are a growing number of resources aimed at training molecular ecologists

588    and evolutionary biologists in computational large-scale data analysis (Andrews and Luikart

589    2014; Belcaid and Toonen 2015; Benestan et al. 2016). We provide a limited an indicative list of

590    the current, most commonly used analysis pipelines that are specific to each data generation

591    method in Supplemental Table 12. Here, we briefly summarize current genomic data analysis

592    pipelines and discuss considerations that are likely to be similar across multiple data generation

593    methods.

594

595    Genomic data analysis often involves multiple steps, and the choice of analysis tool for each step

596    can greatly affect the outcome, with different tools producing different (though usually

597    overlapping) sets of results (e.g., Schurch et al. 2015). All analyses begin by evaluating data

598    quality, trimming sequences if necessary to remove erroneous nucleotides (MacManes 2014),

599    and implementing appropriate data quality filters (e.g., phred scores, read length, and/or read

600    depth). Raw reads also need to be demultiplexed based on unique barcodes if pools of

601  individuals were sequenced in a single lane. Analyses then usually proceed in a *de novo* or

602  genome-enabled manner, depending on available resources. Briefly, sequences can be compared

603  (e.g., to identify variants) by mapping all reads to a reference genome or *de novo* assembling

604  stacks of sequences putatively derived from the same locus, based on sequence similarity. *De*

605  *novo* methods are sensitive to sequencing error, as well as true genetic variation, and therefore

606  can erroneously assemble polymorphic sequences as separate loci or transcripts, requiring further

607  filtering to remove redundancy. The opposite problem can also occur in both *de novo* and

608  reference mapping approaches, where two distinct loci (e.g., paralogous loci) may assemble as a

609  single locus or map to the same reference location. Researchers should therefore recognize the

610  inherent trade-offs when carefully selecting their thresholds for acceptable levels of variation

611  within and among loci.

612

613  Considerations relevant to the selection of subsequent downstream analyses are specific to the

614  type of data generated and the research objective. For example, RADseq analysis pipelines differ

615  in the algorithms used to genotype variants (Table 1S2). Similarly, there are several gene

616  expression analysis pipelines for RNAseq data that compare transcript abundance between

617  samples (Table 1S2). Analysis of TSC data usually uses standard *de novo* assemblers (e.g.,

618  Trinity, Velvet); these assemblers can be run using packages such as PHYLUCE (Faircloth

619  2015), which is designed specifically for use with ultraconserved elements. Unfortunately, for

620  most analyses, there are no unifying recommendations currently available and researchers must

621  evaluate several approaches, each with their own advantages and disadvantages, in order to

622  select the most appropriate tool for their particular experiment and system. Furthermore, we can

623  expect that the recommendations for analysis tools will continue to evolve as new programs

624  become available in the future.

625

626  Guidelines for data quality control and sharing

627  With rapid growth in sequencing platforms and bioinformatics analysis pipelines comes the need

628  to extend existing principles (e.g., Bonin et al. 2004) on quality control, analysis, and

629  transparency. General recommendations for sample and data handling, library preparation, and

630  sequencing have been discussed elsewhere (Paszkiewicz et al. 2014). We therefore focus on the

631  need to produce guidelines on data quality evaluation and reporting for genomic data (e.g.,

632    Morin et al. 2010b). A primary challenge in this area is that quality metrics vary widely across

633    sequencing technologies. Yet, regardless of sequencing platform, the quality of sequencing reads

634    must be evaluated (e.g., using FastQC; Andrews 2010) and reported.

635

636    Best practices guidelines for ~~high-coverage whole~~reference genome sequencing and RNAseq

637    data generation, analysis, and reporting are available from the human-centric ENCODE

638    consortium (www.encodeproject.org). These include minimum depth of sequencing and number

639    and reproducibility of biological replicates. For RNAseq experiments, evaluation of *de novo*

640    assembly quality remains a challenge. Suggested quality metrics include percentage of raw reads

641    mapping back to the assembly and number of assembled transcripts with homology to known

642    proteins (MacManes 2016). Emerging tools such as Transrate (Smith-Unna et al. 2015) attempt

643    to integrate these and other metrics into a comprehensive assembly quality score.

644

645    In contrast, there is not yet any standard way to estimate or report error rates with RADseq or

646    ~~low-coverage~~ genome re-sequencing methods (but see Mastretta-Yanes et al. 2015; Fountain et

647    al. 2016). Recommendations to improve confidence in genotyping include using methods that

648    account for population--level allele frequencies when calling individual genotypes, mapping

649    reads to reference genomes rather than *de novo* assembly (Nadeau et al. 2014; Fountain et al.

650    2016), filtering out PCR duplicates (Andrews et al. 2014), identifying and removing markers in

651    possible repeat regions, and filtering data to include only those with high read depth (>10-20x

652    per locus per individual) (Nielsen et al. 2011). Other analysis methods, such as robust Bayesian

653    methods and likelihood--based approaches that account for read quality in calculations of

654    posterior probabilities of genotypes and per-site allele frequencies utilizing the sample mean site

655    frequency spectrum as a prior (Fumagalli et al. 2013), can account for uncertainty and/or error in

656    the data, and are therefore suitable for use with low to moderate read depths (2-20x per locus;

657    e.g., Han et al. 2015; O'Rawe et al. 2015).

658

659    Due to the large number of analysis tools that are available, data quality and reproducibility

660    ultimately depend on methods and data transparency. All raw sequencing reads should be

661    publicly archived, for example deposited in the NCBI Sequence Read Archive. Many journals,

662    including the *Journal of Heredity* (Baker 2013), now also require that primary data supporting

663   the published results and conclusions (e.g., SNP genotypes, assemblies) be publicly archived in

664   online data repositories (e.g., Dryad). We further recommend making public the analysis

665   pipelines, scripts (e.g., using GitHub), and additional outputs, as appropriate, in order for

666   analyses to be fully reproducible and transparent, which is the cornerstone of the scientific

667   method (Nosek et al. 2015).

668

669   **Future directions**

670   As demonstrated here for one group of mammalian taxa, the rapid growth of the field of non-

671   model genomics has been both impressive and empowering. As we approach a point of relative

672   saturation in reference genomes, we anticipate an increase in population-scale genomic studies

673   that produce lower depth or coverage datasets per individual but across larger sample sizes

674   relative to high-coverage sequencing of a few individuals of each species. In addition (or

675   alternatively), we hope to see increasing efforts to sequence reference transcriptomes and

676   improve NMO genome annotation in ways beyond the inherently limited approach of

677   comparison to gene lists from a few model organisms. Population-scale genomic studies will

678   facilitate greater ecological understanding of natural populations, while efforts to improve

679   annotation will address persistent limitations in our understanding of gene function for NMOs.

680   Ultimately, improving our understanding of local adaptation, adaptive potential, and

681   demographic history through the use of genomic toolkits such as those described here is likely to

682   have important implications for the future conservation of these populations.

683

684   Advances in sequencing technologies and analytical tools will no doubt continue, in some cases

685   drawing on established techniques in model organisms, posing both new opportunities and new

686   challenges for researchers in NMO genomics. Likely the most persistent challenge will remain

687   selecting the data generation and experimental design that is most appropriate for the respective

688   research objective. Our review identified few cases that exhibit relative dominance of a single

689   methodology and analytical pipeline (e.g., RADseq and STACKS, RNAseq and Trinity); rather,

690   more often we found a diversity of approaches even within each category of data generation. In

691   fact, such diversity of approaches has its benefits, with each approach promoting its own

692   advantages (and limitations). Overall, our reflections on lessons learned from the past decade of

693   NMO genomics in one well-studied group of mammalian taxa clearly demonstrate the value,

694 increased ease, and future promise of applying genomic techniques across a wide range of non-

695 model species to gain previously unavailable insights into evolution, population biology, and

696 physiology on a genome-wide scale.

697

**Acknowledgements**

706

**Funding**

719

**References**

721 Albrechtsen A, Nielsen FC, Nielsen R. 2010. Ascertainment biases in SNP chips affect measures
722         of population divergence. *Mol Biol Evol*. 27:2534-2547.
723 Alexander A, Steel D, Hoekzema K, Mesnick S, Engelhaupt D, Kerr I, Payne R, Baker CS.
724         2016. What influences the worldwide genetic structure of sperm whales (*Physeter*
725         *macrocephalus*)? *Mol Ecol*.

726    Allentoft ME, Sikora M, Sjögren K-G, Rasmussen S, Rasmussen M, Stenderup J, Damgaard PB,
727            Schroeder H, Ahlstrom T, Vinner L*, et al.* 2015. Population genomics of Bronze Age
728            Eurasia. *Nature*. 522:167-172.
729    Alvarez M, Schrey AW, Richards CL. 2015. Ten years of transcriptomics in wild populations:
730            what have we learned about their ecology and evolution? *Mol Ecol*. 24:710-725.
731    Anders S, Huber W. 2010. Differential expression analysis for sequence count data. *Genome
732            Biol*. 11:R106.
733    Andrews K, Good JM, Miller MR, Luikart G, Hohenlohe PA. 2016. Harnessing the power of
734            RADseq for ecological and evolutionary genomics. *Nat Rev Genet*. 17:81-92.
735    Andrews KR, Hohenlohe PA, Miller MR, Hand BK, Seeb JE, Luikart G. 2014. Trade-offs and
736            utility of alternative RADseq methods: Reply to Puritz *et al*. 2014. *Mol Ecol*. 23:5943-
737            5946.
738    Andrews KR, Luikart G. 2014. Recent novel approaches for population genomics data analysis.
739            *Mol Ecol*. 23:1661-1667.
740    Andrews S. 2010. FastQC: a quality control tool for high throughput sequence data. Available
741            online at: http://www.bioinformatics.babraham.ac.uk/projects/fastqc
742    Ankeny RA, Leonelli S. 2011. What's so special about model organisms? *Studies in History and
743            Philosophy of Science*. 42:313-323.
744    Armengaud J, Trapp J, Pible O, Geffard O, Chaumot A, Hartmann EM. 2014. Non-model
745            organisms, a species endangered by proteogenomics. *J Proteomics*. 105:5-18.
746    Arnason U, Adegoke JA, Bodin K, Born EW, Esa YB, Gullberg A, Nilsson M, Short RV, Xu X,
747            Janke A. 2002. Mammalian mitogenomic relationships and the root of the eutherian tree.
748            *Proc Natl Acad Sci USA*. 99:8151-8156.
749    Arnason U, Gullberg A, Widegren B. 1991. The complete nucleotide sequence of the
750            mitochondrial DNA of the fin whale, *Balaenoptera physalus*. *J Mol Evol*. 33:556-568.
751    Ávila-Arcos M, Cappellini E, Romero-Navarro JA, Wales N, Moreno-Mayar JV, Rasmussen M,
752            Fordyce SL, Montiel R, Vielle-Calzada J-P, Willerslev E*, et al.* 2011. Application and
753            comparison of large-scale solution-based DNA capture-enrichment methods on ancient
754            DNA. *Sci Rep*. 1:74.
755    Baird NA, Etter PD, Atwood TS, Currey MC, Shiver AL, Lewis ZA, Selker EU, Cresko WA,
756            Johnson EA. 2008. Rapid SNP discovery and genetic mapping using sequenced RAD
757            markers. *PLoS One*. 3:e3376.
758    Baker CS. 2013. *Journal of Heredity* adopts Joint Data Archiving Policy. *J Hered*. 104:1.
759    Barrett RDH, Rogers SM, Schluter D. 2008. Natural selection on a major armor gene in
760            threespine stickleback. *Science*. 322:255-257.
761    Bashiardes S, Veile R, Helms C, Mardis ER, Bowcock AM, Lovett M. 2005. Direct genomic
762            selection. *Nat Methods*. 2:63-69.
763    Belcaid M, Toonen RJ. 2015. Demystifying computer science for molecular ecologists. *Mol
764            Ecol*. 24:2619-2640.
765    Benestan LM, Ferchaud A-L, Hohenlohe PA, Garner BA, Naylor GJP, Baums IB, Schwartz MK,
766            Kelley JL, Luikart G. 2016. Conservation genomics of natural and managed populations:
767            building a conceptual and practical framework. *Mol Ecol*.
768    Bolger AM, Lohse M, Usadel B. 2014. Trimmomatic: a flexible trimmer for Illumina Sequence
769            Data. *Bioinformatics*. 30:2114-2120.

770  Bonin A, Bellemain E, Bronken Eidesen P, Pompanon F, Brochmann C, Taberlet P. 2004. How
771          to track and assess genotyping errors in population genetics studies. *Mol Ecol*. 13:3261-
772          3273.
773  Brown CT, Howe A, Zhang Q, Pyrkosz AB, Brom TH. 2012. A reference-free algorithm for
774          computational normalization of shotgun sequencing data. *arXive*. 1203:4802.
775  Cammen KM, Schultz TF, Rosel PE, Wells RS, Read AJ. 2015. Genomewide investigation of
776          adaptation to harmful algal blooms in common bottlenose dolphins (*Tursiops truncatus*).
777          *Mol Ecol*. 24:4697-4710.
778  Campbell NR, Harmon SA, Narum SR. 2015. Genotyping-in-Thousands by sequencing (GT-
779          seq): a cost effective SNP genotyping method based on custom amplicon sequencing.
780          *Mol Ecol Resour*. 15:855-867.
781  Carroll EL, Baker CS, Watson M, Alderman R, Bannister J, Gaggiotti OE, Gröcke DR,
782          Patenaude N, Harcourt R. 2015. Cultural traditions across a migratory network shape the
783          genetic structure of southern right whales around Australia and New Zealand. *Sci Rep*.
784          5:16182.
785  Catchen JM, Amores A, Hohenlohe PA, Cresko WA, Postlethwait JH. 2011. *Stacks*: building
786          and genotyping loci *de novo* from short-read sequences. *G3*. 1:171-182.
787  Catchen JM, Hohenlohe PA, Bassham S, Amores A, Cresko WA. 2013. Stacks: an analysis tool
788          set for population genomics. *Mol Ecol*. 22:3124-2140.
789  Chancerel E, Lepoittevin C, Le Provost G, Lin Y-C, Jaramillo-Correa JP, Eckert AJ, Wegrzyn
790          JL, Zelenika D, Boland A, Frigerio J-M, *et al*. 2011. Development and implementation of
791          a highly-multiplexed SNP array for genetic mapping in maritime pine and comparative
792          mapping with loblolly pine. *BMC Genomics*. 12:368.
793  Chen H, Patterson N, Reich D. 2010. Population differentiation as a test for selective sweeps.
794          *Genome Res*. 20:393-402.
795  Colosimo PF, Hosemann KE, Balabhadra S, Villarreal G, Grimwood J, Schmutz J, Myers RM,
796          Schluter D, Kingsley DM. 2005. Widespread parallel evolution in sticklebacks by
797          repeated fixation of ectodysplasin alleles. *Science*. 307:1928-1933.
798  Conesa A, Gotz S, Garcia-Gomez JM, Terol J, Talon M, Robles M. 2005. Blast2GO: a universal
799          tool for annotation, visualization and analysis in functional genomics research.
800          *Bioinformatics*. 21:3674-3676.
801  Corander J, Majander KK, Cheng L, Merilä J. 2013. High degree of cryptic population
802          differentiation in the Baltic Sea herring *Clupea harengus*. *Mol Ecol*. 22:2931-2940.
803  Cummings N, King R, Rickers A, Kaspi A, Lunke S, Haviv I, Jowett JBM. 2010. Combining
804          target enrichment with barcode multiplexing for high throughput SNP discovery. *BMC
805          Genomics*. 11:641.
806  Davey JW, Hohenlohe PA, Etter PD, Boone JQ, Catchen JM, Blaxter ML. 2011. Genome-wide
807          genetic marker discovery and genotyping using next-generation sequencing. *Nat Rev
808          Genet*. 12:499-510.
809  De Mita S, Thuillet A-C, Gay L, Ahmadi N, Manel S, Ronfort J, Vigouroux Y. 2013. Detecting
810          selection along environmental gradients: analysis of eight methods and their effectiveness
811          for outbreeding and selfing populations. *Mol Ecol*. 22:1383-1399.
812  De Wit P, Pespeni MH, Palumbi SR. 2015. SNP genotyping and population genomics from
813          expressed sequences - current advances and future possibilities. *Mol Ecol*. 24:2310-2323.
814  Deagle BE, Kirkwood R, Jarman SN. 2009. Analysis of Australian fur seal diet by
815          pyrosequencing prey DNA in faeces. *Mol Ecol*. 18:2022-2038.

26

816  Deméré TA, McGowen MR, Berta A, Gatesy J. 2008. Morphological and molecular evidence for
817      a stepwise evolutionary transition from teeth to baleen in mysticete whales. *Syst Biol*.
818      57:15-37.
819  DePristo MA, Banks E, Poplin R, Garimella KV, Maguire JR, Hartl C, Philippakis AA, del
820      Angel G, Rivas MA, Hanna M*, et al.* 2011. A framework for variation discovery and
821      genotyping using next-generation DNA sequencing data. *Nat Genet*. 43:491-498.
822  Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, Batut P, Chaisson M, Gingeras
823      TR. 2013. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics*. 29:15-21.
824  Eaton DAR. 2014. PyRAD: assembly of *de novo* RADseq loci for phylogenetic analysis.
825      *Bioinformatics*. 30:1844-1849.
826  Ekblom R, Galindo J. 2011. Applications of next generation sequencing in molecular ecology of
827      non-model organisms. *Heredity*. 107:1-15.
828  Ekblom R, Wolf JBW. 2014. A field guide to whole-genome sequencing, assembly and
829      annotation. *Evolutionary Applications*. 7:1026-1042.
830  Ellegren H. 2014. Genome sequencing and population genomics in non-model organisms.
831      *Trends Ecol Evol*. 29:51-63.
832  Ellegren H, Smeds L, Burri R, Olason PI, Backström N, Kawakami T, Künstner A, Mäkinen H,
833      Nadachowska-Brzyska K, Qvarnström A*, et al.* 2012. The genomic landscape of species
834      divergence in *Ficedula* flycatchers. *Nature*. 491:756-760.
835  Elshire RJ, Glaubitz JC, Sun Q, Poland JA, Kawamoto K, Buckler ES, Mitchell SE. 2011. A
836      robust, simple genotyping-by-sequencing (GBS) approach for high diversity species.
837      *PLoS One*. 6:e19379.
838  Enk J, Devault A, Kuch M, Murgha Y, Rouillard J-M, Poinar H. 2014. Ancient whole genome
839      enrichment using baits built from modern DNA. *Mol Biol Evol*. 31:1292-1294.
840  Evans TG. 2015. Considerations for the use of transcriptomics in identifying the 'genes that
841      matter' for environmental adaptation. *J Exp Biol*. 218:1925-1935.
842  Excoffier L, Dupanloup I, Huerta-Sánchez E, Sousa VC, Foll M. 2013. Robust demographic
843      inference from genomic and SNP data. *PLoS Genetics*. 9:e1003905.
844  Faircloth BC. 2015. PHYLUCE is a software package for the analysis of conserved genomic
845      loci. *Bioinformatics*. 32:786-788.
846  Faircloth BC, McCormack JE, Crawford NG, Harvey MG, Brumfield RT, Glenn TC. 2012.
847      Ultraconserved elements anchor thousands of genetic markers spanning multiple
848      evolutionary timescales. *Syst Biol*. 61:717-726.
849  Ferrer-Admetlla A, Liang M, Korneliussen T, Nielsen R. 2014. On detecting incomplete soft or
850      hard selective sweeps using haplotype structure. *Mol Biol Evol*. 31:1275-1291.
851  Flicek P, Birney E. 2009. Sense from sequence reads: methods for alignment and assembly. *Nat
852      Methods*. 6:S6-S12.
853  Foote AD, Liu Y, Thomas GWC, Vinař Ts, Alföldi J, Deng J, Dugan S, van Elk CE, Hunter ME,
854      Joshi V*, et al.* 2015. Convergent evolution of the genomes of marine mammals. *Nat
855      Genet*. 47:272-275.
856  Foote AD, Newton J, Ávila-Arcos MC, Kampmann M-L, Samaniego JA, Post K, Rosing-Asvid
857      A, Sinding M-HS, Gilbert MTP. 2013. Tracking niche variation over millennial
858      timescales in sympatric killer whale lineages. *Proc R Soc Lond B Biol Sci*. 280:20131481.
859  Foote AD, Thomsen PF, Sveegaard S, Wahlberg M, Kielgast J, Kyhn LA, Salling AB, Galatius
860      A, Orlando L, Gilbert MTP. 2012. Investigating the potential use of environmental DNA
861      (eDNA) for genetic monitoring of marine mammals. *PLoS One*. 7:e41781.

862 Foote AD, Vijay N, Ávila-Arcos M, Baird RW, Durban JW, Fumagalli M, Gibbs RA, Hanson
863     MB, Korneliussen TS, Martin MD*, et al.* 2016. Genome-culture coevolution promotes
864     rapid divergence of killer whale ecotypes. *Nat Commun*. 7:11693.
865 Fountain ED, Pauli JN, Reid BN, Palsbøll PJ, Peery MZ. 2016. Finding the right coverage: the
866     impact of coverage and sequence quality on single nucleotide polymorphism genotyping
867     error rates. *Mol Ecol Resour*.
868 Fumagalli M, Vieira FG, Korneliussen TS, Linderoth T, Huerta-Sánchez E, Albrechtsen A,
869     Nielsen R. 2013. Quantifying population genetic differentiation from next-generation
870     sequencing data. *Genetics*. 195:979-992.
871 Fumagalli M, Vieira FG, Linderoth T, Nielsen R. 2014. *ngsTools*: methods for population
872     genetics analyses from Next-Generation Sequencing data. *Bioinformatics*. 30:1486-1487.
873 Gao X, Han J, Lu Z, Li Y, He C. 2013. *De novo* assembly and characterization of spotted seal
874     *Phoca largha* transcriptome using Illumina paired-end sequencing. *Comp Biochem*
875     *Physiol D Genom Proteom*. 8:103-110.
876 Garner BA, Hand BK, Amish SJ, Bernatchez L, Foster JT, Miller KM, Morin PA, Narum SR,
877     O'Brien SJ, Roffler G*, et al.* 2016. Genomics in conservation: case studies and bridging
878     the gap between data and application. *Trends Ecol Evol*. 31:81-83.
879 Glaubitz JC, Casstevens TM, Lu F, Harriman J, Elshire RJ, Sun Q, Buckler ES. 2014. TASSEL-
880     GBS: a high capacity genotyping by sequencing analysis pipeline. *PLoS One*. 9:e90346.
881 Gnerre S, MacCallum I, Przybylski D, Ribeiro FJ, Burton JN, Walker BJ, Sharpe T, Hall G, Shea
882     TP, Sykes S*, et al.* 2011. High-quality draft assemblies of mammalian genomes from
883     massively parallel sequence data. *Proc Natl Acad Sci USA*. 108:1513-1518.
884 Goecks J, Nekrutenko A, Taylor J, The Galaxy Team. 2010. Galaxy: a comprehensive approach
885     for supporting accessible, reproducible, and transparent computational research in the life
886     sciences. *Genome Biol*. 11:R86.
887 Gu H, Smith ZD, Bock C, Boyle P, Gnirke A, Meissner A. 2011. Preparation of reduced
888     representation bisulfite sequencing libraries for genome-scale DNA methylation
889     profiling. *Nat Protoc*. 6:468-481.
890 Gui D, Jia K, Xia J, Yang L, Chen J, Wu Y, Yi M. 2013. *De novo* assembly of the Indo-Pacific
891     humpback dolphin leucocyte transcriptome to identify putative genes involved in the
892     aquatic adaptation and immune response. *PLoS One*. 8:e72417.
893 Gutenkunst RN, Hernandez RD, Williamson SH, Bustamante CD. 2009. Inferring the joint
894     demographic history of multiple populations from multidimensional SNP frequency data.
895     *PLoS Genetics*. 5:e1000695.
896 Haas BJ, Papanicolaou A, Yassour M, Grabherr M, Blood PD, Bowden J, Couger MB, Eccles D,
897     Li B, Lieber M*, et al.* 2013. *De novo* transcript sequence reconstruction from RNA-seq
898     using the Trinity platform for reference generation and analysis. *Nat Protoc*. 8:1494-
899     1512.
900 Han E, Sinsheimer JS, Novembre J. 2015. Fast and accurate site frequency spectrum estimation
901     from low coverage sequence data. *Bioinformatics*. 31:720-727.
902 Hancock-Hanser BL, Frey A, Leslie MS, Dutton PH, Archer FI, Morin PA. 2013. Targeted
903     multiplex next-generation sequencing: advances in techniques of mitochondrial and
904     nuclear DNA sequencing for population genomics. *Mol Ecol Resour*. 13:254-268.
905 Harris K, Nielsen R. 2013. Inferring demographic history from a spectrum of shared haplotype
906     lengths. *PLoS Genetics*. 9:e1003521.

907 Hayden S, Bekaert M, Crider TA, Mariani S, Murphy WJ, Teeling EC. 2010. Ecological
908     adaptation determines functional mammalian olfactory subgenomes. *Genome Res*. 20:1-
909     9.
910 Hedrick PW. 2000 *Genetics of Populations*. Jones and Bartlett Publishers, Sudbury, MA.
911 Helyar SJ, Hemmer-Hansen J, Bekkevold D, Taylor MI, Ogden R, Limborg MT, Cariani A,
912     Maes GE, Diopere E, Carvalho GR*, et al.* 2011. Application of SNPs for population
913     genetics of nonmodel organisms: new opportunities and challenges. *Mol Ecol Resour*.
914     11:123-136.
915 Higdon JW, Bininda-Emonds ORP, Beck RMD, Ferguson SH. 2007. Phylogeny and divergence
916     of the pinnipeds (Carnivora: Mammalia) assessed using a multigene dataset. *BMC Evol*
917     *Biol*. 7:216.
918 Hodges E, Rooks M, Xuan Z, Bhattacharjee A, Gordon DB, Brizuela L, McCombie WR,
919     Hannon GJ. 2009. Hybrid selection of discrete genomic intervals on custom-designed
920     microarrays for massively parallel sequencing. *Nat Protoc*. 4:960-974.
921 Hoffman JI. 2011. Gene discovery in the Antarctic fur seal (*Arctocephalus gazella*) skin
922     transcriptome. *Mol Ecol Resour*. 11:703-710.
923 Hoffman JI, Nicholas HJ. 2011. A novel approach for mining polymorphic microsatellite
924     markers *in silico*. *PLoS One*. 6:e23283.
925 Hoffman JI, Simpson F, David P, Rijks JM, Kuiken T, Thorne MAS, Lacy RC, Dasmahapatra
926     KK. 2014. High-throughput sequencing reveals inbreeding depression in a natural
927     population. *Proc Natl Acad Sci USA*. 111:3775-3780.
928 Hoffman JI, Thorne MAS, Trathan PN, Forcada J. 2013. Transcriptome of the dead:
929     characterisation of immune genes and marker development from necropsy samples in a
930     free-ranging marine mammal. *BMC Genomics*. 14:52.
931 Hoffman JI, Tucker R, Bridgett SJ, Clark MS, Forcada J, Slate J. 2012. Rates of assay success
932     and genotyping error when single nucleotide polymorphism genotyping in non-model
933     organisms: a case study in the Antarctic fur seal. *Mol Ecol Resour*. 12:861-872.
934 Hohenlohe PA, Bassham S, Etter PD, Stiffler N, Johnson EA, Cresko WA. 2010. Population
935     genomics of parallel adaptation in threespine stickleback using sequenced RAD tags.
936     *PLoS Genet*. 6:e1000862.
937 Holt C, Yandell M. 2011. MAKER2: an annotation pipeline and genome-database management
938     tool for second-generation genome projects. *BMC Bioinformatics*. 12:491.
939 Humble E, Martinez-Barrio A, Forcada J, Trathan PN, Thorne MAS, Hoffmann M, Wolf JBW,
940     Hoffman JI. 2016. A draft fur seal genome provides insights into factors affecting SNP
941     validation and how to mitigate them. *Mol Ecol Resour*.
942 Jackson JA, Baker CS, Vant M, Steel DJ, Medrano-González L, Palumbi SR. 2009. Big and
943     slow: phylogenetic estimates of molecular evolution in baleen whales (suborder
944     Mysticeti). *Mol Biol Evol*. 26:2427-2440.
945 Jones FC, Grabherr MG, Chan YF, Russell P, Mauceli E, Johnson J, Swofford R, Pirun M, Zody
946     MC, White S*, et al.* 2012. The genomic basis of adaptive evolution in threespine
947     sticklebacks. *Nature*. 484:55-61.
948 Kajitani R, Toshimoto K, Noguchi H, Toyoda A, Ogura Y, Okuno M, Yabana M, Harada M,
949     Nagayasu E, Maruyama H*, et al.* 2014. Efficient *de novo* assembly of highly
950     heterozygous genomes from whole-genome shotgun short reads. *Genome Res*. 24:1384-
951     1395.

952 Keane M, Semeiks J, Webb AE, Li YI, Quesada V, Craig T, Madsen LB, van Dam S, Brawand
953     D, Marques PI, *et al.* 2015. Insights into the evolution of longevity from the bowhead
954     whale genome. *Cell Reports*. 10:112-122.
955 Khudyakov JI, Champagne CD, Preeyanon L, Ortiz RM, Crocker DE. 2015a. Muscle
956     transcriptome response to ACTH administration in a free-ranging marine mammal.
957     *Physiol Genomics*. 47:318-330.
958 Khudyakov JI, Preeyanon L, Champagne CD, Ortiz RM, Crocker DE. 2015b. Transcriptome
959     analysis of northern elephant seal (*Mirounga angustirostris*) muscle tissue provides a
960     novel molecular resource and physiological insights. *BMC Genomics*. 16:64.
961 Kishida T, Thewissen JGM, Hayakawa T, Imai H, Agata K. 2015. Aquatic adaptation and the
962     evolution of smell and taste in whales. *Zoolog Lett*. 1:9.
963 Koepfli K-P, Paten B, Genome 10K Community of Scientists, O'Brien SJ. 2015. The Genome
964     10K Project: a way forward. *Annu Rev Anim Biosci*. 3:57-111.
965 Korneliussen TS, Albrechtsen A, Nielsen R. 2014. ANGSD: Analysis of Next Generation
966     Sequencing Data. *BMC Bioinformatics*. 15:356.
967 Künstner A, Wolf JBW, Backström N, Whitney O, Balakrishnan CN, Day L, Edwards SV, Janes
968     DE, Schlinger BA, Wilson RK, *et al.* 2010. Comparative genomics based on massive
969     parallel transcriptome sequencing reveals patterns of substitution and selection across 10
970     bird species. *Mol Ecol*. 19:266-276.
971 Lamichhaney S, Berglund J, Almén MS, Maqbool K, Grabherr M, Martinez-Barrio A,
972     Promerová M, Rubin C-J, Wang C, Zamani N, *et al.* 2015. Evolution of Darwin's finches
973     and their beaks revealed by genome sequencing. *Nature*. 518:371-375.
974 Langmead B, Trapnell C, Pop M, Salzberg SL. 2009. Ultrafast and memory-efficient alignment
975     of short DNA sequences to the human genome. *Genome Biol*. 10:R25.
976 Lemmon AR, Emme SA, Lemmon EM. 2012. Anchored hybrid enrichment for massively high-
977     throughput phylogenomics. *Syst Biol*. 61:727-744.
978 Li B, Dewey CN. 2011. RSEM: accurate transcript quantification from RNA-Seq data with or
979     without a reference genome. *BMC Bioinformatics*. 12:323.
980 Li H, Durbin R. 2009. Fast and accurate short read alignment with Burrows-Wheeler transform.
981     *Bioinformatics*. 25:1754-1760.
982 Li H, Durbin R. 2011. Inference of human population history from individual whole-genome
983     sequences. *Nature*. 475:493-496.
984 Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R,
985     1000 Genome Project Data Processing Subgroup. 2009. The Sequence Alignment/Map
986     form and SAMtools. *Bioinformatics*. 25:2078-2079.
987 Li S, Jakobsson M. 2012. Estimating demographic paramaters from large-scale population
988     genomic data using Approximate Bayesian Computation. *BMC Genet*. 13:22.
989 Li Y, Hu Y, Bolund L, Wang J. 2010. State of the art *de novo* assembly of human genomes from
990     massively parallel sequencing data. *Human Genomics* 4:271-277.
991 Lindblad-Toh K, Garber M, Zuk O, Lin MF, Parker BJ, Washietl S, Kheradpour P, Ernst J,
992     Jordan G, Mauceli E, *et al.* 2011. A high-resolution map of human evolutionary
993     constraint using 29 mammals. *Nature*. 478:476-482.
994 Lindqvist C, Schuster SC, San Y, Talbot SL, Qi J, Ratan A, Tomsho LP, Kasson L, Zeyl E, Aars
995     J, *et al.* 2010. Complete mitochondrial genome of a Pleistocene jawbown unveils the
996     origin of polar bear. *Proc Natl Acad Sci USA*. 107:5053-5057.

997   Liu S, Lorenzen ED, Fumagalli M, Li B, Harris K, Xiong Z, Zhou L, Korneliussen TS, Somel M,
998        Babbitt C, *et al.* 2014a. Population genomics reveal recent speciation and rapid
999        evolutionary adaptation in polar bears. *Cell*. 157:785-794.
1000  Liu X, Fun Y-X. 2015. Exploring population size changes using SNP frequency spectra. *Nat
1001       Genet*. 47:555-559.
1002  Liu Y, Zhou J, White KP. 2014b. RNA-seq differential expression studies: more sequence or
1003       more replication? *Bioinformatics*. 30:301-304.
1004  Lotterhos KE, Whitlock MC. 2014. Evaluation of demographic history and neutral
1005       parameterization on the performance of $F_{ST}$ outlier tests. *Mol Ecol*. 23:2178-2192.
1006  Louis M, Viricel A, Lucas T, Peltier H, Alfonsi E, Berrow S, Brownlow A, Covelo P, Dabin W,
1007       Deaville R, *et al*. 2014. Habitat-driven population structure of bottlenose dolphins,
1008       *Tursiops truncatus*, in the North-east Atlantic. *Mol Ecol*. 23:857-874.
1009  Love MI, Huber W, Anders S. 2014. Moderated estimation of fold change and dispersion for
1010       RNA-seq data with DESeq2. *Genome Biol*. 15:550.
1011  MacManes MD. 2014. On the optimal trimming of high-throughput mRNA sequence data. *Front
1012       Genet*. 5:13.
1013  MacManes MD. 2016. Establishing evidence-based best practice for the *de novo* assembly and
1014       evaluation of transcriptomes from non-model organisms. *bioRxiv*.  doi:
1015       http://dx.doi.org/10.1101/035642.
1016  Magera AM, Mills Flemming JE, Kaschner K, Christensen LB, Lotze HK. 2013. Recovery
1017       trends in marine mammal populations. *PLoS One*. 8:e77908.
1018  Malenfant RM, Coltman DW, Davis CS. 2015. Design of a 9K Illumina BeadChip for polar
1019       bears (*Ursus maritimus*) from RAD and transcriptome sequencing. *Mol Ecol Resour*.
1020       15:587-600.
1021  Mamanova L, Coffey AJ, Scott CE, Kozarewa I, Turner EH, Kumar A, Howard E, Shendure J,
1022       Turner DJ. 2010. Target-enrichment strategies for next-generation sequencing. *Nat
1023       Methods*. 7:111-118.
1024  Mancia A, Abelli L, Kucklick JR, Rowles TK, Wells RS, Balmer BC, Hohn AA, Baatz JE, Ryan
1025       JC. 2015. Microarray applications to understand the impact of exposure to environmental
1026       contaminants in wild dolphins (*Tursiops truncatus*). *Mar Genomics*. 19:47-57.
1027  Mancia A, Lundqvist ML, Romano TA, Peden-Adams MM, Fair PA, Kindy MS, Ellis BC,
1028       Gattoni-Celli S, McKillen DJ, Trent HF, *et al.* 2007. A dolphin peripheral blood
1029       leukocyte cDNA microarray for studies of immune function and stress reactions. *Dev
1030       Comp Immunol*. 31:520-529.
1031  Mancia A, Ryan JC, Chapman RW, Wu Q, Warr GW, Gulland FMD, Van Dolah FM. 2012.
1032       Health status, infection and disease in California sea lions (*Zalophus californianus*)
1033       studied using a canine microarray platform and machine-learning approaches. *Dev Comp
1034       Immunol*. 36:629-637.
1035  Mancia A, Warr GW, Chapman RW. 2008. A transcriptomic analysis of the stress induced by
1036       capture-release health assessment studies in wild dolphins (*Tursiops truncatus*). *Mol
1037       Ecol*. 17:2581-2589.
1038  Mastretta-Yanes A, Arrigo N, Alvarez N, Jorgensen TH, Piñero D, Emerson BC. 2015.
1039       Restriction site-associated DNA sequencing, genotyping error estimation and *de novo*
1040       assembly optimization for population genetic inference. *Mol Ecol Resour*. 15:28-41.

1041  McCormack JE, Faircloth BC, Crawford NG, Gowaty PA, Brumfield RT, Glenn TC. 2012.
1042      Ultraconserved elements are novel phylogenomic markers that resolve placental mammal
1043      phylogeny when combined with species-tree analysis. *Genome Res*. 22:746-754.
1044  McGowen MR. 2011. Toward the resolution of an explosive radiation - a multilocus phylogeny
1045      of oceanic dolphins (Delphinidae). *Mol Phylogenet Evol*. 60:345-357.
1046  McGowen MR, Clark C, Gatesy J. 2008. The vestigial olfactory receptor subgenome of
1047      odontocete whales: phylogenetic congruence between gene-tree reconciliation and
1048      supermatrix methods. *Syst Biol*. 57:574-590.
1049  McGowen MR, Gatesy J, Wildman DE. 2014. Molecular evolution tracks macroevolutionary
1050      transitions in Cetacea. *Trends Ecol Evol*. 29:336-346.
1051  McGowen MR, Grossman LI, Wildman DE. 2012. Dolphin genome provides evidence for
1052      adaptive evolution of nervous system genes and a molecular rate slowdown. *Proc R Soc*
1053      *Lond B Biol Sci*. 279:3643-3651.
1054  McGowen MR, Spaulding M, Gatesy J. 2009. Divergence date estimation and a comprehensive
1055      molecular tree of extant cetaceans. *Mol Phylogenet Evol*. 53:891-906.
1056  McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, Garimella K,
1057      Altshuler D, Gabriel S, Daly M*, et al.* 2010. The Genome Analysis Toolkit: A
1058      MapReduce framework for analyzing next-generation DNA sequencing data. *Genome*
1059      *Res*. 20:1297-1303.
1060  McTavish EJ, Hillis DM. 2015. How do SNP ascertainment schemes and population
1061      demographics affect inferences about population history? *BMC Genomics*. 16:266.
1062  Meredith RW, Gatesy J, Emerling CA, York VM, Springer MS. 2013. Rod monochromacy and
1063      the coevolution of cetacean retinal opsins. *PLoS Genetics*. 9:e1003432.
1064  Meyer M, Kircher M, Gansauge M-T, Li H, Racimo F, Mallick S, Schraiber JG, Jay F, Prüfer K,
1065      de Filippo C*, et al.* 2012. A high-coverage genome sequence from an archaic Denisovan
1066      individual. *Science*. 338:222-226.
1067  Miller MR, Dunham JP, Amores A, Cresko WA, Johnson EA. 2007. Rapid and cost-effective
1068      polymorphism identification and genotyping using restriction site associated DNA
1069      (RAD) markers. *Genome Res*. 17:240-248.
1070  Miller W, Schuster SC, Welch AJ, Ratan A, Bedoya-Reina OC, Zhao F, Kim HL, Burhans RC,
1071      Drautz DI, Wittekindt NE*, et al.* 2012. Polar and brown bear genomes reveal ancient
1072      admixture and demographic footprints of past climate change. *Proc Natl Acad Sci USA*.
1073      109:E2382-E2390.
1074  Mirceta S, Signore AV, Burns JM, Cossins AR, Campbell KL, Berenbrink M. 2013. Evolution
1075      of mammalian diving capacity traced by myoglobin net surface charge. *Science*.
1076      340:1234192.
1077  Morin PA, Archer FI, Foote AD, Vilstrup J, Allen EE, Wade P, Durban JW, Parsons K, Pitman
1078      R, Li L*, et al.* 2010a. Complete mitochondrial genome phylogeographic analysis of killer
1079      whales (*Orcinus orca*) indicates multiple species. *Genome Res*. 20:908-916.
1080  Morin PA, Luikart G, Wayne RK, SNP workshop group. 2004. SNPs in ecology, evolution and
1081      conservation. *Trends Ecol Evol*. 19:208-216.
1082  Morin PA, Martien KK, Archer FI, Cipriano F, Steel D, Jackson J, Taylor BL. 2010b. Applied
1083      conservation genetics and the need for quality control and reporting of genetic data used
1084      in fisheries and wildlife management. *J Hered*. 101:1-10.

1085 Morin PA, Parsons KM, Archer FI, Ávila-Arcos M, Barrett-Lennard LG, Dalla Rosa L, Duchêne
1086           S, Durban JW, Ellis GM, Ferguson SH, *et al.* 2015. Geographic and temporal dynamics
1087           of a global radiation and diversification in the killer whale. *Mol Ecol*. 24:3964-3979.
1088 Moura AE, Kenny JG, Chaudhuri R, Hughes MA, Welch AJ, Reisinger RR, de Bruyn PJN,
1089           Dahlheim ME, Hall N, Hoelzel AR. 2014a. Population genomics of the killer whale
1090           indicates ecotype evolution in sympatry involving both selection and drift. *Mol Ecol*.
1091           23:5179-5192.
1092 Moura AE, Nielsen SCA, Vilstrup JT, Moreno-Mayar JV, Gilbert MTP, Gray HWI, Natoli A,
1093           Möller L, Hoelzel AR. 2013. Recent diversification of a marine genus (*Tursiops* spp.)
1094           tracks habitat preference and environmental change. *Syst Biol*. 62:865-877.
1095 Moura AE, van Rensburg CJ, Pilot M, Tehrani A, Best PB, Thornton M, Plön S, de Bruyn PJN,
1096           Worley KC, Gibbs RA, *et al.* 2014b. Killer whale nuclear genome and mtDNA reveal
1097           widespread population bottleneck during the last glacial maximum. *Mol Biol Evol*.
1098           31:1121-1131.
1099 Nadeau NJ, Ruiz M, Salazar P, Counterman B, Alejandro Medina J, Ortiz-Zuazaga H, Morrison
1100           A, McMillan WO, Jiggins CD, Papa R. 2014. Population genomics of parallel hybrid
1101           zones in the mimetic butterflies, *H. melpomene* and *H. erato*. *Genome Res*. 24:1316-
1102           1333.
1103 Narum SR, Buerkle CA, Davey JW, Miller MR, Hohenlohe PA. 2013. Genotyping-by-
1104           sequencing in ecological and conservation genomics. *Mol Ecol*. 22:2841-2847.
1105 Narum SR, Hess JE. 2011. Comparison of $F_{ST}$ outlier tests for SNP loci under selection. *Mol
1106           Ecol Resour*. 11:184-194.
1107 Nelson TM, Apprill A, Mann J, Rogers TL, Brown MV. 2015. The marine mammal microbiome:
1108           current knowledge and future directions. *Microbiology Australia*. 36:8-13.
1109 Ng SB, Turner EH, Robertson PD, Flygare SD, Bigham AW, Lee C, Shaffer T, Wong M,
1110           Bhattacharjee A, Eichler EE, *et al.* 2009. Targeted capture and massively parallel
1111           sequencing of twelve human exomes. *Nature*. 461:272-276.
1112 Nielsen R, Paul JS, Anders A, Song YS. 2011. Genotype and SNP calling from next-generation
1113           sequencing data. *Nat Rev Genet*. 12:433-451.
1114 Noonan JP, Coop G, Kudaravalli S, Smith D, Krause J, Alessi J, Chen F, Platt D, Pääbo S,
1115           Pritchard JK, *et al.* 2006. Sequencing and analysis of Neanderthal genomic DNA.
1116           *Science*. 314:1113-1118.
1117 Nosek BA, Alter G, Banks GC, Borsboom D, Bowman SD, Breckler SJ, Buck S, Chambers CD,
1118           Chin G, Christensen G, *et al.* 2015. Promoting an open research culture: Author
1119           guidelines for journals could help to promote transparency, openness, and reproducibility.
1120           *Science*. 348:1422-1425.
1121 O'Rawe JA, Ferson S, Lyon GJ. 2015. Accounting for uncertainty in DNA sequencing data.
1122           *Trends Genet*. 31:61-66.
1123 Olsen MT, Volny VH, Bérubé M, Dietz R, Lydersen C, Kovacs KM, Dodd RS, Palsbøll PJ.
1124           2011. A simple route to single-nucleotide polymorphisms in a nonmodel species:
1125           identification and characterization of SNPs in the Arctic ringed seal (*Pusa hispida
1126           hispida*). *Mol Ecol Resour*. 11:9-19.
1127 Orlando L, Ginolhac A, Zhang G, Froese D, Albrechtsen A, Stiller M, Schubert M, Cappellini E,
1128           Petersen B, Moltke I, *et al.* 2013. Recalibrating *Equus* evolution using the genome
1129           sequence of an early Middle Pleistocene horse. *Nature*. 499:74-78.

1130 Pabuwal V, Boswell M, Pasquali A, Wise SS, Kumar S, Shen Y, Garcia T, Lacerte C, Wise JP,
1131    Jr., Wise JP, Sr*., et al.* 2013. Transcriptomic analysis of cultured whale skin cells exposed
1132    to hexavalent chromium [Cr(VI)]. *Aquat Toxicol*. 134-135:74-81.

1133 Parker J, Tsagkogeorga G, Cotton JA, Liu Y, Provero P, Stupka E, Rossiter SJ. 2013. Genome-
1134    wide signatures of convergent evolution in echolocating mammals. *Nature*. 502:228-231.

1135 Paszkiewicz KH, Farbox A, O'Neill P, Moore K. 2014. Quality control on the frontier. *Front
1136    Genet*. 5:157.

1137 Patro R, Duggal G, Kingsford C. 2015. Accurate, fast, and model-aware transcript expression
1138    quantification with Salmon. *bioRxiv*.  doi: http://dx.doi.org/10.1101/021592.

1139 Peterson BK, Weber JN, Kay EH, Fisher HS, Hoekstra HE. 2012. Double digest RADseq: an
1140    inexpensive method for *de novo* SNP discovery and genotyping in model and non-model
1141    species. *PLoS One*. 7:e37135.

1142 Poh Y-P, Domingues VS, Hoekstra HE, Jensen JD. 2014. On the prospect of identifying adaptive
1143    loci in recently bottlenecked populations. *PLoS One*. 9:e110579.

1144 Poland JA, Brown PJ, Sorrells ME, Jannink J-L. 2012. Development of high-density genetic
1145    maps for barley and wheat using a novel two-enzyme genotyping-by-sequencing
1146    approach. *PLoS One*. 7:e32253.

1147 Polanowski AM, Robbins J, Chandler D, Jarman SN. 2014. Epigenetic estimation of age in
1148    humpback whales. *Mol Ecol Resour*. 14:976-987.

1149 Puritz JB, Hollenbeck CM, Gold JR. 2014. *dDocent*: a RADseq, variant-calling pipeline
1150    designed for population genomics of non-model organisms. *PeerJ*. 2:e431.

1151 Rasmussen M, Li Y, Lindgreen S, Pedersen JS, Albrechtsen A, Moltke I, Metspalu M, Metspalu
1152    E, Kivisild T, Gupta R*, et al*. 2010. Ancient human genome sequence of an extinct
1153    Palaeo-Eskimo. *Nature*. 463:757-762.

1154 Riesch R, Barrett-Lennard LG, Ellis GM, Ford JKB, Deecke VB. 2012. Cultural traditions and
1155    the evolution of reproductive isolation: ecological speciation in killer whales? *Biol J Linn
1156    Soc Lond*. 2012:1-17.

1157 Robinson JD, Coffman AJ, Hickerson MJ, Gutenkunst RN. 2014. Sampling strategies for
1158    frequency spectrum-based population genomic inference. *BMC Evol Biol*. 14:254.

1159 Robinson MD, McCarthy DJ, Smyth GK. 2010. edgeR: a Bioconductor package for differential
1160    expression analysis of digital gene expression data. *Bioinformatics*. 26:139-140.

1161 Ruan R, Guo A-H, Hao Y-J, Zheng J-S, Wang D. 2015. *De novo* assembly and characterization
1162    of narrow-ridged finless porpoise renal transcriptome and identification of candidate
1163    genes involved in osmoregulation. *Int J Mol Sci*. 16:2220-2238.

1164 Ruegg K, Rosenbaum HC, Anderson EC, Engel M, Rothschild A, Baker CS, Palumbi SR. 2013.
1165    Long-term population size of the North Atlantic humpback whale within the context of
1166    worldwide population structure. *Cons Gen*. 14:103-114.

1167 Schiffels S, Durbin R. 2014. Inferring human population size and separation history from
1168    multiple genome sequences. *Nat Genet*. 46:919-925.

1169 Schubert M, Lindgreen S, Orlando L. 2016. AdapterRemoval v2: rapid adapter trimming,
1170    identification, and read merging. *BMC Res Notes*. 9:88.

1171 Schurch NJ, Schofield P, Gierlinski M, Cole C, Sherstnev A, Singh V, Wrobel N, Gharbi K,
1172    Simpson GG, Owen-Hughes T*, et al*. 2015. Evaluation of tools for differential gene
1173    expression analysis by RNA-seq on a 48 biological replicate experiment. *arXive*.
1174    1505:02017.

1175  Seim I, Ma S, Zhou X, Gerashchenko MV, Lee SG, Suydam R, George JC, Bickham JW,
1176        Gladyshev VN. 2014. The transcriptome of the bowhead whale *Balaena mysticetus*
1177        reveals adaptations of the longest-lived mammal. *Aging*. 6:879-899.
1178  Shafer ABA, Cullingham CI, Côté SD, Coltman DW. 2010. Of glaciers and refugia: a decade of
1179        study sheds new light on the phylogeographic patterns of northwestern North America.
1180        *Mol Ecol*. 19:4589-4621.
1181  Shafer ABA, Davis CS, Coltman DW, Stewart REA. 2014. Microsatellite assessment of walrus
1182        (*Odobenus rosmarus rosmarus*) stocks in Canada. *NAMMCO Scientific Publications*. 9.
1183  Shafer ABA, Gattepaille LM, Stewart REA, Wolf JBW. 2015. Demographic inferences using
1184        short-read genomic data in an approximate Bayesian computation framework: *in silico*
1185        evaluation of power, biases and proof of concept in Atlantic walrus. *Mol Ecol*. 24:328-
1186        345.
1187  Shen Y-Y, Zhou W-P, Zhou T-C, Zeng Y-N, Li G-M, Irwin DM, Zhang Y-P. 2012. Genome-
1188        wide scan for bats and dolphin to detect their genetic basis for new locomotive styles.
1189        *PLoS One*. 7:e46455.
1190  Smith-Unna RD, Boursnell C, Patro R, Hibberd JM, Kelly S. 2015. TransRate: reference free
1191        quality assessment of *de-novo* transcriptome assemblies. *bioRxiv*.
1192  Spies D, Ciaudo C. 2015. Dynamics in transcriptomics: advancements in RNA-seq time course
1193        and downstream analysis. *Comput Struct Biotechnol J*. 13:469-477.
1194  Springer MS, Signore AV, Paijmans JLA, Vélez-Juarbe J, Domning DP, Bauer CE, He K, Crerar
1195        L, Campos PF, Murphy WJ, *et al.* 2015. Interordinal gene capture, the phylogenetic
1196        position of Steller's sea cow based on molecular and morphological data, and the
1197        macroevolutionary history of Sirenia. *Mol Phylogenet Evol*. 91:178-193.
1198  Springer MS, Starrett J, Morin PA, Lanzetti A, Hayashi C, Gatesy J. 2016. Inactivation of
1199        *C4orf26* in toothless placental mammals. *Mol Phylogenet Evol*. 95:34-45.
1200  Sremba AL, Martin AR, Baker CS. 2015. Species identification and likely catch time preiod of
1201        whale bones from South Georgia. *Mar Mamm Sci*. 31:122-132.
1202  Stanke M, Keller O, Gunduz I, Hayes A, Waack S, Morgenstern B. 2006. AUGUSTUS: *ab initio*
1203        prediction of alternative transcripts. *Nucleic Acids Res*. 34:W435-W439.
1204  Stein LD. 2010. The case for cloud computing in genome informatics. *Genome Biol*. 11:207.
1205  Stinchcombe JR, Hoekstra HE. 2008. Combining population genomics and quantitative genetics:
1206        finding genes underlying ecologically important traits. *Heredity*. 100:158-170.
1207  Tabuchi M, Veldhoen N, Dangerfield N, Jeffries S, Helbing CC, Ross PS. 2006. PCB-related
1208        alteration of thyroid hormones and thyroid hormone receptor gene expression in free-
1209        ranging harbor seals (*Phoca vitulina*). *Environ Health Perspect*. 114:1024-1031.
1210  Taylor BL, Gemmell NJ. 2016. Emerging technologies to conserve biodiversity: further
1211        opportunities via genomics. Response to Pimm *et al. Trends Ecol Evol*. 31:171-172.
1212  The *Heliconius* Genome Consortium. 2012. Butterfly genome reveals promiscuous exchange of
1213        mimicry adaptations among species. *Nature*. 487:94-98.
1214  Thomsen PF, Kielgast J, Iversen LL, Møller PR, Rasmussen M, Willerslev E. 2012. Detection of
1215        a diverse marine fish fauna using environmental DNA from seawater samples. *PLoS One*.
1216        7:e41732.
1217  Towns J, Cockerill T, Dahan M, Foster I, Gaither K, Grimshaw A, Hazlewood V, Lathrop S,
1218        Lifka D, Peterson GD, *et al.* 2014. XSEDE: accelerating scientific discovery. *Computing
1219        in Science and Engineering*. 16:62-74.

Tsagkogeorga G, McGowen MR, Davies KT, Jarman S, Polanowski A, Bertelsen MF, Rossiter SJ. 2015. A phylogenomic analysis of the role and timing of molecular adaptation in the aquatic transition of cetartiodactyl mammals. *R Soc Open Sci*. 2:150156.

van Dijk EL, Auger H, Jaszczyzyn Y, Thermes C. 2014. Ten years of next-generation sequencing technology. *Trends Genet*. 30:418-426.

VanRaden PM, Sun C, O'Connell JR. 2015. Fast imputation using medium or low-coverage sequence data. *BMC Genet*. 16:82.

Villar D, Berthelot C, Aldridge S, Rayner TF, Lukk M, Pignatelli M, Park TJ, Deaville R, Erichsen JT, Jasinska AJ, *et al*. 2015. Enhancer evolution across 20 mammalian species. *Cell*. 160:554-566.

Viricel A, Pante E, Dabin W, Simon-Bouhet B. 2014. Applicability of RAD-tag genotyping for interfamilial comparisons: empirical data from two cetaceans. *Mol Ecol Resour*. 14:597-605.

Viricel A, Rosel PE. 2014. Hierarchical population structure and habitat differences in a highly mobile marine species: the Atlantic spotted dolphin. *Mol Ecol*. 23:5018-5035.

Wolf JB. 2013. Principles of transcriptome analysis and gene expression quantification: an RNA-seq tutorial. *Mol Ecol Resour*. 13:559-572.

Xiong Y, Brandley MC, Xu S, Zhou K, Yang G. 2009. Seven new dolphin mitochondrial genomes and a time-calibrated phylogeny of whales. *BMC Evol Biol*. 9:20.

Yandell M, Ence D. 2012. A beginner's guide to eukaryotic genome annotation. *Nat Rev Genet*. 13:329-342.

Yeh R-F, Lim LP, Burge CB. 2001. Computational inference of homologous gene structures in the human genome. *Genome Res*. 11:803-816.

Yim H-S, Cho YS, Guang X, Kang SG, Jeong J-Y, Cha S-S, Oh H-M, Lee J-H, Yang EC, Kwon KK, *et al*. 2014. Minke whale genome and aquatic adaptation in cetaceans. *Nat Genet*. 46:88-92.

Zhao Q-Y, Wang Y, Kong Y-M, Luo D, Li X, Hao P. 2011. Optimizing *de novo* transcriptome assembly from short-read RNA-Seq data: a comparative study. *BMC Bioinformatics*. 12:S2.

Zhou X, Sun F, Xu S, Fan G, Zhu K, Liu X, Chen Y, Shi C, Yang Y, Huang Z, *et al*. 2013. Baiji genomes reveal low genetic variability and new insights into secondary aquatic adaptations. *Nat Commun*. 4:2708.

Zou Z, Zhang J. 2015. No genome-wide protein sequence convergence for echolocation. *Mol Biol Evol*. 32:1237-1241.

1255 Table 1. Current and commonly used tools for analysis of genomic data generated in non-model organisms. Please note that this list is
1256 not exhaustive and new computational tools are continuously being developed.
1257

| **Computational Tool** | **Purpose** | **Strengths/Weaknesses** | **Reference** |
|---|---|---|---|
| *RADseq\** | | | |
| STACKS | quality filtering, *de novo* assembly or reference-aligned read mapping, variant genotyping | scalable (new data can be compared against existing locus catalog); flexible filtering and export options; recently implemented a gapped alignment algorithm to process insertion-deletion (indel) mutations; secondary algorithm adjusts SNP calls using population-level allele frequencies; compatible with input data from multiple RADseq methods | Catchen et al. (2011; 2013), http://catchenlab.life.illinois.edu/stacks/ |
| PyRAD | quality filtering, *de novo* assembly, read mapping, variant genotyping | efficiently processes indel mutations, thus optimal for analysis of highly divergent species; high speed and quality of paired-end library assemblies; compatible with input data from multiple RADseq methods | Eaton (2014) |
| TASSEL-GBS | quality filtering, reference-aligned read mapping, variant genotyping | optimized for single-end data from large sample sizes (tens of thousands of individuals) with a reference genome; performs genome-wide association studies | Glaubitz et al. (2014) |
| dDocent | quality trimming, *de novo* assembly, read mapping, variant genotyping | beneficial in analysis of paired-end data; identifies both SNP and indel variants; most appropriate for ezRAD and ddRAD data | Puritz et al. (2014) |
| AftrRAD | quality filtering, *de novo* assembly, read mapping, variant genotyping | identifies both SNP and indel variants; computationally faster than STACKS and PyRAD | Sovic et al. (2015) |
| *Array-based high-throughput sequencing* | | | |
| Affymetrix Axiom™ Analysis Suite, Illumina® GenomeStudio | genotype scoring | visualization of genotype clusters; quality scores assigned to genotype calls allow user-specific filtering; manual editing possible | |
| *Whole genome sequencing* | | | |
| AdapterRemoval v2, Trimmomatic | trim raw sequences | remove adapter sequences and low-quality bases prior to assembly | Bolger et al. (2014), Schubert et al. (2016) |
| ALLPATHS-LG, PLATANUS, SOAPdenovo | *de novo* genome assembly | designed for short-read sequences of large heterozygous genomes | Li et al. (2010), Gnerre et al. (2011), Kajitani et al. (2014) |
| AUGUSTUS, GenomeScan, MAKER2 | gene annotation | highly accurate evidence-driven or BLASTX-guided gene prediction (Yandell and Ence 2012) | Yeh et al. (2001), Stanke et al. (2006), Holt and Yandell (2011) |

| | | | |
|---|---|---|---|
| Bowtie, bwa | read mapping | rapid short-read alignment with compressed reference genome index, but limited number of acceptable mismatches per alignment (Flicek and Birney 2009) | Langmead et al. (2009), Li and Durbin (2009) |
| SAMtools | data processing, variant calling ~~(SNP and indel discovery)~~ | multi-purpose tool that conducts file conversion, alignment sorting, PCR duplicate removal, and variant (SNP and indel) calling for SAM/BAM/CRAM files | Li et al. (2009) |
| GATK | data processing and quality control, variant calling | suitable for ~~processing and analyses of~~ data with low to high mean read depth across the genome~~coverage data~~; initially optimized for large human datasets, then modified for use with non-model organisms | McKenna et al. (2010), DePristo et al. (2011) |
| ANGSD/NGStools | data processing, variant calling, estimation of diversity metrics, population genomic analyses | suitable for ~~processing and analyses of~~ data with low mean read depth, including ~~coverage and~~ palaeogenomic data; allow downstream analyses such as D-statistics and SFS estimation | Fumagalli et al. (2014), Korneliussen et al. (2014) |
| *RNAseq* | | | |
| Fastx Toolkit, Trimmomatic | trim raw sequences | remove erroneous nucleotides from reads prior to assembly | MacManes (2014) |
| khmer diginorm, Trinity normalization | *in silico* read normalization | reduce~~s~~ memory requirements for assembly, but can result in fragmented assemblies and collapse heterozygosity | Brown et al. (2012); Haas et al. (2013) |
| Trinity | *de novo* and genome-guided transcriptome assembly | accurate assembly across conditions, but requires long runtime if normalization is not used (Zhao et al. 2011) | Haas et al. (2013) |
| bowtie, bowtie2, STAR | read alignment to genome or transcriptome assembly | required for many downstream analyses, but bowtie is computationally intensive and all produce very large output BAM files | Langmead et al. (2009), Dobin et al. (2013) |
| eXpress, kallisto, RSEM, Sailfish, Salmon | estimation of transcript abundance | RSEM requires computationally intensive read mapping back to the assembly; the others are faster streaming alignment, quasi-alignment, or alignment-free algorithms | Li and Dewey (2011), Patro et al. (2015) |
| DESeq, DESeq2, edgeR | differential expression analysis | exhibit highest true positive and lowest false positive rates in experiments with smaller sample sizes (Schurch et al. 2015) | Anders and Huber (2010), Robinson et al. (2010), Love et al. (2014) |
| blast2GO, Trinotate | functional annotation of assembled transcripts | complete annotation pipelines including gene ontology and pathway enrichment analyses | Conesa et al. (2005), Haas et al. (2013) |

1258  \* This is a non-exhaustive list of software that ~~include~~ <u>focuses on</u> *de novo* loci assembly and genotype calling for RADseq data, as many practitioners working on
1259  NMOs will not have access to a reference genome. Other programs (e.g., GATK and ANGSD) that undertake genotype calling using reference-aligned loci ~~only~~
1260  are described in the whole genome sequencing section.

**Figure 1**



Figure 1. Phylogenetic tree showing current genomic resources available for (A) cetaceans and (B) pinnipeds; relationships and branch lengths are based on molecular dating estimates from McGowen et al. (2009), McGowen (2011), and Higdon et al. (2007). Scale is in millions of years ago (MYA). Red circles indicate species with high-~~coverage~~ quality ~~whole~~ reference genomes; green stars indicate ~~low-coverage~~ whole genome re-sequencing data; blue triangles indicate transcriptomes (generated by microarray or RNAseq); and black squares indicate RADseq data.

**Figure 2**



Figure 2. Number of marine mammal genomics publications from 1990 to 2015, categorized by primary methodology and research aim. Genomic methodologies include high-throughput single nucleotide polymorphism (SNP) genotyping and sequencing of mitogenomes, whole genomes (WGS), transcriptomes (generated by microarray or RNAseq), and reduced-representation genomic libraries (RRL). The "Other" category includes studies of microbiomes, BAC libraries, and large (~100) gene sets.

**Figure 3**



Figure 3. Number of BioProjects (~~shaded~~ gray bars) related to marine mammal genomics submitted from 2006 to 2015 to an online public database maintained by NCBI. Early BioProjects were largely microarray datasets. The number of projects created each year, as well as the yearly average (black dots ± SE) and maximum (×) size of data submitted in each BioProject, increased dramatically after 2011, reflecting advances in high-throughput sequencing technologies that facilitated their use in non-model systems.

**Figure 4**

| | 2002 --- 2009 | 2010 | 2011 | 2012 | 2013 | 2014 | 2015 | 2016 |
|---|---|---|---|---|---|---|---|---|
| **Common bottlenose dolphin** (*Tursiops truncatus*) | Sequence mitogenome[1] | | Sequence genome (2.8x)[2] | | Population mitogenomics[3] | | RADseq[4] / Improve genome (3.5x 454, 30x Illumina)[5] | Shotgun sequencing for SNP discovery[6] |
| **Killer whale** (*Orcinus orca*) | | Population mitogenomics[7] | | | | RADseq[8] / Sequence genome (20x)[9] | Sequence genome (200x)[5] | Population genome re-sequencing (avg 2x, N = 48)[10] |
| **Antarctic fur seal** (*Arctocephalus gazella*) | | | Sequence transcriptome[11] / Microsat. discovery[12] | SNP discovery[13] | | | | Sequence genome (200x) to validate SNPs[14] |
| **Polar bear** (*Ursus maritimus*) | Sequence mitogenome[15] | Ancient DNA mitogenome[16] | | Sequence genome (100x) & transcriptome[17] | | Population genome re-sequencing (avg 3.5x, N=61)[18] | RADseq & transcriptome sequencing for SNP discovery[19] | |

1) Xiong et al. 2009; 2) Lindblad-Toh et al. 2011; 3) Moura et al. 2013; 4) Cammen et al. 2015; 5) Foote et al. 2015; 6) Louis et al. unpubl. data; 7) Morin et al. 2010; 8) Moura et al. 2014a; 9) Moura et al. 2014b; 10) Foote et al. 2016; 11) Hoffman 2011; 12) Hoffman and Nicholas 2011; 13) Hoffman et al. 2012; 14) Humble et al. 2016; 15) Arnason et al. 2002; 16) Lindqvist et al. 2010; 17) Miller et al. 2012; 18) Liu et al. 2014; 19) Malenfant et al. 2015
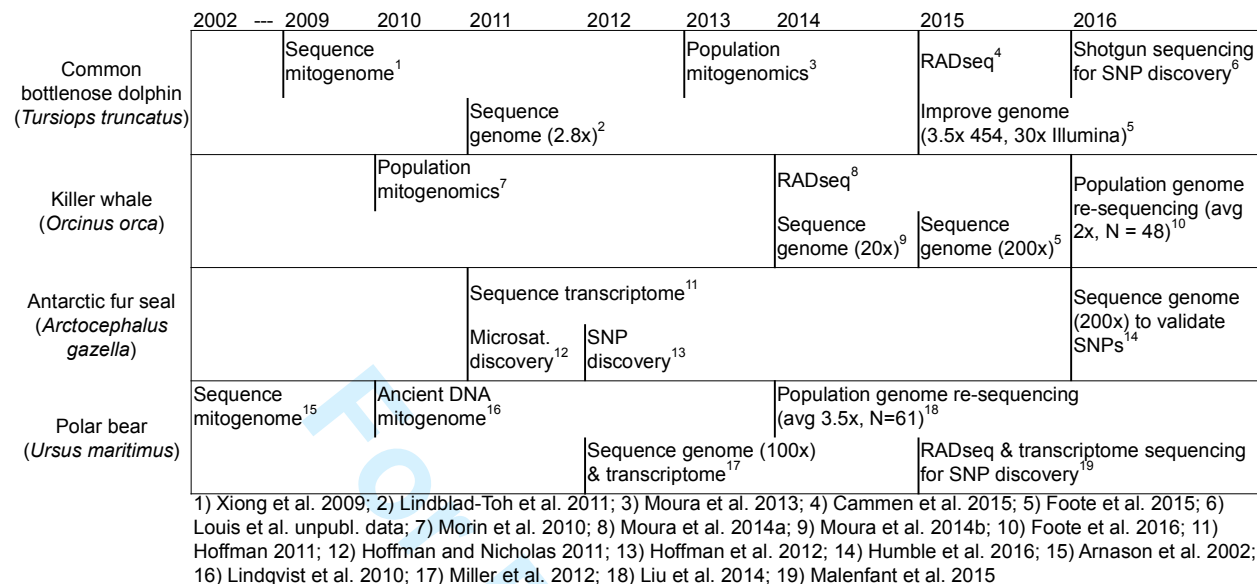
Figure 4. Timelines depicting the independent progression of genomic studies for four representative marine mammal species. Trajectories show the common progression for non-model species from mitogenome sequencing to whole genome sequencing, as well as from sequencing reference specimens to population-scale genomic sequencing. In addition, the timelines reveal the utility of genomic and transcriptomic sequencing for subsequent genetic marker development.