1

2    DR. YING ZHEN (Orcid ID : 0000-0001-8635-9687)

3    DR. ERIC C ANDERSON (Orcid ID : 0000-0003-1326-0840)

4    MISS SIRENA LAO (Orcid ID : 0000-0003-2907-7266)

5

6

7    Article type    : Original Article

8

9

10    **Title:** Genomic divergence across ecological gradients in the Central African rainforest

11    songbird (*Andropadus virens*)

12

13    **Authors:**

14    *Ying Zhen*[a,b*]*, Ryan J. Harrigan*[b]*, Kristen C. Ruegg*[b,c]*, Eric C. Anderson*[d]*, Thomas C. Ng*[e]*,*

15    *Sirena Lao*[b]*, Kirk E. Lohmueller*[a,f] *and Thomas B. Smith* [a,b]

16

17    **Affiliations:**

18    [a] Department of Ecology and Evolutionary Biology, University of California, Los Angeles,

19    610 Charles E Young Drive East, Los Angeles, CA 90095, USA.

20    [b] Center for Tropical Research, Institute of Environment and Sustainability, University of

21    California, Los Angeles, La Kretz Hall, Suite 300, Los Angeles, CA 90095, USA.

22    [c] Department of Ecology and Evolutionary Biology, University of California, Santa Cruz,

23    Santa Cruz, CA 95060, USA.

24    [d] Fisheries Ecology Division, Southwest Fisheries Science Center, National Marine

25    Fisheries Service, 110 McAllister Way, Santa Cruz, CA 95060, USA.

26    [e] Department of Biomolecular Engineering, University of California, Santa Cruz, CA

27    95060, USA.

28  [f]Department of Human Genetics, David Geffen School of Medicine, University of

29  California, Los Angeles, Los Angeles, CA, 90095, USA.

30

31

34

35  [*]**Corresponding author:**

36  Ying Zhen

37  Address: Center for Tropical Research, Institute of Environment and Sustainability,

38  University of California, Los Angeles, La Kretz Hall, Suite 300, Los Angeles, CA 90095,

39  USA.

40  Phone: 310-206-6234

41  Fax: 310-825-5446

42  Email: zhen@g.ucla.edu

43

44  **Running title:** Genomic divergence across ecological gradients

45

46  **Abstract**

47  The little greenbul, a common rainforest passerine from sub-Saharan Africa, has been the

48  subject of long-term evolutionary studies to understand the mechanisms leading to

49  rainforest speciation. Previous research found morphological and behavioral divergence

50  across rainforest-savanna transition zones (ecotones), and a pattern of divergence with

51  gene flow suggesting divergent natural selection has contributed to adaptive divergence

52  and ecotones could be important areas for rainforests speciation. Recent advances in

53  genomics and environmental modeling make it possible to examine patterns of genetic

54  divergence in a more comprehensive fashion. To assess the extent to which natural

55  selection may drive patterns of differentiation, here we investigate patterns of genomic

56  differentiation among populations across environmental gradients and regions. We find

57  compelling evidence that individuals form discrete genetic clusters corresponding to

58  distinctive environmental characteristics and habitat types. Pairwise $F_{ST}$ between

59  populations in different habitats is significantly higher than within habitats, and this

60  differentiation is greater than what is expected from geographic distance alone. Moreover,

61  we identified 140 SNPs that showed extreme differentiation among populations through a

62  genome-wide selection scan. These outliers were significantly enriched in exonic and

63  coding regions, suggesting their functional importance. Environmental association analysis

64  of SNP variation indicates that several environmental variables, including temperature and

65  elevation, play important roles in driving the pattern of genomic diversification. Results

66  lend important new genomic evidence for environmental gradients being important in

67  population differentiation.

68

69  **Introduction**

70      Rainforests are heralded for their exceptionally high biological diversity, yet the

71  evolutionary mechanisms for the generation and maintenance of this diversity have been

72  debated for decades (Haffer 1969; Mayr & O'Hara 1986; Martin 1991; Smith *et al.* 1997;

73  Schneider *et al.* 1999; Moritz *et al.* 2000; Ogden & Thorpe 2002; Price 2008; Schluter

74  2009; Hoorn *et al.* 2010; Ribas *et al.* 2011; Smith *et al.* 2014; Beheregaray *et al.* 2015).

75  Models of rainforest speciation abound. Some emphasize the importance of neutral

76  processes, for example genetic drift in allopatric populations isolated by historical refugia

77  (Haffer 1969), some favor processes such as landscape change (Hoorn *et al.* 2010; Ribas *et*

78  *al.* 2011) or dispersal (Smith *et al.* 2014), while others point toward a dominant role of

79  divergent natural selection across ecological gradients and ecotones (Smith *et al.* 1997;

80  Schneider *et al.* 1999; Ogden & Thorpe 2002; Smith *et al.* 2005; Schluter 2009; Smith *et*

81  *al.* 2011; Beheregaray *et al.* 2015). Each process is expected to shape the genomes of

82  natural populations in different ways, leaving a signal that provides insights into the

83  evolutionary mechanisms that may have led to divergence. Such information is of

84  importance not only to evolutionary geneticists interested in understanding the processes

85  involved in speciation, but also to conservation decision makers, who are interested in

86  preserving biodiversity and prioritizing new regions for protection in the face of rapid

87  anthropogenic change and climate change.

88    In this study, we explore the roles that population-level processes play in shaping

89    biodiversity in Central Africa by examining the genomic diversity in a common songbird,

90    the little greenbul (*Andropadus virens*). The little greenbul provides a particularly useful

91    taxon for this inquiry because it has a broad geographic distribution across sub-Saharan

92    Africa where it occurs in ecologically diverse habitats, and has been the subject of long-

93    term studies of intra-specific diversity and speciation. In the case of *A. virens*, as well as

94    some other rainforest taxa, the rainforest-savanna transition zones (ecotones) have been

95    shown to drive phenotypic divergence and likely speciation (Smith *et al.* 1997, 2005;

96    Kirschel *et al.* 2009; Freedman *et al.* 2010; Mitchell *et al.* 2015; Nadis 2016). Compared

97    to the central rainforest, ecotone habitats differ dramatically in numerous ways. For

98    example, ecotones have less tree cover, lower levels of precipitation, and greater intra-

99    annual variation in environmental variables. These ecological differences may lead to

100   distinctive food resources, pathogens, acoustic environments and predation levels

101   (Slabbekoorn & Smith 2002; Smith *et al.* 2005, 2013). Consequently, these differences in

102   both abiotic and biotic environments are hypothesized to result in divergent selection in

103   ecotone and rainforest populations, leading to locally adapted populations (Smith *et al.*

104   1997, 2005; Kirschel *et al.* 2009; Freedman *et al.* 2010; Sehgal *et al.* 2011; Kirschel *et al.*

105   2011). This hypothesis is supported by the fact that parapatric *A. virens* populations across

106   rainforest-ecotone gradients have undergone significant divergence in morphological (i.e.

107   body mass, wing, tail, tarsus and beak length) and vocal characteristics despite significant

108   levels of gene flow (Smith *et al.* 1997; Slabbekoorn & Smith 2002; Smith *et al.* 2005;

109   Kirschel *et al.* 2011; Smith *et al.* 2013). This pattern of divergence with gene flow and the

110   role of ecotones in driving adaptive divergence is further supported by the fact that

111   allopatric rainforest populations of *A. virens* that were geographically isolated in West and

112   Central Africa for two million years, had much lower levels of phenotypic divergence in

113   these traits compared to the level of divergence observed across a narrow (often 100km)

114   rainforest-ecotone gradient (Smith *et al.* 2005). Together, results for *A. virens* and those

115   from other species suggest that strong divergent natural selection across the rainforest-

116   savanna ecotone transition contributes to adaptive phenotypic divergence despite high

117   levels of ongoing gene-flow (Smith *et al.* 1997, 2001, 2005). Evidence for divergence with

118   gene flow in *A. virens* is also consistent with models of ecological speciation where

119    natural selection caused by shifts in ecology or invasions of new habitats can result in

120    divergence in fitness related traits and might play a prominent role in speciation (Orr &

121    Smith 1998; Schneider *et al.* 1999; Ogden & Thorpe 2002; Rundle & Nosil 2005; Schluter

122    2009; Beheregaray *et al.* 2015). Opportunities for this kind of divergence are possible

123    across the little greenbul range, as they occur across a wide diversity of habitats, including

124    mountains and islands, which are also known as hotspots of diversification and speciation

125    (Darwin 1859; Myers *et al.* 2000; Orme *et al.* 2005). Previous research has found that,

126    compared to *A. virens* populations in mainland rainforests, mountain populations and

127    island populations also show significant divergence in morphological traits typically

128    related to fitness in birds, including body mass, wing length, tail length, tarsus length and

129    bill size (Smith *et al.* 2005). Moreover, both habitats have considerable gene flow with

130    mainland rainforest populations in Lower Guinea (Smith *et al.* 2005), suggesting natural

131    selection may play an important role in divergence of mountain and island populations in

132    *A. virens*.

133          To date, the paucity of high-resolution genomic data for rainforest species such as

134    *A. virens* hinders a full exploration of the evolutionary mechanisms that may be important

135    for diversification. Previous genetic studies on *A. virens* population structure utilized a

136    handful of mtDNA markers (Smith *et al.* 2001) and microsatellite loci (Smith *et al.* 2005).

137    These limited resources were unable to differentiate ecotone and forest populations at

138    genetic level, therefore debates still remain whether the observed phenotypic divergence

139    might be the results of plasticity in traits in response to varying environmental conditions,

140    or strictly genomic divergence between populations in ecotone and rainforest. Recent

141    development of next generation sequencing techniques (NGS), especially restriction site

142    associated DNA (RAD) sequencing, enables one to *de novo* assemble hundreds of

143    thousands of RAD loci across the genome in hundreds of samples without a reference

144    genome. This cost effective method to produce genomic-wide population data provides

145    unprecedented opportunities to assess the patterns of diversity with much greater

146    resolution, to find potential population structure and to identify candidate loci under local

147    selection in non-model species such as *A. virens*.

148          Here we take a population genomic approach leveraging single nucleotide

149    polymorphism (SNP) data generated from RAD sequencing to survey the genome-wide

150 diversity of *A. virens* across multiple ecological habitats in Cameroon and Equatorial

151 Guinea, including rainforests, ectones, mountains, as well as island. Our specific objectives

152 for this approach were to: 1) estimate overall levels of genetic diversity in *A. virens*; 2)

153 determine population structure and differentiation across habitats; 3) identify candidate loci

154 that are potential targets of selection; 4) understand the biological functions of these

155 candidate loci using transcriptome data; and 5) characterize genetic turnover across

156 environmental gradients.

157

158 **Materials and methods**

159 *Sampling, DNA extraction and RADseq library construction*

160     For RAD sequencing, blood samples from adult *A. virens* were collected in Central

161 Africa and stored in Queens Lysis Buffer (Smith *et al.* 1997, 2005). Overall, 217 samples

162 were collected from 15 geographically distant sampling sites (Figure 1A), representing 15

163 populations. Sampling sites were classified into one of four habitat types by researchers in

164 the field and had previously been confirmed using remote sensing data (Slabbekoorn &

165 Smith 2002; Smith *et al.* 2005, 2013). Low quality samples were removed by filtering,

166 resulting in a total of 182 samples included in the final analysis (see *RADseq data*

167 *bioinformatics processing* below). This included seven rainforest populations (83 samples),

168 five ecotone populations (59 samples), two mountain populations (18 samples) and one

169 population from the island of Bioko (12 samples). Each population was represented by 5-

170 22 samples, with a mean representation of 12 samples (Table S1) (Willing *et al.* 2012;

171 Nazareno *et al.* 2017).

172     RAD library preparation followed the methods for traditional RAD as described in

173 Ali et al (2016) that were slightly modified from the original RAD protocol as described in

174 Baird et al (2008). In short, genomic DNA (50 ng) for each sample was digested with 2.4

175 units of SbfI-HF (New England Biolabs, NEB, R3642L) at 37 °C for 1 hr in a 12 μl

176 reaction volume buffered with 1X NEBuffer 4 (NEB, B7004S). Samples were heated to 65

177 °C for 20 minutes, and then 2 μl indexed SbfI P1 RAD adapter (10 nM) was added to each

178 sample. Ligation of inline barcoded P1 adaptors was performed by combining 2 μl of the

179 ligation mix (1.28 μl water, 0.4 μl NEBuffer 4, 0.16 μl rATP (100 mM, Fermentas R0441)

180 with 0.16 μl T4 DNA Ligase (NEB, M0202M) and heating at 20 °C for 1 hr followed by

181    incubation at 65 °C for 20 min. Following the ligation, half the per sample volume or 5 μl

182    of each of the 96 samples were pooled into a single tube and cleaned using 1X Agencourt

183    AMPure XP beads (Beckman Coulter, A63881); the remainder of the sample was stored

184    for use in an additional library preparation if needed. The pooled DNA was then re-

185    suspended in 100 μl low TE and sheared to an average fragment size of 500 base pairs

186    using a Bioruptor NGS sonicator (Diagenode). Sheared DNA was then concentrated to

187    55.5 μl using Ampure XP beads, and used as the template in the NEBNext Ultra DNA

188    Library Prep Kit for Illumina (NEB E7370L; version 1.2). The NEBNext protocol for

189    library prep was followed apart from the fact that we used custom P2 adaptors which were

190    created by annealing an NEBNext Multiplex Oligo for Illumina (NEB, E7335L) to the

191    oligo GATCGGAAGAGCACACGTCTGAACTCCAGTCACIIIIIIATCAGAACA*A (the

192    * represents a Phosphorothioated DNA base). In addition, instead of the USER enzyme

193    step, we used a universal P1 RAD primer

194    (AATGATACGGCGACCACCGAGATCTACACTCTTTCCCTACACGAC*G) and a

195    universal P2 RAD primer (CAAGCAGAAGACGGCATACG*A) during final

196    amplification. The final RAD library was cleaned using AMPure XP beads and sequenced

197    at the UC Berkeley QB3 Vincent J Coates Genome Sequencing Laboratory (GSL) on an

198    Illumina HiSeq2000 (Illumina, San Diego, CA) using single-end 100 bp sequence reads.

199

200    *RADseq data bioinformatics processing*

201         We used FastQC (http://www.bioinformatics.babraham.ac.uk/projects/fastqc/) to

202    assess overall data quality of each RADseq sequencing run. To remove the lowest quality

203    bases, we trimmed all raw sequencing reads (100bp) by 14 bp at the 3' end using *seqtk*

204    (https://github.com/lh3/seqtk). We processed RADseq reads using the *Stacks* pipeline

205    version 1.32 (Catchen *et al.* 2011, 2013) in the following manner. First, we demultiplexed

206    the trimmed data by P1 barcodes and removed low-quality reads and those containing

207    adapter sequences using *process_radtags.* After demultiplexing, reads were 80 bp in length

208    (without barcodes) and data from different runs were combined together. These reads were

209    used to *de novo* assemble RAD loci using *denovo_map.pl* (parameter settings: m=3 M=4

210    n=4). The parameters for *de novo* assembly were determined empirically to limit the

211    impact of over-splitting of loci following methods described in Ilut *et al.* (2014) and

212   Harvey *et al.* (2015). Specifically, we chose one sample that had a depth of coverage close

213   to the median depth coverage of all samples, and ran the *de novo* assembly over a wide

214   range of values of M (1-8; n=M) using *ustacks*. The percentage of homozygous and

215   heterozygous loci plateaued at M=4, suggesting this value appropriately minimized over-

216   splitting of alleles (Figure S1). Thus, we used M=4 for the final run on all samples. *Stacks*

217   implements a multinomial-based likelihood model for SNP calling, by estimating the

218   likelihood of two most frequently observed genotypes at each site and performing a

219   standard likelihood ratio test using a chi-square distribution (Hohenlohe *et al.* 2010;

220   Catchen *et al.* 2011). We used the default alpha (chi-square significance level) of 0.05.

221   Paralogous loci that were stacked together were identified and removed by later quality

222   control steps built into *Stacks* (e.g. max number of stacks per loci = 3; Ilut *et al.* 2014;

223   Harvey *et al.* 2015). After the first round of assembly using *denovo_map.pl*, we ran stacks'

224   correction mode (*rxstacks-cstacks-sstacks*) using the bounded SNP model with a 0.05

225   upper bound for the error rate (bound_high = 0.05). The *rxstacks* program made

226   corrections to genotype and haplotype calls based on population information, rebuilt the

227   catalog loci, and filtered out loci with average log likelihood less than -8.0

228   (http://catchenlab.life.illinois.edu/stacks/).

229          We then identified a set of high-quality RAD loci for downstream population

230   genetic analysis using the following steps: 1) We only kept RAD loci if they were present

231   in at least 80% of all samples, because loci that only assembled in small subset of samples

232   had limited utility in downstream analyses as well as possibly higher error rates. 2) We

233   filtered out RAD loci that had more than 40 SNPs along the 80 bp RAD loci sequence, as

234   these likely represented sequencing errors or over-clustering of paralogous loci. In the final

235   dataset, the RAD locus that has the most SNPs possessed 25 SNPs. Because the alignments

236   look reasonable for the RAD loci that have higher number of SNPs, we did not apply any

237   additional filters to avoid introducing additional biases. 3) We mapped the RAD loci

238   sequences from *A. virens* to the closest reference genome available, the zebra finch genome

239   (version 3.24), using BLAT and removed RAD loci that mapped to multiple positions in

240   the zebra finch genome. 4) We used BLAT to compare RAD loci sequences against each

241   other, and removed ones that had a match. This step further removes over-splitting RAD

242   loci.

243    Following these filters, we obtained our final consensus set of RAD loci (Table S2).

244    Samples that were missing more than 20% of the final consensus RAD loci were identified

245    in a preliminary run and were removed from final analysis because they likely had low

246    quality DNA, low quality libraries, or low sequencing coverage. 182 samples were

247    included in the final dataset (see above). Genotypes were called and filtered using methods

248    implemented in the *Stacks* pipeline (Hohenlohe *et al.* 2010). We exported genotypes for

249    the final consensus RAD loci in VCF format using stacks *populations* program (only bi-

250    allelic SNPs). Additional filters based on genotype calls were performed in *vcftools*

251    (https://vcftools.github.io/index.html) or using custom scripts, which includes removing

252    SNPs from the last 7 bp of the RAD loci as this part of the locus was enriched for

253    erroneous SNPs due to the lower sequencing quality at the 3' end of reads, and filtering

254    sites that have genotyping rate less than 80% of all samples.

255    We used the resulting full SNP dataset with SNPs from all frequencies to estimate

256    genetic diversity statistics such as number of segregating sites (S), average pairwise

257    differences ($\pi$) and Waterson's $\theta$ ($\theta_w$) in each population (Table S1). Rare SNPs that had a

258    minor allele frequency (MAF) less than 2% in the whole sample set were subsequently

259    removed using *vcftools*, and the remaining SNPs were used for downstream analyses such

260    as PCA, pairwise $F_{ST}$ calculations, *Bayescan* outlier analysis, *gradientForest* analysis.

261

262    *RNA extraction, RNAseq library preparation,* and *transcriptome de novo assembly*

263    *A. virens* lacks a reference genome. In order to help determine which of the RAD

264    loci are transcribed, we collected RNAseq data and made a *de novo* assembly of *A. virens*

265    transcriptome. Fresh tissue samples were collected from 10 live individuals in the field

266    (five tissue types: blood, brain, breast tissue, heart and liver). Tissue samples were stored

267    in either PAXgene (Blood RNA Tubes, PreAnalytiX/Qiagen, Switzerland) or Allprotect

268    (Tissue Reagent, Qiagen, Germany) buffer and shipped to laboratory facilities at UCLA.

269    RNA was extracted from each sample and tissue type separately using an RNeasy kit

270    (Qiagen, Germany), and based on quality of extractions (both overall concentration and

271    260/280 ratio), three RNA samples from three tissue types (brain, heart and liver) were

272    chosen to perform library preparations. RNAseq libraries were prepared using Illumina

273    TruSeq RNA Library Prep Kit v.2 (Illumina, San Diego, California) following the

274    manufacture's protocol, and libraries were indexed, normalized, pooled, and sequenced on

275    a single lane on Illumina HiSeq 2500 (paired end 100 bp reads, Rapid Run mode) at GSL.

276        We obtained one lane of paired-end RNAseq data pooled from three tissue types.

277    We first removed bases with quality scores lower than 20 and minimum sequence length of

278    30bp using *trim_galore* (http://www.bioinformatics.babraham.ac.uk/projects/trim_galore/).

279    We then pooled the remaining paired-end reads from different tissues together for a *de*

280    *novo* assembly of the transcriptome using the *Trinity* pipeline (Grabherr *et al.* 2011). We

281    assessed the quality of the assembly using scripts provided in the *Trinity* package, and

282    predicted the coding regions in the assembled transcriptome using *TransDecoder* in *Trinity*.

283

284    *Detecting population structure using genomic data*

285        To detect the underlying population structure among samples, we performed a

286    Principal Component Analysis (PCA) using the bioconductor package *SNPrelate* (Zheng *et*

287    *al.* 2012). 47,482 SNPs with MAF >=2% were used in PCA. The first six principal

288    components were visually examined to identify clustering patterns of samples and to

289    determine whether these genetic clusters tend to segregate with ecological factors or

290    geography. We also used the program ADMIXTURE (Alexander *et al.* 2009) to estimate

291    the ancestry of individual genotypes, using only the first SNP of each RAD loci to limit the

292    impact of linkage disequilibrium. The analysis was run for K = 1-15.

293        To quantify pairwise population differentiation, we calculated pairwise $F_{ST}$

294    between populations using *SNPrelate*. The correlation of population genetic differentiation

295    (pairwise $F_{ST}$) and geographic distance, in other words, the presence of isolation by

296    distance (IBD), was estimated by a simple Mantel test with 999,999 permutations using

297    vegan2.2-1 in R (Mantel 1967; Oksanen *et al.* 2017). Mantel tests are reported using both

298    raw $F_{ST}$ and $F_{ST}/(1- F_{ST})$, as well as both raw Euclidian geographic distance and log-

299    transformed distances (Slatkin 1995; Rousset 1997).

300        Moreover, we found that pairwise $F_{ST}$ between populations from different habitats

301    were higher than pairwise $F_{ST}$ computed between populations from the same habitat (see

302    Results). In principle, this pattern could solely be driven by isolation by distance as

303    populations from the same habitat tend to be located closer geographically to each other

304    compared to population from different habitats. To determine whether the elevated $F_{ST}$

305  between populations from different habitats (compared to $F_{ST}$ between populations from

306  the same habitat) could be explained simply by the differences in geographic distance, we

307  performed permutation tests that accounted for the fact that populations from different

308  habitats tend to be further apart. Specifically, we divided the population pairs into five bins

309  based upon their geographic distances from each other (i.e. <200km, 200-400km, 400-

310  600km, 600-800km, >800km). Then within each bin, we permutated whether the pairwise

311  $F_{ST}$ values are from a within-habitat comparison or a between-habitat comparison. We

312  generated 10,000 such permutations and, for each permutation, performed a t-test on

313  whether the $F_{ST}$ values for between-habitat comparisons were higher than those for within-

314  habitat comparisons. From the permuted data, we built a null distribution of t-statics,

315  which accounted for the effect of geography. Our final empirical p-value for the observed

316  data was calculated as the percentage of permutated datasets that had a t-statistic as large

317  or larger than the one seen in the original data. Similar permutation analyses were applied

318  to the dataset including all habitats as well as to a dataset that only considered rainforest

319  and ecotone populations. In the null distribution of t-statistics for the test of whether $F_{ST}$ is

320  higher between than within habitats, we found that none of the 10,000 permutated datasets

321  had a t-statistic of $F_{ST}$ as large or larger than the one seen in the original data, suggesting a

322  p-value < 1e-04. However, for the null distribution of t-statistics for the test of whether

323  distance is higher between or within habitats, 1581 of the 10,000 permutated datasets had a

324  t-statistic as large or larger than the one seen in the original data, suggesting a p-value =

325  0.158. This suggests that our null distribution of t-statistics accounts for the fact that

326  populations from similar habitats tend to be closer together geographically.

327  As an alternative method to test whether habitat contributed to the observed pattern

328  of population differentiation above and beyond geographic distance, we created a binary

329  matrix that indicated whether a pair of populations was from the same habitat or not. We

330  tested the correlation of genetic distance matrix and this matrix while controlling for

331  geographic distance using a partial Mantel test using vegan 2.2-1 in R (Mantel 1967;

332  Oksanen *et al.* 2017). Partial Mantel tests are performed using both raw $F_{ST}$ and $F_{ST}/(1-$

333  $F_{ST})$, as well as both raw Euclidian geographic distance and log-transformed distances

334  (Slatkin 1995; Rousset 1997).

335

336    *Identifying outlier SNPs under selection*

337    We used *BayeScan2.1* (Foll & Gaggiotti 2008) to identify highly differentiated

338    SNPs that are candidates to be under natural selection. This program takes a Bayesian

339    approach to search for SNPs exhibiting extreme $F_{ST}$ values that could be due to local

340    adaptation. Outlier SNPs were identified using SNPs with MAF >= 2% across all samples,

341    specifying all 15 populations or four habitats (see supplemental Notes). We ran 5000

342    iterations using prior odds of 10 and assessed the statistical significance of a locus being an

343    outlier using a false discovery rate (FDR) of 5%.

344    To explore the spatial patterns of population differentiation across chromosomes,

345    we mapped the consensus RAD loci to the zebra finch genome using BLAT with default

346    parameters. For the uniquely mapped RAD loci, we plotted the $F_{ST}$ of each SNP by

347    genome coordinates to examine spatial patterns of outlier SNPs. To interpret the potential

348    biological function of the outlier SNPs identified by Bayescan analysis, we used a zebra

349    finch genome annotation (v3.2.4.78) to identify outlier SNPs mapped to annotated genic

350    regions.

351    We further examined whether candidate loci under selection were enriched in

352    exonic (transcribed) or coding regions. To do this, we mapped RAD loci to the *de novo*

353    assembled *A. virens* transcriptome using BLAT with default parameters. Any RAD locus

354    that mapped to the transcriptome was considered from exonic regions of the genome, and

355    the remaining RAD loci were labeled "non-transcribed" regions of the genome. Similarly,

356    we mapped RAD loci to predicted coding sequences and categorized them into coding and

357    non-coding RAD loci. We then used a one-sided Fisher's exact test to examine whether

358    there was significant enrichment of outlier loci in exon or coding regions of the genome.

359    Finally, we cross-checked these outliers to see if there was any significant associations

360    with environment using Latent Factor Mixed Models (Frichot *et al.* 2013) (see

361    Supplemental Notes for more details).

362

363    *Detecting environmental drivers of genomic variation*

364    In addition to population structure, we also tested whether allele frequencies in

365    different populations were associated with environmental variables across the range of *A.*

366    *virens* using the package *gradientForest* (Ellis *et al.* 2012) in the R statistical framework

367 (R Working Group, 2014). Gradient forests are an extension of random forests (Breiman

368 2001) that treat response variables (in this case, individual SNP minor allele frequencies

369 within each population) as members of a larger community (the total genome), and

370 provides summary statistics based on ensembled forest runs to indicate an overall

371 association of changes in allele frequency to particular environmental variables (Ellis *et al.*

372 2012; Fitzpatrick & Keller 2015). Gradient forests were run using the following changes to

373 the default settings: number of trees run for each environmental variable = 500, number of

374 SNPs included in each bin = 1000. Allelic frequencies across the genome were predicted

375 for unsampled geographic locations by generating a random set of 100,000 points across

376 the range of *A. virens*. Then we used our final gradient forest model to predict allele

377 frequencies at each of those points, given their environmental characteristics. Ordinary

378 Kriging (Oliver & Webster 1990) was then used within ArcMap (ESRI, Redlands, CA) to

379 generate a continuous surface across the known range of *A. virens* in Cameroon.

380　　　　We used a suite of 17 environment variables (Table S6), including bioclim

381 measures of temperature and precipitation (n=9; any variables showing a Pearson's

382 correlation coefficient > 0.7 were removed) downloaded from the Worldclim database

383 (www.worldclim.org), measures of vegetation and tree cover captured using the NASA

384 MODerate-resolution Imaging Spectroradiometer (MODIS, n=4), elevation and slope

385 captured via the Shuttle Radar Topography Mission (n=2), and surface moisture estimates

386 measures using the Quick Scatterometer (QuikSCAT, n=2). In addition to these variables,

387 and to account for purely geographic associations, we also included Euclidean distances

388 (measured as Latitude and Longitude) as predictor variables in all model.

389

390 **Results**

391 *SNP discovery and overall genetic diversity*

392　　　　We used RAD sequencing to survey the genome-wide diversity of *A. virens*. The

393 final sample set included 15 populations from four different habitats, including rainforests,

394 ecotones, mountains, and an island (Figure 1A; Table S1). After removing low quality

395 reads and samples, we obtained a total of 916 million reads for 182 *A. virens* samples

396 (SRA:xxxxxx). The number of raw sequence reads per sample ranged from 1.60 to 20.73

397 million. On average, 99.2% of these reads were utilized in the *de novo* assembly of the

398   RAD loci. The mean coverage depth ranged from 16× to 136× per sample (mean = 38×,

399   median = 32×, Figure S2). Using this dataset, we assembled and identified 34,657 high

400   quality RAD loci that passed our quality filters and were genotyped in more than 80% of

401   all final samples. On these 34,657 consensus RAD loci, there were a total of 255,290

402   SNPs. The median number of SNPs per RAD locus is seven. With a minimum minor allele

403   frequency filter of 2%, we retained 47,482 SNPs that were present on 23,882 RAD tags

404   (Table S2; Supplemental Notes).

405       The number of segregating sites ranges from 25,936 to 70,598 per population.

406   Waterson's $\theta$ ($\theta_w$) was estimated to be 0.0049 - 0.0076/bp (mean = 0.0064/bp), and $\pi$

407   ranges from 0.0034 - 0.0037/bp (mean = 0.0036/bp) (Table S1), which is comparable to $\pi$

408   estimated from other bird species (Nadachowska-Brzyska *et al.* 2013; Romiguier *et al.*

409   2014). Overall levels of genetic diversity are comparable in each habitat, including the

410   island population (Table S1). The finding that $\theta_w$ is larger than $\pi$ indicates an excess of

411   low-frequency variants relative to the standard neutral model which could be driven by

412   recent population expansions.

413   *Transcriptome assembly and annotation*

414       Transcriptome assembly was performed using 169 million paired-end RNAseq data

415   from three different tissue types. The assembled transcriptome had a GC content of 45%.

416   The average contig length was 815 bp and N50 was 1,619 bp. In total, trinity produced

417   237,226 genes and 286,494 transcripts, and predicted 81,018 coding sequences from these

418   transcripts (Table S5). Of the 34,657 RAD loci we genotyped, 8412 RAD loci (24.2%)

419   were mapped to the *de novo* assembled *A. virens* transcriptome, and 3,618 RAD loci

420   (10.4%) were mapped to the predicted coding sequences (Figure S3). The RAD tags

421   overlapping coding sequence tend to have fewer SNPs than those that do not overlap with

422   coding sequences (Figure S4), consistent with the fact that the coding regions are likely

423   under stronger selective constraint.

424

425   *Population structure*

426       We used PCA to identify population structure in little greenbuls. The first two PCs

427   revealed a clear clustering pattern of individuals from the same habitats (Figure 1B).

428   Populations from island, mountains, rainforests, and ecotones formed four discrete clusters,

429   suggesting genomic divergence across ecological gradients and habitats. Island and

430   mountain populations were most distinct (Figure 1B), however samples from all four

431   habitats separated on PC1, including those from ecotone and rainforest habitats. PC2

432   further separated island and mountain samples from rainforest and ecotone samples.

433   Remarkably, results suggest that, within rainforest and ecotone habitats, individual

434   populations could be distinguished solely on the basis of genomic markers, mostly by PC1,

435   with individuals from the same sampling sites clustering together (Figure 1C). The level of

436   separation of ecotone populations from rainforest populations along PC1 roughly followed

437   a latitudinal gradient, corresponding to environmental and rainfall gradients that

438   distinguish ecotone in the north from rainforest in the south of Cameroon (see

439   environmental analyses below). Specifically, samples collected from sites Wakwa and

440   Ngaoundaba Ranch, toward the more extreme edge of the ecotone habitat and had the most

441   extreme ecotone environmental conditions, formed clusters that were most distant from the

442   rainforest samples, while samples collected from Betare Oya, at lower ecotone that was

443   closest to the central rainforests, formed a cluster closest to the rainforest (Figure 1C). The

444   pattern of genomic differentiation across habitats was confirmed using the program

445   *Admixture* (Figure S5).

446         Pairwise $F_{ST}$ between populations ranged from 0.017 to 0.078 (mean = 0.038;

447   Figure 2A-B; Table S3), indicating low overall levels of genomic differentiation across

448   populations. There was significant correlation between pairwise $F_{ST}$ and geographic

449   distance between the populations (Mantel r = 0.34; mantel simulated p-value = 0.003),

450   suggesting isolation by distance contributes to population differentiation. However,

451   pairwise $F_{ST}$ between populations from different habitats was significantly higher than

452   pairwise $F_{ST}$ between populations within the same habitat (one tailed t test, p-value =

453   1.36e-10; Figure 2A). Pairwise geographic distances between populations from different

454   habitats were also significantly higher than pairwise distances between populations within

455   the same habitat (one tailed t-test, p-value = 1.015e-06). To account for the fact that

456   populations from different habitats were also geographically further apart, we performed a

457   permutation test, where we randomized whether a population pair was from the same or

458   different habitats in different bins stratified by their geographic distance. Using permutated

459   datasets, we built a null distribution of these t-statistics (that already includes the effect of

460  geographic distance), which we used to evaluate the significance of our observed value.

461  The higher $F_{ST}$ value for between-habitat comparison was highly significant when

462  compared to this improved null distribution (p-value < 1e-04, Figure 2C and Figure S6),

463  suggesting that isolation by distance alone cannot explain the higher $F_{ST}$ between habitats

464  than within habitats. Similarly, only considering rainforest and ecotone populations,

465  pairwise $F_{ST}$ was significantly higher between habitats as compared to within habitats (one

466  tailed t-test, p-value = 9.793e-05). Application of the same permutation test shows the

467  higher $F_{ST}$ between ecotone and rainforest populations (p-value=0.0055) cannot be

468  explained by geographic distance alone (Figure 2D and Figure S6).

469        To confirm this finding using an alternative statistical approach, we used partial

470  Mantel tests to determine the contribution of habitat types of population pairs to their

471  genetic differentiation ($F_{ST}$), controlling for geographic distance. We found a highly

472  significant and positive correlation between genetic distance and whether a pair of

473  population comes from the same habitat, and greater genetic differentiation (higher $F_{ST}$)

474  from between-habitat populations compared to within-habitat populations, while

475  controlling for geographic distance (Table 1). Taken together these results suggest that

476  factors other than geographic location, such as local adaptation significantly contribute to

477  population differentiation between habitats.

478        In addition, mountain and island populations were more diverged from other

479  populations (Figure 2B). Interestingly, $F_{ST}$ between two mountain populations were

480  exceptionally high ($F_{ST}$= 0.060) compared to other within habitat pairwise $F_{ST}$ (ranging

481  from 0.017 to 0.040, Figure 2A), despite the fact that the two mountain populations were

482  geographically very close to each other. The $F_{ST}$ values between mountain populations and

483  lowland forest/ecotone populations were larger than pairwise $F_{ST}$ values between lowland

484  populations, suggesting mountain populations are highly differentiated both from one

485  another and from lowland populations.

486

487  *Candidate loci under selection*

488        To further explore potential candidate loci under selection, we identified SNPs with

489  extreme allele frequency differences across populations, which should be enriched by

490  targets of local adaptation. We identified 140 outlier SNPs across all populations with a

491   False Discovery Rate of 5% using *Bayescan*. These candidate SNPs are potential targets of

492   divergent selection across different sampling sites (Figure S7). The 140 outlier SNPs reside

493   in 119 loci, and 40 of these loci mapped to the zebra finch genome (Figure S9). Of these,

494   36 mapped to main scaffolds of known chromosomes and four mapped to the Z

495   chromosome. Only 13 of these outlier loci mapped to annotated genic regions on the zebra

496   finch genome and nine mapped to genes with functional annotations (Table S4).

497       In order to uncover the functional significance of outlier loci, we used the *de novo*

498   assembly of greenbul transcriptome to partition the RAD loci and SNPs into different

499   categories depending on whether they mapped to coding regions or transcribed (exonic)

500   regions (Figure S10). This enabled us to test for enrichment of outlier SNPs in putatively

501   functional regions. Of the 47,482 SNPs, 9677 mapped to the transcriptome, and 42 were

502   outliers based on a Bayescan analysis. Using a one-sided Fisher's exact test, we detected

503   significant enrichment of outlier loci in exonic regions of the genome (p=0.0044; Table S2;

504   Figure S10). Using the predicted coding sequence from the transcriptome, 3,602 SNPs

505   mapped to the predicted coding sequences and 21 of these were outliers. We again detected

506   a significant enrichment of outlier SNPs in protein-coding sequences (p=0.002; Table S2;

507   Figure S10). Taken together, these enrichment results provide additional confidence that

508   the outlier loci found using Bayescan captured functionally important, biologically relevant

509   genetic variants, which were not merely loci that fell within the tail of a neutral

510   distribution.

511

512   *Genomic Turnover Across Environments*

513       Because some environmental adaptation may involve shifts in allele frequency at

514   many loci across the genome (e.g., polygenic selection involving many genes of small

515   effect), we used the *gradientForest* approach to look for correlations in allele frequencies

516   associated with environmental variables. A total of 7238 SNPs, ~15% of all SNPs, had $R^2$

517   values above 0 (0.0073-0.83) when testing for a correlation between frequency and an

518   environmental variable. Of the 19 environmental and geographical variables included in

519   models (Table S6), variables capturing temperature variation (Min Temp: minimum

520   temperature of the coldest month, Temp Range: mean diurnal temperature range, and Mean

521   temp: mean annual temperature) and elevation were most important in explaining

522   environmently associated variation in SNPs (Figure 3A). In some cases, measures of

523   surface moisture or tree cover were also important, but axes for these variables largely

524   overlapped with temperature or elevation measures along PC plot (likely the result of co-

525   linearity in environmental variables) (Figure 3B). Results from LFMM analyses indicated

526   these same variables were also associated with differentiation observed at hundreds of

527   individual loci, although exact functions of these regions remain unknown (see

528   Supplemental Notes).

529         Geographic variables alone were not as important in explaining variation in allele

530   frequency, again suggesting that geographic distance cannot fully account for all variation

531   in SNP frequencies across the range of little greenbuls. Above and beyond neutral

532   processes, selective pressures imposed by differences in these environments best explains

533   the observed genomic patterns of variation. Predictions across Cameroon suggest strong

534   genomic turnover (defined as coordinated shifts in allele frequencies across the genome)

535   throughout the forest, savannah, and ecotone regions, with diagnostic genomic variation

536   occurring in each of these habitats (Figure 3). Distinct SNP frequencies at high elevations

537   (Figure 3B-C), and the fact that elevation explains a large proportion of variation of allele

538   frequencies in the greenbul genome (largely allied with PC1, Figure 3B) also suggest a

539   unique genetic signature in populations at elevation.

540

541   **Discussion**

542         In this study, we used genome-wide RADseq SNPs to characterize the overall level

543   of genetic diversity in *A. virens* populations across four different habitats. We found

544   evidence of population structure of *A. virens* consistent with habitat type and previously

545   observed phenotypic divergence. We demonstrated that population differentiation across

546   habitats cannot be explained solely by isolation-by-distance, suggesting local adaptation

547   further contributes to genomic divergence among habitats. We identified 140 outlier SNPs

548   that are potential targets of selection and the fact that they are significantly enriched in

549   exonic and coding regions suggests they are functionally important. Environmental

550   association analysis further supports this conclusion and shows environmental variables,

551   including temperature and elevation, are highly associated with patterns of genomic

552   variation across the range of the little greenbul.

553        In addition to the differences between rainforest and ecotone populations, other

554    habitats were found to harbor distinct patterns of genetic variation. The population from

555    Bioko island formed a distinctive genetic cluster based on PCA, and also was identified as

556    distinct in environmental association models (Figure 1B and 3), consistent with previous

557    studies (Smith *et al.* 2005). Bioko island is 32 km off the coast of Africa, separated from

558    mainland ~10,000 years ago and has an area of 2,017 km$^2$. Island populations and species

559    may have smaller effective population sizes than mainland populations or sister taxa

560    (Robinson *et al.* 2016), due to possible population bottleneck and considerably smaller

561    ranges. As a result, island populations may have lower genetic diversity compared to their

562    mainland counterparts (Frankham 1997). However, several recent empirical studies

563    suggested this may not always be the case, particularly in birds (Francisco *et al.* 2015;

564    James *et al.* 2016). In our study, the estimates of genetic variation using genome wide SNP

565    markers in greenbul population on Bioko island are comparable to those from mainland

566    populations (Table S1). This is consistent with the recent findings that island populations

567    do not always have lower genetic diversity (Francisco *et al.* 2015; James *et al.* 2016), and

568    the fact that Bioko island is a large island that only recently separated from the mainland.

569        Tropical mountains are well known to support disproportionally high biodiversity

570    and are thought to be hotspots for avian speciation (Roy 1997; Myers *et al.* 2000; Smith *et*

571    *al.* 2000; Orme *et al.* 2005; Fjeldså *et al.* 2007, 2012; Drovetski *et al.* 2013). Little

572    greenbuls are found at elevations up to 2400 m, where environmental variables,

573    particularly temperature and vegetation, change rapidly along altitudinal gradients. Our

574    two mountain populations have high $F_{ST}$ despite being geographically close and from same

575    habitat type. Although the Euclidean distance between these two mountain populations is

576    short, the environmental changes along altitudinal gradients are steep, causing isolation

577    between populations from different mountains and forming "sky islands", between which

578    the level of gene flow probably is much lower than among lowland populations. Moreover,

579    we found that the two different mountain populations exhibited the lowest within-

580    population genetic variation among all sampled populations (Table S1). While this

581    difference was not statistically significant (likely due to small sample sizes), this decreased

582    level variation can inflate $F_{ST,}$ the relative measurement of population differentiation. It

583    also suggests that mountain populations may have overall smaller effective population

584    sizes (consistent with presumably smaller suitable habitat size for mountain populations)

585    and/or have experienced serial bottleneck/founder effects as range expansions occurred.

586    These processes can further contributed to divergence due to stronger genetic drift within

587    each subpopulation leading to faster changes in allele frequency. The idea that elevation

588    can drive genomic changes is supported by previous estimates of morphological

589    divergence (Smith *et al.* 2005), and emphasizes the importance of preserving elevational

590    gradients in tropical ecosystems in general (Thomassen *et al.* 2011).

591            Most of the genes containing outlier SNPs only have annotations predicted from

592    human homologs, except two that have annotations from bird species. Both of these two

593    genes are of particular interest. One outlier locus mapped to the 5UTR/coding junction of

594    *EDIL3*, a calcium-binding protein that has been found to function in avian eggshell

595    biomineralization (Marie *et al.* 2015). Avian eggshells protect the developing embryo and

596    keep the egg free from pathogens. Environmental factors such as temperature, humidity,

597    and partial oxygen pressure have been reported to affect avian eggshell structure, and

598    previous studies documented rapid evolution of eggshell structure in response to

599    colonization of novel environments in the house finch (Stein & Badyaev 2011). The

600    second outlier locus mapped to *MLXIPL* (MLX interacting protein-like), which is co-

601    activator of the carbohydrate response element binding protein that has been correlated

602    with fat deposition in caged chicken (Proszkowiec-Weglarz *et al.* 2008; Li *et al.* 2015).

603    Interestingly, seven more genes that contain outlier SNPs have annotation linked with

604    metabolic traits or diseases in humans. For example, outlier SNPs were found in *RGS6*

605    (Sibbel *et al.* 2011 p. 6), *CSAD* (Comuzzie *et al.* 2012) and *UNC13B* intron (Trégouet *et*

606    *al.* 2008), which were associated with dietary fat intake, food preference, adiposity/obesity

607    and diabetes in humans. Although metabolic traits were not measured, adult little

608    greenbuls from the rainforest have significantly smaller body mass and body size

609    compared to ecotone, mountain, and island populations (Smith *et al.* 1997, 2005), which

610    could be the result of divergent selection of these genes associated with metabolic traits.

611    Several recent studies have discussed the limitations of identifying $F_{ST}$ outlier as loci under

612    divergent selection, and suggest results should be interpreted carefully, because many other

613    factors, including demographic history, recombination rate heterogeneity, and background

614    selection may also create $F_{ST}$ outliers (Roesti *et al.* 2012; Lotterhos & Whitlock 2014;

615    Cruickshank & Hahn 2014). Current work to model the demographic history of *A. virens*

616    should help examine these various possibilities in greater detail.

617         Numerous hypotheses have been proposed for how biodiversity is generated in

618    rainforests (Mayr & O'Hara 1986; Moritz *et al.* 2000). With the rapid advances in

619    genomics and environmental modeling in the last decade, it is now possible to examine

620    these mechanisms in greater depth. Using more powerful genome-wide data, we have

621    shown, for the first time, strong patterns of population structure and genomic

622    differentiation between rainforest and ecotone habitats in *A. virens* (Figure 1). Previously,

623    no genetic differentiation was found between morphologically divergent populations in

624    rainforest and ecotone habitats, leaving open the possibility that the observed

625    morphological difference could simply be the result of a homogenized meta-population

626    that differentially responds to environmental gradients. Although identifying the

627    underlying genetic basis of morphological traits that differ between rainforest and ecotone

628    populations was beyond the scope of this study, our results complement previous work by

629    demonstrating that populations along the rainforest - ecotone gradient are diverging at the

630    genomic level, and raise the possibility that local adaptation could account for the patterns

631    of morphological variation previously observed across ecotone-rainforest gradients.

632    Results also complement past research on reproductive behavior, which found differences

633    in song characteristics along the forest-ecotone gradient, and showed experimentally that

634    singing males respond more aggressively to male songs from their own habitat, suggesting

635    incipient reproductive isolation driven by habitat (Slabbekoorn & Smith 2002; Kirschel *et*

636    *al.* 2011; Smith *et al.* 2013). These patterns of differentiation are consistent with models of

637    ecological speciation, where natural selection caused by shifts in ecology can promote

638    speciation (Orr & Smith 1998; Schneider *et al.* 1999; Schluter 2000; Ogden & Thorpe

639    2002; Rundle & Nosil 2005; Price 2008; Räsänen & Hendry 2008; Schluter 2009;

640    Beheregaray *et al.* 2015; Hanson *et al.* 2016). However, further research is necessary to

641    more fully understand the evolutionary significance of divergence across ecological

642    gradients and ecotones. In particular, studies investigating the underlying genetic basis of

643    phenotypic differentiation and mate choice experiments would provide additional insights

644    into their importance in divergence and speciation.

645    **Author Contributions**

646  T.B.S. and K.E.L conceived of and supervised the project. S.L. and R.J.H conducted the

647  laboratory work. Sequence assemblies, population structure and outlier analysis was

648  primarily carried out by Y.Z. with assistance from T.N., K.R., E.C.A. and K.E.L.

649  Environmental association analysis was performed by R.J.H. The manuscript was written

650  by Y.Z., R.J.H., K.R., K.E.L., and T.B.S., with input from all authors.

651

660

661  **References**

662  Alexander DH, Novembre J, Lange K (2009) Fast model-based estimation of ancestry in

663      unrelated individuals. *Genome Research*, **19**, 1655–1664.

664  Beheregaray LB, Cooke GM, Chao NL, Landguth EL (2015) Ecological speciation in the

665      tropics: insights from comparative genetic studies in Amazonia. *Evolutionary and*

666      *Population Genetics*, **5**, 477.

667  Breiman L (2001) Random Forests. *Machine Learning*, **45**, 5–32.

668  Catchen JM, Amores A, Hohenlohe P, Cresko W, Postlethwait JH (2011) Stacks: Building

669      and Genotyping Loci De Novo From Short-Read Sequences. *G3:*

670      *Genes|Genomes|Genetics*, **1**, 171–182.

671  Catchen J, Hohenloh PA, Bassham SL, Amores A, Cresko WA (2013) Stacks: an analysis

672      tool set for population genomics. *Molecular ecology*, **22**, 3124–3140.

673  Comuzzie AG, Cole SA, Laston SL *et al.* (2012) Novel Genetic Loci Identified for the

674      Pathophysiology of Childhood Obesity in the Hispanic Population. *PLoS ONE*, **7**.

675    Cruickshank TE, Hahn MW (2014) Reanalysis suggests that genomic islands of speciation
676          are due to reduced diversity, not reduced gene flow. *Molecular Ecology*, **23**, 3133–
677          3157.
678    Darwin C (1859) *On the Origin of Species by Means of Natural Selection, Or, The*
679          *Preservation of Favoured Races in the Struggle for Life*. J. Murray.
680    Drovetski SV, Semenov G, Drovetskaya SS *et al.* (2013) Geographic mode of speciation in
681          a mountain specialist Avian family endemic to the Palearctic. *Ecology and*
682          *Evolution*, **3**, 1518–1528.
683    Ellis N, Smith SJ, Pitcher CR (2012) Gradient forests: calculating importance gradients on
684          physical predictors. *Ecology*, **93**, 156–168.
685    Fitzpatrick MC, Keller SR (2015) Ecological genomics meets community-level modelling
686          of biodiversity: mapping the genomic landscape of current and future
687          environmental adaptation. *Ecology Letters*, **18**, 1–16.
688    Fjeldså J, Bowie RCK, Rahbek C (2012) The Role of Mountain Ranges in the
689          Diversification of Birds. *Annual Review of Ecology, Evolution, and Systematics*,
690          **43**, 249–265.
691    Fjeldså J, Johansson US, Lokugalappatti LGS, Bowie RCK (2007) Diversification of
692          African greenbuls in space and time: linking ecological and historical processes.
693          *Journal of Ornithology*, **148**, 359–367.
694    Foll M, Gaggiotti O (2008) A Genome-Scan Method to Identify Selected Loci Appropriate
695          for Both Dominant and Codominant Markers: A Bayesian Perspective. *Genetics*,
696          **180**, 977–993.
697    Francisco FO, Santiago LR, Mizusawa YM, Oldroyd BP, Arias MC (2015) Genetic
698          structure of island and mainland populations of a Neotropical bumble bee species.
699          *bioRxiv*, 27813.
700    Frankham R (1997) Heredity - Abstract of article: Do island populations have less genetic
701          variation than mainland populations? *Heredity*, **78**, 311–327.
702    Freedman AH, Thomassen HA, Buermann W, Smith TB (2010) Genomic signals of
703          diversification along ecological gradients in a tropical lizard. *Molecular Ecology*,
704          **19**, 3773–3788.

705 Frichot E, Schoville SD, Bouchard G, François O (2013) Testing for Associations between
706      Loci and Environmental Gradients Using Latent Factor Mixed Models. *Molecular*
707      *Biology and Evolution*, **30**, 1687–1699.
708 Grabherr MG, Haas BJ, Yassour M *et al.* (2011) Trinity: reconstructing a full-length
709      transcriptome without a genome from RNA-Seq data. *Nature biotechnology*, **29**,
710      644–652.
711 Haffer J (1969) Speciation in Amazonian Forest Birds. *Science*, **165**, 131–137.
712 Hanson D, Moore J-S, Taylor EB, Barrett RDH, Hendry AP (2016) Assessing reproductive
713      isolation using a contact zone between parapatric lake-stream stickleback ecotypes.
714      *Journal of Evolutionary Biology*, **29**, 2491–2501.
715 Harvey MG, Judy CD, Seeholzer GF *et al.* (2015) Similarity thresholds used in DNA
716      sequence assembly from short reads can reduce the comparability of population
717      histories across species. *PeerJ*, **3**, e895.
718 Hohenlohe PA, Bassham S, Etter PD *et al.* (2010) Population Genomics of Parallel
719      Adaptation in Threespine Stickleback using Sequenced RAD Tags. *PLoS Genet*, **6**,
720      e1000862.
721 Hoorn C, Wesselingh FP, Steege H ter *et al.* (2010) Amazonia Through Time: Andean
722      Uplift, Climate Change, Landscape Evolution, and Biodiversity. *Science*, **330**, 927–
723      931.
724 Ilut DC, Nydam ML, Hare MP *et al.* (2014) Defining Loci in Restriction-Based Reduced
725      Representation Genomic Data from Nonmodel Species: Sources of Bias and
726      Diagnostics for Optimal Clustering. *BioMed Research International, BioMed*
727      *Research International*, **2014**, **2014**, e675158.
728 James JE, Lanfear R, Eyre-Walker A (2016) Molecular Evolutionary Consequences of
729      Island Colonization. *Genome Biology and Evolution*, **8**, 1876–1888.
730 Kirschel ANG, Blumstein DT, Smith TB (2009) Character displacement of song and
731      morphology in African tinkerbirds. *Proceedings of the National Academy of*
732      *Sciences*, **106**, 8256–8261.
733 Kirschel ANG, Slabbekoorn H, Blumstein DT *et al.* (2011) Testing Alternative
734      Hypotheses for Evolutionary Diversification in an African Songbird: Rainforest
735      Refugia Versus Ecological Gradients. *Evolution*, **65**, 3162–3174.

736 Li Q, Zhao XL, Gilbert ER *et al.* (2015) Confined housing system increased abdominal

737 and subcutaneous fat deposition and gene expressions of carbohydrate response

738 element-binding protein and sterol regulatory element-binding protein 1 in chicken.

739 *Genetics and molecular research: GMR*, **14**, 1220–1228.

740 Lotterhos KE, Whitlock MC (2014) Evaluation of demographic history and neutral

741 parameterization on the performance of FST outlier tests. *Molecular Ecology*, **23**,

742 2178–2192.

743 Mantel N (1967) The Detection of Disease Clustering and a Generalized Regression

744 Approach. *Cancer Research*, **27**, 209–220.

745 Marie P, Labas V, Brionne A *et al.* (2015) Quantitative proteomics and bioinformatic

746 analysis provide new insight into protein function during avian eggshell

747 biomineralization. *Journal of Proteomics*, **113**, 178–193.

748 Martin C (1991) *The Rainforests of West Africa: Ecology, Threats, and Protection*.

749 Birkhauser, Basel ; Boston.

750 Mayr E, O'Hara RJ (1986) The biogeographic evidence supporting the Pleistocene forest

751 refuge hypothesis. *Evolution*, **40**, 55–67.

752 Mitchell MW, Locatelli S, Sesink Clee PR, Thomassen HA, Gonder MK (2015)

753 Environmental variation and rivers govern the structure of chimpanzee genetic

754 diversity in a biodiversity hotspot. *BMC Evolutionary Biology*, **15**, 1.

755 Moritz C, Patton JL, Schneider CJ, Smith TB (2000) Diversification of Rainforest Faunas:

756 An Integrated Molecular Approach. *Annual Review of Ecology and Systematics*, **31**,

757 533–563.

758 Myers N, Mittermeier RA, Mittermeier CG, da Fonseca GAB, Kent J (2000) Biodiversity

759 hotspots for conservation priorities. *Nature*, **403**, 853–858.

760 Nadachowska-Brzyska K, Burri R, Olason PI *et al.* (2013) Demographic Divergence

761 History of Pied Flycatcher and Collared Flycatcher Inferred from Whole-Genome

762 Re-sequencing Data. *PLoS Genet*, **9**, e1003942.

763 Nadis S (2016) Life on the edge: Saving the world's hotbeds of evolution | New Scientist.

764 Nazareno AG, Bemmels JB, Dick CW, Lohmann LG (2017) Minimum sample sizes for

765 population genomics: an empirical study from an Amazonian plant species.

766 *Molecular Ecology Resources*, n/a-n/a.

767    Ogden R, Thorpe RS (2002) Molecular evidence for ecological speciation in tropical
768        habitats. *Proceedings of the National Academy of Sciences*, **99**, 13612–13615.

769    Oksanen J, Blanchet FG, Friendly M *et al.* (2017) *vegan: Community Ecology Package*.

770    Oliver MA, Webster R (1990) Kriging: a method of interpolation for geographical
771        information systems. *International Journal of Geographical Information Systems*,
772        **4**, 313–332.

773    Orme CDL, Davies RG, Burgess M *et al.* (2005) Global hotspots of species richness are
774        not congruent with endemism or threat. *Nature*, **436**, 1016–1019.

775    Orr MR, Smith TB (1998) Ecology and speciation. *Trends in Ecology & Evolution*, **13**,
776        502–506.

777    Price T (2008) *Speciation in birds*. Roberts and Co., Greenwood Village, Colo.

778    Proszkowiec-Weglarz M, Humphrey BD, Richards MP (2008) Molecular cloning and
779        expression of chicken carbohydrate response element binding protein and Max-like
780        protein X gene homologues. *Molecular and Cellular Biochemistry*, **312**, 167–184.

781    Räsänen K, Hendry AP (2008) Disentangling interactions between adaptive divergence
782        and gene flow when ecology drives diversification. *Ecology Letters*, **11**, 624–636.

783    Ribas CC, Aleixo A, Nogueira ACR, Miyaki CY, Cracraft J (2011) A palaeobiogeographic
784        model for biotic diversification within Amazonia over the past three million years.
785        *Proceedings of the Royal Society of London B: Biological Sciences*, rspb20111120.

786    Robinson JA, Vecchyo DO-D, Fan Z *et al.* (2016) Genomic Flatlining in the Endangered
787        Island Fox. *Current Biology*, **26**, 1183–1189.

788    Roesti M, Salzburger W, Berner D (2012) Uninformative polymorphisms bias genome
789        scans for signatures of selection. *BMC Evolutionary Biology*, **12**, 94.

790    Romiguier J, Gayral P, Ballenghien M *et al.* (2014) Comparative population genomics in
791        animals uncovers the determinants of genetic diversity. *Nature*, **515**, 261–263.

792    Rousset F (1997) Genetic Differentiation and Estimation of Gene Flow from F-Statistics
793        Under Isolation by Distance. *Genetics*, **145**, 1219–1228.

794    Roy MS (1997) Recent diversification in African greenbuls (Pycnonotidae: Andropadus)
795        supports a montane speciation model. *Proceedings of the Royal Society B:*
796        *Biological Sciences*, **264**, 1337–1344.

797    Rundle HD, Nosil P (2005) Ecological speciation. *Ecology Letters*, **8**, 336–352.

798     Schluter D (2000) *The ecology of adaptive radiation*. Oxford University Press, Oxford.

799     Schluter D (2009) Evidence for Ecological Speciation and Its Alternative. *Science*, **323**,
800             737–741.

801     Schneider CJ, Smith TB, Larison B, Moritz C (1999) A test of alternative models of
802             diversification in tropical rainforests: Ecological gradients vs. rainforest refugia.
803             *Proceedings of the National Academy of Sciences of the United States of America*,
804             **96**, 13869–13873.

805     Sehgal RNM, Buermann W, Harrigan RJ *et al.* (2011) Spatially explicit predictions of
806             blood parasites in a widely distributed African rainforest bird. *Proceedings of the*
807             *Royal Society of London B: Biological Sciences*, **278**, 1025–1033.

808     Sibbel SP, Talbert ME, Bowden DW *et al.* (2011) RGS6 variants are associated with
809             dietary fat intake in Hispanics: the IRAS Family Study. *Obesity (Silver Spring,*
810             *Md.)*, **19**, 1433–1438.

811     Slabbekoorn H, Smith TB (2002) Habitat-Dependent Song Divergence in the Little
812             Greenbul: An Analysis of Environmental Selection Pressures on Acoustic Signals.
813             *Evolution*, **56**, 1849–1858.

814     Slatkin M (1995) A measure of population subdivision based on microsatellite allele
815             frequencies. *Genetics*, **139**, 457–462.

816     Smith TB, Calsbeek R, Wayne RK *et al.* (2005) Testing alternative mechanisms of
817             evolutionary divergence in an African rain forest passerine bird. *Journal of*
818             *Evolutionary Biology*, **18**, 257–268.

819     Smith TB, Harrigan RJ, Kirschel ANG *et al.* (2013) Predicting bird song from space.
820             *Evolutionary Applications*, **6**, 865–874.

821     Smith TB, Holder K, Girman D *et al.* (2000) Comparative avian phylogeography of
822             Cameroon and Equatorial Guinea mountains: implications for conservation.
823             *Molecular Ecology*, **9**, 1505–1516.

824     Smith BT, McCormack JE, Cuervo AM *et al.* (2014) The drivers of tropical speciation.
825             *Nature*, **515**, 406–409.

826     Smith TB, Schneider CJ, Holder K (2001) Refugial isolation versus ecological gradients.
827             *Genetica*, **112**–**113**, 383–398.

828 Smith TB, Thomassen HA, Freedman AH *et al.* (2011) Patterns of divergence in the olive

829 sunbird Cyanomitra olivacea (Aves: Nectariniidae) across the African rainforest–

830 savanna ecotone. *Biological Journal of the Linnean Society*, **103**, 821–835.

831 Smith TB, Wayne RK, Girman DJ, Bruford MW (1997) A Role for Ecotones in

832 Generating Rainforest Biodiversity. *Science*, **276**, 1855–1857.

833 Stein LR, Badyaev AV (2011) Evolution of eggshell structure during rapid range

834 expansion in a passerine bird. *Functional Ecology*, **25**, 1215–1222.

835 Thomassen HA, Fuller T, Buermann W *et al.* (2011) Mapping evolutionary process: a

836 multi-taxa approach to conservation prioritization. *Evolutionary Applications*, **4**,

837 397–413.

838 Trégouet D-A, Groop P-H, McGinn S *et al.* (2008) G/T Substitution in Intron 1 of the

839 UNC13B Gene Is Associated With Increased Risk of Nephropathy in Patients With

840 Type 1 Diabetes. *Diabetes*, **57**, 2843–2850.

841 Willing E-M, Dreyer C, Oosterhout C van (2012) Estimates of Genetic Differentiation

842 Measured by FST Do Not Necessarily Require Large Sample Sizes When Using

843 Many SNP Markers. *PLOS ONE*, **7**, e42649.

844 Zheng X, Levine D, Shen J *et al.* (2012) A high-performance computing toolset for

845 relatedness and principal component analysis of SNP data. *Bioinformatics*, **28**,

846 3326–3328.

847

848

849 **Data Accessibility:**

850 - RADseq data: NCBI SRA database BioProject ID PRJNA390986

851 - RNAseq data: NCBI SRA database BioProject ID PRJNA390772

852 - Data files including RAD loci consensus sequences, VCF file and sample information

853 available at Dryad doi:10.5061/dryad.8n8t0

854

855 **Table 1.** Simple Mantel test for IBD (isolation-by-distance) and partial Mantel test for the

856 effect of habitat.

Simple Mantel test: test for IBD

| correlation between | Mantel r | p |
|---|---|---|

| $F_{ST}$ | non-transformed distance | 0.34 | 0.003 |
|---|---|---|---|
| $F_{ST}$ | log-transformed distance | 0.29 | 0.008 |
| $F_{ST}/(1- F_{ST})$ | log-transformed distance | 0.28 | 0.007 |

857

Partial Mantel test: test for the effect of habitat while controlling for IBD

| correlation between | | control for | Mantel r | p |
|---|---|---|---|---|
| $F_{ST}$ | same habitat or not | non-transformed distance | 0.48 | 9.00E-06 |
| $F_{ST}$ | same habitat or not | log-transformed distance | 0.50 | 1.00E-06 |
| $F_{ST}/(1- F_{ST})$ | same habitat or not | log-transformed distance | 0.50 | 3.00E-06 |

858　P-values were generated by 999,999 permutations.

859　Here "distance" refers to the geographic distance separating the pair of populations on

860　which the $F_{ST}$ value was computed.**Figure Legends**

861　**Figure 1. Sampling and population structure.** (A), Sampling locations. Each point is a

862　sampling location, and habitat types are indicated by color same to (B).  (B-C), PCA using

863　SNPs that have a minor allele frequency higher than 2%. Each point presents a sample, and

864　samples are colored by their habitat types (B) and by populations (C).

865



866

867

868　**Figure 2. Pairwise population differentiation.** (A), Pairwise $F_{ST}$ between populations

869　correlates with pairwise geographic distance between populations. Empty circles denote

870　pairs of populations from the same type of habitat (shown by the color of the circle). Solid

871　circles are pairs of populations from different types of habitats (shown by colors of the

872　circle and inside). (B), Heat map of pairwise $F_{ST}$. Sampling locations are grouped by

873　habitat type in both axes. (C) and (D), The pairwise $F_{ST}$ of populations from different

874    habitats are greater than the pairwise $F_{ST}$ of populations from the same habitat, even at the

875    same geographic distance. (C) includes all populations from four habitats, and (D) includes

876    only rainforest and ecotone populations. Histogram shows the null distribution of t-

877    statistics generated by 10000 permutations of habitats within different bins of geographic

878    distance (see Methods). Red dot shows the observed value.

879

880

881    **Figure 3**. **Environmental drivers of genomic variation.** (A), Environmental and

882    geographical variables ranked by their importance in explaining SNP allele frequency

883    variation. (B), PC plot indicates the contribution of the environmental variables to the

884    predicted patterns of frequency differentiation, with labeled vectors indicating the direction

885    and magnitude of environmental gradients with greatest contribution. (C), Predicted spatial

886    variation in population-level genetic composition from SNPs. Red points in (C) are

887    locations where actual samples were collected in this study. Colors in (B) and (C) represent

888    gradients in genomic turnover derived from transformed environmental predictors.

889  Locations with similar colors are expected to harbor populations with similar genetic

890  composition.

891

Author Manuscript

**Table 1.** Simple Mantel test for IBD (isolation-by-distance) and partial Mantel test for the effect of habitat.

Simple Mantel test: test for IBD

| correlation between | | Mantel r | p |
|---|---|---|---|
| $F_{ST}$ | non-transformed distance | 0.34 | 0.003 |
| $F_{ST}$ | log-transformed distance | 0.29 | 0.008 |
| $F_{ST}/(1-F_{ST})$ | log-transformed distance | 0.28 | 0.007 |

Partial Mantel test: test for the effect of habitat while controlling for IBD

| correlation between | | control for | Mantel r | p |
|---|---|---|---|---|
| $F_{ST}$ | same habitat or not | non-transformed distance | 0.48 | 9.00E-06 |
| $F_{ST}$ | same habitat or not | log-transformed distance | 0.50 | 1.00E-06 |
| $F_{ST}/(1-F_{ST})$ | same habitat or not | log-transformed distance | 0.50 | 3.00E-06 |

P-values were generated by 999,999 permutations.

Here "distance" refers to the geographic distance separating the pair of populations on which the $F_{ST}$ value was computed.

mec_14270_f1.pdf

**B** legend:
- Ecotone
- Forest
- Island
- Mountain

**C** legend:
- Betare_Oya
- Bioko
- Kribi
- Lac_Lobeke
- Meiganga
- Mt_Tchabal_Gandaba
- Mt_Tchabal_Mbabo
- Ndibi
- Ngaoundaba_Ranch
- Nkwouak
- No_Ayong
- Sakbayeme
- Tibati
- Wakwa
- Zoebefame

mec_14270_f3.jpg