

1 **Supporting Information**

2 **Genome divergence and diversification within a geographic mosaic of coevolution**

3 Thomas L. Parchman, C. Alex Buerkle, Victor Soria-Carrasco, and Craig W. Benkman

4

5 **Methods**

6 **DNA sequencing, assembly, and variant calling**

7 We sequenced DNA from 219 red crossbills representing nine morphologically and vocally
8 differentiated ecotypes (call types), as well as 12 white-winged crossbills (*L. l. leucoptera*;
9 Table S1, Fig. S1). Blood samples were taken from live-caught birds, or tissue samples from
10 vouchered specimens at the University of California at Berkeley Museum of Vertebrate of
11 Zoology. Red crossbills were assigned to call type based on analysis of recorded contact calls
12 (Groth, 1993; Benkman *et al.*, 2009; Irwin, 2010). We used samples of multiple populations
13 from geographically separate regions within red crossbill ecotypes 2, 3, 4, 5, and 7 (Table S1,
14 Fig. S1). DNA was extracted from blood or tissue samples using Qiagen DNeasy 96 Blood and
15 Tissue Kits (Qiagen Inc., Valencia, CA, USA), and quantified with a NanoDrop
16 spectrophotometer (Thermo Sci. Inc., Waltham, MA, USA).

17 We generated reduced representation genomic libraries for each individual using a
18 genotyping by sequencing protocol that we have used previously for a diversity of taxa (Gompert
19 2012; Nosil *et al.*, 2012; Parchman *et al.*, 2012, 2013). We digested the genomic DNA of each
20 individual with two restriction endonucleases (*EcoRI* and *MseI*; NEB Inc., Ipswich, MA, USA),
21 and then ligated double-stranded adaptor oligonucleotides to the digested fragments. Uniquely
22 barcoded *EcoRI* adaptors (containing 8, 9, or 10 bp barcodes) were ligated to fragments from
23 each bird, and the adaptor-ligated fragments from all individuals were subsequently pooled. We

24 then used 288 replicate PCR reactions to amplify the pooled, adaptor-ligated fragments using a
25 high fidelity proofreading polymerase (Iproof polymerase, BioRad Inc., Hercules, CA, USA) and
26 Illumina-based PCR primers. To reduce the proportion of the genome targeted for sequencing,
27 we ran libraries for 3.5 hours at 90V on a 2% agarose gel and excised gel regions containing
28 fragments between approximately 350 and 450 bp in length. We purified these fragments using
29 Qiaquick Gel Extraction Kits (Qiagen Inc., Valencia, CA, USA). A detailed description of this
30 method and associated oligos is deposited at the Dryad digital repository
31 (<http://dx.doi.org/10.5061/dryad.65d97>). For quality control, fragment size distributions were
32 quantified with a BioAnalyzer (Agilent Inc., Santa Clara, CA, USA), and libraries were
33 evaluated for sequencability using qPCR. Sequencing of the libraries was performed by the
34 National Center for Genome Research (Santa Fe, NM, USA) using three lanes of single-end 100-
35 base sequencing on the Illumina HiSeq 2000.

36 Using custom Perl scripts, we removed potential contaminant DNA (PhiX, *E. coli*) and
37 low quality reads from the raw sequence data. Our sequencing library protocol resulted in reads
38 that began with individual barcodes followed by the *EcoRI* restriction site. We trimmed barcodes
39 and restriction-cut-site-associated bases from sequences, and matched barcodes from each read
40 to the correct individual using a custom Perl script that allows for correction of sequencing errors
41 in the barcode sequences. As a high quality *Loxia* genome is unavailable, we used a two-step
42 assembly procedure to assemble reads into homologous genetic regions. We first used SeqMan
43 NGen 3.0.4 (DNASTAR) to perform a *de novo* assembly for a subset of 30 million randomly
44 sampled sequences. This assembly used a minimum match percentage of 92%, a gap penalty of
45 50, a match size of 50 bp, and a mismatch penalty of 15 (full details of assembly
46 parameterization are available from TLP upon request). This assembly placed 24,352,918 reads

47 into 403,678 contigs. After removing low-quality contigs and those with consensus sequences
48 <82 or >88 bases in length, we generated a partial artificial reference genome from the consensus
49 sequences. This reference genome contained 349,865 contig consensus sequences from the
50 original *de novo* assembly. We next assembled the full collection of sequencing reads for each
51 individual onto the reference genome using the ‘aln’ and ‘samse’ algorithms in bwa 0.7.12
52 (Burrows-Wheeler Aligner; Li & Durbin, 2009).

53 We used `samtools` 1.19 and `bcftools` 1.19 (Li *et al.*, 2009) to identify variant sites in
54 the assembled sequence data and to call bi-allelic single nucleotide variants and estimate
55 genotype likelihoods. We first identified single nucleotide polymorphisms (SNPs) in the
56 assembled red crossbill reads for use in population genetic analyses. We disregarded insertions
57 and deletions, and only called variants when 98% of the individuals had at least one read at that
58 locus. Genotype likelihoods were calculated with `bcftools`, stored in variant call format, and
59 converted to a convenient composite genotype likelihood format for downstream analyses. We
60 removed variants for which more than one alternate allele was observed, and randomly selected a
61 single variant from contigs containing multiple variants to increase the independence of loci. In
62 addition, we discarded variants where the observed allele counts from apparent heterozygous
63 individuals were very unlikely ($P < 0.001$) given a binomial distribution with $P = 0.5$. Finally,
64 we excluded low frequency variants and limited analyses to loci having a SNP and minor allele
65 frequencies >0.03.

66 After removing barcodes from the raw reads, and discarding reads with contaminants, we
67 retained 321,627,388 reads representing 219 individuals. Initial *de novo* assembly placed
68 24,352,918 reads into 403,678 contigs; the 349,865 highest quality contigs from this assembly
69 were used as an artificial reference genome. Assembling the reads for each individual onto this

70 reference created alignments containing a total of 253,736,350 reads, with an average of
71 1,071,000 reads aligned per individual. After using `samtools` and `bcftools` to call variant
72 sites, discarding loci with minor allele frequency <0.03 , and randomly sampling a single SNP per
73 contig, we retained a final set of 18,385 SNPs (mean coverage per individual per locus of 7.2x)
74 for population genetic analyses across the red crossbill complex.

75

76

77 **References**

- 78 Benkman CW, Smith JW, Keenan PC, Parchman TL, Santisteban L (2009) A new species
79 of the red crossbill (Fringillidae: *Loxia*) from Idaho. *Condor*, **111**, 169-176.
- 80 Gompert Z, Lucas LK, Nice CC, Fordyce JA, Forister ML, Buerkle CA (2012) Genomic
81 regions with a history of divergent selection affect fitness of hybrids between two
82 butterfly species. *Evolution*, **66**, 2167-2181.
- 83 Groth JG (1993) *Evolutionary differentiation in morphology, vocalizations, and allozymes*
84 *among nomadic sibling species in the North American red crossbill (*Loxia curvirostra*)*
85 *complex*, vol. 127. University of California Publications in Zoology.
- 86 Irwin K (2010) A new and cryptic call type of the red crossbill. *Western Birds*, **41**, 10-25.
- 87 Li H, Durbin R (2009) Fast and accurate short read alignment with Burrows-Wheeler
88 transform. *Bioinformatics*, **25**, 1754-1760.
- 89 Li H, Handsaker B, Wysoker A, et al. (2009) The Sequence Alignment/Map format and
90 SAMtools. *Bioinformatics*, **25**, 2078-2079.
- 91 Nosil P, Gompert Z, Farkas TE, et al. (2012) Genomic consequences of multiple speciation
92 processes in a stick insect. *Proceedings of the Royal Society B: Biological Sciences*, **279**,
93 5058-5065.
- 94 Parchman TL, Gompert Z, Braun MJ, et al. (2013) The genomic consequences of adaptive
95 divergence and reproductive isolation between species of manakins. *Molecular Ecology*,
96 **22**, 3304-3317.
- 97 Parchman TL, Gompert Z, Mudge J, Schilkey F, Benkman CW, Buerkle CA (2012)
98 Genome-wide association genetics of an adaptive trait in lodgepole pine. *Molecular*
99 *Ecology*, **21**, 2991-3005.

100 Supporting Tables and Figures:

101

102

103 Table S1. The number of individuals sampled of *Loxia curvirostra* ecotypes, including
 104 geographic locations of separate samples within each ecotype, as well as white-winged crossbills
 105 *Loxia leucoptera*.

Taxa/species/location and abbreviation	Sample Size	Source
<i>Loxia curvirostra</i>		
Type 1	21	
Brush Mt., VA (1_A)	6	MVZ
Poverty Hollow, VA (1_B)	10	MVZ
Roan Mt., NC (1_C)	5	MVZ
Type 2	74	
Bears Paw Mts., MT (2_D)	4	CWB/PE
Black Hills, SD (2_E)	7	PE
Grand Marais, MI (2_F)	7	MVZ
Little Rocky Mts., MT (2_G)	8	CWB
Sagehen Creek, CA (2_H)	28	MVZ
Sandia Mts., NM (2_I)	12	CWB
Tennessee Pass, CO (2_J)	8	MVZ
Type 3	18	
Berkeley, CA (3_K)	15	MVZ
Sea Lion Point, OR (3_L)	3	MVZ
Type 4	16	
Sea Lion Point, OR (4_M)	9	MVZ
Thompson Plateau, BC (4_N)	7	MVZ
Type 5	22	
San Juan Mts., CO (5_O)	4	MVZ
Santa Catalina Mts., AZ (5_P)	4	MVZ
Tennessee Pass, CO (5_Q)	14	MVZ
Type 6	13	
Chiricahua Mts., AZ (6)	13	MVZ
Type 7	11	
Odell Creek, OR (7_R)	4	MVZ
Thompson Plateau, BC (7_S)	7	MVZ
Type 9/South Hills crossbill	37	
South Hills, ID (9)	37	CWB
Type 10	7	
Eureka, CA (10)	7	MVZ
<i>Loxia leucoptera</i>	12	
Alaska, Colorado	5	MVZ

106 MVZ: Museum of Vertebrate Zoology, UC Berkeley; CWB: Craig W. Benkman field caught; PE: Pim
 107 Edelaar field caught.
 108

109 Table S2. Pairwise estimates of Nei's D (upper triangle) and F_{ST} (lower triangle) among *Loxia*
 110 *curvirostra* ecotypes (call types).

111

	Type 1	Type 2	Type 3	Type 4	Type 5	Type 6	Type 7	Type 9	Type 10
Type 1	-	0.0059	0.0167	0.0098	0.0096	0.0139	0.0122	0.0134	0.0170
Type 2	0.0099	-	0.0068	0.0068	0.0058	0.0101	0.0088	0.0098	0.0141
Type 3	0.0151	0.0111	-	0.0167	0.0109	0.0150	0.0109	0.0145	0.0167
Type 4	0.0169	0.0115	0.0162	-	0.0104	0.0150	0.0132	0.0149	0.0171
Type 5	0.0168	0.0096	0.0178	0.0178	-	0.0137	0.0119	0.0124	0.0177
Type 6	0.0239	0.0182	0.0250	0.0258	0.0230	-	0.0167	0.0170	0.0216
Type 7	0.0203	0.0162	0.0209	0.0242	0.0192	0.0285	-	0.0156	0.0201
Type 9	0.0219	0.0155	0.0231	0.0227	0.0190	0.0259	0.0241	-	0.0215
Type 10	0.0257	0.0257	0.0250	0.0272	0.0278	0.0330	0.0349	0.0331	-

112

113

114 Table S3. Deviance Information Criterion (DIC) estimates for entropy models run for $k = 2$
115 through $k = 9$. Lower estimates of DIC reflect better model fit.

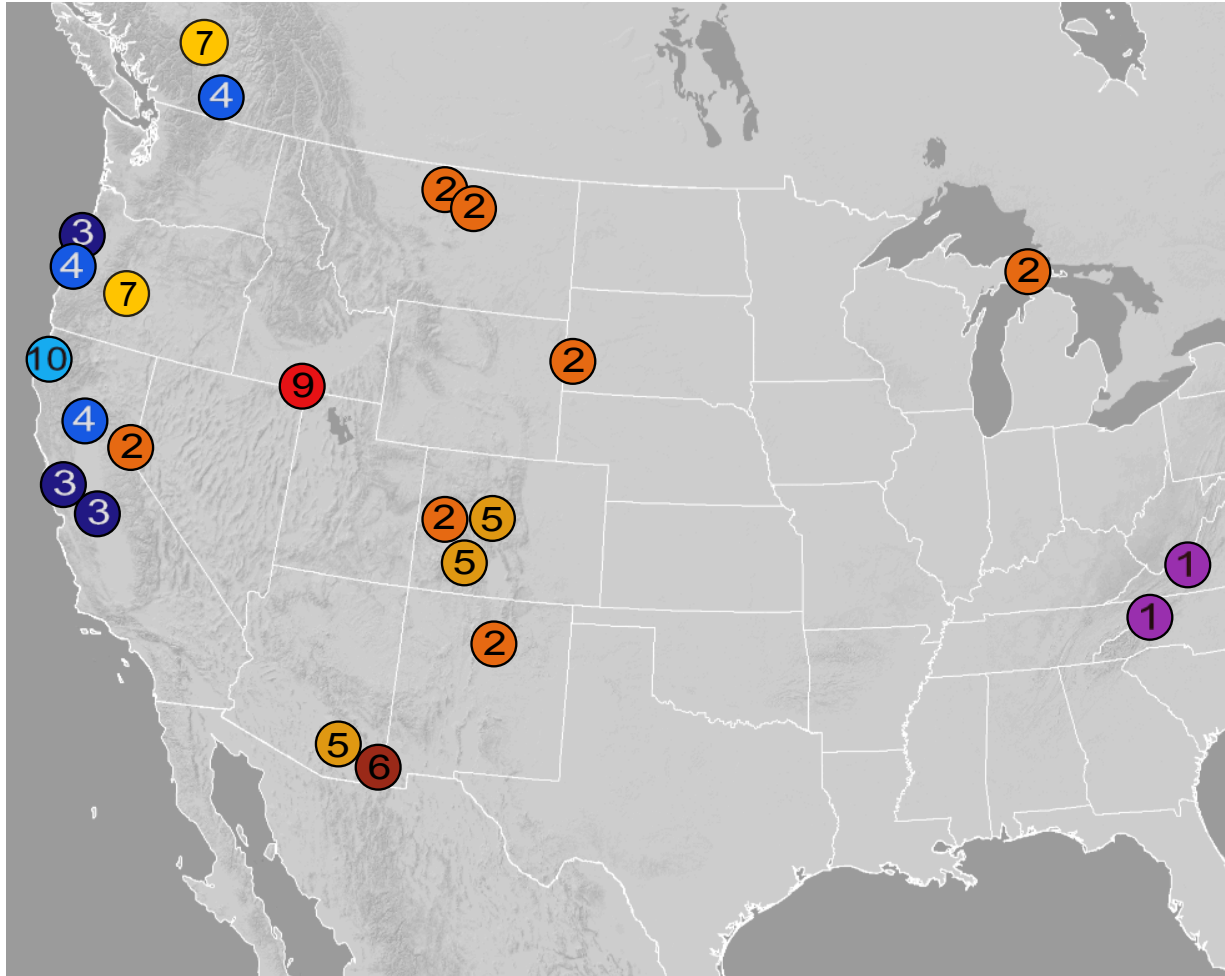
116
117

118	k	DIC
119	2	8,606,675
120	3	8,651,678
121	4	8,502,511
122	5	8,514,392
123	6	8,573,614
124	7	13,860,352
125	8	117,161,960
126	9	151,557,301

127

128 Figure S1. The map illustrates sampling localities for red crossbill (*Loxia curvirostra* complex)
129 ecotypes (call types). Individual points refer to geographically separate collection locations and
130 correspond with Table S1. Circles of the same number and color represent different geographical
131 samples from a given ecotype.

132



133

134