

An annotated chromosome-scale reference genome for Eastern black-eared wheatear (*Oenanthe melanoleuca*)

Valentina Peona [†], Octavio Manuel Palacios-Gimenez [†], Dave Lutgen [†], Remi André Olsen, Niloofar Alaei Kakhki, Pavlos Andriopoulos , Vasileios Bontzorlos , Manuel Schweizer , Alexander Suh , Reto Burri *

Department of Organismal Biology—Systematic Biology, Science for Life Laboratory, Evolutionary Biology Centre, Uppsala University, 75236 Uppsala, Sweden

Department of Population Ecology, Institute of Ecology and Evolution, Friedrich Schiller University Jena, 07743 Jena, Germany

German Centre for Integrative Biodiversity Research (iDiv) Halle-Jena-Leipzig, 04103 Leipzig, Germany

Department of Biology, Institute of Ecology and Evolution, University of Bern, 3012 Bern, Switzerland

Swiss Ornithological Institute, CH-6204 Sempach, Switzerland

Science for Life Laboratory, Department of Biochemistry and Biophysics, Stockholm University, 17165 Solna, Sweden

Section of Ecology and Systematics, Department of Biology, National and Kapodistrian University of Athens, 15772 Athens, Greece

TYTO—Association for the Management and Conservation of Biodiversity in Agricultural Ecosystems, 41335 Larisa, Greece

Natural History Museum Bern, 3005 Bern, Switzerland

School of Biological Sciences, University of East Anglia, NR4 7TU Norwich, UK

*Corresponding author: Swiss Ornithological Institute, Seerose 1, CH-6204 Sempach, Switzerland. Email: reto.burri@vogelwarte.ch

[†]These authors contributed equally to the work.

Abstract

Pervasive convergent evolution and in part high incidences of hybridization distinguish wheatears (songbirds of the genus *Oenanthe*) as a versatile system to address questions at the forefront of research on the molecular bases of phenotypic and species diversification. To prepare the genomic resources for this venture, we here generated and annotated a chromosome-scale assembly of the Eastern black-eared wheatear (*Oenanthe melanoleuca*). This species is part of the *Oenanthe hispanica* complex that is characterized by convergent evolution of plumage coloration and high rates of hybridization. The long-read-based male nuclear genome assembly comprises 1.04 Gb in 32 autosomes, the Z chromosome, and the mitogenome. The assembly is highly contiguous (contig N50, 12.6 Mb; scaffold N50, 70 Mb), with 96% of the genome assembled at the chromosome level and 95.5% benchmarking universal single-copy orthologs (BUSCO) completeness. The nuclear genome was annotated with 18,143 protein-coding genes and 31,333 mRNAs (annotation BUSCO completeness, 98.0%), and about 10% of the genome consists of repetitive DNA. The annotated chromosome-scale reference genome of Eastern black-eared wheatear provides a crucial resource for research into the genomics of adaptation and speciation in an intriguing group of passerines.

Keywords: birds, open-habitat chats, *Oenanthe melanoleuca*, *Oenanthe hispanica* complex, transcriptome, repeat content, transposable elements

Introduction

Wheatears of the genus *Oenanthe* and their relatives—together referred to as “open-habitat chats”—are a group of songbirds that display several remarkable characteristics distinguishing them as a versatile system to address key questions on the evolution of phenotypes and formation of species. Many phenotypes, including multiple conspicuous color ornaments, seasonal migration, and sexual dimorphism, appear independently in multiple branches within open-habitat chats, suggesting a high incidence of convergent evolution (Aliabadian et al. 2012; Alaei Kakhki et al. 2013; Schweizer et al. 2019a, 2019b). Furthermore, hybridization is observed in several species complexes and occurs at notably high rates in the *Oenanthe hispanica* complex that consists of 4 currently recognized taxa (Schweizer et al. 2019a, 2019b): Western black-eared wheatear (*O. hispanica*), pied wheatear (*Oenanthe pleschanka*), cyprus wheatear (*Oenanthe cyprica*), and Eastern black-eared wheatear (*Oenanthe melanoleuca*; Fig. 1). Pied

and Eastern black-eared wheatear hybridize pervasively at the western shores of the Black Sea, in the Caucasus, and in the Alborz mountains of northern Iran (Haffer 1977; Panov 2005). The resulting introgression reaches beyond the hybrid zones (Schweizer et al. 2019a, 2019b), and hybrid zones themselves sport admixed phenotypes that display combinations of plumage color phenotypes divergent between species (mantle and neck-side coloration) (Haffer 1977; Panov 2005). Finally, a phenotype divergently expressed between many wheatear species, black-or-white throat coloration, segregates as polymorphisms in 3 species of the *O. hispanica* complex. Once a high-quality reference genome is available, this polymorphism and the recombination of mantle and neck-side coloration in hybrids provide an excellent opportunity to map these phenotypes to the genome (Buerkle and Lexer 2008) and study their convergent evolution across open-habitat chats. Furthermore, hybridization in several geographic regions enables insights into common or idiosyncratic patterns of evolution under hybridization (Gompert et al. 2017).

Received: December 22, 2022. Accepted: March 8, 2023

© The Author(s) 2023. Published by Oxford University Press on behalf of The Genetics Society of America.

This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<https://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

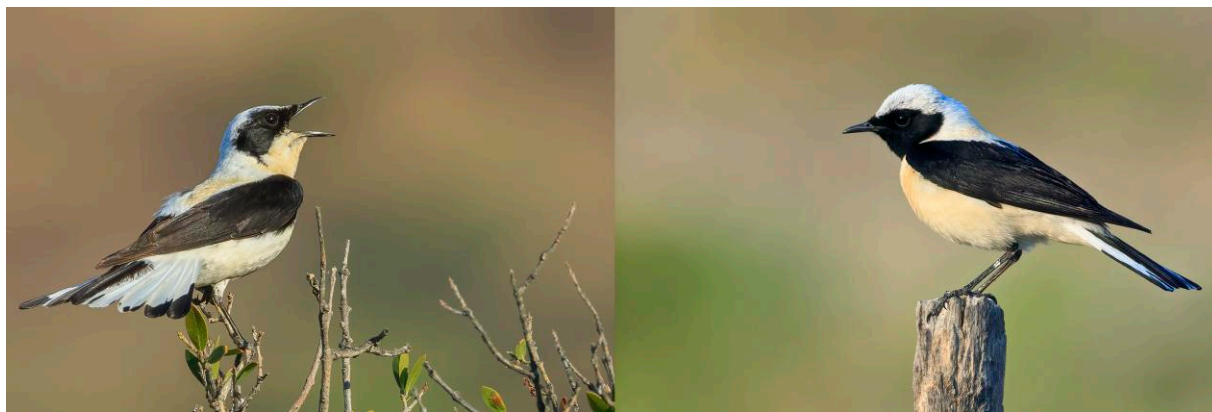


Fig. 1. Eastern black-eared wheatear (*O. melanoleuca*). The species sports a white-throated (left; Agii Pantess, Greece, June 2022) and a black-throated phenotype (right; Lesvos, Greece, May 2017) in males. © Reto Burri.

Here, we describe the de novo assembly and annotation of a chromosome-scale reference genome for the Eastern black-eared wheatear (*O. melanoleuca*). The assembly includes models for 32 autosomes, the Z chromosome, and the mitogenome that together cover 90% of the k-mer-based genome size estimate (94% with unplaced scaffolds included); it is highly contiguous with a scaffold N50 of 70 Mb and benchmarking universal single-copy orthologs (BUSCO) completeness score of 95.5%. This reference genome enables genomic research into the evolutionary history of phenotypic and species diversification in wheatears and their close relatives.

Material and methods

Sampling, tissue preservation, and nucleic acid extraction

To obtain optimal starting material for a reference individual, we freshly sampled a male Eastern black-eared wheatear (*O. melanoleuca*) well outside known hybrid zones (Haffer 1977; Panov 2005) in Galaxidi, Greece (sampling permit no. 181968/989, issued by the Ministry of Environment and Energy, General Secretariat of Environment, General Directorate of Forests and Forest Environment, Directorate of Forest Management, and Department of Wildlife and Game Management; export permit no. 55980/1575, Regional CITES Management Authority Attika). For this purpose, we sampled about 100 μ L of blood from the brachial vein, and, after euthanizing the bird, we extracted all tissues possible. Tissues were immediately snap-frozen in liquid nitrogen. Throughout transportation and storage preceding DNA extraction, the samples were kept at a temperature below -80°C .

To obtain ultra-high molecular weight (UHMW) DNA from the reference individual, NGI Uppsala (Sweden) extracted DNA from the blood sample using the Bionano Prep Blood and Cell Culture DNA Isolation Kit (Bionano, San Diego, USA). Electrophoresis on a Femto Pulse instrument showed a mean DNA fragment length of about 200 kb, with fragments reaching up to 800 kb.

To prepare muscle tissue for Hi-C sequencing library preparation, we pulverized breast muscle tissue from the reference individual in a mortar. To avoid unfreezing of the tissue powder, the procedure was carried out in a climate chamber at 4°C under regular addition of liquid nitrogen.

To prepare RNA for full-length transcript sequencing, we extracted total RNA from 8 snap-frozen tissues kept at -80°C (brain, breast muscle, heart, kidney, liver, lung, spleen, and testis) using

the RNeasy Mini Kit (Qiagen; Hombrechtikon, Switzerland) according to the manufacturer's instructions. RNA quality was assessed with a Fragment Analyzer (Agilent). RNA from spleen showed considerable degradation and was excluded from further analyses.

De novo genome sequencing and reference genome assembly and annotation

Assembly strategy and data acquisition

To obtain a chromosome-scale reference genome, our strategy largely followed the multiplatform approach recommended by Peona et al. (2021a). In brief, it consisted of (1) a phased primary assembly based on long reads, (2) polishing and scaffolding of the primary assembly with linked-read sequencing data, and (3) scaffolding of the secondary assembly with proximity ligation (Hi-C) information.

To this end, we obtained a total of 215-Gb (unique coverage 151 Gb) Pacific Biosciences (PacBio) long-read sequence data, 54-Gb linked-read sequence data, and 83-Gb Hi-C data. NGI Uppsala (Sweden) prepared a PacBio library from UHMW DNA using the SMRTbell Template Prep Kit 1.0 and sequenced this library on 18 SMRT Cells 1M v3 on a PacBio Sequel instrument (Sequel Binding Kit 3.0, Sequel Sequencing Plate 3.0). PacBio long-read data was initially processed using SMRT Link v6. A linked-read sequencing library was prepared using the 10x Genomics Chromium Genomic Kit (from the same DNA extraction as used for PacBio sequencing; 10x Genomics, Inc., Pleasanton, CA, USA; Cat No. 120215), and a Hi-C library was prepared with the Dovetail Omni-C kit (Scotts Valley, CA, USA; Cat No. 21005). The linked-read and Hi-C libraries were prepared and sequenced on a NovaSeq 6000 instrument (S4 lane, 150-bp paired-end reads) at the facilities of NGI Stockholm (Sweden).

Genome size estimation

We estimated genome size by counting k-mer frequency of the quality-checked 10x Genomics linked reads. To this end, we first trimmed 22 bp from all 10x Genomics linked reads using fastp (Chen et al. 2018) to remove indices from R1 reads and keep symmetric read lengths for the R2 reads. We then counted k-mers of size 21 using jellyfish 2.2.10 (Marçais and Kingsford 2011) and used GenomeScope (Vurtture et al. 2017) to estimate genome size from k-mer count histograms.

De novo genome assembly

We assembled the PacBio long reads into the phased primary assembly using the Falcon Unzip 0.5 assembler (Chin et al. 2016), followed by polishing with Arrow 1.9.0. Before assembly polishing, we masked repeat regions of the phased primary assembly with RepeatMasker 4.1.0 (Smit et al. 1996–2010) using a custom repeat library (Suh et al. 2018; Boman et al. 2019; Weissensteiner et al. 2020; Peona et al. 2021a, 2021b) to make accurate assembly corrections without overcorrecting large repeats. We then polished the masked assembly with 2 rounds of Pilon v1.22 (Walker et al. 2014) with the parameter “–fix indels” using the reference individual’s linked-read data. To purge duplicate scaffolds from the assembly, we ran `purge_dups` 1.2.5 (Guan et al. 2020) on the polished assembly. Prior to scaffolding with linked-read data, we split potential mis-assemblies with reference–individual linked-read data using Tigmint 1.2.4 (Jackman et al. 2018). With the aim to scaffold the polished remaining contigs, we applied ARCS 1.2.2 and LINKS 2.0.0 using the reference individual’s linked-read data using default parameters (Warren et al. 2015; Yeo et al. 2018).

To further scaffold the assembly, we applied the 3D DNA pipeline (Dudchenko et al. 2017) to join the sequences into chromosomes. We first used Juicer v.1.6 (Durand et al. 2016) to map Hi-C data against the contigs and to filter reads and then ran the `asm-pipeline` v.180922 to generate a draft scaffolding.

Finally, we corrected mis-assemblies based on the visual inspection of the proximity map using Juicebox 2.13.06 (Robinson et al. 2018). The final chromosome-level assembly was polished with 2 additional rounds of Pilon as described above.

To assess homology of the assembled scaffolds with bird chromosomes, we aligned the final genome assembly to the genomes of collared flycatcher (*Ficedula albicollis*) (FicAlb1.5) (Kawakami et al. 2014a, 2014b), zebra finch (*taeGut3.2.4*) (Warren et al. 2010), and chicken (GRCg6a) (Bellott et al. 2017) using D-GENIES (Cabanettes and Klopp 2018). Chromosomes were named according to homology with these 3 genomes. In cases, such as chicken chromosomes 1 and 4 that are split to multiple chromosomes in songbirds, the nomenclature in the wheatear genome was adapted to the species whose homologous chromosome matched closest.

Mitogenome assembly

To assemble the mitochondrial genome, we used the MitoFinder 1.4 (Allio et al. 2020) and mitoVGP 2.2 (Formenti et al. 2021) pipelines with the published *Oenanthe isabellina* mitochondrial genome (GenBank accession number: NC_040290.1) as reference. We ran MitoFinder with the reference individual’s short-read data (linked-read data but without making use of the linked-read haplotype information), and with mitoVGP, we made joint use of the linked-read and long-read data. From MitoFinder, we extracted the longest contig containing all 13 protein-coding genes, 2 rRNA genes, and 22 tRNAs annotated by MitoFinder as mitogenome assembly. We annotated both assemblies using the MITOS WebServer (<http://mitos2.bioinf.uni-leipzig.de/index.py>).

We then aligned both resulting assemblies with the mitogenomes of isabelline wheatear (*O. isabellina*, NC_040290.1) and northern wheatear (*O. oenanthe*, MN356231.1) using MUSCLE (Edgar 2004) in MEGA X (Stecher et al. 2020) and generated a circular mitogenome map using CGView (Stothard and Wishart 2005).

Assembly quality evaluation

To evaluate assembly quality at each assembly step, we estimated basic assembly statistics using QUAST 5.0.2 (Gurevich et al. 2013)

and evaluated the completeness of expected gene content in the assembly based on BUSCO (Simão et al. 2015) with the avian data set `aves_odb10` (8,338 BUSCO) in BUSCO 5.0.0.

Repeat annotation

The final version of the genome assembly was used to de novo characterize both interspersed and tandem repeats. For interspersed repeats, we used RepeatModeler2 (Flynn et al. 2020) with the option “-LTR_struct” to obtain an improved characterization of LTR retrotransposons which are commonly found in avian genomes (Kapusta and Suh 2017; Boman et al. 2019; Peona et al. 2021a, 2021b). The resulting library of raw consensus sequences was filtered from consensus sequences of tandem repeats (for which we ran a specific analysis; see below) and from protein-coding genes using the Snakemake pipeline `repeatlib_filtering_workflow` v0.1.0 (https://github.com/NBISweden/repeatlib_filtering_workflow).

For tandem repeats, we used RepeatExplorer2 (Novák et al. 2020) to search for satellite DNA (satDNA) sequences using the reference individual’s 10x Genomics linked reads. Prior to RepeatExplorer2 graph-based clustering analysis, sequencing reads were preprocessed and checked by quality with FastQC (Babraham Bioinformatics: Cambridge 2012) using the public online platform at <https://galaxy-elixir.cerit-sc.cz/>. We processed the reads with the “quality trimming”, “FASTQ interlacer on the paired end reads,” and “FASTQ to FASTA converter”, followed by “RepeatExplorer2 clustering” tools with default parameters. Each reference sequence assembled by RepeatExplorer2 consisted of a monomer of the satDNA consensus sequence. The relative genomic abundance and nucleotide divergence (Kimura 2-parameter distance) of each detected satDNA were estimated by sampling 4 million read pairs and aligning them to the satDNA library with RepeatMasker 4.1.0 (Smit et al. 1996–2010). The sampled reads were mapped to dimers of satDNA consensus sequences, and for smaller satDNAs, several monomers were concatenated until reaching roughly 150-bp array length. The resulting `RepeatMasker.align` file was then parsed to the script `calcDivergenceFromAlign.pl` from RepeatMasker utils. The relative abundance of each satDNA sequence was then estimated as the proportion of nucleotides aligned with the reference sequence with respect to the total Illumina library size.

The RepeatModeler2 library was then merged with the satDNA library produced here and with known avian consensus sequences of transposable elements (TEs) from Repbase (Bao et al. 2015), Dfam (Storer et al. 2021, 2021), flycatcher, blue-capped cordon-bleu, hooded crow, and paradise crow (Suh et al. 2018; Boman et al. 2019; Weissensteiner et al. 2020; Peona et al. 2021a, 2021b). This library was then used to annotate the genome assembly with RepeatMasker (Smit et al. 1996–2010). The annotation produced was processed with the script `calcDivergenceFromAlign.pl` from RepeatMasker utils to calculate the divergence between repeats and their consensus sequences using the Kimura 2-parameter distance corrected for the presence of CpG sites.

Full-length transcript sequencing and genome annotation

We aimed to establish a high-quality genome annotation based on full-length transcripts. To this end, for each of the abovementioned 7 tissues, the NGS platform of the University of Bern, Switzerland, prepared an Iso-Seq library using the SMRTbell Express Template Prep Kit 2.0 (Pacific Biosciences). These 7 libraries were then sequenced on 3 separate SMRT cells 8M, sequencing twice 5 tissues (brain and testis, lung, muscle, and heart) and once 2 tissues (liver and kidney) per SMRT cell. Sequencing of these SMRT cells was conducted on a Pacific Biosciences Sequel II instrument at the Genomic Technologies Facility in Lausanne,

Table 1. Assembly statistics for different versions of the *O. melanoleuca* genome.

		Falcon unzip, Arrow	+ Pilon, purge_dups	+ Tigmint	+ 3D DNA (all)	+ 3D DNA (chrom)
Basic stats	No. contigs/scaffolds ^a	1,681	381	383	588 ^a	32 ^a
	No. contigs/scaffolds ^a >50 kb	1,610	347	348	143 ^a	31 ^a
	Assembly length (Gb)	1.29	1.04	1.04	1.04 ^a	1.00 ^a
	Contig/scaffold ^a N50 (Mb)	8.6	13.5	12.6	69.6 ^a	69.7 ^a
	Contig/scaffold ^a L50	35	23	24	6 ^a	5 ^a
	Largest contig/scaffold ^a (Mb)	45.3	45.3	45.3	148.4 ^a	148.4 ^a
BUSCO	Complete (%)	96.9	96.4	96.4	96.2	95.5
	Complete single-copy (%)	90.6	95.9	95.9	95.7	95.1
	Complete duplicated (%)	6.3	0.5	0.5	0.5	0.4
	Fragmented (%)	0.7	0.7	0.7	0.9	0.9
	Missing (%)	2.4	2.9	2.9	2.9	3.6

^a Numbers concerning scaffolds instead of contigs.

Switzerland. As the libraries underloaded, 5 libraries (all but liver and kidney) were jointly sequenced on an additional SMRT cell 8M on a Pacific Biosciences Sequel IIe at the NGS platform of the University of Bern.

Circular consensus sequences (CCS), full-length nonchimeric transcripts, and polished high- and low-quality transcripts were obtained by the NGS platform at the University of Bern separately for each run using the IsoSeq 3 pipeline (ICS v10.1). Polished full-length isoforms for each sequencing run were merged by tissue and then separately mapped to the reference genome using Minimap v2.2 (-ax splice) (Li 2018, 2021). Transcriptome annotations were generated by first collapsing redundant transcripts using TAMA collapse (-x no_cap), before generating open reading frame (ORF) and nonsense-mediated mRNA decay (NMD) predictions using the scripts implemented in TAMA-GO (Kuo et al. 2020) for each of the 7 tissues. We then evaluated tissue-specific transcriptome completeness using BUSCO (Simão et al. 2015) with the avian data set *aves_odb10* (8,338 BUSCO) in BUSCO 5.0.0. Additional transcriptome annotation statistics were obtained using the *agat_sp_statistics.pl* script implemented in the AGAT toolkit (Dainat 2019).

We annotated the repeat soft-masked genome using GeMoMa 1.9 (Keilwagen et al. 2018; Keilwagen et al. 2019), a homology-based gene prediction tool. This tool is based on the annotation of protein-coding genes and intron position conservation in a reference genome to predict the annotation of protein-coding genes in the target genome. We used the genomes of chicken (GCA_016699485.1; International Chicken Genome Sequencing Consortium 2004), zebra finch (GCA_003957565.2; Warren et al. 2010), silvereye (GCA_001281735.1; Cornetti et al. 2015), and collared flycatcher (GCA_000247815.2; Ellegren et al. 2012; Kawakami et al. 2014a, 2014b) as references for the homology-based gene prediction, along with the reference individual's transcriptome obtained from Iso-Seq data to incorporate RNA evidence for the splice prediction. Using the Extract RNA-seq Evidence tool implemented in GeMoMa, we obtained intron position and coverage. This information was fed into the GeMoMa pipeline (GeMoMa.m = 200,000, AnnotationFinalizer.r = SIMPLE, pc = true, and o = true) to obtain predicted protein-coding gene models. To account for redundancies/duplicates resulting from the predicted protein-coding genes potentially stemming from each of the 4 reference species, genome annotation completeness was assessed by recomputing BUSCO using the BUSCOrecomputer tool in GeMoMa.

Functional annotation of protein-coding genes was obtained with InterProScan 5.59 (Jones et al. 2014; Paysan-Lafosse et al. 2022). InterProScan ran with the following settings: *-goterms -iprlookup*

-appl CDD, COILS, Gene3D, HAMAP, MobiDBLite, PANTHER, Pfam, PIRSF, PRINTS, PROSITEPATTERNS, PROSITEPROFILES, SFLD, SMART, SUPERFAMILY, and TIGRFAM. Predicted protein-coding genes were further annotated through a protein Blast search (-evalue 0.000001, -seg yes, -soft_masking true, and -lcase_masking) against the Swiss-Prot database (Uniprot Consortium 2019). We then merged the predicted protein-coding gene models and the functional annotation using the *agat_sp_manage_functional_annotation.pl* script, obtained summary statistics using *agat_sp_statistics.pl* and *agat_sp_functional_statistics.pl*, both implemented in the AGAT toolkit. Gene ontology (GO terms) were visualized with WEGO 2.0 (wego.genomics.cn).

Results and discussion

Nuclear genome assembly

The polished, unzipped primary assembly contained a total of 1,681 contigs, of which all were >25 kb long and 1,610 were >50 kb long (Table 1). Total assembly length was 1.29 Gb, with the longest contig spanning 45.3 Mb, contig N50 of 8.6 Mb, and half of the assembly placed in 35 contigs. Avian BUSCO were 96.9% complete, with 90.6% being single-copy genes (Table 1).

Purging duplicated contigs resulted in an assembly comprised of 381 contigs with a total assembly length of 1.04 Gb, contig N50 of 13.5 Mb, and half of the assembly placed in 23 contigs (Table 1). After this step, BUSCO completeness remained at 96.4%, but an improvement to nearly 96% single-copy BUSCOs was achieved (Table 1).

Starting from an already highly contiguous assembly, the linked-read data did not yield any scaffolding improvement. Still, Tigmint detected several supposed mis-assemblies and split the assembly into 451 scaffolds. However, an alignment of the original contigs in D-GENIES (Cabanettes and Klopp 2018) showed that all but one of the original contigs (see below) were collinear with the collared flycatcher genome. Given this result and that the proximity ligation data would correct mis-assemblies in subsequent steps, we decided to keep the original contigs except for one aligning to flycatcher chromosomes 2 and 3. For the latter contig, we used the output of Tigmint that split the contig in line with the alignment. The 2 split parts covered all but 12,527 bp of the original contig. Visual inspection of the missing sequence showed that it almost entirely consisted of repeats. We left this sequence in the assembly as a separate contig.

The proximity ligation information obtained through Hi-C scaffolding corrected a number of scaffolds, resulting in a higher number of scaffolds (588) than the number of contigs it started from

Table 2. Comparison of genome assembly and annotation summary statistics of *O. melanoleuca* with other songbird species (*J. hyemalis*, *F. coelebs*, *M. melodia*, *T. guttata*, *F. albicollis*, *M. vitellinus*, and *G. fortis*). Modified from Friis et al. (2022).

	<i>Oenanthe</i>	<i>Junco</i>	<i>Fringilla</i>	<i>Melospiza</i>	<i>Taeniopygia</i>	<i>Ficedula</i>	<i>Manacus</i>	<i>Geospiza</i>
Genome assembly length (Gb)	1.04	0.99	0.99	1.36	1.22	1.1	1.17	1.04
Genome contig N50 (kb)	7,700	75	67	8,300	38	410	194	30
Genome BUSCO scores (%)	C	95.5	95.4	94.1	87.9	93.8	96.5	96.0
	S	95.1	95.2	93.8	87.3	91.9	96	95.6
	D	0.4	0.2	0.3	0.6	1.9	0.5	1.5
	F	0.9	1.6	2.0	7.2	2.3	0.8	1
	M	3.6	3.0	4.0	5.0	3.9	2.7	2.9
No. of genes	18,143	19,026	17,703	15,086	17,561	16,763	18,976	14,399
Mean gene length (bp)	28,23,218	15,402	15,818	14,457	26,458	31,394	27,847	30,164
No. of CDS	31,333	23,245	17,703	15,086	17,561	16,763	18,976	14,399
Mean CDS length (bp)	1,682	1,647	1,679	1,325	1,677	1,942	1,929	1,766
No. of exons	320,754	229,210	221,872	131,940	171,767	189,043	190,390	164,721
Mean exon length (bp)	164	167	165	153	255	253	264	195
Mean no. exons/gene	102	9.9	10.2	8.7	10.3	12.2	11.5	11.4
No. of introns	289,421	205,965	200,041	116,724	153,909	171,236	171,089	149,563

BUSCO parameters are C, complete genes; S, complete and single-copy genes; D, complete and duplicated genes; F, fragmented genes; M, missing genes.

(383). However, the scaffolding yielded a highly contiguous chromosome-scale assembly (N50, 69.6 Mb; L50, 6) with BUSCO completeness of still >96% and almost all BUSCOs in single copy (Table 1). This final assembly contained all macrochromosomes and the majority of microchromosomes usually found in the latest generation of avian genome assemblies (Kapusta et al. 2017; Rhie et al. 2021; Peona et al. 2021a, 2021b). A total of 96% of the assembly was placed into chromosome models, and the chromosome-only assembly covered still 95.5% of BUSCO (Table 1).

The final assembly length closely matched the one of previous linked-read-based assemblies of the same species and closely related ones (Schweizer et al. 2019a, 2019b; Lutgen et al. 2020). The genome size estimated from the k-mer distribution of linked-read sequence was between 1.105 and 1.106 Gb, with 0.925–0.926 Gb of unique and 0.179–0.180 Gb (16%) repeat sequence and 0.75–0.76% heterozygosity (GenomeScope model fit 98–99%). The full final reference genome assembly thus covered 94% of the genome size estimate, with 90% of the estimated genome size placed in chromosomes. A total of 96% of the assembly were placed in 33 chromosomes with homologs in collared flycatcher, zebra finch, and chicken, according to which we adapted the chromosome nomenclature. The differences in genome size estimates based on the k-mer approach and the genome assembly length are likely the result of highly repetitive sequences (e.g. centromeres, telomeres, and satDNAs) that collapsed during the assembly process (Peona et al. 2018). Assembly contiguity and completeness (as judged by BUSCO scores) of the *O. melanoleuca* assembly compared favorably with other songbird genome assemblies (Table 2).

Mitogenome assembly

MitoFinder and MitoVGP assembled mitogenomes of 16,944 bp and 18,631 bp length, respectively. The mitochondrial contigs assembled by the 2 pipelines were congruent, except for 9 single base pair mismatches, for a 1,827-bp-long insert in the MitoVGP assembly and of a 141-bp-long insert in the MitoFinder assembly. We decided to not consider either of these inserts in the final mitogenome assembly for the following reasons. First, neither of the inserts was observed in the mitogenomes of isabelline and northern wheatear. For the long insert in the MitoVGP assembly, moreover, the coverage of short reads mapped to the MitoVGP assembly was strongly reduced (Supplementary Fig. 1), and the insertion constituted a partial duplication of *nd6*, duplications of 2 tRNAs (Glu and Pro), and a partial duplication of the control region

likely caused by an assembly artifact. The short insert in the MitoFinder assembly was not observed in the other wheatear mitogenomes, and if real, we would expect long reads to cover this insert. Because base calling based on short reads is expected to have higher quality, we retained the MitoFinder assembly, but without the 141-bp insert as final mitogenome.

The final mitogenome (as also both original assemblies) contained all 13 protein-coding genes, 2 rRNAs, and 22 tRNAs (Fig. 2). All genes, except 8 tRNAs and *nd6*, were located on the heavy DNA strand. Both gene order and strandedness were concordant with those observed in northern wheatear (*O. oenanthe*) (Wang et al. 2020).

Repetitive element annotation

The de novo identification of repetitive elements resulted in the characterization of 572 raw consensus sequences from RepeatModeler2 and 16 satellite DNA consensus sequences from RepeatExplorer2. The consensus sequences from RepeatModeler2 were filtered from tandem repeats and protein-coding genes. This resulted in a final library of 477 consensus sequences (Supplementary File 1). Among these consensus sequences, RepeatModeler2 classified 226 sequences as LTR retrotransposons, 98 as LINE retrotransposons, 21 as DNA transposons, and 5 as SINE retrotransposons, and 112 sequences were unclassified (“unknown”).

The genome assembly annotation run with RepeatMasker using the repeat library produced here and merged with already known avian repeats showed that ~10% of the assembled genome is repetitive (Fig. 3a and Supplementary Table 1 and File 2). This finding indicates that many repeats collapsed during the genome assembly process. An example of this were satDNAs that represented ~0.8% of the sequenced reads but only <0.3% of the genome assembly, suggesting that satDNA repeats [such as in (peri-)centromeric and (sub-)telomeric regions] are the most collapsed repeats. Most of the repeats annotated were LTR and LINE retrotransposons (Fig. 3a). While it is common to find LINES as most abundant TEs in avian genomes (Kapusta and Suh 2017; Manthey et al. 2018; Galbraith et al. 2021; Peona, Blom et al. 2021a), it is less common to find so similar percentages of LINE and LTR retrotransposons. This is especially true for a male genome assembly such as the present one here that does not include the W chromosome which is highly enriched in LTRs and acts as a refugium for most of the full-length genomic LTR elements in birds (Peona et al. 2021a, 2021b; Warmuth et al. 2022). The TE

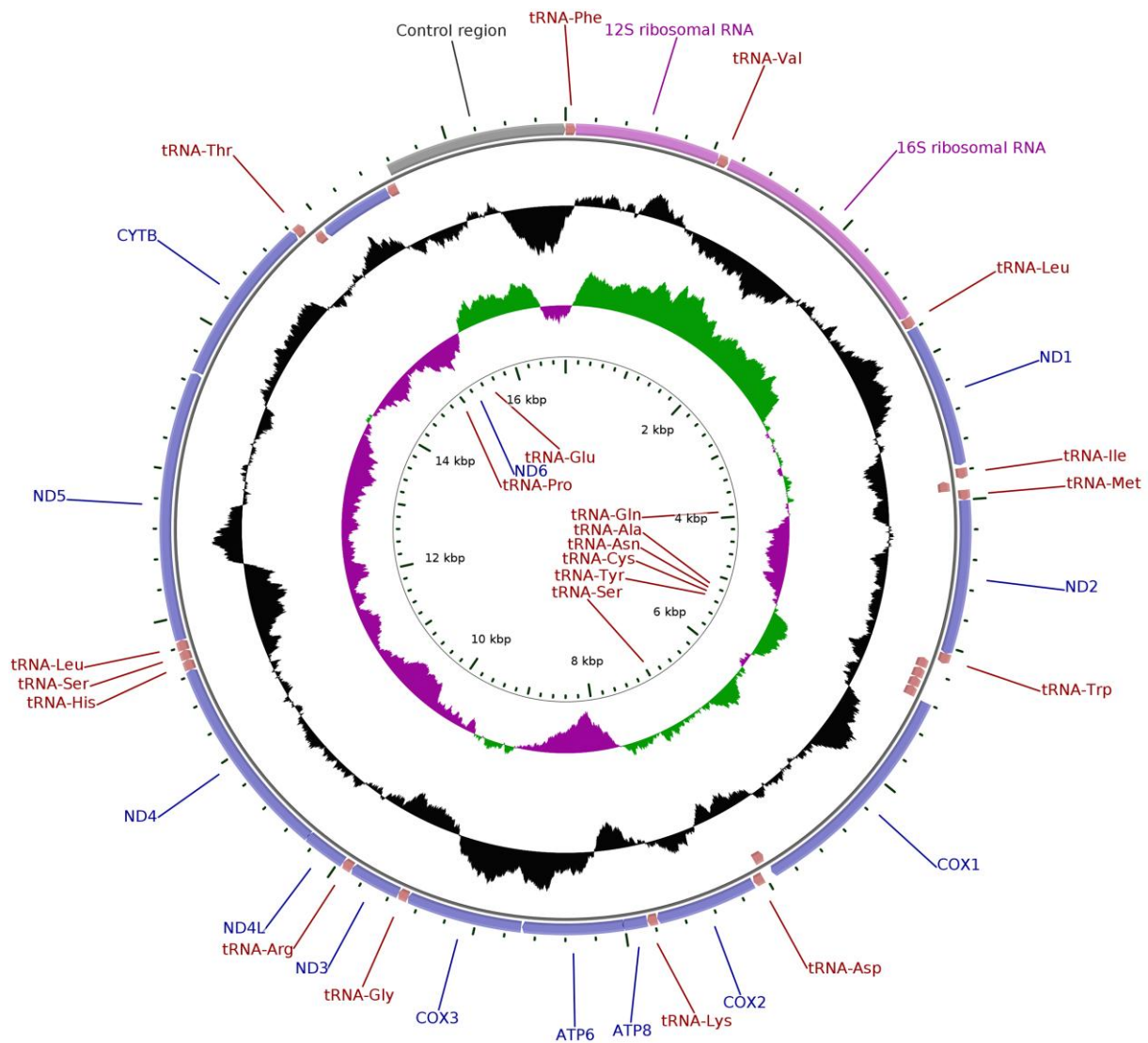


Fig. 2. Circular sketch map of the *O. melanoleuca* mitogenome assembly. The outer circle shows coding sequences, rRNAs, and tRNAs. The black track on the middle circle indicates GC content. On the inner circle, positive and negative GC skews in nucleotide composition are indicated.

landscape (Fig. 3b) suggests that LINE retrotransposons experienced a drop in their genomic accumulation in recent times (0–5% divergence; Fig. 3b), whereas LTR retrotransposons kept accumulating at the same rate. Such a recent replacement of LINE retrotransposon activity with a diversity of LTR retrotransposons has been noted in other songbirds and seems to have occurred independently in the so far analyzed passerine families, i.e. estrildid finches (Warren et al. 2010, Boman et al. 2019), flycatchers (Suh et al. 2018), crows (Weissensteiner et al. 2020), and birds-of-paradise (Peona et al. 2021a, 2021b). Finally, the satDNA landscape (Fig. 3b) shows that satDNA arrays experienced differential amplification in copy numbers in recent times (0–10% divergence), implying fast evolution of this genomic fraction in the genome (Peona et al. 2022).

Transcriptome sequencing, genome annotation, and gene function prediction

Iso-Seq sequencing yielded a total of 4,627,382 CCS reads (125,633–1,087,892 reads per tissue, Table 3). This resulted in numbers of high-quality isoforms ranging from 16,078 to 80,600

per tissue. On average 8,833 genes were predicted per tissue, ranging from 4,772 in muscle to 10,924 in the liver. Transcriptome completeness evaluated through BUSCO ranged from 31.2 to 57.5% complete BUSCO per tissue (Table 3).

The Iso-Seq transcriptomes were then used as splice evidence in GeMoMa to perform a predominantly homology-based annotation of the reference genome. We predicted 18,143 protein-coding genes with a total of 320,754 exons and 289,421 introns. The number of exons, CDS, and introns was higher for our *O. melanoleuca* annotation compared with the annotations of other songbirds, such as *Junco hyemalis*, *Fringilla coelebs*, *Melospiza melodia*, *Taeniopygia guttata*, *Ficedula albicollis*, *Manacus vitellinus*, and *Geospiza fortis* (Table 2). Mean gene length, CDS length, exon length, and number of exons per gene, on the other hand, were in the range of values obtained for the abovementioned songbird annotations (Table 3). Of the 18,143 predicted genes, 17,553 (96.7%) were annotated with protein families or function assignment, and 12,472 (68.7%) genes obtained a GO term assignment through InterProScan. The most abundant GO terms were associated with “cell part,” “cell” and “membrane” in the cellular

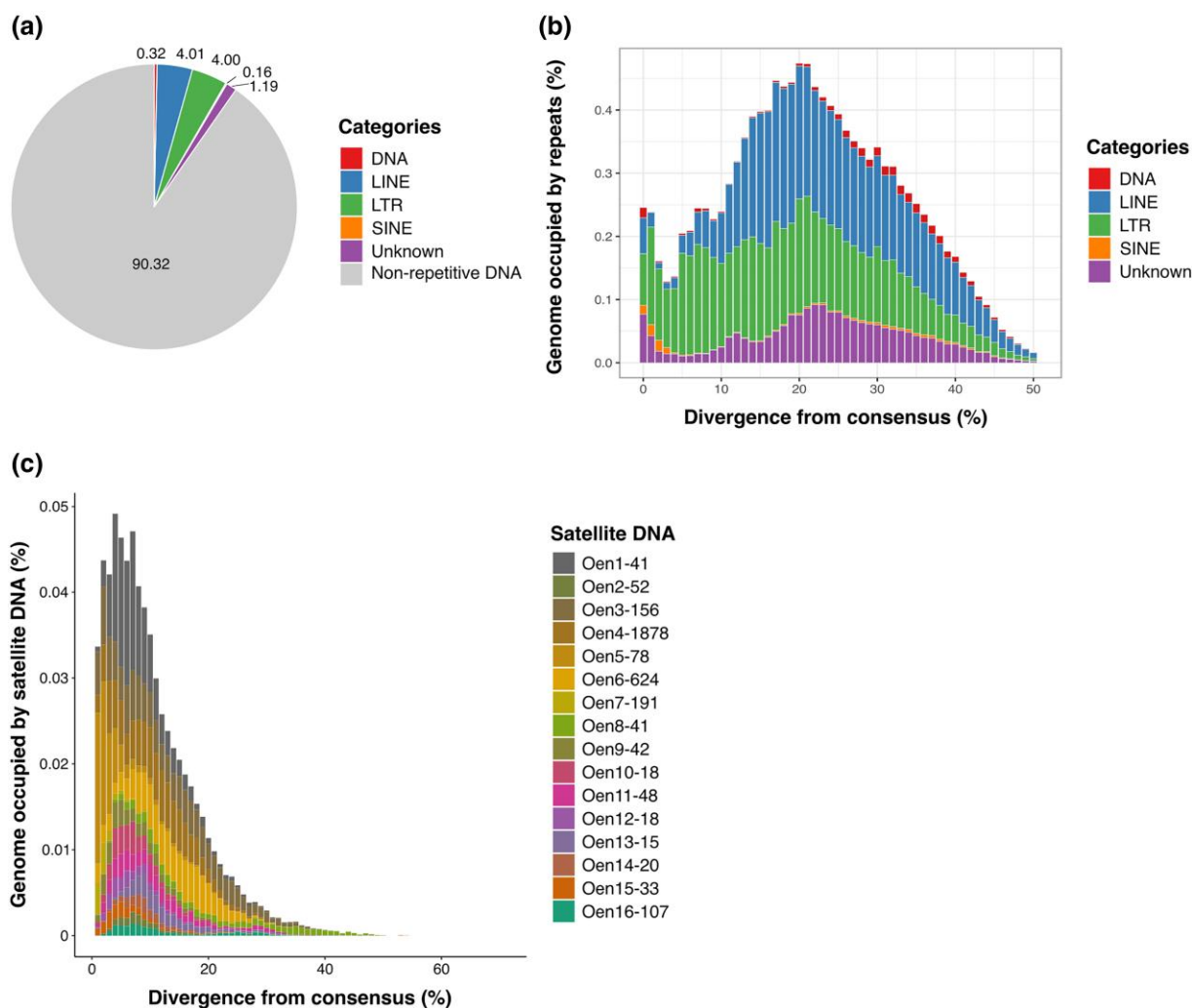


Fig. 3. Repeat annotation landscapes. a) Pie chart summarizing the TE content annotated in the genome assembly. b) TE landscape. The divergence between interspersed repeat copies and their consensus sequences is shown on the X-axis as genetic distance calculated using the Kimura 2-parameter distance. The percentage of the genome assembly occupied by transposable elements is shown on the Y-axis. c) Satellite DNA landscape. The divergence between the satellite DNA consensus sequences and sequences annotated in the short-read library is shown on the X-axis as genetic distance calculated using the Kimura 2-parameter distance. The percentage of the genome (short reads) annotated as satellite DNA is shown on the Y-axis.

Table 3. Iso-Seq data characterization and transcriptome completeness.

		Brain	Heart	Kidney	Liver	Lung	Muscle	Testis
Transcriptome	No. of CCS reads	847,617	253,468	723,158	1,087,892	1,061,936	125,633	527,678
	High-quality isoforms	73,422	80,600	45,097	47,491	28,508	16,078	44,605
	Low-quality isoforms	734	844	616	384	151	94	284
	No. of genes	10,449	10,448	9,063	10,924	6,564	4,772	9,613
	Mean gene length (bp)	24,193	20,119	16,350	15,125	18,528	17,397	17,415
	No. of CDS	27,449	28,747	25,790	27,202	13,551	8,447	23,009
	Mean CDS length (bp)	972	985	932	823	894	980	960
	No. of exons	231,169	222,791	235,989	194,325	108,084	69,859	184,794
	Mean exon length (bp)	246	248	223	225	221	224	209
BUSCO	Mean no. of exons/mRNA	8.4	8.2	7.9	7.1	8.0	8.3	8.0
	Complete (%)	56.80	57.50	48.30	49.40	38.30	31.20	49.3
	Single-copy (%)	40.30	39.50	33.60	34.70	31.1	27.00	34.6
	Duplicated (%)	16.50	18.00	14.70	14.70	7.20	4.20	14.70
	Fragmented (%)	2.90	2.10	2.60	3.20	2.00	1.10	2.30
Missing (%)	40.30	40.60	49.10	47.40	59.70	67.70	48.40	

component category, “binding” in the molecular function category, and “cellular metabolic process” or “metabolic process” in the biological process category (Supplementary Fig. 2). BUSCO completeness of the final annotation as judged from avian BUSCO ($n=8,338$) was 98.0%, with 97.4% single-copy BUSCO, 0.6% duplicated BUSCO, 0.6% fragmented BUSCO, and 1.5% missing BUSCO. This suggests an accurate and rather complete annotation.

Data availability

All data, including the assembly, its annotation, and the original sequencing, data are available on the European Nucleotide Archive under project accession PRJNA937434. Code for the repeat analysis is available on <https://github.com/ValentinaBoP/WheatearGenomeAnalysis>.

Supplemental material available at figshare: <https://doi.org/10.25387/g3.22209697>.

Acknowledgements

We warmly thank Marta Burri for preparing RNA and Giulio Formenti, Remi Allio, and Lauren Coombe for support with computational questions. We are indebted to NGI Uppsala, namely Mai-Britt Mosbech and Olga Vinnere Pettersson, for UHMW DNA extraction and long-read sequencing and Ignas Bunikis for running the primary assembly, as well as NGI Stockholm for the preparation of linked-read and Hi-C sequencing data. Finally, we thank the NGS platform at the University of Bern, namely Pamela Nicholson and Catia Coito, for the preparation of Iso-Seq data. Computations were performed on resources provided by the Swedish National Infrastructure for Computing (SNIC) through the Uppsala Multidisciplinary Center for Advanced Computational Science (UPPMAX) and the High-Performance Computing Cluster EVE, a joint effort of the Helmholtz Centre for Environmental Research (UFZ) and the German Centre for Integrative Biodiversity Research (iDiv) Halle-Jena-Leipzig. We thank the administration and support staff of EVE: Thomas Schnicke and Ben Langenberg (UFZ) and Christian Krause (iDiv).

Funding

The present project was supported by the Deutsche Forschungsgemeinschaft (DFG), grant number BU3456/3-1 to RB, and the Fonds National de la Recherche (FNR) Luxembourg, grant number 14575729, to DL; OMP-G was supported by Vetenskapsrådet (VR), grant number 2020-03866; VP was supported via grants to AS from Vetenskapsrådet (grant number 2020-04436) and Formas (2017-01597). NAK was supported by a Georg Forster Research Stipend of the Alexander von Humboldt Foundation. Part of the analyses were enabled by resources provided by the Swedish National Infrastructure for Computing (SNIC) at Uppsala partially funded by the Swedish Research Council through grant agreement no. 2018-05973.

Conflicts of interest statement

The authors declare no conflict of interest.

Literature cited

Alaei Kakhki NA, Schweizer M, Lutgen D, Bowie RCK, Shirihai H, Suh A, Schielzeth H, Burri R. A phylogenomic assessment of processes

underpinning convergent evolution in open-habitat chats. *Mol Biol Evol.* 2023;40(1):msac278. doi:10.1093/molbev/msac278.

Aliabadian M, Kaboli M, Förschler MI, Nijman V, Chamani A, Tillier A, Prodon R, Pasquet E, Ericson PGP, Zuccon D. Convergent evolution of morphological and ecological traits in the open-habitat chat complex (Aves, Muscicapidae: Saxicolinae). *Mol Phylogenet Evol.* 2012;65(1):35–45. doi:10.1016/j.ympev.2012.05.011.

Allio R, Schomaker-Bastos A, Romiguier J, Prosdociami F, Nabholz B, Delsuc F. Mitofinder: efficient automated large-scale extraction of mitogenomic data in target enrichment phylogenomics. *Mol Ecol Resour.* 2020;20(4):892–905. doi:10.1111/1755-0998.13160.

Babraham Bioinformatics: Cambridge. FastQC; version 0.10.1: a quality control tool for high throughput sequence data. Cambridge: Babraham Bioinformatics; 2012.

Bao W, Kojima KK, Kohany O. Repbase Update, a database of repetitive elements in eukaryotic genomes. *Mob DNA.* 2015;6(11). 10.1186/s13100-015-0041-9

Bellott DW, Skaletsky H, Cho T-J, Brown L, Locke D, Chen N, Galkina S, Pyntikova T, Koutseva N, Graves T, et al. Avian W and mammalian Y chromosomes convergently retained dosage-sensitive regulators. *Nat Genet.* 2017;49(3):387–394. doi:10.1038/ng.3778.

Boman J, Frankl-Vilches C, da Silva dos Santos M, de Oliveira EHC, Gahr M, Suh A. The genome of blue-capped cordon-bleu uncovers hidden diversity of LTR retrotransposons in zebra finch. *Genes (Basel).* 2019;10(4):301. doi:10.3390/genes10040301.

Buerkle CA, Lexer C. Admixture as the basis for genetic mapping. *Trends Ecol Evol.* 2008;23(12):686–694. doi:10.1016/j.tree.2008.07.008.

Cabanettes F, Klopp C. D-GENIES: dot plot large genomes in an interactive, efficient and simple way. *PeerJ* 2018;6:e4958. doi:10.7717/peerj.4958.

Chen S, Zhou Y, Chen Y, Gu J. Fastp: an ultra-fast all-in-one FASTQ preprocessor. *Bioinformatics* 2018;34(17):i884–i890. doi:10.1093/bioinformatics/bty560.

Chin C-S, Peluso P, Sedlazeck FJ, Nattestad M, Concepcion GT, Clum A, Dunn C, O'Malley R, Figueroa-Balderas R, Morales-Cruz A, et al. Phased diploid genome assembly with single-molecule real-time sequencing. *Nat Methods.* 2016;13(12):1050–1054. doi:10.1038/nmeth.4035.

Cornetti L, Valente LM, Dunning LT, Quan X, Black RA, Hébert O, Savolainen V. The genome of the “Great Speciator” provides insights into bird diversification. *Genome Biol Evol.* 2015;7(9):2680–2691. doi:10.1093/gbe/evv168.

Dainat J. AGAT: another Gff analysis toolkit to handle annotations in any GTF/GFF format. (Version v0.7.0). Zenodo. 2019. doi:10.5281/zenodo.3552717.

Dudchenko O, Batra SS, Omer AD, Nyquist SK, Hoeger M, Durand NC, Shamim MS, Machol I, Lander ES, Aiden AP, et al. De novo assembly of the *Aedes aegypti* genome using Hi-C yields chromosome-length scaffolds. *Science* 2017;356(6333):92–95. doi:10.1126/science.aal3327.

Durand NC, Shamim MS, Machol I, Rao SSP, Huntley MH, Lander ES, Aiden EL. Juicer provides a one-click system for analyzing loop-resolution Hi-C experiments. *Cell Syst.* 2016;3(1):95–98. doi:10.1016/j.cels.2016.07.002.

Edgar RC. MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Res.* 2004;32(5):1792–1797. doi:10.1093/nar/gkh340.

Ellegren H, Smeds L, Burri R, Olason PI, Backström N, Kawakami T, Kunstner A, Makinen H, Nadachowska-Brzyska K, Qvarnstrom A, et al. The genomic landscape of species divergence in *Ficedula* flycatchers. *Nature* 2012;491(7426):756–760. doi:10.1038/nature11584.

- Flynn JM, Hubley R, Goubert C, Rosen J, Clark AG, Feschotte C, Smit AF. Repeatmodeler2 for automated genomic discovery of transposable element families. *Proc Natl Acad Sci U S A*. 2020;117(17):9451–9457. doi:10.1073/pnas.1921046117.
- Formenti G, Rhie A, Balacco J, Haase B, Mountcastle J, Fedrigo O, Brown S, Capodiferro MR, Al-Ajli FO, Ambrosini R, et al. Complete vertebrate mitogenomes reveal widespread repeats and gene duplications. *Genome Biol*. 2021;22(1):120. doi:10.1186/s13059-021-02336-9.
- Friis G, Vizueta J, Ketterson ED, Milá B. A high-quality genome assembly and annotation of the dark-eyed junco *Junco hyemalis*, a recently diversified songbird. *G3 (Bethesda)* 2022;12(6):jkac083. doi:10.1093/g3journal/jkac083.
- Galbraith JD, Kortschak RD, Suh A, Adelson DL. Genome stability is in the eye of the beholder: CR1 retrotransposon activity varies significantly across avian diversity. *Genome Biol Evol*. 2021;13(12):evab259. doi:10.1093/gbe/evab259.
- Gompert Z, Mandeville EG, Buerkle CA. Analysis of population genomic data from hybrid zones. *Annu Rev Ecol Evol Syst*. 2017;48:207–229. doi:10.1146/annurev-ecolsys-110316-022652.
- Guan D, McCarthy SA, Wood J, Howe K, Wang Y, Durbin R. Identifying and removing haplotypic duplication in primary genome assemblies. *Bioinformatics* 2020;36(9):2896–2898. doi:10.1093/bioinformatics/btaa025.
- Gurevich A, Saveliev V, Vyahhi N, Tesler G. QUAST: quality assessment tool for genome assemblies. *Bioinformatics* 2013;29(8):1072–1075. doi:10.1093/bioinformatics/btt086.
- Haffer J. Secondary contact zones of birds in northern Iran. Bonn, Germany, Bonner Zoologische Monographien; 1977.
- International Chicken Genome Sequencing Consortium. Sequence and comparative analysis of the chicken genome provide unique perspectives on vertebrate evolution. *Nature* 2004;432(7018):695–716. doi:10.1038/nature03154.
- Jackman SD, Coombe L, Chu J, Warren RL, Vandervalk BP, Yeo S, Xue Z, Mohamadi H, Bohlmann J, Jones SJM, et al. Tigmint: correcting assembly errors using linked reads from large molecules. *BMC Bioinformatics* 2018;19(1):393. doi:10.1186/s12859-018-2425-6.
- Jones P, Binns D, Chang H-Y, Fraser M, Li W, McAnulla C, McWilliam H, Maslen J, Mitchell A, Nuka G, et al. Interproscan 5: genome-scale protein function classification. *Bioinformatics* 2014;30(9):1236–1240. doi:10.1093/bioinformatics/btu031.
- Kapusta A, Suh A. Evolution of bird genomes—a transposon’s-eye view. *Ann N Y Acad Sci*. 2017;1389(1):164–185. doi:10.1111/nyas.13295.
- Kapusta A, Suh A, Feschotte C. Dynamics of genome size evolution in birds and mammals. *Proc Natl Acad Sci U S A*. 2017;114(8):E1460–E1469.
- Kawakami T, Backström N, Burri R, Husby A, Olason P, Rice AM, Ålund M, Qvarnström A, Ellegren H. Estimation of linkage disequilibrium and interspecific gene flow in *Ficedula* flycatchers by a newly developed 50k single-nucleotide polymorphism array. *Mol Ecol Resour*. 2014a;14(6):1248–1260. doi:10.1111/1755-0998.12270.
- Kawakami T, Smeds L, Backström N, Husby A, Qvarnström A, Mugal CF, Olason P, Ellegren H. A high-density linkage map enables a second-generation collared flycatcher genome assembly and reveals the patterns of avian recombination rate variation and chromosomal evolution. *Mol Ecol*. 2014b;23(16):4035–4058. doi:10.1111/mec.12810.
- Keilwagen J, Hartung F, Grau J. Gemoma: homology-based gene prediction utilizing intron position conservation and RNA-seq data. *Methods Mol Biol*. 2019;1962:161–177. doi:10.1007/978-1-4939-9173-0_9.
- Keilwagen J, Hartung F, Paulini M, Twardziok SO, Grau J. Combining RNA-seq data and homology-based gene prediction for plants, animals and fungi. *BMC Bioinformatics* 2018;19(1):189. doi:10.1186/s12859-018-2203-5.
- Kuo RI, Cheng Y, Zhang R, Brown JWS, Smith J, Archibald AL, Burt DW. Illuminating the dark side of the human transcriptome with long read transcript sequencing. *BMC Genomics*. 2020;21(1):751. doi:10.1186/s12864-020-07123-7.
- Li H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* 2018;34(18):3094–3100. doi:10.1093/bioinformatics/bty191.
- Li H. New strategies to improve minimap2 alignment accuracy. *Bioinformatics* 2021;37(23):4572–4574. doi:10.1093/bioinformatics/btab705.
- Lutgen D, Ritter R, Olsen R-A, Schielzeth H, Gruselius J, Ewels P, García JT, Shirihai H, Schweizer M, Suh A, et al. Linked-read sequencing enables haplotype-resolved resequencing at population scale. *Mol Ecol Resour*. 2020;20(5):1311–1322. doi:10.1111/1755-0998.13192.
- Manthey JD, Moyle RG, Boissinot S. Multiple and independent phases of transposable element amplification in the genomes of piciformes (woodpeckers and allies). *Genome Biol Evol*. 2018;10(6):1445–1456. doi:10.1093/gbe/evy105.
- Marçais G, Kingsford C. A fast, lock-free approach for efficient parallel counting of occurrences of k-mers. *Bioinformatics* 2011;27(6):764–770. doi:10.1093/bioinformatics/btr011.
- Novák P, Neumann P, Macas J. Global analysis of repetitive DNA from unassembled sequence reads using RepeatExplorer2. *Nat Protoc*. 2020;15(11):3745–3776. doi:10.1038/s41596-020-0400-y.
- Panov EN. Wheatears of the palearctic. In: Wilson MG, editor. *Ecology, Behaviour and Evolution of the Genus Oenanthe*. Sofia-Moscow: Pensoft 2005. (Series Faunistica).
- Paysan-Lafosse T, Blum M, Chuguransky S, Grego T, Pinto BL, Salazar GA, Bileschi ML, Bork P, Bridge A, Colwell L, et al. Interpro in 2022. *Nucleic Acids Res*. 2022;51(D1):D418–D427. doi:10.1093/nar/gkac993.
- Peona V, Blom MPK, Xu L, Burri R, Sullivan S, Bunikis I, Liachko I, Haryoko T, Jönsson KA, Zhou Q, et al. Identifying the causes and consequences of assembly gaps using a multiplatform genome assembly of a bird-of-paradise. *Mol Ecol Resour*. 2021a;21(1):263–286. doi:10.1111/1755-0998.13252.
- Peona V, Kutschera VE, Blom MPK, Irestedt M, Suh A. Satellite DNA evolution in Corvoidea inferred from short and long reads. *Mol Ecol*. 2022;32(6):1288–1305. doi:10.1111/mec.16484.
- Peona V, Palacios-Gimenez OM, Blommaert J, Liu J, Haryoko T, Jönsson KA, Irestedt M, Zhou Q, Jern P, Suh A. The avian W chromosome is a refugium for endogenous retroviruses with likely effects on female-biased mutational load and genetic incompatibilities. *Philos Trans R Soc Lond B Biol Sci*. 2021b;376(1833):20200186. doi:10.1098/rstb.2020.0186.
- Peona V, Weissensteiner MH, Suh A. How complete are “complete” genome assemblies? : —an avian perspective. *Mol Ecol Resour*. 2018;18(6):1188–1195. doi:10.1111/1755-0998.12933.
- Rhie A, McCarthy SA, Fedrigo O, Damas J, Formenti G, Koren S, Uliano-Silva M, Chow W, Fungtammasan A, Kim J, et al. Towards complete and error-free genome assemblies of all vertebrate species. *Nature* 2021;592(7856):737–746. doi:10.1038/s41586-021-03451-0.
- Robinson JT, Turner D, Durand NC, Thorvaldsdóttir H, Mesirov JP, Aiden EL. Juicebox.js provides a cloud-based visualization system for Hi-C data. *Cell Syst*. 2018;6(2):256–258.e1. doi:10.1016/j.cels.2018.01.001.
- Schweizer M, Warmuth VM, Alaei Kakhki N, Aliabadian M, Förschler M, Shirihai H, Ewels P, Gruselius J, Olsen R-A, Schielzeth H, et al.

- Genome-wide evidence supports mitochondrial relationships and pervasive parallel phenotypic evolution in open-habitat chats. *Mol Phylogenet Evol.* 2019a;139:106568. doi:10.1016/j.ympev.2019.106568.
- Schweizer M, Warmuth V, Alaei Kakhki N, Aliabadian M, Förschler M, Shirihai H, Suh A, Burri R. Parallel plumage color evolution and pervasive hybridization in wheatears. *J Evol Biol.* 2019b; 32(1):100–110. doi:10.1111/jeb.13401
- Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM. BUSCO: assessing genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics* 2015;31(19): 3210–3212. doi:10.1093/bioinformatics/btv351.
- Smit AFA, Hubley R, Green P. RepeatMasker Open 3.3.0. 1996–2010.
- Stecher G, Tamura K, Kumar S. Molecular evolutionary genetics analysis (MEGA) for macOS. *Mol Biol Evol.* 2020;37(4):1237–1239. doi: 10.1093/molbev/msz312.
- Storer J, Hubley R, Rosen J, Wheeler TJ, Smit AF. The Dfam community resource of transposable element families, sequence models, and genome annotations. *Mob DNA.* 2021;12(1):2. doi:10.1186/s13100-020-00230-y.
- Stothard P, Wishart DS. Circular genome visualization and exploration using CGView. *Bioinformatics* 2005;21(4):537–539. doi:10.1093/bioinformatics/bti054.
- Suh A, Smeds L, Ellegren H. Abundant recent activity of retrovirus-like retrotransposons within and among flycatcher species implies a rich source of structural variation in songbird genomes. *Mol Ecol.* 2018;27(1):99–111. doi:10.1111/mec.14439.
- UniProt Consortium. Uniprot: a worldwide hub of protein knowledge. *Nucleic Acids Res.* 2019;47(D1):D506–D515. doi:10.1093/nar/gky1049.
- Vurtture GW, Sedlazeck FJ, Nattestad M, Underwood CJ, Fang H, Gurtowski J, Schatz MC. Genomescope: fast reference-free genome profiling from short reads. *Bioinformatics* 2017;33(14): 2202–2204. doi:10.1093/bioinformatics/btx153.
- Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, Cuomo CA, Zeng Q, Wortman J, Young SK, et al. Pilon: an integrated tool for comprehensive microbial variant detection and genome assembly improvement. *PLoS One* 2014;9(11):e112963. doi:10.1371/journal.pone.0112963.
- Wang E, Zhang D, Braun MS, Hotz-Wagenblatt A, Pärt T, Arlt D, Schmaljohann H, Bairlein F, Lei F, Wink M. Can mitogenomes of the northern wheatear (*Oenanthe oenanthe*) reconstruct its phylogeography and reveal the origin of migrant birds? *Sci Rep.* 2020; 10(1):9290. doi:10.1038/s41598-020-66287-0.
- Warmuth VM, Weissensteiner MH, Wolf JBW. Accumulation and ineffective silencing of transposable elements on an avian W chromosome. *Genome Res.* 2022;32(4):671–681. doi:10.1101/gr.275465.121.
- Warren WC, Clayton DF, Ellegren H, Arnold AP, Hillier LW, Kunster A, Searle S, White S, Vilella AJ, Fairly S, et al. The genome of a songbird. *Nature* 2010;464(7289):757–762. doi:10.1038/nature08819.
- Warren RL, Yang C, Vandervalk BP, Behsaz B, Lagman A, Jones SJM, Birol I. LINKS: scalable, alignment-free scaffolding of draft genomes with long reads. *GigaScience.* 2015;4(1):35. doi:10.1186/s13742-015-0076-3.
- Weissensteiner MH, Bunikis I, Catalán A, Francoijs K-J, Knief U, Heim W, Peona V, Pophaly SD, Sedlazeck FJ, Suh A, et al. Discovery and population genomics of structural variation in a songbird genus. *Nat Commun.* 2020;11(1):3403. doi:10.1038/s41467-020-17195-4.
- Yeo S, Coombe L, Warren RL, Chu J, Birol I. ARCS: scaffolding genome drafts with linked reads. *Bioinformatics* 2018;34(5):725–731. doi: 10.1093/bioinformatics/btx675.

Editor: A. Sethuraman