1 **Evaluating the role of reference-genome phylogenetic distance on evolutionary inference**

2

3 Aparna Prasad[1], Eline D Lorenzen[1], Michael V Westbury[1*]

4

5    1. GLOBE Institute, University of Copenhagen, Øster Voldgade 5-7, Copenhagen,
6        Denmark

7

8 Corresponding author: m.westbury@sund.ku.dk

9

10 **Abstract**

11

12      When a high-quality genome assembly of a target species is unavailable, an option to

13 avoid the costly *de novo* assembly process is a mapping-based assembly. However, mapping

14 shotgun data to a distant relative may lead to biased or erroneous evolutionary inference.

15 Here, we used short-read data from a mammal and a bird species (beluga and rowi kiwi) to

16 evaluate whether reference genome phylogenetic distance can impact downstream

17 demographic (PSMC) and genetic diversity (heterozygosity, runs of homozygosity) analyses.

18 We mapped to assemblies of species of varying phylogenetic distance (conspecific to

19 genome-wide divergence of >7%), and *de novo* assemblies created using cross-species

20 scaffolding. We show that while reference genome phylogenetic distance has an impact on

21 demographic analyses, it is not pronounced until using a reference genome with >3%

22 divergence from the target species. When mapping to cross-species scaffolded assemblies, we

23 are unable to replicate the original beluga demographic analyses, but can with the rowi kiwi,

24 presumably reflecting the more fragmented nature of the beluga assemblies. As for genetic

25 diversity estimates, we find that increased phylogenetic distance has a pronounced impact;

26 heterozygosity estimates deviate incrementally as phylogenetic distance increases. Moreover,

27 runs of homozygosity are removed when mapping to any non-conspecific assembly.

28 However, these biases can be reduced when mapping to a cross-species scaffolded assembly.

29    Taken together, our results show that caution should be exercised when selecting the

30    reference genome for mapping assemblies. Cross-species scaffolding may offer a way to

31    avoid a costly, traditional *de novo* assembly, while still producing robust, evolutionary

32    inference.

33

34    **Introduction**

35

36        The large extent of genetic information within the nuclear genome enables powerful

37    evolutionary inferences using just a single individual. Two options are available for genome

38    assembly: mapping-based assemblies using a closely-related species as reference, or *de novo*

39    assemblies. In the former approach, relatively little time and monetary expense is invested in

40    sequencing one individual to high coverage (>20x). After assembly, it is possible to make

41    population-wide evolutionary inferences of the target species, including levels of genetic

42    diversity and inbreeding, adaptive genomic changes, and demographic history (Barnett et al.,

43    2020; Lord et al., 2020; Michael V. Westbury, Petersen, Garde, Heide-Jørgensen, &

44    Lorenzen, 2019).

45

46        Although mapping-based assemblies are less costly than *de novo* assemblies, there are

47    some caveats. Biases towards the reference genome allele may be an issue when analysing

48    population-level datasets. Such errors can arise during variant calling, when the alternative

49    allele fails to be called altogether, or when heterozygous sites are incorrectly called as

50    homozygous for the reference allele (Brandt et al., 2015; Ros-Freixedes et al., 2018).

51    Although such issues are known to occur when mapping to a conspecific from a different

52    population, biases caused by mapping to phylogenetically more distant taxa have only

53    somewhat been addressed (Armstrong et al., 2020; M. V. Westbury et al., 2021). Problems

54    with correctly identifying variants may arise due to decreased mapping efficiency as

55    reference-genome phylogenetic distance increases (Shapiro & Hofreiter, 2014). However, the

56    consequences of this on downstream analyses have yet to be comprehensively assessed. This

57    leads to the question of whether the potentially costly *de novo* assembly process can be

58    avoided when assemblies from phylogenetically more distant species are available. Insights

59    into this will be especially important for the study of extinct species, where a conspecific

60    reference genome is unlikely to be available (Barnett et al., 2020; Palkopoulou et al., 2018).

61

62    Here, we investigate the influence of the reference genome's phylogenetic distance to

63    the target species on downstream evolutionary analyses. Specifically, we focused on estimates

64    of (i) demographic history (reconstructed using PSMC), (ii) genetic diversity (genome-wide

65    heterozygosity estimated using ANGSD, ROHan, SAMtools/BCFtools), (iii) inbreeding/runs

66    of homozygosity (using ROHan). Additionally, we investigated whether biases can be

67    overcome by using cross-species scaffolded con-specific assemblies as mapping reference.

68    We applied our methodology to two taxonomically disparate datasets; one based on mammals

69    (beluga and incrementally divergent cetacean species), and one based on birds (rowi kiwi and

70    incrementally divergent paleognath species). We selected these datasets based on the

71    assumption that mammal and bird genomes may respond differently to reference biases,

72    allowing for more generalised conclusions. Furthermore, while beluga whales are a relatively

73    abundant species (Hobbs et al., 2019), rowi kiwi are threatened with extinction and have

74    much lower population numbers (Robertson et al., 2017), which may play a role.

75

76

77   **Materials and methods**

78

79        A simplified version of the methodologies implemented in this manuscript are

80   presented in figure 1.

81

82   **Data**

83        For the cetacean comparative dataset, we downloaded the raw Illumina reads and an

84   assembled genome of the beluga (*Delphinapterus leucas*, Genbank accession code:

85   GCF_002288925.2, SRA code: SRR5197961). In addition, we downloaded genome

86   assemblies for five other cetacean species with varying genomic distance to the beluga (Table

87   1): narwhal (*Monodon monoceros*, Genbank accession code: GCF_005190385.1), narrow-

88   ridged finless porpoise (*Neophocaena asiaeorientalis*, Genbank accession code:

89   GCF_003031525.1), bottlenose dolphin (*Tursiops truncatus*, Genbank accession code:

90   GCF_001922835.1), sperm whale (*Physeter macrocephalus*, Genbank accession code:

91   GCA_002837175.2), and minke whale (*Balaenoptera acutorostrata*, Genbank accession

92   code: GCF_000493695.1). Assembly length, N50, and level of missing data for each

93   assembly are listed in supplementary table S1.

94

95        For the paleognath comparative dataset, we downloaded raw reads and an assembled

96   genome of the rowi kiwi (Genbank accession: GCF_003343035.1, SRA accession:

97   SRR6918118). We also downloaded published assemblies for three palaeognathae species of

98   varying genomic distance to the rowi kiwi (Table 1): North Island brown kiwi (termed brown

99   kiwi here, *A. mantelli*, Genbank accession: GCF_001039765.1), great spotted kiwi (termed

100  spotted kiwi here, *A. haastii*, Genbank accession: GCA_003342985.1), and emu (*Dromaius*

4

101    *novaehollandiae,* Genbank accession: GCF_003342905.1). Assembly length, N50, and levels

102    of missing data for each assembly are listed in supplementary table S2.

103

104    **De novo and mapping assemblies**

105        All mappings were performed following the same procedure for the beluga/cetacean

106    and rowi kiwi/palaeognathae species datasets. We trimmed adapter sequences and removed

107    reads shorter than 30bp from the raw reads using skewer v0.2.2 (Jiang, Lei, Ding, & Zhu,

108    2014)), and mapped the trimmed reads using BWA v0.7.15 (Li & Durbin, 2009) and the mem

109    algorithm utilising default parameters. We parsed the output and removed duplicates and

110    reads with a mapping quality lower than 30 with SAMtools v1.6 (Li et al., 2009).

111

112        For the beluga, we mapped the beluga short-read data to (i) published assemblies

113    (beluga, narwhal, finless porpoise, bottlenose dolphin, sperm whale, minke whale), (ii) a *de*

114    *novo* contig-level beluga assembly constructed for the purposes of this study, and (iii) five *de*

115    *novo* beluga assemblies produced by cross-species scaffolding the contig-level assembly with

116    each published non-beluga assembly using in-silico mate pair (MP) libraries.

117

118        For the rowi kiwi, we mapped the rowi short-read data to (i) published assemblies

119    (rowi kiwi, brown kiwi, spotted kiwi, emu), (ii) a *de novo* contig-level rowi kiwi assembly

120    constructed for the purposes of this study, and (iii) three *de novo* rowi kiwi assemblies

121    produced by cross-species scaffolding the contig-level assembly with each of the non-rowi

122    kiwi published assembles using in-silico MP libraries.

123

124    We constructed *de novo* contig-level assemblies for both the beluga and the rowi kiwi

125    by first performing an error-correction step on the adapter-trimmed reads in tadpole from the

126    BBtools toolsuite (Bushnell, 2014) and a kmer size of 31. We constructed the *de novo*

127    assembly with the error-corrected reads using SOAPdenovo2 pregraph and contig (Luo et al.,

128    2012), specifying a kmer size of 51, and otherwise using default parameters.

129

130    To construct the cross-species scaffolded assemblies we used the contig-level

131    assemblies from above and scaffolded them either five times independently in the case of the

132    beluga, or three times independently in the case of the rowi kiwi. For this, we constructed in-

133    silico MP libraries using a modified version of the cross-species scaffolding pipeline (Grau,

134    Hackl, Koepfli, & Hofreiter, 2018) and repeat-masked versions of the published non-beluga

135    assemblies (narwhal, finless porpoise, bottlenose dolphin, sperm whale, minke whale) and

136    published non-rowi assemblies (brown kiwi, spotted kiwi, emu). Repeats were masked based

137    on the Genbank annotations.

138

139    In short, we constructed a fasta consensus sequence using a consensus base-call

140    approach (-doFasta 2) in ANGSD v0.931 (Korneliussen, Albrechtsen, & Nielsen, 2014) and a

141    minimum read depth of 3 (-minInddepth 3), minimum mapping quality of 25 (-minmapq 25),

142    minimum base quality of 25 (-minq 25), and only considered reads that mapped uniquely to

143    one location (-uniqueonly 1). We converted this fasta sequence into a pseudo-fastq sequence

144    with a quality score of 40 for all covered bases using BBtools (Bushnell, 2014), as input for

145    the seq-scripts pipeline.

146

147     From the consensus sequences, we constructed in-silico MP libraries with approximate

148     insert sizes of 1kb, 2kb, 3kb, 5kb, 8kb, 10kb, 15kb, and 20kb using seq-scripts

149     (https://github.com/thackl/seq-scripts), specifying a read length of 150bp and read depth of

150     100x. We scaffolded the contig-level beluga and rowi kiwi assemblies using SOAPdenovo2

151     map and scaff. To reduce the chances of mis-assembly by the longer in-silico MP libraries, we

152     specified MP libraries of different insert sizes as different ranks in the SOAP config file. The

153     shortest insert sizes had higher rankings; if a longer-insert library contradicted the shorter

154     inserts, they were not used for scaffolding.

155

156     We assessed the final assembly quality in the form of contiguity (N50) and amount of

157     missing data using QUAST v4.5 (Gurevich, Saveliev, Vyahhi, & Tesler, 2013).

158

159     **Sex-scaffold filtering**

160     To exclude sex-linked scaffolds in downstream demographic and heterozygosity

161     analyses, we determined the scaffolds that likely originated from the sex-chromosomes for

162     each of the scaffolded assemblies (published assemblies and the cross-species scaffolded

163     assemblies) used in this study. We found putative sex-chromosome scaffolds in the cetacean

164     species by aligning all assemblies to the Cow X (Genbank accession: CM008168.2) and

165     Human Y (Genbank accession: NC_000024.10) chromosomes, and putative sex-chromosome

166     scaffolds in the palaeognath species by aligning all assemblies to the chicken (*Gallus gallus*)

167     W (Genbank accession: CM000121.5) and Z (Genbank accession: CM000122.5)

168     chromosomes. We did this using satsuma synteny v2.1 (Grabherr et al., 2010) and default

169     parameters. We also removed scaffolds smaller than 10 kilobase pairs (kb) from all

170     downstream analyses.

7

171 **Divergence estimates**

172       To ensure comparability between the divergence estimates of our datasets, we

173 calculated the autosome-wide divergence of our species of interest, either beluga or rowi kiwi,

174 to the other species included in the study. To do this, we downloaded the raw reads for all the

175 species (Supplementary table S3) and mapped them back to either the published beluga

176 assembly or rowi kiwi assembly. We calculated the pairwise distance between species from

177 the resultant mapped bam files and a consensus base call approach in ANGSD (-doIBS 2),

178 and specifying the following parameters: -minq 25 -minmapq 25 -minind 4 -setMinDepthInd

179 5 -uniqueonly 1 -docounts 1 make a distance matrix (-makematrix 1), and only including

180 autosomal scaffolds over 10kb in length (-rf).

181

182 **Demographic reconstruction**

183       To determine the influence of (a) phylogenetic distance of the reference genome to the

184 target species, (b) reference genome contiguity, and (c) the utility of cross-species scaffolded

185 reference genomes on demographic reconstruction, we ran a Pairwise Sequentially Markovian

186 Coalescent model (PSMC) (Li & Durbin, 2011) on each diploid genome, resulting in a total of

187 twelve replicates for the beluga dataset and eight for the rowi kiwi dataset. We called diploid

188 genome sequences using SAMtools and BCFtools v1.6 (Narasimhan et al., 2016), specifying

189 a minimum quality score of 20 and minimum coverage of 10.

190

191       We ran PSMC specifying atomic intervals 4+25*2+4+6. Beluga PSMC outputs were

192 plotted using a generation time of 32 years (Garde et al., 2015) and mutation rate of 1.65e-08

193 (Michael V. Westbury et al., 2019). To plot the rowi kiwi PSMC outputs, we calculated a

194 mutation rate using the pairwise distance of the rowi kiwi to the brown kiwi (0.003123) and

195 the formula pairwise distance x 2 / divergence time. We used a divergence time ~3.8 Ma (De

196 Cahsan & Westbury, 2020) which resulted in a mutation rate of $1.64 \times 10^{-9}$ per year or $4.1 \times 10^{-8}$

197 per generation assuming a generation time of 25 years (Weir, Haddrath, Robertson,

198 Colbourne, & Baker, 2016).

199

200 **Genetic diversity**

201 To determine the influence of (a) phylogenetic distance of the reference genome to the

202 target species, (b) reference genome contiguity, and (c) the utility of cross-species scaffolded

203 reference genomes on genetic diversity estimates, we estimated the autosome-wide

204 heterozygosity of the beluga mapped to all twelve cetacean assemblies and the rowi kiwi

205 mapped to all eight paleognath assemblies.

206

207 We calculated heterozygosity for each of the datasets using ANGSD. We estimated

208 autosomal heterozygosity using allele frequencies (-doSaf 1), taking genotype likelihoods into

209 account with the GATK algorithm (-GL 2), and specifying the following filters: only include

210 sites with a read depth of at least 5 (-mininddepth 5), minimum mapping and base qualities of

211 30 (-minmapq 30, -minq 30), only include reads mapping uniquely to one location (-

212 uniqueonly 1), only include reads where both read pairs map (-only_proper_pairs 1), only

213 include autosomal scaffolds (-rf), and the extended adjust quality scores around indels

214 parameter (-baq 2). Heterozygosity was computed from the output of this using realSFS from

215 the ANGSD toolsuite, specifying 20 megabase pairs (Mb) windows of covered sites (-nSites).

216

217 We subsequently tested whether the results are consistent regardless of parameter and

218 software selection using the beluga dataset mapped to the six published assemblies. We

219  assessed the influence of parameter selection in ANGSD on heterozygosity estimates by

220  computing heterozygosity using the procedure outlined above, but replacing the 'extended

221  adjust quality scores around indels' parameter (-baq 2), with (i) adjust quality scores around

222  insertion/deletions (indels) (-baq 1), (ii) no indel quality score adjustment (-baq 0), or (iii)

223  extended adjust quality scores around indels (-baq 2), and adjust quality for reads with

224  multiple mismatches to the reference (-C 50).

225

226          To assess whether software choice can impact results, we used two additional methods

227  to compute heterozygosity of the beluga mapped to the six published assemblies: ROHan

228  (Renaud, Hanghøj, Korneliussen, Willerslev, & Orlando, 2019), and the PSMC input diploid

229  file (SAMtools/BCFtools).

230

231          In ROHan, we used default parameters to calculate autosome-wide levels of

232  heterozygosity and runs of homozygosity (ROH). The default parameters specify a 1 MB

233  window as being a ROH, if the window has an average heterozygosity of less than 1e-5. To

234  calculate average autosome-wide heterozygosity from the diploid file used for the PSMC

235  analysis, we used seqtk comp (https://github.com/lh3/seqtk).

236

237  **Inbreeding (runs of homozygosity)**

238          As ROHan simultaneously outputs runs of homozygosity as well as autosome-wide

239  levels of heterozygosity, we could evaluate how reference genome phylogenetic distance

240  influences perceived inbreeding estimates using ROH. We did not retrieve any significant

241  ROH in the beluga dataset using ROHan, so we were unable to investigate this further using

242    the beluga data. We repeated the above ROHan analysis using the rowi kiwi mapped to all

243    eight assemblies (both published and cross-species scaffolded).

244

245    **Results**

246

247    **Mapping**

248        Mapping results of the beluga raw reads to each cetacean assembly can be found in

249    supplementary tables S4 and S5. Mapping results of the rowi kiwi raw reads to each

250    paleognath assembly can be found in supplementary tables S6 and S7. As phylogenetic

251    distance to the reference genome increases, there is a general trend of a decreasing number of

252    unique reads mapping. This trend is not seen when mapping to the five and three cross-

253    species scaffolded assemblies for beluga and rowi kiwi, respectively. However, less reads

254    map to these assemblies than to the conspecific assemblies.

255

256    **Cross-species scaffolded *de novo* assemblies**

257        Our contig-level beluga assembly had an N50 of ~3.5 kb. The cross-species scaffolded

258    assemblies were more contiguous, with N50s ranging from 283 kb to 614 kb (Supplementary

259    table S8). However, these assemblies also had a lot of introduced missing data (16% - 18%).

260    In comparison, the original assemblies for each species had N50s of 6.3 Mb - 122.2 Mb, and

261    missing data rates of 0.5% - 6% (Supplementary table S1).

262

263        Our contig-level rowi kiwi assembly had an N50 ~6.6 kb. The cross-species

264    scaffolded assemblies were more contiguous, with N50s ranging from 1.9 Mb - 4.9 Mb

265    (Supplementary table S9). These assemblies also had large amounts of introduced missing

11

266    data (10% - 21%). However, this was comparable to the brown kiwi assembly with 14% data.

267    Assembly contiguities were also more comparable to the published assemblies, which had

268    N50s of 1.4 Mb - 5.7 Mb (Supplementary table S2).

269

270    **Demographic reconstruction**

271        **Beluga -** With increasing phylogenetic distance of the reference genome, we see an

272    incremental increase in deviation from the pattern obtained when mapping to the published

273    beluga assembly (Fig 2A). However, we do not see an incremental change when using our

274    five cross-species scaffolded assemblies as reference. Instead we see that all newly assembled

275    genomes produce the same PSMC output. However, this output differs from the pattern

276    obtained when mapping to the published beluga assembly (Fig 2B).

277

278        When comparing PSMC results produced by mapping to the published beluga

279    assembly, and by mapping to our *de novo* contig-level beluga assembly, we see a pattern of

280    increase in $N_e$ ~500 thousand years ago (kya) followed by a decrease ~150 kya. This is

281    consistent between both assemblies. However, the values of effective population size ($N_e$) are

282    much lower when mapping to the *de novo* contig-level assembly (Supplementary fig S1).

283

284        **Rowi kiwi -** Unlike the beluga, the PSMC results of the rowi kiwi were vastly

285    different when mapping to phylogenetic distant references compared to the published rowi

286    assembly (Supplementary fig S2). However, we do see the incremental change as

287    phylogenetic distance increases when mapping to the non-rowi assemblies. We investigated if

288    there was a problem with the published rowi assembly by reassembling it using the published

289    short-read and 3 kb mate-paired libraries (Sackton et al., 2019) with SOAPdenovo. Our

12

290    reassembled rowi kiwi genome was much less contiguous than the published version (0.3 Mb

291    vs 1.7 Mb) and had more missing data (12.7% vs. 1.6%). However, the PSMC produced when

292    mapping to this assembly was much more consistent with what we would have expected

293    based on the beluga results, and shows a demographic history similar to when mapping to the

294    assemblies from the other three non-rowi kiwi species (Fig 3A). Hence, we only considered

295    this re-assembly when assessing the inference of reference genome on demographic history

296    results in the rowi kiwi. The results produced after mapping to the cross-species scaffolded

297    rowi kiwi assemblies are much more similar to those from the re-assembled published rowi

298    kiwi assembly (Fig 3B).

299

300        When comparing PSMC results produced by mapping to the re-assembled rowi

301    genome, and by mapping to the contig-level assembly, we see similar general trajectories.

302    However, the values of effective population size (Ne) are much lower when mapping to the *de*

303    *novo* contig-level assembly as seen in the beluga (Supplementary fig S3).

304

305    **Genetic diversity**

306        When using ANGSD, ROHan, and SAMtools/BCFtools, we see a general trend of

307    increasing heterozygosity as reference genome phylogenetic distance increases. Which is also

308    consistent when applying the alternative ANGSD parameter sets (i) -baq 1 instead of -baq 2,

309    and (ii) -baq 0 instead of -baq 2 (Fig 4A,B, Supplementary figs S4 and S5, Supplementary

310    tables S10 and S11). In contrast, when using ANGSD parameter set (iii) adjusted for reads

311    with multiple mismatches to the reference (-C 50), we see a general trend of decreasing

312    heterozygosity levels as phylogenetic distance increases (Supplementary fig S6).

313

13

314     When using the cross-species scaffolded assemblies as reference genomes, we obtain

315     results more comparable to those obtained when using the published conspecific assemblies

316     as reference (Fig 4C,D). However, we do not see this when using SAMtools/BCFtools and the

317     beluga dataset. Instead, we observe a decrease in heterozygosity relative to the published

318     conspecific beluga assembly. The decrease is of a similar magnitude regardless of which

319     cetacean species was used for scaffolding (Supplementary table S12).

320

321     The quality of the assembly also appears to play a role; higher genome-wide

322     heterozygosity was estimated when mapping to the *de novo* contig-level beluga assembly

323     compared to a scaffolded assembly (Supplementary fig S7). This same pattern was also seen

324     when comparing the *de novo* contig-level rowi kiwi assembly to our reassembled version, but

325     not compared to the published assembly (Supplementary fig S8).

326

327     **Inbreeding**

328     When mapping the beluga reads to any reference genome (including the published

329     conspecific beluga assembly), we did not uncover any ROH. When running ROHan on the

330     rowi kiwi mapped to a published conspecific rowi kiwi assembly, we uncovered ROH, but not

331     when mapping to any of the non-rowi kiwi assemblies (Table 2).

332

333     **Discussion**

334

335     Through a detailed comparison of results produced after mapping to multiple reference

336     genomes from two unique datasets, we show that the choice of reference genome for mapping

337     of short-read data of a target species can and does impact downstream evolutionary

14

338    inferences. In general, as the phylogenetic distance of the reference genome increases, results

339    become incrementally less reliable with regards to demographic history, genetic diversity, and

340    inbreeding estimates.

341

342         With regards to demographic history analyses using PSMC, phylogenetic distance of

343    the reference genome to the target species did not appear to affect the overall trajectories, but

344    did result in relatively decreased $N_e$ estimates. However, this only became apparent when

345    using a reference genome more than 0.14% different to the target species (e.g. beluga vs

346    finless porpoise) (Figs 2A, 3A). Based on these results, if a conspecific assembly is not

347    available, using an assembly from a relatively closely-related species is unlikely to interfere

348    with the overall demographic trajectory.

349

350         In contrast with the demographic results, the bias that reference-genome selection

351    plays on genetic diversity estimates is more noticeable. Reference bias can cause

352    heterozygous sites to be incorrectly called as homozygous for the reference allele (Brandt et

353    al., 2015; Ros-Freixedes et al., 2018). However, we see a general increase in heterozygosity,

354    as opposed to the expected decrease (Fig 4). Therefore, misalignments may be a larger factor

355    in falsely calling heterozygous alleles as opposed to simply incorrect base calling. When

356    applying a strict filter that corrects for reads with multiple mismatches to the reference

357    genome (-C 50), it may be possible to eliminate increased heterozygosity due to

358    misalignments (Supplementary figure S6). However, this is still associated with issues, as we

359    observed a general decrease in heterozygosity due to putatively incorrect basecalls.

360

15

361     Phylogenetic distance driven reference bias was especially apparent when estimating

362     ROH (Table 2). When mapping to a non-conspecific reference genome, we observed a

363     complete loss of ROH, which would lead to the incorrect inference of no inbreeding in this

364     individual. As we show that global heterozygosity rates increase as phylogenetic distance of

365     the reference increases, this could artificially increase the heterozygosity level in ROH,

366     making the ROH no longer observable.

367

368     One method we investigated for its putative ability to overcome these biases, without

369     performing a traditional conspecific *de novo* assembly, is cross-species scaffolding. The

370     biggest attraction of a cross-species scaffolded assembly over a traditional conspecific *de*

371     *novo* assembly is that it only requires a single lane of Illumina sequencing, and an available

372     assembly from a closely-related species. However, at least in the case of the beluga dataset, it

373     can result in much more fragmented assemblies (Supplementary tables S8), and may therefore

374     not always be applicable, especially when highly contiguous assemblies are required (e.g. for

375     PSMC and ROH analyses). Nevertheless, using cross-species scaffolded assemblies as

376     reference resulted in relatively reliable PSMC and ROH results for the rowi kiwi (Fig 3, Table

377     2), as well as reliable genetic results in all comparisons when using ANGSD (Fig 4). The

378     reliability of these results may reflect that ANGSD uses genotype likelihoods to call

379     heterozygosity (Korneliussen et al., 2014), as opposed to direct genotype calls, and therefore

380     may be more reliable when heterozygous sites do not have the perfect near-50/50 allele ratios.

381

382     Despite the promising results when using cross-species scaffolded assemblies with our

383     rowi kiwi dataset, PSMC results were less reliable using the beluga dataset. This could result

384     from the low quality and highly-fragmented nature of the beluga cross-species scaffolded

16

385    assemblies (Supplementary table S8). The inability to produce contiguous scaffolds like the

386    rowi kiwi, may have arisen due to the highly fragmented nature of the beluga contig

387    assembly, with an N50 of only ~3.5 kb. To create a more contiguous final assembly, the N50

388    would need to be increased. Here, we implemented SOAPdenovo on a single library of

389    random insert sizes. The use of different software (Butler et al., 2008) and lab protocols

390    (Weisenfeld et al., 2014) to ensure the insert sizes are uniform may improve the contiguity of

391    the final assembly, and make results reliant on highly contiguous data, such as PSMC, more

392    reliable.

393

394         Although the assemblies are more fragmented, this does not mean they are completely

395    devoid of information for the PSMC analysis. Results using the cross-species scaffolded

396    assemblies still present the increase in $N_e$ ~500 kya and decrease ~150 kya seen when using

397    the published beluga assembly as reference, but with slightly decreased $N_e$ values (Fig 2).

398    Furthermore, when comparing PSMC results produced via mapping to the scaffold-level and

399    contig-level assemblies, we also see a similar pattern of population size change, but with

400    different values of $N_e$ (Supplementary figs S1 and S3). This suggests that contiguity may not

401    influence the pattern as much as the scale, and may still be useful for investigating relative

402    changes in $N_e$ rather than absolute values of $N_e$ itself.

403

404         Our analyses uncovered a potential problem with the published rowi kiwi assembly.

405    When comparing results mapped to the published assembly against non-conspecific

406    assemblies, cross-species scaffolded assemblies, and a reassembly of the published data, we

407    uncover large discrepancies in the results, especially in the PSMC results (Supplementary figs

408    S2 and S3). As our assemblies all used the same published raw data, we suspect that these

17

409   discrepancies resulted from miss-assemblies during the original *de novo* assembly process in

410   Allpaths-LG (Butler et al., 2008). Although outside of the scope of the present study, these

411   results show that caution should be exercised in reference genome selection for mapping

412   assemblies; if multiple assemblies are available, it may be beneficial to test robustness of

413   results against multiple reference genomes.

414

415   Taken together, our results show that demographic analyses of a single individual

416   mapped to a phylogenetically distant reference genome may be considered reliable with

417   regards to demographic trajectories (as in relative changes in $N_e$, rather than absolute values

418   of $N_e$). However, the phylogenetic distance of the reference genome can lead to

419   overestimation of heterozygosity and, in turn, underestimations of ROH. Finally, if no

420   assembly from a suitably closely related species is available as a mapping reference, cross-

421   species scaffolded assemblies appear to be a valid and likely more suitable option for

422   evolutionary inference.

423

424   **Acknowledgements**

427

428   **References**

429   Armstrong, E. E., Taylor, R. W., Miller, D. E., Kaelin, C. B., Barsh, G. S., Hadly, E. A., &

430   Petrov, D. (2020). Long live the king: chromosome-level assembly of the lion (*Panthera*

431   *leo*) using linked-read, Hi-C, and long-read data. *BMC Biology*, *18*(1), 3.

432    Barnett, R., Westbury, M. V., Sandoval-Velasco, M., Vieira, F. G., Jeon, S., Zazula, G., …

433        Gilbert, M. T. P. (2020). Genomic Adaptations and Evolutionary History of the Extinct

434        Scimitar-Toothed Cat, *Homotherium latidens*. *Current Biology: CB*, *30*, 1–8.

435    Brandt, D. Y. C., Aguiar, V. R. C., Bitarello, B. D., Nunes, K., Goudet, J., & Meyer, D.

436        (2015). Mapping Bias Overestimates Reference Allele Frequencies at the HLA Genes in

437        the 1000 Genomes Project Phase I Data. *G3* , *5*(5), 931–941.

438    Bushnell, B. (2014). BBTools software package. *URL Http://sourceforge.*

439        *Net/projects/bbmap*.

440    Butler, J., MacCallum, I., Kleber, M., Shlyakhter, I. A., Belmonte, M. K., Lander, E. S., …

441        Jaffe, D. B. (2008). ALLPATHS: de novo assembly of whole-genome shotgun

442        microreads. *Genome Research*, *18*(5), 810–820.

443    De Cahsan, B., & Westbury, M. V. (2020). Complete mitochondrial genomes offer insights

444        into the evolutionary relationships and comparative genetic diversity of New Zealand's

445        iconic kiwi (*Apteryx* spp.). *New Zealand Journal of Zoology*, 1–9.

446    Garde, E., Hansen, S. H., Ditlevsen, S., Tvermosegaard, K. B., Hansen, J., Harding, K. C., &

447        Heide-Jørgensen, M. P. (2015). Life history parameters of narwhals (*Monodon*

448        *monoceros*) from Greenland. *Journal of Mammalogy*, *96*(4), 866–879.

449    Grabherr, M. G., Russell, P., Meyer, M., Mauceli, E., Alföldi, J., Di Palma, F., & Lindblad-

450        Toh, K. (2010). Genome-wide synteny through highly sensitive sequence alignment:

451        Satsuma. *Bioinformatics* , *26*(9), 1145–1151.

452    Grau, J. H., Hackl, T., Koepfli, K.-P., & Hofreiter, M. (2018). Improving draft genome

453        contiguity with reference-derived in silico mate-pair libraries. *GigaScience*, *7*(5), giy029.

454    Gurevich, A., Saveliev, V., Vyahhi, N., & Tesler, G. (2013). QUAST: quality assessment tool

455        for genome assemblies. *Bioinformatics* , *29*(8), 1072–1075.

456   Hobbs, R. C., Reeves, R. R., Prewitt, J. S., Desportes, G., Breton-Honeyman, K., Christensen,

457        T., … Watt, C. A. (2019). Global Review of the Conservation Status of Monodontid

458        Stocks. *Marine Fisheries Review*, *81*, 1+.

459   Jiang, H., Lei, R., Ding, S.-W., & Zhu, S. (2014). Skewer: a fast and accurate adapter trimmer

460        for next-generation sequencing paired-end reads. *BMC Bioinformatics*, *15*, 182.

461   Korneliussen, T. S., Albrechtsen, A., & Nielsen, R. (2014). ANGSD: Analysis of Next

462        Generation Sequencing Data. *BMC Bioinformatics*, *15*, 356.

463   Li, H., & Durbin, R. (2009). Fast and accurate short read alignment with Burrows–Wheeler

464        transform. *Bioinformatics* , *25*(14), 1754–1760.

465   Li, H., & Durbin, R. (2011). Inference of human population history from individual whole-

466        genome sequences. *Nature*, *475*(7357), 493–496.

467   Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., … 1000 Genome

468        Project Data Processing Subgroup. (2009). The Sequence Alignment/Map format and

469        SAMtools. *Bioinformatics* , *25*(16), 2078–2079.

470   Lord, E., Dussex, N., Kierczak, M., Díez-Del-Molino, D., Ryder, O. A., Stanton, D. W. G., …

471        Dalén, L. (2020). Pre-extinction Demographic Stability and Genomic Signatures of

472        Adaptation in the Woolly Rhinoceros. *Current Biology: CB*, *30*(19), 3871–3879.e7.

473   Luo, R., Liu, B., Xie, Y., Li, Z., Huang, W., Yuan, J., … Wang, J. (2012). SOAPdenovo2: an

474        empirically improved memory-efficient short-read de novo assembler. *GigaScience*,

475        *1*(1), 18.

476   Narasimhan, V., Danecek, P., Scally, A., Xue, Y., Tyler-Smith, C., & Durbin, R. (2016).

477        BCFtools/RoH: a hidden Markov model approach for detecting autozygosity from next-

478        generation sequencing data. *Bioinformatics* , *32*(11), 1749–1751.

479   Palkopoulou, E., Lipson, M., Mallick, S., Nielsen, S., Rohland, N., Baleka, S., … Reich, D.

480  (2018). A comprehensive genomic history of extinct and living elephants. *Proceedings of*

481  *the National Academy of Sciences of the United States of America*, *115*(11), E2566–

482  E2574.

483  Renaud, G., Hanghøj, K., Korneliussen, T. S., Willerslev, E., & Orlando, L. (2019). Joint

484  Estimates of Heterozygosity and Runs of Homozygosity for Modern and Ancient

485  Samples. *Genetics*, *212*(3), 587–614.

486  Robertson, H. A., Baird, K., Dowding, J. E., Elliott, G. P., Hitchmough, R. A., Miskelly, C.

487  M., … Taylor, G. A. (2017). *Conservation status of New Zealand birds, 2016* (p. 23).

488  Department of Conservation, Wellington: New Zealand Threat Classification Series 19.

489  Ros-Freixedes, R., Battagin, M., Johnsson, M., Gorjanc, G., Mileham, A. J., Rounsley, S. D.,

490  & Hickey, J. M. (2018). Impact of index hopping and bias towards the reference allele on

491  accuracy of genotype calls from low-coverage sequencing. *Genetics, Selection,*

492  *Evolution: GSE*, *50*(1), 64.

493  Sackton, T. B., Grayson, P., Cloutier, A., Hu, Z., Liu, J. S., Wheeler, N. E., … Edwards, S. V.

494  (2019). Convergent regulatory evolution and loss of flight in paleognathous birds.

495  *Science*, *364*(6435), 74–78.

496  Shapiro, B., & Hofreiter, M. (2014). A paleogenomic perspective on evolution and gene

497  function: new insights from ancient DNA. *Science*, *343*(6169), 1236573.

498  Weir, J. T., Haddrath, O., Robertson, H. A., Colbourne, R. M., & Baker, A. J. (2016).

499  Explosive ice age diversification of kiwi. *Proceedings of the National Academy of*

500  *Sciences of the United States of America*, *113*(38), E5580–E5587.

501  Weisenfeld, N. I., Yin, S., Sharpe, T., Lau, B., Hegarty, R., Holmes, L., … Jaffe, D. B.

502  (2014). Comprehensive variation discovery in single human genomes. *Nature Genetics*,

503  *46*(12), 1350–1355.

504 Westbury, M. V., Le Duc, D., Duchêne, D. A., Krishnan, A., Prost, S., Rutschmann, S., …

505    Hofreiter, M. (2021). Ecological Specialisation and Evolutionary Reticulation in Extant

506    Hyaenidae. *Molecular Biology and Evolution*, msab055.

507 Westbury, M. V., Petersen, B., Garde, E., Heide-Jørgensen, M. P., & Lorenzen, E. D. (2019).

508    Narwhal Genome Reveals Long-Term Low Genetic Diversity despite Current Large

509    Abundance Size. *iScience*, *15*, 592–599.

510

**Author contributions**

512 Conceptualization, MVW; Formal analysis, AP, MVW; Writing – Original Draft MVW;

513 Writing – Review & Editing EDL, MVW; Funding Acquisition, EDL; Supervision, EDL,

514 MVW.

515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536

537 **Tables**

538

539 **Table 1:** Genome-wide pairwise divergence estimates of the species used in this study.

| Species | Compared to | Divergence |
|---|---|---|
| Narwhal | Beluga | 0.0050 |
| Finless porpoise | Beluga | 0.0143 |
| Bottlenose dolphin | Beluga | 0.0202 |
| Sperm whale | Beluga | 0.0318 |
| Minke whale | Beluga | 0.0344 |
| Brown kiwi | Rowi kiwi | 0.0031 |
| Spotted kiwi | Rowi kiwi | 0.0079 |
| Emu | Rowi kiwi | 0.0734 |

540

541

542 **Table 2:** Autosomal heterozygosity and runs of homozygosity (ROH) estimates of the rowi
543 kiwi when mapped to a variety of different reference genomes. Reference genomes named
544 'Rowi -' are constructed using cross-species scaffolding and the species depicted after the
545 hyphen.

| Reference genome | Global heterozygosity rate | Lower limit | Upper limit | Segments in ROH (%) | Avg. length of ROH (bp) |
|---|---|---|---|---|---|
| Rowi (published) | 0.00121 | 0.00111 | 0.00126 | 4.42 | 1,644,440 |
| Rowi (re-assembled) | 0.00105 | 0.00086 | 0.00123 | 1.57 | 1,076,920 |
| Brown kiwi | 0.00121 | 0.00107 | 0.00139 | 0.22 | 1,000,000 |
| Spotted kiwi | 0.00119 | 0.00108 | 0.00134 | 0.00 | 0 |
| Emu | 0.00177 | 0.00157 | 0.00193 | 0.00 | 0 |
| Rowi - brown kiwi | 0.00107 | 0.00097 | 0.00116 | 2.74 | 1,833,330 |
| Rowi - spotted kiwi | 0.00106 | 0.00098 | 0.00117 | 2.82 | 1,444,440 |
| Rowi - emu | 0.00099 | 0.00089 | 0.00107 | 3.28 | 1,558,140 |

546

547

548

549

550

551

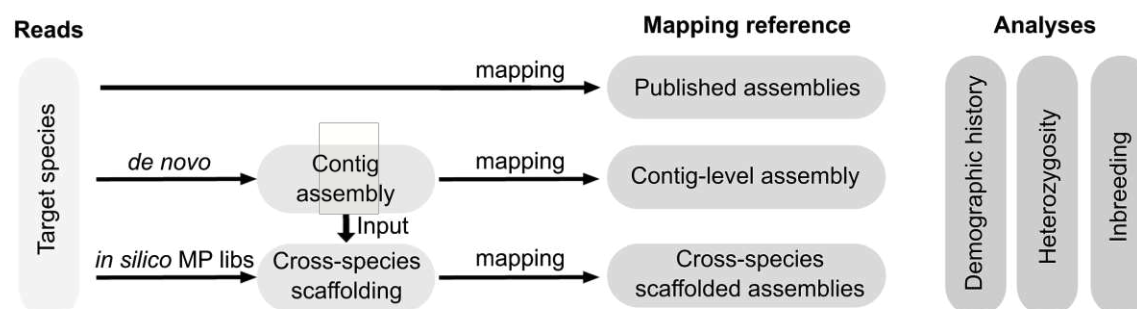552

553 **Figures**

554



555
556
557 **Figure 1:** Overview of the approaches used to investigate the role reference genome plays in
558 downstream demographic history and genetic diversity results. Raw reads are mapped to
559 published assemblies, a *de novo* contig-level assembly, or cross-species scaffolded
560 assemblies. Contig-level assemblies are constructed using the raw reads. Cross-species
561 scaffolded assemblies are made by scaffolding the contig assembly using *in-silico* mate-pair
562 (MP) libraries.

563
564

**Figure 2:** Beluga demographic history over the last 2.5 million years. Demographic trajectories in each panel represent genomes generated by mapping beluga reads to (A) assemblies of six different phylogenetically distant species (including a conspecific), colours indicate species of the reference genome, and (B) *de novo* assemblies constructed using cross-species scaffolding and the published beluga assembly, colours show the species used to scaffold the *de novo* beluga contig-level assembly.
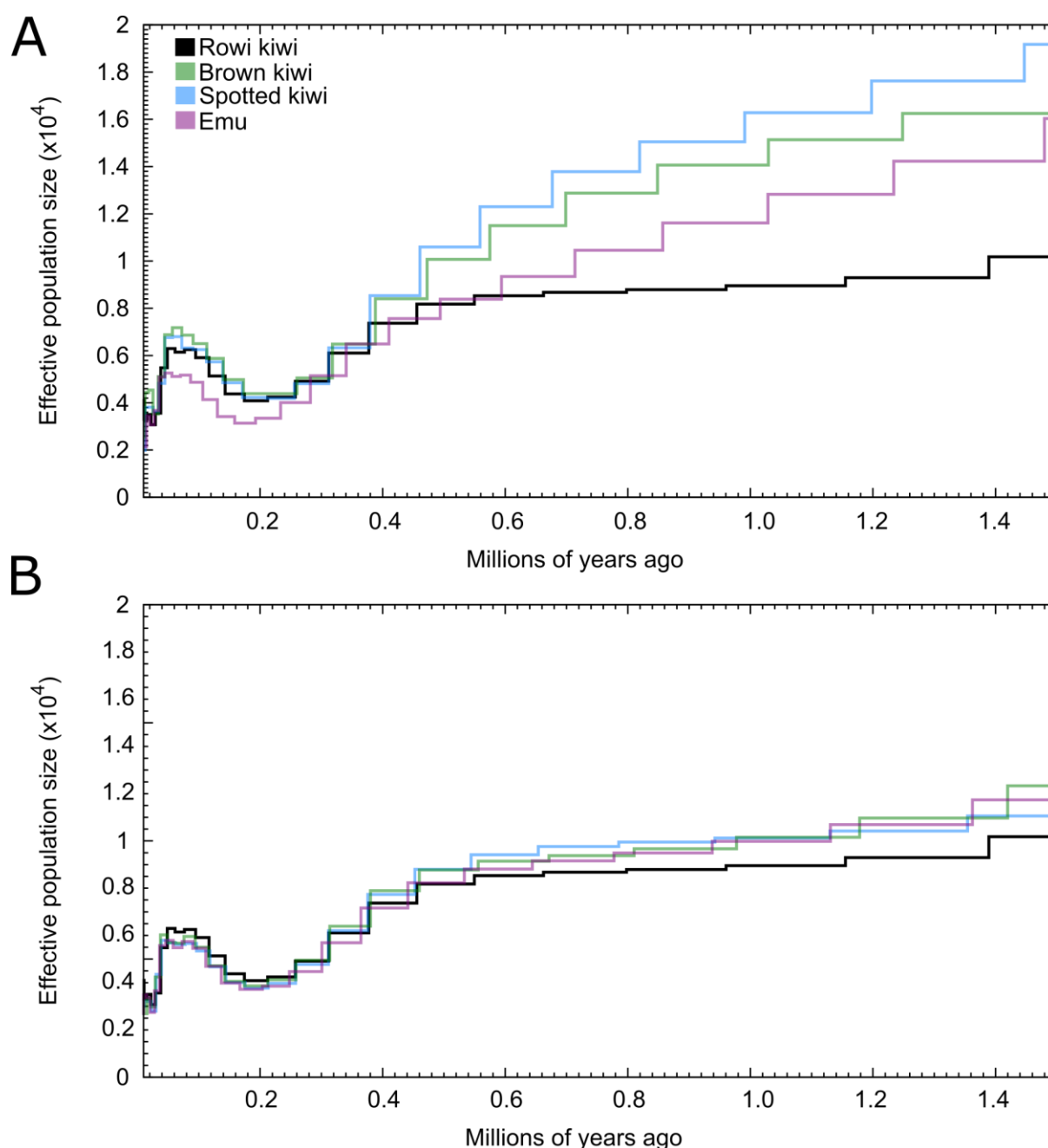
**Figure 3:** Rowi kiwi demographic history over the last 1.5 million years. Demographic trajectories in each panel represent genomes generated by mapping rowi kiwi reads to (A) assemblies of four different phylogenetically distant species (including our re-assembled rowi kiwi assembly), colours indicate species of the reference genome, and (B) *de novo* rowi kiwi assemblies constructed using cross-species scaffolding and our reassembled rowi kiwi assembly - colours show the species used to scaffold the *de novo* rowi kiwi contig-level assembly.
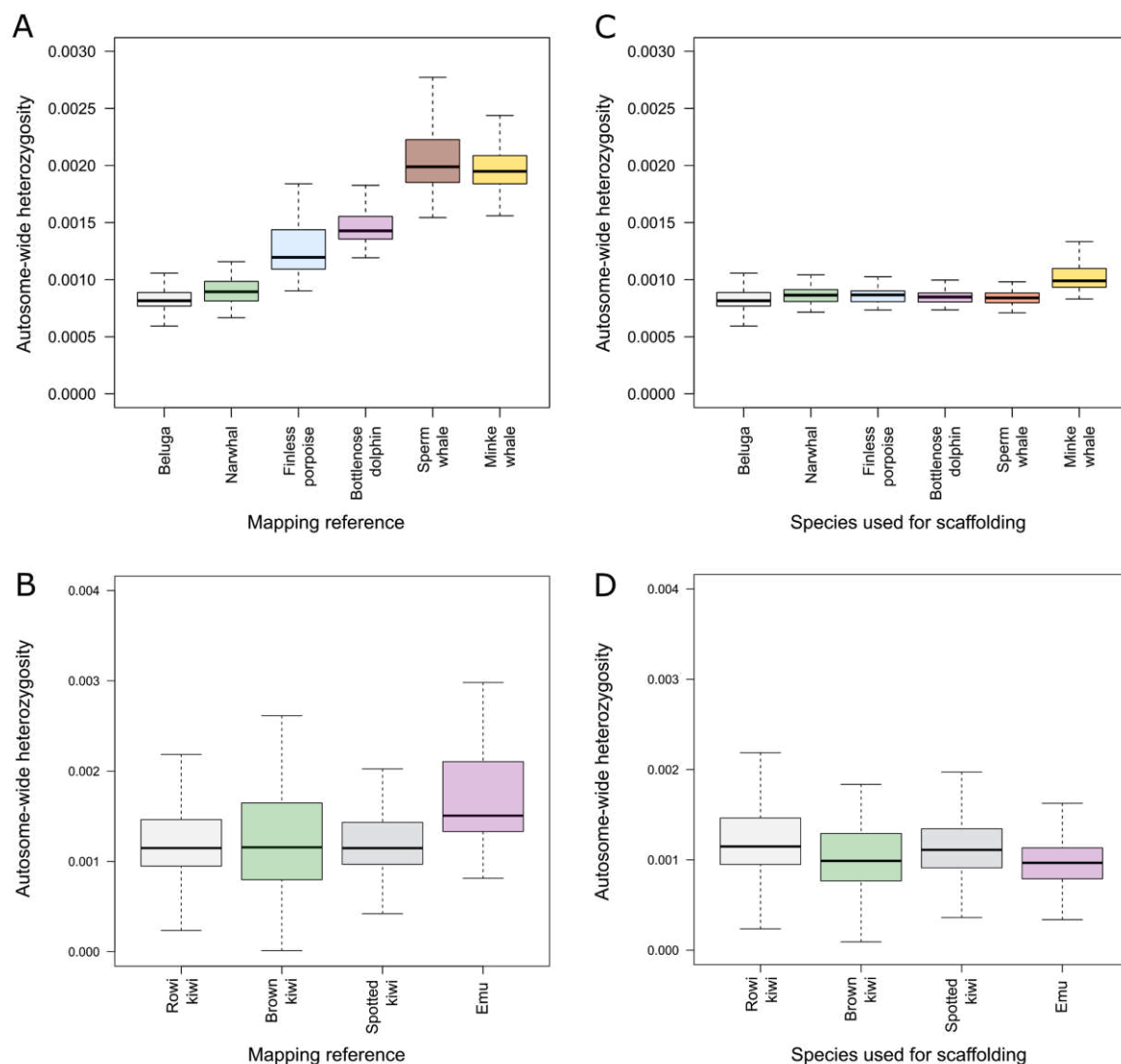
**Figure 4:** Autosome-wide heterozygosity estimates of the beluga and rowi kiwi mapped to different reference genomes. A single beluga individual was mapped to (A) six downloaded assemblies, and (B) a published beluga assembly and *de novo* beluga assemblies constructed using cross-species scaffolding. A single rowi kiwi individual was mapped to (C) our re-assembled rowi kiwi genome and the three downloaded non-rowi kiwi assemblies, and (D) *de novo* rowi kiwi assemblies constructed using either the published rowi kiwi mate-pair libraries or using cross-species scaffolding with mate-pair (MP) libraries constructed from each of the three non-rowi kiwi assemblies.

27