

1 **Transcript- and annotation-guided genome assembly of the European starling**

2

3 Katarina C. Stuart^{1†}, Richard J. Edwards^{2†}, Yuanyuan Cheng³, Wesley C. Warren⁴, David W.
4 Burt⁵, William B. Sherwin¹, Natalie R. Hofmeister^{6,7}, Scott J. Werner⁸, Gregory F. Ball⁹,
5 Melissa Bateson¹⁰, Matthew C. Brandley¹¹, Katherine L. Buchanan¹², Phillip Cassey¹³, David
6 F. Clayton¹⁴, Tim De Meyer¹⁵, Simone L. Meddle¹⁶, Lee A. Rollins¹

7 ¹ Evolution & Ecology Research Centre, School of Biological, Earth and Environmental
8 Sciences, UNSW Sydney, Sydney, New South Wales, Australia

9 ² School of Biotechnology and Biomolecular Sciences, UNSW Sydney, Sydney, New South
10 Wales, Australia

11 ³ School of Life and Environmental Sciences, The University of Sydney, Sydney, New South
12 Wales, Australia

13 ⁴ Department of Animal Sciences, Institute for Data Science and Informatics, The University
14 of Missouri, Columbia, Missouri, USA

15 ⁵ Office of the Deputy Vice-Chancellor (Research and Innovation), The University of
16 Queensland, Brisbane, Australia

17 ⁶ Department of Ecology and Evolutionary Biology, Cornell University, Ithaca, NY 14850

18 ⁷ Fuller Evolutionary Biology Program, Cornell Lab of Ornithology, Ithaca, NY 14850

19 ⁸ United States Department of Agriculture, Animal and Plant Health Inspection Service,
20 Wildlife Services, National Wildlife Research Center, Fort Collins, Colorado, USA.

21 ⁹ Department of Psychology, University of Maryland, College Park, MD 20742 USA

22 ¹⁰ Institute of Neuroscience, Newcastle University, Newcastle upon Tyne, UK

23 ¹¹ Carnegie Museum of Natural History, Pittsburgh, Pennsylvania, USA

24 ¹² School of Life and Environmental Sciences, Deakin University, Waurin Ponds, VIC, 3228,
25 Australia

26 ¹³ Invasion Science & Wildlife Ecology Lab, University of Adelaide, Adelaide SA 5005,
27 Australia

28 ¹⁴ Department of Genetics & Biochemistry, Clemson University, South Carolina 29634

29 ¹⁵ Dept. of Data Analysis & Mathematical Modelling, Faculty of Bioscience Engineering,
30 Ghent University, Ghent, Belgium

31 ¹⁶ The Roslin Institute, The Royal (Dick) School of Veterinary Studies, The University of
32 Edinburgh, Midlothian, EH25 9RG, UK

33 † Joint first authors

34

35 **Abstract (250 words)**

36 The European starling, *Sturnus vulgaris*, is an ecologically significant, globally invasive
37 avian species that is also suffering from a major decline in its native range. Here, we present
38 the genome assembly and long-read transcriptome of an Australian-sourced European starling
39 (*S. vulgaris* vAU), and a second North American genome (*S. vulgaris* vNA), as
40 complementary reference genomes for population genetic and evolutionary characterisation.
41 *S. vulgaris* vAU combined 10x Genomics linked-reads, low-coverage Nanopore sequencing,
42 and PacBio Iso-Seq full-length transcript scaffolding to generate a 1050 Mb assembly on
43 1,628 scaffolds (72.5 Mb scaffold N50). Species-specific transcript mapping and gene
44 annotation revealed high structural and functional completeness (94.6% BUSCO
45 completeness). Further scaffolding against the high-quality zebra finch (*Taeniopygia guttata*)
46 genome assigned 98.6% of the assembly to 32 putative nuclear chromosome scaffolds. Rapid,
47 recent advances in sequencing technologies and bioinformatics software have highlighted the
48 need for evidence-based assessment of assembly decisions on a case-by-case basis. Using *S.*
49 *vulgaris* vAU, we demonstrate how the multifunctional use of PacBio Iso-Seq transcript data
50 and complementary homology-based annotation of sequential assembly steps (assessed using
51 a new tool, SAAGA) can be used to assess, inform, and validate assembly workflow
52 decisions. We also highlight some counter-intuitive behaviour in traditional BUSCO metrics,
53 and present BUSCOMP, a complementary tool for assembly comparison designed to be robust
54 to differences in assembly size and base-calling quality. Finally, we present a second starling
55 assembly, *S. vulgaris* vNA, to facilitate comparative analysis and global genomic research on
56 this ecologically important species.

57 **Keywords:** *Sturnus vulgaris*, genome assembly, genome assessment, genome annotation,
58 full-length transcripts

59

60

61 **1. Introduction**

62 The European starling (*Sturnus vulgaris*) is a globally invasive passerine that was
63 deliberately introduced during early European acclimatisation efforts into North America,
64 Australia, New Zealand, and South Africa during the mid-late 19th century (Feare 1985).
65 More recently, the species was accidentally introduced into South America (Palacio et al.
66 2016). Since these introductions the invasive ranges of the starling have been expanding, with
67 the species now occupying a range in excess of 38,400,000 km² globally (BirdLife
68 International 2020), posing threats to the economics and health of the agriculture industry, as
69 well as local biodiversity (Bomford & Sinclair 2002; Koch et al., 2009; Palacio et al., 2016;
70 Linz et al., 2017). Recent molecular ecology studies of individuals from the invasive ranges
71 of North America, Australia, and South Africa report that these populations are undergoing
72 rapid and independent evolution in response to novel local selection pressures (Phair et al.,
73 2018; Hofmeister et al., 2019; Bodt et al., 2020; Stuart et al., 2020), a common phenomenon
74 in many invasive populations (Prentis et al., 2008). This suggests the starling has a flexible
75 invasion strategy, potentially enabling colonisation of ecosystems vastly different from those
76 in their native range.

77 Despite their invasive range success, European starlings are increasingly of ecological
78 concern within their native range (Rintala et al., 2003; Robinson et al., 2005). High densities
79 of native range starlings have traditionally been supported by cattle farming across Europe,
80 because starlings preferentially feed in open grasslands, and benefit from invertebrates in
81 overturned soil produced by livestock grazing (Coleman 1977). A shift towards modern
82 indoor cattle rearing processes across Europe may contribute to the decline in starling
83 numbers, which has been a concern since the 1980s (Wretenberg et al., 2006). This decline is
84 reflected globally, with starling and other avifauna numbers decreasing sharply over the last
85 few decades (Spooner et al., 2018; Rosenberg et al., 2019), though this may be further

86 amplified for starling populations subjected to control strategies to reduce their economic
87 impact (Linz et al., 2007). The biological and ecological importance of this species is evident
88 from its prolific use in research, as it is the most studied non-domesticated passerine (Bateson
89 & Feenders 2010). It is evident that future research on the European starling will focus on
90 identifying patterns of evolutionary diversification, as well as investigating genes associated
91 with invasion success. Such research provides important information for the improvement of
92 control measures and may also provide insight into recovery and dispersive potential in other
93 species that would benefit global conservation efforts. For this, a high-quality, annotated
94 reference genome is essential.

95 Once reliant on large consortia, assembling high-quality reference genomes for
96 genetic analyses is now commonplace. Nevertheless, *de novo* assembly of non-model
97 organism genomes still poses many challenges. Best practices may have not been established
98 for the study species/data, and basic information such as genome size, repeat landscape, and
99 ploidy may be unknown. Furthermore, high-quality references can be generated in multiple
100 ways, which can serve varied research purposes. Rapid developments in both sequencing
101 technology and bioinformatics methods can quickly outdate benchmarking attempts. Whilst
102 not always documented in final publications, the standard practice for non-model species
103 genomes is to select from multiple assemblies generated using different assembly methods,
104 none of which is universally best (Rhie et al., 2020; Whibley et al., 2020). This complexity
105 can be magnified further when sequencing occurs across multiple technology platforms that
106 may be combined and utilised in different ways (Jayakumar & Sakakibara 2019; Kono &
107 Arakawa 2019). The challenge is then to select the best combination of tools and assembly
108 decisions, based on the quality of the genome assemblies produced.

109 A multitude of tools and approaches are available for genome assembly assessment,
110 though some may not be applicable or feasibly implemented for a particular

111 species/assembly and/or the data available (e.g. Bradnam et al., 2013; Hunt et al., 2013; Yuan
112 et al., 2017). Common approaches employed to guide genome assembly decisions focus on
113 contiguity (how continuous the assembled sequences are) and completeness (whether the
114 assembly contains all the genetic information for that species). Two such approaches are
115 assembly statistics (e.g., contig/scaffold counts, and L50/N50 statistics of the number and
116 shortest length of sequences needed to cover 50% of the assembly) and “Benchmarking
117 Universal Single Copy Orthologs” (BUSCO) estimates of genome completeness (Simão et
118 al., 2015). Assembly statistics are very quick to generate and easy to understand, but
119 interpretation can be challenging due to hidden assembly errors and artefacts, which can
120 create false signals. BUSCO assesses the presence or absence of highly conserved lineage-
121 specific genes but is limited to a set of common single-copy genes that may represent easier
122 regions of the genome to sequence and assemble based on current bioinformatics
123 technologies. Furthermore, BUSCO analysis is vulnerable to unpredictable misreporting of
124 presence/completeness for specific genes as a consequence of assembly differences elsewhere
125 in the genome (Edwards et al., 2018; Edwards 2019; Field et al., 2020; Edwards et al., 2021).
126 In addition to the above drawbacks, these methods do not explicitly test the genome
127 assembly’s ability to perform the role for which it was intended (e.g., to serve as a reference
128 genome for specific genomic analysis).

129 Here, we present the first official release of the European starling draft genome, *S.*
130 *vulgaris* vAU. This assembly represents the first synthesis of species-specific full-length
131 transcripts, together with genomic data for this species. In this paper, we complement genome
132 statistics and BUSCO completeness with transcriptome- and annotation-based assessments
133 that help determine genome assembly quality and completeness in the absence of a reference
134 genome to benchmark against. We demonstrate how full-length transcripts can be utilised in
135 genome assembly scaffolding and assessment, in addition to transcriptome construction and

136 annotation. We show how BUSCOMP (<https://github.com/slimsuite/buscomp>) can help avoid
137 over-interpretation or misinterpretation of small differences in BUSCO completeness. We
138 also explore how lightweight homology-based annotation by GEMOMA (Keilwagen et al.,
139 2018), can be used as an assembly assessment using a new tool, SAAGA
140 (<https://github.com/slimsuite/saaga>). Finally, we compare the Australian *S. vulgaris* vAU
141 assembly (GCF_JAGFZL000000000) with an additional (short read) assembly of a North
142 American bird, *S. vulgaris* vNA (GCF_001447265.1), enabling reference-specific biases to
143 be identified in future starling genomics studies.

144

145 **2. Material and Methods**

146 **2.1 Genome assembly and scaffolding**

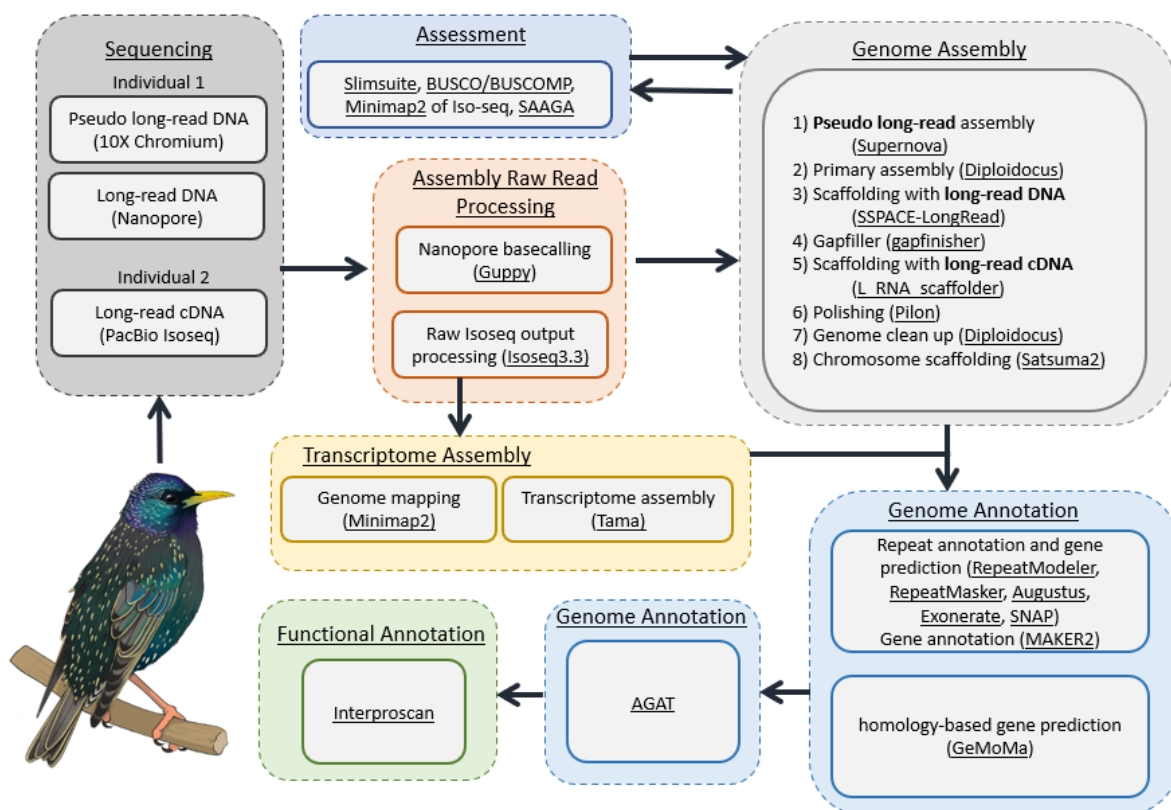
147 The *S. vulgaris* vAU genome assembly used 10x Chromium linked reads and low coverage
148 ONT long reads (Appendix 1: Genomic DNA sample collection, gDNA extraction, and
149 sequencing), and was produced via eight assembly steps (Fig. 1). The 10x reads were
150 assembled into an initial diploid assembly using SUPERNOVA (v2.1.1) (Weisenfeld et al.,
151 2017) with barcode fraction and reads subsample calculated following SUPERNOVA best
152 practices for a genome size based on k-mer counts calculation by JELLYFISH v2.2.10 (Marçais
153 & Kingsford 2011) (parameters: bfrac = 0.8, maxreads = 550 million, Supplementary
154 Materials: Appendix 2, Validation of SUPERNOVA genome size prediction using JELLYFISH,
155 Supplementary Materials: Fig. S1). This assembly was then split into non-redundant primary
156 and alternative haploid assemblies using DIPLOIDOCUS (parameters: runmode= diphapnr)
157 (v0.9.5) (<https://github.com/slimsuite/diploidocus>). First, both SUPERNOVA pseudohap2
158 assemblies were combined and any sequences lacking definitive base calls (100% Ns) were
159 removed. Remaining scaffolds were size-sorted and gaps reduced in size to a maximum of 10
160 Ns then subject to an all-by-all search with MINIMAP2 (v2.17) (Li 2018) (--cs -p 0.0001 -x

161 asm20 -N 250). (Note that gap size reduction is used for MINIMAP2 searching only, and the
162 non-redundant pseudodiploid assembly produced has the same gap sizes as generated by
163 SUPERNOVA.). Any sequences that were 100% contained within another sequence were
164 removed. Where two or more scaffolds had an 100% identical sequence, only one was kept.
165 Scaffolds are then matched into haplotig pairs based on their SUPERNOVA names. Where a
166 single haplotig is found, it is assigned as diploid, under the assumption that the two original
167 haplotigs were identical with one removed, and added to the primary assembly. (Note: it is
168 possible that only one parent had this scaffold, e.g., a sex chromosome scaffold or structural
169 variant.). If two haplotigs are identified, the longest is assigned to the primary assembly and
170 the shorter to the alternative assembly. The primary assembly should therefore contain an
171 entire haploid copy of the genome, whilst the alternative assembly contains the subset of
172 scaffolds with heterozygous haplotigs.

173 The primary haploid assembly produced by DIPLOIDOCUS was scaffolded using the
174 filtered ONT reads using the program SSPACE-LONGREAD (v1-1) (Boetzer & Pirovano
175 2014). The filtered ONT reads were then used to gap-fill the assembly using GAPFINISHER
176 (v1.0) (Kammonen et al., 2019). Clustered high quality Iso-Seq reads (see section 2.2 cDNA
177 analysis) were then used for a secondary round of scaffolding using L_RNA_SCAFFOLDER
178 (Xue et al., 2013). Paired-end 10x linked reads were processed with 10X Genomics LONG
179 RANGER (v2.2) and mapped onto this scaffolded assembly using BWA *mem* before error
180 correction of SNPs and indels using PILON (v1.23) (Walker et al., 2014) (parameters: --
181 diploid --fix all settings). To validate the scaffolds, the assembly was analysed using the
182 BREAK10X toolkit in SCAFF10X (v3.1) (<https://github.com/wtsi-hpag/Scaff10X>). The
183 assembly was further checked for assembly artefacts and contamination using DIPLOIDOCUS
184 (parameters: runmode=purgehaplotig & runmode=vecscreen (ref); runmode=DipCycle was
185 tested yet discarded due to over-pruning, see Supplementary Materials: Fig. S2) (v0.9.5).

186 Avian species are characterised by distinctive and constrained karyotypes, generally
 187 comprised of approximately 10 macrochromosomes and approximately 30 indistinguishable
 188 microchromosomes (Griffin et al., 2007; O'Connor et al., 2019), a pattern to which the *S.*
 189 *vulgaris* genome conforms (Calafati & Capanna 1981). Therefore, we aligned our assembly
 190 to the chromosome scale assembly of zebra finch (*Taeniopygia guttata*) (NCBI=
 191 GCF_008822105.2) (Balakrishnan et al., 2010) using SATSUMA2
 192 (<https://github.com/bioinfologics/satsuma2>) to create putative chromosomes assuming
 193 orthology. This assembly formed the final updated draft genome we present for the species:
 194 *Sturnus vulgaris* vAU.

195



196

197 **Figure 1: Workflow for genome assembly and annotation.** A summary of all the
 198 experimental methods used for sequencing, genome assembly, transcriptome assembly,
 199 genome annotation, and functional annotation, with programs used underlined.
 200

201

202

203

204 **2.2 Transcriptome assembly and analysis**

205 Raw PacBio Iso-Seq whole transcript reads (Appendix 3: Transcriptome sample collection,
206 RNA extraction, and sequencing) were processed using the protocol outlined in SMRT LINK
207 (v9.0) (PacBio, California, United States). Briefly, this involved generating Circular
208 Consensus Sequences (CCS) using CCS (v4.2.0), which were then processed using Lima
209 (v1.11.0) for primer removal and demultiplexing. The reads were further processed (PolyA
210 tail minimum length = 8) and clustered using ISO-SEQ (v3.3). The high quality clustered Iso-
211 Seq reads were then aligned to the reference genome (see section 2.1 Genome assembly and
212 scaffolding) using minimap2 (v2.17) (Li 2018), before further processing using TAMA
213 collapse (Kuo et al., 2020) (settings -a 100 -z 30 -sj sj_priority -lde 5). Both these steps were
214 assessed using BUSCO (v3.0.2b) (Simão et al., 2015) (parameters: aves lineage,
215 transcriptome mode), alongside a short read transcriptome produced from *S. vulgaris* liver
216 RNA (Richardson et al., 2017), as well as other available avian Iso-seq transcriptomes
217 (Workman et al., 2018; Yin et al., 2019).

218 **2.3 Genome annotation and functional annotation**

219 Each stage of genome assembly was annotated using GEMOMA v1.7.1 (Keilwagen et al.,
220 2018) using the 26 avian genome annotations available on Ensembl at the time this analysis
221 was conducted (Supplementary Materials, Table S1) and with the high quality clustered Iso-
222 seq, as RNA evidence. The GEMOMA *GeMoMaPipeline* function was run to complete the full
223 pipeline with a maximum intron size of 200 kb (parameters: tblastn=false
224 GeMoMa.m=200000 GeMoMa.Score=ReAlign AnnotationFinalizer.r=SIMPLE pc=true
225 o=true).

226 The final *S. vulgaris* vAU genome assembly was also annotated with MAKER2 (Holt &
227 Yandell 2011) (BLAST+ v2.9 (Camacho et al., 2009), AUGUSTUS v3.3.2 (Stanke &
228 Morgenstern 2005), EXONERATE v2.2.0 (Gs & E 2005), REPEATMASKER v4.0.7 (Smit et al.,

229 2013), REPEATMODELER v1.0.11 (Flynn et al., 2020), and SNAP v0.15.4 (Korf 2004) using
230 repeat-filtered Swiss-Prot protein sequences (downloaded Aug 2018) (UniProt Consortium
231 2019). A custom AUGUSTUS species database was created by running BUSCO using the
232 Optimization mode Augustus self-training mode (--long), using the aves database for lineage.
233 MAKER2 was run using the recommended protocol, including generation of a repeat library,
234 and with the TAMA-processed Iso-Seq data included as primary species transcript evidence,
235 and the pre-existing short read liver transcript data (Richardson et al., 2017) provided as
236 alternate transcript evidence in the first iteration of the MAKER2 annotation process. We ran
237 MAKER2 for a total of three training runs, using the hidden Markov models (HMMs)
238 produced from SNAP training in each subsequent run. *Ab initio* genes were not retained in the
239 final annotation model to produce high quality and conservative gene predictions. GEMOMA
240 and MAKER2 annotations for the final *S. vulgaris* vAU assembly were combined using the
241 AGAT *agat_sp_merge_annotations* function to produce the final annotation. Functional
242 annotation of protein-coding genes were generated using INTERPROSCAN 5.25–64.0
243 (parameters: -dp -goterms -iprlookup -appl TIGRFAM, SFLD, Phobius, SUPERFAMILY,
244 PANTHER, Gene3D, Hamap, ProSiteProfiles, Coils, SMART, CDD, PRINTS, Pro
245 SitePatterns, SignalP_EUK, Pfam, ProDom, MobiDBLite, PIRSF, TMHMM). BLAST was
246 used to annotate predicted genes using all Swiss-Prot proteins (parameters: -evalue 0.000001
247 -seg yes -soft_masking true -lcase_masking -max_hsps). Annotation summaries were
248 generated using the AGAT *agat_sp_functional_statistics.pl* script, BEDTOOLS was used to
249 calculate gene coverage statistics. Gene ontology terms were assigned using WEGO v2.0 (Ye
250 et al., 2018).

251 **2.4 Annotation assessment using SAAGA: Summarise, Annotate & Assess Genome**

252 **Annotations**

253 SAAGA (Summarise, Annotate & Assess Genome Annotations) (v0.5.3)
254 (<https://github.com/slimsuite/saaga>) was used to assess annotation quality and compare
255 predicted proteins to the repeat- and transposase-filter Swiss-Prot protein sequences used for
256 MAKER2 annotation (above). SAAGA performs a reciprocal MMseqs2 (Steinegger & Söding
257 2017) search of annotated proteins against a (high-quality) reference proteome, identifying
258 best hits for protein identification and employing coverage ratios between query and hit
259 proteins as a means of annotation assessment to generate summary statistics, including:

260 • **Protein length ratio.** The length ratio of the annotated proteins versus its top
261 reference hit

262 • **F1 score.** An annotation consistency metric calculated using the formula:

$$266 \quad (2 \times \text{PROTCOV} \times \text{REFCOV}) / (\text{PROTCOV} + \text{REFCOV})$$

263 where PROTCOV is the proportion of the annotated protein covered by its best
264 reference protein hit, and REFCOV is the proportion of the best reference protein hit
265 covered by the annotated protein.

267 • **Completeness.** The summed percentage coverage of reference proteome.

268 • **Purity.** The summed percentage reference coverage of the annotated proteome.

269 • **Homology.** The percentage of annotated genes with any hit in reference.

270 • **Orthology.** The percentage of annotated genes with reciprocal best hits in reference.

271 • **Duplicity.** The mean number of annotated genes sharing the same best reference hit.

272 • **Compression.** The number of unique annotated genes that were the top hit for
273 reference proteins, divided by the total number of reference proteins with a hit.

274 • **Multiplicity.** The ratio of total number of annotated genes to reference proteins.

275 For protein length ratio and F1 score, values close to 1 means that the query protein closely
276 matches the length of the hit protein, indicating high fidelity of the gene prediction model and
277 underlying assembly. The remaining metrics will be closer to 1 (or 100%) for complete

278 annotations and assemblies without duplications, akin to BUSCO scores. Although the
279 maximum achievable value for these metrics will generally be unknown, comparative values
280 can be used to assess improvement in assembly and/or annotation.

281 SAAGA scores may be used to compare alternate annotations of the same assembly,
282 or to compare alternative assemblies in conjunction with consistent annotation. Low genome
283 contiguity, misassemblies, or frameshifting indels will affect the quality of predicted genes,
284 with poorer assemblies reporting more fragmented or truncated genes. This approach has
285 been facilitated by the rapid homology-based gene prediction program GEMOMA, which uses
286 reference genome annotation to predict protein-coding genes in the target genome. The
287 program can be run from one line of code and may be parallelised to run much faster than
288 other annotation software (e.g., MAKER2). The ease of this annotation tool opens the way for
289 conducting annotations for the purpose of assessment on sequential or even competing
290 genome annotation steps. Assessing the quality of protein-coding region predictions will help
291 ensure the final genome assembly can produce a high-quality annotation. Here, we used the
292 repeat-filtered Swiss-Prot database used in annotation, and the *Gallus gallus* reference
293 proteome (UP000000539_9031), to assess predicted protein quality and annotated proteome
294 completeness.

295 **2.5 Genome assembly completeness assessment**

296 Assembly contiguity and completeness was assessed for sequential genome assembly steps of
297 the *S. vulgaris* vAU assembly and compared to existing passerine chromosome level
298 assemblies available on NCBI, including the *S. vulgaris* vNA assembly (Assembly accession
299 GCF_001447265.1, Supplementary Material: Appendix 4, Assembly and annotation of the *S.*
300 *vulgaris* vNA genome version).

301 **2.5.1 BUSCO and BUSCOMP assembly completeness assessment**

302 Genome completeness was estimated using BUSCO (v3.0.2b, genome mode, aves lineage).
303 BUSCO results were collated across all assemblies using BUSCOMP v0.10.1
304 (<https://github.com/slimsuite/buscomp>). BUSCOMP collated BUSCO outputs across all
305 genome assembly stages and compiled a maximal non-redundant set of 4727 complete
306 BUSCOs found at single copy in at least one assembly. Compiled BUSCO predicted gene
307 sequences were mapped onto each assembly to be rated with MINIMAP2 v2.17 (Li 2018) and
308 re-scored in terms of completeness, thereby providing a robust and consistent means of
309 assessing comparable completeness across assemblies of the same genome.

310 **2.5.2 PacBio Iso-Seq completeness assessment**

311 The PacBio Iso-Seq reads were mapped on to genome assemblies using MINIMAP2
312 (parameters: -ax splice -uf --secondary=no --splice-flank=no -C5 -O6,24 -B4) (Li 2018) and
313 the number of Iso-Seq transcripts mapping on to each assembly, and their corresponding
314 mapping quality, was calculated.

315 **2.5.3 KAT k-mer completeness assessment**

316 The final genome assembly completeness was assessed by examining the read k-mer
317 frequency distribution with different assembly copy numbers based on the 10x Chromium
318 linked reads using K-MER ANALYSIS TOOLKIT (KAT) v2.4.2 (Mapleson et al., 2017) (30 bp
319 trimmed for R1 reads, and 16 bp trimmed for R2 reads).

320 **2.6 Additional genome statistics**

321 The Iso-Seq and final annotation transcript density, final annotation gene density, global SNP
322 variant density (based on a whole genome data set of 24 individuals from United Kingdom,
323 North America, and Australia, N=8 (Hofmeister et al., 2020), and GC-content were
324 calculated in sliding windows of width 1 Mb using BEDTOOLS v 2.27.1 (Quinlan & Hall
325 2010), and plotted across the largest 32 scaffolds in our final genome assembly (representing

326 more than 98% of the total assembly captured on putative chromosomes orthologous to other
327 avian chromosomes) using CIRCLIZE (v 0.4.9) (Gu et al., 2014).

328 **2.7 Genome assembly correction**

329 NCBI VecScreen flagged possible bacterial and adapter contamination in the final *S. vulgaris*
330 vAU assembly, which was missed by earlier contamination screening steps. An updated
331 version of DIPLOIDOCUS (runmode vecscreen) was run to mask shorter adapter sequences and
332 flag additional organism contaminants (screenmode=purge vecmask=27). Four related
333 bacterial strains (Delftia acidovorans SPH-1, Acidovorax sp. JS42, Alicyclophilus
334 denitrificans K601, Paraburkholderia xenovorans LB400) were identified, and so GABLAM
335 v2.30.5 (Davey et al., 2006) was used to search these four genomes against the final
336 assembly, and purge small contigs (<5,000 bp) that contained sequence matches (285 short
337 contigs excluded). For larger scaffolds that contained possible embedded contaminated
338 sequences, the high quality ONT reads were mapped using Minimap2 over the regions. For
339 those contaminated sites that had Nanopore reads spanning the contaminated region, the
340 sequences were masked, and for those lacking nanopore support, the scaffold was split and/or
341 trimmed to remove the contaminating sequence (seq 4 trimmed, seq 12 and 31 split into
342 chromosome and unplaced scaffold). Finally, gaps of unknown size were standardised to 100
343 bp, and mitochondrial genome insertions into the nuclear genome were assessed using
344 NUMTfinder (<https://github.com/slimsuite/numtfinder>) (Edwards et al., 2021) (none located).
345 This paper primarily analyses *S. vulgaris* vAU1.0 (which we refer to as *S. vulgaris* vAU),
346 while the final NCBI release (accession = JAGFZL000000000) is explicitly referred to as *S.*
347 *vulgaris* vAU1.1 when relevant.

348 **2.8 BUSCO versus BUSCOMP performance benchmarking**

349 BUSCO-containing scaffolds from the DIPLOIDOCUS primary haploid SUPERNOVA assembly
350 of *Sturnus vulgaris* vAU were extracted into a reduced genome ‘pribusco’ assembly for

351 additional BUSCO and BUSCOMP benchmarking (Supplementary Materials: Fig. S3).
352 BUSCO v3.0.2b (Simão et al., 2015) (HMMER v3.2.1 (Wheeler & Eddy 2013),
353 AUGUSTUS v3.3.2 (Stanke & Morgenstern 2005), BLAST+ v2.2.31 blast(Camacho et al.,
354 2009), EMBOSS v6.6.0 (Rice et al., 2000)) was run in genome mode with the aves_odb9
355 dataset (n=4915) on: the non-redundant pseudodiploid ('dipnr'), primary ('pri') and
356 alternative ('alt') assemblies; BUSCO-containing scaffolds from the primary assembly
357 ('pribusco'); a reverse-complemented copy ('revcomp'), combined with 'pribusco' to make a
358 100% duplicate assembly ('duplicate'); a direct copy ('copy') combined with 'duplicate' to
359 make a triplicated assembly ('triplicate'); three randomly shuffled versions of 'pribusco'
360 ('shuffle1', 'shuffle2', 'shuffle3'), added in combination to 'pribusco' to generate datasets of
361 increasing assembly size without increasing duplication levels ('2n', '3n' and '4n'); ten
362 straight repeats of the 'pribusco' run ('rep0' to 'rep9'). All BUSCO results were processed
363 with BUSCOMP v0.11.0 (MINIMAP2 v2.17). In addition to the full BUSCOMP analysis of all
364 runs, the following subsets were grouped for analysis (Supplementary File 1, BUSCO v3
365 BUSCOMP output):

- 366 • Pseudodip: 'dipnr', 'pri' and 'alt'. (Haploid versus diploid assemblies.)
- 367 • Core: 'dipnr', 'pri', 'alt', 'pribusco' and 'revcomp'. (Assembly filtering and
368 manipulation.)
- 369 • Duplication: 'copy', 'duplicate', 'triplicate'. (Duplicating scaffolds.)
- 370 • Size: 'shuffle1', 'shuffle2', 'shuffle3', '2n', '3n', '4n'. (Increasing assembly size
371 without duplication.)
- 372 • Replicates: 'rep0' to 'rep9'.

373 The same analysis was repeated with BUSCO v5.0.0 (Simão et al., 2015) (SEPP v4.3.10
374 (Mirarab et al., 2012), BLAST v2.11.0 (Camacho et al., 2009), HMMer v3.3 (Wheeler &
375 Eddy 2013), AUGUSTUS v3.3.2 (Stanke & Morgenstern 2005), PRODIGAL v2.6.3 (Hyatt et

376 al., 2010), METAELK v20200908 (Levy Karin et al., 2020)) and the aves_odb10 dataset
377 (n=8338).

378 **3. Results**

379 **3.1 *Sturnus vulgaris* vAU genome assembly**

380 Genome assembly of *Sturnus vulgaris* vAU combined three different sequencing
381 technologies for *de novo* genome assembly (10x Genomics linked reads, ONT long reads,
382 and PacBio Iso-Seq full length transcripts) (Table 1), before a predicted reference-based
383 scaffolding to the chromosome level using the high-quality reference assembly of *T. guttata*
384 (NCBI REF: GCF_008822105.2). Approximately 109 Gb (97x coverage) of 10x linked read
385 data (subsampling during assembly to 56x based on the estimated genome size of 1.119 Gb,
386 barcode subsampling of 80%) were assembled with SUPERNOVA (v2.1.1) (Weisenfeld et al.,
387 2017) (step 1) and converted to a primary haploid assembly (step2). We generated
388 approximately 8 Gb of raw genomic reads using an ONT minion, which were reduced to 5
389 Gb after stringent filtering (Table 1). These data were used to scaffold the genome (step 3)
390 and gap-fill (step 4), reducing the total number of scaffolds from 18,439 to 7,856, increasing
391 the scaffold N50 from 1.76 Mb to 7.12 Mb, and decreasing the scaffold L50 from 146 to 39
392 (Supplementary Materials: Fig. S4). These measures were further improved after Iso-Seq
393 scaffolding (step 5) (7,776 scaffolds, N50 7.12 Mb, and L50 38), followed by Pilon polishing
394 using 10x linked reads (step 6). Finally, following haplotig removal (step 7), chromosomal
395 alignment against the *T. guttata* reference genome (step 8) reduced the final number of
396 scaffolds to 1,628 (N50 72.5 Mb, and L50 5) (Supplementary Materials: Fig. S4), with 98.6%
397 of the assembly assigned to 32 putative nuclear chromosome scaffolds. While no whole
398 mitochondrial genome insertions were found, 27 smaller mitochondrial pseudogenes
399 (NUMTs) were located in *S. vulgaris* vAU1.1, with scaffold 31 (corresponding to the Z
400 chromosome) containing the highest amount (Table S2).

401 **Table 1: Summary of sequencing data** for *Sturnus vulgaris* vAU genome assembly and
 402 annotation

Genetic Data	Platform	Library	Library length/Mean Insert Size (kb)	Mean raw Read Length (bp)	Number of Reads	Number of bases (Gb)
gDNA	Hiseq X Ten	Paired-end 10x Chromium	51.7kb	150	361,950,449	108.58
gDNA	ONT MinION	Ligation	47kb	6,417	1,225,865	7.865
cDNA	PacBio	Iso-Seq	Full transcripts (brain) (2.6 kb)	12,000	20,558,110	38.650
cDNA	PacBio	Iso-Seq	Full transcripts (heart + testes) (2.0 kb)	10,000	18,985,944	29.496

403

404 ***Improvements to genome assembly completeness during scaffolding***

405 Sequential steps of scaffolding, polishing, and quality control (Fig. 1, Supplementary
 406 Materials: Fig. S2, Table S3) improved the genome assembly statistics considerably from the
 407 initial SUPERNOVA *S. vulgaris* assembly (Supplementary Materials: Fig. S4). BUSCO
 408 completeness was approximately 94.6%, which was largely achieved by the initial assembly
 409 (92.9%), but somewhat improved over the additional assembly steps (Fig. 2a). The final
 410 BUSCO completeness score is comparable to other chromosome-level passerine assemblies
 411 on NCBI (Fig. 3a). BUSCO predictions are susceptible to base calling errors and can also
 412 fluctuate due to changes elsewhere in the genome assembly (Edwards 2019) (see section 2.8
 413 BUSCO versus BUSCOMP performance benchmarking). As a consequence, BUSCO can
 414 under-report the true number of complete BUSCO genes in an assembly (Edwards et al.,
 415 2018; Field et al., 2020; Edwards et al., 2021). We therefore used BUSCOMP to compile
 416 complete BUSCO genes from across all stages of the assembly. Only 70 (1.4%) of the 4,915
 417 Aves BUSCO genes were found to be “Missing” from all assembly versions, with 4,764
 418 (96.9%) rated “Complete” in at least one stage (Fig 3a, BUSCOMP).

419 The final assembly had the fewest unmapped Iso-Seq reads (Fig. 2b), with the largest
 420 improvement seen post gap-filling, followed by chromosome scaffolding. An increase in
 421 missing Iso-Seq transcripts was observed after scaffolding with the Iso-Seq reads themselves,

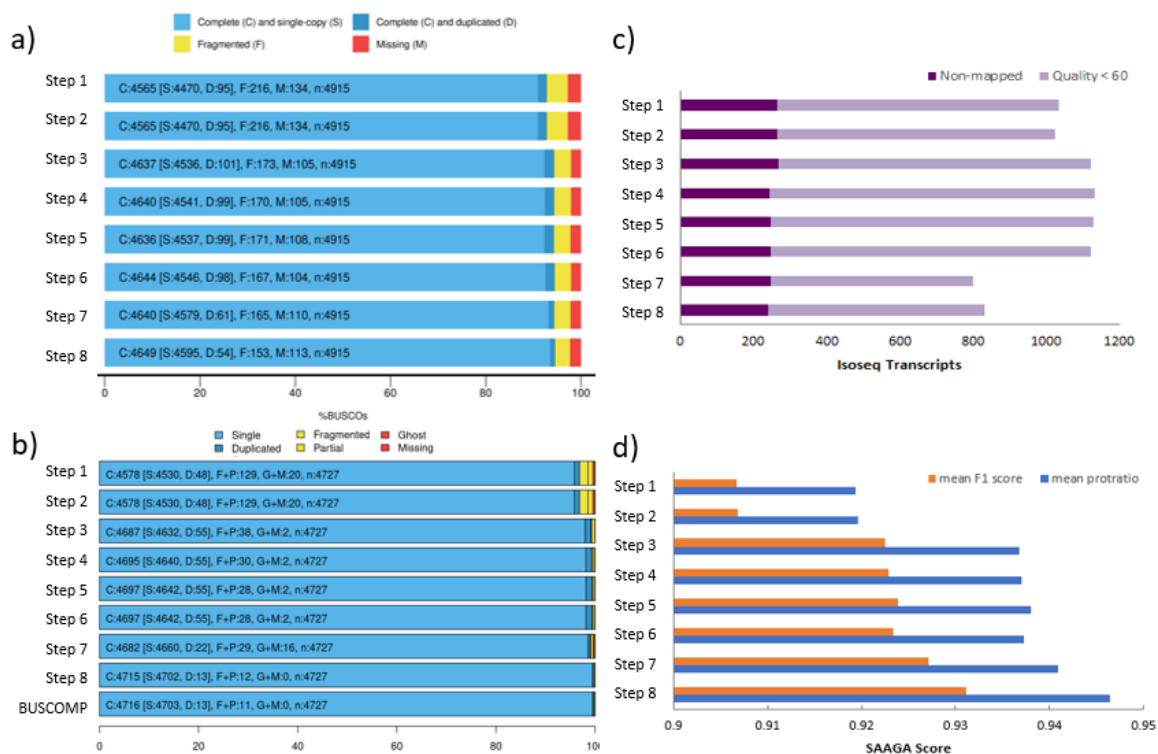
422 and post long-read scaffolding, due to reads no longer partially matching at scaffold ends.
423 Polishing caused a minimal improvement on the total number of mapped Iso-Seq reads, and
424 none were lost during scaffold clean-up with DIPLOIDOCUS (runmode purgehaplotig and
425 vecscreen). Assessment using GEMOMA and SAAGA revealed that across these assembly
426 steps we see a generally consistent increase in the quality of the predicted proteins during
427 annotation (Fig. 2c), with the largest increases occurring post long-read scaffolding, followed
428 by chromosome scaffolding, and then scaffold clean-up.

429 Of the 33,454 high quality isoform transcripts in the PacBio Iso-Seq data, only 241
430 failed to map to the final genome assembly, a 17.2% decrease compared to the 291 that failed
431 to map to *S. vulgaris* vNA (Fig. 3b).

432 ***Final genome assembly size, heterozygosity, and contiguity***

433 The *S. vulgaris* vAU assembly of 1,049,838,585 bp covers approximately 93.78% of
434 the total estimated 1.119 Gb genome size (Supplementary Materials: Appendix 2 Validation
435 of SUPERNOVA genome size prediction using JELLYFISH). A similar estimation of genome
436 completeness was reported by K-MER ANALYSIS TOOLKIT (KAT), with the raw read1s
437 (forward reads) estimating a genome completeness of 96.7% (estimated genome size 1.125
438 Gb, estimated heterozygosity rate 0.57%) and read2s (reverse reads) estimating a genome
439 completeness of 95.92% (estimated genome size 1.135 Gb, estimated heterozygosity rate
440 0.54%) (Supplementary Materials: Fig. S5). Predicted genome sizes based on either read1s or
441 read2s using KAT were slightly larger than the estimation generated by JELLYFISH using all
442 the read data, however the length range was relatively consistent (1.119-1.135 Gb). This
443 assembly reports a scaffold N50 of 72.5 Mb and L50 of 5, with a total of 1,628 scaffolds
444 (Table 2); 98.6% (1,035,260,756 bp) of the sequence length has been assigned to the 32
445 putative nuclear chromosomes (identified via the *T. guttata* v3.2.4 assembly), plus a
446 mitochondrial genome. The final assembly contains 14 macrochromosomes (> 20 Mb, as

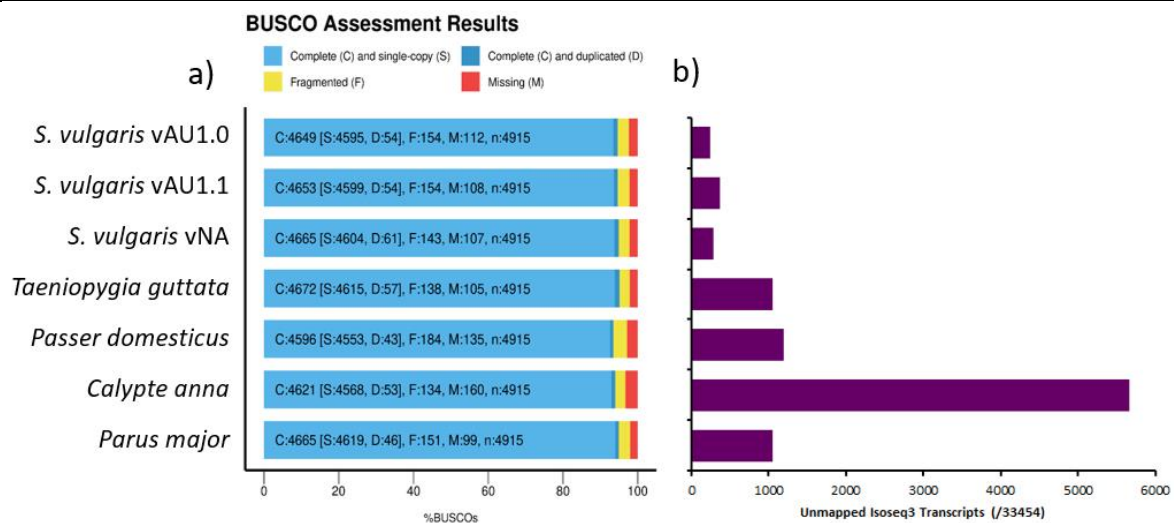
447 described in Backström et al., 2010), with relative sizes appearing in consensus with known
 448 karyotype of *S. vulgaris* (Calafati & Capanna 1981). Macrochromosome scaffolds account
 449 for 81.9% of the total assembly size, with the remainder on microchromosomes (16.9%) or
 450 unplaced scaffolds. While these large scaffolds remain only putative chromosomes assuming
 451 karyotype orthology until they can be validated with further read data, increased
 452 completeness scores post chromosomal alignment across all assembly assessment metrics
 453 (Fig. 2) support the assembly structure.



454
 455 **Figure 2: *Sturnus vulgaris* vAU assembly steps overview.** Quality and completeness
 456 assessments for eight sequential assembly steps: step 1 (SUPERNOVA assembly), step 2
 457 (DIPLODOCUS primary assembly), step 3 (SPACE-LONGREADS scaffolding), step 4
 458 (GAPFINISHER gapfilling), step 5 (L_RNA_SCAFFOLDER), step 6 (PILON polishing), step 7
 459 (DIPLODOCUS clean up), and step 8 (SATSUMA2 Chromosome scaffolding). **a)** BUSCO (Aves,
 460 n=4,915) completeness rating summaries for the sequential steps of *S. vulgaris* genome
 461 assembly. **b)** BUSCOMP completeness results for the 4,727 BUSCO genes identified as
 462 single copy and complete in one or more assembly stages. The final BUSCOMP row
 463 compiles the best rating for each gene across all eight steps. **c)** The number of Iso-Seq reads
 464 that failed to map to each assembly step. **d)** SAAGA annotation scores of mean protein length
 465 ratio (blue) and F1 score (orange) (see Methods for details).
 466

467 **Table 2: *Sturnus vulgaris* overview of assembly statistics for vAU1.0, vAU1.1, and vNA,**
 468 **assessed using BUSCOMP.**

	<i>Sturnus vulgaris</i> vAU1.0	<i>Sturnus vulgaris</i> vAU1.1	<i>Sturnus vulgaris</i> vNA
Total length (bp)	1,049,838,585	1,043,825,671	1,036,755,994
Number of scaffolds	1,628	1,344	2,361
Scaffold N50 (bp)	72,525,610	72,244,370	3,416,708
Scaffold L50	5	5	89
Largest scaffold (bp)	151,927,750	151,503,485	11,828,398
Mean scaffold length (bp)	644,864.0	776,656.01	439,117.3
Median scaffold length (bp)	1,337	1,343	4,856
Number of Contigs	23,815	23,340	22,666
Contig N50 (bp)	145,864	147,322	147,183
Contig L50	2,030	2,010	1,908
Gap (N) length (bp)	13,242,113 (1.26%)	0.74%	23,939,528 (2.31%)
GC (Guanine-Cytosine) content (%)	41.73%	41.72%	41.49%



469 **Figure 3: Assessment of *Sturnus vulgaris* and comparison avian reference assemblies. a)**
 470 **BUSCO (Aves) assessments of assembly completeness of *S. vulgaris* vAU1.0, and the NCBI**
 471 **uploaded genome *S. vulgaris* vAU1.1, presented alongside *S. vulgaris* vNA and four recent**
 472 **high-quality avian reference genomes (*Taeniopygia guttata* assembly accession**
 473 **GCF_008822105.2, *Passer domesticus* assembly accession GCA_001700915.1, *Calypte***
 474 ***anna* assembly accession GCA_003957555.2, *Parus major* assembly accession**
 475 **GCA_001522545.3). b) Total number of Iso-Seq transcripts that failed to map to each**
 476 **assembly.**
 477

478
 479

480 **3.2 *Sturnus vulgaris* vAU whole transcriptome data analysis**

481 We generated approximately 68 Gb of PacBio Iso-Seq whole transcript (39,544,054
482 subreads) (Table 1). This produced a total of 33,454 clustered high-quality (predicted
483 accuracy ≥ 0.99) reads, and 157 clustered low-quality (predicted accuracy < 0.99) reads
484 (Supplementary Materials: Table S4). These high-quality read data were used to improve the
485 scaffold assembly of the genome using L_RNA_SCAFFOLDER (see section 2.1) and assess
486 genome completeness (using count comparison of unmapped Iso-Seq reads, see section
487 2.5.2). After being passed through the TAMA *collapse* pipeline, a total of 28,448 non-
488 redundant transcripts were retained to create the final *S. vulgaris* vAU transcriptome, which
489 was used for gene prediction when completing the annotation of the genome assembly. This
490 final three tissue (brain, gonad, heart) Iso-Seq transcriptome had a moderate level overall
491 BUSCO completeness of around 63% that compares to other avian Iso-Seq transcriptomes
492 (Fig. 4a), with a wide range of gene ontology terms identified in the final Iso-Seq transcript
493 list (Fig. 4b) that resembled other avian Iso-Seq GO term distributions (Yin et al., 2019).

505 **3.3 *Sturnus vulgaris* genome annotation**

506 The initial annotation produced by GEMOMA, informed by the 26 avian genome
507 annotations available at the time on Ensembl (Supplementary Materials, Table S1), predicted
508 21,539 protein coding genes, with 97.2% BUSCO completeness (93.1% complete when
509 longest protein-per-gene extracted with SAAGA) (Fig. 5). The initial MAKER2 annotation
510 reported 13,495 genes, and a BUSCO completeness of 79.5% (Fig. 5). The merged final
511 annotation reported a BUSCO completeness of 98.2% (Fig. 5a), and this annotation predicted
512 a total of 21,863 protein-coding genes and 79,359 mRNAs. There was an average of 10.7
513 exons and 9.7 introns per mRNA, with an average intron length of 3,364 (Table 3). Of these,
514 1,764 are single-exon genes and 2,330 single-exon mRNA. Predicted coding sequences make
515 up 5.4% of the assembly, with 44.77% remaining outside any gene annotation (Fig. 5b).

516 The predicted transcripts were mapped using SAAGA to the Swiss-Prot database,
517 with 66,890 transcripts returning successful hits (84.3%) and 12,469 transcripts remaining
518 unknown (15.7%) for the final annotation (Fig. 6a). The known proteins had an average
519 length of 652 amino acids (aa) and the unknown proteins had an average length of 426 aa
520 (Fig. 6a). Most of the predicted proteins were of high quality, with around 56% of them
521 having an F1 score (see Methods) of greater than 0.95 (Fig. 6b). Similar results were seen
522 when the *Gallus gallus* reference proteome was used, with 69,714 known proteins of average
523 length of 646 aa, and 9,645 known proteins of average length of 401 aa, and the final merged
524 annotation having the same F1 score distribution (Fig 6c & 6d).

525 The GEMOMA annotation had similar protein quality patterns, with 57,026 known
526 proteins (average length 664 aa), and 10,400 unknown proteins (average length 401 aa) (Fig.
527 6e). The MAKER2 displayed much greater similarity in protein length histogram between
528 known and unknown proteins, with shorter proteins with known homologs (average length
529 565 aa), but longer unknown proteins (average length 549 aa) (Fig. 6f). The *S. vulgaris* vNA

530 annotation final merged annotation had extremely similar statistics to the final *S. vulgaris*
531 vAU annotation, with an average known protein length of 650 aa, and an average unknown
532 protein length of 407 aa (Supplementary Materials: Figure S2b).

533 **3.4 *Sturnus vulgaris* genome-wide patterns of genomics features**

534 Transcript density compared between mapped Iso-Seq reads and predicted transcripts
535 in the final annotation displayed similar patterns, with some minor variation in patterns
536 between the two tracks (Fig. 7; track 1). Final predicted gene densities (Fig. 7; track 2) were
537 largely following the patterns seen in transcript densities. Further, patterns of transcript and
538 gene numbers across the genome track relatively consistently to GC content (Fig. 7; track 4).

539 Global whole genome variant data (Fig. 7; track 3) revealed genomic regions where
540 variant density is low or non-existent, indicative of high genetic conservation across the
541 species, and genomic regions where variant density peaks are indicative of variant hotspots.
542 Interestingly, we see regions of high conservation corresponding to peaks in gene and/or
543 transcript numbers (e.g., midway through chromosome 4), which may be indicative of
544 regions of highly conserved genes and possibly centromere locations.

545

546

547

548

549

550

551

552

553

554

555

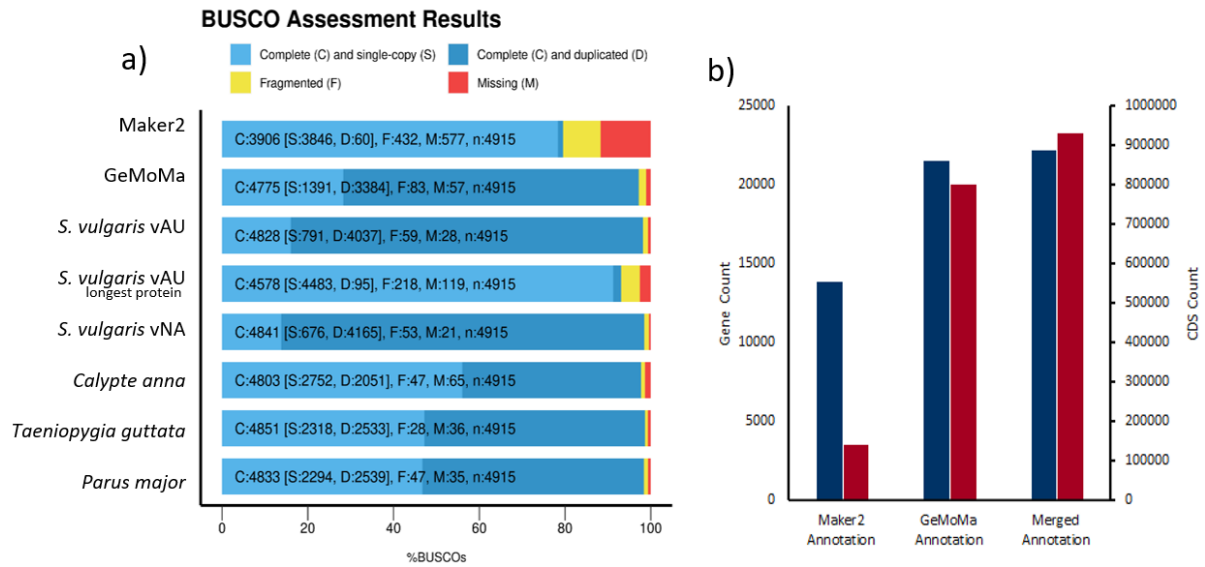
556

557

558 **Table 3: Summary of genome annotation of *Sturnus vulgaris* vAU and vNA assemblies.**
 559 Statistics extracted using AGAT *agat_sp_functional_statistics.pl*.
 560

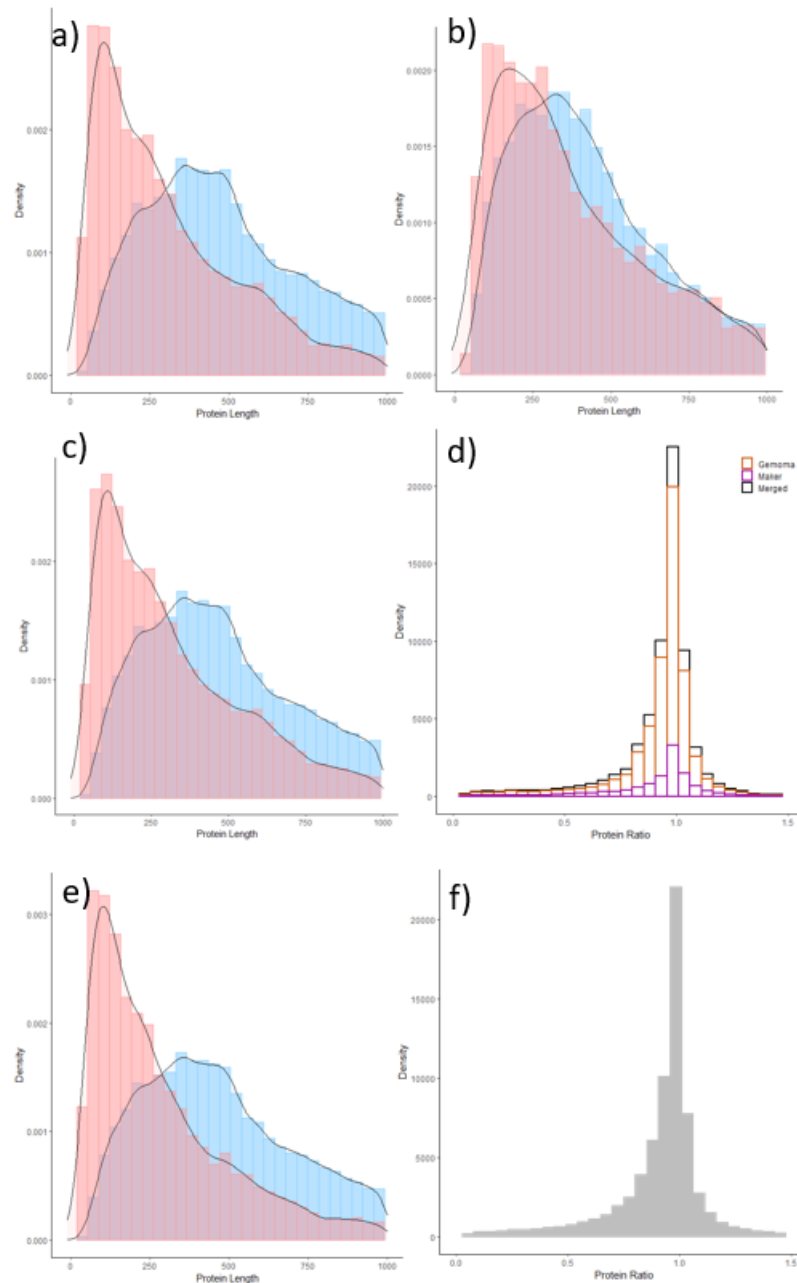
Genome Annotation		<i>S. vulgaris</i> vAU	<i>S. vulgaris</i> vNA
Genes	Total number	21,863	21,944
	Average length	34,699 bp	35,761 bp
	mean mRNAs per gene	3.6	3.7
mRNA	Total number	79,359	81,714
	Average length	38,073 bp	37,857 bp
	mean exons per mRNA	11.8	11.8
CDS	Total number	79,359	81,714
	Average length	1,851	1,836
	Average intron in CDS	3,364	3,343
	length		
Exons	Total number	933,014	962,220
	Mean length	163	158
Gene Function	Ontology Term	60.26% (13174/21863)	59.68% (13097/21944)
	InterPro	78.87% (17244/21863)	77.57% (17022/21944)
	SUPERFAMILY	60.36% (13197/21863)	58.26% (12786/21944)

561
 562
 563
 564
 565
 566



567
568
569
570
571
572
573
574

Figure 5: *Sturnus vulgaris* assessment of annotation. **a)** BUSCO (Aves) assessments of initial MAKER2 and GEMOMA assemblies, the final *S. vulgaris* vAU annotation, the final annotation with the longest protein-per-gene extracted using SAAGA, the final *S. vulgaris* vNA annotation (combined GEMOMA and MAKER2 annotation), and the ensemble annotations of three additional passerines. **b)** The number of genes (blue) and CDS (red) in the MAKER2 annotation, GEMOMA annotation, and merged annotation.



575

576 **Figure 6: Summary of predicted annotated proteins. a)** Protein lengths for known proteins

577 (blue, with a located Swiss-Prot comparison) and unknown proteins (red, those that did not

578 map to Swiss-Prot) for the GEMOMA annotation compared to Swiss-Prot. **b)** Protein lengths

579 of known and unknown proteins for the MAKER2 annotation compared to Swiss-Prot. **c)**

580 Protein lengths of known and unknown proteins for the merged GEMOMA and MAKER2

581 annotation compared to Swiss-Prot. **d)** Protein length ratio between output from SAAGA for

582 all known Swiss-Prot proteins (where a score close to 1 indicates a high-quality gene

583 annotation, protein length ratio calculated as annotated protein length / best Swiss-Prot

584 reference protein length) (merged annotation = black, GEMOMA annotation = orange,

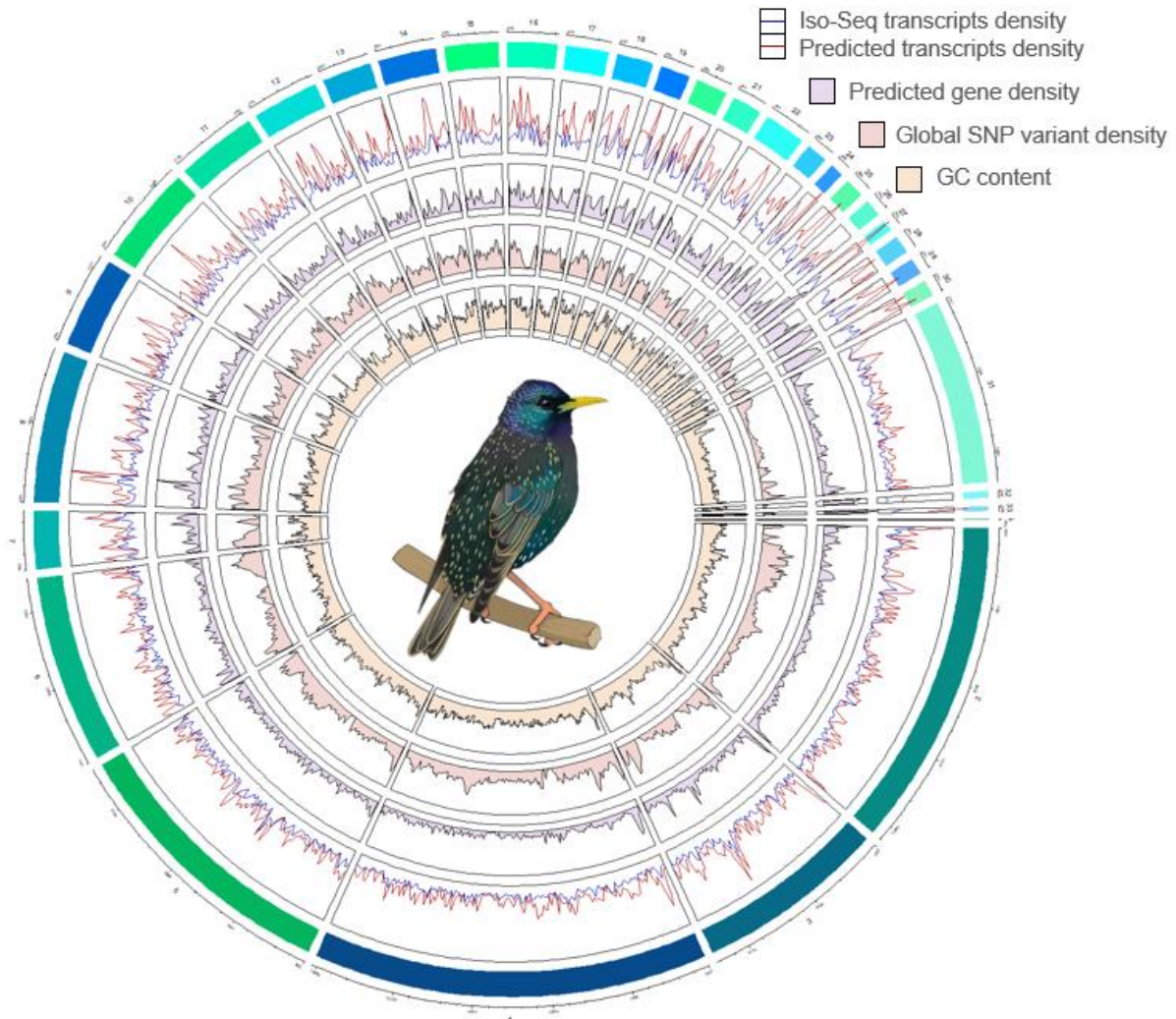
585 MAKER2 annotation = purple). **e)** Protein lengths of known and unknown proteins for the

586 merged GEMOMA and MAKER2 annotation compared to *Gallus gallus* reference proteome

587 (UP000000539_9031). **f)** Protein length ratio between output from SAAGA for the merged

588 annotation against the *Gallus gallus* reference proteome.

589



590

591 **Figure 7: CIRCLIZE plot of the 33 main chromosomal scaffolds** (32 putative autosomes
592 plus mtDNA) in the *Sturnus vulgaris* (*S. vulgaris* vAU) genome assembly (>98% of the total
593 assembly length). The tracks denote variable values in 1,000,000 bp sliding windows. From
594 the outermost track in, the variables displayed are track 1 (Iso-Seq transcripts as blue line,
595 final annotation transcripts as red line), track 2 (final annotation gene counts, purple area),
596 track 3 (variant density, red area), and track 4 GC content (yellow area).

597

598

599

600

601

602

603

604

605

606

607

608

609

610 **3.5 BUSCO versus BUSCOMP performance benchmarking**

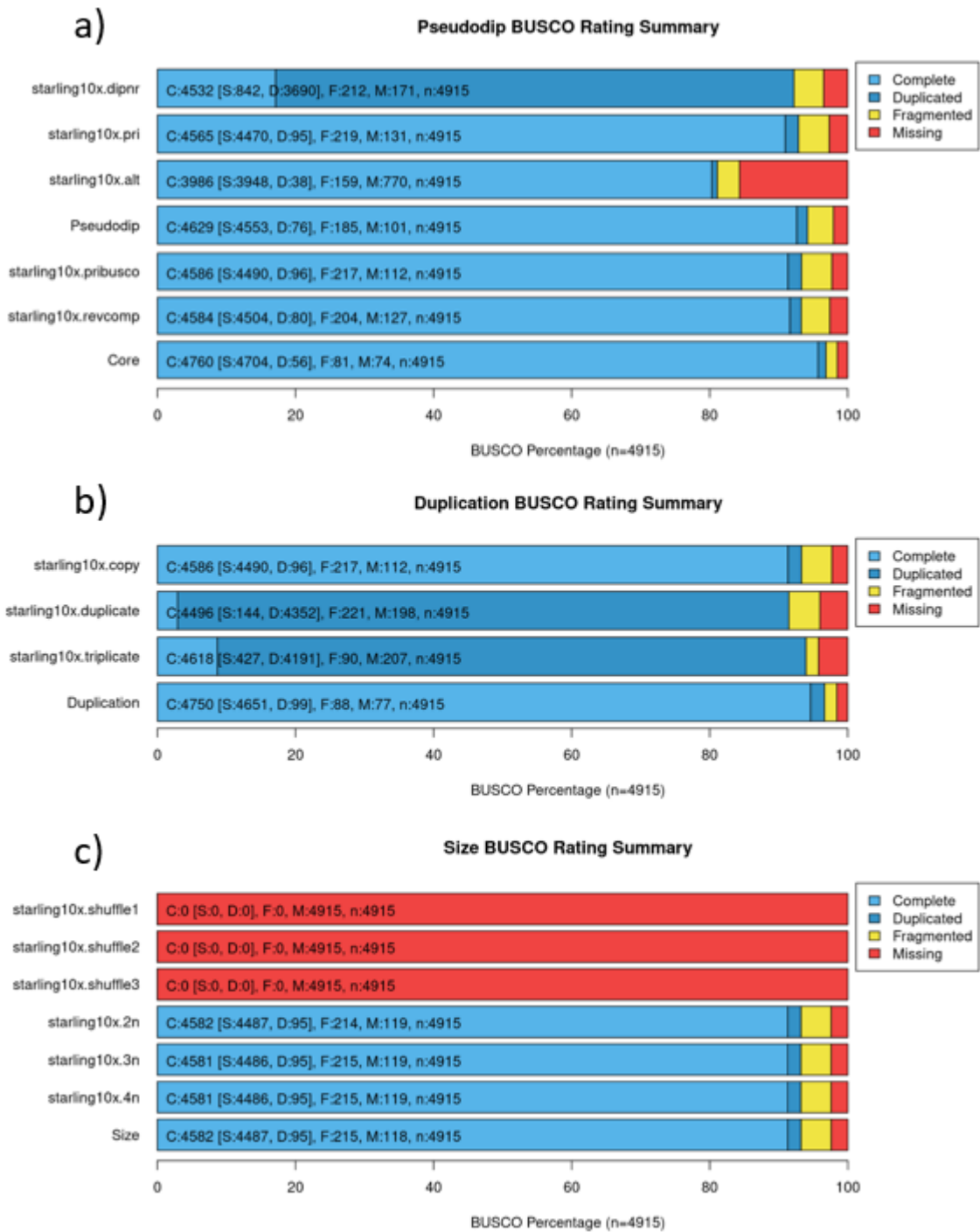
611 For the non-redundant pseudodiploid *S. vulgaris* vAU SUPERNOVA assembly, BUSCOMP
612 revealed differences in the BUSCO ratings of scaffolds dependent on the assembly
613 background (Fig. 8). Despite the primary ('pri') assembly being a subset of the non-
614 redundant pseudodiploid ('dinpr') assembly, it identified more "Complete" BUSCO genes
615 (4,565 versus 4,532) with fewer "Missing" (131 versus 171) (Fig. 8a). The alternative
616 assembly ('alt') subset similarly returned a partially overlapping set of BUSCO genes with
617 'dinpr', including some not found in 'dinpr' or 'pri': in total, only 101 genes were missing
618 from all three assemblies. Reducing the primary assembly to the 968 (of 18,439) scaffolds
619 containing a complete BUSCO gene ('pribusco'), increased the number of complete genes
620 from 4,565 to 4,586 and reduced the number missing from 131 to 112. Most unexpectedly,
621 reverse complementing these scaffolds reduced the complete BUSCO by two genes and
622 increased the number missing by fifteen (Fig. 8a). All five assemblies returned complete
623 BUSCO genes that were fragmented or missing in all the other four assemblies
624 (Supplementary File 1, BUSCOMP v3 results), for a combined total of 4,760 complete and
625 only 74 missing.

626 Adding direct or reverse-complemented copies of the 'pribusco' scaffolds increased
627 the number of "Duplicated" genes, but still returned single copy complete genes (Fig. 8b).
628 Doubling and then tripling the assembly size also increased the number of "Missing" genes
629 from 112 to 198 ('duplicate') and then 207 ('triplicate'). As before, these summary numbers
630 hide some gene gains as well as gene losses; only 77 genes are missing from all three
631 BUSCO runs, with 4,750 returned as complete by at least one. Adding randomly shuffled
632 versions of the 'pribusco' scaffolds only had a marginal effect on BUSCO ratings, with four
633 ('2n') to five ('3n', '4n') fewer complete genes returned and seven additional genes missing

634 following addition of the random sequences (Fig. 8c). Ten replicate analyses of the ‘pribusco’
635 scaffolds returned identical results (Supplementary File 1, BUSCOMP v3 results).

636 In contrast, BUSCOMP completeness is much more consistent across all datasets,
637 with the primary assembly returning the same numbers of complete, partial/fragmented and
638 missing genes as the pseudodiploid assembly (Supplementary File 1, BUSCOMP v3 results).
639 Similarly, reverse complementing scaffolds or increasing genome size gives no difference to
640 the completion statistics. Unlike BUSCO, BUSCOMP rates 100% of complete BUSCO genes
641 for duplicated or triplicated scaffolds as ‘Duplicate’ rather than ‘Single Copy’. Most
642 reassuringly, every complete BUSCO gene returned by a variant or subset of the
643 pseudodiploid assembly is also returned as ‘Complete’ in the pseudodiploid assembly itself.
644 Results using BUSCO v5 and the updated lineage data were qualitatively the same as v3,
645 showing largely identical trends (Supplementary File 2, BUSCOMP v5 results). The
646 exception is that reverse-complementing scaffolds reduced the complete BUSCO genes by
647 one (7,555 to 7,554) and increased the number missing by one (391 to 392). Curiously, this
648 was not reflected by analysis of the duplicated scaffolds, in which all 7,555 ‘pribusco’
649 complete genes were returned as complete and duplicated. It should be noted that the
650 ‘pribusco’ scaffolds for the v5 analysis are missing a greater proportion of the BUSCOMP-
651 compiled single copy complete BUSCO genes because they were still defined from v3 data.

652



653
654 **Figure 8. Compiled BUSCO results for benchmarking data.**

655
656
657
658
659
660
661

662 4. Discussion

663 Here, we present a high-quality, near-complete reference genome for the European
664 starling, *Sturnus vulgaris* vAU, with chromosome-level scaffolding that assigns 98.6% of the
665 genome assembly length to 32 putative nuclear chromosome scaffolds. We demonstrate the
666 utility of both transcripts and gene annotation in validating *S. vulgaris* vAU assembly
667 processes. BUSCOMP, Iso-Seq transcript, and SAAGA annotation assessment were largely
668 in agreement with one another, though each provided additional fine-scale feedback on
669 assembly improvements achieved by each assembly step. These analyses highlight the
670 benefits of these complementary assessment approaches in ensuring that aspects of genome
671 quality are not sacrificed to improve non-specific assembly quality metrics, such as N50. We
672 also present a second, North American, genome assembly, *S. vulgaris* vNA
673 (GCF_001447265.1). Overall, the *S. vulgaris* vAU assembly improved genome assembly
674 statistics over the *S. vulgaris* vNA genome, with a greater percentage of the estimated 1.119
675 Gb genome represented (94% vs 93%), an increase of scaffold N50 from 3.42 Mb to 72.5
676 Mb, and a decrease in scaffold L50 from 89 to 5. The *S. vulgaris* vNA still has good
677 assembly statistics (Table 2, Table S3) and has a marginally higher BUSCO completeness
678 (Fig. 3a) and BUSCOMP completeness (Supplementary Materials: Fig. S6) of approximately
679 20 BUSCO sequences. There is increasing recognition of the importance of pan-genomes
680 (genome assemblies that differentiate between genes/regions shared by all members of the
681 species, and dispensable or rare genes/regions) (Hirsch et al., 2014; Sherman & Salzberg
682 2020), which are essential for many model organisms (Vernikos et al., 2015). Having these
683 two high-quality *de novo* assemblies from different populations will improve future genomic
684 work on the global invasive populations of this species, and facilitate review of structural
685 variation (e.g., inversions) that may exist across different populations. It should be noted,
686 however, that the final scaffolding step for *S. vulgaris* vAU assumed structural conservation

687 between the starling and zebra finch and thus future synteny analyses may want to use the
688 earlier assembly step.

689 **4.1 BUSCO and BUSCOMP assembly completeness assessment**

690 BUSCO (Simão et al., 2015) is an extremely useful and widely-used assembly
691 assessment tool, providing information on which conserved lineage specific genes are
692 present, fragmented, or absent from a genome assembly. The program, however, can suffer
693 from inconsistent BUSCO gene identification, where a particular BUSCO may be dropped
694 from a report due to changes to contigs or scaffolds elsewhere in the assembly (Edwards
695 2019), which can result in under-reporting of assembly completeness (Edwards et al., 2018;
696 Field et al., 2020; Edwards et al., 2021). Here, we confirm this behaviour on benchmarking
697 datasets derived from the *S. vulgaris* vAU pseudodiploid 10x linked read assembly
698 (Supplementary Materials: Fig. S3, 8). Adding and removing scaffolds can both alter the
699 BUSCO ratings for “Complete” genes within the unchanged scaffolds (Fig. 8, Supplementary
700 File 1, BUSCOMP v3 results, Supplementary File 2, BUSCOMP v5 results). Many of these
701 changes are likely to be the consequence of changes in score thresholds and/or gene
702 prediction models. However, we also demonstrate some unexpected behaviours that are
703 harder to explain, such as changes to BUSCO gene ratings when scaffolds are reverse
704 complemented (Fig 8a).

705 This unpredictable variability in the identification of BUSCOs across genome
706 assembly versions poses some obvious challenges when trying to compare alternate versions
707 of the same assembly. This is particularly true when trying to interpret small changes in
708 BUSCO ratings as assemblies near completion. In addition, an important feature of BUSCO
709 is that it incorporates sequence quality in the context of the gene prediction models it
710 generates. This is desirable for assessing final assembly quality, but can present problems
711 when comparing early assembly stages, prior to error-correction by “polishing”. BUSCOMP

712 (<https://github.com/slimsuite/buscomp>) is robust to differences in assembly size and base-
713 calling quality and rates the “completeness potential” of an assembly based on the presence
714 of genes first identified for that species by BUSCO. Here, we used BUSCOMP analysis of
715 sequential assembly steps to gain a more accurate understanding of how assembly decisions
716 affected genome completeness (Fig. 2, Supplementary Materials: Fig. S4). BUSCOMP
717 analysis can then be complemented by other tools, such as KAT (Mapleson et al., 2017),
718 SAAGA (<https://github.com/slimsuite/saaga>), and BUSCO itself to get additional assessment
719 of sequence quality.

720 **4.1. Transcript- and annotation-guided *Sturnus vulgaris* vAU genome assembly**

721 The assembly of the *S. vulgaris* vAU genome was improved by assessing mapped Iso-
722 Seq whole transcripts and quality scores of predicted proteins from homology-based
723 annotation. Mapping of the high quality Iso-Seq reads proved to be an extremely fast method
724 of assessment (33,454 Iso-seq sequences mapped in <5 mins with 16 CPU cores), while the
725 GEMOMA and SAAGA compute time of 12 hrs per assembly was roughly comparable to
726 BUSCO (approximately 50 CPU hrs per assembly on an average machine), though more
727 computationally intensive (GEMOMA ran for approximately 200 CPU hours per assembly,
728 and SAAGA ran for approximately 8 CPU hours per assembly). Over the eight sequential
729 assembly steps, there was a decrease in unmapped Iso-Seq reads, indicating improved
730 sequence representation, primarily due to gap-filling yielding the greatest decrease in
731 unmapped Iso-Seq transcripts across all assembly steps. Similarly, the quality of annotated
732 proteins predicted by GEMOMA, as assessed by SAAGA, demonstrated ongoing
733 improvements through ONT scaffolding, clean-up, and chromosome alignment. It is also
734 noteworthy that increases in large-scale sequence connectivity using the *T. guttata* genome
735 (Peona et al., 2018) improved the assembly’s performance across all metrics, including

736 completeness estimates, although future Hi-C (or similar) analysis will be required to confirm
737 the predicted genome structure.

738 Further, BUSCOMP provided an important means of standardising BUSCO
739 annotation ratings across the multiple assembly steps. This method, together with the mapped
740 Iso-Seq reads, can deal with the unpolished intermediary genome steps, and does not suffer
741 the same sequence identification accuracy issues as the traditional stand-alone BUSCO
742 analysis. Together, the standardised assessment reported by BUSCOMP, and the
743 comprehensive and genome/species specific set of genes provided by Iso-seq and
744 GEMOMA/SAAGA showcase the complementary features of these annotation approaches for
745 assembly assessment.

746 **4.2. Improvements to contiguity and completeness during *Sturnus vulgaris* vAU genome** 747 **assembly**

748 Several alternative assembly pipelines were assessed (Supplementary Materials: Fig.
749 S2) and most of the upstream assembly decisions were based primarily on establishing
750 reasonable base assembly statistics (scaffold N50, scaffold L50, contig numbers). Assembly
751 size increased during scaffolding steps, due to estimated bases in gaps, while a decrease in
752 assembly size was only seen during scaffold clean up using *Purgehaplotigs*. Of all the
753 scaffolding steps, scaffolding with the low coverage ONT long reads resulted in the greatest
754 decrease of scaffold L50 (146 scaffold to 39 scaffold, Supplementary Materials: Fig. S4d)
755 and total scaffold number (18,439 scaffolds to 7,856 scaffolds, Supplementary Materials: Fig.
756 S4a). It has previously been shown that even low coverage of ONT data in conjunction with
757 10x may produce high quality genome assemblies (Ma et al., 2019). This was true for our
758 data, which demonstrates the utility of even low coverage, long read sequencing
759 (approximately 4.5% coverage based on the estimated genome size of 1.119 Gb) in greatly
760 improving the contiguity of scaffolds generated by short read genome assemblers (though Hi-

761 C data may serve this purpose at a lower cost to scaffold ratio and may assist in identifying
762 misassemblies, which is often not a focus of long-read scaffolding tools). While additional
763 scaffolding using the Iso-Seq whole transcripts did not result in a large increase in continuity
764 (Supplementary Materials: Fig. S4), the Iso-Seq reads were nevertheless able to scaffold
765 some sequences that failed low coverage ONT scaffolding, reducing the total scaffold count
766 by approximately 100 (7,856 scaffold to 7,776 scaffolds, Supplementary Materials: Fig. S4a).
767 This long-read transcript scaffolding served to minimise the number of fragmented genes in
768 the final assembly, helping downstream analysis and gene prediction models. The final
769 assembly maintained reasonably short contig N50 and high contig L50, which will only be
770 improved with much more extensive long-read sequencing of the species. Nevertheless,
771 scaffolding the *S. vulgaris* genome against that of *T. guttata* was able to further scaffold the
772 genome to a predicted chromosome level, assigning 98.6% of the assembly to previously
773 characterised chromosomes. In support of the assumed synteny of this step, we saw small
774 increases in assembly quality and completeness metrics.

775 The final two assembly steps (contig clean up and chromosomal alignment) were
776 primarily guided by high BUSCO scores and low missing Iso-Seq transcripts (Supplementary
777 Materials: Figure S3). DIPLOIDOCUS *vecscreen* did not flag any contamination and so did not
778 result in any assembly decreases. Over-pruning of contigs during clean up (using
779 DIPLOIDOCUS *DipCycle* which is stricter than just *Purgehaplotigs*) resulted in too many
780 (>1,000) discarded scaffolds that decreased assembly completeness scores, most likely
781 because of low coverage ONT long read data. While this drastically improved assembly
782 statistics, this came at the cost of dropped BUSCO sequences. Lastly, assembly duplication
783 analysis using KAT agreed with BUSCO results, indicating there was little final assembly
784 sequence duplication when comparing to raw read k-mer counts (Supplementary Materials:
785 Fig. S5).

786 **4.3. *Sturnus vulgaris* vAU transcriptome**

787 When comparing the completeness of this new starling transcriptome data to existing
788 Illumina short read transcript data produced using liver tissue (Richardson et al., 2017), we
789 see an increase of about 20% in BUSCO completeness, with a particularly large increase in
790 the number of duplicated BUSCO, a result of the alternate transcript isoforms captured
791 through the Iso-Seq. Assessing the effect the TAMA pipeline had on BUSCO completeness,
792 we see a small drop in complete BUSCOs (Fig. 2a) that appear to have been lost during the
793 mapping to genome assembly step. Finally, comparing our final transcriptome to two other
794 avian Iso-Seq transcriptomes gives an indication of how much unique transcript information
795 is added by the addition of tissues into pooled Iso-Seq sequencing runs. The single tissue Iso-
796 Seq liver transcriptome of *Calypste anna* (Anna's hummingbird) (Workman et al., 2018)
797 reported similar BUSCO completeness to the short read *S. vulgaris* liver transcriptome. The
798 eight tissue Iso-Seq transcriptome of *Anas platyrhynchos* (mallard) (Yin et al., 2019) yielded
799 an increase of 30% in complete BUSCOs, consistent with the expectation that our three-tissue
800 Iso-Seq library will be missing a number of tissue-specific genes.

801 **4.4. *Sturnus vulgaris* genome annotation**

802 Of the approximately 22,000 genes reported in the final annotation, 65% were from
803 GEMOMA, and 35% from MAKER2, with the source being randomly selected for common
804 annotation. MAKER2 predicted a higher number of genes in *S. vulgaris* vNA versus *S.*
805 *vulgaris* vAU (15,150 vs 13,495), while GEMOMA predicted a higher number of genes in the
806 *S. vulgaris* vAU genome (21,539 vs 20,414). The ratio in predicted MAKER2 and GEMOMA
807 was more biased towards the homology-based predictor, with an approximate ratio of 1:5
808 between MAKER2 and GEMOMA (Fig. 5b). Merging of the MAKER2 annotation to the
809 GEMOMA annotation resulted in an increase in 1.1.% in BUSCO completeness. Duplication
810 levels were much higher in the GEMOMA annotation when compared to MAKER2 (Fig. 5a).

811 This is not unreasonable, as the GEMOMA annotation will be biased toward well-
812 characterised genes and so may contain more transcripts per gene (Fig. 5b), whereas MAKER2
813 will inform the prediction of more taxon or possibly species-specific coding sequences. High
814 congruence between Iso-Seq and predicted transcript numbers indicate regions of accurate
815 annotation predictions (Fig. 7). In contrast, Iso-Seq transcripts that are dissimilar or much
816 lower to the predicted transcript densities, are either genomics regions producing tissue
817 specific transcripts not captured by their brain, testes, or muscle, or more likely annotated
818 transcript overprediction.

819 For the final *S. vulgaris* vAU annotation, the predicted proteins of unknown origin
820 (those that failed to map to Swiss-Prot database or *Gallus gallus* proteome) had a smaller
821 average length than those with known homologs (Fig. 6a & 6c). Similar results were found
822 when this approach was used to assess genes predicted in the *R. marina* genome assembly
823 (Edwards et al., 2018), and are indicative that these ‘unknown’ proteins are fragmented and
824 lower quality predictions that may be due to underlying assembly issues with contiguity or
825 frameshifting indels. The poorer quality could also reflect low stringency MAKER2 gene
826 predictions or homology based GEMOMA annotation of low-quality reference genes. The
827 known proteins predicted by MAKER2 (Fig. 6f) were of apparent lower quality than those
828 reported by GEMOMA as indicated by their shorter lengths and lower protein ratios (Fig. 6e),
829 which may be a result from a combination of incorrect gene predictions, and the high-quality
830 reference homologs inflating quality scores of the GEMOMA annotation in comparison.
831 The known protein lengths were similar across the *S. vulgaris* vAU and vNA annotations
832 (652 vs 650 aa), though there was a slightly larger difference in average unknown protein
833 length (426 vs 407 aa). Although this increase in *S. vulgaris* vAU is very slight, it may
834 indicate increased quality of unknown protein predictions in the vAU annotation, possibly
835 due to the more Iso-Seq data mapping to the vAU genome (Fig. 4b) or the higher contiguity.

836 Predicted genes were more commonly shorter than their closest reference protein hits,
837 indicative there might still be some truncated gene predictions, consistent with the large
838 number of assembly gaps. Nevertheless, the final annotation has a strong protein ratio peak
839 around 1.0 for known proteins (Fig. 6b & 6d), indicating that the bulk of these predicted
840 genes were of lengths similar to their Swiss-Prot homologs and hence deemed high quality.

841 Near identical assembly pipelines were used for the annotation of the two genome
842 assemblies, with the resulting final gene count predictions comparable to other high-quality
843 avian genomes and expected gene counts in eukaryote genomes. Both genome assembly
844 versions reported similar final annotation statistics, with *S. vulgaris* vNA reporting slightly
845 more predicted genes (Table 3), and a larger predicted gene coverage over the genome
846 (59.09% gene coverage vs 55.23%), indicating this increase in predicted genes is not just a
847 result of more overlapped predictions, though it could be a result of smaller assembly size
848 and higher gene duplication (Fig. 5a).

849

850 **5. Conclusion:**

851 This paper highlights the multifunctional use of species-specific transcript data, and
852 the importance of diverse assessment tools in the assembly and assessment of reference
853 genomes and annotations. We present a high-quality, annotated *S. vulgaris* vAU reference
854 genome, scaffolded at the chromosome level. Alongside a second assembly, *S. vulgaris* vNA,
855 these data provide vital resources for characterising the diverse and changing genomic
856 landscape of this globally important avian. In addition to improving the completeness of gene
857 annotation, we demonstrate the utility of long-read transcript data for genome quality
858 assessment and assembly scaffolding. We also reveal some counter-intuitive behaviour of
859 BUSCO genome completeness statistics, and present complementary two tools, BUSCOMP

860 and SAAGA, which can identify and resolve potential artefacts, and inform assembly
861 pipeline decisions.

862 **Author Contributions**

863 Project conception: all authors

864 Sample Collection: KCS, SJW, MCB

865 Lab Work: KCS, YC, LAR, WCW

866 Data Analysis: KCS, RJE, YC, WCW

867 Program Development: RJE

868 Manuscript Writing: KCS, RJE

869 Manuscript Editing: All authors

870 **Acknowledgements:**

871 We thank non-author members of the Starling Genome Consortium for their support of this
872 project including Wim Vanden Berghe. Thank you to Stella Loke, Annabel Whibley, and
873 Mark Richardson for their guidance of Nanopore sequencing and analysis. Thank you to
874 Daniel Selechnik for assistance with RNA extractions. Art credit (Fig. 7 illustration) to
875 Megan Bishop. RJE was funded by the Australian Research Council (LP160100610 and
876 LP18010072). DWB and YC acknowledge grant funding from the Human Sciences Frontier
877 Programme (Grant RGP0030/2015). SM acknowledges Roslin Institute Strategic Grant
878 funding from the UK Biotechnology and Biological Sciences Research Council
879 (BB/P013759/1). LAR was supported by a Scientia Fellowship from UNSW.

880

881 **Programs**

882 BUSCOMP documentation: <https://slimsuite.github.io/buscomp/>

883 Diploidocus documentation: <https://slimsuite.github.io/diploidocus/>

884 SAAGA documentation: <https://slimsuite.github.io/saaga/>

885

886

887 **References:**

- 888 Backström N, Forstmeier W, Schielzeth H, Mellenius H, Nam K, Bolund E, Webster MT, Öst T,
889 Schneider M, Kempenaers B *et al.* 2010 The recombination landscape of the zebra finch
890 *Taeniopygia guttata* genome. *Genome Research* **20** 485–495. (doi:10.1101/gr.101410.109)
- 891 Balakrishnan CN, Edwards SV & Clayton DF 2010 The Zebra Finch genome and avian genomics in the
892 wild. *Emu - Austral Ornithology* **110** 233–241. (doi:10.1071/MU09087)
- 893 Bateson M & Feenders G 2010 The use of passerine bird species in laboratory research: implications
894 of basic biology for husbandry and welfare. *ILAR Journal* **51** 394–408.
895 (doi:10.1093/ilar.51.4.394)
- 896 BirdLife International 2020 Species factsheet: *Sturnus vulgaris*. In *BirdLife International*.
- 897 Bodt LH, Rollins LA & Zichello J 2020 Genetic Diversity of the European Starling (*Sturnus vulgaris*)
898 Compared Across Three Invasive Ranges.
- 899 Boetzer M & Pirovano W 2014 SSPACE-LongRead: scaffolding bacterial draft genomes using long
900 read sequence information. *BMC Bioinformatics* **15** 211. (doi:10.1186/1471-2105-15-211)
- 901 Bomford M & Sinclair R 2002 Australian research on bird pests: impact, management and future
902 directions. *Emu - Austral Ornithology* **102** 29–45. (doi:10.1071/MU01028)
- 903 Bradnam KR, Fass JN, Alexandrov A, Baranay P, Bechner M, Birol I, Boisvert S, Chapman JA, Chapuis
904 G, Chikhi R *et al.* 2013 Assemblathon 2: evaluating de novo methods of genome assembly in
905 three vertebrate species. *GigaScience* **2**. (doi:10.1186/2047-217X-2-10)
- 906 Calafati P & Capanna E 1981 Karyotype analysis in ornithological studies: the chromosomes of six
907 species of oscines (passeriformes). *Avocetta* **5** 1–5.
- 908 Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K & Madden TL 2009 BLAST+:
909 architecture and applications. *BMC Bioinformatics* **10** 421. (doi:10.1186/1471-2105-10-421)
- 910 Coleman JD 1977 The foods and feeding of starlings in Canterbury. *Proceedings (New Zealand*
911 *Ecological Society)* **24** 94–109.
- 912 Davey NE, Shields DC & Edwards RJ 2006 SLiMDisc: short, linear motif discovery, correcting for
913 common evolutionary descent. *Nucleic Acids Research* **34** 3546–3554.
914 (doi:10.1093/nar/gkl486)
- 915 Edwards R 2019 BUSCOMP: BUSCO compilation and comparison – Assessing completeness in
916 multiple genome assemblies. *F1000Research* **8**:995 (slides).
917 (doi:10.7490/f1000research.1116972.1)
- 918 Edwards RJ, Tuipulotu DE, Amos TG, O’Meally D, Richardson MF, Russell TL, Vallinoto M, Carneiro M,
919 Ferrand N, Wilkins MR *et al.* 2018 Draft genome assembly of the invasive cane toad, *Rhinella*
920 *marina*. *GigaScience* **7**. (doi:10.1093/gigascience/giy095)
- 921 Edwards RJ, Field MA, Ferguson JM, Dudchenko O, Keilwagen J, Rosen BD, Johnson GS, Rice ES,
922 Hillier LD, Hammond JM *et al.* 2021 Chromosome-length genome assembly and structural
923 variations of the primal Basenji dog (*Canis lupus familiaris*) genome. *BMC Genomics* **22** 188.
924 (doi:10.1186/s12864-021-07493-6)

- 925 Feare CJ 1985 *The Starling*. Shire.
- 926 Field MA, Rosen BD, Dudchenko O, Chan EKF, Minoche AE, Edwards RJ, Barton K, Lyons RJ, Tuipulotu
927 DE, Hayes VM *et al.* 2020 Canfam_GSD: De novo chromosome-length genome assembly of
928 the German Shepherd Dog (*Canis lupus familiaris*) using a combination of long reads, optical
929 mapping, and Hi-C. *GigaScience* **9**. (doi:10.1093/gigascience/giaa027)
- 930 Flynn JM, Hubley R, Goubert C, Rosen J, Clark AG, Feschotte C & Smit AF 2020 RepeatModeler2 for
931 automated genomic discovery of transposable element families. *Proceedings of the National
932 Academy of Sciences* **117** 9451–9457. (doi:10.1073/pnas.1921046117)
- 933 Griffin DK, Robertson LBW, Tempest HG & Skinner BM 2007 The evolution of the avian genome as
934 revealed by comparative molecular cytogenetics. *Cytogenetic and Genome Research* **117** 64–
935 77. (doi:10.1159/000103166)
- 936 Gs S & E B 2005 Automated generation of heuristics for biological sequence comparison. *BMC
937 Bioinformatics* **6** 31–31. (doi:10.1186/1471-2105-6-31)
- 938 Gu Z, Gu L, Eils R, Schlesner M & Brors B 2014 circlize Implements and enhances circular visualization
939 in R. *Bioinformatics (Oxford, England)* **30** 2811–2812. (doi:10.1093/bioinformatics/btu393)
- 940 Hirsch CN, Foerster JM, Johnson JM, Sekhon RS, Muttoni G, Vaillancourt B, Peñagaricano F, Lindquist
941 E, Pedraza MA, Barry K *et al.* 2014 Insights into the Maize Pan-Genome and Pan-
942 Transcriptome. *The Plant Cell* **26** 121–135. (doi:10.1105/tpc.113.119982)
- 943 Hofmeister NR, Werner SJ & Lovette IJ 2019 Environment but not geography explains genetic
944 variation in the invasive and largely panmictic European starling in North America. *BioRxiv*
945 643858. (doi:10.1101/643858)
- 946 Hofmeister NR, Stuart K, +, Rollins LA & Clayton DF 2020 Replicated invasions by starlings reveal
947 hotspots of natural selection. *Unpublished*.
- 948 Holt C & Yandell M 2011 MAKER2: an annotation pipeline and genome-database management tool
949 for second-generation genome projects. *BMC Bioinformatics* **12** 491. (doi:10.1186/1471-
950 2105-12-491)
- 951 Hunt M, Kikuchi T, Sanders M, Newbold C, Berriman M & Otto TD 2013 REAPR: a universal tool for
952 genome assembly evaluation. *Genome Biology* **14** R47. (doi:10.1186/gb-2013-14-5-r47)
- 953 Hyatt D, Chen G-L, LoCascio PF, Land ML, Larimer FW & Hauser LJ 2010 Prodigal: prokaryotic gene
954 recognition and translation initiation site identification. *BMC Bioinformatics* **11** 119.
955 (doi:10.1186/1471-2105-11-119)
- 956 Jayakumar V & Sakakibara Y 2019 Comprehensive evaluation of non-hybrid genome assembly tools
957 for third-generation PacBio long-read sequence data. *Briefings in Bioinformatics* **20** 866–876.
958 (doi:10.1093/bib/bbx147)
- 959 Keilwagen J, Hartung F, Paulini M, Twardziok SO & Grau J 2018 Combining RNA-seq data and
960 homology-based gene prediction for plants, animals and fungi. *BMC Bioinformatics* **19** 189.
961 (doi:10.1186/s12859-018-2203-5)
- 962 Koch AJ, Munks SA & Spencer C 2009 Bird use of native trees retained in young eucalypt plantations:
963 species richness and use of hollows. *Wildlife Research* **36** 581–591. (doi:10.1071/WR09037)

- 964 Kono N & Arakawa K 2019 Nanopore sequencing: Review of potential applications in functional
965 genomics. *Development, Growth & Differentiation* **61** 316–326.
966 (doi:<https://doi.org/10.1111/dgd.12608>)
- 967 Korf I 2004 Gene finding in novel genomes. *BMC Bioinformatics* **5** 59. (doi:10.1186/1471-2105-5-59)
- 968 Kuo RI, Cheng Y, Zhang R, Brown JWS, Smith J, Archibald AL & Burt DW 2020 Illuminating the dark
969 side of the human transcriptome with long read transcript sequencing. *BMC Genomics* **21**
970 751. (doi:10.1186/s12864-020-07123-7)
- 971 Levy Karin E, Mirdita M & Söding J 2020 MetaEuk—sensitive, high-throughput gene discovery, and
972 annotation for large-scale eukaryotic metagenomics. *Microbiome* **8** 48. (doi:10.1186/s40168-
973 020-00808-x)
- 974 Li H 2018 Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics (Oxford, England)*
975 **34** 3094–3100. (doi:10.1093/bioinformatics/bty191)
- 976 Linz GM, Homan HJ, Gaulker SM, Penry LB & Bleier WJ 2007 European starlings: A review of an
977 invasive species with far-reaching impacts. *Managing Vertebrate Invasive Species*.
- 978 Linz G, Johnson R & Thiele J 2017 European Starlings. In *Ecology and Management of Terrestrial*
979 *Vertebrate Invasive Species in the United States*, 1st ed, pp 311–332. Ed WC Pitt. Boca
980 Raton : Taylor & Francis, 2018. | “A CRC title, part of the Taylor & Francis imprint, a member
981 of the Taylor & Francis Group, the academic division of T&F Informa plc.”: CRC Press.
982 (doi:10.1201/9781315157078-15)
- 983 Ma Z (Sam), Li L, Ye C, Peng M & Zhang Y-P 2019 Hybrid assembly of ultra-long Nanopore reads
984 augmented with 10x-Genomics contigs: Demonstrated with a human genome. *Genomics* **111**
985 1896–1901. (doi:10.1016/j.ygeno.2018.12.013)
- 986 Mapleson D, Garcia Accinelli G, Kettleborough G, Wright J & Clavijo BJ 2017 KAT: a K-mer analysis
987 toolkit to quality control NGS datasets and genome assemblies. *Bioinformatics* **33** 574–576.
988 (doi:10.1093/bioinformatics/btw663)
- 989 Marçais G & Kingsford C 2011 A fast, lock-free approach for efficient parallel counting of occurrences
990 of k-mers. *Bioinformatics* **27** 764–770. (doi:10.1093/bioinformatics/btr011)
- 991 Mirarab S, Nguyen N & Warnow T 2012 SEPP: SATé-enabled phylogenetic placement. *Pacific*
992 *Symposium on Biocomputing. Pacific Symposium on Biocomputing* 247–258.
993 (doi:10.1142/9789814366496_0024)
- 994 O’Connor RE, Kiazim L, Skinner B, Fonseka G, Joseph S, Jennings R, Larkin DM & Griffin DK 2019
995 Patterns of microchromosome organization remain highly conserved throughout avian
996 evolution. *Chromosoma* **128** 21–29. (doi:10.1007/s00412-018-0685-6)
- 997 Palacio FX, Maragliano RE & Montalti D 2016 Functional role of the invasive European Starling,
998 *Sturnus vulgaris*, in Argentina. *Emu* **116** 387–393. (doi:10.1071/MU16021)
- 999 Peona V, Weissensteiner MH & Suh A 2018 How complete are “complete” genome assemblies?—An
1000 avian perspective. *Molecular Ecology Resources* **18** 1188–1195.
1001 (doi:<https://doi.org/10.1111/1755-0998.12933>)

- 1002 Phair DJ, Roux JIL, Berthouly-Salazar C, Visser V, Vuuren BJ van, Cardilini APA & Hui C 2018 Context-
1003 dependent spatial sorting of dispersal-related traits in the invasive starlings (*Sturnus*
1004 *vulgaris*) of South Africa and Australia. *BioRxiv* 342451. (doi:10.1101/342451)
- 1005 Prentis PJ, Wilson JRJ, Dormontt EE, Richardson DM & Lowe AJ 2008 Adaptive evolution in invasive
1006 species. *Trends in Plant Science* **13** 288–294. (doi:10.1016/j.tplants.2008.03.004)
- 1007 Quinlan AR & Hall IM 2010 BEDTools: a flexible suite of utilities for comparing genomic features.
1008 *Bioinformatics* **26** 841–842. (doi:10.1093/bioinformatics/btq033)
- 1009 Rhie A, McCarthy SA, Fedrigo O, Damas J, Formenti G, Koren S, Uliano-Silva M, Chow W,
1010 Functamman A, Gedman GL *et al.* 2020 Towards complete and error-free genome
1011 assemblies of all vertebrate species. *BioRxiv* 2020.05.22.110833.
1012 (doi:10.1101/2020.05.22.110833)
- 1013 Rice P, Longden I & Bleasby A 2000 EMBOS: the European Molecular Biology Open Software Suite.
1014 *Trends in Genetics: TIG* **16** 276–277. (doi:10.1016/s0168-9525(00)02024-2)
- 1015 Richardson MF, Sherwin WB & Rollins LA 2017 De Novo Assembly of the Liver Transcriptome of the
1016 European Starling, *Sturnus vulgaris*. *Journal of Genomics* **5** 54–57. (doi:10.7150/jgen.19504)
- 1017 Rintala J, Tiainen J & Pakkala T 2003 Population trends of the Finnish starling *Sturnus vulgaris*,
1018 1952–1998, as inferred from annual ringing totals. *Annales Zoologici Fennici* **40** 365–385.
- 1019 Robinson RA, Siriwardena GM & Crick HQP 2005 Status and population trends of Starling *Sturnus*
1020 *vulgaris* in Great Britain. *Bird Study* **52** 252–260. (doi:10.1080/00063650509461398)
- 1021 Rosenberg KV, Dokter AM, Blancher PJ, Sauer JR, Smith AC, Smith PA, Stanton JC, Panjabi A, Helft L,
1022 Parr M *et al.* 2019 Decline of the North American avifauna. *Science* **366** 120–124.
1023 (doi:10.1126/science.aaw1313)
- 1024 Sherman RM & Salzberg SL 2020 Pan-genomics in the human genome era. *Nature Reviews Genetics*
1025 **21** 243–254. (doi:10.1038/s41576-020-0210-7)
- 1026 Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV & Zdobnov EM 2015 BUSCO: assessing
1027 genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics*
1028 (*Oxford, England*) **31** 3210–3212. (doi:10.1093/bioinformatics/btv351)
- 1029 Smit A, Hubley R & Green P 2013 *RepeatMasker Open-4.0*.
- 1030 Spooner FEB, Pearson RG & Freeman R 2018 Rapid warming is associated with population decline
1031 among terrestrial birds and mammals globally. *Global Change Biology* **24** 4521–4531.
1032 (doi:https://doi.org/10.1111/gcb.14361)
- 1033 Stanke M & Morgenstern B 2005 AUGUSTUS: a web server for gene prediction in eukaryotes that
1034 allows user-defined constraints. *Nucleic Acids Research* **33** W465–W467.
1035 (doi:10.1093/nar/gki458)
- 1036 Steinegger M & Söding J 2017 MMseqs2 enables sensitive protein sequence searching for the
1037 analysis of massive data sets. *Nature Biotechnology* **35** 1026–1028. (doi:10.1038/nbt.3988)

- 1038 Stuart KC, Cardilini APA, Cassey P, Richardson MF, Sherwin W, Rollins LA & Sherman CDH 2020
1039 Signatures of selection in a recent invasion reveals adaptive divergence in a highly vagile
1040 invasive species. *Molecular Ecology* **n/a**. (doi:10.1111/mec.15601)
- 1041 UniProt Consortium 2019 UniProt: a worldwide hub of protein knowledge. *Nucleic Acids Research* **47**
1042 D506–D515. (doi:10.1093/nar/gky1049)
- 1043 Vernikos G, Medini D, Riley DR & Tettelin H 2015 Ten years of pan-genome analyses. *Current Opinion*
1044 *in Microbiology* **23** 148–154. (doi:10.1016/j.mib.2014.11.016)
- 1045 Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, Cuomo CA, Zeng Q, Wortman J,
1046 Young SK *et al.* 2014 Pilon: An Integrated Tool for Comprehensive Microbial Variant
1047 Detection and Genome Assembly Improvement. *PLOS ONE* **9** e112963.
1048 (doi:10.1371/journal.pone.0112963)
- 1049 Weisenfeld NI, Kumar V, Shah P, Church DM & Jaffe DB 2017 Direct determination of diploid genome
1050 sequences. *Genome Research* **27** 757–767. (doi:10.1101/gr.214874.116)
- 1051 Wheeler TJ & Eddy SR 2013 nhmmer: DNA homology search with profile HMMs. *Bioinformatics* **29**
1052 2487–2489. (doi:10.1093/bioinformatics/btt403)
- 1053 Whibley A, Kelley JL & Narum SR 2020 The changing face of genome assemblies: Guidance on
1054 achieving high-quality reference genomes. *Molecular Ecology Resources*. (doi:10.1111/1755-
1055 0998.13312)
- 1056 Workman RE, Myrka AM, Wong GW, Tseng E, Welch KC & Timp W 2018 Single-molecule, full-length
1057 transcript sequencing provides insight into the extreme metabolism of the ruby-throated
1058 hummingbird *Archilochus colubris*. *GigaScience* **7**. (doi:10.1093/gigascience/giy009)
- 1059 Wretenberg J, Lindström Å, Svensson S, Thierfelder T & Pärt T 2006 Population trends of farmland
1060 birds in Sweden and England: similar trends but different patterns of agricultural
1061 intensification. *Journal of Applied Ecology* **43** 1110–1120. (doi:10.1111/j.1365-
1062 2664.2006.01216.x)
- 1063 Xue W, Li J-T, Zhu Y-P, Hou G-Y, Kong X-F, Kuang Y-Y & Sun X-W 2013 L_RNA_scaffolder: scaffolding
1064 genomes with transcripts. *BMC Genomics* **14** 604. (doi:10.1186/1471-2164-14-604)
- 1065 Ye J, Zhang Y, Cui H, Liu J, Wu Y, Cheng Y, Xu H, Huang X, Li S, Zhou A *et al.* 2018 WEGO 2.0: a web
1066 tool for analyzing and plotting GO annotations, 2018 update. *Nucleic Acids Research* **46**
1067 W71–W75. (doi:10.1093/nar/gky400)
- 1068 Yin Z, Zhang F, Smith J, Kuo R & Hou Z-C 2019 Full-length transcriptome sequencing from multiple
1069 tissues of duck, *Anas platyrhynchos*. *Scientific Data* **6** 275. (doi:10.1038/s41597-019-0293-1)
- 1070 Yuan Y, Bayer PE, Scheben A, Chan C-KK & Edwards D 2017 BioNanoAnalyst: a visualisation tool to
1071 assess genome assembly quality using BioNano data. *BMC Bioinformatics* **18** 323.
1072 (doi:10.1186/s12859-017-1735-4)
- 1073
- 1074
- 1075