



THE UNIVERSITY *of* EDINBURGH

Edinburgh Research Explorer

A chromosome-level genome assembly of a free-living white-crowned sparrow (*Zonotrichia leucophrys gambelii*)

Citation for published version:

Wu, Z, Miedzinska, K, Krause, J, Perez, J, Wingfield, JC, Meddle, S & Smith, J 2024, 'A chromosome-level genome assembly of a free-living white-crowned sparrow (*Zonotrichia leucophrys gambelii*)', *Scientific Data*, vol. 11, no. 1, 86, pp. 1-10. <https://doi.org/10.1038/s41597-024-02929-6>

Digital Object Identifier (DOI):

[10.1038/s41597-024-02929-6](https://doi.org/10.1038/s41597-024-02929-6)

Link:

[Link to publication record in Edinburgh Research Explorer](#)

Document Version:

Peer reviewed version

Published In:

Scientific Data

General rights

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact openaccess@ed.ac.uk providing details, and we will remove access to the work immediately and investigate your claim.



1 **High quality genome assembly of a free-living white-crowned**
2 **sparrow (*Zonotrichia leucophrys gambelii*)**

3
4 **Authors**

5 Zhou Wu¹, Katarzyna Miedzinska¹, Jesse S. Krause^{2,3}, Jonathan H. Pérez⁴, John C. Wingfield²,
6 Simone L. Meddle¹, Jacqueline Smith¹

7
8 **Affiliations**

9 1. The Roslin Institute and Royal (Dick) School of Veterinary Studies R(D)SVS, The University of
10 Edinburgh, Easter Bush, Midlothian EH25 9RG, UK.

11 2. Department of Neurobiology, Physiology, and Behavior, University of California, Davis, CA
12 95616, USA.

13 3. Department of Biology, University of Nevada Reno, Reno, NV 89557, USA.

14 4. Department of Biology, University of South Alabama, Mobile, AL, 36688, USA.

15
16 Corresponding author(s): Zhou Wu (zhou.wu@roslin.ed.ac.uk) and Jacqueline Smith
17 (Jacqueline.smith@roslin.ed.ac.uk)

18
19
20 **Abstract**

21 The white-crowned sparrow, *Zonotrichia leucophrys*, is a passerine bird with large distribution
22 and which is extensively adapted to environmental changes. It has historically acted as a model
23 species in studies on avian ecology, physiology and behaviour. Here, we present a high-quality
24 chromosome-level genome of *Zonotrichia leucophrys* using PacBio and OmniC sequencing
25 data. Gene models were constructed by combing RNA-seq and Iso-seq data from liver,
26 hypothalamus, and ovary. In total a 1.12 Gb genome was generated, including 31
27 chromosomes assembled in complete scaffolds along with other, unplaced scaffolds. This
28 high-quality genome assembly offers an important genomic resource for the research
29 community using the white-crowned sparrow as a model for understanding avian genome
30 biology and development, and provides a genomic basis for future studies, both fundamental
31 and applied.

33 **Background & Summary**

34 The white-crowned sparrow (WCS; *Zonotrichia leucophrys*) is a small passerine bird that is
35 commonly found in North America and has been historically studied to provide understanding
36 of the biology and ecology in wild, free-living birds. There are five recognized sub-species of
37 white-crowned sparrow (*Zonotrichia leucophrys pugetensis*, *gambelii*, *nuttalli*, *oriantha*, and
38 *leucophrys*) with variation in geographic distribution, appearance and migratory behaviour.
39 White-crowned sparrows offer great opportunities to understand the evolution of subspecies
40 through hybridization and introgression that is characterized by the genomic landscape. As a
41 model species for understanding divergence of behavioural and physiological process, genetic
42 methodologies and approaches have been commonly employed to study the underlying
43 mechanisms using genetic markers on mitochondria or across the whole genome [1]. However,
44 to date, a good quality genome assembly for the white-crowned sparrow has not been
45 available. Previous studies investigating the genetics of *Zonotrichia* species often utilize
46 nucleotide polymorphisms in representative segments of the genome, such as microsatellite
47 markers, genotyping-by-sequencing (GBS), SNP arrays developed for closely-related species,
48 and other restriction site-associated DNA sequencing (RADseq) approaches [1–5]. As a high-
49 quality reference assembly was not available for past genetic studies on white-crowned
50 sparrows, assemblies of other bird species were commonly used as a reference, e.g. genomes
51 of the white-throated sparrow (*Zonotrichia albicollis*), zebra finch (*Taeniopygia guttata*),
52 canary (*Serinus canaria*) or chicken (*Gallus gallus*) [6–9]. The compatibility of these type of
53 studies could be greatly improved by using a specific reference genome assembly and gene
54 models of the white-crowned sparrow.

55 To this end, we present a high-quality chromosome level genome assembly for the white-
56 crowned sparrow using the subspecies *Zonotrichia leucophrys gambelii*. Previous studies
57 suggested that the *Zonotrichia leucophrys* karyotype is $2n=82$ [10–12]. This comprises several
58 pairs of micro-chromosomes, characterized by small size and higher gene density, in which is
59 a feature of bird karyotypes [13]. We combined long-read sequencing (PacBio) and
60 information on DNA compartment proximity (Omni-C) and present a genome of 1.12 Gb,
61 including 3,792 scaffolds with a scaffold N50 of 72 Mb. We assembled 31 relatively complete
62 chromosomes, representing all macro-chromosomes (including the Z sex chromosome), most
63 of the intermediate chromosomes and a good number of micro-chromosomes.

64

65 **Methods**

66 **Sample collection**

67 Samples were collected from two wild, free-living female Gambel's white-crowned sparrows
68 (*Zonotrichia leucophrys gambelii*) captured on breeding grounds in the vicinity of Toolik Lake
69 Research Station on the North Slope of Alaska (N 68° 45', W149° 52') in May 2016 (for DNA
70 extraction) and 20th July 2016 (for RNA extraction). There were no severe weather
71 perturbations (e.g., snowstorm) observed on the days of collection. Following capture with a
72 mist net, a blood sample was collected within three minutes of capture by venipuncture of the
73 alar vein with a 26-gauge needle and transferred into heparinized glass microcapillary tubes
74 (VWR: 15401-56). The birds were quickly sedated with isoflurane and euthanized within three
75 minutes. Following euthanasia, the left pectoralis muscle, brain, liver and ovary were
76 dissected, flash frozen on dry ice, wrapped individually in aluminium foil into labelled plastic
77 bags and kept frozen on dry ice until they were stored in a -80°C freezer upon returning to the
78 laboratory.

79 For DNA extraction, a frozen sample of pectoralis muscle from one individual was sent on dry
80 ice to Dovetail Genomics (California, USA). The RNA samples from the other individual were
81 later shipped on dry ice to the Roslin Institute, University of Edinburgh, UK, where they were

82 stored at -80°C. Approximately 100 mg of liver and ovarian tissue was homogenized for RNA
83 extraction and for the hypothalamus we used 150 mg of tissue.
84 The work was approved by the University of California, Davis, USA Institutional Animal Care
85 and Use Committee (AICUC) under protocol 19758, United States Fish and Wildlife Service -
86 Federal MB90026B-0 and The Animal Welfare and Ethical Review Body at the Roslin Institute,
87 The University of Edinburgh, UK.

88 Genome sequencing

89 Pectoralis muscle was used to obtain high molecular weight DNA (50 to 100 Kb), which was
90 subsequently used for PacBio library preparation after satisfactory quality control. The library
91 preparation, sequencing and scaffolding were carried out by Dovetail Genomics (California,
92 USA) according to their standard genome assembly pipeline (<https://dovetailgenomics.com/>).
93 In short, the PacBio SMRTbell library was constructed using SMRTbell Express Template Prep
94 Kit 2.0 (PacBio, Menlo Park, CA, USA). Sequencing of the genome was performed with PacBio
95 Sequel II 8M SMRT cells, yielding 273.6 Gb data. Sequences were then assembled into scaffolds
96 by using Wtdbg2 [14], followed by contamination detection and duplicated haplotig purging
97 using Blobtools (v2.9) [15] and purge_dups (v1.1.2) [16] respectively.

98 A proximity ligation library was generated by the Omni-C technique [17], followed by
99 sequencing on an Illumina HiSeqX platform. Chromatin was fixed in place in the nucleus with
100 formaldehyde before extraction (for technical note, see [https://dovetailgenomics.com/wp-](https://dovetailgenomics.com/wp-content/uploads/2021/09/Omni-C-Tech-Note.pdf)
101 [content/uploads/2021/09/Omni-C-Tech-Note.pdf](https://dovetailgenomics.com/wp-content/uploads/2021/09/Omni-C-Tech-Note.pdf)). Fixed chromatin was digested with DNase
102 I, fragmented chromatin ends were repaired and biotinylated to adapters followed by
103 proximity ligation. Crosslinks were then reversed, the DNA purified and the biotin
104 subsequently removed. The DNA library was prepared and sequenced to produce 2 x 150bp
105 paired-end reads at a coverage of around 30X. The Omni-C technology uses a sequence-
106 independent endonuclease which provides even, unbiased genome coverage. The HiRise
107 pipeline was employed for further scaffolding of the *de novo* assembly [18]. The genome
108 assembly and Omni-C sequences were used as input for the HiRise pipeline, mainly to
109 determine genomic distance between proximity ligation reads to identify the joins and mis-
110 joins within the scaffolds. The interaction matrix was corrected (--filterThreshold -2.5 3) and
111 visualized by HiExplorer (V3.7.2) [19] (supplementary file 1 **Figure S1**). In addition, we used
112 short-read sequences from a WCS individual (the same one used in RNA-sequencing) to
113 perform genome polishing, using POLCA [20] and pilon (v1.24) [21] with default parameters.

114 RNA-seq sample preparation and sequencing

115 In order to generate a gene model for the white-crowned sparrow genome, we used three
116 RNA-sequencing datasets of the brain (specifically the hypothalamus), liver, and ovary from
117 an individual independently. To isolate RNA for RNA-sequencing, RNA samples were
118 homogenized in TRIzol reagent (Invitrogen) and the Direct-zol RNA Miniprep kit (Zymo
119 Research USA) protocol was followed for RNA extraction. After elution of the total RNA in
120 RNase-free water, we ensured a minimum of 500ng RNA with a concentration of >12.5ng/μL
121 for library preparation. The library construction involved PolyA selection and subsequent
122 sequencing on the BGI DNBSEQ platform [22,23], generating 150 bp paired-end reads and
123 around 30 million sequences per read. The reads were mapped to the genome using STAR
124 (version 2.7.8a) [24] with default options. The RNA-seq data were used to assist the gene
125 model annotations and the mapping rate was also used to validate the completeness of the
126 assembly.

127 Iso-seq library preparation and sequencing

128 The same 3 RNA samples (hypothalamus, liver and ovary) were further prepared for long-read
129 isoform sequencing (Iso-seq). We implemented quality control (QC) using three available
130 methods: NanoDrop spectrophotometer (Thermo Fisher, USA), Qubit 3 fluorometer
131 (Invitrogen, US), and the TapeStation 4200 system (Agilent, US). The starting concentration of
132 the samples were 324 ng/ul, 46 ng/ul and 44 ng/ul, respectively, with RIN > 8. To ensure the
133 quantity of RNA for Iso-seq, libraries were prepared in three technical replicates for ovary and
134 in four technical replicates for liver and hypothalamus. The amount of RNA used for a single
135 reaction was: 0.5 µg for ovary and liver, and 2 µg for hypothalamus. The full-length cDNA was
136 produced using the Teloprime full-length cDNA amplification kit (v1) from Lexogen (cat. No
137 013.24) according to manufacturer's protocols. To determine the Optimal Endpoint PCR (OEP)
138 cycle, a qPCR assay was performed on an aliquot of the full-length double-stranded cDNA using
139 a Light Cycler 480 SW 1.5 machine, and the OEP was determined at 20 cycles corresponding
140 to 80% of the maximum fluorescence value (plateau phase) on the amplification curve.
141 Subsequently, the libraries were purified on columns provided by the manufacturer and the
142 technical replicates were then pooled and subjected to QC. The average concentration of each
143 library was 40 ng/µl. The size distribution, as confirmed by the D5000 screen tape on the
144 TapeStation, ranged from 600 to 2500 bp with a significant peak observed around 1500 bp.
145 Full-length cDNA were then used for PacBio SMRT sequencing on the Sequel system (version
146 2.1). In total, PacBio Iso-seq generated 112 GB data, including 47,186,447 subreads with an
147 average length of 1,389 bp. circular consensus sequences (CCSs) were then created, which
148 subsequently produced 12,219 full-length non-chimeric (flnc) reads with poly-A tail.

149 Genome quality assessment and chromosome assignment

150 Thirty-one relatively complete chromosomes have been assembled, including all macro-
151 chromosomes, intermediate chromosomes and most of the micro-chromosomes,
152 representing, 1, 1A, 2-4, 4A, 5-15, 17-29, Z (**Figure 1**). In total, the size of the Gambel's white-
153 crowned sparrow genome is 1,123,996,003 bp, including 3,792 scaffolds and 4,117 contigs
154 (**Table 1**). Chromosome assignment was based on the zebra finch genome assembly
155 (bTaeGut1.4.pri) (**Figure 2**). In case of future amendment, the corresponding scaffold
156 assignment is presented in **Table 2**. In addition, some scaffolds showed shorter alignment to
157 the zebra finch genome. Although we do not have the full confidence to assign them as
158 complete chromosomes, they can tentatively be assumed to represent the chromosomes with
159 complex sequence structure, such as micro-chromosomes 30, 31, 32, 35 and W. These results
160 are separately represented in supplementary file 1 (**Figure S2**). The prospective chromosomes
161 were visualized by a circos plot using the circlize (v0.4.15) [25] package in R with annotation
162 of genome characteristics, including Ns and gaps, repeat distribution, and GC content.
163 Completeness of the assembly was assessed with Benchmarking Universal Single-Copy
164 Orthologs (BUSCO) for both the assembled genome sequences and the annotated
165 transcriptome (**Figure 3**). The genome has an overall BUSCO score of 96.9% when compared
166 with a total 'aves' (odb10) background, with 0.5% duplication, suggesting good completeness
167 and contiguity of the assembly.

168 The assembly was evaluated by computing quality statistics and detecting repeat elements in
169 the final assembly. First, basic features for the assembly were calculated (e.g., N50, N90, GC
170 content etc.) using available scripts (<https://github.com/WenchaoLin/assemblyStatics>) (**Table**
171 **1**). The genome assembly shows good contiguity and completeness, the scaffold N50 is 71.97
172 Mb, contig N50 is 14.73 Mb and the GC content is 42.80%. In particular, 26,361 bp of Ns are
173 seen in the assembly, making up 0.002% of the total sequence. As for repeat sequences,
174 RepeatModeler (v2.0.2) [26] was used to firstly build the repeat models (such as transposable
175 element families) and then repeat sequences were annotated and masked in place using
176 RepeatMasker (v4.1.2) [27] (**Table 3**). In total, 14.97% of sequences were identified as repeats

177 and soft-masked in the final output. The GC content and repeat content for each chromosome
178 show significantly negative correlation with chromosome size (**Figure 4**). This is particularly
179 pronounced in micro-chromosomes, where GC and repeat content are relatively high. Overall,
180 our assembly for the white-crowned sparrow is comparable to previously published genome
181 assemblies of passerine birds in closely-related families (i.e., *Passerellidae* and *Emberizidae*),
182 regarding the genome size (ranging 1.03 – 1.11 Gb), GC content (41.52 - 42.75%), repeat
183 content (8.4% - 12.19%) and BUSCO score (e.g., complete aves BUSCO ranging 91 - 96.2%)
184 [28,29].

185 Gene model annotation

186 To generate a gene model annotation for the white-crowned sparrow assembly, various
187 sources of evidence and different methodological approaches were integrated, and results
188 consolidated to produce a non-redundant prediction. First, we performed an Iso-seq gene
189 model annotation, following the nf-core/isoseq pipeline for Iso-seq data processing
190 (<https://github.com/nf-core/isoseq>) [30]. In short, raw Iso-seq subreads were converted to
191 CCS using default parameters and subsequently to FLNC reads. LIMA was then used to identify
192 and remove barcodes and primer sequences. Given the library preparation kit used in our
193 study, poly-A clean-up was run with primers suggested by TAMA toolkits [31] for optimized
194 retention of transcripts. The sequences were then mapped to the genome assembly using
195 minimap2 [32], followed by processing with TAMA collapse and TAMA merge. Annotations
196 that were created by subreads belonging to the same tissue were then merged, and
197 annotations further merged across tissues.

198 Furthermore, we used the BRAKER (v2.1.6) annotation pipeline [33] with ETP mode using
199 transcriptomic evidence and protein homology evidence that was retrieved from closely-
200 related reference species. The transcriptomic evidence was acquired from the three RNA-seq
201 tissue samples that were mapped to the genome assembly using STAR (version 2.7.8a) with
202 default parameters [12]. The large protein database includes OrthoDB vertebrate as well as
203 chicken (GRCg6a) and zebra finch (bTaeGut1.4.pri). The aligned RNA-seq and protein database
204 was used to support the training of GeneMark-ETP (version 4.71_lic) [34], followed by
205 AUGUSTUS (version 3.4.0) training and prediction with the same extrinsic information.
206 Augustus training was run with "--species chicken" parameters. Using the BRAKER pipeline, an
207 *ab initio* prediction was also generated [35].

208 In addition, the transcript alignments were further utilized to detect splice junctions using
209 portcullis (1.2.4). The results across multi-samples contributed to a unified set of annotation
210 using PsiCLASS (v1.0.3) [36]. We then predicted open reading frames (ORF) using Transdecoder
211 (5.5.0) (<https://github.com/TransDecoder/TransDecoder>) with an additional search for known
212 proteins using Swiss-Prot (uniprot_sprot, retrieved 2023 May) or pfam (3.1b2) using blastp
213 (2.10.0+) [37] or hmmscan (3.3.2) [38]. Gth (GenomeThreader 1.7.1) was also used to gain a
214 protein alignment based gene structure prediction using the predicted protein sequences
215 (<https://genomethreader.org/>).

216 Finally, the results of the above-mentioned predictions were all combined to a consensus
217 annotation using EVM (EvidenceModeler-v2.0.0). We combined different sources of
218 annotations, including the Iso-seq alignment, transcript alignment, protein alignment,
219 GeneMark, and BRAKER predictions (both *ab initio* and with evidence). The BUSCO score for
220 the transcriptome annotation using 'aves' database for assessment) shows 95.1% complete,
221 2.2% fragmented and 2.7% missing BUSCOs (**Figure 3**). In total, the annotation resulted in
222 25,044 genes and 201,833 exons, with an average gene length of 19382.32 bp, an average
223 exon count of 8.06 per gene, and an average exon length of 217.85 bp (**Figure S3**). The overall

224 noncoding features of the annotation were predicted using CPC2 (0.1) [39]. In total, we
225 identified 18,674 coding genes and 6,370 noncoding genes. In addition, 495 tRNA were
226 detected by using tRNAscan-SE and the details of 737 noncoding sequences (e.g. rRNA) were
227 identified with the Rfam library using Infernal (Supplementary file 2) [40]. We show that
228 overall distribution of gene features correlates with chromosome size (**Figure 4**), in other
229 words, the total number of genes is positively correlated with chromosome length, while the
230 gene density is negatively correlated with chromosome length, with micro-chromosomes (e.g.
231 25, 27, 28, 29) exhibiting high density of gene features (**Figure S4**) as has been shown for
232 chicken, turkey (*Meleagris gallopavo*) and barn swallow (*Hirundo rustica*) [41–43].

233 **Data Records**

234 The data presented in this paper were deposited in National Center for Biotechnology
235 Information (NCBI) databases, with all sequences found under project accession number
236 PRJNA889240. The Whole Genome Shotgun project has been deposited at GenBank under the
237 accession JAPPSN000000000 (we have updated the genome file in NCBI, the latest version will
238 be available upon acceptance of the paper). The version described in this paper is version
239 JAPPSN010000000, the GenBank sequence accession is GCA_028769735.1 (an updated
240 version will become public once accepted). The RNA-seq data can be accessed via
241 SRR21858074, SRR21858075 and SRR21858076; the Iso-seq data is available under
242 SRR21856897, SRR21856898 and SRR21856899; the whole genome sequencing data is
243 available under SRR25788565.

244

245 **Technical Validation**

246 In order to assess the quality of *Zonotrichia Leucophrys* genome assembly, we used multiple
247 methods and datasets for validation. Whole genome alignment to some closely related avian
248 species was performed, including zebra finch (*Taeniopygia guttata*, bTaeGut1.4.pri, RefSeq
249 accession: GCF_003957565.2), and white-throated sparrow (*Zonotrichia albicollis*,
250 *Zonotrichia albicollis*-1.0.1, Ensembl 108). NUCmer (NUCleotide MUMmer) aligner built in
251 MUMmer (version 3.1) [44] was used with default parameters. The percentage of total aligned
252 bases to zebra finch and white-throated sparrow is 82.43% and 80.38%, respectively.

253 We then filtered the alignment for the minimum alignment identity at 30%. A DOT plot was
254 used to visualize the cross-species alignment by adapting R code from dotPlotly
255 (<https://github.com/tpoorten/dotPlotly>) with alignment cut off: queries with total alignments
256 > 80000 bp, minimum alignments > 3000 bp.

257 To evaluate the quality of the RNA-seq data, FastQC (v0.11.7) [45] and QualiMap (v.2.2.1) [46]
258 were used to assess the sequence and mapping quality, respectively. As shown in **Figure S5**,
259 the input RNA-seq data has high quality, as demonstrated by the statistics of reads, e.g. base
260 quality. The RNA-seq data was mapped to our assembled genome using STAR (version 2.7.8a)
261 [24]. The input raw reads and mapping quality are summarized in **Table 4**, with an average
262 uniquely mapping rate of 90.98%, indicating good quality and successful alignment to the
263 genome assembly. Similarly, the short-read whole-genome sequencing data were mapped to
264 the final assembly and then assessed for mapping quality. BWA-MEM [47] was used for
265 mapping with recommended parameters, and the percentage of mapped reads was 99.4%
266 with a mean mapping score of 22.07.

267 **Code Availability**

268 The majority of the data analyses were completed using standard bioinformatic tools running
269 on the Linux system. The version and code/parameters of the main software tools are
270 described in text. Additional scripts used to generate the results and the figures can be found
271 in the github repository https://github.com/wzuhou/Genome_assembly_annotation.

272

273 **Acknowledgements**

274 This work was supported by Biotechnology and Biological Sciences Research Council (BBSRC),
275 UK BB/V001647/1 to JS and SLM, Roslin Institute Strategic Grant funding from BBSRC,
276 BBS/E/D/30002276 and BBS/E/RL/230001C, to JS and SLM and National Science Foundation
277 (NSF) Office of Polar Programs ARC 0909133 to JCW and Integrative Organismal Systems IOS
278 1558049 to JWC and SLM. We thank Sebastien Guizard for helpful suggestions on running the
279 nf-core/iseq pipeline and Valerie Bishop for laboratory assistance. We also thank Hannah J.
280 Lau, Helen E. Chmura, Jeffrey Cheah, and Ryan E. Swanson for their help in the field and logistic
281 support from the staff at Toolik Field Station, The University of Alaska Fairbanks, USA as well
282 as Brian Barnes and Jeannette Moore for logistical support while in Fairbanks, Alaska USA. We
283 also thank Dovetail Genomics (California, USA) for preparing the DNA for genome sequencing
284 and for preparing the draft assembly for the genome. For the purpose of open access, the
285 author has applied a Creative Commons Attribution (CC BY) licence to any Author Accepted
286 Manuscript version arising from this submission.

287

288 **Author contributions**

289 SLM, JSK and JHP collected the samples. KM conducted RNA sample preparation and Iso-seq
290 library preparation. JCW provided the genome samples. ZW performed all data analyses and
291 wrote the manuscript. JS, SLM and JCW provided supervision. All authors contributed to the
292 manuscript preparation.

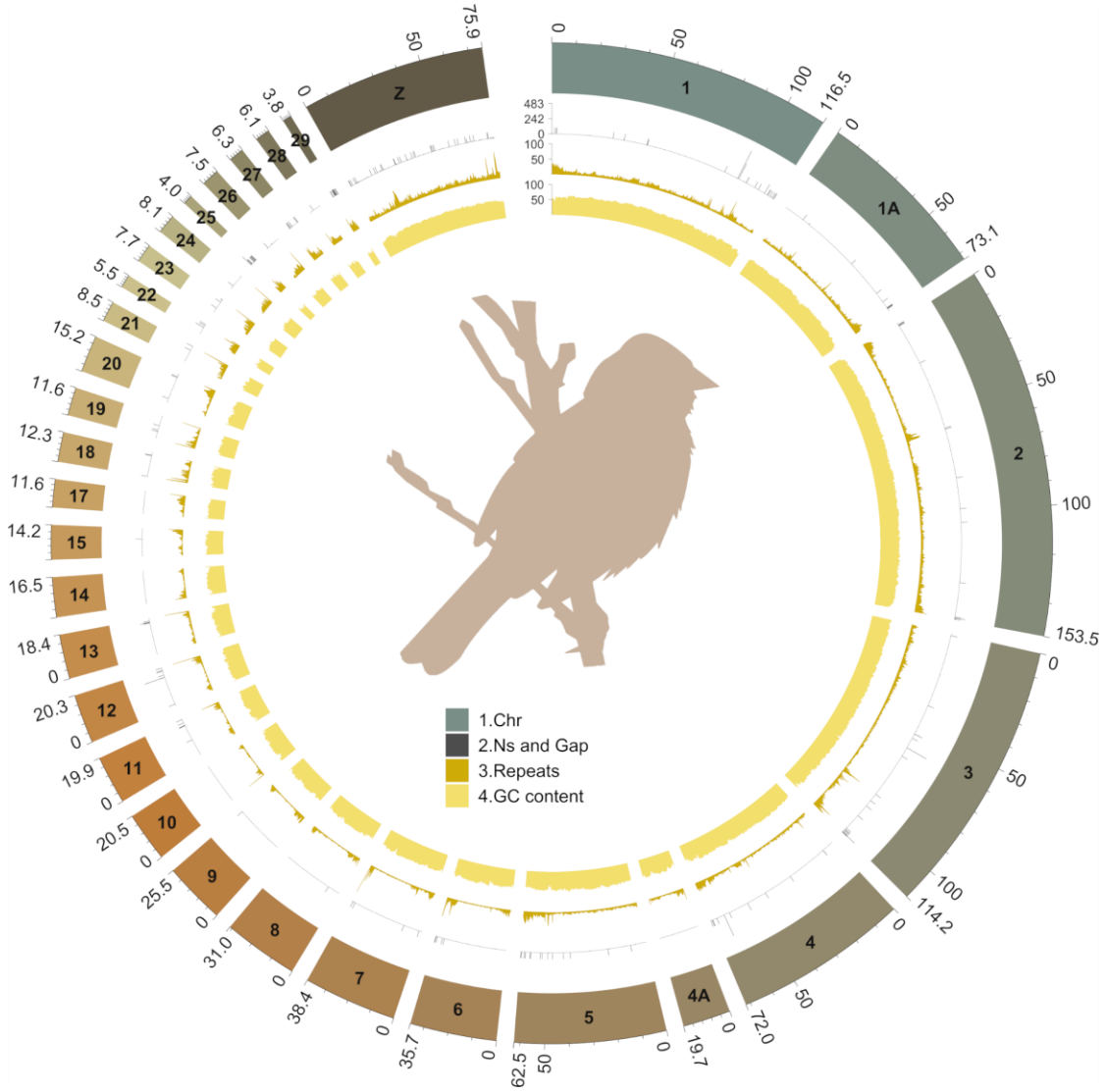
293

294 **Competing interests**

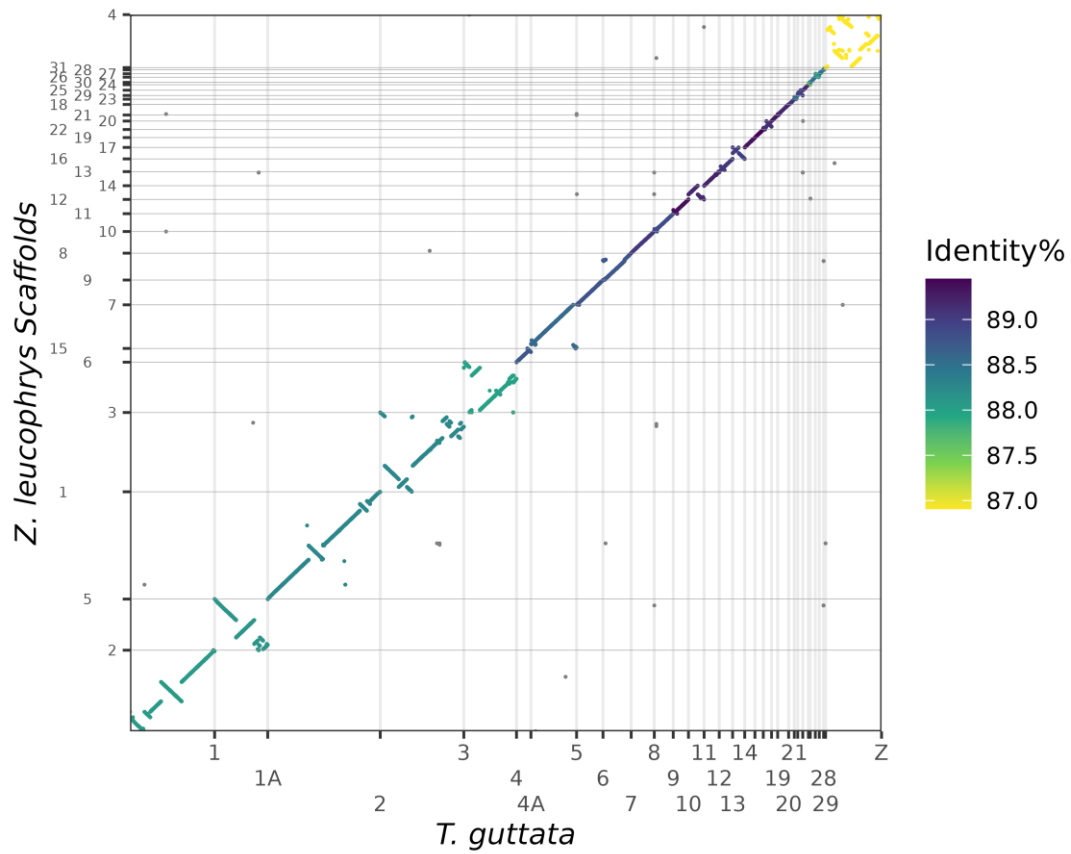
295 The authors declare there is no conflict of interest.

296

297



299 **Figure 1**
 300 Overview of the genome assembly of the white-crowned sparrow (*Zonotrichia leucophrys*
 301 *gambelii*). The size of chromosomes is displayed in Mb, the Ns and Gaps are in bp, while
 302 repeats and GC content are presented as percentages (window size 200k). The bird silhouette
 303 image was downloaded from <https://www.phylopic.org/> (provided 2017 Aug 29, by Matt
 304 Wilkins) under the Creative Commons (CC0) 1.0 Universal Public Domain Dedication License.
 305

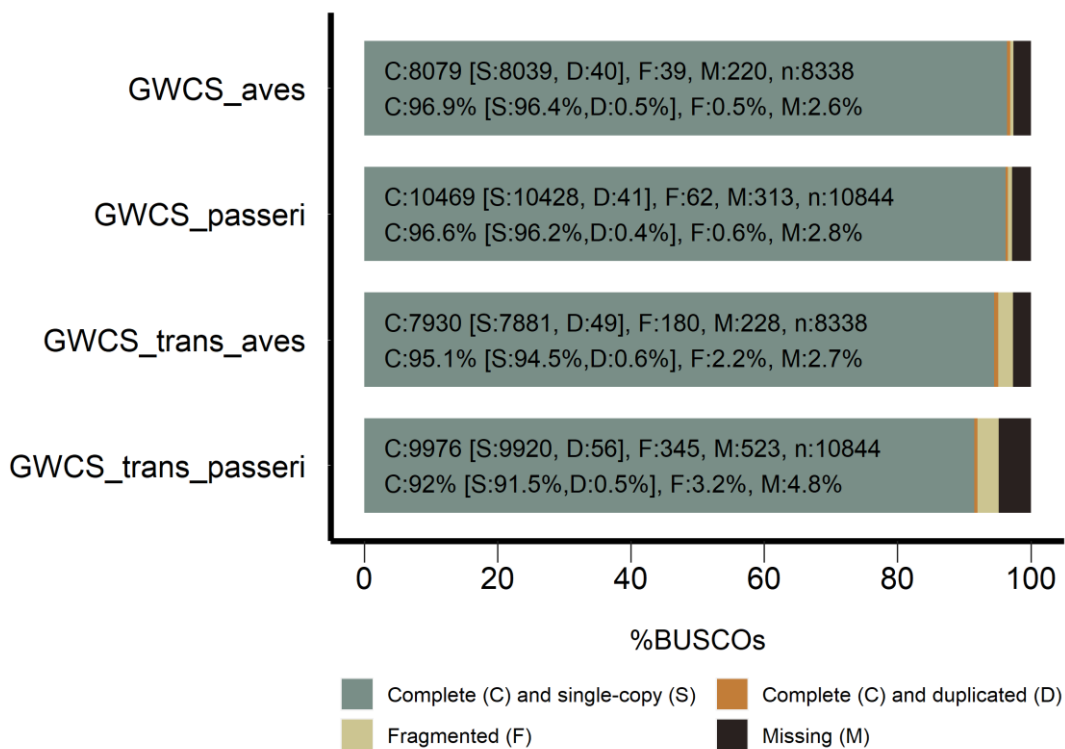


306
307
308
309
310
311

Figure 2

Whole-genome alignment between assemblies of the white-crowned sparrow (*Zonotrichia leucophrys gambelii*) and zebra finch (*Taeniopygia guttata*; version: bTaeGut1.4.pri). The y-axis displays the representative scaffolds of the white-crowned sparrow genome.

BUSCO Assessment Results

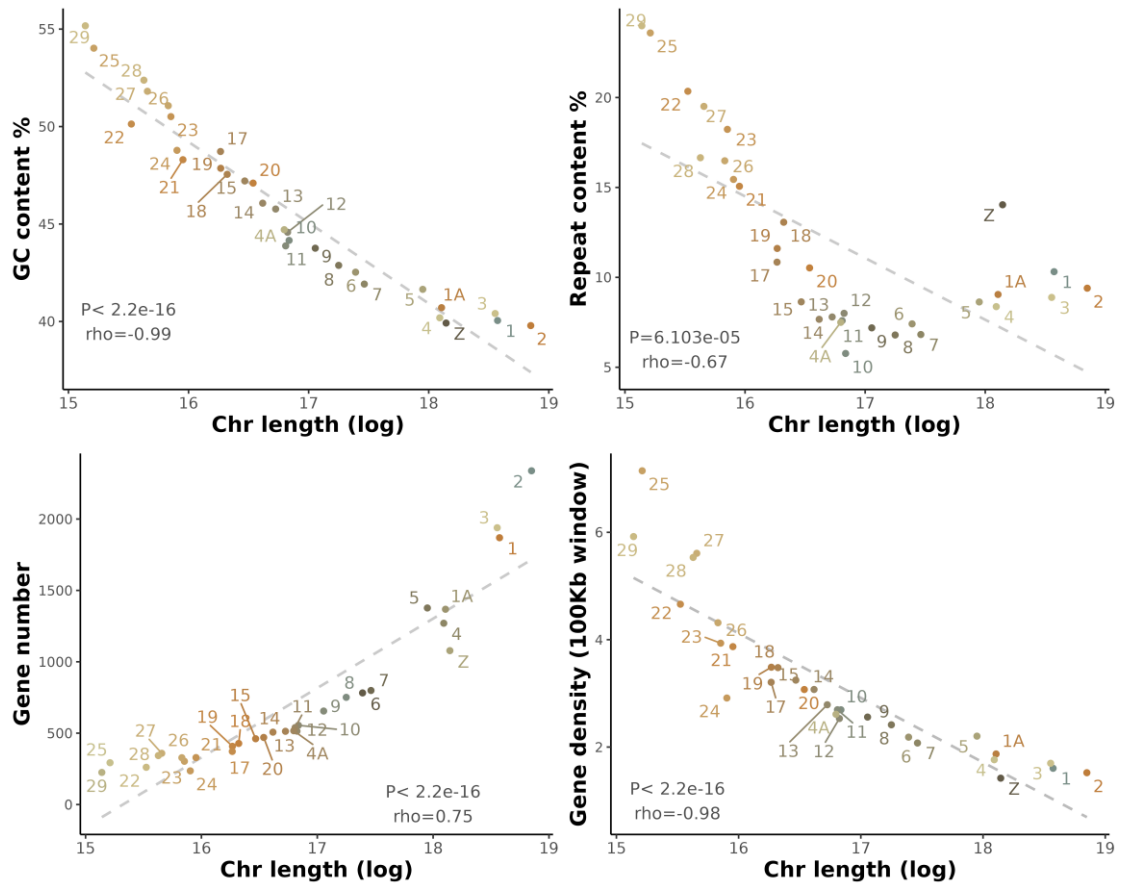


312 **Figure 3**

313 Assessment of Benchmarking Universal Single-Copy Orthologs (BUSCOs) of the white-crowned
 314 sparrow (*Zonotrichia leucophrys gambelii*) genome and transcriptome using aves and
 315 Passeriformes (abbreviated as passeri) (odb10) databases.

316

317



318

319 **Figure 4**

320 Correlation between chromosome size and GC content, repeat elements, number of genes
 321 and gene density of the white-crowned sparrow (*Zonotrichia leucophrys gambelii*) genome.

322 The chromosome size is log transformed and the P value was calculated by Spearman's test.

323

324 **Figure Legends**

325 **Figure 1**

326 Overview of the genome assembly of the white-crowned sparrow (*Zonotrichia leucophrys*
327 *gambelii*). The size of chromosomes is displayed in Mb, the Ns and Gaps are in bp, while
328 repeats and GC content are presented as percentages (window size 200k). The bird silhouette
329 image was downloaded from <https://www.phylopic.org/> (provided 2017 Aug 29, by Matt
330 Wilkins) under the Creative Commons (CC0) 1.0 Universal Public Domain Dedication License.

331

332 **Figure 2**

333 Whole-genome alignment between assemblies of the white-crowned sparrow (*Zonotrichia*
334 *leucophrys gambelii*) and zebra finch (*Taeniopygia guttata*; version: bTaeGut1.4.pri). The y-
335 axis displays the representative scaffolds of the white-crowned sparrow genome.

336

337 **Figure 3**

338 Assessment of Benchmarking Universal Single-Copy Orthologs (BUSCOs) of the white-crowned
339 sparrow (*Zonotrichia leucophrys gambelii*) genome and transcriptome using aves and
340 passeriformes (odb10) databases.

341

342 **Figure 4**

343 Correlation between chromosome size and GC content, repeat elements, number of genes
344 and gene density of the white-crowned sparrow (*Zonotrichia leucophrys gambelii*) genome.
345 The chromosome size is log transformed and the P value was calculated by Spearman's test.

346

347 **Tables**

348 **Table 1**

349 Assessment of the white-crowned sparrow genome assembly.

Assembly features	Gambels_ncbi_update
Counts of scaffold sequences	3,792
Length of scaffold sequences	1,123,996,003
Largest scaffold name	Scaffold_1_153547327
Largest scaffold length	153,547,327
Scaffold N50	71,969,017
Counts of N50	6
Scaffold N90	6,309,133
Counts of N90	27
GC content (%)	42.80%
N Length	26,361
N content (%)	0.002%
Counts of contigs	4,117
Maximum length of contigs	40,609,704
contig N50	14,729,340
Counts of contig N50	25
contig N90	546,537
Counts of contig N90	179

350

351 **Table 2** Chromosome assignment for the white-crowned sparrow assembly.

Scaffold name	Chromosome
Scaffold_2_116484495	1
Scaffold_5_73051372	1A

Scaffold_1_153547327	2
Scaffold_3_114162194	3
Scaffold_6_71969017	4
Scaffold_15_19713544	4A
Scaffold_7_62472784	5
Scaffold_9_35708988	6
Scaffold_8_38401667	7
Scaffold_10_31016323	8
Scaffold_11_25524209	9
Scaffold_12_20527583	10
Scaffold_14_19948824	11
Scaffold_13_20270949	12
Scaffold_16_18355265	13
Scaffold_17_16474596	14
Scaffold_19_14189122	15
Scaffold_22_11597714	17
Scaffold_20_12261182	18
Scaffold_21_11615082	19
Scaffold_18_15211132	20
Scaffold_23_8480127	21
Scaffold_29_5518869	22
Scaffold_25_7671743	23
Scaffold_24_8071077	24
Scaffold_30_4037257	25
Scaffold_26_7504969	26
Scaffold_27_6309133	27
Scaffold_28_6126336	28
Scaffold_31_3761913	29
Scaffold_4_75875312	Z

352

353

354

Table 3

Repeat elements identified in the assembly.

Repeats	Count	Length (bp)	Percentage (%)
Retroelements	234,891	96,498,034	8.59
SINEs	2,311	291,064	0.03
LINEs	133,634	37,252,295	3.31
LTR elements	98,946	58,954,675	5.25
DNA transposons	7,445	1,092,740	0.10
Rolling-circles	1,858	1,015,043	0.09
Unclassified	89,799	46,879,085	4.17
Total interspersed repeats		144,469,859	12.85
Small RNA	749	82,339	0.01
Satellites	7,681	5,697,135	0.51
Simple repeats	235,850	13,986,714	1.24
Low complexity	49,091	3,115,412	0.28

Bases masked		168,298,524	14.97
--------------	--	-------------	-------

355

356

Table 4

357

Validation of the white-crowned (*Zonotrichia leucophrys gambelii*) RNA-seq dataset.

Sample type	Number of input reads (pairs)	Uniquely mapped reads	Number of total splices	Mismatch rate per base
Gonad	33,585,925	92.25%	31,206,515	0.67%
Hypothalamus	34,035,354	89.73%	20,131,958	0.62%
Liver	34,085,391	90.97%	29,982,801	0.57%

358

359

References

360

1. Taylor RS, Bramwell AC, Clemente-Carvalho R, Cairns NA, Bonier F, Dares K, et al.

361

Cytonuclear discordance in the crowned-sparrows, *Zonotrichia atricapilla* and *Zonotrichia*

362

leucophrys. *Mol Phylogenet Evol.* 2021;162:107216.

363

2. Mccallum Q, Askelson K, Fogarty F, Natola L, Nikelski E, Huang A, et al. Extreme sex

364

chromosome differentiation, likely driven by inversion, contrasts with mitochondrial

365

paraphyly between species of crowned sparrows. *bioRxiv preprint.* 2022;

366

3. Cheviron ZA, Whitehead A, Brumfield RT. Transcriptomic variation and plasticity in rufous-

367

collared sparrows (*Zonotrichia capensis*) along an altitudinal gradient. *Mol Ecol.*

368

2008;17:4556–69.

369

4. Lipshutz SE, Overcast IA, Hickerson MJ, Brumfield RT, Derryberry EP. Behavioural response

370

to song and genetic divergence in two subspecies of white-crowned sparrows (*Zonotrichia*

371

leucophrys). *Mol Ecol.* 2017;26:3011–27.

372

5. Weckstein JD, Zink RM, Blackwell-Rago RC, Nelson DA. Anomalous variation in

373

mitochondrial genomes of White-crowned (*Zonotrichia leucophrys*) and Golden-crowned (*Z.*

374

atricapilla) Sparrows: Pseudogenes, hybridization, or incomplete lineage sorting? *Auk.*

375

2001;118:231–6.

376

6. Krause JS, McGuigan MA, Bishop VR, Wingfield JC, Meddle SL. Decreases in

377

Mineralocorticoid but not Glucocorticoid Receptor mRNA Expression During the Short Arctic

378

Breeding Season in Free-Living Gambel's White-Crowned Sparrow (*Zonotrichia leucophrys*

379

gambelii). *J Neuroendocrinol.* 2015;27:66–75.

380

7. Krause JS, Watkins T, Reid AMA, Cheah JC, Pérez JH, Bishop VR, et al. Gene expression of

381

sex steroid metabolizing enzymes and receptors in the skeletal muscle of migrant and

382

resident subspecies of white-crowned sparrow (*Zonotrichia leucophrys*). *Oecologia.*

383

2022;199:549–62.

384

8. Krause JS, Pérez JH, Reid AMA, Cheah J, Bishop V, Wingfield JC, et al. Acute restraint stress

385

does not alter corticosteroid receptors or 11 β -hydroxysteroid dehydrogenase gene

386

expression at hypothalamic–pituitary–adrenal axis regulatory sites in captive male white-

387

crowned sparrows (*Zonotrichia leucophrys gambelii*). *Gen Comp Endocrinol.* 2021;303.

388

9. Jones S, Pfister-Genskow M, Cirelli C, Benca RM. Changes in brain gene expression during

389

migration in the white-crowned sparrow. *Brain Res Bull.* 2008;76:536–44.

390

10. Shields GF. Comparative Avian Cytogenetics: A Review [Internet]. *Condor.* 1982 [cited

391

2023 Feb 13]. p. 45. Available from: [https://www-jstor-](https://www-jstor-org.ezproxy.is.ed.ac.uk/stable/1367820?sid=primo)

392

[org.ezproxy.is.ed.ac.uk/stable/1367820?sid=primo](https://www-jstor-org.ezproxy.is.ed.ac.uk/stable/1367820?sid=primo)

393

11. Shields GF. Bird chromosomes. *Current ornithology Vol1.* 1983;189–209.

394

12. Degrandi TM, Barcellos SA, Costa AL, Garnero AD V, Hass I, Gunski RJ. Introducing the

395

Bird Chromosome Database: An Overview of Cytogenetic Studies in Birds. *Cytogenet*

396

Genome Res. 2020;160:199–205.

397

13. Degrandi TM, Barcellos SA, Costa AL, Garnero AD V, Hass I, Gunski RJ. Introducing the

398

Bird Chromosome Database: An Overview of Cytogenetic Studies in Birds. *Cytogenet*

399 Genome Res [Internet]. 2020 [cited 2023 Feb 13];160:199–205. Available from:
400 www.karger.com/cgr

401 14. Ruan J, Li H. Fast and accurate long-read assembly with wtdbg2. *Nat Methods*.
402 2020;17:155–8.

403 15. Laetsch DR, Blaxter ML. BlobTools: Interrogation of genome assemblies. *F1000Res*.
404 2017;6:1287.

405 16. Guan D, Guan D, McCarthy SA, Wood J, Howe K, Wang Y, et al. Identifying and removing
406 haplotypic duplication in primary genome assemblies. *Bioinformatics*. 2020;36:2896–8.

407 17. Lieberman-Aiden E, van Berkum NL, Williams L, Imakaev M, Ragozcy T, Telling A, et al.
408 Comprehensive Mapping of Long-Range Interactions Reveals Folding Principles of the Human
409 Genome. *Science (1979)*. 2009;326:289–93.

410 18. Putnam NH, O’Connell BL, Stites JC, Rice BJ, Blanchette M, Calef R, et al. Chromosome-
411 scale shotgun assembly using an in vitro method for long-range linkage. *Genome Res*.
412 2016;26:342–50.

413 19. Wolff J, Rabbani L, Gilsbach R, Richard G, Manke T, Backofen R, et al. Galaxy HiCExplorer
414 3: A web server for reproducible Hi-C, capture Hi-C and single-cell Hi-C data analysis, quality
415 control and visualization. *Nucleic Acids Res*. 2020;48:W177–84.

416 20. Zimin A V., Salzberg SL. The genome polishing tool POLCA makes fast and accurate
417 corrections in genome assemblies. *PLoS Comput Biol*. 2020;16.

418 21. Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, et al. Pilon: An
419 integrated tool for comprehensive microbial variant detection and genome assembly
420 improvement. *PLoS One*. 2014;9:112963.

421 22. Jeon SA, Park JL, Park S-J, Kim JH, Goh S-H, Han J-Y, et al. Comparison between MGI and
422 Illumina sequencing platforms for whole genome sequencing. *Genes Genomics*.
423 2021;43:713–24.

424 23. Patterson J, Carpenter EJ, Zhu Z, An D, Liang X, Geng C, et al. Impact of sequencing depth
425 and technology on de novo RNA-Seq assembly. *BMC Genomics*. 2019;20.

426 24. Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, et al. STAR: ultrafast
427 universal RNA-seq aligner. *Bioinformatics*. 2013;29:15–21.

428 25. Gu Z, Gu L, Eils R, Schlesner M, Brors B. Circlize implements and enhances circular
429 visualization in R. *Bioinformatics*. 2014;30:2811–2.

430 26. Flynn JM, Hubley R, Goubert C, Rosen J, Clark AG, Feschotte C, et al. RepeatModeler2 for
431 automated genomic discovery of transposable element families. *Proc Natl Acad Sci U S A*.
432 2020;117:9451–7.

433 27. Chen N. Using RepeatMasker to identify repetitive elements in genomic sequences.
434 *Current protocols in bioinformatics / editorial board, Andreas D Baxevanis . [et al]*.
435 2004;Chapter 4:4.10.1-4.10.14.

436 28. Friis G, Vizueta J, Ketterson ED, Milá B. A high-quality genome assembly and annotation
437 of the dark-eyed junco *Junco hyemalis* , a recently diversified songbird. Dunlap J, editor. G3
438 Genes|Genomes|Genetics [Internet]. 2022 [cited 2022 May 25];12. Available from:
439 <https://academic.oup.com/g3journal/article/doi/10.1093/g3journal/jkac083/6566302>

440 29. Feng S, Stiller J, Deng Y, Armstrong J, Fang Q, Reeve AH, et al. Dense sampling of bird
441 diversity increases power of comparative genomics. *Nature [Internet]*. 2020 [cited 2021 Mar
442 30];587:252–7. Available from: <https://doi.org/10.1038/s41586-020-2873-9>

443 30. Guizard S, Miedzinska K, Smith J, Smith J, Kuo R, Davey M, et al. nf-core/iseq: Simple
444 gene and isoform annotation with PacBio Iso-Seq long-read sequencing. Robinson P, editor.
445 *Bioinformatics*. 2023;

446 31. Kuo RI, Cheng Y, Zhang R, Brown JWS, Smith J, Archibald AL, et al. Illuminating the dark
447 side of the human transcriptome with long read transcript sequencing. *BMC Genomics*.
448 2020;21.

449 32. Coster W De, Rijk P De, Roeck A De, Pooter T De, D'Hert S, Strazisar M, et al. Structural
450 variants identified by Oxford Nanopore PromethION sequencing of the human genome.
451 *Genome Res.* 2019;29:1178–87.

452 33. Brůna T, Hoff KJ, Lomsadze A, Stanke M, Borodovsky M. BRAKER2: automatic eukaryotic
453 genome annotation with GeneMark-EP+ and AUGUSTUS supported by a protein database.
454 *NAR Genom Bioinform.* 2021;3:1–11.

455 34. Brůna T, Lomsadze A, Borodovsky M. GeneMark-EP+: eukaryotic gene prediction with
456 self-training in the space of genes and proteins. *NAR Genom Bioinform.* 2020;2.

457 35. Ter-Hovhannisyan V, Lomsadze A, Chernoff YO, Borodovsky M. Gene prediction in novel
458 fungal genomes using an ab initio algorithm with unsupervised training. *Genome Res.*
459 2008;18:1979–90.

460 36. Song L, Sabunciyany S, Yang G, Florea L. A multi-sample approach increases the accuracy
461 of transcript assembly. *Nat Commun.* 2019;10:1–7.

462 37. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment search tool. *J*
463 *Mol Biol.* 1990;215:403–10.

464 38. Eddy SR. A probabilistic model of local sequence alignment that simplifies statistical
465 significance estimation. *PLoS Comput Biol.* 2008;4:e1000069.

466 39. Kang YJ, Yang DC, Kong L, Hou M, Meng YQ, Wei L, et al. CPC2: A fast and accurate coding
467 potential calculator based on sequence intrinsic features. *Nucleic Acids Res.* 2017;45:W12–6.

468 40. Nawrocki EP, Eddy SR. Infernal 1.1: 100-fold faster RNA homology searches.
469 *Bioinformatics.* 2013;29:2933–5.

470 41. Smith J, Bruley CK, Paton IR, Dunn I, Jones CT, Windsor D, et al. Differences in gene
471 density on chicken macrochromosomes and microchromosomes. *Anim Genet.* 2000;31:96–
472 103.

473 42. City NY, City NY, City NY, City NY, Cedex O. Pangenomics provides insights into the role of
474 synanthropy in barn swallow evolution. 2022;

475 43. Barros CP, Derks MFL, Mohr J, Wood BJ, Crooijmans RPMA, Megens H-J, et al. A new
476 haplotype-resolved turkey genome to enable turkey genetics and genomics research.
477 *Gigascience.* 2022;12.

478 44. Kurtz S, Phillippy A, Delcher AL, Smoot M, Shumway M, Antonescu C, et al. Versatile and
479 open software for comparing large genomes. 2004;5:12.

480 45. Andrews S. FastQC A quality control tool for high throughput sequence data. FastQC A
481 quality control tool for high throughput sequence data.
482 2010;<http://www.bioinformatics.babraham.ac.uk/projects/>.

483 46. Okonechnikov K, Conesa A, García-Alcalde F. Qualimap 2: Advanced multi-sample quality
484 control for high-throughput sequencing data. *Bioinformatics.* 2016;32:292–4.

485 47. Li H, Durbin R. Fast and accurate long-read alignment with Burrows-Wheeler transform.
486 *Bioinformatics.* 2010;26:589–95.

487