1    **Title**: High quality whole genome sequence of an abundant Holarctic odontocete, the harbour

2    porpoise (*Phocoena phocoena*)

3

4    **Authors**: Marijke Autenrieth[1], Stefanie Hartmann[2], Ljerka Lah[1,3], Anna Roos[4], Alice B.

5    Dennis[1], Ralph Tiedemann[1]*

6

7    [1] University of Potsdam, Institute of Biochemistry and Biology, Evolutionary Biology &

8    Systematic Zoology, 14476 Potsdam, Germany

9    [2] University of Potsdam, Institute of Biochemistry and Biology, Evolution and Adaptive

10   Genomics, 14476 Potsdam, Germany

11   [3] current address: Novartis BTDM Mengeš, Kolodvorska 27, SI-1234 Mengeš, Slovenia (All

12   work in regards to this manuscript was performed at University of Potsdam[1])

13   [4] Swedish Museum of Natural History, SE-104 05 Stockholm, Sweden

14

15   **\*Corresponding author**: Prof. Dr. Ralph Tiedemann tiedeman@uni-potsdam.de

16

17

18  **Abstract**

19

20  The harbour porpoise (*Phocoena phocoena*) is a highly mobile cetacean found in waters

21  across the Northern hemisphere. It occurs in coastal water and inhabits water basins that vary

22  broadly in salinity, temperature, and food availability. These diverse habitats could drive

23  differentiation among populations; population structure within the north Atlantic (north of 51°

24  latitude) is not fully resolved, particularly in relation to Baltic Sea populations. Here we report

25  the first harbour porpoise genome, assembled *de novo* from a Swedish Kattegat individual.

26  The genome is one of the most complete cetacean genomes currently available, with a total

27  size of 2.7 Gb, and 50% of the total length found in just 34 scaffolds. Using the largest

28  scaffolds, we were able to examine chromosome-level rearrangements relative to the genome

29  of the closest related species available, domestic cattle (*Bos taurus*). The draft annotation

30  comprises 22,154 predicted gene models, which we further annotated through matches to

31  NCBI nucleotide database, GO categorization, and motif prediction. To infer the adaptive

32  abilities of this species, as well as their population history, we performed Bayesian skyline

33  analysis of the genome, which is concordant with the demographic history of this species,

34  including expansion and fragmentation events. Overall, this genome assembly, together with

35  the draft annotation, represents a crucial addition to the limited genetic markers currently

36  available for the study of porpoise and cetacean conservation, phylogeny, and evolution.

2

## Introduction

As an apex predator, the harbour porpoise (*Phocoena phocoena*) is a key indicator for conservation and biodiversity measurements in the Nordic Seas (Hooker & Gerber, 2004; Lawrence et al., 2016; Sergio et al., 2008). Marine mammals in particular face many threats from their environment (Fietz et al., 2013; Godard-Codding et al., 2011) including noise pollution (Dyndo et al., 2015; Nabe-Nielsen et al., 2014), marine debris and by-catch (Scheidat et al., 2008; Unger et al., 2017), predation by grey seals (Leopold et al., 2014), and infectious diseases (Siebert et al., 2001; van Beurden et al., 2017). These threats impact structure, boundaries, and stability of populations. This is especially true in the Kattegat/Baltic Sea area, where broad ecological shifts have occurred on a relatively short time scale. Since forming 15,000 years ago, the Baltic has undergone periods of brackish, marine, and completely fresh water, and encountered increasing and continuous humans impacts including eutrophication, pollution and overharvesting (Korpinen et al., 2012; Paasche et al., 2015; Ukkonen et al., 2014; Varjopuro et al., 2014). This geological history has created a series of challenges to marine species, and has likely fostered local adaption and population differentiation.

Harbour porpoises are the most abundant costal cetaceans across their wide distribution from sub-polar to temperate waters in the Northern hemisphere (Fontaine et al., 2017; Gaskin, 1984). As one of the smallest marine mammals, they belong to the *Delphinoida* and are the sister group to the *Monodontidae* (Gatesy et al., 2013; Geisler et al., 2011; Hassanin et al., 2012). Three subspecies of harbour porpoise, *P. p. vomerina* (North Pacific), *P. p. relicta* (Black Sea) and *P. p. phocoena* (North Atlantic), can be differentiated genetically (Rosel et al., 1999), but also by morphological traits including body size and diet (Fontaine et al., 2017; Galatius et al., 2012).

3

61  The population size of harbour porpoises in European Atlantic Shelf waters is estimated to be

62  375,000 with shifts across the last decade in the exact regions they occupy (e.g. in the North

63  Sea; (Hammond et al., 2013). Estimates of population size in the western Baltic Sea are

64  smaller, approximately 40,000 animals (Benke et al., 2014; Scheidat et al., 2008; Viquerat et

65  al., 2014). The Baltic Sea proper population, which is not included in the former surveys, has

66  very low estimates (below 500 individuals; Amundin, 2016) and is considered critically

67  endangered (Benke et al., 2014; Hammond et al., 2008; Scheidat et al., 2008).

68      As with other marine mammals in the Northern Atlantic, e.g. grey and harbour seals,

69  subpopulations of the harbor porpoise arose during the end of the last glacial period as North

70  Sea populations recolonized the Baltic Sea (Fietz et al., 2016). Now these different

71  populations show shifts in habitat use based largely on food availability (Hammond et al.,

72  2013) and activity patterns (Nuuttila et al., 2017), and display fine scale morphological and

73  genetic differences (Fontaine et al., 2012, 2014; Wiemann et al., 2010) and significant

74  isolation by distance (Lah et al., 2016). Recent studies based on morphometric and genetic

75  data suggest that different ecotypes of harbour porpoise in the North Atlantic and Baltic Sea

76  exist and may need further conservation measures (Fontaine et al., 2014, 2017; Galatius et al.,

77  2012).

78      These fine scale differences in morphology and behavior may constitute local

79  adaptation, yet the genes underlying such a potentially adaptive differentiation are still

80  unknown and would be best investigated on a whole-genome scale. To examine this, there is a

81  need for high quality genomic resources for this species. A genome will also allow for a

82  broader investigation of population structure, demographic history, functional, and

83  evolutionary questions, as has been shown for other cetacean species in recent studies (Foote

84  et al., 2016; Keane et al., 2015; Nery et al., 2013; Sun et al., 2013; Yim et al., 2013; Zhou et

85    al., 2013). To this end, a full genome will enable mapping of so far anonymous nuclear

86    microsatellite (Wiemann et al. 2010) and SNP (Lah et al. 2016) loci, thus facilitating

87    population genomic inference.

88        We present here the first *de novo* assembly of the full genome of the harbor porpoise,

89    scaffolded with *in vitro* proximity ligation data (hereafter "Chicago" library), and draft-

90    annotated to predict its coding proteins and their functions (Deposited at NCBI as BioProject:

91    PRJNA417595 with BioSample-ID: SAMN08000480). We demonstrate chromosome-level

92    homology with other Cetartiodactyla (Gatesy et al., 2013), and insight into past population

93    dynamics using a Bayesian skyline plot (Li & Durbin, 2011).

94

95

96    **Materials and Methods**

97    *DNA sampling*

98    Tissue for whole genome sequencing came from a single individual from the Kattegat

99    (Glommen - Falkenberg), Sweden (ID: C2009/02665). Muscle tissue was sampled in July

100   2009 from a by-caught female of probably young age (22.4kg, 110.5m), frozen, and

101   transported to Potsdam, Germany for DNA extraction. Sample preparation and Genomic DNA

102   isolation were performed following the Quiagen DNeasy Blood & Tissue Kit (Cat 69506,

103   Hilden). Successful high molecular weight DNA-isolation was confirmed by Sanger

104   sequencing of the mitochondrial control region, and visualization of fragment sizes of the

105   entire extraction using the Tape Station (Agilent 2200, Santa Clara, CA 95051). By mtDNA

106   sequencing, we verified that the analyzed specimen carried haplotype PHO7 (Tiedemann et

107   al., 1996), indicative of the separated Beltsea population of the Kattegat/Western Baltic Sea

108   region (Lah et al., 2016; Wiemann et al., 2010).

5

109   *Genome sequencing and assembly*

110   The draft *de novo* assembly was constructed from two libraries (insert sizes ca. 300 and ca.

111   500bp); sequenced in 125bp PE on the Illumina HiSeq 2500 at EUROFINS Genomics. Reads

112   were trimmed using CUTADAPT v1.10 (Martin, 2011) and an initial assembly was made using

113   SOAPDENOVO2 (Luo et al., 2015). DNA from the same sample was used by Dovetail

114   Genomics for construction of a Chicago library (Putnam et al., 2016), and sequenced in 150bp

115   PE reads on an Illumina NextSeq500 at the University of Potsdam. The draft assembly was

116   then scaffolded with the Chicago library results for the final HiRise assembly, performed by

117   Dovetail Genomics.

118        Presence of core, single copy, and orthologous genes was measured using CEGMA and

119   BUSCO, run in the genome mode for the Laurasiatheria database (Simão et al., 2015).

120   BLOBTOOLS was run to examine potential contaminants, based on divergence in GC-content

121   and read coverage variation across the assembly (Laetsch & Blaxter, 2017).

122

123   *Genome annotation*

124   Genome annotation was performed by MAKER2 (Holt & Yandell, 2011) in two steps. MAKER2

125   makes use of different programs and draws from several lines of evidence. Prior to

126   annotation, repetitive elements were soft-masked with REPEATMASKER (Smit et al., 2013-

127   2015) using the te_protein repeat database (Smith et al., 2007). In the first MAKER2 run, three

128   gene predictors were used: SNAP (Bromberg et al., 2008) was *ab initio* trained with the

129   CEGMA results (Parra et al., 2007), GENEMARK-ES (Ter-Hovhannisyan et al., 2008) was run

130   using an HMM produced by *ab initio* training on the whole *P. phocoena* genome, and

131   AUGUSTUS was run using the presets for human, as is recommended for vertebrates (Stanke et

132   al., 2004). Protein sequences, supplied as evidence were obtained from the complete

6

133    SwissProt database (553,941 Proteins) plus NCBI entries of 184,527 proteins from eight

134    different cetacean groups (On 20 March 2017, all hits to following keywords:

135    "*Balaenopteridae*", "*Lipotes vexillifer*", "*Neophocaena*", "*Orcinus orca*", "*Phocoena*",

136    "*Physeter catodon*", "*Pontoporia blainvillei*", "*Tursiops truncatus*").

137         For the second MAKER2 run, we created a new SNAP-HMM based on the first MAKER2

138    output, and ran it with the same parameters as the first run, exchanging only the SNAP HMM

139    and excluding the protein evidence. The resulting CDS predictions were extracted from the

140    final gff file, which was created by *fathom* implemented in SNAP (Bromberg et al., 2008).

141    These gene predictions were further verified by a BLASTN search against the entire GENBANK

142    non-redundant nucleotide sequence database (date downloaded 21.07.2017). Summary

143    statistics were generated using GENOME ANNOTATION GENERATOR (Hall et al., 2014). We then

144    used all CDS and their BLAST results in BLAST2GO (Goetz et al., 2008) to identify conserved

145    protein domains with INTERPROSCAN (including a Pfam comparison). We functional

146    annotated the CDS with GO terms, which are a controlled vocabulary to describe gene

147    function constantly actualized by the Gene Ontology Consortium (Ashburner et al., 2000;

148    Carbon et al., 2017).

149

150    *Comparative genomics*

151    The closest relative with a chromosome-level assembly currently available is the domestic

152    cattle, *Bos taurus*. To validate our assembly, we compared our scaffolds to the *B. taurus*

153    chromosomes (assembly UMD 3.1.1 downloaded from NCBI, ACCESSION

154    DAAA00000000). Specifically, the 122 *P. phocoena* scaffolds of at least 1Mbp were aligned

155    to the *B. taurus* chromosomes using the nucmer software of the MUMMER package v. 3.23

156    (Kurtz et al., 2004). From the coordinates of these alignments, runs of ten or more

7

157   consecutive matches of each at least 250bp between a given *P. phocoena* scaffold and a *B.*

158   *taurus* chromosome were extracted using custom perl scripts. Their start and end positions

159   were used to generate a CIRCOS (http://circos.ca/) plot that shows regions of collinearity as

160   well as rearrangements. For the CIRCOS plot, separate ribbons are displayed between a *B.*

161   *taurus* chromosome and a *P. phocoena* scaffold for consecutive hits that were each no more

162   than 20,000 bp apart. If a hit is more than 20,000bp from the next run of consecutive hits, a

163   new ribbon was started; in total 24,394 separate ribbons were constructed (Figure 1).

164

165   *Population genomics*

166   In using genome-wide diploid sequence data it is possible to reconstruct the population

167   history in estimating population sizes through the past (Li et al., 2011). To estimate the

168   demographic history of the individual sequenced, we used the SNP Frequency spectra based

169   on our genome assembly, which is a haploid sequence, and the PE reads used to construct the

170   *de novo* assembly, prior to Chicago scaffolding (described above, we used both insert sizes).

171   These reads were first mapped back to the final assembly using BWA (Li & Durbin, 2009).

172   SNP data was extracted from the resulting bam files, and variants were extracted using

173   SAMTOOLS VS.1.6. (Li, Handsaker, et al., 2009), and BCFTOOLS (Li, Handsaker, et al., 2009),

174   implemented with the script vcfutils.pl (Li, Handsaker, et al., 2009). This generated a final

175   *.fq.gz file, which was then used to generate the final Bayesian skyline plot in the PSMC

176   package, using perl scripts psmc2history.pl and psmc_plot.pl (Li et al., 2011). The parameters

177   of the PSMC analysis were set following the recommendation from the authors (Li & Durbin,

178   2011, https://github.com/lh3/psmc) and we applied a generation time of 10 years (Birkun Jr. &

179   Frantzis, 2008) and a mutation rate of $2.2 \times 10^{-9}$ year/site (Taylor et al., 2007).

180

8

181 **Results**

182

183 De novo *assembly of the* P. phocoena *genome*

184 Shotgun sequencing produced a total number of 1,268M reads (Table 1), these were used to

185 generate a draft assembly with 2.4M scaffolds and an N50 of 33.1kb. This assembly was

186 combined with the Chicago library data (556M read) for final scaffolding by Dovetail

187 Genomics (Putnam et al., 2016). The final HiRise assembly from Dovetail contains ca. 2M

188 scaffolds (Table 2) and has a total length of 2.7Gb (N50 of 23.8Mb). The greatest

189 improvements from the addition of the Chicago libraries is in building up the 34 longest

190 scaffolds, which make up approximately half of the entire assembly (Table 2). The CIRCOS

191 plot illustrates the near-completeness of these long scaffolds. We observe almost complete

192 coverage of the cow chromosomes by scaffolds bigger than 1Mb in our assembly (Figure 1).

193 The BUSCO and CEGMA analyses also suggests that we have largely reconstructed the entire

194 genome, and identified 96.9% (91.3% complete) of the 2,586 Eukaryotic and 94.2 (88.7%

195 complete) of the 6,253 Laurasiatheria BUSCO core genes and 90% of the 248 ultra-conserved

196 CEGs (54% complete).

197

198 *Genome completeness and annotation*

199 The MAKER2 annotation resulted in the prediction of 22,154 coding genes (Table 3). In total

200 21,750 CDS had a BLAST hit against the nucleotide database, which accounts for 98% of the

201 total CDSs. Of these BLAST hits, 99% account for vertebrate, and these were dominated

202 (90%) by hits to Cetacea (thereof 59% *Tursiops truncatus*, 27% *Orcinus orca*). Further

203 annotation with INTERPROSCAN revealed 250,126 features of these predicted proteins. These

204 comprise hits in several protein domain databases, e.g. 23,319 PFAM protein domains, 37,046

9

205  PANTHER gene families, 24,538 SUPERFAMILY annotations and 31,114 GENE3D domains.

206  Assignment of the BLAST results to Gene Ontology (GO) categories resulted in 55,143 hits

207  across the GO categories (

208  Figure *2*).

209

210  *Inference of Kattegat/Baltic population history*

211  We inferred the population history of the harbour porpoise *P. phocoena* based on one single

212  individual (Li et al., 2011) using the PSMC algorithm, which combines all generated PE read

213  data generated. Between eight and four million years ago the inferred population size ($N_e$) was

214  low, around 10,000 individuals (Figure 3). It began to increase slightly at 3Myr, and rose

215  more rapidly around 2Myr, reaching an $N_e$ of 45,000 during the following 1.5 Myr. The

216  estimated population size peaked approximates 400kyra before it dropped to a quarter of the

217  original size around 100kyrs ago, leading to a very low $N_e$, similar to that seen in present day

218  populations (Hammond et al., 2013).

219

220  **Discussion**

221  We present here a high quality *de novo* genome assembly for the harbour porpoise *Pho-*

222  *coena phocoena*. With a GC-content of 41.4% and a total length of 2.7 GB, this assembly is

223  comparable to other high quality genomes (Groenen et al., 2013; Zimin et al., 2009). BUSCO

224  and CEGMA gene scans support a near completeness of core genes in the assembly, and sup-

225  port that we have largely reconstructed the entire genome. For almost completely covering the

226  chromosomes of the *B. taurus* genome (Figure 1), only 122 scaffolds are needed, including

227  the 34 largest scaffolds representing 50% of the whole genome. Of these largest scaffolds

228  some completely match single *B. taurus* chromosomes, e.g., chromosome 25. Other *B. taurus*

10

229    chromosomes are in only 2-3 pieces in our scaffolds, e.g. chromosomes 12, 24. Based on this

230    comparison, we infer that our assembly represents a nearly complete genome of *P. phocoena*,

231    and that our largest scaffolds are nearly-complete chromosomes. The CIRCOS plot also illus-

232    trates chromosomal rearrangements between domestic cattle and the harbour porpoise, two

233    species diverged approximately 60Myrs ago within the Cetartiodactyla (Gatesy et al., 2013).

234    These chromosomal rearrangements are seen several times among distinct lineage of Cetarti-

235    odactyla (Avila et al., 2015; Kulemzina et al., 2009, 2011; Pauciullo et al., 2014), e.g., com-

236    parison between camel, pig and domestic cattle (Balmus et al., 2007).

237        The number of annotated genes (22,154) is comparable to other published cetacean ge-

238    nomes: 21,459 bottlenose dolphin (Lindblad-Toh et al., 2011), 20,605 minke whale (Yim et

239    al., 2013), 22,711 grey whale (DeWoody et al., 2017). They appear to broadly span key func-

240    tional gene categories, e.g. biological processes, cellular components and molecular function,

241    both across the annotated GO terms and the INTERPROSCAN analysis. With this information

242    we can directly search for known, respectively key genes, for further investigations, e.g. se-

243    lection or adaptive traits.

244        The harbor porpoise is estimated to have split from is closest relative ca. 5Myr ago

245    (Gatesy et al., 2013). Interestingly our Bayesian skyline plot (Figure 3) coincides with this

246    date by starting a population expansion around that time point. Around 4.5 Myr ago an

247    expansion occurred, during which time the North Atlantic is known to have cooled, leading to

248    an extinction of 65% of the marine organisms (Stanley, 1995). The harbour porpoise is well

249    known in subarctic regions and some populations (e.g. Greenland) occur in areas which freeze

250    to a large extent during winter (Tolley & Rosel, 2006). Therefore, an extinction of other

251    marine species during a cold water period does not preclude that the harbour porpoise could

252    increase its population size and expand through the Atlantic. During the last interglacial

11

253   period, Eemian, the inferred $N_e$ remained relatively high at around 50,000 individuals before,

254   dropping dramatically with the beginning of the last glacial period 100kya. When comparing

255   this pattern to the demographic history of other cetaceans, it is most similar to the bottlenose

256   dolphin (*Tursiops truncatus*), a related species with a similar North Atlantic distribution

257   (Brüniche-Olsen et al., n.d.; Foote et al., 2016; Yim et al., 2013; Zhou et al., 2013). The newly

258   forming sea ice areas, around 400kya ago, could have led to fragmentation of different

259   populations, and therefore lead to a drop in regional total effective population size in regards

260   to our sample. A potential low population size we see postulated for today would fit to the

261   history of the Baltic Sea and the population status of *P. phocoena* (Johannesson et al., 2011;

262   Johannesson & André, 2006; Ukkonen et al., 2014). Specifically, there is strong evidence for a

263   Western Baltic/Kattegat (i.e., Beltsea) population separated from the North Sea/North Atlantic

264   (Hammond et al., 2013; Lah et al., 2016), which currently counts approximately 40,000

265   animals (Benke et al., 2014; Scheidat et al., 2008; Viquerat et al., 2014). Our sequenced

266   specimen was assigned with high likelihood to this Beltsea population by mtDNA analysis

267   (exhibiting haplotype PHO 7; cf. Tiedemann et al., 1996; Wiemann et al., 2010).

268   In this study we present the first whole genome assembly and annotation of the harbour

269   porpoise, at this point the most complete assembly for the Family *Phocoenidae*. This genome

270   adds to the Cetacean genome collection by supplying important resources for further

271   investigation within the *Odontoceti* as well as outside the *Cetacea*. This will provide an

272   invaluable resource for further genetic studies within the harbour porpoise itself, both as a

273   resource for whole-genome investigations into population structure and to identify key genes

274   associated with local adaptation. This genome represents a crucial genetic resource for further

275   investigation in the population genetics and phylogeny on other species of the *Phocoenidae*

276   including the currently most rare marine mammal, the almost extinct Vaquita (*Phocoena*

277   *sinus*) (Taylor et al., 2017), and is hence especially important for conservation efforts.

278   **Acknowledgments**

286

287   **References**

288
289   Amundin, M. (2016). SAMBAH - Static Acoustic Monitoring of the Baltic Sea Harbour
290      porpoise. LIFE Project Number LIFE 08 NAT/S/000261 European Commission, 77pp.
291      available at http://www.sambah.org/SAMBAH-Final-Report-FINAL-for-website-April-
292      2017.pdf.
293   Ashburner, M., Ball, C. A., Blake, J. A., Botstein, D., Butler, H., Cherry, J. M., … Sherlock,
294      G. (2000). Gene ontologie: Tool for the unification of biology. *Nature Genetics*, *25*(1),
295      25–29. doi:10.1038/75556.Gene
296   Avila, F., Baily, M. P., Merriwether, D. A., Trifonov, V. A., Rubes, J., Kutzler, M. A., …
297      Raudsepp, T. (2015). A cytogenetic and comparative map of camelid chromosome 36
298      and the minute in alpacas. *Chromosome Research*, *23*(2), 237–251. doi:10.1007/s10577-
299      014-9463-3
300   Balmus, G., Trifonov, V. A., Biltueva, L. S., O'Brien, P. C. M., Alkalaeva, E. S., Fu, B., …
301      Ferguson-Smith, M. A. (2007). Cross-species chromosome painting among camel, cattle,
302      pig and human: Further insights into the putative Cetartiodactyla ancestral karyotype.
303      *Chromosome Research*, *15*(4), 499–515. doi:10.1007/s10577-007-1154-x
304   Benke, H., Bräger, S., Dähne, M., Gallus, A., Hansen, S., Honnef, C. G., … Verfuß, U. K.
305      (2014). Baltic Sea harbour porpoise populations: Status and conservation needs derived
306      from recent survey results. *Marine Ecology Progress Series*, *495*, 275–290.
307      doi:10.3354/meps10538
308   Birkun Jr., A. A., & Frantzis, A. (2008). *Phocoena phocoena ssp. relicta*. The IUCN Red List
309      of Threatened Species 2008: e.T17030A6737111. Downloaded on 12 October 2017.
310   Bromberg, Y., Yachdav, G., & Rost, B. (2008). SNAP predicts effect of mutations on protein
311      function. *Bioinformatics*, *24*(20), 2397–2398. doi:10.1093/bioinformatics/btn435
312   Brüniche-Olsen, A., Westerman, R., Kazmierczyk, Z., Vertyankin, V. V., Godard-Codding, C.,

313   Bickham, J. W., & DeWoody, J. A. (n.d.). The inference of gray whale (Eschrichtius
314       robustus) population attributes from whole-genome sequences.
315   Carbon, S., Dietze, H., Lewis, S. E., Mungall, C. J., Munoz-Torres, M. C., Basu, S., …
316       Westerfield, M. (2017). Expansion of the gene ontology knowledgebase and resources:
317       The gene ontology consortium. *Nucleic Acids Research*, *45*(D1), D331–D338.
318       doi:10.1093/nar/gkw1108
319   DeWoody, J. A., Fernandez, N. B., Brüniche-Olsen, A., Antonides, J. D., Doyle, J. M., San
320       Miguel, P., … Bickham, J. (2017). Characterization of the gray whale (Eschrichtius
321       robustus) genome and a genotyping array based on single nucleotide polymorphisms in
322       candidate genes. *Biological Bulletin*, *232*(June), 186–197.
323   Dyndo, M., Wiśniewska, D. M., Rojano-Doñate, L., & Madsen, P. T. (2015). Harbour
324       porpoises react to low levels of high frequency vessel noise. *Scientific Reports*, *5*, 11083.
325       doi:10.1038/srep11083
326   Fietz, K., Galatius, A., Teilmann, J., Dietz, R., Frie, A. K., Klimova, A., … Olsen, M. T.
327       (2016). Shift of grey seal subspecies boundaries in response to climate, culling and
328       conservation. *Molecular Ecology*, *25*(17), 4097–4112. doi:10.1111/mec.13748
329   Fietz, K., Graves, J. A., & Olsen, M. T. (2013). Control Control Control: A Reassessment and
330       Comparison of GenBank and Chromatogram mtDNA Sequence Variation in Baltic Grey
331       Seals (*Halichoerus grypus*). *PLoS ONE*, *8*(8), 1–7. doi:10.1371/journal.pone.0072853
332   Fontaine, M. C., Roland, K., Calves, I., Austerlitz, F., Palstra, F. P., Tolley, K. A., … Aguilar,
333       A. (2014). Postglacial climate changes and rise of three ecotypes of harbour porpoises,
334       *Phocoena phocoena*, in western Palearctic waters. *Molecular Ecology*, *23*(13), 3306–
335       3321. doi:10.1111/mec.12817
336   Fontaine, M. C., Snirc, A., Frantzis, A., Koutrakis, E., Oztürk, B., Oztürk, A. a, & Austerlitz,
337       F. (2012). History of expansion and anthropogenic collapse in a top marine predator of
338       the Black Sea estimated from genetic data. *Proceedings of the National Academy of
339       Sciences of the USA*, *109*(38), E2569-76. doi:10.1073/pnas.1201258109
340   Fontaine, M. C., Thatcher, O., Ray, N., Piry, S., Brownlow, A., Davison, N. J., … Goodman,
341       S. J. (2017). Mixing of porpoise ecotypes in southwestern UK waters revealed by genetic
342       profiling. *Royal Society Open Science*, *4*(3), 160992. doi:10.1098/rsos.160992
343   Foote, A. D., Vijay, N., Ávila-Arcos, M. C., Baird, R. W., Durban, J. W., Fumagalli, M., …
344       Wolf, J. B. W. (2016). Genome-culture coevolution promotes rapid divergence of killer
345       whale ecotypes. *Nature Communications*, *7*(May), 11693. doi:10.1038/ncomms11693
346   Galatius, A., Kinze, C. C., & Teilmann, J. (2012). Population structure of harbour porpoises in
347       the Baltic region: evidence of separation based on geometric morphometric comparisons.
348       *Journal of the Marine Biological Association of the United Kingdom*, *92*(8), 1669–1676.
349       doi:10.1017/S0025315412000513
350   Gaskin, D. (1984). The harbour porpoise *Phocoena phocoena* (L.): regional populations,
351       status, and infromation on direct and indirect catches. *Reports of the International
352       Whaling Commission*, *34*, 569–586.
353   Gatesy, J., Geisler, J. H., Chang, J., Buell, C., Berta, A., Meredith, R. W., … McGowen, M. R.
354       (2013). A phylogenetic blueprint for a modern whale. *Molecular Phylogenetics and
355       Evolution*, *66*(2), 479–506. doi:10.1016/j.ympev.2012.10.012
356   Geisler, J. H., McGowen, M. R., Yang, G., & Gatesy, J. (2011). A supermatrix analysis of
357       genomic, morphological, and paleontological data from crown Cetacea. *BMC
358       Evolutionary Biology*, *11*(1), 112. doi:10.1186/1471-2148-11-112
359   Godard-Codding, C. A. J., Clark, R., Fossi, M. C., Marsili, L., Maltese, S., West, A. G., …
360       Stegeman, J. J. (2011). Pacific ocean-wide profile of CYP1A1 expression, stable carbon

14

361   and nitrogen isotope ratios, and organic contaminant burden in sperm whale skin
362   biopsies. *Environmental Health Perspectives*, *119*(3), 337–343.
363   doi:10.1289/ehp.0901809
364   Goetz, S., Garccia-Gomez, M. J., Terol, J., Williams, T. D., Nagaraj, S. H., Nueda, M. J., …
365   Conesa, A. (2008). High-throughput functional annotation and data mining with the
366   Blast2GO suite. *Nucleic Acids Research*, *36*(10), 3420–3435. doi:10.1093/nar/gkn176
367   Groenen, M. A. M., Archibald, A. L., Uenishi, H., Tuggle, C. K., Takeuchi, Y., Rothschild, M.
368   F., … Hunt, T. (2013). Analyses of pig genomes provide insight into porcine demography
369   and evolution. *Nature*, *491*(7424), 393–398. doi:10.1038/nature11622
370   Hall, B., DeRego, T., & Geib, S. (2014). GAG: the Genome Annotation Generator (Version
371   1.0) [Software]. doi:Available from http://genomeannotation.github.io/GAG
372   Hammond, P. S., Bearzi, G., Bjørge, A., Forney, K. A., Karczmarski, L., Kasuya, T., …
373   Wilson, B. (2008). *Phocoena phocoena* (Baltic Sea subpopulation). (errata version
374   published in 2016) The IUCN Red List of Threatened Species 2008:
375   e.T17031A98831650. Downloaded on 10 October 2017. Retrieved from
376   http://www.iucnredlist.org/details/17031/0
377   Hammond, P. S., Macleod, K., Berggren, P., Borchers, D. L., Burt, L., Canadas, A., …
378   Vazquez, J. A. (2013). Cetacean abundance and distribution in European Atlantic shelf
379   waters to inform conservation and management. *Biological Conservation*, *164*, 107–122.
380   doi:10.1016/j.biocon.2013.04.010
381   Hassanin, A., Delsuc, F., Ropiquet, A., Hammer, C., Jansen Van Vuuren, B., Matthee, C., …
382   Couloux, A. (2012). Pattern and timing of diversification of Cetartiodactyla (Mammalia,
383   Laurasiatheria), as revealed by a comprehensive analysis of mitochondrial genomes.
384   *Comptes Rendus - Biologies*, *335*(1), 32–50. doi:10.1016/j.crvi.2011.11.002
385   Holt, C., & Yandell, M. (2011). MAKER2: an annotation pipeline and genome-database
386   management tool for second-generation genome projects. *BMC Bioinformatics*, *12*(1),
387   491. doi:10.1186/1471-2105-12-491
388   Hooker, S. K., & Gerber, L. (2004). Marine Reserves as a Tool for Ecosystem-Based
389   Management : The Potential Importance of Megafauna. *BioScience*, *54*(1), 27–39.
390   doi:10.1641/0006-3568(2004)054[0027:MRAATF]2.0.CO;2
391   Johannesson, K., & André, C. (2006). Life on the margin: Genetic isolation and diversity loss
392   in a peripheral marine ecosystem, the Baltic Sea. *Molecular Ecology*, *15*(8), 2013–2029.
393   doi:10.1111/j.1365-294X.2006.02919.x
394   Johannesson, K., Smolarz, K., Grahn, M., & André, C. (2011). The future of baltic sea
395   populations: Local extinction or evolutionary rescue? *Ambio*, *40*(2), 179–190.
396   doi:10.1007/s13280-010-0129-x
397   Keane, M., Semeiks, J., Webb, A. E., Li, Y. I., Quesada, V., Craig, T., … deMagalhães, J. P.
398   (2015). Insights into the evolution of longevity from the bowhead whale genome. *Cell*
399   *Reports*, *10*(1), 112–122. doi:10.1016/j.celrep.2014.12.008
400   Korpinen, S., Meski, L., Andersen, J. H., & Laamanen, M. (2012). Human pressures and their
401   potential impact on the Baltic Sea ecosystem. *Ecological Indicators*, *15*(1), 105–114.
402   doi:10.1016/j.ecolind.2011.09.023
403   Kulemzina, A. I., Trifonov, V. A., Perelman, P. L., Rubtsova, N. V., Volobuev, V., Ferguson-
404   Smith, M. A., … Graphodatsky, A. S. (2009). Cross-species chromosome painting in
405   Cetartiodactyla: Reconstructing the karyotype evolution in key phylogenetic lineages.
406   *Chromosome Research*, *17*(3), 419–436. doi:10.1007/s10577-009-9032-3
407   Kulemzina, A. I., Yang, F., Trifonov, V. A., Ryder, O. A., Ferguson-Smith, M. A., &
408   Graphodatsky, A. S. (2011). Chromosome painting in Tragulidae facilitates the

15

409    reconstruction of Ruminantia ancestral karyotype. *Chromosome Research*, *19*(4), 531–
410    539. doi:10.1007/s10577-011-9201-z
411 Kurtz, S., Phillippy, A., Delcher, A. L., Smoot, M., Shumway, M., Antonescu, C., & Salzberg,
412    S. L. (2004). Versatile and open software for comparing large genomes. *Genome*
413    *Biology*, *5*(2), R12. doi:10.1186/gb-2004-5-2-r12
414 Laetsch, D. R., & Blaxter, M. L. (2017). BlobTools: Interrogation of genome assemblies.
415    *F1000Research*, *6*(1287), 1–15. doi:10.12688/f1000research.12232.1 and
416    doi:10.5281/zenodo.845347
417 Lah, L., Trense, D., Benke, H., Berggren, P., Gunnlaugsson, Þ., Lockyer, C., … Tiedemann,
418    R. (2016). Spatially Explicit Analysis of Genome-Wide SNPs Detects Subtle Population
419    Structure in a Mobile Marine Mammal, the Harbor Porpoise. *PLoS ONE*, *11*(10), 1–23.
420    doi:10.1371/journal.pone.0162792
421 Lawrence, J. M., Armstrong, E., Gordon, J., Lusseau, S. M., & Fernandes, P. G. (2016).
422    Passive and active, predator and prey: using acoustics to study interactions between
423    cetaceans and forage fis. *ICES Journal of Marine Science*, *73*(8), 2075–2084.
424    doi:10.1093/icesjms/fsw013
425 Leopold, M. F., Begeman, L., van Bleijswijk, J. D. L., IJsseldijk, L. L., Witte, H. J., & Gröne,
426    A. (2014). Exposing the grey seal as a major predator of harbour porpoises. *Proc. R. Soc.*
427    *B*, *282*, 20142429. doi:10.1098/rspb.2014.2429
428 Li, H., & Durbin, R. (2009). Fast and accurate short read alignment with Burrows-Wheeler
429    transform. *Bioinformatics*, *25*(14), 1754–1760. doi:10.1093/bioinformatics/btp324
430 Li, H., & Durbin, R. (2011). Inference of human population history from individual whole-
431    genome sequences. *Nature*, *475*(7357), 493–496. doi:10.1038/nature10231
432 Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., … Durbin, R. (2009).
433    The Sequence Alignment/Map format and SAMtools. *Bioinformatics*, *25*(16), 2078–
434    2079. doi:10.1093/bioinformatics/btp352
435 Lindblad-Toh, K., Garber, M., Zuk, O., Lin, M. F., Parker, B. J., Washietl, S., … Kellis, M.
436    (2011). A high-resolution map of human evolutionary constraint using 29 mammals.
437    *Nature*, *478*(7370), 476–482. doi:10.1038/nature10530
438 Luo, R., Liu, B., Xie, Y., Li, Z., Huang, W., Yuan, J., … Wang, J. (2015). Erratum:
439    SOAPdenovo2: an empirically improved memory-efficient short-read de novo assembler.
440    *GigaScience*, *4*(1), 30. doi:10.1186/s13742-015-0069-2
441 Martin, M. (2011). Cutadapt removes adapter sequences from high-throughput sequencing
442    reads. *EMBnet.journal*, *17*(1), 10. doi:10.14806/ej.17.1.200
443 Nabe-Nielsen, J., Sibly, R. M., Tougaard, J., Teilmann, J., & Sveegaard, S. (2014). Effects of
444    noise and by-catch on a Danish harbour porpoise population. *Ecological Modelling*, *272*,
445    242–251. doi:10.1016/j.ecolmodel.2013.09.025
446 Nery, M. F., Gonzalez, D. J., & Opazo, J. C. (2013). How to Make a Dolphin: Molecular
447    Signature of Positive Selection in Cetacean Genome. *PLoS ONE*, *8*(6), 2–8.
448    doi:10.1371/journal.pone.0065491
449 Nuuttila, H. K., Courtene-Jones, W., Baulch, S., Simon, M., & Evans, P. G. H. (2017). Don't
450    forget the porpoise: acoustic monitoring reveals fine scale temporal variation between
451    bottlenose dolphin and harbour porpoise in Cardigan Bay SAC. *Marine Biology*, *164*(3),
452    1–16. doi:10.1007/s00227-017-3081-5
453 Paasche, Ø., Österblom, H., Neuenfeldt, S., Bonsdorff, E., Brander, K., Conley, D. J., …
454    Stenseth, N. C. (2015). Connecting the Seas of Norden. *Nature Climate Change*, *5*(2),
455    89–92. doi:10.1038/nclimate2471
456 Parra, G., Bradnam, K., & Korf, I. (2007). CEGMA: A pipeline to accurately annotate core

457          genes in eukaryotic genomes. *Bioinformatics*, *23*(9), 1061–1067.
458          doi:10.1093/bioinformatics/btm071

459 Pauciullo, A., Perucatti, A., Cosenza, G., Iannuzzi, A., Incarnato, D., Genualdo, V., …
460          Iannuzzi, L. (2014). Sequential cross-species chromosome painting among river buffalo,
461          cattle, sheep and goat: A useful tool for chromosome abnormalities diagnosis within the
462          family bovidae. *PLoS ONE*, *9*(10). doi:10.1371/journal.pone.0110297

463 Putnam, N. H., Connell, B. O., Stites, J. C., Rice, B. J., Hartley, P. D., Sugnet, C. W., …
464          Rokhsar, D. S. (2016). Chromosome-scale shotgun assembly using an in vitro method for
465          long-range linkage. *Genome Research*, *26*, 342–350. doi:10.1101/gr.193474.115

466 Rosel, P. E., Tiedemann, R., & Walton, M. (1999). Genetic evidence for limited trans-Atlantic
467          movements of the harbor porpoise, *Phocoena phocoena*. *Marine Biology (Berlin)*,
468          *133*(4), 583–591.

469 Scheidat, M., Gilles, A., Kock, K. H., & Siebert, U. (2008). Harbour porpoise *Phocoena*
470          *phocoena* abundance in the southwestern Baltic Sea. *Endangered Species Research*,
471          *5*(2–3), 215–223. doi:10.3354/esr00161

472 Sergio, F., Caro, T., Brown, D., Clucas, B., Hunter, J., Ketchum, J., … Hiraldo, F. (2008). Top
473          Predators as Conservation Tools: Ecological Rationale, Assumptions, and Efficacy.
474          *Annual Review of Ecology, Evolution, and Systematics*, *39*(1), 1–19.
475          doi:10.1146/annurev.ecolsys.39.110707.173545

476 Siebert, U., Wünschmann, A., Weiss, R., Frank, H., Benke, H., & Frese, K. (2001). Post-
477          mortem findings in harbour porpoises (*phocoena phocoena*) from the German North and
478          Baltic Seas. *Journal of Comparative Pathology*, *124*(2–3), 102–114.
479          doi:10.1053/jcpa.2000.0436

480 Simão, F. A., Waterhouse, R. M., Ioannidis, P., Kriventseva, E. V., & Zdobnov, E. M. (2015).
481          BUSCO: Assessing genome assembly and annotation completeness with single-copy
482          orthologs. *Bioinformatics*, *31*(19), 3210–3212. doi:10.1093/bioinformatics/btv351

483 Smit, A. F. A., Hubley, R., & Green, P. (n.d.). RepeatMasker Open-4.0.
484          <http://www.repeatmasker.org>.

485 Smith, C. D., Edgar, R. C., Yandell, M. D., Smith, D. R., Celniker, S. E., Myers, E. W., &
486          Karpen, G. H. (2007). Improved repeat identification and masking in Dipterans. *Gene*,
487          *389*(1), 1–9. doi:10.1016/j.gene.2006.09.011

488 Stanke, M., Steinkamp, R., Waack, S., & Morgenstern, B. (2004). AUGUSTUS: A web server
489          for gene finding in eukaryotes. *Nucleic Acids Research*, *32*(Web Server issue), W309–
490          W312. doi:10.1093/nar/gkh379

491 Stanley, S. M. (1995). 7 Neogene Ice Age in the North Atlantic Region: Climatic Changes,
492          Biotic Effects, and Forcing Factors. In *Effects of Past Global Change on Life*.
493          Washington (DC): National Academies: National Research Council (US) Panel on
494          Effects of Past Global Change on Life.

495 Sun, Y. B., Zhou, W. P., Liu, H. Q., Irwin, D. M., Shen, Y. Y., & Zhang, Y. P. (2013). Genome-
496          wide scans for candidate genes involved in the aquatic adaptation of dolphins. *Genome*
497          *Biology and Evolution*, *5*(1), 130–139. doi:10.1093/gbe/evs123

498 Taylor, B. L., Chivers, S. J., Larese, J., & Perrin, W. F. (2007). Generation length and percent
499          mature estimates for IUCN assessments of cetaceans. *Administrative Report LJ-07-01*
500          *National Marine Fisheries*, 24. doi:10.1.1.530.4789

501 Taylor, B. L., Rojas-Bracho, L., Moore, J., Jaramillo-Legorreta, A., Ver Hoef, J. M.,
502          Cardenas-Hinojosa, G., … Hammond, P. S. (2017). Extinction is Imminent for Mexico's
503          Endemic Porpoise Unless Fishery Bycatch is Eliminated. *Conservation Letters*, *10*(5),
504          588–595. doi:10.1111/conl.12331

505 Ter-Hovhannisyan, V., Lomsadze, A., Chernoff, Y. O., & Borodovsky, M. (2008). Gene
506     prediction in novel fungal genomes using an ab initio algorithm with unsupervised
507     training, 1979–1990. doi:10.1101/gr.081612.108

508 Tiedemann, R., Harder, J., Gmeiner, C., & Haase, E. (1996). Mitochondrial DNA sequence
509     patterns of Harbour porpoises (*Phocoena phocoena*) from the North and the Baltic Sea.
510     *Zeitschrift Für Säugetierkunde*, *61*, 104–111.

511 Tolley, K., & Rosel, P. (2006). Population structure and historical demography of eastern
512     North Atlantic harbour porpoises inferred through mtDNA sequences. *Marine Ecology
513     Progress Series*, *327*, 297–308. doi:10.3354/meps327297

514 Ukkonen, P., Aaris-Sorensen, K., Arppe, L., Daugnora, L., Halkka, A., Lougas, L., … Stora, J.
515     (2014). An Arctic seal in temperate waters: History of the ringed seal (*Pusa hispida*) in
516     the Baltic Sea and its adaptation to the changing environment. *The Holocene*, *24*(12),
517     1694–1706. doi:10.1177/0959683614551226

518 Unger, B., Herr, H., Benke, H., Böhmert, M., Burkhardt-Holm, P., Dähne, M., … Siebert, U.
519     (2017). Marine debris in harbour porpoises and seals from German waters. *Marine
520     Environmental Research*, 1–8. doi:10.1016/j.marenvres.2017.07.009

521 van Beurden, S. J., Ijsseldijk, L. L., van de Bildt, M. W. G., Begeman, L., Wellehan, J. F. X.,
522     Waltzek, T. B., … Penzes, J. J. (2017). A novel cetacean adenovirus in stranded harbour
523     porpoises from the North Sea: detection and molecular characterization. *Archives of
524     Virology*, *162*(7), 2035–2040. doi:10.1007/s00705-017-3310-8

525 Varjopuro, R., Andrulewicz, E., Blenckner, T., Dolch, T., Heiskanen, A. S., Pihlajamäki,
526     M., … Psuty, I. (2014). Coping with persistent environmental problems: Systemic delays
527     in reducing eutrophication of the Baltic Sea. *Ecology and Society*, *19*(4).
528     doi:10.5751/ES-06938-190448

529 Viquerat, S., Herr, H., Gilles, A., Peschko, V., Siebert, U., Sveegaard, S., & Teilmann, J.
530     (2014). Abundance of harbour porpoises (*Phocoena phocoena*) in the western Baltic,
531     Belt Seas and Kattegat. *Marine Biology*, *161*(4), 745–754. doi:10.1007/s00227-013-
532     2374-6

533 Wiemann, A., Andersen, L. W., Berggren, P., Siebert, U., Benke, H., Teilmann, J., …
534     Tiedemann, R. (2010). Mitochondrial Control Region and microsatellite analyses on
535     harbour porpoise (*Phocoena phocoena*) unravel population differentiation in the Baltic
536     Sea and adjacent waters. *Conservation Genetics*, *11*(1), 195–211. doi:10.1007/s10592-
537     009-0023-x

538 Yim, H.-S., Cho, Y. S., Guang, X., Kang, S. G., Jeong, J.-Y., Cha, S.-S., … Lee, J.-H. (2013).
539     Minke whale genome and aquatic adaptation in cetaceans. *Nature Genetics*, *46*(1), 88–
540     92. doi:10.1038/ng.2835

541 Zhou, X., Sun, F., Xu, S., Fan, G., Zhu, K., Liu, X., … Yang, G. (2013). Baiji genomes reveal
542     low genetic variability and new insights into secondary aquatic adaptations. *Nature
543     Communications*, *4*(2708), 1–6. doi:10.1038/ncomms3708

544 Zimin, A. V, Delcher, A. L., Florea, L., Kelley, D. R., Schatz, M. C., Puiu, D., … Salzberg, S.
545     L. (2009). A whole-genome assembly of the domestic cow, *Bos taurus*. *Genome
546     Biology*, *10*(4), R42. doi:10.1186/gb-2009-10-4-r42

547
548

549 **Data Accessibility**

550 The genome assembly, finale genome sequence and the draft annotation are deposit on NCBI

551 under BioProject-ID: PRJNA417595 and BioSample-ID: SAMN08000480).

552

553 **Authors Contributions**

554 R.T. and L.L. designed the study; A.R. provided the sample and associated biological

555 information. L.L. performed molecular lab work, S.H. performed initial *de novo* assembly,

556 M.A. executed all genome annotations and analyses, M.A., S.H., and, A.B.D. analyzed and

557 interpreted the results, M.A. wrote the manuscript. All authors edited and approved the final

558 manuscript.

559

560

**Tables and Figures**

**Table 1** Sequencing statistics of libraries used for the two assemblies.

|  | Insert (bp) | n reads | Coverage |
|---|---|---|---|
| Illumina Library 300 | 300 | 794 M | 74.1 X |
| Illumina Library 500 | 500 | 473 M | 44.2 X |
| Chicago library | 1-50kb | 528 M | 87.2 X |

**Table 2** Assembly statistics of the harbour porpoise genome.

|  | SOAPdenovo assembly | Dovetail HiRise Assembly |
|---|---|---|
| Total length | 2,669.6 Mb | 2,681.2 Mb |
| *Scaffolds* | | |
| N50 (number/length) | 23,685 / 0.032Mb | 34 / 23.8Mb |
| N90 (number) | 159,889 | 43,146 |
| Longest | 304,733 | 67,078,619 |
| Number | 2,139,681 | 2,025,248 |

**Table 3** Genome annotation statistics

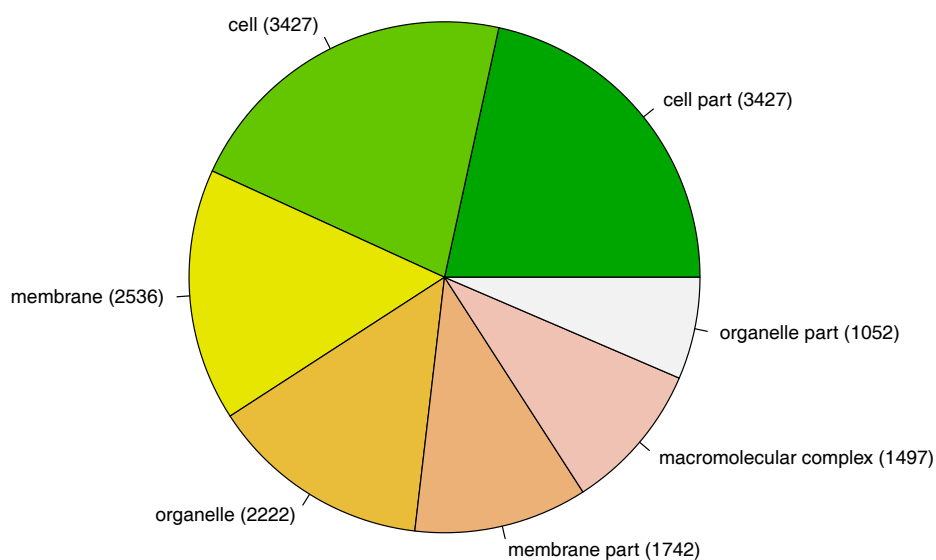|  | Exons | Introns | Genes | CDS |
|---|---|---|---|---|
| Number | 171,735 | 149,581 | 22,154 | 22,154 |
| Longest | 10,332 | 1,532,135 | 2,299,565 | 20,613 |
| Mean length | 166 | 3,966 | 28,051 | 1,282 |
| % genome covered by | - | - | 23,2% | 1.1% |
| GC% in CDS | - | - | - | 54.48% |

20

**Figure 1** CIRCOS plot of harbour porpoise scaffolds (at least 1Mbp, black on the left outer rim) mapped against the cow (*Bos taurus*) chromosomes (colored bars, labeled X and 1-29).

567   The *B. taurus* autosomal chromosomes (1-29) as well as the X chromosome (x) are shown in

568   different colors. The 122 largest *P. phocoena* scaffolds with consecutive MUMMER hits be-

569   tween a *B. taurus* chromosome of at least 250bp that are no more than 20,000bp apart are

570   shown in black. Matches between *B. taurus* chromosomes and *P. phocoena* scaffolds are

571   shown in the color of the *B. taurus* chromosomes.


**Figure 2** GO-Terms; GO annotation level 2, separated by gene ontology terms: "Cellular Component", "Biological Process", and "Molecular Function". Separate categories are listed, with the number of hits in parentheses.
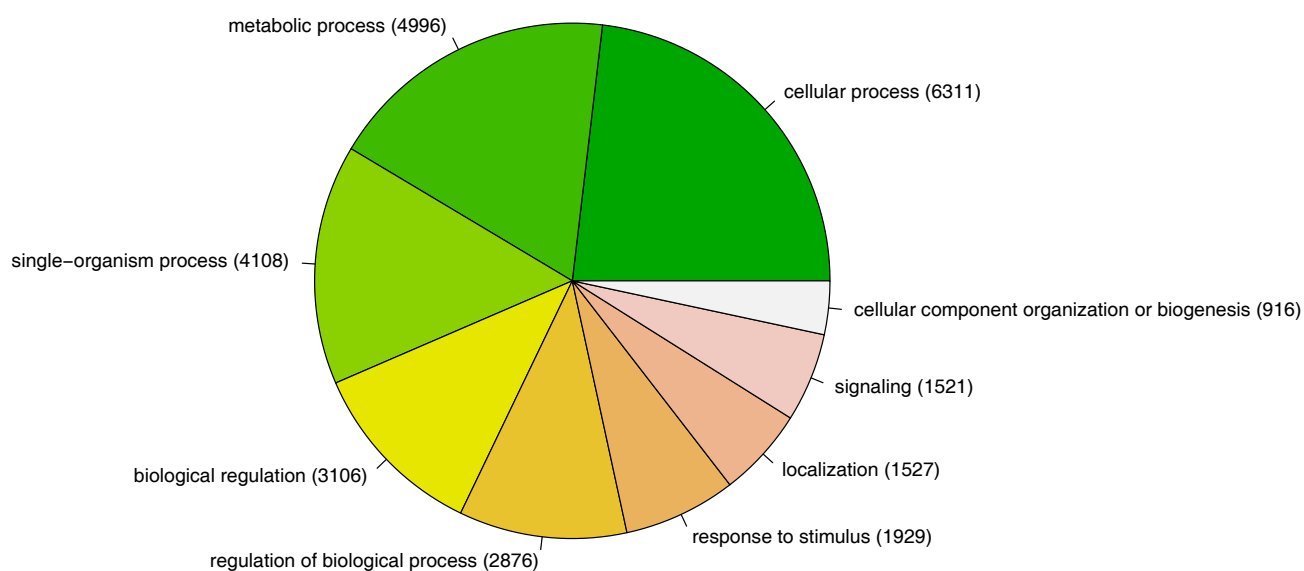

**Figure 3** PSMC estimated harbour porpoise population size changes over time for the Baltic Sea. g = generation time; μ = mutation rate (per site, per year). Porpoise data generated on the basis of mapping PE reads to whole genome scaffolds during SNP calling.

**Cellular Component**



**Biological Process**



**Molecular Function**