

1 **Genomic islands of speciation harbor genes underlying coloration differences in a pair of**  
2 **Neotropical seedeaters**

3

4 Tram N. Nguyen<sup>1</sup>, Márcio Repenning<sup>2</sup>, Carla Suertegaray Fontana<sup>3,4</sup>, Leonardo Campagna<sup>1,5,\*</sup>

5

6 1. Department of Ecology and Evolutionary Biology, Cornell University, 215 Tower Road, Ithaca, NY  
7 14853, USA.

8 2. Universidade Federal do Rio Grande (FURG), Laboratório de Aves Aquáticas e Tartarugas Marinhas  
9 (LAATM). Av. Itália km 8, Campus Carreiros, 96203-900, Rio Grande, Rio Grande do Sul, Brazil.

10 3. Universidade Federal do Rio Grande do Sul (UFRGS). Laboratório de Ecologia de Comunidades e  
11 Populações, Instituto de Biociências, Avenida Bento Gonçalves, 9500, 91501-970, Porto Alegre, Rio  
12 Grande do Sul, Brazil.

13 4. Programa de Pós-graduação em Biodiversidade Animal, Universidade Federal de Santa Maria.

14 5. Fuller Evolutionary Biology Program, Cornell Lab of Ornithology, 159 Sapsucker Woods Road, Ithaca,  
15 NY, 14850, USA.

16 \*Correspondence: [lc736@cornell.edu](mailto:lc736@cornell.edu)

17

18 Short title: Islands of speciation in two Neotropical seedeaters

19 **Abstract**

20

21 Incomplete speciation can be leveraged to associate phenotypes with genotypes, thus providing insights  
22 into the traits relevant to the reproductive isolation of diverging taxa. We investigate the genetic  
23 underpinnings of the phenotypic differences between *Sporophila plumbea* and *S. beltoni*. *S. beltoni* has  
24 only recently been described based, most notably, on differences in bill coloration (yellow vs. black in *S.*  
25 *plumbea*). Both species are indistinguishable through mtDNA or reduced-representation genomic data,  
26 and even whole-genome sequencing revealed low genetic differentiation. Demographic reconstructions  
27 attribute this genetic homogeneity to gene flow, despite divergence in the order of millions of generations.  
28 We found a narrow hybrid zone in southern Brazil where genetically, yet not phenotypically, admixed  
29 individuals appear to be prevalent. Despite the overall low genetic differentiation, we identified three  
30 narrow peaks along the genome with highly differentiated SNPs. These regions harbor six genes, one of  
31 which is involved in pigmentation (*EDN3*) and is a candidate for controlling bill color. Within the outlier  
32 peaks we found signatures of resistance to gene flow, as expected for islands of speciation. Our study  
33 shows how genes related to coloration traits are likely involved in generating prezygotic isolation and  
34 establishing species boundaries early in speciation.

35

36 **Key words**

37

38 Ancestral Recombination Graph, Genome scan, Hybridization, Islands of speciation, *Sporophila*.

39

40 **Introduction**

41

42 Increasing knowledge about the genetic, ecological, and biogeographic contexts in which new  
43 species are formed is central to understanding the processes which lead to speciation. General patterns  
44 can be found by accumulating examples of how speciation has taken place throughout the tree of life [1-  
45 4]. The genetic architecture of reproductive barriers can be studied by leveraging non-model systems in  
46 the early stages of speciation and/or by focusing on hybridizing species where these barriers break down.

47 For recently diverged taxa, genetic differentiation tends to be localized around areas relevant to  
48 speciation, but with time, genomes will continue to diverge as other differences start to follow [5]. In such  
49 cases, natural hybridization or experimental crossing can help break up genetic linkage through  
50 recombination, and phenotypic traits can still be associated statistically to genetic changes [6]. Whereas  
51 studies following these types of designs have revealed how genomic landscapes are clearly  
52 heterogeneously differentiated [5], the processes driving such differentiation remain debated. Areas of  
53 high differentiation (or outlier regions) can be shaped by processes such as demographic range  
54 expansions [7], background selection combined with low levels of recombination [8], selective sweeps [9]  
55 or resistance to gene flow [10,11]. These processes can leave different signatures in the outlier genomic  
56 regions which researchers can examine with various traditional summary statistics, such as  $F_{ST}$ ,  $D_{XY}$ ,  
57 Tajima's D, and  $\Pi$ , to uncover why genomes are heterogeneously differentiated [11]. However, clarity on  
58 the processes which have shaped many of these genomic regions is still generally lacking [8]. With the  
59 increasing feasibility of whole-genome resequencing, the toolkit for inferring evolutionary processes from  
60 genomic patterns continues to grow, offering researchers greater power to pinpoint regions of the  
61 genome under selection that mediate phenotypes of interest [12].

62 In birds, genome scans among recently diverged taxa have repeatedly uncovered coloration  
63 genes in areas of high differentiation in species that otherwise show shallow overall levels of genomic  
64 differentiation [10,13-17]. Collectively, these findings point to the importance of coloration differences in  
65 the early stages of speciation, and their likely effectiveness as mechanisms of prezygotic isolation [18-  
66 20]. Notably, *de novo* mutations in coloration genes, or the reassembly of different variants from standing  
67 variation across permeable species limits [21], may promote speciation. However, many taxa in the early  
68 stages of speciation may be ephemeral [22,23], with their persistence dependent on the accumulation of  
69 additional mutations over time, perhaps facilitated by the biogeographic context in which they were  
70 formed. Nevertheless, recently diverged species, whether they persist or not, will be informative of the  
71 types and location of genetic changes which can lead to initiating speciation.

72 The Neotropical genus *Sporophila* contains over 40 species and has the highest speciation rates  
73 of the species-rich Tanager family, to which it belongs [24,25]. This high speciation rate is in part the  
74 consequence of several small groups of rapidly radiating (recent and closely related) taxa in the early

75 stages of speciation, making the genus a valuable study system to understand the origin of reproductive  
76 isolation and the process of speciation. Many of these groups have been studied from a genetic  
77 perspective, and consistently show differentiation in male sexual traits (including coloration patterns)  
78 despite little genome-wide genetic differentiation [14,26,27]. The largest of such groups is known as the  
79 capuchino seedeaters, which contains 12 species differing most notably in the plumage patterns of adult  
80 males and in their vocalizations (a primarily cultural trait) [28,29]. Despite their phenotypic diversity,  
81 capuchino seedeaters show extremely low genetic differentiation, except for a few areas of the genome  
82 which are enriched for pigmentation genes. These genes likely mediate the different coloration  
83 phenotypes observed across males from different species [14,20,30], which in turn, impact species  
84 recognition [20]. Similarly, the Variable Seedeater superspecies complex (*S. corvina*, *S. intermedia*, *S.*  
85 *murallae*, and *S. americana*) and Morelet's Seedeater (*S. moreletti*) represent additional cases of  
86 plumage diversity, yet shallow genetic differentiation within *Sporophila* [26,27,31,32]. Overall, these  
87 examples suggest that localized genetic differences may mediate divergence in coloration traits,  
88 contributing to reproductive isolation and speciation [20].

89 Another *Sporophila* species that is not part of the previously mentioned groups, the Plumbeous  
90 Seedeater (*S. plumbea*), was traditionally thought to contain black and yellow-billed individuals. However,  
91 this presumed intraspecific variation was not a common polymorphism and did not correlate with  
92 seasonality or age (Figure 1A and 1B, [33]). In 2013, the yellow-billed individuals were shown to  
93 constitute a new species, the Tropeiro seedeater (*S. beltoni*), that bred in the grassland highlands of  
94 southern Brazil [34,35]. *S. beltoni* is differentiated from *S. plumbea* primarily through bill coloration, but  
95 also through male size and bill shape, vocalizations, migration phenology and subtle aspects of male  
96 plumage coloration. Furthermore, both species also differ in their breeding ranges, habitat, and phenology  
97 [34,35]. *S. beltoni* is sexually dimorphic and shares a narrow contact zone (~50-100 km) in the north of its  
98 distribution with the southernmost breeding population of *S. plumbea* (Figure 1A). In this contact zone,  
99 individuals are mostly segregated by habitat and elevation, though they can breed in close proximity  
100 where their respective habitats are available [34]. Although hybrids have not been definitively identified in  
101 this region, a few individuals with irregularly colored and streaky bills were observed, raising the  
102 possibility that hybridization occurs in this contact zone. However, a preliminary study using reduced-

103 representation genomic data failed to distinguish these two species [35], suggesting they were of very  
104 recent origin and/or continued to experience high levels of gene flow. Together, the phenotypic difference  
105 and low genetic differentiation between this species pair provides an opportunity to further investigate the  
106 process of speciation.

107 Here we assessed the degree of genetic differentiation between *S. beltoni* and *S. plumbea* using  
108 whole-genome resequencing, and conducted the first study to search for genomic regions associated with  
109 their phenotypic differences. Using the genomes of individuals from each species across both their  
110 allopatric ranges (i.e., outside of the contact zone) and within the contact zone, we conducted a  
111 demographic reconstruction to understand the demographic patterns in these different geographic  
112 regions. We then performed a genome scan which highlighted three narrow regions containing genes that  
113 likely mediate differences in beak coloration. We subsequently used genealogy-based statistics derived  
114 from topologies extracted from the ancestral recombination graph (ARG) to infer the processes that have  
115 shaped genomic outlier regions, and found signatures of reduced gene flow, as is expected for islands of  
116 speciation. Our study contributes to the emerging body of literature showing differentiation in genes  
117 related to bird coloration, which may mediate prezygotic isolation, before speciation is complete.

118

## 119 **Results**

120

### 121 **Low genomic differentiation between *S. beltoni* and *S. plumbea***

122

123 While *S. beltoni* and *S. plumbea* individuals sampled outside of the contact zone (n=11 for each  
124 species) clustered separately in a PCA derived from ~15 million genome-wide SNPs, those sampled  
125 within the contact zone resembled each other, overlapping in the PCA (n=7 for each species; Figure 1C,  
126 Figure S1). Moreover, individuals could not be assigned to species using mtDNA (cytochrome c oxidase  
127 subunit 1 sequences), irrespective of sampling location (Figure 1D). Although the overall level of genome-  
128 wide differentiation was consistently low across groups, it was higher when comparing between allopatric  
129 individuals versus between sympatric individuals ( $F_{ST} = 0.0031$  for allopatric, ~0 for sympatric, and on  
130 average 0.0015 for all samples combined). Using these genomic data, we then reconstructed the history

131 of divergence between the two species. We found similarly large current effective population sizes for  
132 both taxa (in the order of millions of individuals), which were much larger than the size of the inferred  
133 ancestral population (~5 fold; Figure 2A). The divergence time between *S. beltoni* and *S. plumbea* was  
134 estimated to be between ~3.2 and 3.6 million generations (Figure 2B), and the genetic similarity between  
135 the species could be partly explained by high levels of gene flow since their split (between 2 and 4  
136 migrants per generation; Figure 2C). There was some variation in the parameter values obtained from  
137 models using birds sampled from different parts of the range (allopatric, sympatric, or a combination of  
138 both; Figure 2). However, the most notable difference was the direction of gene flow (Figure 2C), which  
139 we interpret with caution as it can be particularly challenging to infer under certain scenarios [36]. While  
140 migration was inferred from *S. beltoni* into *S. plumbea* within the contact zone, the opposite was true for  
141 allopatric individuals (or when all samples were combined). It is possible that the patterns of hybridization  
142 may be different within and outside of the contact zone, with the former being representative of more  
143 recent gene flow (see Discussion).

144

#### 145 **Differentiated SNPs are concentrated in three outlier peaks**

146

147 When comparing all samples, a small number of SNPs showed differentiation above  $F_{ST}=0.20$   
148 (1.1%) and only 40 SNPs (out of ~15 million) had an  $F_{ST}$  value above 0.7 (Figure 3A). Of the 40 variants  
149 showing the highest  $F_{ST}$  values, 39 were clustered in three divergence peaks: on contigs 404 (24 elevated  
150 SNPs in ~100kb on the Z chromosome), 33 (9 elevated SNPs in ~66 kb on chromosome 20), and 382 (6  
151 elevated SNPs in ~2 kb on chromosome 11) (Figure 3A). Taken together, the variants from these peaks  
152 could be combined to cluster individuals by species in a PCA (Figure S1), and the same was true when  
153 distinguishing allopatric individuals using the variants from each peak separately (Figure 3B). However,  
154 with the exception of contig 33, sympatric individuals were genetically similar to each other and could not  
155 always be assigned to species based on their genotypes in these genomic regions (Figure 3B, Figure  
156 S2). The 39 variants showing the highest  $F_{ST}$  values in these peaks were in high linkage-disequilibrium,  
157 both within but also between peaks, suggesting they are co-inherited despite being on different  
158 chromosomes (Figure S3). Moreover, each species showed a common haplotype on each peak, with less

159 intraspecific variation compared to interspecific differences (Figures S4-S6). Only a few sympatric  
160 individuals were heterozygotes or sometimes homozygotes for haplotypes from the other species. We did  
161 not find individuals that showed admixture at all three peaks simultaneously (Figures S4-S6).

162

### 163 ***EDN3* and *PRLR* are candidate genes implicated in the phenotypic differences between species**

164

165 The divergence peaks we identified involved a total of six genes (Figure 3B, Table 1). Two of  
166 these genes, *EDN3* which encodes the protein Endothelin-3, and *PRLR* which encodes the Prolactin  
167 receptor, may mediate the difference in beak coloration between *S. beltoni* and *S. plumbea* (Table 1; see  
168 Discussion). The remaining genes could be related to differences between these taxa that are not  
169 immediately obvious to us (Table 1). The variants showing the highest  $F_{ST}$  values in these peak regions  
170 were mostly non-coding (74% fell in intergenic regions and 13% in introns; Figure 3B), however two were  
171 responsible for non-synonymous mutations in the *HSDL1* gene on contig 382. Three additional mutations  
172 fell within *EDN3* but our annotation of that region was not of sufficient quality to distinguish if they fell  
173 within introns or in exons.

174

### 175 **Outlier genomic regions are speciation islands**

176

177 We explored the divergence peaks to search for signatures that could reveal the processes  
178 involved in shaping these regions. We first looked at Tajima's D and nucleotide diversity, which tend to be  
179 reduced after selective sweeps. Although Tajima's D was on average negative and nucleotide diversity  
180 was reduced in the three peaks, more extreme values could be found outside of the peaks within the  
181 same contigs (Figure S7). Therefore, based on these statistics alone, it is not conclusive whether  
182 selective sweeps have occurred within these outlier peaks. We next looked at H12 and H2/H1 which can  
183 be used to distinguish hard and soft selective sweeps, but again did not find conclusive patterns (Figure  
184 S7). We also compared  $F_{ST}$  and  $D_{XY}$  values and observed slightly elevated  $D_{XY}$  in some windows within  
185 the  $F_{ST}$  divergence peaks (compared to both windows outside of the  $F_{ST}$  peaks or in contigs without  $F_{ST}$   
186 peaks), suggesting these regions may have accumulated some absolute sequence differences between

187 the species (Figure S8). This prompted us to search for signatures of selection in more depth using three  
188 statistics derived from topologies extracted from the ancestral recombination graph (ARG) (Figure 4A).  
189 The species enrichment score measures the probability of observing clades containing individuals from a  
190 certain species [9]. The relative TMRCA half-life,  $RTH'$  [37], is a normalized version of the time to most  
191 recent common ancestry (TMRCA). Finally, the inter-species cross-coalescence measures the average  
192 age of the most recent coalescence events between tips in the tree which belong to the different species.  
193 We followed the reasoning outlined by Hejase et al. [9], where areas under selection are expected to  
194 show species enrichment, but these could be the product of either species-specific and recent selective  
195 sweeps, or they could be older islands of speciation that resist gene flow. These two extreme scenarios  
196 can be distinguished by using  $RTH'$ , which we expect to be low when selective sweeps produce shallow  
197 clades, and the inter-species cross-coalescence, which we expect to show larger values when gene flow  
198 is selected against (i.e., delayed cross-coalescence [9]). Following this logic, we obtained values for the  
199 three statistics in the divergence peaks and for a control set of contigs which represented ~10% of the  
200 genome, and we used the latter to establish species-specific thresholds of significance for each statistic.  
201 We also conducted this analysis for all samples together, and for sympatric and allopatric samples  
202 separately. The analysis shows genealogies within the peak regions which are significantly enriched for  
203 each species (Figure 4B, Table 2, Table S2, Figures S9-S12). These clades show low  $RTH'$  values,  
204 sometimes for one species but in some cases for both, and for the peaks on contigs 404 and 382 we also  
205 observed delayed cross-coalescence (Table 2, Table S2). These patterns were accentuated when  
206 conducting the analysis with samples from outside of the contact zone and generally weaker or lost when  
207 using sympatric birds from within the contact zone. Taken together, these regions show evidence of  
208 having undergone selection and experiencing reduced migration compared to other areas in the genome,  
209 especially when analyzing allopatric samples which may be more representative of historical processes  
210 (see Discussion).

211

## 212 **Discussion**

213

214 Whole-genome sequencing of *Sporophila beltoni* and *Sporophila plumbea* revealed little genetic  
215 differentiation between this species pair despite their marked phenotypic differences, similar to what has  
216 been found in other *Sporophila* species [14,26,32]. The highest levels of differentiation were concentrated  
217 in three narrow peaks, containing only six genes, which are candidates for mediating the phenotypic  
218 differences between our focal taxa. However, unlike the other *Sporophila* species which have shown this  
219 pattern of genetic homogeneity [9,14], we inferred *S. beltoni* and *S. plumbea* to be comparatively old, in  
220 the order of 3 million generations (~1 Myr). For comparison, the 10 species of southern capuchino  
221 seedeaters (those found south of the Amazon River [38]) began to diverge within a narrow window over  
222 the last million years [39]. Our demographic modelling suggests high levels of gene flow, and not recent  
223 origin, as a possible explanation for the genetic homogeneity between *S. beltoni* and *S. plumbea*. This  
224 scenario of relatively deep divergence, high overall gene flow, and a small number of highly differentiated  
225 regions suggests that these divergence peaks could be islands of speciation [40], i.e., regions of the  
226 genome that resist homogenization through gene flow. In contrast, in capuchino seedeaters most  
227 divergence peaks, which are also enriched for coloration genes, were found to be shaped through recent  
228 species-specific selective sweeps [9]. Therefore, divergence in regions of the genome containing  
229 pigmentation genes has likely evolved through different processes in *S. beltoni/plumbea* and capuchino  
230 seedeaters [9,14].

231 The two outlier peaks with the highest levels of differentiation and encompassing most of the  
232 highly differentiated SNPs (on contig 404 and contig 33), contained only three genes, two of which have  
233 functions that may mediate the phenotypic differences between our focal taxa. *EDN3* functions in  
234 melanoblast proliferation and differentiation, and mutations (including duplications) in either the *EDN3*  
235 gene itself or its receptors (*EDNRs*) can lead to hypo or hypermelanization in quails and chickens [41],  
236 and to several piebalding phenotypes in domestic pigeons [42]. These differences in pigmentation have  
237 been shown to affect not only feathers, but also other body parts, like the chicken comb [43] or even  
238 internal organs [41]. The Prolactin receptor (*PRLR*) is implicated in keratinization and feather formation  
239 [44] and Prolactin itself is known to stimulate molt [45] and is involved in coloration phenotypes in other  
240 taxa [46]. It is therefore likely that these genes are involved in mediating the coloration differences  
241 between *S. beltoni* and *S. plumbea*, the most notable of which is in the beak. These differences in beak

242 color may contribute to promoting prezygotic reproductive isolation, as has been found between  
243 subspecies of Long-tailed Finches (*Poephila acuticauda*) [47,48]. Within the genus *Sporophila*, species  
244 have either gray, black, yellow or orange bills [24]. The clade containing *S. plumbea* and *S. beltoni* has a  
245 total of six species (also including *S. albogularis*, *S. falcirostris*, *S. collaris* and *S. schistacea*) [49], and  
246 four of these six species have orange bills, suggesting that the black melanized bill is possibly a derived  
247 trait in *S. plumbea*. More work is needed to understand if *EDN3* and *PRLR* mediate beak coloration  
248 across multiple *Sporophila* species, as is the case with the *BCO2* gene and carotenoid-based coloration  
249 in *Setophaga* warblers [50].

250         It was previously unclear whether *S. beltoni* and *S. plumbea* hybridized in their contact zone in  
251 Southern Brazil (Figure 1A, [34]), mainly due to the lack of individuals with a clear mixed phenotype in the  
252 wild. However, our study shows extensive admixture in this region, likely due to ongoing hybridization.  
253 Therefore, this narrow area of contact likely constitutes a hybrid zone. This hybrid zone coincides with a  
254 region where the breeding habitat of both species is degraded due to human activities, an aspect which  
255 may have contributed to increasing the overlap between the breeding ranges of both species and resulted  
256 in high levels of recent gene flow [51]. It remains to be determined whether this habitat alteration, and the  
257 possibly associated increase in gene flow, remains restricted to this narrow contact zone or if it may have  
258 consequences for the completion of speciation. Within the divergence peaks, sympatric individuals were  
259 more similar to each other than allopatric individuals, and in some analyses (e.g., PCA or admixture  
260 plots), some were indistinguishable (Figure 3B, Figure S1-S2). The exception was the peak on contig 33  
261 which contained the *EDN3* gene that may be involved in controlling beak coloration, allowing the correct  
262 species identification in the contact zone (Figure 3B). Heterogeneous genomic landscapes tend to show a  
263 larger number of peaks outside of contact zones, where gene flow is reduced [52-54]. This increased  
264 level of divergence could be due to selection on traits in allopatry, unrelated to reproductive isolation and  
265 where gene flow does not occur [52]. If this is true for our focal taxa, the lack of admixture in *EDN3* in  
266 both allopatry and sympatry constitutes further evidence that this gene may mediate phenotypic  
267 differences relevant to reproductive isolation. However, we note that the genome-wide level of  
268 differentiation between *S. beltoni* and *S. plumbea* is very low in allopatry ( $F_{ST} = 0.0031$ ), and that we infer  
269 high levels of genome-wide gene flow when analyzing sympatric or allopatric samples (Figure 2C). It is

270 therefore possible that the three peaks of differentiation that we have identified have resisted gene-flow  
271 during speciation. This prompted us to search for signatures of different evolutionary processes which  
272 could have shaped our candidate genomic regions.

273 We did not find clear patterns using traditional summary statistics (e.g., Tajima's D, nucleotide  
274 diversity,  $D_{XY}$ ) within the peak regions that would allow us to draw robust conclusions about the processes  
275 that have shaped them. We therefore leveraged the ARG and genealogy-based statistics that would  
276 potentially give us greater resolution and allow us to distinguish between two different extreme scenarios  
277 which could elevate  $F_{ST}$ : recent and species-specific selective sweeps and resistance to gene flow [9,17].  
278 Regions which have undergone recent species-specific selective sweeps are expected to show clades  
279 containing most individuals from the relevant species that are significantly shallower (i.e., showing  
280 reduced diversity) than other areas of the genome. Islands of speciation, on the other hand, should show  
281 signatures of reduced gene flow (i.e., delayed cross-coalescence) when compared to other regions of the  
282 genome. We note that these two extreme models are not mutually exclusive, and a locus could undergo a  
283 selective-sweep after which gene flow between species would be reduced [9]. This latter scenario is  
284 supported by our data in the peak regions, where we tend to see comparatively shallow clades for one or  
285 both species, as expected for loci that have undergone selective sweeps. Moreover, we also see  
286 evidence for delayed cross-coalescence (primarily outside of the hybrid zone, but also more generally  
287 when the p-value thresholds are not as stringent; Table 1 vs. Table S2). Taken together, our findings  
288 support the islands of speciation model [11,40], where a few loci that may be relevant to generating  
289 coloration differences and prezygotic isolation, have undergone selection and resisted gene flow.

290 Our study represents the first investigation into the phenotypic differences between *S. beltoni* and  
291 *S. plumbea* and shows how coloration genes can diverge early in speciation and likely mediate prezygotic  
292 isolation. We also exemplify how the ARG can be leveraged to distinguish between competing processes  
293 which shape genome evolution, finding evidence of reduced gene flow in speciation islands compared to  
294 the rest of the genome. Taken together, our genome scan, demographic reconstructions, and ARG-based  
295 analyses reveal the complicated interplay between selection and gene flow in shaping species  
296 boundaries.

297

298 **Materials and methods**

299

300 ***Sampling and dataset***

301

302 We sampled a total of 36 individuals for whole genome resequencing, 18 *S. beltoni* and 18 *S. plumbea*.  
303 Samples originated from allopatric portions of the ranges of both species (n=11 for each species), as well  
304 as from the contact zone (n=7 for each species) (Figure 1; Table S1). For this study we only used adult  
305 males so that identification would be unambiguous, and so that our depth of coverage for autosomes and  
306 sex chromosomes would be equivalent. Birds were captured during breeding seasons (austral  
307 spring/summer) between 2008 and 2014 using mist nest, and subsequently banded and bled before  
308 being released. To conduct fieldwork in Brazil and collect samples, we secured licenses from the National  
309 Center for Research on the Conservation of Wild Birds (CEMAVE, licenses 3711 and 3778) and the  
310 Instituto Chico Mendes of Conservation of Brazilian Biodiversity (ICMBio), through the Brazilian System  
311 and Information on Biodiversity (SISBIO, license numbers: 13310, 36881, 35434, and 52714).

312

313 ***Sequencing and variant discovery***

314

315 We used the DNEasy blood and tissue kit (Qiagen, CA, USA) to obtain DNA from blood and  
316 prepared libraries following the protocol recommended for the TruSeq Nano DNA library preparation kit  
317 protocol (550 bp inset size). Libraries were pooled into two groups of 18 samples, using concentrations of  
318 adapter-ligated DNA (determined through digital polymerase chain reaction), and each group was  
319 sequenced on its own Illumina NextSeq 500 lane at the Cornell Institute for Biotechnology core facility.  
320 This produced a total of ~1.6 billion 151 bp paired-end reads, and consequently based on the number of  
321 raw reads obtained for each individual, we expected an average depth of coverage of 5.6 +/- 1.1. We  
322 assessed sequence quality and performed filtering and adapter removal as described previously [17]. We  
323 first trimmed low-quality bases from individual reads and collapsed overlapping paired reads using  
324 AdapterRemoval version 2.1.1 (options: --trimns --trimqualities --minquality 10) [55]. We then aligned the  
325 filtered sequences from each individual to the *Sporophila hypoxantha* reference genome [14] using

326 Bowtie2 version 2.4.3 [56] with the very sensitive local option. Despite using a reference genome from a  
327 different, yet closely related species, the average alignment rate was 96% (Table S1). The average depth  
328 of coverage across all samples and sites, calculated using Qualimap v2.2.1 [57], was 5.4 +/- 1.1x (Table  
329 S1). We followed the genotyping pipeline described in detail in [17] and filtered the resulting SNP dataset  
330 with VCFtools version 0.1.16 [58]. Briefly, we first marked PCR duplicates and realigned around indels  
331 using “MarkDuplicates” and “IndelRealigner”, respectively. We subsequently used the “Haplotypecaller”  
332 and “GenotypeGVCFs” modules from GATK version 3.8.1 [59] to identify variants and conduct joint  
333 genotyping across samples, resulting in a single variant file for the entire dataset. We then selected SNPs  
334 with the “SelectVariants” module of GATK (option: --selectType SNP) and filtered out those that did not  
335 satisfy the following filters: QD < 2, FS > 60.0, MQ < 30.0, ReadPosRankSum < -8.0. Finally, we  
336 performed additional post-hoc filtering using VCFtools to retain 15,663,387 variants present in at least  
337 85% of individuals, with mean depth of coverage between 2 and 50 and a minor allele count of at least  
338 six. We explored the sensitivity of our dataset to a range of filtering parameters in different combinations,  
339 including up to a minimum depth of coverage of 5, a minor allele count of at least 10 and less than 5%  
340 missing data. These different filters produced congruent patterns in our PCAs and Manhattan plots, yet  
341 reduced the number of total SNPs retained. We therefore decided to retain a larger number of SNPs and,  
342 unless otherwise stated, use the dataset described above for downstream analyses.

343

#### 344 ***Analysis of mitochondrial DNA***

345

346 Because we did not recover high-quality mitochondrial genomes for every individual from the  
347 whole genome data, possibly due to miss-assemblies related to nuclear sequences of mitochondrial  
348 origin, we decided to sequence the cytochrome c oxidase I (*COI*) gene, which is commonly used for  
349 species identification [60]. We used primers BirdF1 and COIBirdR2 and followed the procedures  
350 described by Kerr et al. [61]. We Sanger sequenced PCR-amplified DNA using both the forward and  
351 reverse primers at the Cornell Institute for Biotechnology core facility. Forward and reverse sequences  
352 were combined into a consensus sequence for each individual and aligned using Geneious version 10.2.6  
353 [62]. In total we obtained 27 *COI* sequences, 8 allopatric and 6 sympatric *S. beltoni* and 7 allopatric and 6

354 sympatric *S. plumbea* (Table S1). We used these sequences to construct a *COI* haplotype network using  
355 the R version 4.0.2 [63] package *pegas* v1.1 [64].

356

### 357 ***Genetic differentiation, admixture and summary statistics***

358

359 We first assessed genome-wide patterns of differentiation among species and sampling locations  
360 (i.e., allopatric or sympatric) by conducting a Principal Component Analysis (PCA) in the R package  
361 *SNPRelate* version 3.3 [65]. Two samples originated from individuals that were later found to be related  
362 (TNN26 and TNN27, relatedness of 0.538) and separated from the remaining samples in a PCA (Figure  
363 S1A). We therefore excluded these samples from the PCA in Figure 1C, which produced the same  
364 pattern as when the samples were retained, and PC2 vs. PC3 were plotted (Figure S1B). We used the  
365 same programs to produce PCAs from the SNPs found within the divergence peaks, after subsetting the  
366 dataset by genomic location in *VCFtools*. We also calculated  $F_{ST}$  values for individual SNPs using  
367 *VCFtools* and displayed these in a histogram produced in R or as a Manhattan plot using the R package  
368 *qqman* version 0.1.8 [66]. For the latter, we simplified the plot by only visualizing SNPs with  $F_{ST} > 0.2$ . We  
369 also used *VCFtools* to subset the dataset by species and contig to calculate statistics base on the  
370 frequency of the most common haplotypes (H1, H2, H12 and H2/H1) in *SelectionHapStats* [67]. We used  
371 a window size (non-overlapping) of 25 SNPs and merged haplotypes differing in one position (--  
372 distanceThreshold 1). We also used *VCFtools* to calculate  $R^2$  values and estimate linkage disequilibrium  
373 (LD) among all the different sites showing the highest differentiation ( $F_{ST} > 0.7$ ) in the divergence peaks.  
374 For these same variants (24 on contig 404, 9 on contig 33 and 6 on contig 382), within each divergence  
375 peak, we explored local haplotypes by first phasing and imputing missing data using default parameters  
376 in *BEAGLE* version 3.3.2 [68]. We plotted the two haplotypes for each peak region from every individual  
377 using the function *phylo.heatmap* from the R package *phytools* [69], and clustered them according to  
378 similarity by producing a distance matrix in the R package *vegan* [70]. Finally, we explored ancestry along  
379 the contigs of interest by running *Admixture* version 1.3.0 [71] in 50 kb (non-overlapping) sliding windows  
380 using a K value of two. As the peak on contig 382 was narrow (~2 kb), we used a sliding window of 10 kb  
381 in this case, as a compromise between not averaging out across too many SNPs while maintaining

382 sufficient resolution. We plotted these values separately for allopatric and sympatric samples using a  
383 smoothing line in ggplot2 [72]. To calculate Tajima's D and nucleotide diversity ( $\pi$ ) in the contigs of  
384 interest, we used a vcf file that followed the same filtering steps described above except for that it was not  
385 filtered by minor allele frequency (the same file used in the demographic reconstruction), as this would  
386 bias both summary statistics that are sensitive to low frequency alleles. We calculated these statistics in  
387 non-overlapping, 5 kb windows. Finally, we calculated  $D_{XY}$  using the program pixy version 1.2.7 [73]. For  
388 this specific analysis we included invariant sites by running GATK's "GenotypeGVCFs" module with the  
389 option "--includeNonVariantSite". Following filtering recommendations [73], we then separated our data  
390 into variant and invariant sites and filtered them independently, as applying population genetic based  
391 filters such as minor-allele frequency (MAF) filtering would result in the loss of our invariant sites. We  
392 omitted the MAF filter to obtain invariant sites, but all other filtering steps were implemented as described  
393 previously (e.g., missing data or coverage filters).

394

### 395 ***Demographic reconstruction***

396

397 We estimated current and ancestral effective population sizes for both species, the splitting time  
398 between the taxa and migration rates (total of six demographic parameters) using an isolation with  
399 migration model as implemented in G-PhoCS version 1.3 [74]. Because this analysis is computationally  
400 intensive, we subsampled individuals from the dataset. We conducted three different analyses using  
401 individuals from different parts of the ranges of both species: allopatric (eight from each species),  
402 sympatric (seven from each species) or a combination of both types of samples (4 allopatric and 4  
403 sympatric from each species). To avoid biasing our analyses by excluding low frequency alleles, we used  
404 a vcf file that was not filtered for minor allele frequency (~61M SNPs), but otherwise had the same filters  
405 applied as described previously. We subsequently discarded contigs that were smaller than 1 Mb, Z-  
406 linked, or had outlier SNPs (contig 404, 33, 382 and 910). We used the "FastaAlternateReferenceMaker"  
407 module in GATK to generate sequence files for each individual for the remaining contigs, and sampled  
408 and aligned 2,500 loci, each 1 kb in length and at least 100 kb apart. We ran G-PhoCS for 1 million  
409 generations (for the sympatric and allopatric analyses) or for 300,000 iterations (for the combined

410 analysis) and discarded the initial 50,000 as burn-in. We used the coda package in R [75] to check for  
411 convergence and subsequently converted median and 95% Bayesian credible intervals to generations or  
412 individuals (from mutation scale) by using an approximate mutation rate estimate of  $10^{-9}$  per bp per  
413 generation [76]. We have focused our interpretations primarily on relative comparisons between the  
414 species which are independent of our assumption of mutation rate, as is the number of migrants per  
415 generation [74].

416

### 417 ***Identification of genes in outlier genomic regions***

418

419 We defined  $F_{ST}$  outlier regions as the genomic coordinates encompassed by the SNPs showing  
420  $F_{ST} > 0.7$ . By using this approach, we aim to focus on the regions showing the highest levels of  
421 differentiation between species, however this does not imply that other regions, showing more subtle  
422 patterns, are not relevant to shaping phenotypic differences. We had previously mapped every contig in  
423 the reference genome to its corresponding chromosome in the Zebra Finch [14]. Most of these outlier  
424 SNPs were clustered together and located in three peaks (contig 404, contig 33 and contig 382), yet there  
425 was a single additional elevated SNP on scaffold 910 (near the gene LRRRC15), corresponding to  
426 chromosome 9. We did not consider this a divergence peak as it included a single position, and don't  
427 discuss it further. We searched for the genes within each divergence peak ( $\pm 5$  kb) by inspecting our  
428 reference genome annotation using Geneious version 10.2.6 [62]. We used the NCBI database  
429 (<http://www.ncbi.nlm.nih.gov/>) to find information on the function of these different genes in other  
430 organisms. Finally, the outlier genomic regions were similar to the rest of the genome in average missing  
431 data (0.05 vs. 0.05), variant quality (QUAL; 1683.19 vs 1678.24) and depth of coverage (4.97 vs 5.00).

432

### 433 ***Statistics derived from the Ancestral Recombination Graph (ARG)***

434

435 We estimated local trees by inferring ARGs using the arg-sample module in ARGweaver version  
436 1 [37]. We ran ARGweaver on unphased data in 5-7 mb intervals which included the areas of high  
437 differentiation in the three contigs with outlier peaks, plus a control set consisting of 10 contigs totaling

438 ~110 megabases. For the control set we broke large contigs into 5 mb fragments and ran the software  
439 following the settings described previously [17]. Briefly, we set the mutation and the recombination rates  
440 to 10–9/bp/gen, the effective population size to 500,000 individuals and the remaining parameters as  
441 follows: -c 5 -ntimes 20 -maxtime 1e7 -delta 0.005 -resample-window-iters 1 -resample-window 10000 -n  
442 1000. We retained the last of 1,000 MCMC iterations, discarded the first and last 50 kb of each ARG  
443 block, and extracted trees every 500 bp. We used these trees to calculate three statistics as described in  
444 detail in [9] (Figure 4A). The species enrichment score measures the probability of observing subtrees of  
445 different sizes containing individuals from a certain species, and is calculated for both species and every  
446 subtree. For each tree we retained the maximum enrichment score obtained for a given species across all  
447 subtrees. A high enrichment score is obtained when a local tree has a clade with most or all individuals  
448 from that species. We note that when only two species are present in a tree like in this study, when one  
449 species enrichment score is high it means that individuals from the second species may also group  
450 together on the tree, and therefore both species enrichment scores will tend to be elevated. The relative  
451 TMRCA half-life [9,37], RTH', is a normalized version of the time to most recent common ancestry  
452 (TMRCA). To calculate this statistic, we divided the time to the most recent common ancestor of half of  
453 the haploid samples for a species by the age of the youngest subtree containing at least half of all the  
454 haploid samples (irrespective of species). This normalization controls for the variation in coalescence  
455 times in different trees across the genome [37]. Finally, the inter-species cross coalescence measures the  
456 average age of the most recent coalescence events between tips in the tree which belong to the different  
457 species. We added the ages of the ten most recent cross-coalescence events in each tree and  
458 normalized them in the same way we did to calculate RTH' [9]. We obtained these statistics for each tree  
459 (sampled every 500 bp) and then averaged those values across 5 kb non-overlapping windows. We  
460 generated empirical distributions for these statistics from the ~22,000 windows (~110 mb) obtained from  
461 the control contigs, and used these distributions to assess statistical significance as described previously  
462 [9,17]. The analysis was conducted separately for allopatric and sympatric samples, and a third time for  
463 all samples combined. We exported trees which show extreme enrichment scores for each species for  
464 illustration purposes. The code to calculate these three statistics and conduct these statistical tests is  
465 deposited in GitHub ([https://github.com/CshISiepelLab/bird\\_capuchino\\_analysis](https://github.com/CshISiepelLab/bird_capuchino_analysis); see [9]).

466

## 467 **Acknowledgements**

468 We thank Bronwyn Butcher for help with library preparation and sequencing. This project was funded  
469 through the Athena grant from the Cornell Lab of Ornithology (to T.N.N.) and NSF DEB-2232929 (to  
470 L.C.). The Brazilian Federal Agencies (CNPq and CAPES) provided scholarships to M.R. The Neotropical  
471 Grassland Conservancy (NGC) supported field work. The Brazilian National Research Council (CNPq  
472 grants (303318/2013, 309438/2016) supported C.S.F. We would like to acknowledge the property owners  
473 who kindly authorized our fieldwork. We thank Renata Grieco for permission to reproduce the bird  
474 illustrations in Figure 1. The sequence data generated for this project is archived on GenBank (Bio  
475 Project PRJNA382416).

476

## 477 **Author contributions**

478 Designed project: T.N.N., M.R. and L.C. Performed fieldwork: M.R. and C.S.F. Obtained genomic data:  
479 T.N.N. Analyzed genomic data: T.N.N. and L.C. Supervised research: L.C. Wrote manuscript: T.N.N. and  
480 L.C., with edits from all authors.

481

## 482 **Declaration of interests**

483 The authors declare no competing interests.

484

## 485 **Figure titles and legends**

486

487 **Figure 1: Genome-wide differentiation between *S. beltoni* and *S. plumbea*.** (A) Range map for both  
488 species showing sampling localities in Brazil, and the area of contact (encompassed in the rectangle on  
489 the map and also in the topographic map in the inset) following [34]. Samples are color-coded by species  
490 and geographic origin (as shown in C). (B) Examples of adult male *S. beltoni* and *S. plumbea* individuals;  
491 photos by Márcio Repenning. (C) PCA derived from ~15 million genome-wide SNPs. (D) Haplotype  
492 network derived from sequences of the mitochondrial gene *COI*. Different haplotypes are connected

493 either by solid (based on the minimum spanning tree) or by dashed lines (indicating alternative paths).

494 The number of mutational steps between haplotypes is shown as solid bars.

495

496 **Figure 2: Demographic reconstruction.** (A) Current and ancestral population sizes, (B) divergence time  
497 and (C) migrants per generation. Each set of parameters was estimated three times, once with a random  
498 subset of all samples combined and also after subsampling either exclusively allopatric or sympatric  
499 individuals. Bars indicate median values and 95% credible intervals.

500

501 **Figure 3: Outlier peaks between *S. beltoni* and *S. plumbea*.** (A) Manhattan plot showing  $F_{ST}$  values  
502 calculated using all individuals from both species. Note that for simplicity the plot displays the 1.1% of  
503 SNPs with  $F_{ST} > 0.2$ . The inset shows a histogram of the distribution of  $F_{ST}$  values from all genome-wide  
504 SNPs, with those showing values above 0.7 in red. The plots in (B) represent a zoom into each peak  
505 shown in (A) and display the genes that are annotated in that region. The insets show PCAs derived from  
506 the variants found in each peak. Samples are color-coded by species and the geographic region to which  
507 they belong (as in Figure 1).

508

509 **Figure 4: ARG-based statistics.** (A) Graphical representation of how species enrichments,  $RTH'$  and  
510 cross-coalescence values are calculated (modified from [17]). The species enrichment score is calculated  
511 for each subtree (example shaded in gray) in a topology and represents the probability of observing that  
512 number of samples of a particular species (the one represented in red in the example) in that subtree  
513 under a hypergeometric distribution. The score for a given tree and species is the maximum score  
514 associated with the full tree.  $RTH'$  is the ratio between the TMRCA of half of the samples from a species  
515 and the age of the youngest subtree that contains at least half of all samples (irrespective of species).  
516 The cross-coalescence time is calculated by adding the time of the ten most recent cross-coalescence  
517 events in the tree and dividing it by the same value used to normalize TMRCA to obtain  $RTH'$ . In the  
518 example the color of the leaves indicates two hypothetical species. (B) Plots showing these statistics for  
519 contig 404 in 5 kb windows, with the peak region shown by dashed vertical blue lines. These particular  
520 plots show statistics derived from individuals sampled outside of the hybrid zone. Species enrichment

521 values for *S. beltoni* and *S. plumbea* are higher than the genome-wide threshold of statistical significance,  
522 while RTH' values are significantly lower. Therefore, subtrees from topologies in this region contain most  
523 individuals from each species and are comparatively younger than in the rest of the genome. The cross-  
524 coalescence times in the peak region are also significantly delayed, meaning that there is less recent  
525 gene flow than across other areas of the genome.

526 **Tables**

527

528 **Table 1. Regions of high divergences between species.** Details for each outlier peak, including the genes they contain and a summary of their  
 529 functions.

530

531

Contig	Peak size (kb)	Highest SNP $F_{ST}$ in window	Number of SNPs with $F_{ST} > 0.7$	Number of genes	Genes	Gene function
404	102	0.794	24	2	<i>PRLR</i> , <i>AGXT2</i>	<p><i>PRLR</i> encodes for the prolactin hormone receptor affecting breeding behavior and parental care in birds. <i>PRLR</i> also plays a role in keratinization.</p> <p><i>AGXT2</i> encodes for a protein which catalyzes the conversion of glyoxylate to glycine.</p>
33	66	0.901	9	1	<i>EDN3</i>	<p>Encodes a protein affecting the development of melanocytes, resulting in both hypo and hyperpigmentation in chickens and is also related to egg production.</p>

---

						<i>DNAAF1</i> encodes a protein required for cilia stability.
382	2	0.795	6	3	<i>DNAAF1</i> , <i>HSDL1</i> , <i>MBTPS1</i>	<i>HSDL1</i> enables oxidoreductase activity.
						<i>MBTPS1</i> encodes a member of the subtilisin-like proprotein convertase family.

---

532

533 **Table 2: Results from statistical tests conducted on ARG-based statistics.** Tests were conducted separately for allopatric, sympatric, or all  
 534 samples combined. Species enrichment and RTH' scores were compared to the distribution of these values across a set of control contigs without  
 535 outlier peaks. Cross-coalescence (CC) values were compared between the peaks and their adjacent regions, and for control regions, between  
 536 focal and adjacent areas. The distribution of the difference in CC values between focal and adjacent areas was assessed across all control  
 537 regions and used to determine the statistical significance of the difference in CC values observed between the peak regions and those areas  
 538 immediately adjacent to the peaks. Statistically significant results are shown by asterisks, which also denote the level of significance (\*0.005;  
 539 \*\*0.001). Note that we ran species enrichment and RTH' tests for each species, while the cross-coalescence test is run for the species pair and  
 540 therefore only presented once in the top of the table.  
 541

		peak on contig 404			peak on contig 33			peak on contig 382		
		Species enrichment	RTH'	CC	Species enrichment	RTH'	CC	Species enrichment	RTH'	CC
<i>S. beltoni</i>	all samples	8.95**	0.246**	0.15	7.4**	0.367*	-0.05	4.15**	0.845	0.3*
	allopatric	9.16**	0.239*	0.32*	6.39**	0.346*	0.17	4.39**	0.772	0.5**
	sympatric	1.83	0.545	0.07	2.44**	0.49*	0.06	1.25	0.72	0.06
<i>S. plumbea</i>	all samples	9.28**	0.235**		4.76**	0.709		5.66**	0.753	
	allopatric	10.34**	0.193*		4.65**	0.742		6.24**	0.469*	
	sympatric	1.76	0.44*		1.7	0.693		1.33	0.828	

542 **References**

543

544 1. Schluter D. The ecology of adaptive radiation. Oxford University Press; 2000.

545 2. Price T. Speciation in birds. Roberts and Co.; 2008.

546 3. Coyne JA, Orr A. Speciation. Oxford University Press; 2004.

547 4. Grant PR, Grant BR. How and why species multiply: the radiation of Darwin's finches. Princeton  
548 University Press; 2007.

549 5. Seehausen O, Butlin RK, Keller I, Wagner CE, Boughman JW, Hohenlohe PA, et al. Genomics and  
550 the origin of species. *Nature Reviews Genetics*. 2014;15: 176–192. doi:10.1038/nrg3644

551 6. Barton NH, Gale KS, Harrison R. Genetic analysis of hybrid zones. Hybrid zones and the  
552 evolutionary process. Oxford University Press; 1993: 13–45.

553 7. Klopstein S, Currat M, Excoffier L. The fate of mutations surfing on the wave of a range expansion.  
554 *Molecular Biology and Evolution*. 2006;23: 482–490. doi:10.1093/molbev/msj057

555 8. Cruickshank TE, Hahn MW. Reanalysis suggests that genomic islands of speciation are due to  
556 reduced diversity, not reduced gene flow. *Molecular ecology*. 2014;23: 3133–3157.

557 9. Hejase HA, Salman-Minkov A, Campagna L, Hubisz MJ, Lovette IJ, Gronau I, et al. Genomic islands  
558 of differentiation in a rapid avian radiation have been driven by recent selective sweeps.  
559 *Proceedings of the National Academy of Sciences of the United States of America*. 2020;48: 30554-  
560 30565 doi:10.1101/2020.03.07.977694

561 10. Toews DPL, Taylor SA, Vallender R, Brelsford A, Butcher BG, Messer PW, et al. Plumage genes and  
562 little else distinguish the genomes of hybridizing warblers. *Current Biology*. 2016;26: 2313–2318.  
563 doi:10.1016/j.cub.2016.06.034

564 11. Irwin DE, Alcaide M, Delmore KE, Irwin JH, Owens GL. Recurrent selection explains parallel evolution  
565 of genomic regions of high relative but low absolute differentiation in a ring species. *Molecular  
566 Ecology*. 2016;25(18): 4488-4507.  
567

568 12. Lewanski AL, Grundler MC, Bradburd GS. The era of the ARG: an empiricist's guide to ancestral  
569 recombination graphs. *ArXiv [Preprint]*. 2023 Oct 18:arXiv:2310.12070v1. PMID: 37904740; PMCID:  
570 PMC10614969.

571 13. Poelstra JW, Vijay N, Bossu CM, Lantz H, Ryll B, Muller I, et al. The genomic landscape underlying  
572 phenotypic integrity in the face of gene flow in crows. *Science*. 2014;344: 1410–1414.  
573 doi:10.1126/science.1253226

574 14. Campagna L, Repenning M, Silveira LF, Fontana CS, Tubaro PL, Lovette IJ. Repeated divergent  
575 selection on pigmentation genes in a rapid finch radiation. *Science Advances*. 2017;3: e1602404.  
576 doi:10.1126/sciadv.1602404

577 15. Stryjewski KF, Sorenson MD. Mosaic genome evolution in a recent and rapid avian radiation. *Nature  
578 Ecology and Evolution*. 2017;1: 1912–1922. doi:10.1038/s41559-017-0364-7

- 579 16. Wang S, Rohwer S, Zwaan DR, Toews DPL, Lovette IJ, Mackenzie J, et al. Selection on a small  
580 genomic region underpins differentiation in multiple color traits between two warbler species.  
581 *Evolution Letters*. 2020;4: 502–515. doi:10.1002/evl3.198
- 582 17. Campagna L, Mo Z, Siepel A, Uy JAC. Selective sweeps on different pigmentation genes mediate  
583 convergent evolution of island melanism in two incipient bird species. Bierné N, editor. *PLoS*  
584 *Genetics* 2022;18: e1010474. doi:10.1371/journal.pgen.1010474
- 585 18. Uy JAC, Irwin DE, Webster MS. Behavioral isolation and incipient speciation in birds. *Annual Review*  
586 *of Ecology, Evolution, and Systematics*. 2018;49: 1–24. doi:10.1146/annurev-ecolsys-110617-  
587 062646
- 588 19. Campagna L, Toews DPL. The genomics of adaptation in birds. *Current Biology*. 2022;32: R1173–  
589 R1186. doi:10.1016/j.cub.2022.07.076
- 590 20. Turbek SP, Browne M, Di Giacomo AS, Kopuchian C, Hochachka WM, Estalles C, et al. Rapid  
591 speciation via the evolution of pre-mating isolation in the Iberá Seedeater. *Science*. 2021;371:  
592 eabc0256. doi:10.1126/science.abc0256
- 593 21. Marques DA, Meier JI, Seehausen O. A combinatorial view on speciation and adaptive radiation.  
594 *Trends in Ecology & Evolution*. 2019; S0169534719300552. doi:10.1016/j.tree.2019.02.008
- 595 22. Rosenblum EB, Sarver BAJ, Brown JW, Des Roches S, Hardwick KM, Hether TD, et al. Goldilocks  
596 meets Santa Rosalia: An ephemeral speciation model explains patterns of diversification across  
597 time scales. *Evolutionary Biology* 2012;39: 255–261. doi:10.1007/s11692-012-9171-x
- 598 23. Kautt AF, Kratochwil CF, Nater A, Machado-Schiaffino G, Olave M, Henning F, et al. Contrasting  
599 signatures of genomic divergence during sympatric speciation. *Nature*. 2020;588: 106–111.  
600 doi:10.1038/s41586-020-2845-0
- 601 24. Winkler DW, Billerman SM, and Lovette IJ. Tanagers and Allies (Thraupidae), version 1.0. In *Birds of*  
602 *the World* (SM Billerman, BK Keeney, PG Rodewald, and TS Schulenberg, Editors). Cornell Lab of  
603 Ornithology, Ithaca, NY, USA. 2020. <https://doi.org/10.2173/bow.thraup2.01>.
- 604 25. Burns KJ, Shultz AJ, Title PO, Mason NA, Barker FK, Klicka J, et al. Phylogenetics and diversification  
605 of tanagers (Passeriformes: Thraupidae), the largest radiation of Neotropical songbirds. *Molecular*  
606 *Phylogenetics and Evolution*. 2014;75: 41–77. doi:10.1016/j.ympev.2014.02.006
- 607 26. Mason NA, Olvera-Vital A, Lovette IJ, Navarro-Sigüenza AG. Hidden endemism, deep polyphyly, and  
608 repeated dispersal across the Isthmus of Tehuantepec: Diversification of the White-collared  
609 seedeater complex (Thraupidae: *Sporophila torqueola*). *Ecology and Evolution* 2018;8: 1867–1881.  
610 doi:10.1002/ece3.3799
- 611 27. Ocampo D, Winker K, Miller MJ, Sandoval L, Uy JAC. Rapid diversification of the Variable Seedeater  
612 superspecies complex despite widespread gene flow. *Molecular Phylogenetics and Evolution*.  
613 2022;173: 107510. doi:10.1016/j.ympev.2022.107510
- 614 28. Campagna L, Benites P, Loughheed SC, Lijtmaer DA, Di Giacomo AS, Eaton MD, et al. Rapid  
615 phenotypic evolution during incipient speciation in a continental avian radiation. *Proceedings of the*  
616 *Royal Society of London. Series B: Biological Sciences*. 2012;279: 1847–1856.  
617 doi:10.1098/rspb.2011.2170
- 618 29. Repenning M, Fontana CS. Distinguishing females of capuchino seedeaters: call repertoires provide  
619 evidence for species-level diagnosis. *Revista Brasileira de Ornitologia*. 2019;27: 70–78.  
620 doi:10.1007/BF03544451

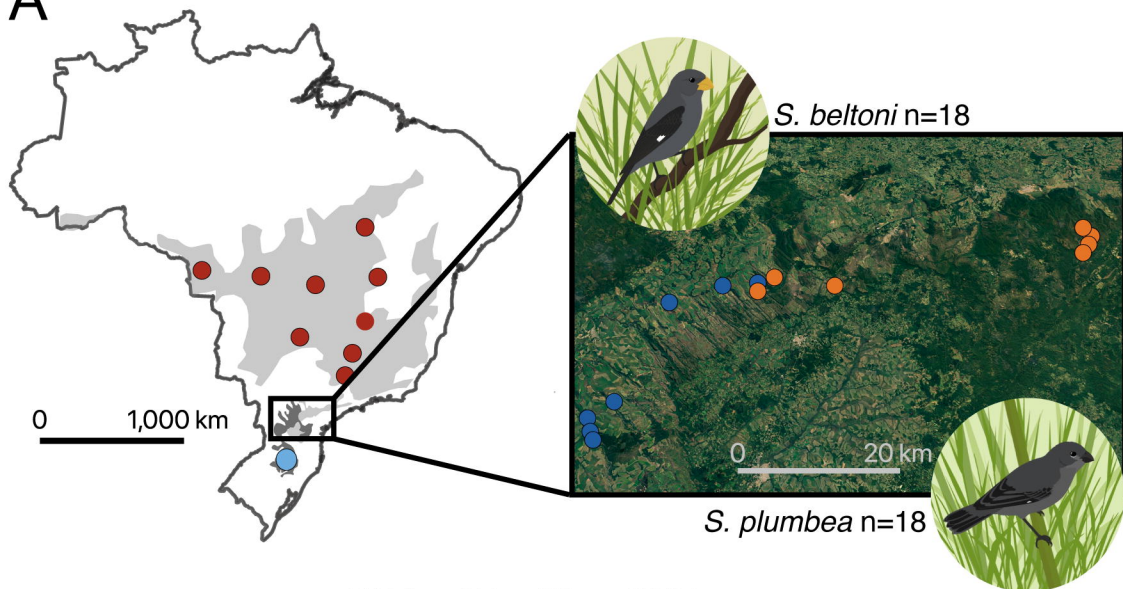
- 621 30. Estalles C, Turbek SP, José Rodríguez-Cajarville M, Silveira LF, Wakamatsu K, Ito S, et al.  
622 Concerted variation in melanogenesis genes underlies emergent patterning of plumage in  
623 capuchino seedeaters. *Proceedings of the Royal Society of London. Series B: Biological Sciences*.  
624 2022;289: 20212277. doi:10.1098/rspb.2021.2277
- 625 31. Stiles FG. When black plus white equals gray: The nature of variation in the variable seedeater  
626 complex (Emberizinae: *Sporophila*). *Ornitologia Neotropical*. 1996;7: 75–107.
- 627 32. Ocampo D, Winker K, Miller MJ, Sandoval L, Uy JAC. Replicate contact zones suggest a limited role  
628 of plumage in reproductive isolation among subspecies of the variable seedeater (*Sporophila*  
629 *corvina*). *Molecular Ecology*. 2023; mec.16938. doi:10.1111/mec.16938
- 630 33. Hellmayr CE. Catalogue of birds of the Americas and the adjacent islands in Field Museum of Natural  
631 History. Chicago; 1918. doi:10.5962/bhl.title.5570
- 632 34. Repenning M, Fontana CS. A new species of gray seedeater (Emberizidae: *Sporophila*) from upland  
633 grasslands of southern Brazil. *The Auk*. 2013;130: 791–803. doi:10.1525/auk.2013.12167
- 634 35. Repenning M. Variacao geografica em *Sporophila* (Aves: Thraupidae) com base em evidencias  
635 fenotipicas, ecologicas e geneticas. PhD Thesis, Pontificia Universidade Católica do Rio Grande do  
636 Sul. 2017. Available: <http://tede2.pucrs.br/tede2/handle/tede/7598>
- 637 36. Thawornwattana Y, Huang J, Flouri T š, Mallet J, Yang Z. Inferring the direction of introgression using  
638 genomic sequence data. *Molecular Biology and Evolution*. 2023; msad178.  
639 doi:10.1093/molbev/msad178
- 640 37. Rasmussen MD, Hubisz MJ, Gronau I, Siepel A. Genome-Wide inference of ancestral recombination  
641 graphs. *PLoS Genetics*. 2014;10: e1004342. doi:10.1371/journal.pgen.1004342
- 642 38. Lijtmaer DA, Sharpe NM, Tubaro PL, Loughheed SC. Molecular phylogenetics and diversification of the  
643 genus *Sporophila* (Aves: Passeriformes). *Molecular Phylogenetics and Evolution*. 2004;33(3): 562-  
644 579. doi: 10.1016/j.ympev.2004.07.011
- 645 39. Campagna L, Silveira LF, Tubaro PL, Loughheed SC. Identifying the sister species to the rapid  
646 capuchino seedeater radiation (Passeriformes: *Sporophila*). *The Auk*. 2013;130: 645–655.  
647 doi:10.1525/auk.2013.13064
- 648 40. Wolf JBW, Ellegren H. Making sense of genomic islands of differentiation in light of speciation.  
649 *Nature Reviews Genetics*. 2017;18: 87–100. doi:10.1038/nrg.2016.133
- 650 41. Akiyama T, Kinoshita K. Body color expression in birds. *Pigments, Pigment Cells and Pigment*  
651 *Patterns*. 2021; 91–126.
- 652 42. Maclary ET, Wauer R, Phillips B, Brown A, Boer EF, Samani AM, et al. An allelic series at the  
653 *EDNRB2* locus controls diverse piebalding patterns in the domestic pigeon. *PLoS genetics*. 2023;  
654 19(10), p.e1010880. <https://doi.org/10.1371/journal.pgen.1010880>
- 655 43. Dong X, Li J, Zhang Y, Han D, Hua G, Wang J, et al. Genomic analysis reveals pleiotropic alleles at  
656 *EDN3* and *BMP7* involved in chicken comb color and egg production. *Frontiers in Genetics*.  
657 2019;10: 612. doi:10.3389/fgene.2019.00612
- 658 44. Luo C, Shen X, Rao Y, Xu H, Tang J, Sun L, et al. Differences of Z chromosome and genomic  
659 expression between early- and late-feathering chickens. *Molecular Biology Reports*. 2012;39: 6283–  
660 6288. doi:10.1007/s11033-012-1449-7

- 661 45. Kuenzel W. Neurobiology of molt in avian species. *Poultry Science*. 2003;82: 981–991.  
662 doi:10.1093/ps/82.6.981
- 663 46. Skold H, Amundsen T, Svensson P, Mayer I, Bjelvenmark J, Forsgren E. Hormonal regulation of  
664 female nuptial coloration in a fish. *Hormones and Behavior*. 2008;54: 549–556.  
665 doi:10.1016/j.yhbeh.2008.05.018
- 666 47. Hooper DM, Griffith SC, Price TD. Sex chromosome inversions enforce reproductive isolation across  
667 an avian hybrid zone. *Molecular ecology*. 2019; 28(6): 1246–1262.  
668 <https://doi.org/10.1111/mec.14874>
- 669 48. McDiarmid CS, Finch F, Peso M, van Rooij E, Hooper DM, Rowe M, Griffith SC. Experimentally  
670 testing mate preference in an avian system with unidirectional bill color introgression. *Ecology and*  
671 *Evolution*. 2023;13(2): e9812. <https://doi.org/10.1002/ece3.9812>
- 672 49. Mason NA, Burns KJ. Molecular phylogenetics of the Neotropical seedeaters and seed-finches  
673 (*Sporophila*, *Oryzoborus*, *Dolospingus*). *Ornitologia Neotropical*. 2013;24: 139–155.
- 674 50. Baiz MD, Wood AW, Brelsford A, Lovette IJ, Toews DPL. Pigmentation genes show evidence of  
675 repeated divergence and multiple bouts of introgression in *Setophaga* warblers. *Current Biology*.  
676 2021;31: 643–649.e3. doi:10.1016/j.cub.2020.10.094
- 677 51. Grabenstein KC, Taylor SA. Breaking barriers: Causes, consequences, and experimental utility of  
678 human-mediated hybridization. *Trends in Ecology & Evolution*. 2018;33: 198–212.  
679 doi:10.1016/j.tree.2017.12.008
- 680 52. Vijay N, Bossu CM, Poelstra JW, Weissensteiner MH, Suh A, Kryukov AP, et al. Evolution of  
681 heterogeneous genome differentiation across multiple contact zones in a crow species complex.  
682 *Nature Communications*. 2016;7: 13195. doi:10.1038/ncomms13195
- 683 53. Semenov GA, Linck E, Enbody ED, Harris RB, Khaydarov DR, Alström P, et al. Asymmetric  
684 introgression reveals the genetic architecture of a plumage trait. *Nature Communications*. 2021;12:  
685 1019. doi:10.1038/s41467-021-21340-y
- 686 54. Walsh J, Kovach AI, Olsen BJ, Shriver WG, Lovette IJ. Bidirectional adaptive introgression between  
687 two ecologically divergent sparrow species. *Evolution*. 2018;72: 2076–2089. doi:10.1111/evo.13581
- 688 55. Schubert M, Lindgreen S, Orlando L. AdapterRemoval v2: rapid adapter trimming, identification, and  
689 read merging. *BMC Research Notes*. 2016; 9: 88. <https://doi.org/10.1186/s13104-016-1900-2>
- 690 56. Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. *Nature Methods*. 2012;9:  
691 357–359. doi:10.1038/nmeth.1923
- 692 57. García-Alcalde F, Okonechnikov K, Carbonell J, Cruz LM, Götz S, Tarazona S, et al. Qualimap:  
693 evaluating next-generation sequencing alignment data. *Bioinformatics*. 2012;28: 2678–2679.  
694 doi:10.1093/bioinformatics/bts503
- 695 58. Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA, et al. The variant call format  
696 and VCFtools. *Bioinformatics*. 2011;27: 2156–2158. doi:10.1093/bioinformatics/btr330
- 697 59. Van der Auwera GA, O'Connor BD. Genomics in the cloud: Using docker, GATK, and WDL in Terra.  
698 O'Reilly Media, Incorporated; 2020.

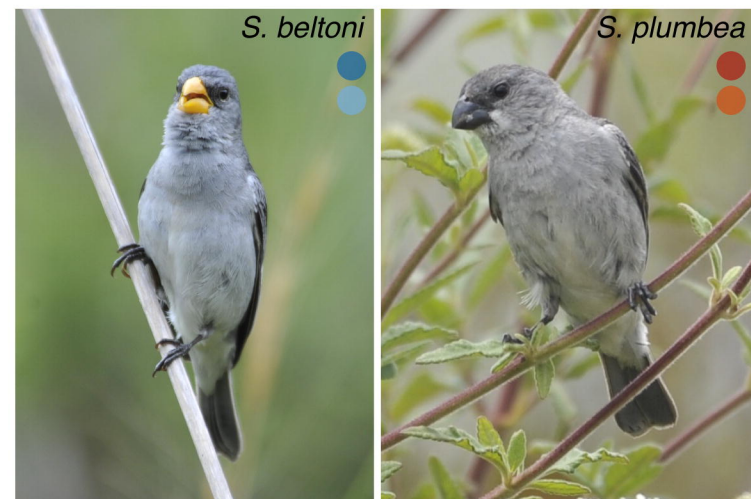
- 699 60. Hebert PDN, Ratnasingham S, de Waard JR. Barcoding animal life: cytochrome c oxidase subunit 1  
700 divergences among closely related species. *Proceedings of the Royal Society of London. Series B:*  
701 *Biological Sciences*. 2003;270. doi:10.1098/rsbl.2003.0025
- 702 61. Kerr KCR, Lijtmaer DA, Barreira AS, Hebert PDN, Tubaro PL. Probing evolutionary patterns in  
703 Neotropical birds through DNA barcodes. DeSalle R, editor. *PLoS ONE*. 2009;4: e4379.  
704 doi:10.1371/journal.pone.0004379
- 705 62. Kearse M, Moir R, Wilson A, Stones-Havas S, Cheung M, Sturrock S, et al. Geneious Basic: An  
706 integrated and extendable desktop software platform for the organization and analysis of sequence  
707 data. *Bioinformatics*. 2012;28: 1647–1649. doi:10.1093/bioinformatics/bts199
- 708 63. R Core Team. *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R  
709 Foundation for Statistical Computing; 2021. Available: <https://www.R-project.org/>
- 710 64. Paradis E. pegas: an R package for population genetics with an integrated–modular approach.  
711 *Bioinformatics*. 2010;26: 419–420. doi:10.1093/bioinformatics/btp696
- 712 65. Zheng X, Levine D, Shen J, Gogarten S, Laurie C, Weir B. AhHigh-performance computing toolset for  
713 relatedness and principal component analysis of SNP data. *Bioinformatics*. 2012;28: 3326–3328.  
714 doi:10.1093/bioinformatics/bts606
- 715 66. Turner SD. qqman: an R package for visualizing GWAS results using Q-Q and manhattan plots.  
716 *Journal of Open Source Software*. 2018;3: 731–732.
- 717 67. Garud NR, Messer PW, Buzbas EO, Petrov DA. Recent selective sweeps in North American  
718 *Drosophila melanogaster* show signatures of soft sweeps. *PLoS Genetics*. 2015;11: e1005004.
- 719 68. Browning SR, Browning BL. Rapid and accurate haplotype phasing and missing-data inference for  
720 whole-genome association studies by use of localized haplotype clustering. *The American Journal*  
721 *of Human Genetics*. 2007/09/21 ed. 2007;81: 1084–1097. doi:10.1086/521987
- 722 69. Revell LJ. phytools: an R package for phylogenetic comparative biology (and other things): phytools:  
723 R package. *Methods in Ecology and Evolution*. 2012;3: 217–223. doi:10.1111/j.2041-  
724 210X.2011.00169.x
- 725 70. Oksanen J, Blanchet FG, Friendly M, Kindt R, Legendre P, McGlenn D, et al. vegan: Community  
726 ecology package. R package version 2.5-6. <https://CRAN.R-project.org/package=vegan>; 2019.  
727 Available: <https://CRAN.R-project.org/package=vegan>
- 728
- 729 71. Alexander DH, Novembre J, Lange K. Fast model-based estimation of ancestry in unrelated  
730 individuals. *Genome Research*. 2009;19: 1655–1664.
- 731
- 732 72. Wickham H. ggplot2: Elegant Graphics for Data Analysis. Springer-Verlag New York; 2016. Available:  
733 <https://ggplot2.tidyverse.org>
- 734
- 735 73. Korunes KL, Samuk K. pixy: Unbiased estimation of nucleotide diversity and divergence in the  
736 presence of missing data. *Molecular Ecology Resources*. 2021;21: 1359–1368.
- 737 74. Gronau I, Hubisz MJ, Gulko B, Danko CG, Siepel A. Bayesian inference of ancient human  
738 demography from individual genome sequences. *Nature Genetics*. 2011;43: 1031–1034.  
739 doi:10.1038/ng.937
- 740 75. Plummer M, Best N, Cowles K, Vines K. CODA: Convergence diagnosis and output analysis for  
741 MCMC. *R News*. 2006;6: 7–11.

- 742 76. Smeds L, Qvarnström A, Ellegren H. Direct estimate of the rate of germline mutation in a bird.  
743 Genome Research. 2016;26: 1211–1218. doi:10.1101/gr.204669.116

A

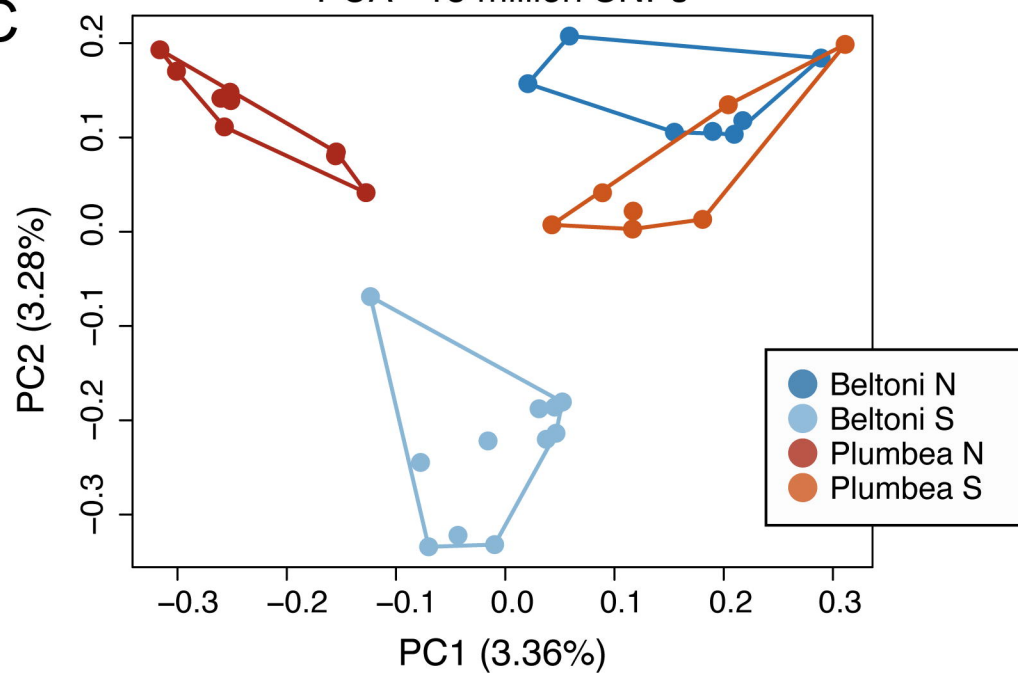


B

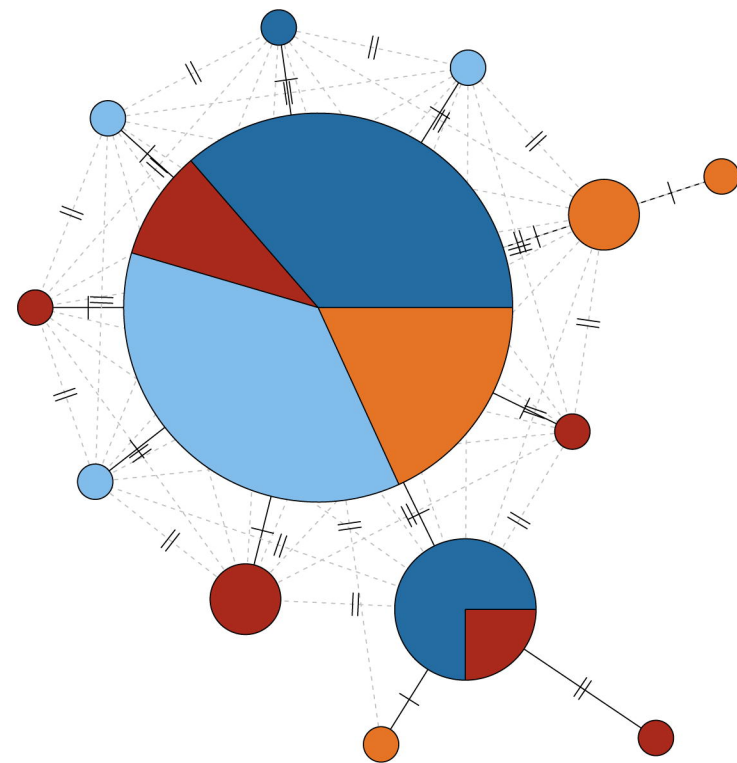


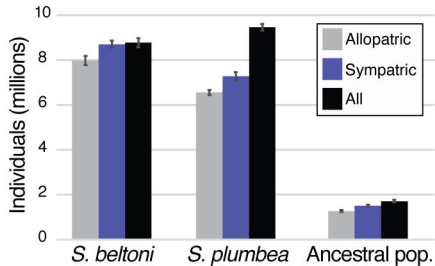
C

PCA ~15 million SNPs

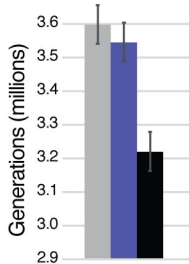


D

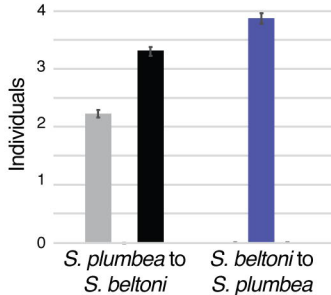


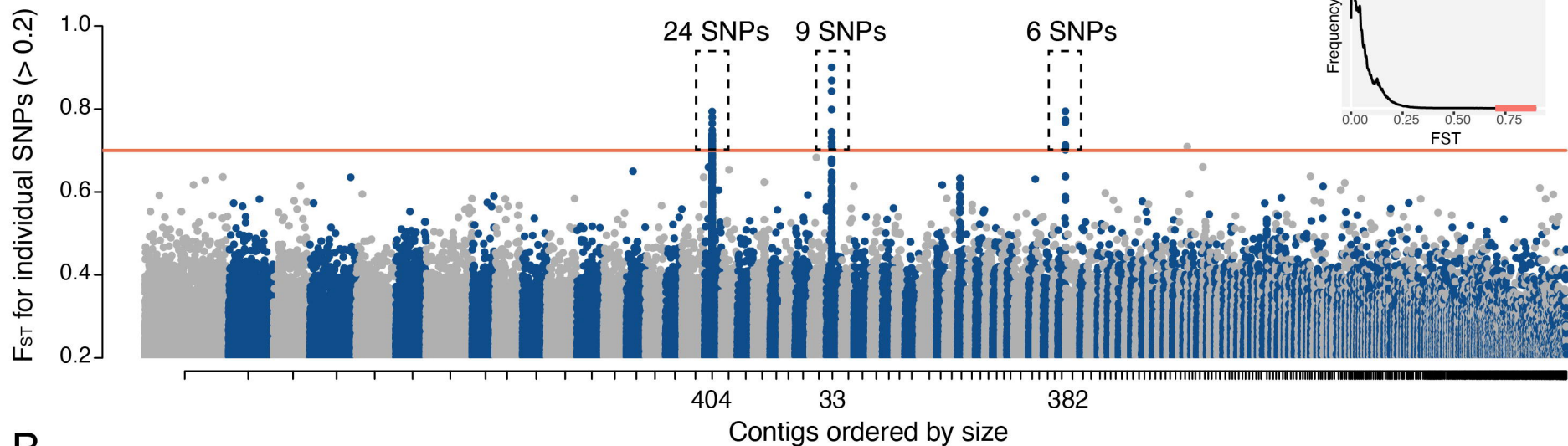
**A**Effective population sizes ( $N_e$ )**B**

## Divergence time

**C**

## Migrants per generation



**A****B**