



21 <sup>10</sup>School of Biological Sciences, University of Auckland, Auckland, New Zealand.

22 \*Correspondence to: d.lambert@griffith.edu.au.

23

24 **Abstract:** Microsatellites are widely used in population genetics, but their evolutionary  
25 dynamics remain poorly understood. It is unclear whether microsatellite loci drift in length over  
26 time. We identify more than 27 million microsatellites using a novel and unique dataset of  
27 modern and ancient Adélie penguin genomes along with data from 63 published chordate  
28 genomes. We investigate microsatellite evolutionary dynamics over two time scales: one based  
29 on the Adélie penguin samples dating to approximately 46.5 kya, the other dating to the  
30 diversification of chordates more than 500 Mya. We show that the process of microsatellite allele  
31 length evolution is at dynamic equilibrium; while there is length polymorphism among  
32 individuals, the length distribution for a given locus remains stable. Many microsatellites persist  
33 over very long time scales, particularly in exons and regulatory sequence. These often retain  
34 length variability, suggesting that they may play a role in the maintenance of evolutionary  
35 plasticity.

36

## 37 **Introduction**

38           Microsatellites or short tandem repeats (STRs), consisting of tandem repeats of two to six  
39 base pair motifs, are prevalent in both prokaryotic and eukaryotic genomes. Some microsatellites  
40 have been shown to be functionally important<sup>1-3</sup>, but most are assumed to evolve neutrally, and  
41 for this reason, along with their abundance and high variability, they have been used extensively  
42 in population genetics studies<sup>4</sup>. However, their evolutionary dynamics remain poorly understood,  
43 and it is unclear whether microsatellite loci are in dynamic equilibrium with respect to the length  
44 of alleles, or whether alleles experience directional drift in length. This is important because the  
45 mutation processes that underlie these important genetic markers are central to the evolutionary  
46 models that employ microsatellites.

47           In this study, when describing microsatellites, we consider both the number of base pairs  
48 in the underlying motif *period* and, for each allele, how many times the motif appears (the *repeat*  
49 *number*). We refer to microsatellites that contain only exact copies of the motif as *pure*. The total  
50 length of a microsatellite allele (in nucleotides) is the product of the period and repeat number.  
51 Repeat number is thought to change through a process of replication slippage<sup>5,6</sup>, by which  
52 strands may transiently dissociate during DNA replication and then mispair with a different copy  
53 of the repeat, resulting in the insertion or deletion of one or more repeat units. Microsatellites are  
54 highly plastic in evolutionary terms, with mutation rates due to replication slippage generally  
55 several orders of magnitude higher than for point mutation<sup>7</sup>.

56           Much still remains to be learned about the mutational processes involved in microsatellite  
57 evolution. The overall process can be thought of as a birth–death process (increase or decrease in  
58 length of microsatellite by the birth or death of individual repeat units) embedded within a  
59 second birth–death process (microsatellite loci appear and disappear) all happening along a

60 branching process (the population history). Slippage during DNA replication is thought to be the  
61 main cause of changes in length, with mismatch repair reducing the mutation rate<sup>8</sup>, but  
62 recombination may also play a role, and point mutations must be taken into account. The  
63 processes by which new microsatellites appear, and by which they eventually degenerate and  
64 disappear, are particularly poorly understood<sup>9</sup>.

65 Existing models are highly simplified and only take into account changes in length (and  
66 occasionally purity) in existing microsatellites, ignoring the processes of ‘birth’ and ‘death’ by  
67 which microsatellite loci appear and disappear (and perhaps reappear)<sup>10</sup>. Some of these models  
68 have been designed so that they have a stationary distribution (for example those of Kruglyak *et*  
69 *al.*<sup>11</sup>, Calabrese *et al.*<sup>12</sup> and Amos *et al.*<sup>13</sup>), but it is not clear whether this is biologically realistic.  
70 It may be that an individual microsatellite locus is never at equilibrium, tending instead to  
71 increase in length throughout its life, but that the birth-death process causes the genome-wide  
72 distribution of allele lengths at all microsatellite loci to be at equilibrium.

73 A key open question is thus whether the alleles at a microsatellite locus increase or  
74 decrease in average length over time, or whether each locus is maintained at an equilibrium  
75 length. While some pedigree studies have shown a bias in favor of gain of repeats<sup>14</sup>, suggesting  
76 that microsatellites should rapidly increase in size<sup>15</sup>, others have found that slippage has a  
77 length-dependent bias<sup>16-18</sup>, supporting earlier suggestions that constraints exist on repeat number  
78 at microsatellite loci<sup>19</sup>. On the basis of the former observation, it has been suggested that  
79 microsatellites increase in length until the accumulation of point mutations hinders slippage and  
80 ultimately leads to the degeneration of the microsatellite locus<sup>10,11</sup>. Alternatively, Amos *et al.*<sup>13</sup>  
81 recently proposed a model consistent with the latter observations, in which inter-allelic  
82 interactions in heterozygous individuals may drive the process whereby longer-than-average

83 alleles tend to get shorter and shorter-than-average alleles tend to get longer (which they call the  
84 centrally directed mutation model).

85 Here we make use of exceptionally well-preserved ancient DNA from a unique set of  
86 Adélie penguin samples reported here for the first time, and genotype 177,974 microsatellites in  
87 both modern and ancient genomes, including some dating to approximately 46.5 kya. In addition,  
88 we are able to time the evolutionary origin of many of these loci by aligning them with more  
89 than 27 million microsatellites from a large set of published chordate genomes and mapping  
90 them onto a recent phylogeny<sup>20</sup>. Our data include microsatellites that date to the diversification  
91 of chordates more than 500 Mya. We show that allele lengths at microsatellite loci are in  
92 dynamic equilibrium, and these have remained stable over hundreds of millions of years and  
93 through many speciation events. While there is length polymorphism among individuals, the  
94 overall length distribution for a given locus does not change appreciably over time. We show that  
95 microsatellites can persist over very long time scales, particularly those in exons and regulatory  
96 sequence, while retaining length variability. This suggests that microsatellites may play a role in  
97 the maintenance of evolutionary plasticity.

98

## 99 **Results**

### 100 Microsatellite dynamics in Adélie penguin samples

101 Genomes obtained from ancient biological remains allow us to observe changes in  
102 sequence variation that cannot be observed using only contemporary sequences. Here we have  
103 used whole-genome sequence data of ancient Adélie penguin remains from 23 individuals dated  
104 at up to 46,587 years old, as well as from 26 modern individuals, to identify 177,974

105 microsatellite loci in an Adélie penguin reference genome. Most loci are close to the minimum  
106 length detectable for each period (especially in the case of pure loci), with very small numbers of  
107 longer loci up to thousands of base pairs in length. The length distributions of these  
108 microsatellite loci are shown in Supplementary Fig. 1. We determined the genotype of these loci  
109 in each of the ancient and modern Adélie samples, and allele length distributions for each sample  
110 are shown in Supplementary Fig. 2. These genotype data enable us to obtain length distributions  
111 for microsatellites at different time points, and hence to test whether there is any evidence for  
112 directional drift in microsatellite length.

113         To test whether microsatellite allele length is stationary, or whether the average allele  
114 lengths of individual loci increase over time, we used BayesFactor<sup>21</sup> to compare generalized  
115 linear mixed models in which allele length is treated as dependent on different combinations of  
116 possible explanatory variables. The explanatory variables considered were: the motif of the  
117 allele, the surrounding sequence type (exon, intron, regulatory, or intergenic), and sample age.  
118 We also tested for an interaction between surrounding sequence type and sample age. In addition  
119 to these fixed effects, which are assumed to be the same for all genomes, we also treated the  
120 sample, i.e., the particular Adélie genome, as a random effect; this is equivalent to allowing a  
121 different intercept in the regression model for each genome. Impurity affects the length at which  
122 microsatellites can be detected, so models were fit separately for pure and impure microsatellites.  
123 Similarly, models were fit separately for loci of different periods because different alignment  
124 score thresholds were used to detect them, so that their allele lengths cannot be compared  
125 directly. Bayes factors for all models tested are given in Supplementary Table 1, and posterior  
126 estimates of effect sizes in Supplementary Table 2. For both pure and impure microsatellites of  
127 each period, the best-supported model is that in which allele length is dependent on the motif and

128 surrounding sequence type. Our data provide positive evidence for this model, being at least  
129 seven times more likely to be observed under this model than under a model in which length  
130 depends on sample age. Since length does not depend on sample age in this model, we infer that  
131 the process of expansion and contraction of microsatellite alleles is effectively stationary, or  
132 nearly stationary, over a time-scale of tens of thousands of years.

133

#### 134 Microsatellite locus age inference

135 To investigate microsatellite dynamics over a much longer timescale and across a broad  
136 range of species, we used whole-genome alignments of 48 avian species from Zhang *et al.*<sup>22</sup>  
137 along with the genomes of fifteen non-avian vertebrate species that span the chordate tree. We  
138 identified a total of over 27 million microsatellites in the 63 genomes, and a breakdown of the  
139 numbers of loci of each period detected in each genome is given in Supplementary Table 3. For  
140 each of these species, we used a whole-genome alignment to chicken to generate a standard set  
141 of coordinates for all microsatellites present in the alignment. The number of microsatellite loci  
142 in any species that can be aligned to the chicken genome, and the overall number of bases  
143 aligned to the chicken genome, are negatively correlated with the time since the most recent  
144 common ancestor of that species and chicken (see Supplementary Fig. 3). We were able to map  
145 approximately 5.4 million microsatellites across the 63 species to almost 2.9 million loci in the  
146 chicken genome. Of these, almost 2.2 million microsatellite loci were found in only a single  
147 species and 680,804 loci had microsatellites conserved across two or more species. Exact  
148 numbers of microsatellites detected and aligned are given in Supplementary Table 4.

149 We used the dated avian whole-genome phylogeny published by Jarvis *et al.*<sup>20</sup>, to which  
150 we added 15 non-avian species with estimated divergence times taken from the Timetree of Life

151 ([www.timetree.org](http://www.timetree.org))<sup>23</sup>. To infer gains and losses of microsatellite loci in different lineages, we  
152 carried out ancestral state reconstruction on a subtree whose topology is relatively  
153 uncontroversial, agreeing with the trees published by Jarvis *et al.*<sup>20</sup> and Prum *et al.*<sup>24</sup>, and on  
154 which we expect incomplete lineage sorting events to be rare<sup>25</sup>. This allows us to infer the edge  
155 on which any locus present in Adélie penguin was gained, and hence to estimate the ages of  
156 these loci. Distributions of estimated ages for loci in intergenic, intronic, exonic, and regulatory  
157 sequence are shown in Supplementary Fig. 4. Supplementary Fig. 5 shows the numbers of  
158 inferred gains and losses of microsatellites on each edge of the subtree, scaled according to both  
159 the length of the edge and the amount of sequence that can be aligned to the chicken genome.  
160 The total numbers of extant microsatellite loci whose origins were inferred to pre-date selected  
161 ancestral nodes are given in Supplementary Table 5.

162         The relative densities of microsatellite loci (including both pure and impure  
163 microsatellites) in different types of sequence for different age brackets are shown in Fig. 1.  
164 While the older age brackets contain fewer loci overall, those loci are much more likely to be  
165 found in regulatory or coding sequences. The percentage of loci found in regulatory or coding  
166 sequence for each bracket is shown in Supplementary Table 6. Microsatellite loci in regulatory  
167 or coding sequences thus appear to be conserved over longer periods on average than those in  
168 intergenic sequence or introns. This suggests that they are maintained by selection, be it directly  
169 for the presence of a microsatellite or for the surrounding sequence. Total numbers of loci  
170 genotyped in the Adélie penguin samples for each age bracket are given in Supplementary Table  
171 7, along with the percentages of loci at which we observe multiple genotypes, showing that these  
172 loci retain length variability in Adélie penguins.

173



174 Microsatellite dynamics through deep time

175           To test whether the process of microsatellite mutation results in allele length distributions  
176 that are stationary over evolutionary time-scales (millions of years), we used BayesFactor as  
177 described above, replacing the sample age parameter with the locus age estimate. Bayes factors  
178 for all models tested are given in Supplementary Table 8. For all subsets of the data comprising  
179 pure and impure microsatellites of each period, the best-fitting model for allele length is  
180 dependent on motif, surrounding sequence type, locus age, and an interaction between  
181 surrounding sequence type and locus age. In all cases, the data provide very strong evidence for  
182 this model, being more likely under this model than under any other by a factor of at least  $10^{12}$ .  
183 We sampled from the posterior distribution of the full model for each subset to obtain posterior  
184 estimates of effect sizes, and these are shown in Supplementary Table 9. The effect of locus age  
185 is shown separately in Table 1. The effect sizes are very small (on the order of one nucleotide per  
186 hundred million years). For loci of periods 2 and 3, we also tested interactions between motif and  
187 locus age, and between motif and surrounding sequence type for subsets of our data, and found  
188 strong evidence for these interactions. We were unable to test these interactions for loci of longer  
189 periods because of the rapid increase in numbers of motifs as period increases.

190           Distributions of allele lengths for loci of different ages in different types of surrounding  
191 sequence are shown in Fig. 2. In agreement with the results of the linear mixed-model, a very  
192 slow increase in mean allele length over time can be seen for microsatellites in intron and  
193 intergenic sequence. Overall, pure di- and tetranucleotide loci in protein-coding sequence have  
194 significantly shorter mean allele lengths than those in non-coding sequence, while impure tri-  
195 and hexanucleotide microsatellites in protein-coding sequence have significantly longer mean  
196 allele lengths than those in non-coding sequence (see Table 2). It is likely that selection against

197 frameshift mutations in coding sequence limits microsatellite expansion when the period is not a  
198 multiple of three<sup>26</sup>.

199

## 200 **Discussion**

201 To summarize our results, the mean allele length at any given microsatellite locus  
202 changes very little, on scales ranging from a few thousand to hundreds of millions of years, with  
203 estimated effect sizes on the order of one nucleotide per hundred million years. There is a  
204 gradual increase in allele length variation over time, as can be seen in Fig. 2. This suggests that  
205 the replication slippage process that generates length polymorphism is in a dynamic equilibrium,  
206 such that increases and decreases in length remain approximately balanced. These results are  
207 consistent with the findings of Sun *et al.*<sup>18</sup> that longer alleles tend to decrease in length and  
208 shorter alleles tend to increase. We recommend that population geneticists and ecologists use  
209 models of microsatellite evolution that have stationary distributions, such as those of Kruglyak *et*  
210 *al.*<sup>11</sup>, Calabrese *et al.*<sup>12</sup> or Amos *et al.*<sup>13</sup>, rather than those, such as the stepwise mutation model<sup>27</sup>,  
211 that allow allele lengths to drift upwards indefinitely.

212 We have also shown that microsatellites can persist, and remain variable, over very long  
213 periods of evolutionary time, with 257 extant microsatellite loci dating from before the origin of  
214 chordates, and 3,938 pre-dating the divergence of mammals and reptiles. Although we observe a  
215 slight decrease in heterozygosity with locus age (Supplementary Fig. 6), nevertheless, we  
216 observe multiple alleles in the Adélie samples for many ancient loci (Supplementary Table 7).  
217 The microsatellite loci that persist over very long periods are more often found in coding  
218 sequence and in regulatory regions. A disproportionate number of these variable ancient loci are

219 trimer repeats located in protein-coding genes, which must code for a homopolymer run of  
220 amino acids. These trimer repeats in coding sequences make up only 0.55% of all loci that are  
221 variable in our Adélie samples, but 5.67% of variable loci that pre-date the divergence of extant  
222 birds, and 9.86% of those that pre-date the divergence of mammals and reptiles. It seems likely  
223 that selection is acting to maintain variability at these loci, which could act as mediators of rapid  
224 phenotypic change<sup>2</sup>.

225         A limitation of using short read data is that longer alleles are effectively censored from  
226 our data; however, as can be seen in Supplementary Fig. 1, the overwhelming majority of loci  
227 are much shorter (in the reference genome) than the read length. In addition, the reads from  
228 ancient Adélie samples are shorter than those from modern samples. This means that longer  
229 alleles are less likely to be genotyped in the ancient samples, and therefore we would expect this  
230 to give a signal for increasing allele length over time. However, we do not observe any such  
231 signal despite this potential bias, presumably because any such signal is swamped by the much  
232 larger number of shorter loci. As long-read sequencing becomes more common, and as methods  
233 for genotyping microsatellites in long-read data are developed, it may become feasible to verify  
234 our results for a more complete data set.

235

## 236 **Materials and Methods**

### 237 Contemporary Adélie penguin samples

238         Blood samples from Adélie penguins were collected from individuals at active breeding  
239 colonies, using methods as described in Millar *et al.*<sup>28</sup>, in six locations around Antarctica:  
240 Tongerson Island (AP samples) the Mawson region (B samples), Cape Adare (CA), Cape Bird

241 (CB), Coulman Island (CI), and Inexpressible Island (II). Collection and sequencing information  
242 is given in Supplementary Table 10.

243

#### 244 Ancient Adélie penguin samples

245 Sub-fossil bones were collected in abandoned nests discovered along coastal ice-free areas  
246 both in the vicinity of presently occupied colonies and in relict colonies discovered in sites where  
247 penguins do not breed at present<sup>29-31</sup> (Supplementary Table 11). Ornithogenic soils were  
248 stratigraphically excavated to find penguin bones and other remains as described previously<sup>32,33</sup>.

249 Radiocarbon AMS dates were supplied by NOSAMS, Woods Hole Oceanographic Institute,  
250 the New Zealand Institute of Geological and Nuclear Sciences, Lower Hutt, New Zealand, and  
251 Institut for Fysik og Astronomi, Aarhus Universitet, Denmark. Radiocarbon dates were  
252 calibrated with CALIB 7.1 (<http://calib.qub.ac.uk/calib/>)<sup>34</sup> using the Marine Reservoir Correction  
253 Database 2013 and applying a delta-*R* of  $791 \pm 121$  [<sup>35</sup>]. Mean ages and 2 delta standard  
254 deviation values were considered.

255

#### 256 Modern DNA extraction

257 For the 26 modern Adélie penguin samples, genomic libraries were prepared by first  
258 extracting DNA from Seutin-preserved blood or ethanol-preserved soft tissue samples. DNA was  
259 then purified using Qiagen DNEasy spin-columns according to the manufacturer's protocol  
260 (Qiagen, Valencia, CA, USA) and eluted in 100  $\mu$ L UltraPure™ water (Life Technologies,  
261 Grand Island, New York, USA).

262

#### 263 Ancient DNA extraction

264 All laboratory work with ancient Adélie penguin samples prior to PCR-amplification of  
265 genomic libraries (see below) was carried out in a physically isolated laboratory used only for  
266 ancient DNA work, following strict guidelines to minimize external contamination. Designated  
267 blank samples consisting originally of 200  $\mu$ L digestion buffer were carried sequentially through  
268 all DNA extraction and library building procedures at a minimum ratio of one blank for every  
269 eight samples.

270 DNA was extracted from ancient bone or muscle tissue samples by first digesting ca. 0.1 g  
271 bone/tissue shavings in 200  $\mu$ L digestion buffer (consisting of 180  $\mu$ L of 0.5 M EDTA, 10  $\mu$ L of  
272 10% N-lauryl sarcosine, 10  $\mu$ L of 20mg/mL proteinase K) for 12–18 hours at 55°C with  
273 rotational mixing (ca. 10 rpm). This was followed by 2–5 rounds of organic extraction with 1–  
274 1.5 mL ultra-pure buffer-saturated phenol and one round of extraction with 1–1.5 mL chloroform  
275 (Sigma-Aldrich, St. Louis, MO, USA). Extracts were purified with Qiagen MinElute or PCR  
276 Purification columns using high concentration buffer PB or PE (10:1 buffer:sample volume ratio)  
277 to improve retention of small fragments, and 2 $\times$  spin-through centrifugation for sample  
278 application and elution stages to further maximize yield. Final elutions were completed in a  
279 volume of 22  $\mu$ L UltraPure™ water or NEB buffer EB.

280

### 281 DNA library construction

282 Purified extracted DNA of modern samples was quantified with a Qubit 2.0 Fluorometer  
283 and dsDNA HS Assay kit (Life Technologies, Grand Island, NY, USA) and ca. 0.5–1.5  $\mu$ g of  
284 DNA was sheared with a Covaris sonicator (Covaris, Woburn, MA, USA) to an average size of  
285 300–600 bp (base pairs). Sheared extracts were adapter-ligated and enriched using the standard  
286 NEBNext or NEBNext Ultra protocol (catalog #E6040 and #E7370) and NEBNext multiplex

287 Illumina primers (catalog #E7335) (NEB, Ipswich, MA, USA) in ½-size recommended reaction  
288 volumes for end-preparation, adapter ligation, and enrichment reactions. Enrichments were  
289 performed under recommended cycling conditions, with 10–14 cycles of enrichment for each  
290 sample and using Phusion High-Fidelity Master Mix (NEB catalog #M0531). Samples were  
291 submitted for 101 bp paired-end sequencing on an Illumina HiSeq2000 at BGI-Hong Kong,  
292 using one or ca. 1.33 lanes for each sample.

293 Genomic libraries for ancient samples were built following two strategies. Library building  
294 for all Holocene samples and initial attempts for two late Pleistocene samples (CB070121.08,  
295 CB070121.16) were completed based on Meyer *et al.*<sup>36</sup> with minor adjustments. Based on low  
296 endogenous yields for the two late Pleistocene samples, a second attempt at library building was  
297 made for all three late Pleistocene samples (CB070121.08, CB070121.13, CB070121.16)  
298 following the NEBNext Ultra protocol, and using ½-size reactions for end-preparation, adapter  
299 ligation and enrichment reactions. For all ancient samples, enrichment reactions were completed  
300 by mixing ca. 11.5 µL of the heat-inactivated adapter-ligation reaction, 0.5 µL each of 25 µM  
301 NEBNext index and universal primers, and 12.5 µL 2× Phusion Hi-Fidelity Master Mix.  
302 Enrichment reactions were carried out under recommended cycling conditions, with 12–22  
303 cycles of enrichment for each sample.

304 Finished ancient libraries were purified using Axygen MAG-PCR SPRI beads (Corning  
305 Life Sciences, Tewksbury, MA, USA) at a ratio of 0.7-1.1:1 Axygen:sample volume to minimize  
306 concentration of potential adapter dimers<sup>37</sup> and quantified with a Qubit 2.0 Fluorometer.  
307 Libraries were submitted for 101bp single-end (SE) sequencing on an Illumina HiSeq2000 to  
308 either BGI-Hong Kong or the National High-Throughput Sequencing Center (University of  
309 Denmark, <http://seqcenter.ku.dk/>), using from between two and 10.5 lanes of sequencing for

310 each sample with resultant genome-wide average sequencing depths of ca. 22× and 8× for  
311 modern and ancient samples, respectively (Supplementary Tables 10 and 11).

312

### 313 Alignment

314 For all sequence pools, adapter sequences were trimmed from reads using Cutadapt<sup>38</sup> v. 1.1  
315 under default parameters. Low-quality reads were filtered with Trimmomatic<sup>39</sup> v. 0.22, with  
316 minimum trailing and leading quality of 20, average quality over 20bp sliding windows of 20,  
317 and minimum lengths of 80bp for modern reads and 30bp for ancient reads. Trimmed and  
318 filtered Illumina reads for each Adélie penguin sample were mapped to the Adélie reference  
319 genome<sup>22</sup> using Bowtie2<sup>40</sup> with the ‘--very-sensitive’ preset option.

320

### 321 Genomes

322 In this study we use the 48 avian genomes reported by Jarvis *et al.*<sup>20</sup>. We also use the  
323 pairwise alignments to the chicken genome that were used by Jarvis *et al.* in generating their  
324 whole-genome multiple alignment. This consists of a set of pairwise alignments for each species  
325 with each individual chromosome of the chicken genome as reference.

326 In addition, we use genomes of the fifteen non-avian species for which whole-genome  
327 alignments to the chicken galGal3 assembly are available from the UCSC genome browser.  
328 These are: human (hg19), chimpanzee (panTro3), orangutan (ponAbe2), mouse (mm9), rat (rn4),  
329 guinea pig (cavPor3), horse (equCab2), opossum (monDom5), platypus (ornAna1), lizard  
330 (anoCar2), frog (xenTro3), zebrafish (danRer4), fugu (fr2), lamprey (petMar1), and lancelet  
331 (braFlo1). All genomes used are listed in Supplementary Table 12.

332

333 Microsatellite detection

334 Microsatellite loci were identified in all 63 genomes using Tandem Repeats Finder (TRF)<sup>41</sup>  
335 with the following parameters: match weight 2; mismatch weight 7; indel weight 7; matching  
336 probability 80; indel probability 10; minimum alignment score 18; maximum period size 6. The  
337 results were then filtered using the alignment score thresholds shown in Supplementary Table 13,  
338 taken from Willems *et al.*<sup>42</sup>. This gave us five sets of microsatellites for each species: for dimer,  
339 trimer, tetramer, pentamer and hexamer repeats, with their respective score thresholds.

340 Microsatellite loci were compared against the annotations for all the avian genomes, to  
341 determine which loci fall within protein coding sequences or introns. Putative regulatory regions  
342 were identified by extracting the set of conserved nonexonic elements identified in the chicken  
343 genome by Lowe *et al.*<sup>43</sup> from each of the avian genome alignments. All remaining sequences  
344 were assumed to be intergenic.

345 Microsatellites identified in the Adélie penguin reference genome using TRF were  
346 genotyped in the Adélie penguin samples using RepeatSeq<sup>44</sup> (which requires a list of pre-  
347 identified loci and sequence reads as input), and the output formatted as tables for analysis in  
348 R<sup>45</sup>. Tables of genotype calls were imported into R and summary statistics calculated for each  
349 locus, including the mode, mean and standard deviation of the allele lengths observed in the  
350 samples, and the number of alleles observed. These were combined with the TRF output  
351 containing the motif, purity and nucleotide composition of the locus in the reference genome.

352

353 Homology matching

354 First, for each species and period, we coded the microsatellite loci detected above as  
355 features in a general feature format file. Next, we used MafFilter<sup>46</sup> to extract these features from



356 the pairwise alignment between the species in question and each chicken chromosome, and  
357 output the coordinates that each feature aligns to in the chicken genome. A custom R script was  
358 used to produce a table matching each set of chicken coordinates to the corresponding  
359 microsatellite locus. Motifs were standardized by calculating the lexicographically minimal  
360 rotation to allow for loci to begin at different positions within a repeat unit (e.g., the motif TGA  
361 was standardized as ATG).

362 For each chicken chromosome and period, we combined the motif tables for all 63 species,  
363 and used a custom Java program to assign similarity scores to pairs of loci based on the distance  
364 between them (in terms of chicken coordinates) and the similarity of their motifs. Loci were  
365 scored if they were no more than 60 bp apart and their motifs differed by no more than one  
366 substitution. Testing different values of the length threshold showed that larger values did not  
367 increase the numbers of homologous loci detected. In addition, we manually checked a small  
368 sample of loci to verify that the loci detected were indeed homologous. Similarity was calculated  
369 as

$$377 \quad sim = \frac{p}{1 + d}$$

370 where  $p$  is the proportion of sites in the motif that are identical, and  $d$  is the distance in base pairs  
371 between the loci (zero if the loci overlap). We then used the Markov Cluster Algorithm (MCL)<sup>47</sup>  
372 with the `--abc` input option and default settings to identify clusters of putatively homologous  
373 loci (across all 63 genomes). These clusters were converted into a matrix with the 48 species as  
374 columns and locations as rows, containing the motif for each species where a microsatellite is  
375 present. The matrix was also output as a presence/absence matrix, with ones where a  
376 microsatellite is present and zeroes otherwise.

378 To avoid any false negatives where a given region is not represented in the alignment for  
379 some species, we checked the local region of the alignment for any species missing from a given  
380 cluster, and recoded them as unknown ('?', as opposed to '0' for absent) in the presence/absence  
381 matrix if the region was not covered in the alignment.

382

### 383 Ancestral state reconstruction

384 We used the R package phangorn<sup>48</sup> to perform ancestral state reconstruction on the timetree  
385 reported in Jarvis *et al.*<sup>20</sup> using our presence/absence matrices. The maximum likelihood  
386 reconstructions available do not allow non-reversible models (i.e. the rates of gains and losses  
387 are assumed to be equal), so we used the “ACCTRAN” parsimony method. Numbers of gains  
388 and losses of homologous microsatellites inferred for each edge were then counted, ignoring any  
389 changes from a known state to unknown.

390 We also calculated the numbers of microsatellite losses required under a Dollo process,  
391 where any microsatellite locus only ever arises once, but may be lost in multiple lineages.  
392 However, the results were not appreciably different to those obtained under parsimony.

393

### 394 Adélie locus age determination

395 Minimum ages were calculated for loci present in the Adélie penguin genome by using the  
396 ancestral state reconstruction results to identify the most recent gain of the locus on the path from  
397 the root to Adélie. This allows for loci being gained independently in different lineages, or lost  
398 and re-gained. These locus ages were combined with the genotype statistics calculated above,  
399 allowing us to examine the relationship between locus age, length, purity, and surrounding  
400 sequence type.

401

## 402 Model fitting

403 The ‘**generalTestBF**’ function of the R package BayesFactor<sup>21</sup> was used to fit generalized  
404 linear mixed models to the ancient and modern Adélie genotype data. A Bayes Factor (BF) is a  
405 measure that quantifies the evidence for a hypothesis compared to an alternative hypothesis  
406 given the data. The following thresholds have been suggested to quantify the evidence for one  
407 hypothesis over another as reported by BFs: BF < 3: insignificant, BF 3–20: positive, BF 20–  
408 150: strong, BF > 150 very strong<sup>49</sup>.

409 We tested the dependence of microsatellite allele length on sample age, motif, sample,  
410 surrounding sequence type, and an interaction between sample age and surrounding sequence  
411 type for all loci genotyped in Adélie. For those loci for which we were able to estimate the age  
412 (i.e., those that were alignable to the chicken genome), we tested the dependence of allele length  
413 on estimated locus age, motif, sample, surrounding sequence type, and an interaction between  
414 locus age and surrounding sequence type. In both cases, the sample was treated as a random  
415 effect, and all other variables as fixed effects. Sample age and Locus age variables were centred.  
416 Models were fit separately for pure and impure microsatellite loci of each period (2 to 5). To  
417 obtain estimates of effect sizes, we used the ‘posterior’ function of BayesFactor to generate  
418 samples from the posterior distributions of the full models. We also tested interactions between  
419 motif and locus age, and between motif and surrounding sequence type for subsets of our data  
420 for loci of periods 2 and 3.

421

422 Our workflow for detecting homologous microsatellite loci and estimating their ages, starting  
423 from genome sequences and pairwise alignments, is given in Supplementary Fig. 7.

424

425 **Data and code availability:** The datasets generated and analysed during the current study, and  
426 the code used for analysis, are available in the Dryad repository,  
427 <https://doi.org/10.5061/dryad.7gt3rg2>. The Adélie penguin sequence read data have been  
428 deposited with links to BioProject accession number PRJNA210803 in the NCBI BioProject  
429 database (<https://www.ncbi.nlm.nih.gov/bioproject/>).

430

## 431 **References**

- 432 1 Kashi, Y. & King, D. G. Simple sequence repeats as advantageous mutators in evolution.  
433 *Trends Genet.* **22**, 253-259, doi:10.1016/j.tig.2006.03.005 (2006).
- 434 2 Gemayel, R., Vences, M. D., Legendre, M. & Verstrepen, K. J. Variable tandem repeats  
435 accelerate evolution of coding and regulatory sequences. *Annu. Rev. Genet.* **44**, 445-477,  
436 doi:10.1146/annurev-genet-072610-155046 (2010).
- 437 3 Mirkin, S. M. Expandable DNA repeats and human disease. *Nature* **447**, 932-940,  
438 doi:10.1038/nature05977 (2007).
- 439 4 Schlötterer, C. The evolution of molecular markers—just a matter of fashion? *Nat. Rev.*  
440 *Genet.* **5**, 63-69, doi:10.1038/nrg1249 (2004).
- 441 5 Ellegren, H. Microsatellite mutations in the germline: Implications for evolutionary  
442 inference. *Trends Genet.* **16**, 551-558, doi:10.1016/S0168-9525(00)02139-9 (2000).
- 443 6 Levinson, G. & Gutman, G. A. Slipped-strand mispairing: A major mechanism for DNA  
444 sequence evolution. *Mol. Biol. Evol.* **4**, 203-221 (1987).
- 445 7 Bhargava, A. & Fuentes, F. F. Mutational dynamics of microsatellites. *Mol. Biotechnol.*  
446 **44**, 250-266, doi:10.1007/s12033-009-9230-4 (2010).

- 447 8 Schlötterer, C. Evolutionary dynamics of microsatellite DNA. *Chromosoma* **109**, 365-  
448 371, doi:10.1007/s004120000089 (2000).
- 449 9 Kelkar, Y. D., Eckert, K. A., Chiaromonte, F. & Makova, K. D. A matter of life or death:  
450 How microsatellites emerge in and vanish from the human genome. *Genome Res.* **21**,  
451 2038-2048, doi:10.1101/gr.122937.111 (2011).
- 452 10 Buschiazzo, E. & Gemmell, N. J. The rise, fall and renaissance of microsatellites in  
453 eukaryotic genomes. *Bioessays* **28**, 1040-1050, doi:10.1002/bies.20470 (2006).
- 454 11 Kruglyak, S., Durrett, R. T., Schug, M. D. & Aquadro, C. F. Equilibrium distributions of  
455 microsatellite repeat length resulting from a balance between slippage events and point  
456 mutations. *Proc. Natl. Acad. Sci. USA* **95**, 10774-10778, doi:10.1073/pnas.95.18.10774  
457 (1998).
- 458 12 Calabrese, P. P., Durrett, R. T. & Aquadro, C. F. Dynamics of microsatellite divergence  
459 under stepwise mutation and proportional slippage/point mutation models. *Genetics* **159**,  
460 839-852, doi:melanogaster species complex drosophila-melanogaster saccharomyces-  
461 cerevisiae genetic distances range constraints point mutations tandem repeat allele size  
462 loci population (2001).
- 463 13 Amos, W., Kosanović, D. & Eriksson, A. Inter-allelic interactions play a major role in  
464 microsatellite evolution. *Proc. R. Soc. Lond., Ser. B: Biol. Sci.* **282**, 20152125,  
465 doi:10.1098/rspb.2015.2125 (2015).
- 466 14 Weber, J. L. & Wong, C. Mutation of human short tandem repeats. *Hum. Mol. Genet.* **2**,  
467 1123-1128 (1993).
- 468 15 Rose, O. & Falush, D. A threshold size for microsatellite expansion. *Mol. Biol. Evol.* **15**,  
469 613-615, doi:10.1093/oxfordjournals.molbev.a025964 (1998).

- 470 16 Xu, X., Peng, M., Fang, Z. & Xu, X. The direction of microsatellite mutations is  
471 dependent upon allele length. *Nat. Genet.* **24**, 396-399, doi:10.1038/74238 (2000).
- 472 17 Huang, Q.-Y. *et al.* Mutation patterns at dinucleotide microsatellite loci in humans. *Am.*  
473 *J. Hum. Genet.* **70**, 625-634, doi:10.1086/338997 (2002).
- 474 18 Sun, J. X. *et al.* A direct characterization of human mutation based on microsatellites.  
475 *Nat. Genet.* **44**, 1161-1165, doi:10.1038/ng.2398 (2012).
- 476 19 Garza, J. C., Slatkin, M. & Freimer, N. B. Microsatellite allele frequencies in humans and  
477 chimpanzees, with implications for constraints on allele size. *Mol. Biol. Evol.* **12**, 594-  
478 603 (1995).
- 479 20 Jarvis, E. D. *et al.* Whole-genome analyses resolve early branches in the tree of life of  
480 modern birds. *Science* **346**, 1320-1331, doi:10.1126/science.1253451 (2014).
- 481 21 Morey, R. D. & Rouder, J. N. Package "BayesFactor". (2015). <[https://cran.r-](https://cran.r-project.org/web/packages/BayesFactor/BayesFactor.pdf)  
482 [project.org/web/packages/BayesFactor/BayesFactor.pdf](https://cran.r-project.org/web/packages/BayesFactor/BayesFactor.pdf)>.
- 483 22 Zhang, G. *et al.* Comparative genomics reveals insights into avian genome evolution and  
484 adaptation. *Science* **346**, 1311-1320, doi:10.1126/science.1251385 (2014).
- 485 23 Hedges, S. B., Marin, J., Suleski, M., Paymer, M. & Kumar, S. Tree of life reveals clock-  
486 like speciation and diversification. *Mol. Biol. Evol.* **32**, 835-845,  
487 doi:10.1093/molbev/msv037 (2015).
- 488 24 Prum, R. O. *et al.* A comprehensive phylogeny of birds (Aves) using targeted next-  
489 generation DNA sequencing. *Nature*, 3-11, doi:10.1038/nature15697 (2015).
- 490 25 Suh, A., Smeds, L. & Ellegren, H. The dynamics of incomplete lineage sorting across the  
491 ancient adaptive radiation of neoavian birds. *PLoS Biol.* **13**, e1002224,  
492 doi:10.1371/journal.pbio.1002224 (2015).

- 493 26 Metzgar, D., Bytof, J. & Wills, C. Selection against frameshift mutations limits  
494 microsatellite expansion in coding DNA. *Genome Res.* **10**, 72-80, doi:10.1101/gr.10.1.72  
495 (2000).
- 496 27 Ohta, T. & Kimura, M. A model of mutation appropriate to estimate the number of  
497 electrophoretically detectable alleles in a finite population. *Genet. Res.* **22**, 201-204,  
498 doi:10.1017/S0016672300012994 (1973).
- 499 28 Millar, C. D. *et al.* Mutation and evolutionary rates in Adélie penguins from the  
500 Antarctic. *PLoS Genet.* **4**, e1000209, doi:10.1371/journal.pgen.1000209 (2008).
- 501 29 Baroni, C. & Orombelli, G. Abandoned penguin rookeries as Holocene paleoclimatic  
502 indicators in Antarctica. *Geology* **22**, 23-26, doi:10.1130/0091-  
503 7613(1994)022<0023:APRAHP>2.3.CO;2 (1994).
- 504 30 Baroni, C. in *Treatise on Geomorphology* Vol. 8 (eds J. F. Schroeder, R. Giardino, & J.  
505 Harbor) 430-459 (Elsevier, 2013).
- 506 31 Lorenzini, S. *et al.* Adélie penguin dietary remains reveal Holocene environmental  
507 changes in the western Ross Sea (Antarctica). *Palaeogeogr., Palaeoclimatol.,*  
508 *Palaeoecol.* **395**, 21-28, doi:10.1016/j.palaeo.2013.12.014 (2014).
- 509 32 Lambert, D. M. *et al.* Rates of evolution in ancient DNA from Adélie penguins. *Science*  
510 **295**, 2270-2273, doi:10.1126/science.1068105 (2002).
- 511 33 Ritchie, P. A., Millar, C. D., Gibb, G. C., Baroni, C. & Lambert, D. M. Ancient DNA  
512 enables timing of the Pleistocene origin and Holocene expansion of two Adélie penguin  
513 lineages in Antarctica. *Mol. Biol. Evol.* **21**, 240-248, doi:10.1093/molbev/msh012 (2004).
- 514 34 Reimer, P. J. *et al.* IntCal13 and Marine13 radiocarbon age calibration curves 0–50,000  
515 years cal BP. *Radiocarbon* **55**, 1869-1887, doi:10.2458/azu\_js\_rc.55.16947 (2013).

- 516 35 Hall, B. L., Henderson, G. M., Baroni, C. & Kellogg, T. B. Constant Holocene Southern-  
517 Ocean 14C reservoir ages and ice-shelf flow rates. *Earth Planet. Sci. Lett.* **296**, 115-123,  
518 doi:10.1016/j.epsl.2010.04.054 (2010).
- 519 36 Meyer, M. & Kircher, M. Illumina sequencing library preparation for highly multiplexed  
520 target capture and sequencing. *Cold Spring Harbor Protocols* **5**, pdb-prot5448,  
521 doi:10.1101/pdb.prot5448 (2010).
- 522 37 Quail, M. A., Swerdlow, H. & Turner, D. J. Improved protocols for the Illumina genome  
523 analyzer sequencing system. *Curr. Protoc. Hum. Genet.* **62**, 18.12.11-18.12.27,  
524 doi:10.1002/0471142905.hg1802s62 (2009).
- 525 38 Martin, M. Cutadapt removes adapter sequences from high-throughput sequencing reads.  
526 *EMBnet.journal* **17**, 10-12, doi:10.14806/ej.17.1.200 (2011).
- 527 39 Bolger, A. M., Lohse, M. & Usadel, B. Trimmomatic: A flexible trimmer for Illumina  
528 sequence data. *Bioinformatics* **30**, 2114-2120, doi:10.1093/bioinformatics/btu170 (2014).
- 529 40 Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat. Methods*  
530 **9**, 357-359, doi:10.1038/nmeth.1923 (2012).
- 531 41 Benson, G. Tandem repeats finder: A program to analyze DNA sequences. *Nucleic Acids*  
532 *Res.* **27**, 573-580, doi:10.1093/nar/27.2.573 (1999).
- 533 42 Willems, T. F., Gymrek, M., Highnam, G., Mittelman, D. & Erlich, Y. The landscape of  
534 human STR variation. *Genome Res.*, 1894-1904, doi:10.1101/gr.177774.114 (2014).
- 535 43 Lowe, C. B., Clarke, J. A., Baker, A. J., Haussler, D. & Edwards, S. V. Feather  
536 development genes and associated regulatory innovation predate the origin of Dinosauria.  
537 *Mol. Biol. Evol.* **32**, 23-28, doi:10.1093/molbev/msu309 (2015).



- 538 44 Highnam, G. *et al.* Accurate human microsatellite genotypes from high-throughput  
539 resequencing data using informed error profiles. *Nucleic Acids Res.* **41**, e32,  
540 doi:10.1093/nar/gks981 (2013).
- 541 45 R: A Language and Environment for Statistical Computing (R Foundation for Statistical  
542 Computing, Vienna, Austria, 2011).
- 543 46 Dutheil, J. Y., Gaillard, S. & Stukenbrock, E. H. MafFilter: A highly flexible and  
544 extensible multiple genome alignment files processor. *BMC Genomics* **15**, 53,  
545 doi:10.1186/1471-2164-15-53 (2014).
- 546 47 van Dongen, S. *Graph clustering by flow simulation* PhD thesis, University of Utrecht,  
547 (2000).
- 548 48 Schliep, K. P. phangorn: Phylogenetic analysis in R. *Bioinformatics (Oxford, England)*  
549 **27**, 592-593, doi:10.1093/bioinformatics/btq706 (2011).
- 550 49 Kass, R. E. & Raftery, A. E. Bayes Factors. *Journal of the American Statistical*  
551 *Association* **90**, 773-795, doi:10.1080/01621459.1995.10476572 (1995).

552

553 **Acknowledgments:** We thank John Macdonald and Peter Ritchie for assistance with collection  
554 of contemporary Adélie penguin samples. **Funding:** This research was supported by a Human  
555 Frontier Science Program grant (RGP0036/2011) and an Australian Research Council Linkage  
556 grant (2157200). Preliminary studies were funded by the Australia–India Strategic Research  
557 Fund to D.M.L. In addition, we thank Griffith University and the University of Tasmania for  
558 support and the BGI for sequencing of contemporary Adélie penguins and the Copenhagen DNA  
559 Sequencing Facility for ancient DNA sequencing. We are grateful to the Italian National  
560 Program on Antarctic Research (PNRA- 4.2/2004) and Antarctica New Zealand for support for

561 Antarctic fieldwork. **Author contributions:** B.J.M., B.R.H., C.D.M. and D.M.L. conceived and,  
562 together with M.A.C., designed the study. D.M.L., C.D.M., and B.R.H. acquired funding. C.B. &  
563 M.C.S conducted geomorphologic field survey and discovered relict penguin colonies, sampled  
564 and dated in collaboration with M.P. and C.D.M. M.P., C.D.M. and D.M.L participated in  
565 collection of contemporary Adélie penguin samples. M.P. carried out DNA library construction.  
566 R.L. and G.Z. provided genome alignments. B.J.M., M.A.C. and B.R.H. analyzed the data.  
567 B.J.M., M.A.C., M.P., B.R.H. and D.M.L. wrote and revised the manuscript, with contributions  
568 from the other authors. **Competing interests:** All authors declare that they have no competing  
569 interests.

570

571 **Fig. 1. Relative densities of microsatellite loci in intergenic, intron, exon, and regulatory**  
572 **sequences in the Adélie penguin genome.** Relative densities for loci inferred to have arisen on  
573 the branches shown on the tree. Those with periods two and three are displayed above and below  
574 the branches, respectively. Each plot shows relative rather than absolute densities, because the  
575 densities decrease rapidly with increasing locus age. Edge lengths are not drawn to scale.

576

577 **Fig. 2. Distributions of allele lengths for loci of different ages in different types of**  
578 **surrounding sequence.** Distributions of mean allele lengths (in nucleotides) of pure (A) and  
579 impure (B) microsatellite loci present in Adélie penguin and conserved across six age brackets,  
580 for loci with periods two to six in intergenic, intron, exon, and regulatory sequences. The six age  
581 brackets in each cluster correspond to loci that arose most recently on the branch leading to  
582 Adélie penguin; on the branch leading to penguins; within neoaves or on the branch leading to  
583 neoaves; on the branch leading to neognathae; on the branch leading to birds; outside sauria.  
584 (Note that no pure exonic microsatellites of period 5 are inferred to have arisen outside sauria.)  
585 Each box extends from the lower to upper quartiles of the length distribution, and the interior  
586 line indicates the median. The whiskers extend to the most extreme points within ( $1.5 \times$   
587 interquartile range) of the quartiles. Total numbers of loci are shown below each box.

588

589

590 **Table 1: Posterior mean effect of locus age on length.**

591

Period	Pure	Impure
2	0.0051 [0.0049–0.0052]	0.0100 [0.0095–0.0105]
3	0.0056 [0.0055–0.0058]	0.0140 [0.0136–0.0145]
4	0.0038 [0.0034–0.0040]	0.0125 [0.0113–0.0137]
5	0.0058 [0.0048–0.0067]	0.0113 [0.0097–0.0128]
6	0.0010 [0.0006–0.0014]	0.0120 [0.0111–0.0132]

592

593 Mean inferred rate of increase in microsatellite length, in nucleotides per million years. Numbers

594 in brackets represent 95% highest posterior density intervals.

595

596 **Table 2: Mean and standard error of allele lengths at microsatellite loci in different types of surrounding sequence**

597

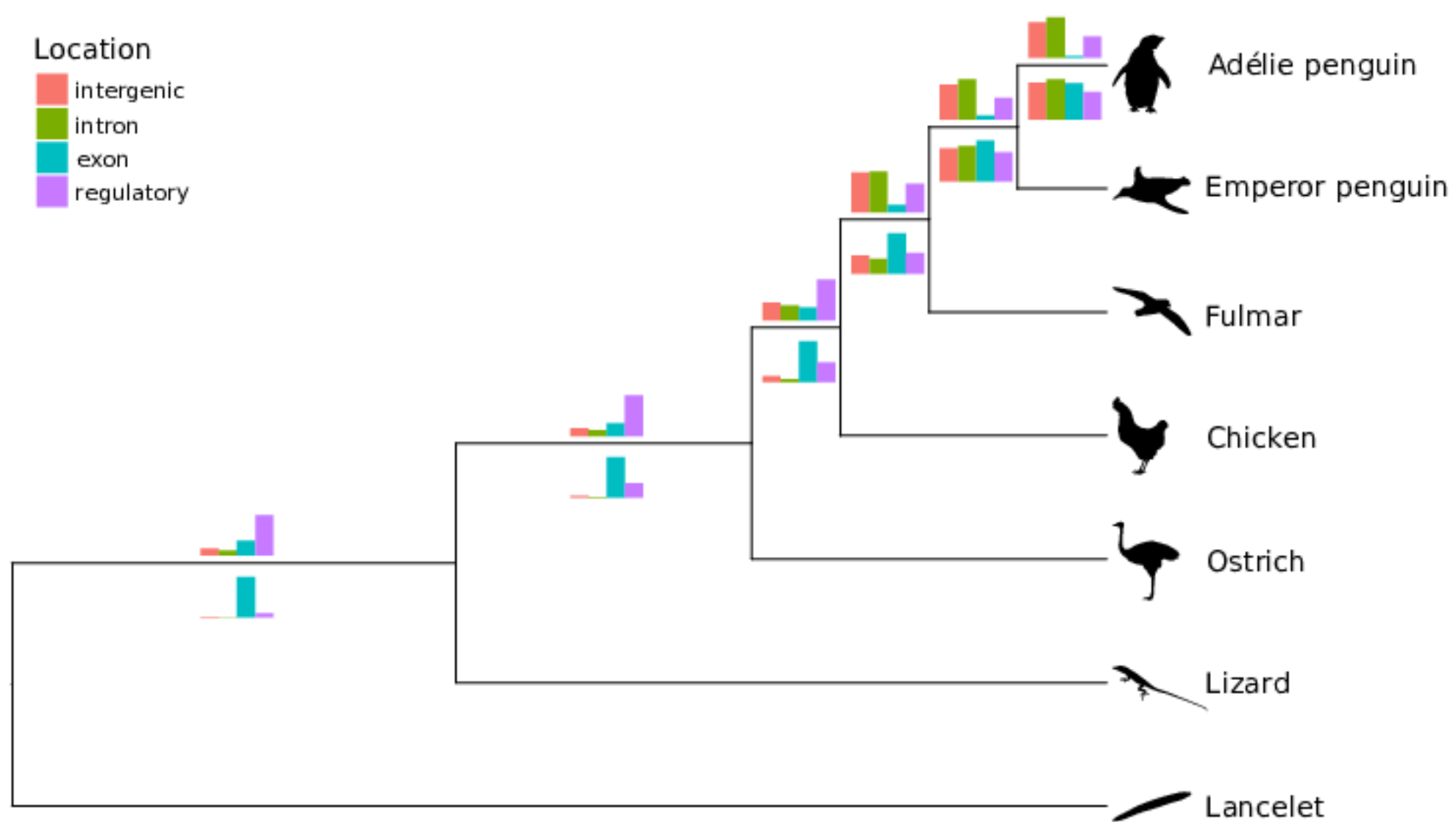
Period	Pure				Impure			
	Intergenic	Intron	Exon	Regulatory	Intergenic	Intron	Exon	Regulatory
2	13.06 [0.0027]	12.93 [0.0039]	11.79 [0.0121]	12.79 [0.0100]	20.53 [0.0080]	20.11 [0.0129]	19.66 [0.0718]	21.03 [0.0303]
3	15.99 [0.0048]	15.70 [0.0071]	16.06 [0.0154]	16.03 [0.0158]	24.02 [0.0162]	23.68 [0.0302]	26.73 [0.0521]	22.79 [0.0402]
4	16.40 [0.0043]	16.03 [0.0058]	14.83 [0.0258]	15.61 [0.0103]	24.59 [0.0111]	24.12 [0.0174]	24.45 [0.1324]	23.15 [0.0360]
5	18.93 [0.0091]	18.35 [0.0127]	17.35 [0.0468]	17.70 [0.0226]	27.79 [0.0156]	26.84 [0.0244]	24.44 [0.1561]	25.86 [0.0524]
6	18.44 [0.0086]	18.02 [0.0124]	17.78 [0.0137]	17.95 [0.0280]	27.73 [0.0187]	26.74 [0.0274]	29.71 [0.1027]	26.75 [0.0870]

598

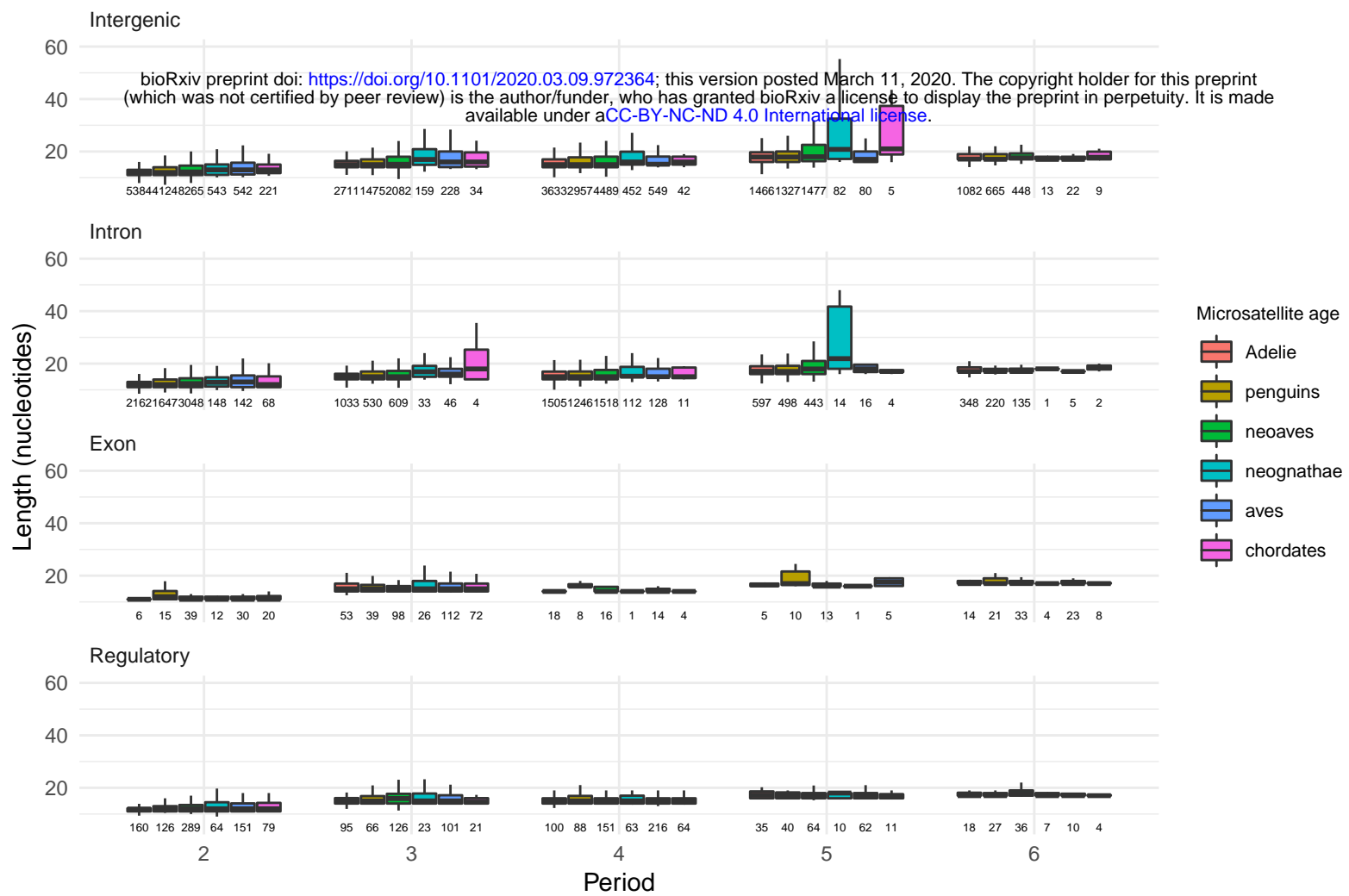
599

### Location

- intergenic
- intron
- exon
- regulatory



A



B

