



Joint identification of sex and sex-linked scaffolds in non-model organisms using low depth sequencing data

Nursyifa, Casia; Brüniche-Olsen, Anna; Garcia-Erill, Genis ; Heller, Rasmus; Albrechtsen, Anders

Published in:
Molecular Ecology Resources

DOI:
[10.1111/1755-0998.13491](https://doi.org/10.1111/1755-0998.13491)

Publication date:
2021

Document version
Peer reviewed version

Citation for published version (APA):
Nursyifa, C., Brüniche-Olsen, A., Garcia-Erill, G., Heller, R., & Albrechtsen, A. (2021). Joint identification of sex and sex-linked scaffolds in non-model organisms using low depth sequencing data. *Molecular Ecology Resources*. <https://doi.org/10.1111/1755-0998.13491>

1 **Joint identification of sex and sex-linked scaffolds in non-model organisms using low depth**
2 **sequencing data**

3 Casia Nursyifa^{1*}, Anna Brüniche-Olsen^{1*}, Genis Garcia Erill¹, Rasmus Heller^{1§}, Anders
4 Albrechtsen^{1§},

5
6 *¹Section for Computational and RNA Biology, Department of Biology, University of Copenhagen,*
7 *Denmark.*

8 * equal contribution

9 § corresponding authors

10 Anders Albrechtsen: albrecht@binf.ku.dk

11 Rasmus Heller: rheller@bio.ku.dk

12

13 **Keywords**

14 Autosomes, bioinformatics, resequencing, scaffold-level assembly

15

16 **Abstract**

17 Being able to assign sex to individuals and identify autosomal and sex-linked scaffolds are essential
18 in most population genomic analyses. Non-model organisms often have genome assemblies at
19 scaffold-level and lack characterization of sex-linked scaffolds. Previous methods to identify sex
20 and sex-linked scaffolds have relied on synteny between the non-model organism and a closely
21 related species or prior knowledge about the sex of the samples to identify sex-linked scaffolds. In
22 the latter case, the difference in depth of coverage between the autosomes and the sex
23 chromosomes are used. Here we present ‘Sex Assignment Through Coverage’ (SATC), a method

24 to assign sex to samples and identify sex-linked scaffolds from next generation sequencing (NGS)
25 data. The method works for species with a homogametic/heterogametic sex determination system
26 and only requires a scaffold-level reference assembly and sampling of both sexes with whole
27 genome sequencing (WGS) data. We use the sequencing depth distribution across scaffolds to
28 jointly identify: i) male and female individuals and ii) sex-linked scaffolds. This is achieved
29 through projecting the scaffold depths into a low-dimensional space using principal component
30 analysis (PCA) and subsequent Gaussian mixture clustering. We demonstrate the applicability of
31 our method using data from five mammal species and a bird species complex. The method is freely
32 available at <https://github.com/popgenDK/SATC> as R code and a graphical user interface (GUI).

33

34 **Introduction**

35 The increasing number of non-model organism scaffold-level genome assemblies provides new
36 information on biodiversity and how evolutionary processes have shaped it (Ellegren 2014). An
37 essential part of genome assembly and annotation is the identification of autosomes and sex
38 chromosomes. Vertebrate species are generally diploid with the majority of their genome
39 represented by a variable number of autosomes and two sex chromosomes (Graves 2008).

40 In mammals, the homogametic sex is the female (XX) and the heterogametic sex is the male (XY).

41 This is opposite in birds, where males are homogametic (ZZ) and females heterogametic (ZW).

42 Due to their inheritance, sex chromosomes differ from autosomes in several aspects of their
43 population genetics and molecular evolution, e.g. by having a smaller effective population size
44 (N_e) than autosomes (Ellegren 2009) and by having different patterns of population differentiation,
45 especially under incipient or complete speciation (Presgraves 2018). Therefore, it is often
46 preferable to separate them from autosomes in population genetic analyses.

47 Ideally, whole genome assemblies should be at chromosome-level and fully annotated, but
48 due to the high cost and challenges associated with complete genome assembly this is often not
49 prioritized for the first generation of a reference genome (Ellegren 2014). Consequently, several
50 approaches to identify sex chromosomes in scaffold-level assemblies have been developed (for a
51 review see Palmer *et al.*, (2019)). One of them is whole genome synteny alignment (Grabherr *et*
52 *al.*2010), where the scaffold-level genome assembly is aligned to a chromosome-level assembly
53 from a closely related species, and sex-linked scaffolds are identified based on sequence similarity
54 to the reference. There are several obstacles to this approach, most importantly the availability of
55 a chromosome-level assembly of a closely related species, but also the accelerated evolution of sex
56 chromosomes in many lineages which causes a high degree of divergence even for closely related
57 species (Charlesworth *et al.*2018; Irwin 2018; Meisel & Connallon 2013; Presgraves 2018), and
58 computational time (Pennell *et al.*2018). Moreover, neo-sex chromosomes are likely to be missed
59 due to their partial or complete synteny with autosomes in the closely related species genome
60 (Graves 2008).

61 An alternative approach is to use genome depth of coverage (DoC) based methods. These
62 are represented by two groups of methods, Y-linked and X-linked. Both methods require prior
63 sexing information of individuals, which can be challenging to obtain for species with limited
64 sexual dimorphism, for cryptic species and in non-invasive sampling situations. Y-linked scaffolds
65 can be identified based on their low DoC compared to the autosomes, due to the expectation of
66 $\frac{1}{2}$ ×autosomal DoC in samples of the heterogametic sex (Hall *et al.*2013). X-linked scaffolds can
67 be identified when sequencing reads from both sexes are mapped to a reference, with the
68 expectation that the homogametic-linked scaffolds will have 1×autosomal DoC in the
69 homogametic sex and $\frac{1}{2}$ ×autosomal DoC in the heterogametic sex (Ellegren 2009, Graves 2008).

70 Due to noise in next generation sequencing (NGS) data, scaffold DoC distributions often deviate
71 from this expectation. For example, autosomal and sex-linked scaffolds can have overlapping DoC
72 distributions, making it challenging to clearly identify individual X and Y scaffolds (Malde *et al.*,
73 2019). The DoC approaches are furthermore highly sensitive to parameters used for the read
74 mapping, pre- and post-mapping filtering steps and data quality, e.g. whether or not repetitive
75 regions are removed, the average genome-wide DoC, etc. (Smeds *et al.*, 2015).

76 Here we present ‘Sex Assignment Through Coverage’ (SATC): a method and software to
77 jointly identify sex-linked scaffolds and determine the sex of each sample mapped to a scaffold-
78 level assembly. The method requires sequencing depth information from whole genome
79 resequencing of male and female samples and applies principal component analysis (PCA) and
80 Gaussian mixture clustering to group the dataset into males and females. Hence, the method
81 harnesses the systematic—but noisy—difference in DoC of the heterogametic sex-linked scaffolds
82 compared to the autosomal scaffolds. To illustrate how the method works we applied it to five
83 mammal species and a bird species complex with different levels of DoC, assembly quality (N50)
84 and sample sizes. Our method is very fast with computational time being less than a minute for
85 100 samples. We also evaluated the quality of input data required for SATC to work by
86 subsampling our original data in various ways. We anticipate that it will be widely useful for
87 inferring individual sex and identification of sex-linked scaffolds for non-model organisms.

88

89 **Materials and methods**

90 *Method*

91 The inputs for SATC are scaffold lengths and mapping statistics (e.g., number of reads mapped to
92 each scaffold for each sample). These are quickly generated with the ‘idxstats’ command in

93 SAMTOOLS (Li *et al.*, 2009) from indexed bam files. The SATC method works by: i) normalizing
94 the DoC of each scaffold within each sample, ii) reducing the dimensionality of the normalized
95 DoC using PCA, iii) clustering the samples using Gaussian mixture clustering on the top PCs, and
96 iv) identifying the sample sex and sex-linked scaffolds from the clustering and the DoC

97 In the first step, we calculate for each sample the average DoC per scaffold and normalize
98 it by the mean DoC of the M longest scaffolds in the reference assembly. Suppose we have $s= 1,$
99 $2, \dots, S$ scaffolds, $n= 1, \dots, N$ samples and a matrix $R \in \mathbb{N}^{S \times N}$ containing the number of reads
100 mapped to each scaffold for each sample. For simplicity we here assume that the scaffolds are
101 ordered by length with scaffold 1 being the largest; however, this is not a requirement when using
102 the method. From this we generate a matrix of normalized mean DoC for each scaffold and sample

$$103 \quad D_{sn} = \frac{R_{sn}}{l_s} \times \frac{\sum_{i=1}^M l_i}{\sum_{i=1}^M R_{in}}$$

104 where R_{sn} is the number of mapped reads on scaffold s , sample n , l_s is the length in bases of scaffold
105 s , R_{in} is the number of mapped reads on scaffold i and l_i is the length in bases of scaffold i , $i=1, 2,$
106 \dots, M longest scaffolds. In our default implementation we set M equal to 5 as we found it to give
107 robust results, but this number can be changed by the user. We choose to use more than one scaffold
108 to allow SATC to give warnings if one of the normalization scaffolds has a substantially different
109 DoC than the others, which can alert the user that the scaffold should not be used in the
110 normalization—it could for example be a sex-linked scaffold. Using multiple scaffolds for the
111 normalization also provides the option to normalize by the median instead of the mean DoC if
112 there is a depth difference between scaffolds. After this normalization, most scaffolds will have a
113 normalized DoC close to 1. Following preliminary analyses, we found that filtering out scaffolds
114 $< 100\text{kb}$ in length and those with a mean normalized DoC outside the range of 0.3 – 2.0 improved

115 the performance of the method, but these thresholds can be changed by the user. We then center
116 the matrix by subtracting the mean normalized DoC from that of each scaffold

$$117 \quad \widetilde{D}_{sn} = D_{sn} - (\sum_{n=1}^N D_{sn} / N)$$

118 and perform PCA on the \widetilde{D} matrix. We proceed by identifying two clusters of samples that
119 tentatively represent different sex groups (e.g. male and female), on the assumption that a large
120 proportion of DoC variance is attributable to differences between males and females in the DoC of
121 reads mapping to sex-linked scaffolds. For this we use the first two principal components as input
122 for a clustering analysis using Gaussian finite mixture models in *mclust* (Scrucca *et al.*, 2016).
123 Mclust has several models and the default option chooses the model that gives the lowest Bayesian
124 Information Criterion (BIC). We chose a model with ellipsoidal distribution, equal volume and
125 variable shape and distribution (EVV), that assumes equal variance in the two groups. Sometimes,
126 this model fails to identify groups and in that case SATC will choose the model with the lowest
127 BIC.

128 To identify the sex-linked scaffolds we use the two groups identified by the above
129 clustering. We first apply a *t*-test assuming unequal variance for each scaffold to test for significant
130 differences in mean DoC between the groups. We use a Bonferroni corrected *p*-value cut off of
131 $0.05/S$ to identify the sex-linked scaffolds, and scaffolds identified with this test are broadly
132 referred to as sex-linked in the following. In the *t*-test we assume normally distributed normalized
133 DoC within each group, which might be violated for some datasets. We then calculate the average
134 DoC for each scaffold in each group. The heterogametic (XY/ZW) and homogametic (XX/ZZ)
135 sexes are expected to show a mean normalized DoC ratio of 0.5:1 between the two groups for
136 scaffolds situated on the X/Z chromosome. Therefore, to identify scaffolds that are highly likely
137 to be situated on the X/Z chromosome, we retain sex-linked scaffolds identified with the above

138 method for which the mean difference of normalized DoC between two sex groups is between 0.4
139 and 0.6. We refer to these scaffolds as X/Z-linked.

140 For the datasets analyzed in our study SATC performed well without any filtering of the
141 mapped reads or scaffold contents. However, if the user wishes to add additional filters, for
142 example to exclude repetitive regions of the genome or remove N's in the genome assembly, then
143 the input file for the method is easily customizable. All analyses were done in R (R Core Team
144 2019) and are freely available at <https://github.com/popgenDK/SATC> as R code and a graphical
145 user interface (GUI)..

146

147 *Application to empirical datasets*

148 To test our method we used six mammal and bird whole genome sequencing (WGS) datasets with
149 low to medium DoC and mapped to scaffold-level assemblies. We used unpublished WGS for
150 impala (*Aepyceros melampus*), muskox (*Ovibos moschatus*), waterbuck (*Kobus ellipsiprymnus*)
151 and gray whale (*Eschrichtius robustus*). WGS for leopard (*Panthera pardus*) were from Pečnerová
152 *et al.*, (2021) and we used only the 49 individuals that passed their QC. The scaffold assemblies
153 for the waterbuck, impala and muskox were from the ruminant genome project (Chen *et al.*, 2019),
154 the leopard were from Kim *et al.*, (2016) while the gray whale assembly is unpublished. WGS data
155 from Darwin's finches species complex represents 15 species (Lamichhaney *et al.*, 2015;
156 Lamichhaney *et al.*, 2016)—the mangrove finch (*Camarhynchus heliobates*), the woodpecker
157 finch (*Camarhynchus pallidus*), the small tree finch (*Camarhynchus parvulus*), the medium tree
158 finch (*Camarhynchus pauper*), the large tree finch (*Camarhynchus psittacula*), the gray warbler-
159 finch (*Certhidea fusca*), the green warbler-finch (*Certhidea olivacea*), the Española cactus finch
160 (*Geospiza conirostris*), the sharp-beaked ground finch (*Geospiza difficilis*), the medium ground

161 finch (*Geospiza fortis*), the small ground finch (*Geospiza fuliginosa*), the large ground finch
162 (*Geospiza magnirostris*), the common cactus finch (*Geospiza scandens*), the Cocos finch
163 (*Pinaroloxias inornate*), and the vegetarian finch (*Platyspiza crassirostris*). We also included two
164 related tanager species; the black-faced grass-quit (*Tiaris bicolor*) and the lesser Antillean
165 bullfinch (*Loxigilla noctis*). All finches species were mapped to the medium ground finch assembly
166 (Zhang *et al* 2014).

167 The sequencing data from the six different datasets were generated as part of separate
168 studies and were therefore filtered and pre-processed in different ways (for a detailed description
169 see Supplementary Information Text S1). Hence, their heterogeneity allows us to assess whether
170 our method is broadly applicable across a range of data treatment regimes, representing the variety
171 of pipelines used in practice for non-model sequencing analysis.

172

173 ***Validation of sexing and sex-linked scaffolds***

174 To evaluate the sensitivity of the sample sex assignment we mapped the read data from the five
175 mammal species to a closely related, well-annotated chromosome-level reference genome that
176 included an annotated X chromosome. We did not include the Y chromosome because it was not
177 available in all reference genomes and is in general much harder to assemble. For the impala we
178 mapped to the goat (ARS1), for the leopard to the domestic cat (*Felis_catus_9.0*), for the waterbuck
179 we used the cow (*bosTau8*), for the muskox we used the sheep (*Oar_rambouillet_v1.0*), and for
180 the gray whale we used the blue whale (*mBalMus1.pri.v3*). Based on this cross-species mapping
181 we calculated the normalized DoC of X-mapping reads for each individual by the DoC of reads
182 mapping to the five largest autosomal chromosomes of the same external reference genome as an
183 external evaluation of the SATC accuracy in assigning the correct sample sex. As 16 of the 17

184 species in the Darwin's finches dataset were mapped to an outgroup (the medium ground finch)
185 we did not explore mapping this dataset to further outgroups. To evaluate if evolutionary
186 divergence among taxa would influence our method we ran SATC on the finch data with and
187 without the more distantly related species, the black-faced grass-quit and the lesser Antillean
188 bullfinch.

189 In addition, we also used an independent verification of sexing for the gray whale dataset
190 where we had sex information from field observations ($n = 37$) and from single nucleotide
191 polymorphism (SNP) genotyping of two markers in the zinc finger proteins of the *Zfx* and *Zfy*
192 genes ($n = 70$) (DeWoody *et al.*, 2017; Brüniche-Olsen *et al.*, 2018).

193 Finally, we used LASTZ (Harris 2007) to explore the degree of synteny between the SATC
194 inferred sex-linked scaffolds and sex chromosomes from the closest relative with a chromosome-
195 level assembly. For this we focused on the gray whale and Darwin's finches datasets, as these were
196 the taxa for which the closely related, chromosome-level assembly was heterogametic (blue whale
197 (mBalMus1.pri.v3) and chicken (bGalGal1.mat.broiler.GRCg7b), respectively). We ran LASTZ
198 for each putative sex-linked scaffold using a low-sensitivity alignment with the following settings
199 --notransition --step=20 --nogapped.

200

201 *Testing the limitations of SATC performance under extreme data conditions*

202 Three main factors are likely to influence the performance of SATC: i) overall sequencing data
203 DoC, ii) the number and sex distribution of samples, and iii) reference genome assembly quality.
204 We therefore tested the robustness of the SATC results to factors i and ii, while the impact of factor
205 iii was implicitly evaluated by using species with highly variable reference genome qualities (Table
206 1).

207 To evaluate the robustness of our method for low DoC data we downsampled the
208 sequencing reads to a DoC of 1X, 0.5X, 0.1X and 0.01X. Downsampling for each sample in each
209 scaffold was done by randomly subsampling reads to achieve the desired DoC directly from the
210 idxstats files. We then compared the inferred sex from these lower DoC scenarios with their
211 externally validated sex assignment found by using the full dataset. We also compared the inferred
212 sex-linked scaffolds and their combined length to the sex-linked scaffolds identified using the full,
213 non-downsampled dataset.

214 To assess the accuracy of sexing in scenarios with large sex ratio imbalances, we performed
215 stratified downsampling for each sex group, reducing each group to 75%, 50%, 25% and 10% of
216 their original sample size and also downsampled each data set to only 2 samples for each sex. To
217 evaluate the inferred sex, we visualized the normalized DoC of sex-linked scaffolds for each
218 sample using boxplots where the median should be a half for the heterogametic and one for the
219 homogametic. We also visualize the PCA plot with uncertainty profiles calculated from *mclust* in
220 which uncertainty is measured as probability belonging to its group and the two sex groups clusters
221 are visualized by two ellipsoids as scaled variances of inferred Gaussian mixture components.

222

223 **Results**

224 We analyzed whole-genome resequencing data from five mammal species and one bird species
225 complex. The WGS data was mapped to the scaffold-level genome assembly from the same
226 species. The datasets varied both in terms of quality of assembly and in DoC, which ranged from
227 3.13–13.76X. As shown in Table S1 the quality of the genome assembly varied between species,
228 ranging from a scaffold N50 of 344kb in impala to 46.8Mb in muskox and each genome assembly
229 contained 2,796–88,935 scaffolds. These assemblies are representative for many low to medium

230 quality draft genomes. Even after removing scaffolds <100kb, a high number of scaffolds remained
231 for some species i.e., up to 7,717 for the impala (Table S1).

232 The normalized DoC was very noisy across scaffolds for most species (Figure S1).
233 However, when we performed a PCA on the DoC matrix we observed a clear separation into two
234 groups for all taxa (Figure 1, left column). All taxa—except the impala—separated in two distinct
235 groups based on PC1; for impala this partitioning was on PC2. After applying Gaussian finite
236 mixture models clustering on the two first axes of the PCA we could clearly group the samples
237 from each taxa into two groups with characteristic and distinct normalized DoC patterns. We
238 interpret these two groups as the homogametic and heterogametic sex, with the homogametic sex
239 being the one with the highest DoC of the sex-linked scaffolds (Figure S2).

240 To identify sex-linked scaffolds we performed a *t*-test for each scaffold testing differences
241 in mean DoC between the two sex groups identified above. This test identified 54-589 sex-linked
242 scaffolds across the six taxa, with by far the highest number found in the two most challenging
243 data sets from the impala and waterbuck (Table 1 and Figure S2). These sex-linked scaffolds might
244 not be exclusively from a sex chromosome, and we therefore define them loosely as sex-linked.
245 We furthermore identify X/Z-linked scaffolds by retaining only those sex-linked scaffolds that had
246 a difference of 0.4–0.6 between the mean normalized DoC in the heterogametic and homogametic
247 individuals, respectively. This yielded between 11 and 113 X/Z-linked scaffolds in each species
248 (Table 1 and Figure 1, right column). The total length of sex-linked scaffolds was 126–187 Mb
249 across the five mammals and 79 Mb in the bird species complex, which is close to the length of
250 the assembled X/Z chromosomes for the close relative of each species (Table 1).

251 Many of the sex-linked scaffolds that do not conform to the X/Z-linked criterion were
252 scaffolds with DoC that are correlated with sex, but do not adhere to the 0.5:1 ratio between the

253 heterogametic and homogametic sexes. Some of these scaffolds show large and unexplained
254 deviations from the autosomal DoC, while some show an expected Y/W ratio of 0.5:0 (Figure S2,
255 muskox and Darwin's finches). Most of the taxa had a small number of scaffolds that were assigned
256 as sex-linked but not X/Z-linked, however the two species with the most fragmented assemblies,
257 impala and waterbuck, showed a much larger difference between the cumulative length of sex-
258 linked and X-linked scaffolds (Table 1). For the impala, the cumulative inferred X-linked scaffold
259 length was just 1.4% of the sex-linked scaffold length, and for the waterbuck this was 24.8%,
260 whereas the corresponding numbers were 79.4-98.4% for the remaining species (Table 1). We
261 therefore conclude that our sex-linked scaffold identification method likely identifies the majority
262 of the X/Z-linked scaffolds while our stricter ratio-based criterion will not identify most of the
263 X/Z-linked scaffolds when the reference genome quality and/or resequencing DoC is low.

264

265 *Verification of SATC sex assignment*

266 To validate the inferred sex for the samples, we mapped the five mammal species reads to a closely
267 related reference genome containing an annotated chromosome-level assembly, and used the
268 normalized DoC of reads mapping to the X chromosome to classify samples as heterogametic or
269 homogametic, respectively. This validation showed 100% agreement with the SATC inferred sex
270 in all cases and across all species (Figure S3). The Darwin's finches dataset represented 15 closely
271 related species and two more distantly related species. We found that even with such a
272 heterogeneous dataset, SATC was able to identify the sample sex and the sex-linked scaffolds
273 regardless of whether we included the most distantly related species in the complex (Figure S4).

274 We further verified the SATC sample sex assignment using SNP-genotyping and field
275 observations. The SNP-genotyping and field observations for the gray whale dataset showed high

276 concordance (Table S2). We found a 100% agreement between the SATC results and field
277 observations. We identified one misclassification in the SNP-genotyping, a female identified by
278 SATC and field observations, but misidentified as a male with genotyping, probably due to a
279 sample mixup in the SNP-genotyping or sample contamination.

280 We evaluated the synteny of the inferred sex-linked scaffolds with known sex
281 chromosomes by aligning the inferred sex-linked scaffolds to XY and ZW chromosomes from
282 closely related species and visualizing the similarity across the sex chromosomes (Figure S5 and
283 S6). For both the gray whale and Darwin's finches we found a higher percentage of scaffold
284 sequence aligning to the X/Z chromosomes (X mean aligned sequence % = 62.9 SD = 15.5 and Z
285 mean aligned sequence % = 6.0; SD = 3.4) than the sex-limiting Y/W chromosomes (Y mean
286 aligned sequence % = 1.2; SD = 0.6 and W mean aligned sequence % 1.9; SD = 2.9). Generally,
287 the percentage of aligned sequence was higher in the gray whales than in Darwin's finches, as
288 illustrated by the low amount of aligned sequence % in Darwin's finches (Figure S5+S6). However,
289 as seen in Fig. S5 and Fig. S6 the majority of the X/Z chromosomes of the closely related reference
290 genomes had clear synteny with one of the sex-linked scaffolds identified by SATC, confirming
291 that we have successfully identified the majority of the X/Z-linked scaffolds in our study taxa.

292 *SATC performance in challenging data situations*

293 We evaluated how the method performed for: i) low DoC data, ii) data with a low number of
294 samples, and iii) data with imbalanced sex ratios. By downsampling sequencing reads we showed
295 that SATC still correctly assigned sample sex for all taxa down to 0.1X. For five of the datasets,
296 SATC also worked at a very low DoC of 0.01X. However, at this extremely low DoC the method
297 failed to assign the correct sex for the impala dataset. (Figure S7). Moreover, we inferred the same
298 set of sex-linked scaffolds at lower DoC for most datasets except the impala and the waterbuck, in

299 which SATC missed many of the sex-linked scaffolds at 0.1X and 0.01X. (Table S3). SATC also
300 remains highly accurate in assigning sex with either a low total number of samples or with highly
301 unbalanced sex ratios among the input samples. For the gray whales and leopard the method
302 correctly identified the sexes both at very low samples sizes ($n_{XX}=2$, $n_{XY}=2$) and with extreme sex
303 ratio imbalance (e.g. $n_{XX}=47$, $n_{XY}=3$). For the impala, Darwin's finches, muskox and waterbuck
304 SATC performed well for a range of scenarios, but failed at sex ratio imbalances of 81:8, 5:126,
305 5:56, and 5:20 respectively (Figure S8). For all of the failed scenarios we detected issues with both
306 sample DoC clustering and with the DoC distribution on the inferred sex-linked scaffolds, both of
307 which are clear indications of poor performance (Figure S9). The PCA plots for these failed SATC
308 analyses illustrate how the clustering can diagnose overlapping or low separation of the two
309 ellipsoidals representing the two putative sex groups. We also observed that some individuals could
310 not be assigned to a cluster with high certainty, which is illustrated in the diagnostic PCA plots by
311 large point sizes. The corresponding boxplots of normalized DoC in all the SATC inferred sex-
312 linked scaffolds clearly do not follow the expected normalized DoC of 0.5 for the heterogametic
313 and 1.0 for the homogametic sex (Fig. S9). For comparison, the corresponding evaluation plots in
314 the situations where SATC works well are shown in Figure S10.

315 We also evaluated the effect of using sex-linked scaffolds in the DoC normalization using
316 the muskox data as an example. We did not see any effect of including the largest inferred sex-
317 linked scaffold along with the four largest non-sex-linked scaffolds in the DoC normalization. We
318 also tested the more extreme case of using the largest non-sex-linked scaffold and the 20 largest
319 sex-linked scaffolds. This resulted in correct clustering and assignment of sex (Figure S11).
320 However, the method now identified the autosomal scaffolds as sex-linked scaffolds, with most of
321 them being categorized as having abnormal DoC ratios. We believe that this latter case is extremely

322 unlikely to occur in a real analysis and it can be avoided easily by inspecting the DoC plots given
323 by the software. In addition the software gives several warnings to the user that the chosen
324 normalizing scaffolds have very different DoCs and suggests to the user to change the
325 normalization scaffolds or use the median DoC instead.

326

327 **Discussion**

328 The increasing amount of whole genome resequencing data presents new avenues for population
329 genomic analyses. Here we add to the analytical toolset by introducing SATC, a method for joint
330 individual sex assignment and identification of sex-linked scaffolds from NGS data in species with
331 a homogametic/heterogametic sex-chromosome system. Our method is automated,
332 computationally light, robust to pre-mapping filtering and has a high accuracy even with
333 challenging data. We anticipate this will be particularly useful for non-model organisms and for
334 samples collected in the field, where information on the sex of individuals and a chromosome-level
335 assembly is often lacking.

336 The benefits of identifying sex-linked scaffolds when carrying out population genetic
337 studies are many. First, sex chromosomal sites may be desirable for specific analyses, such as
338 association (Lee *et al.*2017; Luciano *et al.*2019; Zuo *et al.*2013), gene expression (Grath & Parsch
339 2016), or any evolutionary genetic studies on X/Y or Z/W chromosomes (Gottipati *et al.*2011).
340 Second, if sex-linked scaffolds are not flagged and treated separately they may bias analyses such
341 as demographic history inference (Li & Durbin 2011), genome scans or genome-wide values of
342 summary statistics, including F_{ST} (Lambert *et al.*2010), genetic diversity (Hammer *et al.*2010) and
343 allele frequency distribution (Clayton 2008). Third, analyzing males and females separately can

344 elucidate patterns of sex-biased dispersal (Bidon *et al.*2014) or unequal contributions to offspring
345 diversity (Pérez-González *et al.*2014).

346 We show here that PCA on normalized scaffold DoC is a robust approach to identify
347 individual sex for a range of data situations, including having only low-depth resequencing data
348 and a low quality draft assembly as a reference genome. Having assigned sample sex we can easily
349 reverse the perspective and utilize this information to identify which scaffolds are sex-linked by
350 exploiting the expected 0.5:1 ratio between heterogametic and homogametic DoC for each X/Z
351 linked scaffold. Our SATC approach does not rely on prior knowledge of sample sex or sex-linked
352 scaffolds, and is, to our knowledge, the only automated software that accomplishes this without
353 any external information. The recommended usage of SATC would be to flag all scaffolds
354 identified as sex-linked and remove them from further analyses that assume autosomal
355 chromosome data. Conversely, if X or Z-linked sites are desired, we recommend to include only
356 those that are flagged as X/Z-linked i.e., approximately follow the expected 0.5:1 ratio between
357 DoC in the heterogametic and homogametic sex when compared between the two inferred groups
358 of same-sex samples. Note that it is much harder to identify Y/W-linked scaffolds due to the
359 smaller size of these chromosomes and their highly repetitive sequence content. This leads to a
360 higher occurrence of misassembled Y/W scaffolds, which can distort the DoC. Despite this, some
361 of the sex-linked scaffolds we identify show DoC patterns consistent with being Y/W-linked,
362 having approximately 0.5 normalized DoC in males and very low (virtually zero) DoC in females
363 (Figure S2).

364 We show that highly fragmented genome assemblies can be used in SATC, albeit with
365 reduced performance of the sex-linked scaffold identification. The two examined species with the
366 lowest reference genome scaffold N50s, impala and waterbuck, showed deviating patterns from

367 the rest by having very noisy scaffold DoC (Figure S1). The most obvious explanation for this is
368 poor assembly quality, e.g. regions from both sex chromosomes and autosomes could be
369 erroneously joined in the same scaffold. The impala had the lowest-quality genome assembly (N50
370 = 344kb), and for this species we found that the grouping of sexes occurred in PC2 rather than PC1
371 (Figure 1). Despite this, our clustering approach was still able to assign sample sex with 100%
372 accuracy. In addition to lower assembly quality, biological factors could also influence the ability
373 of SATC to correctly identify sex-linked scaffolds. For example, impalas are known to have
374 segregating karyotypic polymorphisms (Pagacova *et al.*2011), which could potentially influence
375 the depth across scaffolds and exacerbate the noise in DoC. The waterbuck had a higher than
376 expected amount of sex-linked scaffold content with about 40Mb more content than the X
377 chromosome of the cow reference genome. This could again be influenced by karyotypic
378 polymorphisms, which are known to occur both within and between different subspecies of
379 waterbuck (Kingswood *et al.*1998). Autosome-to-X translocations are known from several species
380 of bovids (Efron *et al.*1976; Gallagher Jr & Womack 1992; Kumamoto *et al.*1996), and if such
381 are segregating within our samples they would complicate the DoC-based identification of sex-
382 linked scaffolds. We also observed a large difference between the total amount of sex-linked
383 scaffold content identified by the DoC ratio for these two species, whereas this was much smaller
384 for the other species (Table 1), confirming that excessive noise in scaffold DoC can challenge the
385 use of hard thresholds for identifying sex-linked scaffolds.

386 Our SATC method also works well with heterogeneous datasets. Darwin's finches
387 encompass around 18 species of passerine birds (Grant & Grant 2020). We analyzed 15 of these
388 species, which diverged during the last 150,000-900,000 years but still have some degree of
389 interspecies gene flow (Lamichhaney *et al.*2015). Despite the heterogeneous data we were able to

390 assign sex and to identify sex-linked scaffolds in the medium ground finch reference genome
391 assembly when read data from across these species were analyzed in SATC together. We extended
392 the Darwin's finches dataset with two more distantly related (>900ky) tanager species, the black-
393 faced grassquit and lesser Antillean bullfinch (Lamichhaney *et al.*2015), and show that the PCA
394 clustering method was still able to reliably assign sample sex as well as identify sex-linked
395 scaffolds (Figure S4).

396 Finally, we found that there was a high degree of synteny between the identified sex-linked
397 scaffolds and regions on the sex chromosomes in the closest chromosomal assembly for both
398 mammals and birds (Figure S5 and S6). The synteny was apparent despite a relatively low mean
399 sequence identity (< 6%) between the chicken sex chromosomes and the medium ground finch
400 sex-linked scaffolds. This illustrates that in cases where there is high divergence (i.e., low sequence
401 similarity, chromosome degradation), synteny based approaches alone might have challenges.
402 SATC has an advantage over synteny-based approaches because it can differentiate between sex-
403 linked scaffolds on actual sex chromosomes and on neo-sex chromosomes. The fusion between a
404 sex chromosome and an autosome generates a neo-sex chromosome, and these can not be
405 distinguished from the true sex chromosomes through synteny. Therefore, SATC is expected to
406 perform better at identifying sex-linked scaffolds than synteny-based methods in taxa with more
407 dynamic sex-chromosomes (i.e. non-homeotherms) where neo-sex chromosomes are fairly
408 common.

409 Under extreme data conditions we showed that SATC still works well in both sample sex
410 assignment and sex-linked scaffolds identification. Ultra-low DoC (0.01X) is not an issue for
411 SATC when genome assembly quality is reasonably good, in our analyses represented by the
412 muskox, leopard, gray whale, and Darwin's finches cases (Figure S7). However, SATC did fail in

413 some scenarios with large sex imbalance or with a very low number of individuals (Figure S8).
414 The amount of sex imbalance at which SATC no longer performed well varied according to
415 genome assembly quality. In the better-quality scaffold-level assemblies from the four taxa
416 mentioned above, SATC performed well up to a sex imbalance of 1:10, while for the highly
417 fragmented assemblies from impala and waterbuck SATC failed at 1:4. In failed scenarios the
418 parameters in the model can be tweaked for better performance, as we have for now focused on
419 identifying generic settings that work well across many cases. Importantly, the user can easily
420 evaluate SATC performance using a PCA plot with uncertainty profile and individual DoC plot
421 (Figure S9 and S10). If these diagnostic plots show unsatisfactory results, the user should assume
422 that SATC has not performed well and the results should not be trusted. In this case, the user can
423 try either to change the scaffolds used for normalization, normalization method (e.g. using median
424 instead of mean) or change the clustering method where we also have an option to perform
425 hierarchical clustering instead of Gaussian mixture model-based clustering.

426 We explored SATC's performance on a range of genomic datasets from mammals and
427 birds. However, we did not test how well the method works for taxa with highly degenerated sex-
428 chromosomes or taxa with more homomorphic sex chromosomes. We also did not perform
429 extensive testing of the methods on reduced-representation-sequencing (RRS) datasets, but
430 preliminary analyses showed that the noisy DoC distributions across individual RRS loci and
431 possible allelic dropout (Heller et al., 2021) make such data challenging for SATC without further
432 attention to pre-filtering, e.g. on missingness. We consider it outside the scope of this study to
433 investigate how SATC could be improved to take the peculiarities of RRS data into account. We
434 emphasize that our method works on WGS data without any need for sophisticated filtering. For
435 example, we did not exclude repeat annotated regions or remove regions without mapped reads

436 prior to calculating the DoC in any of the species datasets. It is possible that additional filtering of
437 the data could improve the identification of sex-linked scaffolds in some cases, but we focused
438 instead on demonstrating the robustness of the method by showing that it works in challenging
439 data situations. We found that a single set of settings for the different cutoff values—minimum
440 scaffold length, maximum DoC, ratio of homogametic/heterogametic scaffold DoC—yielded
441 usable results for all the species analyzed here. However, the SATC software allows the user to
442 modify these settings if needed. We encourage users to try different cutoff settings to assess the
443 sensitivity of the analyses.

444

445 **Acknowledgements**

446 ABO was supported by a Carlsberg Foundation Reintegration Fellowship (CF19-0427). RH, GGE
447 and CN were supported by a DFF Sapere Aude research grant (DFF8049-00098B), and RH was
448 furthermore supported by an ERC Starting Grant (No 853442). AA and GGE are supported by the
449 Lundbeck foundation (R215-2015-4174). AA is supported by the Novo Nordisk Foundation
450 (NNF20OC0061343). We thank Jonas Meisner for clustering input suggestions. We thank John
451 W. Bickham and J. Andrew DeWoody for providing gray whale data, Peter van der Wolf for gray
452 whales photograph, and Patrícia Chrzanová Pečnerová for providing muskox summary data. We
453 also thank Kristian Ebbesen Hanghøj for helping with writing R code implementation and the
454 PopGen group at University of Copenhagen for helpful comments on previous versions of the
455 manuscript.

456

457

458

459 **Data accession**

460 The datasets analyzed in this study are available at the European Nucleotide Archive under the
461 BioProject accession codes: leopard (PRJEB41230), waterbuck (PRJEB28089) and Darwin's
462 Finches is on short read archive (PRJNA263122 and PRJNA301892). Idxstats files are available
463 for all taxa (leopard, waterbuck, Darwin's finches, gray whale, impala and muskox) at
464 <https://github.com/popgenDK/SATC/tree/main/examples/idxstats>. The genome assemblies used
465 were downloaded from NCBI for goat (*Capra hircus*, ARS1, GCA_001704415.1), domestic cat
466 (*Felis catus*, Felis_catus_9.0, GCA_000181335.4), cow (*Bos taurus*, bosTau8,
467 GCA_000003055.4), sheep (*Ovis aries*, Oar_rambouillet_v1.0, GCA_002742125.1), blue whale
468 (*Balaenoptera musculus*, mBalMus1.pri.v3, GenBank assembly accession GCA_009873245.2),
469 the medium ground finch (*G. fortis*, GCF_000277835.1_GeoFor_1.0), impala (*Aepyceros*
470 *melampus*, IMP GCA_006408695.1), leopard (*Panthera pardus*, PanPar1.0, GCA_001857705.1)
471 and waterbuck (*Kobus ellipsiprymnus*, DFW, GCA_006410655.1). Rasmus Heller contributed the
472 unpublished muskox and impala data. The unpublished gray whale data is contributed by Anna
473 Brüniche-Olsen. The software framework is freely available at
474 <https://github.com/popgenDK/SATC>.

475

476

477

478

479

480

481

482 **Author contributions**

483 The work was conceived by CN and AA. The research design was planned by CN, ABO, GGE,
484 RH and AA. The data was analyzed by CN and GGE. CN, ABO and RH wrote the manuscript
485 with input from GGE and AA. All authors read and approved the final version of the manuscript.

486

487 **References**

488 Bidon T *et al.*, (2014) Brown and polar bear Y chromosomes reveal extensive male-biased
489 gene flow within brother lineages. *Molecular Biology and Evolution* **31**, 1353-1363.

490 Brüniche-Olsen A *et al.*, (2018) Genetic data reveal mixed-stock aggregations of gray
491 whales in the North Pacific Ocean. *Biol. Lett.* 14: 20180399.

492 Charlesworth B, Campos JL, & Jackson BC (2018) Faster-X evolution: Theory and
493 evidence from *Drosophila*. *Molecular Ecology* **27**, 3753-3771.

494 Chen L *et al.*, (2019) Large-scale ruminant genome sequencing provides insights into their
495 evolution and distinct traits. *Science* 364.6446.

496 Clayton D (2008) Testing for association on the X chromosome. *Biostatistics* **9**, 593-600.

497 Dewoody JA *et al.*, (2017) "Characterization of the gray whale *Eschrichtius robustus*
498 genome and a genotyping array based on single-nucleotide polymorphisms in candidate
499 genes." *The Biological Bulletin* 232.3: 186-197.

500 Effron M, Bogart M, Kumamoto A, & Benirschke K (1976) Chromosome studies in the
501 mammalian subfamily Antilopinae. *Genetica* **46**, 419-444.

502 Ellegren H (2009) The different levels of genetic diversity in sex chromosomes and
503 autosomes. *Trends in Genetics* **25**, 278-284.

504 Ellegren H (2014) Genome sequencing and population genomics in non-model organisms.
505 *Trends in Ecology & Evolution* **29**, 51-63.

506 Gallagher Jr D & Womack J (1992) Chromosome conservation in the Bovidae. *Journal of*
507 *Heredity* **83**, 287-298.

508 Gottipati S, Arbiza L, Siepel A, Clark AG, & Keinan A (2011) Analyses of X-linked and
509 autosomal genetic variation in population-scale whole genome sequencing. *Nature*
510 *Genetics* **43**, 741.

511 Grabherr MG *et al.*, (2010) Genome-wide synteny through highly sensitive sequence
512 alignment: Satsuma. *Bioinformatics* **26**, 1145-1151.

513 Grant PR & Grant BR (2020) *How and why species multiply* Princeton University Press.

514 Grath S & Parsch J (2016) Sex-biased gene expression. *Annual Review of Genetics* **50**, 29-
515 44.

516 Graves JAM (2008) Weird animal genomes and the evolution of vertebrate sex and sex
517 chromosomes. *Annual review of genetics* 42: 565-586.

518 Hall AB *et al.*(2013) Six novel Y chromosome genes in Anophelesmosquitoes discovered
519 by independently sequencing males and females. *BMC Genomics* **14**, 273.

520 Hammer MF *et al.*, (2010) The ratio of human X chromosome to autosome diversity is
521 positively correlated with genetic distance from genes. *Nature Genetics* **42**, 830.

522 Harris R.S. (2007) Improved pairwise alignment of genomic DNA. Ph.D. Thesis, The
523 Pennsylvania State University.

524 Heller R *et al.*, (2021). A reference-free approach to analyse RADseq data using standard
525 next generation sequencing toolkits. *Molecular Ecology Resources*, 21(4), 1085–1097.

526 Irwin DE (2018) Sex chromosomes and speciation in birds and other ZW systems.
527 *Molecular Ecology* **27**, 3831-3851.

528 Kim S *et al.*, (2016). Comparison of carnivore, omni- vore, and herbivore mammalian
529 genomes with a new leopard assembly. *Genome Biol.* *17*, 211.

530 Kingswood S, Kumamoto A, Charter S, Aman R, & Ryder O (1998) Brief communication.
531 Centric fusion polymorphisms in waterbuck (*Kobus ellipsiprymnus*). *Journal of Heredity*
532 **89**, 96-100.

533 Kumamoto A, Charter S, Houck M, & Frahm M (1996) Chromosomes of *Damaliscus*
534 (*Artiodactyla*, *Bovidae*): simple and complex centric fusion rearrangements. *Chromosome*
535 *Research* **4**, 614-621.

536 Lambert CA *et al.*, (2010) Highly punctuated patterns of population structure on the X
537 chromosome and implications for African evolutionary history. *The American Journal of*
538 *Human Genetics* **86**, 34-44.

539 Lamichhaney S *et al.*, (2015) Evolution of Darwin's finches and their beaks revealed by
540 genome sequencing. *Nature* **518**, 371-375.

541 Lamichhaney S *et al.*(2017) X chromosome-wide association study identifies a
542 susceptibility locus for inflammatory bowel disease in Koreans. *Journal of Crohn's and*
543 *Colitis* **11**, 820-830.

544 Li H & Durbin R (2011) Inference of human population history from individual whole-
545 genome sequences. *Nature* **475**, 493-U484.

546 Li H *et al.*, (2009) The sequence alignment/map format and SAMtools. *Bioinformatics* **25**,
547 2078 - 2079.

548 Luciano M *et al.* (2019) The influence of X chromosome variants on trait neuroticism.
549 *Molecular Psychiatry*, 1-9.

550 Malde K, Skern R, & Glover KA (2019) Using sequencing coverage statistics to identify
551 sex chromosomes in minke whales. *arXiv preprint arXiv:1902.06654*.

552 Meisel RP & Connallon T (2013) The faster-X effect: integrating theory and data. *Trends*
553 *in Genetics* **29**, 537-544.

554 Pagacova E, Cernohorska H, Kubickova S, Vahala J, & Rubes J (2011) Centric fusion
555 polymorphism in captive animals of family Bovidae. *Conservation Genetics* **12**, 71-77.

556 Palmer DH, Rogers TF, Dean R, & Wright AE (2019) How to identify sex chromosomes
557 and their turnover. *Molecular Ecology* **28**, 4709-4724.

558 Pečnerová P *et al.*, (2021) High Genetic Diversity and Low Differentiation Reflect the
559 Ecological Versatility of the African Leopard. *Current Biology: CB* 31 (9): 1862–71.e5.

560 Pennell MW, Mank JE, & Peichel CL (2018) Transitions in sex determination and sex
561 chromosomes across vertebrate species. *Molecular Ecology* **27**, 3950-3963.

562 Pérez-González J *et al.*, (2014) Males and females contribute unequally to offspring genetic
563 diversity in the polygynandrous mating system of wild boar. *Plos One* **9**, e115394.

564 Presgraves DC (2018) Evaluating genomic signatures of “the large X-effect” during
565 complex speciation. *Molecular Ecology* **27**, 3822-3830.

566 Scrucca L, Fop M, Murphy TB, & Raftery AE (2016) mclust 5: clustering, classification
567 and density estimation using Gaussian finite mixture models. *The R journal* **8**, 289.

568 Smeds L *et al.*, (2015) Evolutionary analysis of the female-specific avian W chromosome.
569 *Nature Communications* **6**.

570 Team RDC (2019) R: A language and environment for statistical computing, Vienna,
571 Austria.

572 Zhang G *et al.*, (2014) Comparative genomics reveals insights into avian genome evolution
573 and adaptation. *Science* 346(6215):1311-20.

574 Zuo L *et al.*, (2013) Sex chromosome-wide association analysis suggested male-specific
575 risk genes for alcohol dependence. *Psychiatric Genetics* **23**, 233.

576

577

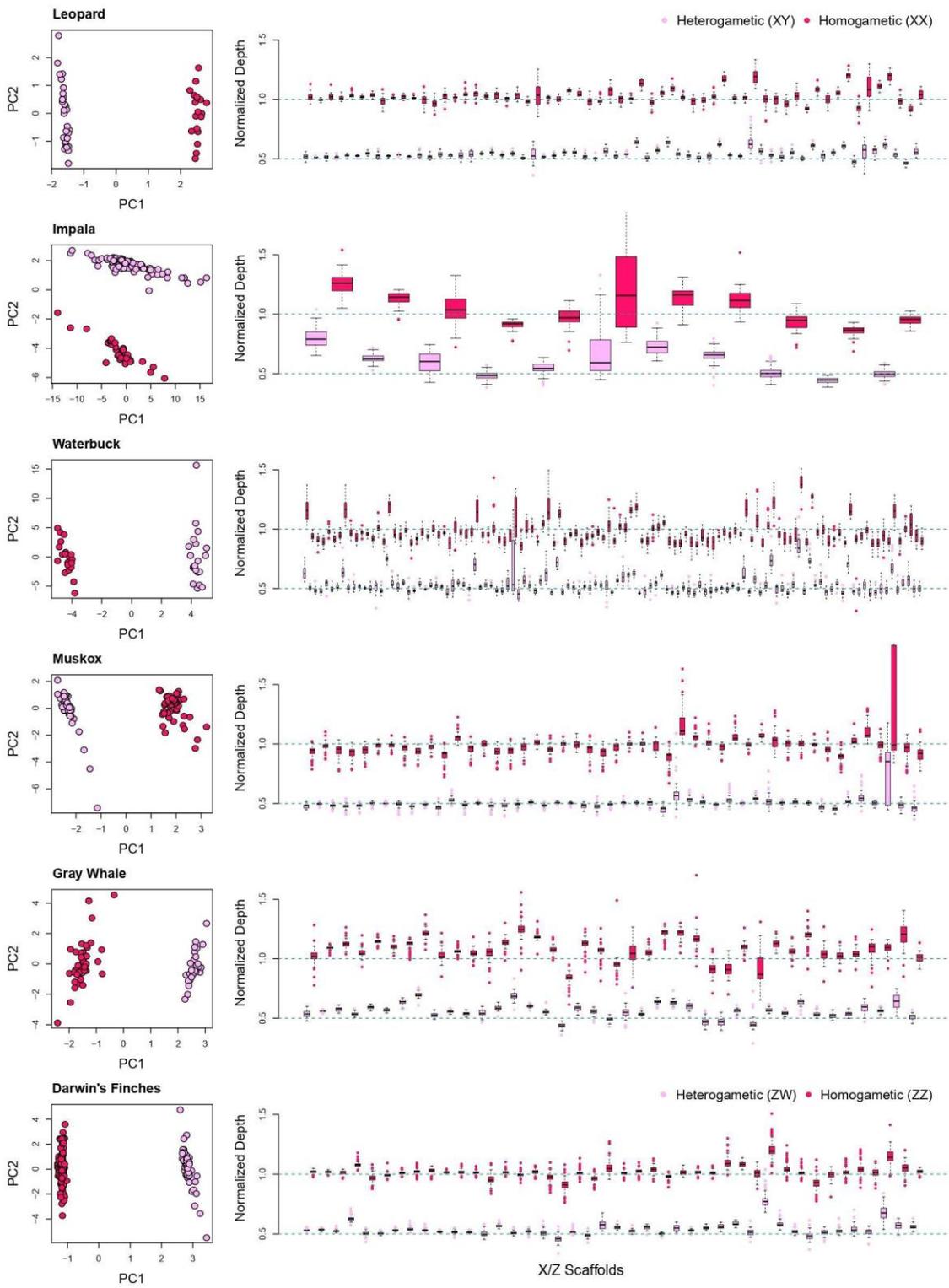
578 **Table 1.** Basic properties of the species analyzed with our SATC method. For each species the
579 number of samples (*N*), depth of coverage (DoC), total number of scaffolds (#scaffolds), number
580 of inferred X or Z scaffolds based on mean difference of normalized DoC being between 0.4 and
581 0.6 (#scaffolds X/Z), total numbers of sex-linked scaffolds based on t-test (#scaffolds sex-linked),
582 length of inferred X/Z scaffolds (X/Z (Mb), total length of all identified sex-linked scaffolds (Mb),
583 and the sex ratio for the samples. Inferred sex was estimated based on Gaussian mixture clustering
584 from the top 2 PCs inferred from closely related species with a chromosomal level assembly. DoC
585 is the average number of reads on each position in the genome calculated by summing all mapped
586 reads multiplied by average read length and divided by total length of scaffolds.

Species	<i>N</i>	DoC	#scaffolds	#scaffolds X/Z	#scaffolds sex-linked	X/Z (Mb)	Sex-linked scaffolds (Mb)	sex ratio (hetero/homo)
Leopard	49	4.2X	50,378	60	66	113	126	30/19
Impala	11 3	3.1X	24,159	11	589	2.1	143	81/32
Muskox	10 3	11.4X	7,072	47	54	126	128	47/56
Gray whales	73	5.8X	2,796	39	62	104	131	46/27

Waterbuck	40	3.4X	88,935	113	400	46.3	187	20/20
Darwin's finches	17 2	13.8X	27,240	42	62	72.8	78.5	122/50

587

588



589

590

591

592 **Figure 1: PCA plots with Gaussian mixture clustering and boxplots.**

593 Left column: PCA plots of normalized DoC across all scaffolds and samples from 5 mammalian
594 species and the 15 species making up Darwin's finches species complex. Two clusters are inferred
595 homogametic (dark pink) and heterogametic (light pink). Right column: Boxplot of normalized
596 DoC from inferred X/Z scaffolds based on mean difference of two sex clusters within range of 0.4
597 and 0.6. Scaffolds are sorted based on their length (x-axis). Each scaffold is represented by two
598 boxplots from homogametic and heterogametic groups. Expected median values for each group
599 are shown by horizontal green dashed lines of 0.5 (heterogametic) and 1.0 (homogametic).

