# Identifying the causes and consequences of assembly gaps using a multiplatform genome assembly of a bird-of-paradise

Valentina Peona[1,2], Mozes P.K. Blom[3,4], Luohao Xu[5,6], Reto Burri[7], Shawn Sullivan[8], Ignas Bunikis[9], Ivan Liachko[8], Knud A. Jønsson[10], Qi Zhou[5,6,11], Martin Irestedt[3], Alexander Suh[1,2]

## Affiliation

[1] Department of Ecology and Genetics – Evolutionary Biology, Uppsala University, Science for Life Laboratories, Norbyvägen 18D, SE-752 36, Uppsala, Sweden
[2] Department of Organismal Biology – Systematic Biology, Uppsala University, Norbyvägen 18D, SE-752 36, Uppsala, Sweden
[3] Department of Bioinformatics and Genetics, Swedish Museum of Natural History, SE-104 05, Stockholm, Sweden
[4] Museum für Naturkunde, Leibniz Institut für Evolutions- und Biodiversitätsforschung, Berlin, Germany
[5] MOE Laboratory of Biosystems Homeostasis & Protection, Life Sciences Institute, Zhejiang University, Hangzhou, China
[6] Department of Molecular Evolution and Development, University of Vienna, Vienna, Austria
[7] Department of Population Ecology, Institute of Ecology and Evolution, Friedrich-Schiller-University Jena, Dornburger Strasse 159, D-07743 Jena, Germany
[8] Phase Genomics, Inc. 1617 8th Ave N, Seattle, WA 98109 USA
[9] Uppsala Genome Center, Science for Life Laboratory, Dept. of Immunology, Genetics and Pathology, Uppsala University, SE‑752 37, Uppsala, Sweden
[10] Natural History Museum of Denmark, University of Copenhagen, Universitetsparken 15, DK-2100 Copenhagen, Denmark
[11] Center for Reproductive Medicine, The 2nd Affiliated Hospital, School of Medicine, Zhejiang University

## Correspondence

V.P. (valentina.peona@ebc.uu.se), A.S. (alexander.suh@ebc.uu.se)

## Abstract

Genome assemblies are currently being produced at an impressive rate by consortia and individual laboratories. The low costs and increasing efficiency of sequencing technologies have opened up a whole new world of genomic biodiversity. Although these technologies generate high-quality genome assemblies, there are still genomic regions difficult to assemble, like repetitive elements and GC-rich regions (genomic "dark matter"). In this study, we compare the efficiency of currently used

45  sequencing technologies (short/linked/long reads and proximity ligation maps) and combinations

46  thereof in assembling genomic dark matter starting from the same sample. By adopting different *de-*

47  *novo* assembly strategies, we were able to compare each individual draft assembly to a curated

48  multiplatform one and identify the nature of the previously missing dark matter with a particular focus

49  on transposable elements, multi-copy MHC genes, and GC-rich regions. Thanks to this multiplatform

50  approach, we demonstrate the feasibility of producing a high-quality chromosome-level assembly for

51  a non-model organism (paradise crow) for which only suboptimal samples are available. Our

52  approach was able to reconstruct complex chromosomes like the repeat-rich W sex chromosome and

53  several GC-rich microchromosomes. Telomere-to-telomere assemblies are not a reality yet for most

54  organisms, but by leveraging technology choice it is possible to minimize genome assembly gaps for

55  downstream analysis. We provide a roadmap to tailor sequencing projects around the completeness

56  of both the coding and non-coding parts of the genomes.

57

## Introduction

59
60  With the advent of Next Generation Sequencing (NGS) technologies, the field of genomics has grown

61  exponentially and during the last 10 years the genomes of almost 10,000 species of prokaryotes and

62  eukaryotes have been sequenced (from NCBI Assembly database, O'Leary et al. (2015)). Traditional

63  NGS technologies rely on DNA amplification and generation of millions of short reads (few hundreds

64  of bp long) that subsequently have to be assembled into contiguous sequences (contigs; Goodwin et

65  al. (2016)). Although the technique has been revolutionary, the short-read length together with

66  difficulties to sequence regions with extreme base composition poses serious limitations to genome

67  assembly (Chaisson et al. 2015; Peona et al. 2018). Technological biases are therefore impeding the

68  complete reconstruction of genomes and substantial regions are systematically missing from genome

69  assemblies. These missing regions are often referred to as the genomic "dark matter" (Johnson et al.

70   2005). It is key now for the genomics field to overcome these limitations and investigate this dark

71   matter.

72

73   Repetitive elements represent an important and prevalent part of the genomic dark matter of many

74   genomes, given that their abundance and repetitive nature makes it difficult to fully and confidently

75   assemble their sequences. This is particularly problematic when the read length is significantly shorter

76   than the repetitive element, in which case it is impossible to anchor the reads to unique genomic

77   regions. To what extent repeats can hamper genome assemblies depends on whether they are

78   interspersed or arranged in tandem. Highly similar interspersed repeats, like for example transposable

79   elements (TEs), may introduce ambiguity in the assembly process and cause assembly (contig)

80   fragmentation. On the other hand, tandem repeats are repetitive sequences arranged head-to-tail or

81   head-to-head such as microsatellites and some multi-copy genes (e.g., ribosomal DNA and genes of

82   the Major Histocompatibility Complex, MHC). Reads shorter than the tandem repeat array will not

83   resolve the exact number of the repeat unit, resulting in the collapse of the region into fewer copies.

84   Some particular genomic regions enriched for repeats tend to be systematically missing or

85   underrepresented in traditional genome assemblies. These regions include: 1) telomeres at the

86   chromosome ends that are usually composed of microsatellites; 2) centromeres, essential for

87   chromosome segregation often specified by satellites that can be arranged in higher-order structures

88   like the alpha satellite in humans (Willard and Waye 1987) or by transposable elements in flies

89   (Chang et al. 2019); 3) multi-copy genes like MHC genes (Shiina et al. 2009); d) non-recombining

90   and highly heterochromatic chromosomes like the Y and W sex chromosomes (Chalopin et al. 2015;

91   Smeds et al. 2015; Hobza et al. 2017). As these regions play an essential role in the functioning and

92   evolution of genomes, the need to successfully assemble them is a pressing matter.

93

94   The other main limitation of traditional NGS methods is the shortcoming in reading regions with

95   extreme base composition (an enrichment of either A+T or G+C nucleotides), thus representing

3

96  another source of genomic dark matter. Extreme base composition mainly affects the last step of the

97  standard library preparation for Illumina sequencers that involves PCR amplification (Dohm et al.

98  2008; Aird et al. 2011). GC-rich regions tend to have higher melting temperatures than the rest of the

99  genome and are thus not as accessible with standard PCR protocols. On the other side of the spectrum,

100  AT-rich regions are also challenging to be amplified with standard PCR conditions and polymerases

101  (Oyola et al. 2012) because they require lower melting and extension temperatures (Su et al. 1996).

102  Several protocols have been developed to help minimize the phenomenon of GC-skewed coverage

103  (uneven representation of GC-rich regions), including PCR-free library preparation (Kozarewa et al.

104  2009) and isolation of the GC-rich genomic fraction prior to sequencing (Tilak et al. 2018).

105  Nonetheless, there is no single method that entirely solves base composition biases of short-read

106  sequencing and gives a homogeneous representation of the genome (Tilak et al. 2018). As a result,

107  extremely GC-rich or AT-rich regions may not be assembled at all.

108

109  It is essential to be aware of technological biases and genome assembly incompleteness during project

110  design since these can affect downstream analysis and mislead biological interpretations (Thomma et

111  al. 2016; Weissensteiner et al. 2017; Domanska et al. 2018; Peona et al. 2018). For example, GC-

112  skewed coverage is particularly important in birds, where 15% of genes are so GC-rich that they are

113  often not represented in Illumina-based genome assemblies (Hron et al. 2015; Botero-Castro et al.

114  2017). Whether these genes are truly missing or mostly hiding due to technological limitations is still

115  debated (Lovell et al. 2014, Botero-Castro 2017). However the "missing gene paradox" in birds is a

116  clear example of how sequencing technologies can shape our view of genome evolution. Furthermore,

117  some GC-rich sequences can form non-B DNA structures, i.e., alternative DNA conformations to the

118  canonical double helix such as G-quadruplexes (G4). G4 structures are a four-stranded DNA/RNA

119  topologies that seem to be involved into numerous cellular processes, such as regulation of gene

120  expression (Du et al. 2008; Du et al. 2009; Raiber et al. 2011), genetic and epigenetic stability

121  (Schiavone et al. 2014), and telomere maintenance (Biffi et al. 2012). On the repetitive element side,

122  for example, transposable elements are a major target of epigenetic silencing (Law and Jacobsen

123  2010) that may influence the epigenetic regulation of nearby genes (Cowley and Oakey 2013; Chuong

124  et al. 2016; Tanaka et al. 2019). The epigenetic effect of transposable elements may be beneficial or

125  deleterious, but in either case it is important to acknowledge their potential involvement in the

126  evolution of gene expression (Lerat et al. 2019). More generally, repetitive elements can play

127  important roles in many molecular and cellular mechanisms, and as a source of genetic variability

128  (Bourque et al. 2018). They have contributed to evolutionary novelty in many organismal groups, by

129  giving rise to important evolutionary features like the mammalian placenta (Emera and Wagner

130  2012), the vertebrate adaptive immune system (Kapitonov and Koonin 2015; Zhang et al. 2019) and

131  other telomere repair systems (Levis et al. 1993; McGurk et al. 2019). Thus, having genome

132  assemblies that are as complete as possible facilitates research into a multitude of molecular

133  phenomena (Slotkin 2018).

134

135  To achieve more complete genomes, we need new technologies. Recently, long-read single-molecule

136  sequencing technologies with virtually no systematic error profile (Eid et al. 2009) have led to more

137  complete and contiguous assemblies (English et al. 2012; Loomis et al. 2013; Pettersson et al. 2019;

138  Smith et al. 2019). To date two sequencing strategies have been developed that produce very long

139  reads from single-molecules: 1) Pacific Biosciences (PacBio) SMRT sequencing, in which the

140  polymerases incorporate fluorescently labelled nucleotides and the luminous signals are captured in

141  real time by a camera; 2) Oxford Nanopore Technologies, which sequences by recording the electrical

142  changes caused by the passage of the different nucleotides through voltage sensitive synthetic pores.

143  These new sequencing techniques have already yielded numerous highly contiguous *de-novo*

144  assemblies (Faino et al. 2015; Gordon et al. 2016; Seo et al. 2016; Bickhart et al. 2017; Weissensteiner

145  et al. 2017; Michael et al. 2018; Yoshimura et al. 2019) and helped improving the completeness of

146  existing ones (Chaisson et al. 2014; Jain et al. 2018), as well as characterizing complex genomic

147    regions like the human Y centromere and MHC gene clusters (Rhoads and Au 2015; Westbrook et

148    al. 2015; Jain et al. 2018; Sedlazeck et al. 2018).

149

150    However, resolving entire chromosomes remains a difficult endeavour even with single-molecule

151    sequencing (except for small fungal and bacterial genomes (Ribeiro et al. 2012; Thomma et al. 2016)).

152    Even though no single technology is able to yield telomere-to-telomere assemblies, it is still possible

153    to bridge separate contigs into scaffolds using long-range physical data and obtain chromosome-level

154    assemblies. Scaffolding technologies are becoming more and more commonly used (Vertebrate

155    Genome Project ; Dudchenko et al. 2017; Belser et al. 2018; Deschamps et al. 2018; Li et al. 2019;

156    Wallberg et al. 2019). The two most common ones are linked-reads (Weisenfeld et al. 2017) and

157    proximity ligation techniques (reviewed in Sedlazeck et al. (2018)). Linked-read libraries are based

158    on a system of labelling reads belonging to a single input DNA molecule with the same barcode

159    (Weisenfeld et al. 2017). In this way, using high molecular weight DNA allows to connect different

160    genomic portions (contigs) that may be distantly located but physically part of the same molecule.

161    High-throughput proximity ligation techniques as Hi-C and CHiCAGO are able to span very distant

162    DNA regions by sequencing the extremities of chromatin loops that could be up to Megabases apart

163    in a linear fashion (for more details see Lieberman-Aiden et al. (2009)). While Hi-C is applied directly

164    on intact nuclei, the CHiCAGO protocol reconstructs chromatin loops *in-vitro* from extracted DNA.

165    All these libraries are then sequenced on an Illumina platform. As linked reads and proximity ligation

166    techniques are becoming more and more popular used nowadays, we also implement and test them

167    in the present study.

168

169    Although a plethora of new sequencing technologies and assembly methods are currently being

170    successfully implemented, it remains unclear how they complement each other in the assembly

171    process. Here we address these assembly and knowledge gaps using a bird as a model. Bird genomes

172    represent a promising target to investigate that as their genomic features make it relatively easy to

173    assemble most parts with the exception of few complex regions per chromosome. In fact, the typical

174    avian genome is characterized by a small genome size (mean of ~1 Gb Kapusta and Suh (2017);

175    Gregory (2019)) and low overall repeat content (about 10% overall, with the exception of

176    woodpeckers that have 20% (Kapusta and Suh 2017). However, there are gene-rich and GC-rich

177    microchromosomes (Burt 2002; Griffin and Burt 2014; Miller and Taylor 2016) as well as a highly

178    repetitive W chromosomes (at least in non-ratite birds Zhou et al. (2014); Smeds et al. (2015); Bellott

179    et al. (2017)) that are still difficult to assemble.

180

181    In this study, to understand which genomic sequences are missing in regular draft genome assemblies

182    with respect to a high-quality and curated assembly, we generated several draft *de-novo* genomes and

183    a reference genome for the same sample of the paradise crow (*Lycocorax pyrrhopterus*, 'lycPyr').

184    The paradise crow is a member of the birds-of-paradise family (Paradisaeidae), one of the most

185    prominent examples of an extreme phenotypic radiation driven by strong sexual selection, and as

186    such, a valuable system for the study of speciation, hybridization, phenotypic evolution and sexual

187    selection (Shedlock et al. 2004; Irestedt et al. 2009; Ligon et al. 2018; Prost et al. 2019; Xu et al.

188    2019). We sequenced one female paradise crow individual with all the technologies that worked with

189    a DNA sample of mean 50 kb molecule length. We combined short, linked, and long-read libraries

190    together with Hi-C and CHiCAGO proximity ligation maps into a multiplatform reference assembly.

191    All these technologies permitted us to curate the resulting assembly by controlling for consistency

192    between multiple independent data types and make majority rule decision in conflicting cases. The

193    curated assembly enabled us to: 1) demonstrate the feasibility of obtaining a high-quality assembly

194    of a non-model organism with limited sample amount and non-optimal sample quality (a situation

195    that empiricists commonly face); 2) identify which genomic regions are actually gained from

196    combining technologies compared to draft assemblies of each individual technology; 3) assess the

197    strengths and weaknesses of the implemented technologies regarding the efficiency of assembling

198    difficult repeats and GC-rich regions; and 4) quantify how technologies can widen or limit the study

7

199    of specific genomic features (e.g., TEs, satellite repeats, MHC genes, non-B DNA structures), thus

200    providing a roadmap to investigate them.

201

## 202    Results

203

204    We leveraged the power of data generated from multiple sequencing approaches for the same sample

205    of paradise crow to generate a gold-quality assembly and to assess limitations of regular draft

206    genomes based on any single technology. Briefly, we combined short, linked and long reads with

207    proximity-ligation data to obtain a high-quality assembly despite the limitations of a non-model

208    organism such as limited sample amount and non-optimal quality. For each sequencing technology,

209    we produced an independent *de-novo* assembly. These assemblies were compared using majority-

210    rule decisions by manually curating the final assembly. Finally, the multiplatform assembly was

211    compared to each *de-novo* version to assess the amount of repeats and other complex regions

212    previously missing from the individual assemblies. We then evaluated the completeness of each

213    assembly using a variety of different metrics, including established scores such as BUSCO,

214    contig/scaffold N50, LTR Assembly Index and new metrics like overall repeat content, number of

215    MHC IIB exons, GC and G4 content, as well as number and nature of gaps.

216

217    **Long and short read *de-novo* assemblies**

218    In order to compare the efficiency of short, linked, and long reads, we produced independent draft

219    assemblies for each of the different sequence libraries. One draft genome assembly of *L. pyrrhopterus*

220    based on short reads (Illumina) is already available from Prost et al. (2019) ('lycPyrIL'; **Table 1**).

221    For the present study, we produced two linked-read libraries (10X Genomics Chromium) from which

222    we assembled two draft genomes ('lycPyrSN1' and lycPyrSN2'; where 'SN' stands for Supernova)

223    and a PacBio library from the same paradise crow sample that generated the primary assembly

224    'lycPyrPB' (**Table 1** and **Methods** section). In total, four independent *de-novo* assemblies were

225    generated.

226    We first evaluated the completeness of these assemblies by assessing their fragmentation, contig and

227    scaffold N50 and by counting the number of core genes present with BUSCO (Nishimura et al. 2017;

228    Waterhouse et al. 2017). In terms of fragmentation, the PacBio primary assembly ('lycPyrPB')

229    consisted of about 3,000 contigs, while lycPyrIL had ~3,000 scaffolds, and the 10XGenomics

230    assemblies had about ~14,000 scaffolds (**Table 1**). The short and linked-read assemblies all had a

231    scaffold N50 of about 4 Mb while the PacBio assembly had a contig N50 of 6 Mb (**Table 1,**

232    **Supplementary Table S1**). Notably, there is a 10-times higher of contig N50 in lycPyrPB relative to

233    the lycPyrIL assembly, indicating significant improvement in assembly continuity in the PacBio vs.

234    Illumina assembly. Next, we used the BUSCO tool (Nishimura et al. 2017) to identify correctly

235    assembled core genes (percentage of only single-copy and complete genes follow): lycPyrIL 93.8%,

236    lycPyrSN1 92.5%, lycPyrSN2 91.5%, lycPyrPB 84.8% prior to any assembly polishing

237    (**Supplementary Table S2**). Similarly, we estimated genome completeness and quality of the

238    intergenic and repetitive sequences with the LTR Assembly Index (LAI, Ou et al. (2018)). This index

239    is calculated as the proportion of full-length LTR retrotransposons over the total length of full-length

240    LTR retrotransposons plus their fragments. LAI could only be calculated for lycPyrPB since the other

241    *de-novo* assemblies did not have enough complete LTR elements for the algorithm to work. lycPyrPB

242    has an LAI score of 11.89, which is typical of a reference-quality assembly (Ou et al. 2018), and

243    higher than chicken (galGal5, RefSeq accession number GCF_000002315.6; Bellott et al. (2017))

244    with an LAI score of 7.54. We cannot exclude that the higher score in paradise crow is caused by

245    biological differences in LTR load between the species. More details about the LAI score distribution

246    across chromosomes and genomes are found in **Supplementary Table S3, Supplementary Figure**

247    **S1** and **Supplementary Figure S2**.

248

**Table 1.** Draft and multiplatform assemblies generated for the paradise crow. For each assembly the sequencing technology and software used to produce them are shown together with contig N50, scaffold N50 and the number of gaps.

| Assembly | Technology | Software | Contig N50 (bp) | N contigs | Scaffold N50 (bp) | N scaffolds | N gaps[a] | Missing assembly[b] (%) |
|---|---|---|---|---|---|---|---|---|
| **lycPyrIL** | Illumina HiSeq2500 (PE + MP)[c] | ALLPATHS-LG | 620,719 | 10,766 | 4,227,710 | 3,216 | 14,573 | 3.82 |
| **lycPyrPB** | PacBio RSII C6-P4 | Falcon | 6,644,420 | 3,422 | - | - | - | 0.45 |
| **lycPyrSN1** | 10X Genomics Chromium HiSeqX | Supernova2 | 144,856 | 29,791 | 4,360,585 | 13,934 | 21,550 | 4.53 |
| **lycPyrSN2** | 10X Genomics Chromium HiSeqX | Supernova2 | 149,640 | 27,366 | 4,748,626 | 14,217 | 20,131 | 2.62 |
| **lycPyrHiC** | PacBio + Phase Genomics Hi-C | Proximo | 6,644,420 | 3,422 | 70,588,898 | 2,927 | 533 | 0.45 |
| **lycPyrILPB** | lycPyrIL + gap-filling with PacBio | PBJelly | 1,982,606 | 6,895 | 4,229,628 | 3,216 | 10,422 | 3.03 |
| **lycPyr2** | PacBio + Dovetail CHiCAGO | HiRise | 6,294,665 | 3,463 | 6,644,037 | 3,227 | 282 | 0.45 |
| **lycPyr3** | lycPyr2 + 10X Genomics | ARCS + LINKS | 6,294,665 | 3,463 | 8,009,555 | 3,121 | 345 | 0.27 |
| **lycPyr4** | lycPyr3 + Phase Genomics Hi-C | Proximo | 6,294,665 | 3,463 | 69,071,023 | 1,713 | 1,791 | 0.27 |
| **lycPyr5** | lycPyr4 + manual curation with alignments + gap filling | PBJelly | 7,540,011 | 3,269 | 74,173,823 | 1,700 | 1,631 | 0.001 |
| **lycPyr6** | lycPyr5 + manual curation with Hi-C | Juicer | 7,540,011 | 3,271 | 74,173,823 | 1,700 | 1,635 | - |

[a] The number of gaps is estimated as the count of stretches of N nucleotides within a scaffold.
[b] The percentage of incompleteness is relative to the final version of the multiplatform assembly: (1 − (assembly size / final assembly size)) * 100. The N nucleotides are excluded from the calculation.
[c] PE: paired end reads; MR: mated reads.

249

10

**The multiplatform reference assembly**

To generate a high-quality genome assembly, we combined five technologies (short, linked, and long reads in addition to a CHiCAGO and Hi-C proximity ligation maps) into one multiplatform assembly. This process was divided into 9 steps (**Figure 1**), described in further detail in the **Methods** section.

First, we assembled the PacBio long reads into the primary assembly (lycPyrPB; 3,442 contigs) and it was scaffolded and corrected for misassemblies with the Dovetail CHiCAGO map ('lycPyr2'; **Figure 1a-b**). The scaffolding software HiRise introduced 98 breaks and made 293 joins of scaffolds (gaps of 100 bp were introduced at this stage), as well as closed 11 gaps between contigs and resulted into an assembly of 3,227 scaffolds (**Table 1** and **Supplementary Table S1**). Subsequently we polished the assembly with long reads (two rounds of Arrow; Chin et al. (2016)) and short reads (two rounds of Pilon; Walker et al. (2014); **Figure 1c**).

We then continued to scaffold lycPyr2 with two types of long-range information in order to get a chromosome-level assembly. First, we used 10X Genomics linked reads (SN1 library; 24 kb mean molecule length; **Figure 1d**) that encode medium-range spatial information that placed 235 contigs into 131 new scaffolds. Of these new scaffolds we kept only 88 and discarded potential chimeric scaffolds, which were identified by being composed of sex-linked contigs and autosomal ones (based on male/female short-read coverage; see **Methods**). We then confirmed the chimeric nature of such scaffolds by constructing an additional assembly based on scaffolding lycPyrPB with the Hi-C map ('lycPyrHiC'; **Table 1**). Phase Genomics Hi-C, i.e., 3D chromatin conformation data, can bridge sequences megabases apart (Burton et al. 2013) and theoretically reconstruct entire chromosomes (Hi-C super-scaffolds). In this way lycPyrHiC represented a second independent verification of the collinearity or chimeric nature of the contigs. Accordingly, we checked whether the contigs resided on different Hi-C super-scaffolds. Once we removed the chimeric contigs, we obtained 'lycPyr3' that contained a total of 3,121 scaffolds. Secondly, we scaffolded lycPyr3 with Phase Genomics Hi-C and

11

276 obtained 38 super-scaffolds ('lycPyr4'; **Figure 4e**) that harboured 1,446 contigs/scaffolds and

277 accounted for 97% of the assembly, while 1,675 contigs/scaffolds remained unplaced (3%). As most

278 of these super-scaffolds (32 out of 38) correspond to entire chromosomes of other avian species, we

279 call them "chromosome models". Examining the post-scaffolding Hi-C heatmap, we found that

280 chromosomes 1 and 2 were split into two Hi-C super-scaffolds, respectively. Therefore, following

281 the high level of Hi-C interaction between these super-scaffold pairs in the heatmap (**Supplementary**

282 **Figure S3**), we manually combined the respective super-scaffold pair into one chromosome model

283 (see **Methods**); the assembly thus resulted in 36 chromosome models.

284

285 We proceeded to further manually curate the chromosome models by looking for misassemblies

286 (**Figure 1f**) and used long reads for gap-filling (**Figure 1g**). We corrected fine scale orientation issues

287 of contigs within scaffolds through whole genome alignments (see **Figure 2** and **Methods**) and

288 corrected more orientation, order issues and erroneous chromosomal translocations through the

289 inspection of Hi-C heatmaps (see **Figure 1i** and **Methods**). We first corrected 43 misassemblies by

290 aligning the draft genomes and three outgroups to lycPyr4 (**Figure 2** and **Methods**). Next, we

291 extended contig ends and filled scaffold gaps with long reads using PBJelly ('lycPyr5'). PBJelly filled

292 106 gaps, extended 56 gaps on both ends and extended only one end of 292 gaps (**Supplementary**

293 **Table S4**). Finally, we further checked for misassemblies with the help of the Hi-C data. We

294 generated a Hi-C heatmap of lycPyr5 with Juicer (Durand et al. 2016) and detected misassemblies

295 though the visual inspection of such a map with JuiceBox (Dudchenko et al. 2018) following the

296 indications given by (Lajoie et al. 2015) and (Dudchenko et al. 2018). The Hi-C heatmap showed

297 mostly orientation and ordering problems within lycPyr5 (**Supplementary Figure S4**) that can be

298 identified from the ribbon-like patterns in the interaction map (Dudchenko et al. 2018). Finally, the

299 map highlighted the misplacement of two contigs between chromosome models (**Supplementary**

300 **Figure S4**). In total 76 misassemblies were corrected to generate the final assembly ('lycPyr6') with

301 a scaffold N50 of ~75 Mb (**Table 1**).

12

302

303   In parallel to the assembly of lycPyr6, we also generated a simpler multiplatform assembly by gap-

304   filling the Illumina primary assembly (lycPyrIL) with PacBio reads ('lycPyrILPB'). PBJelly was used

305   to gap-fill the Illumina assembly and successfully closed 4,151 gaps, reducing the total number of

306   gaps from 14,573 to 10,422. It also double extended 418 gaps and single extended 2,597 gaps

307   (**Supplementary Table S4**). The numbers of scaffolds and scaffold N50 did not significantly change

308   from lycPyrIL (**Table 1**).

309

310   **Chromosome models: macrochromosomes, microchromosomes and sex chromosomes**

311   We obtained 36 chromosome models comprised of 16 macrochromosome models, 18

312   microchromosome models and two sex chromosome models. All the macrochromosome models

313   showed homology to chicken chromosomes and were named after their homologous counterparts.

314   The same applies for 12 of 18 microchromosomes, while the remaining 6 showed no homology with

315   chicken chromosomes and therefore were tentatively named as unknown chromosomes "chrUN1-6".

316   The chromosomes homologous to chicken are mostly syntenic with respect to chicken with few

317   exceptions. In fact, chicken chromosome 1 and 4 are split in two in Passeriformes and correspond,

318   respectively, to chromosome 1 and 1A, and chromosome 4 and 4A (Kapusta and Suh 2017).

319   The Z and W sex chromosome models had an assembled size of 73.5 Mb and 21.4 Mb, respectively,

320   and were comparable to chicken (82 Mb and 7 Mb, galGal6a, RefSeq accession number

321   GCF_000002315.6; Bellott et al. (2017)). Z and W models were also largely consistent with the sex-

322   linked contigs previously identified using male/female coverage comparisons (**Supplementary**

323   **Table S5** and **Methods**), only 3.11 Mb of the W and 3.99 Mb of the Z chromosome were contigs not

324   previously identified as sex-linked. Finally, the pseudoautosomal region (PAR) seemed to be

325   fragmented into two parts. We identified two contigs that are homologous to the PAR of flycatcher;

326   one of them was placed by Hi-C onto the Z while the other was placed onto the W chromosome model

327   (**Supplementary Table S5**). While the Z chromosome showed a repetitive content similar to the

13

328    autosomes (~10%), the W was extremely repeat-rich (~70%, **Figure 3a, Supplementary Table S6**).

329    The dotplots of the alignments of the paradise crow sex chromosomes with the chicken sex

330    chromosomes (**Supplementary Figure S5 and Supplementary Figure S6**) showed that the two Z

331    chromosomes had a high level of synteny and collinearity while the repetitiveness of the two W

332    chromosomes made it difficult to identify shared single-copy regions other than very small ones. The

333    sex chromosomes were also easily identified in the post-clustering Hi-C heatmap (**Supplementary**

334    **Figure S3**), as their hemizygosity can be expected to result in roughly half of the amount of Hi-C

335    interactions (calculated as the frequency of shared paired-end reads between contigs/scaffolds) within

336    each chromosome model and with the other chromosome models.

337    Finally, the LTR Assembly Index calculated on the single chromosomes yielded high scores (min 0

338    on chromosome 10, mean 13.14, max 21.41 on chromosome W) that have been suggested to be

339    indicative of reference and gold-quality assemblies (Ou et al. (2018), **Supplementary Figure S1** and

340    **Supplementary Table S3**).

341

342    **GC content and G4 motif prediction**

343    GC-rich regions are commonly underrepresented in traditional NGS assemblies because of the

344    aforementioned GC-skewed coverage phenomenon (see **Introduction**). Comparing the different *de-*

345    *novo* assemblies, we noticed that indeed lycPyrPB showed more GC-rich regions (54,532 windows

346    of 1 kb size with GC > 58.8%) with respect to lycPyrIL, SN1 and SN2 (45,966, 45,720 and 52,080

347    such windows, **Figure 3b**, **Supplementary Table S7**, **Supplementary Figure S7**). Thus, lycPyrSN1

348    shared a similar number of GC-rich regions while lycPyrSN2 was closer to lycPyrPB

349    (**Supplementary Figure S7**, **Supplementary Table S7**).
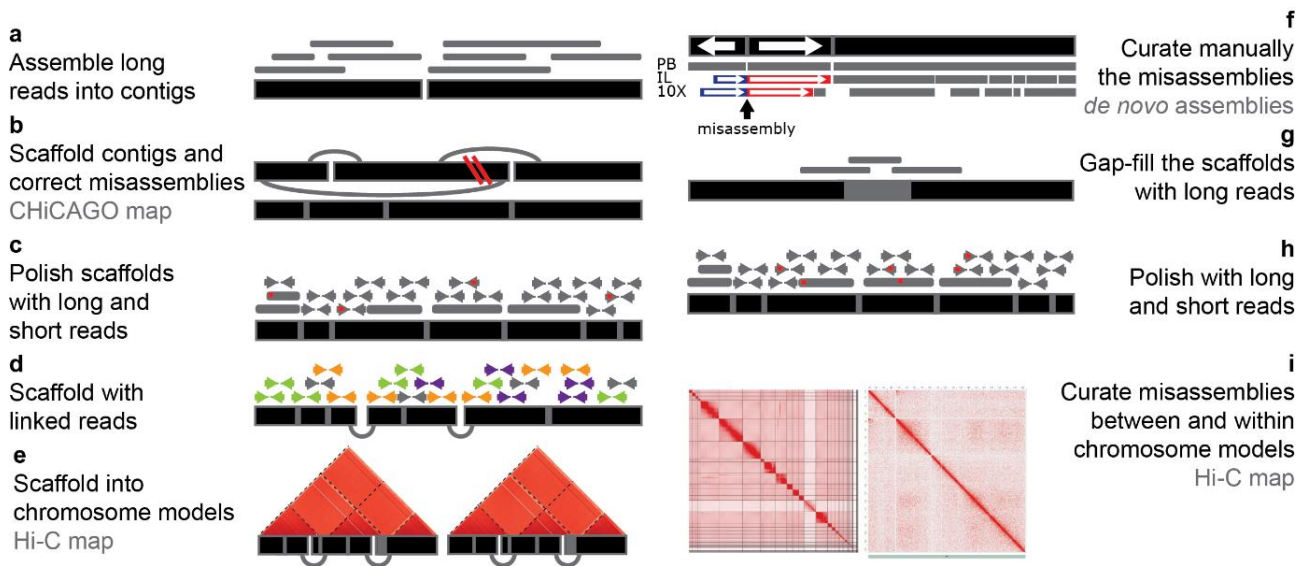
350

351    Since GC-rich regions may form G-quadruplexes motifs and structures (G4), we expected the

352    depletion of GC-rich short reads to limit the representation of G4 motifs in short read assemblies.

353    Conversely, we expected G4 motifs to be more abundant in long read assemblies, since these have
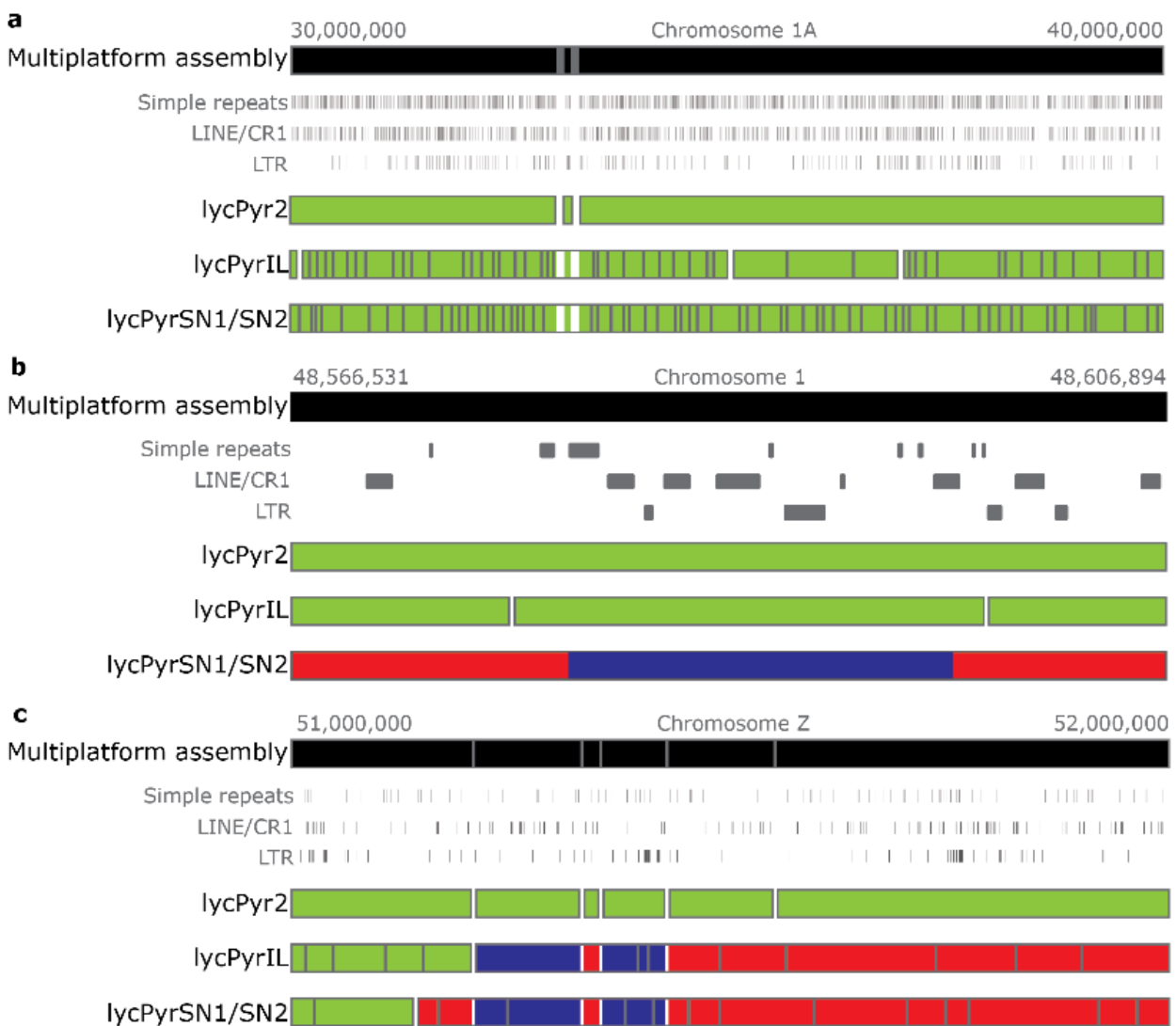
14

354 been suggested to be virtually free from sequence-based biases (Eid et al. 2009). To test this, we

355 predicted the presence of G4 motifs using Quadron (Sahakyan et al. 2017) in all the different

356 assemblies. All the *de-novo* Illumina-based assemblies had fewer predicted G4 sites the PacBio

357 assemblies (**Figure 3c** and **Supplementary Table S8**). lycPyrSN2 and lycPyrIL had 7.3 and 7.5 Mb

358 (169,214 and 166,602 motifs) occupied by G4 sequences and about 1.6 Mb or 24,000 motifs less than

359 lycPyr6 (9.1 Mb, 193,248 motifs). lycPyrSN1 was the assembly with the fewest G4 motifs predicted

360 (6.5 Mb, 149,275 motifs). The PacBio primary assembly lycPyrPB had 8.42 Mb of predicted G4,

361 which was slightly higher in lycPyr2 after the correction with Dovetail CHiCAGO (8.43 Mb; **Figure**

362 **3c** and **Supplementary Table S8**). In the final assembly lycPyr6, G4 motifs were more present on

363 microchromosomes than on macrochromosomes (**Figure 3d**).

364



365 **Figure 1**. Overview of the multiplatform assembly process. (**a**) Long reads were assembled into
366 contigs. (**b**) The primary assembly was corrected and scaffolded using long-range information
367 provided by the CHiCAGO proximity ligation map. (**c**) The assembly was then polished from base-
368 calling errors with both short and long reads and (**d**) further scaffolded with linked-reads. (**e**) The
369 scaffolds are ordered and oriented into chromosome models according to the Hi-C proximity
370 ligation map. (**f**) The chromosome models were aligned to the *de-novo* assemblies based only on
371 one single technology and then manually inspected to find misassemblies and correct them
372 following the majority rule (more details in **Figure 2** and **Methods**). PB: PacBio long-read
373 assembly; IL: Illumina short-read assembly; 10X: 10XGenomics linked-read assemblies (**g**) Long
374 reads were used to gap-fill the assembly and (**h**) to polish the final version together with short reads.
375 (**i**) Hi-C heatmaps were used to identify and correct misassemblies between and within chromosome
376 models.

15

377



**Figure 2**. Examples of the manual curation of the assembly (**step f** in **Figure 1**). The multiplatform assembly is aligned to the other *de-novo* assemblies from the same sample. The grey lines within the assemblies represent gaps between different contigs or scaffolds while the white lines represent gaps within the same scaffold. Green means that the contigs/scaffolds align to the reference in the same orientation for their entire length while red and blue highlight contigs/scaffolds that partially align in the forward (red) and reverse (blue) direction to the reference. **(a)** Here 10 Mb of chromosome 1A are shown that are in accordance with all the *de-novo* assemblies. Nonetheless, short-read based technologies yielded much more fragmented scaffolds. **(b)** Example of a scaffold orientation misassembly in the 10XGenomics assembly. The other two assemblies span the inverted region and both agree with the multiplatform assembly. **(c)** Example of how two different assemblies could help to identify which contigs have to be re-oriented and re-ordered in the final assembly. In lycPyrIL, lycPyrSN1 and lycPyrSN2 we had scaffolds than span the misoriented (blue) region and bridge it to contigs that showed concordant orientation with the multiplatform assembly. This indicated that we have only a small local inversion of two PacBio contigs
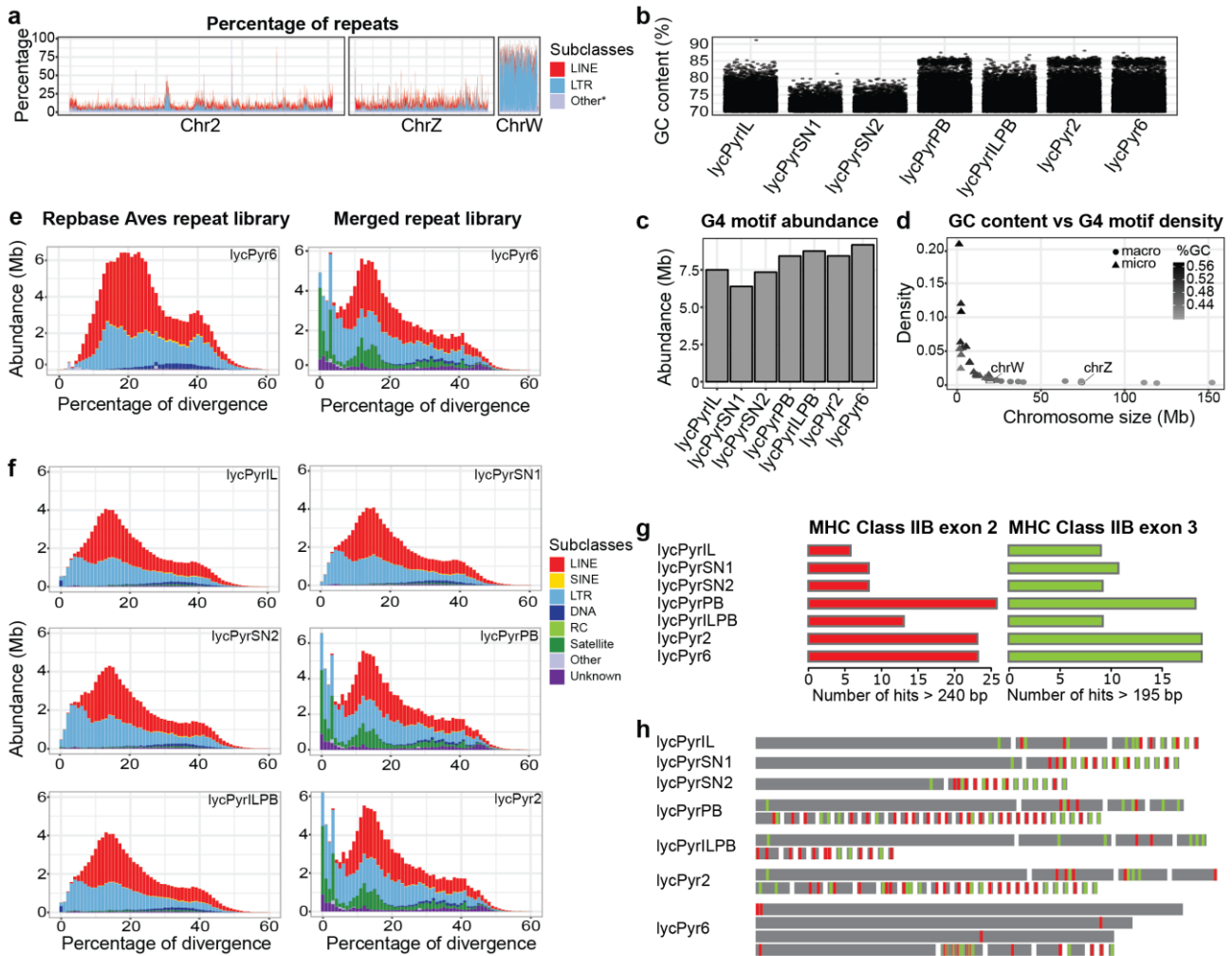
392

**Repeat library**

16

394   To obtain an in-depth annotation of interspersed and tandem repeats, the *de-novo* characterization of

395   repetitive elements and manual curation thereof are essential (Platt et al. 2016). We manually curated

396   a total of 183 consensus repeat sequences generated from lycPyrIL and lycPyrPB to have an optimal

397   repeat characterisation. In Prost et al. (2019) a total of 112 raw consensus sequences were produced

398   using RepeatModeler on three Illumina-based birds-of-paradise (*Astrapia rothschildii*, *L.*

399   *pyrrhopterus* and *Ptiloris paradiseus*; including lycPyrIL) but only the 37 most abundant from

400   lycPyrIL were manually curated. We then curated the remaining 75 and added 71 more *de-novo*

401   consensus sequences based on curated raw consensus sequences from RepeatModeler run on

402   lycPyrPB. Our new bird-of-paradise specific repeat library is now composed of the following

403   numbers of consensus sequences: 56 ERVK, 56 ERVL, 37 ERV1, 5 CR1, 4 LTR, 9 satellites, 2 DNA

404   transposons, 1 SINE/MIR, and 13 unknown repeats. All the consensus sequences curated for the three

405   species of birds-of-paradise (*L. pyrrhopterus, A. rothschildii, P. paradiseus*) are given in

406   **Supplementary Table S9**. Eventually, we merged birds-of-paradise consensus sequences together

407   with the Repbase Aves library, the flycatcher (Suh et al. 2018), the blue-capped cordon blue (Boman

408   et al. 2019) and the hooded crow libraries (Weissensteiner et al. 2019).

409

410   Custom and *de-novo* repeat libraries substantially improve the identification and masking of repeats

411   in genome assemblies (Platt et al. 2016). To quantify this effect for our assemblies, we compared a

412   general avian repeat library with our curated one. The custom library resulted in masking a higher

413   fraction of the genome in every assembly (**Figure 3e-f**). When comparing the masked fraction with

414   the custom library to the fraction masked with the Repbase library, we see that lycPyrIL, lycPyrILPB,

415   and lycPyrSN1 have 20% more masked repeats (from 78 Mb to 94 Mb), while lycPyrSN2 has 21.68%

416   (from 83 to 101 Mb), lycPyrPB 38% (from 87 Mb to 120 Mb), and lycPyr6 38% (from 88 Mb to 122

417   Mb; see **Figure 3e**, **Supplementary Table S10**). In particular, with the new library we were able to

418   identify 9.4 Mb of satellite DNA in the PacBio-based assemblies, while the standard Repbase avian

419   library identified only 1 Mb (**Figure 3e-f**, **Supplementary Table S10**). Relative to lycPyr6, most of

17

420    the satellites and unknown repeats remain unassembled in the short-read and linked-read assemblies

421    (**Figure 3f** and **Figure4b**).

422
423



424    **Figure 3**. **(a)** Comparison of the repeat content across chromosome 2 (representative of autosomes),
425    Z and W calculated as the percentage of repeats per window of 50 kb. Here LINE and LTR are
426    shown as major components of the mobile element repertoire and all the other types of repeats are
427    merged into the "Other*" category. **(b)** Distribution of GC-content per window (10 kb) across
428    assemblies on the left side of the violin plots. GC-content distribution of the windows containing
429    G4 motifs on the right side of the violin plots. **(c)** G4 motif abundance across different paradise
430    crow assemblies. **(d)** G4 motif density across the chromosome models of the final assembly; the
431    chromosomes are arranged by size; macrochromosomes are coloured in light grey while
432    microchromosomes (smaller than 20 Mb) are shown in dark grey. The density distribution of G4 in
433    micro and macro chromosomes was statistically different (t-test p-value: 0.01). **(e)** Repeat landscape
434    of lycPyr6 masked with the Repbase Aves repeat library (on the left) and masked with the custom
435    library produced in this study which also included the Repbase Aves library (on the right). **(f)**
436    Repeat landscapes of the four *de-novo* assemblies of the paradise crow masked with the custom
437    repeat library. **(g)** Abundance of MHC class IIB exon 2 and exon3 in the different paradise crow
438    assemblies. **(h)** Schematic visualization of the instances of MHC class IIB exon 2 (red) and 3
439    (green). Each black rectangle represents a different contig or scaffold.

440

**MHC class IIB analysis**

441

442 In birds, the multi-copy gene family of the major histocompatibility complex (MHC) is arranged as

443 a megabase long tandem repeat array (Miller and Taylor 2016). Since we expect it to be even more

444 difficult to correctly assemble than the aforementioned interspersed repeats (O'Connor et al. 2019),

445 it represents a prime candidate region for measuring the quality of an assembly.

446

447 We used the presence of entire copies of the second (most variable) and third (more conserved) exons

448 of the MHC class IIB as proxies of assembly quality (Hughes and Yeager 1998). Overall, we found

449 that short-read assemblies had fewer MHC gene copies than long-read assemblies (**Figure 3g-h**),

450 while linked-read assemblies performed better than Illumina alone. Regarding exon 2 (**Figure 3g**),

451 PacBio retrieved 26 copies while Illumina and 10XGenomics assembly only hold 6-8. However, it is

452 worth noting that after correcting lycPyrPB with the Dovetail CHiCAGO map, 3 copies were lost

453 (not detectable as full-length exons anymore) and were not restored by the subsequent steps of

454 sequence corrections and curation. The results were similar for exon 3 (**Figure 3g**): PacBio

455 assemblies retrieved 18-19 copies while the other technologies retrieved only 9-11 copies. In this case

456 we see that the molecule input length of 10XGenomics library has an effect on the assembly of these

457 genes, where the library with shorter molecule length had assembled more copies than the longer one

458 (11 vs 9 exon 2 copies; **Figure3g-h**). On the other hand, while Dovetail CHiCAGO prevented the

459 identification of some exon 2 copies, it increased the number of assembled copies of exon 3.
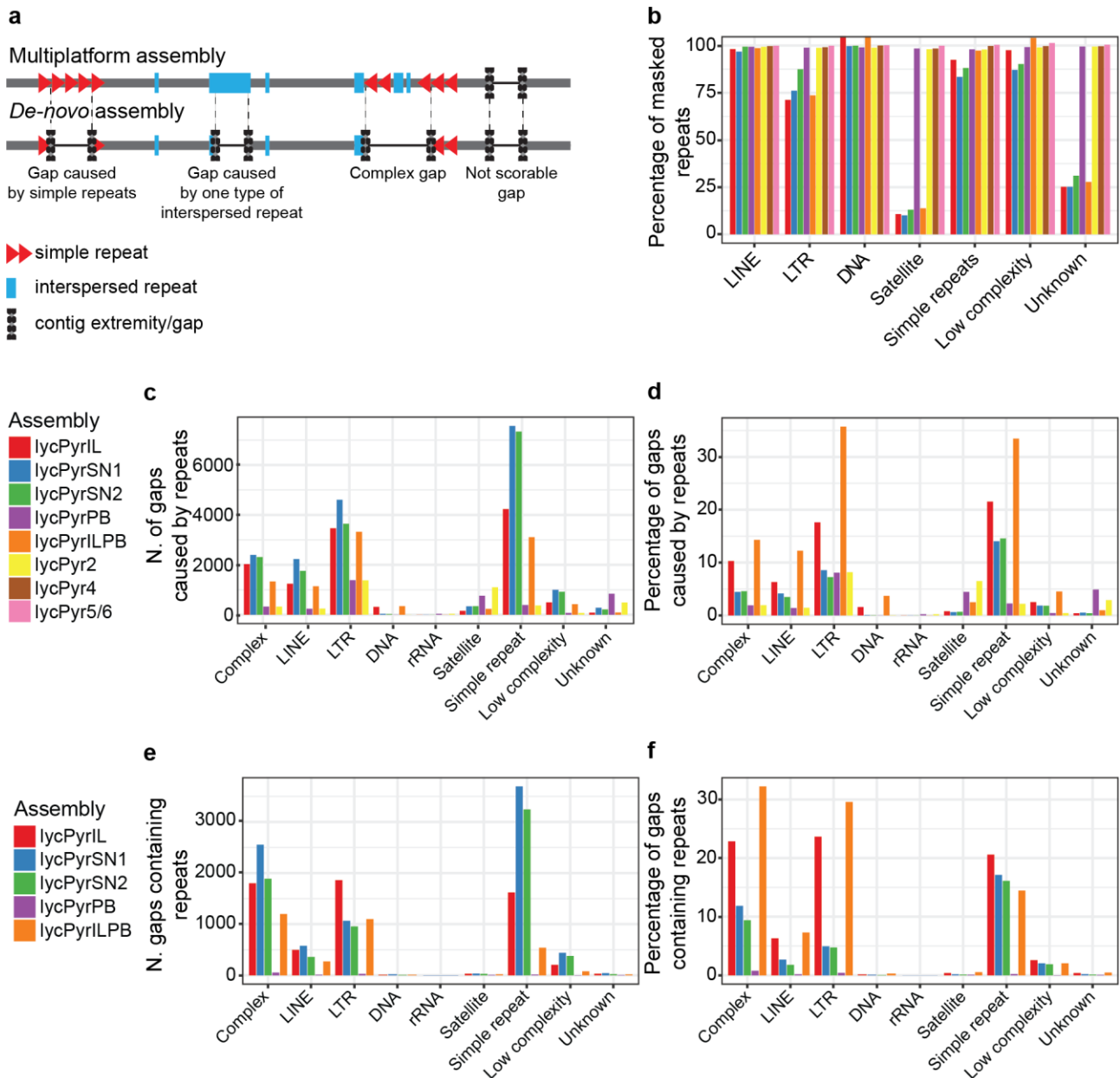
460

461 **Gap analysis**

462 The process of scaffolding links contigs together without adding any information about the missing

463 DNA between them, but it is possible to use long reads to fill those gaps. For this we utilized PBJelly

464 (English et al. 2012) to extend and bridge contigs in the assembly by locally assembling PacBio reads

465 to the contig extremities. Once the software finds reads aligned to the contig extremities, the

466 extremities can be: 1) extended on one or both sides to reduce the gap length, 2) extended and bridged

19

467    to fill the entire gap, 3) extended over the length of the gap without being ultimately bridged

468    (overfilled). PBJelly extended the extremities of 348 gaps, closed 116 gaps and overfilled 236 gaps

469    (**Supplementary Table S4**). This gap-filling step added a total of 2.96 Mb to the assembly. All the

470    sequences that were extended or gap-filled were more GC-rich (40%-89%, mean 58%) than the

471    average GC content of 40% and 2865 G4 motifs were added for a total of 171 kb. Only 800 kb of the

472    2.96 Mb added were repetitive elements; specifically, ~400 kb of LTR elements were added, 120 kb

473    of LINE, 142 kb of satellite DNA and 90 kb of simple and low complexity repeats (**Supplementary**

474    **Table S4**).

475

476    Furthermore, we investigated the causes of assembly fragmentation in several assemblies by

477    analysing the immediate adjacency of repetitive elements to the gaps (lower part of **Figure 4a**). We

478    found that simple repeats were the major fragmentation cause in Illumina and 10XGenomics

479    assemblies, followed by LTR and LINE elements (**Figure 4c-d**). In contrast, PacBio gaps (lycPyrPB

480    and lycPyr2) seemed to be mainly caused by LTR elements and secondarily by satellites (**Figure 4c-**

481    **d**).

482

**Figure 4** Overview of the causes and content of gaps in the paradise crow assemblies by comparing all the assembly versions to the final version. **(a)** Schematic representation of how gaps were categorized based on the flanking regions and content. **(b)** Proportion of repeats present in each assembly version respect to the reference (lycPyr6). **(c)** Number of gaps caused by the major repeat groups. **(d)** Proportion of gaps caused by the major repeat groups. **(e)** Number of gaps that contain (map to) repeats. **(f)** Proportion of gaps that contain (map to) repeats.

Finally, we quantitatively and qualitatively assessed which repeats in the final multiplatform assembly lycPyr6 were collapsed as gaps in the draft assemblies (**Figure 4e**-**f**). Many gaps in the Illumina and 10XGenomics draft assemblies corresponded to complex regions consisting of multiple types of repetitive elements (**Figure 4e-f**). Among draft assembly gaps containing only a single type

21

495    of repeat in lycPyr6, most were caused by simple repeats, LTR retrotransposons, and LINE

496    retrotransposons in short-read and linked-read assemblies (**Figure 4e-f**).

497

498    **Discussion**

499

500    Assembling complete eukaryotic genomes is a complex and demanding endeavour often limited by

501    technological biases and assembly algorithms (Alkan et al. 2010; Sedlazeck et al. 2018). In the last

502    decade, NGS technologies defined the standard of genome assemblies. Although they provided an

503    unprecedented view on the structure and evolution of many coding regions (Zhang et al. 2014), short

504    reads hardly inform on the entire complexity of a genome (Thomma et al. 2016). Indeed, the

505    systematic absence from genome assemblies and the difficulty to characterize the nature of many

506    such genomic regions (e.g. centromeres, telomeres, other repeats and highly heterochromatic regions)

507    gave these "unassemblable" sequences the evocative name of genomic "dark matter" (Johnson et al.

508    2005; Weissensteiner and Suh 2019).

509

510    In this study, we demonstrated that a combined effort involving multiple state-of-the-art methods for

511    long-read sequencing and scaffolding yielded a high-quality reference for a non-model organism. We

512    showed that a multiplatform approach was highly successful in resolving elevated quantities of

513    genomic dark matter in respect to single-technology assemblies (regular draft assemblies) and thus

514    resulted in a much more complete assembly. In order to assess genome completeness we focused

515    mostly on the quantification and characterization of previously inaccessible regions within genomic

516    dark matter, such as large transposable elements, GC-rich regions, and the high-copy MHC locus.

517

518    We generated a *de-novo* multiplatform assembly of a female bird-of-paradise genome by combining

519    the cutting-edge technologies that are now being implemented in many assembly projects (Faino et

520    al. 2015; Gordon et al. 2016; Seo et al. 2016; Bickhart et al. 2017; Weissensteiner et al. 2017; Michael

22

521   et al. 2018; Yoshimura et al. 2019), namely Illumina short reads, 10XGenomics linked reads, PacBio

522   long reads and two proximity ligation maps with Dovetail CHiCAGO and Phase Genomics Hi-C.

523   The choice of using a bird-of-paradise is manifold. First, avian genomes are small among amniotes

524   and have an overall repeat content of 10%, which make most genomic regions relatively "easy" to

525   assemble. This has made it possible to focus on regions that are challenging to assemble in eukaryotic

526   genomes of any size and complexity, like the repeat-rich W sex chromosome, and the GC-rich

527   microchromosomes. Second, birds-of-paradise is a highly promising system for the study of

528   speciation, hybridization and sexual selection (Irestedt et al. 2009; Prost et al. 2019; Xu et al. 2019).

529   A gold standard genome for this family will consequently expose new possibilities for more in-depth

530   studies of the genomic evolution behind the spectacular radiation of birds-of-paradise.

531

532   By employing a multiplatform approach, we 1) could assemble a chromosome-level genome which

533   includes the W chromosome and several previously inaccessible microchromosomes (i.e.,

534   comparable to the chicken genome, so far the best avian genome available); 2) report that a substantial

535   proportion (up to 90%) of repeat categories like satellites and LTR retrotransposons are missing from

536   most types of *de-novo* assemblies (**Figure 3e-f**, **Figure 4b**); and 3) identify simple repeats and LTR

537   retrotransposons as the major causes of assembly fragmentation (**Figure 4c-d**).

538

539   **A chromosome-level assembly for a non-model organism**

540   Our final assembly comprises 36 chromosome models. This assembled chromosome number is

541   similar to the known karyotype of another bird-of-paradise species *Ptiloris intercedens* (36-38

542   chromosome pairs; Les Christidis, personal communication). Among these models, there are 16

543   macrochromosomes, 12 microchromosomes, and the Z and W sex chromosomes showing homology

544   to chicken chromosomes (galGal6a). The remaining 6 models do not share homology with known

545   chicken chromosomes (galGal6a) and they might be putatively uncharacterized microchromosomes.

546   Microchromosomes are known to be very GC-rich (Burt 2002) and indeed this trend is present in our

23

547   data as well (**Figure 3d**). Base composition can create biases during the sequencing process especially

548   when a PCR step is required for the library preparation (Dohm et al. 2008; Aird et al. 2011) thus

549   limiting the representation of GC-rich and AT-rich reads in the data. Although, long read sequencing

550   technologies like PacBio have reduced amplification-based biases to a minimum (Schadt et al. (2010)

551   but see Guiblet et al. (2018)), we could not assemble contiguous sequences for all

552   microchromosomes. Among the unknowns and unassembled chromosomes, chromosome 16 which

553   is one of the most complex avian chromosomes and also holds the MHC (Miller and Taylor 2016).

554   The absence of these chromosomes is likely explained by that they are by far the densest in G4 motifs

555   of all chromosomes (**Figure 3d**). Given that DNA polymerase tends to introduce sequencing errors

556   in the presence of G4 structures (Guiblet et al. 2018), it is tempting to think that the depletion of

557   microchromosomes from assemblies is not only due to GC content per se but also due to the potential

558   presence of non-B structures (like G4) that elevated GC content appears to correlate with.

559   Nonetheless, even with the extensive use of cytogenetics the last chicken assembly (galGal5; Warren

560   et al. (2017)) completely lacks 5 microchromosomes. It thus seems plausible that these chromosomes

561   need special efforts to be recovered.

562

563   One of the most surprising outcomes of this multiplatform approach is the successful assembly of the

564   highly repetitive W chromosome which turned out to be larger (assembly size 21 Mb) and more

565   repetitive than the chicken equivalent (assembly size 9 Mb; Bellott et al. (2017)). In both species, it

566   is likely that the assembled sequences cover the euchromatic portions of the W. Birds have a ZW sex

567   chromosome system where the female is the heterogametic sex and the female-specific W is

568   analogous to the mammalian male-specific Y chromosome. Comparable to the mammalian Y

569   (Charlesworth et al. 2000), the W chromosome is highly repetitive and difficult to assemble

570   (Weissensteiner and Suh 2019). Previous studies focusing on the repetitive content of the avian W in

571   chicken (Bellott et al. 2017) and collared flycatcher (Smeds et al. 2015) showed in both cases a repeat

572   density of about 50%. In our assembly of the paradise crow, we found the W chromosome to be even

24

573     more repetitive with a repeat density of ~70% and highly enriched for LTR retrotransposons (**Figure**

574     **3a** and **Supplementary Table S6**). Having assembled chromosomes is key to improve any genomic

575     analysis but studies on sex chromosome evolution in birds has so far been heavily biased towards Z

576     (Zhou et al. 2014; Yazdi and Ellegren 2018; Xu et al. 2019). With genome assemblies like the present,

577     it will be possible to improve reconstructions how the two sex chromosomes diverged. We can

578     already see that the W chromosome evolves rapidly (**Supplementary Figure S5**) via accumulation

579     of transposable elements and only few regions appear syntenic between paradise crow and chicken

580     W.

581

582     **How complete are genome assemblies?**

583     Previous studies (see for example Etherington et al. (2019); Paajanen et al. (2019)) have assessed the

584     efficiency of available sequencing technologies in genome assembly and genome completeness

585     mainly through summary statistics like scaffold N50 and BUSCO. Scaffold N50 indicates the

586     minimum scaffold size among the largest scaffolds making up half of the assembly, while BUSCO

587     values measure the number of complete/incomplete/missing core genes in the assembly. However,

588     genome completeness goes beyond scaffold N50 and gene presence (Thomma et al. 2016; Domanska

589     et al. 2018; Sedlazeck et al. 2018). Genes usually occupy a small fraction of genomes and new

590     sequencing technologies commonly yield high N50 values. Therefore, these statistics have a very

591     limited scope in perspective of what the new sequencing technologies can achieve.

592

593     Although often being used as proxy of assembly quality, scaffold N50 is hardly meaningful in this

594     regard since it does not inform about the completeness and correctness of the assembled sequences.

595     If we order the scaffolds by decreasing size, scaffold N50 value can only reflect the fragmentation

596     level of the first half of the assembly regardless of whether the second half is made up of shorter

597     sequences. Finally, contig N50 should be used as a measure of contiguity, rather than scaffold N50,

598     as contig length measures sequences not interrupted by gaps.

599

600     Most of the currently available avian genomes score more than 94% of BUSCO gene completeness

601     (Peñalba et al. 2019) with various degrees of fragmentation, suggesting that it has become

602     straightforward to generate short-read assemblies with high BUSCO values. On the other hand,

603     BUSCO seems to be limited by the sequencing errors introduced by PacBio in the identification of

604     gene models (Watson and Warr 2019). Even with multiple rounds of error correction, BUSCO fails

605     to recognize genes that are actually present, at least partially, in the assembly (Watson and Warr

606     2019). Moreover, BUSCO seems to be trained and based on a set of core genes identified from Sanger

607     and Illumina assemblies. As such, BUSCO does not quantify genes in PacBio assemblies that were

608     previously missing in Illumina genomes, which would be needed for a fair genome completeness

609     comparison. This tendency is also evident from our results: for example after gap-filling lycPyrIL

610     with long reads, 10 genes were not detectable anymore in the resulting assembly lycPyrILPB

611     (**Supplementary Table S2**). A similar dynamic was observed also during the assembly process of

612     the superb fairy-wren *Malarus cyaneus* (Peñalba et al. 2019) where BUSCO values dropped with

613     long-read gap-filling but were restored after sequence polishing.

614

615     The new technologies have the potential to assemble very repetitive regions (e.g. MHC) and elusive

616     chromosomes (e.g., W and microchromosomes). For this reason, quality assessment should rely

617     upon measuring the efficiency in assembling difficult regions and not on those regions that we

618     already obtain with previous technologies. We therefore decided to measure genome completeness

619     and quality by characterising and quantifying repetitive regions.

620

621     Long reads were instrumental, not only to find and mask more repeats, but also to assemble and

622     discover previously overlooked repetitive sequences. In fact, by adding PacBio sequence data we

623     were able to significantly increase the number of predicted repeat subfamilies compared to the

624     repeat library previously built on three birds-of-paradise species (from 112 to 183 consensus

26

625    sequences; Prost et al. (2019)). These 71 new consensus sequences were only predicted by

626    RepeatModeler using the PacBio assembly, probably because the respective repeats were too

627    fragmented or assembled in too few copies in Illumina assemblies. A clear example is given by the

628    satellite DNA repeats that are severely depleted from both the lycPyrIL assembly (**Figure 3e-f**,

629    **Figure4b**) and from the previous repeat library. With our new repeat library we could increase the

630    base pairs masked by RepeatMasker by up to 38 % within the same assembly (lycPyr6). This

631    indicates that while longer read lengths are important for assembling repeats, only with a

632    comprehensive repeat library we can quantify their actual efficiency.

633

634    Repetitive elements are not only made up of transposable elements and satellite repeats, but also of

635    multi-copy genes. One of the most repetitive gene family is the Major Histocompatibility Complex

636    (MHC) involved in the adaptive immune response. In birds, MHC genes are located on one of the

637    most difficult chromosomes to assemble, namely chromosome 16 (Miller and Taylor 2016). We

638    recovered several scaffolds from this chromosome for which the only, though fragmented, assembly

639    exists from chicken (Warren et al. 2017). We counted how many MHC IIB copies we could retrieve

640    in the different assemblies, using BLAST hits to exon 2 and 3 sequences as proxy. We found the

641    maximum number of copies in lycPyrPB (**Figure 3g-h**) followed by lycPyr6, suggesting that the

642    misassembly correction with the CHiCAGO map affected the MHC genes, with the number of hits

643    of exon 2 decreasing and for exon 3 increasing. Short-read assemblies harbour fewer MHC IIB exon

644    copies but we note that 10XGenomics could assemble a couple more copies compared to standard

645    Illumina data. Moreover, lycPyrSN1 contained slightly more MHC genes than lycPyrSN2 assembled

646    with longer input molecule length.

647

648    As a further use of repetitive elements as quality measures, we tested the LTR Assembly Index

649    (LAI; Ou et al. (2018)) that assesses the quality of an assembly from the completeness of the LTR

650    retrotransposons present. It was not possible to obtain values for the Illumina and 10XGenomics

651 assemblies because the tool requires a certain baseline quantity of the full-length LTR assembled to

652 run as initial requirements. Nonetheless, both lycPyrPB and lycPyr6 show LAI scores (respectively

653 11.89 and 13.59, **Supplementary Table S3**, **Supplementary Figure S1**) typical for high-quality

654 reference genomes (as indicated in Ou et al. (2018)) and higher than those of chicken

655 (**Supplementary Figure S2**). The increase in LAI value from lycPyrPB and lycPyr6 indicates that

656 the assembly curation process, mostly gap-filling and polishing, improved the quality of the primary

657 assembly.

658

659 In addition to repetitive elements, base composition is the other main factor that limits completing

660 genome assemblies. We thus assessed the GC-content per window for each assembly (**Figure 3b**,

661 **Supplementary Figure S7**) and as expected, found more GC-rich windows in lycPyrPB compared

662 to the other *de-novo* assemblies (**Supplementary Figure S7**). High GC-content is often associated

663 with non-B DNA structures like G4 that have been shown to introduce sequencing errors during

664 polymerisation (Guiblet et al. 2018). We predicted the presence of G4 motifs in our assemblies

665 (**Figure 3c**) and Illumina and 10XGenomics assemblies have about 1.6-2.6 Mb less of G4 compared

666 to lycPyrPB. In this case, linked reads did not help to get a more complete overview of this genomic

667 feature respect to regular Illumina libraries. On the other hand, the overall curation from lycPyrPB to

668 lycPyr6 improved G4 prediction. G4 structures influence various molecular mechanisms such as

669 alternative splicing and recombination, therefore more complete assemblies make these regions

670 accessible for comparative genomic analysis.

671

672 **Strengths and limitations of sequencing technologies**

673 Nowadays, we have a plethora of sequencing technologies to choose from, each with their own

674 advantages and limitations. On top of that, the large number of assembly tools available and

675 hundreds of parameters to tweak makes it inevitable to produce numerous different assembly

676 versions. For example, we generated 15 different assemblies only for the parameter optimization of

28

677  the linked-read scaffolding (**Figure 1d**) and there are studies generating even 400 assemblies in

678  total (Montoliu-Nerin et al. 2019). In such a situation, it might seem difficult to decide how to

679  choose the "best" assembly among dozens. Here we present what we learned from the different

680  technologies and how they help in resolving the genomic regions that are most difficult to assemble.

681

682  We used two types of *de-novo* assemblies based on Illumina sequencing. The first, lycPyrIL is an

683  Illumina assembly made from multiple insert size libraries of paired end and mate pair reads (Prost

684  et al. 2019); the second on 10XGenomics linked reads (lycPyrSN1 and SN2). It is notable that

685  lycPyrIL is much more contiguous than lycPyrSN1 and SN2 (contig N50 of 620 kb vs 145-150 kb;

686  **Table 1**) and has much fewer gaps. Although lycPyrIL is a less fragmented assembly, lycPyrSN2

687  has a better resolution for repeats since 7 Mb more repeats are masked and a larger number of MHC

688  IIB exons are present (**Figure3g-h**) as well as G4 motifs (**Figure 3g**). Nonetheless, the contiguity

689  reached in lycPyrPB for the same sample at contig level (contig N50 of 6 Mb) is ten-fold higher

690  than in lycPyrIL and even outscores lycPyrIL scaffold N50 of 4 Mb. 10XGenomics linked reads

691  bring long-range information through the barcode system that is useful for local phasing, detection

692  of structural variations (Zheng et al. 2016; Marks et al. 2019), scaffolding (Yeo et al. 2017) and

693  construction of recombination maps (Dréau et al. 2019; Sun et al. 2019). We used the barcode

694  information to scaffold the PacBio assembly (lycPyr3, **Table 1**) without obtaining many new

695  scaffolds but this could be due to the already high contiguity of the input lycPyrPB assembly.

696  Finally, we note that the molecule input length for the 10XGenomics libraries have different effects

697  on the assembly and BUSCO scores. That is, lycPyrSN1 (24 kb mean molecule length library)

698  outscores lycPyrSN2 (26.1 kb mean molecule length library) in the number of complete BUSCO

699  genes (**Supplementary Table S2**). Even though 10XGenomics linked reads consist of short reads,

700  both lycPyrSN1 and lycPyrSN2 have more missing genes compared to lycPyrIL (**Supplementary**

701  **Table S2**).

702

703    Long reads together with proximity ligation maps are game changers in genomics. Their

704    combination yielded a very high-quality assembly for a non-model bird with suboptimal sample

705    quality (see mean molecule lengths for 10XGenomics assemblies above). The PacBio assembly is

706    by far the most contiguous and a suitable genomic backbone to obtain chromosome models

707    including the W chromosome and several microchromosomes. The main weakness linked to PacBio

708    is the introduction of sequencing errors (mostly short indels) that must be corrected with accurate

709    short reads. As mentioned before, the sequencing errors hinder the identification of gene models

710    (BUSCO) and protein prediction (Watson and Warr 2019). Moreover, the PacBio assembly is likely

711    not free of misassemblies (e.g., chimeric contigs). Thus a second type of independent data is

712    necessary to detect such errors; e.g., ~100 potential misassemblies were identified by the

713    CHiCAGO proximity map. The CHiCAGO map was very useful to correct the assembly and make

714    a first scaffolding, but neither alone nor with 10XGenomics scaffolding yielded a chromosome-

715    level assembly. The only type of data implemented here that allowed the generation of chromosome

716    models was the Hi-C map. The latter does not rely on extracted DNA quality or library insert size,

717    but instead on *in-situ* proximity within the nuclei of the fixed sample. As such, Hi-C data is an

718    effective replacement of linkage maps for scaffolding purposes (Dudchenko et al. 2017) and can be

719    used to manually curate assemblies.

720

721    A direct way to identify the limits of sequencing data is to investigate where assemblers fail to

722    resolve sequences, i.e. where contig fragmentation occurs. Therefore, we characterized what causes

723    contig fragmentation in each assembly by analysing sequences directly adjacent to gaps and

724    inferring the gap content of draft assemblies by aligning their flanks to the final multiplatform

725    version lycPyr6 (**Figure 4a**). In general, we found that long and/or homogeneous repeats such as

726    LTR retrotransposons, satellites, and simple repeats are the main fragmentation causes in every

727    assembly, though the specific repeat type changed with the technology. Short-read and linked-read

728    contigs mostly break at simple repeats. Even though the percentage of simple repeats assembled in

30

729    lycPyrIL, lycPyrSN1 and lycPyrSN2 ranges between 80-90% relative to lycPyr6 (**Figure 4b**),

730    simple repeats also caused most of the assembly gaps, indicating that insert size and linked read

731    methods are not sufficient to unambiguously solve those regions (**Figure 4c-d**). At the same time,

732    the gaps of these three assemblies, when compared to the final multiplatform assembly, mainly

733    contain LTR retrotransposons, simple repeats and complex repeats (defined as arrays of different

734    types of repeats; **Figure 4e-f**). LTR retrotransposons are the second most abundant retrotransposons

735    in the paradise crow assembly and several kilobases long. These features make LTR

736    retrotransposons the major cause of fragmentation in the PacBio assembly and the second in the

737    short-read ones. This partially unexpected trend is likely because LTR retrotransposons are

738    underrepresented in lycPyrIL, lycPyr SN1 and lycPyrSN2 (as indicated by their lack of part of the

739    recent LTR activity; **Figure 3e-f**). The same pattern can be observed for the multicopy rRNA genes:

740    the only assemblies showing gaps caused by rRNA genes are the PacBio-based and this is likely

741    because PacBio was the only technology able to (partially) solve those repeats (**Figure 4c-d**). It is

742    interesting that linked reads appear to better distinguish long repeats like LTR retrotransposons than

743    short-read libraries based on insert size (**Figure 4b**). The satellite portion of the genome was

744    significantly better assembled with PacBio long reads (~9 Mb), while neither multiple Illumina

745    libraries nor linked reads could assemble more than 1 Mb of satellites. This is probably due to the

746    highly homogeneous nature of long stretches of satellites that make satellite arrays collapse during

747    assembly (Hartley and O'Neill 2019). Similar to LTR retrotransposons and rRNA genes, satellites

748    are barely assembled in lycPyrIL, lycPyrSN1 and lycPyrSN2. Therefore satellites are not a major

749    cause of contig fragmentation in Illumina-based assemblies. LINEs are usually short

750    retrotransposons due to 5' truncation during integration (Levin and Moran 2011)and in the paradise

751    crow and other songbirds they seem to be mostly present in old copies (**Figure 3e**; Suh et al.

752    (2018); Weissensteiner et al. (2019)). Therefore they likely are less homogeneous elements, with

753    more diagnostic mutations and hence easier to assemble. In fact, both Illumina and 10XGenomics

754    assemblies have 96-98% of LINEs assembled and LINEs represent only the fourth causative factor

31

755  of fragmentation. Finally, we noticed a disproportion of DNA transposons annotated in the Illumina

756  assemblies (lycPyrIL and lycPyrILPB) compared to the other assemblies. This phenomenon might

757  be explained by annotation issues linked to the fragmentation of those regions or by the presence of

758  unsolved haplotypes. DNA transposons have been inactive in songbirds for even longer than LINEs

759  (Kapusta and Suh 2017) and should thus be rather straightforward to assemble.

760

## 761 Conclusions

762  Thanks to a manually curated multiplatform assembly and three *de-novo* draft assemblies for the

763  same sample, we were able to characterise and measure genome completeness across sequencing

764  technologies. As expected, long-read assemblies are more complete than short-read assemblies but

765  completeness has been usually measured with statistics that are optimized for short reads rather than

766  for long reads. Scaffold N50 and BUSCO values do not reflect the entire potential and strengths of

767  new sequence technologies, therefore we measured completeness focusing on the most difficult-to-

768  assemble genomic regions. By doing so, we traced the essential steps for generating a high-quality

769  assembly for a non-model organism while optimizing costs and efforts.

770

771  Based on our assembly comparisons, the essential elements to make a chromosome-level assembly

772  are a contiguous primary assembly based on long reads, an independent set of data for correcting

773  misassemblies (CHiCAGO map or linked reads) and polish sequencing errors (short or linked

774  reads), and a Hi-C map for chromosome-level scaffolding. PacBio needs error correction both at the

775  nucleotide level (base calling errors and short indels) and at the assembly level (e.g., chimeric

776  contigs). For both scopes it is possible to use Illumina data but a note of caution is due. First, when

777  polishing the assembly for base calling errors and short indels, short reads could over-homogenize

778  repetitive sequences and thus it would be advisable to correct only outside repeats. In addition,

779  10XGenomics linked reads can also be used to correct both sequencing errors and misassemblies

780  (e.g., Tigmint, Jackman et al. (2018)) and to scaffold the genome (ARCS, Yeo et al. (2017), ARKS,

32

781    Coombe et al. (2018), fragScaff, Adey et al. (2014)). In general, the spatial information brought by

782    linked reads seems to be very versatile (e.g., assembly correction, scaffolding, structural variation

783    inference, haplotype phasing) and able to better avoid over-collapsing of repetitive elements and

784    genes (**Figure 3** and **4**). Therefore, if budgets and sample material are limited, this technology may

785    be suitable to obtain a better genomic overview than short reads alone. Nevertheless, long reads

786    provide the most detailed look into difficult-to-assemble genomic regions. We summarized the

787    strengths and limitations of the implemented technologies in **Figure 5** that can be used as a guide

788    for choosing technologies and ranking assemblies.

789

790    We have shown that recent technological developments have led to enormous improvements in

791    assembly quality and completeness, paving the way to more complete comparative genomic analyses,

792    including regions that were previously inaccessible within genomic dark matter. At the same time,

793    awareness of technological strengths and weaknesses in resolving repeat-rich and GC-rich regions is

794    fundamental for choosing the most suitable technology when designing sequencing projects, and will

795    help in a dilemma many genome scientists face these days: choosing the best assembly among many.



796

797    **Figure 5**. Summary of the relative efficiency of the different technologies over
798    quality/completeness parameters. Green: most effective; red: least effective.

799

800

# Methods

**Samples**

We used pectoral muscle samples from three vouchered specimens of *Lycocorax pyrrhopterus* ssp. *obiensis* collected on Obi Island (Moluccas, Indonesia) in 2013, from the Museum Zoologicum Bogoriense (MZB) in Bogor, Indonesia, temporarily on loan at the Natural History Museum of Denmark. One female (voucher: MZB 34.073) sample preserved in DMSO was used for PacBio, Illumina and 10XGenomics sequencing and for the Dovetail CHiCAGO library, one female sample (voucher: MZB 34.070) preserved in RNAlater was used for the Hi-C library with Phase Genomics, and one male sample preserved in DMSO (voucher: MZB 34.075) was used for Illumina sequencing.

**Sequencing technologies and *de-novo* assemblies**

We sequenced the female sample MZB 34.073 using a) PacBio RSII C6-P4 (mean of 11 kb and N50 of 16 kb for read length) for a total coverage of 72X; b) 10XGenomics with a HiSeqX Illumina machine (24 kb mean molecule length, 280 bp library insert size, 150 bp read length, net coverage 39.7X); c) 10XGenomics with HiSeqX Illumina machine (26.1 kb mean molecule length, 280 bp library insert size, 150 bp read length, net coverage 37.9X). DNA was extracted using magnetic beads on a Kingfisher robot, except for library c) which was based on DNA extracted with agarose gel plugs as in (Weissensteiner et al. 2017). In addition to these libraries, we also used the Illumina libraries and assembly produced in (Prost et al. 2019): Illumina HiSeq 2500 TruSeq paired-end libraries (180 bp and 550 bp insert sizes) and Nextera mate pair libraries (5 kb and 8 kb insert sizes) for a total coverage of 90X. Furthermore, two paired-end libraries (125 bp read length) of chromatin-chromatin interactions from CHiCAGO and Hi-C techniques were produced using a HiSeq 2500 by Dovetail Genomics (Putnam et al. 2016) and Phase Genomics (more details below), respectively. Finally, we generated a paired-end library with insert size of 650 bp on an Illumina HiSeqX machine for the male sample.

34

826    For each library/technology (namely Illumina, 10XGenomics and PacBio) we made independent *de-*

827    *novo* assemblies. (Prost et al. 2019) used ALLPATHS-LG (Butler et al. 2008) for Illumina data while

828    we used Falcon (Chin et al. 2016) for PacBio data and Supernova2 (Weisenfeld et al. 2017) for

829    10XGenomics data (**Table 1**). All the basic genome statistics of the assemblies (**Supplementary**

830    **Table S1**)    were    calculated    using    the    Perl    script    assemblathon_stats.pl    from

831    https://github.com/KorfLab/Assemblathon/blob/master/assemblathon_stats.pl.

832

833    **Identification of sex-linked contigs and PAR**

834    Given the extreme conservation of the Z chromosomes of songbird (Xu et al. 2019), we used the Z-

835    chromosome sequence of great tit as a query to search for homologous Z-linked contigs in paradise

836    crow. The aligner nucmer was used to perform the one-to-one alignment of the great tit genome and

837    lycPyrPB. Those contigs with more than 60 percent sequence aligned to great tit Z chromosome were

838    identified as Z-linked. We further calculated the sequencing coverage using the female Illumina

839    paired-end libraries to confirm the half-coverage pattern of candidate Z-linked contigs relative to

840    autosomal contigs. We used BWA-MEM to map the reads and the samtools depth function to estimate

841    contig coverage. To identify candidate W-linked contigs, we calculated the re-sequencing coverage

842    of the male individual, because W-linked contigs are female-specific and are not expected to be

843    mapped by male reads while the coverage of female reads should be half of that of autosomes. We

844    used the known PAR sequences of collared flycatcher (Smeds et al. 2014) to identify the homologous

845    PAR contigs in paradise crow. As expected, the PAR contigs were found to show similar re-

846    sequencing coverage in both the male and the female as on the autosomes.

847

848    **Multiplatform approach**

849    We created three types of multiplatform assemblies, one that combines only Illumina and PacBio data

850    (lycPyrILPB, see **Table 1**), a second one combining PacBio and Hi-C data, and a third more

35

851    comprehensive one that combines three types of sequencing data and two types of proximity ligation

852    data (lycPyr6).

853

854    For the first type of assembly (lycPyrILPB), we used the Illumina assembly lycPyrIL (Prost et al.

855    2019) as genomic backbone and gap-filled it with PacBio long reads using the software PBJelly

856    (PBSuite v. 15.8.24) maintaining the all the default options but -min 10 to consider only gaps of at

857    least 10 base pairs length. The second multiplatform assembly lycPyrHiC was built by scaffolding

858    the PacBio primary assembly (lycPyrPB) with Hi-C data.

859

860    For the most comprehensive assembly (lycPyr6), we combined PacBio, Illumina, 10XGenomics,

861    CHiCAGO and Hi-C data (**Figure 1**). The first step was to assemble the PacBio reads into a primary

862    assembly with the Falcon software (Chin et al. (2016); **Figure 1a**). The primary contigs were

863    corrected and scaffolded with the Dovetail CHiCAGO map generating lycPyr2 (**Figure 1b**) using the

864    software HiRise (Putnam et al. 2016). lycPyr2 then was polished with long reads (two runs of Arrow;

865    Chin et al. (2016)) and short reads (three runs of Pilon 1.22; Walker et al. (2014); **Figure 1c**). Since

866    PacBio sequencing is prone to introduce short indels in the reads (Eid et al. 2009), we addressed

867    specifically these sequencing problems with Pilon while we did not correct single nucleotide variants.

868    Furthermore, in order to not over-polish repetitive regions (i.e., homogenising them with short reads),

869    we excluded Pilon corrections falling within repeats identified by RepeatMasker 4.0.7 using our

870    custom repeat library.

871

872    We then scaffolded lycPyr2 using the long-range information given by 10XGenomics linked reads

873    with the software ARCS 1.0.1 (Yeo et al. (2017); parameters -s 95 -e 1000 -m 20-100000) and LINKS

874    1.8.5 (Warren et al. (2015); parameters -a 0.2) generating lycPyr3 (**Figure 1d**). The parameters for

875    ARCS and LINKS have been chosen after generating 15 assemblies with different values for -m -e -

876    a (**Supplementary Table S11**). The optimal parameter combination was established by minimising

877    a) the number of "private" scaffolds belonging only to one combination of parameters, b) the number

878    of scaffolds containing putative in-silico chromosomal translocations.

879

880    lycPyr3 was scaffolded into chromosome models (clusters of contigs and scaffolds) with the Phase

881    Genomics Hi-C data and the Proximo Hi-C scaffolding pipeline (lycPyr4; **Figure 1e**). Hi-C data were

882    generated using a Phase Genomics (Seattle, WA) Proximo Hi-C Animal Kit. Following the

883    manufacturer's instructions for the kit, intact cells from two samples were crosslinked using a

884    formaldehyde solution, digested using the Sau3AI restriction enzyme, and proximity ligated with

885    biotinylated nucleotides to create chimeric molecules composed of fragments from different regions

886    of the genome that were physically proximal *in-vivo*, but not necessarily linearly proximal.

887    Continuing with the manufacturer's protocol, molecules were pulled down with streptavidin beads

888    and processed into an Illumina-compatible sequencing library.

889

890    Reads were aligned to the draft assembly lycPyr3 following the manufacturer's recommendations.

891    Briefly, reads were aligned using BWA-MEM (Li and Durbin (2010); v. 0.7.15-r1144-dirty) with the

892    -5SP and -t 8 options specified, while keeping the other parameters as default. SAMBLASTER (Faust

893    and Hall 2014) was used to flag PCR duplicates, which were later excluded from analysis. Alignments

894    were then filtered with samtools (Li et al. (2009); v1.5, with htslib 1.5) using the -F 2304 filtering

895    flag to remove non-primary and secondary alignments, as well as read pairs in which one or more

896    mates were unmapped. Phase Genomics' Proximo Hi-C genome scaffolding platform was used to

897    create chromosome-scale scaffolds from the draft assembly as described in (Bickhart et al. 2017). As

898    in the LACHESIS method (Burton et al. 2013), this process computes a contact frequency matrix

899    from the aligned Hi-C read pairs, normalised by the number of Sau3AI restriction sites (GATC) on

900    each contig, and constructs scaffolds in such a way as to optimise expected contact frequency and

901    other statistical patterns in Hi-C data. Approximately 286,000 separate Proximo runs were performed

37

902 to optimise the number of scaffolds and scaffold construction in order to make the scaffolds as

903 concordant with the observed Hi-C data as possible.

904

905 Two chromosomes (chr1 and chr2) appeared to be split into two different super-scaffolds (or clusters)

906 respectively, thus they were manually put together following the orientation suggested by the Hi-C

907 interaction heatmap (**Supplementary Figure S3**). We then manually inspected the assembly lycPyr4

908 for misassemblies (**Figure 1f** and **Figure 2**) by aligning the four *de-novo* assemblies (lycPyrIL,

909 lycPyrPB, lycPyrSN1 and lycPyrSN2) to it using Satsuma2 (Grabherr et al. 2010) and chromosome

910 models from three songbird outgroups (*Ficedula albicollis*, *Taeniopygia guttata* and *Parus major*)

911 with LASTZ 1.04.00 (Harris 2007). We identified misassemblies by looking for regions in which the

912 different *de-novo* assemblies were in conflict with the final assembly (schematically showed in

913 **Figure 2**). We applied the majority rule for each scaffolding or orientation conflict found between

914 lycPyr4 and the four draft assemblies. To make any decisions against the scaffold configuration in

915 lycPyr4, three of the four *de-novo* assemblies needed to be in discordance with lycPyr4 and show the

916 same pattern of discordance. In cases where only two *de-novo* assemblies showed the same pattern

917 of discordance and the other were not informative, we used the information provided by the outgroups

918 to decide whether to keep the lycPyr4 scaffold configuration or correct it. With this approach we were

919 able to identify 45 intra-scaffold misassemblies at a fine scale, all of them being orientation issues of

920 PacBio contigs within scaffolds.

921

922 Then, we gap-filled the assembly using PBJelly (PBSuite 15.8.24; English et al. (2012)) with the

923 default options except for the parameter -min 10 in order to consider the gaps longer than 10 bps

924 (**Figure 1g**). After the gap-filling step that used the PacBio reads, we ultimately polished the genome

925 with long reads using Arrow (one run; PacBio library) and with short reads using Pilon (two runs;

926 Illumina library; **Figure 1h**).

38

927   The last step of assembly curation involved the generation of Hi-C heatmaps on lycPyr5 by mapping

928   the Hi-C library to the assembly using Juicer 1.5 (Durand et al. (2016); **Figure 1i**). We manually

929   inspected the Hi-C maps for misassemblies using Juicebox 1.9.8

930   (https://github.com/aidenlab/Juicebox) and corrected lycPyr5 accordingly (**Supplementary Figure**

931   **S3**). This way, we manually solved remaining assembly issues regarding the orientation and order of

932   some contigs or scaffolds within the chromosome models, as well as corrected *in-silico* chromosomal

933   translocations.

934

935   The completeness of the assemblies was assessed with gVolante (Nishimura et al. 2017) using

936   BUSCO v3 for avian genomes (**Supplementary Table S2**) and with LTR Assembly Index (Ou et al.

937   (2018); **Supplementary Table S3**).

938

939   The mitochondrial genome was identified as one PacBio contig by aligning the mtDNA of *Corvus*

940   *corax* (GenBank accession number KX245138.1) to lycPyrPB. It was annotated using DOGMA

941   (Wyman et al. 2004) and tRNAscan-SE 1.3.1 (Lowe and Eddy (1997); **Supplementary Table S12**).

942

943   **Chromosome nomenclature**

944   Since the chicken genome is the best avian genome assembled so far with reliable chromosome

945   information (Warren et al. 2017), we named and oriented our chromosome models according to

946   homology with galGal5 (RefSeq accession number GCF_000002315.6). In the case that our

947   chromosome models were not completely collinear with chicken, we oriented them following the

948   orientation of the majority of the model respect to chicken. Finally, if the chromosome models did

949   not share any homology with chicken, their orientation was not changed.

950

951   **Repeat library**

39

952     We produced a *de-novo* repeat library for paradise crow by running the RepeatMasker 4.0.7 and

953     RepeatModeler 1.0.8 software on the PacBio *de-novo* assembly. We hard-masked lycPyrPB with the

954     Aves repeat library from Repbase (version 20170127; https://www.girinst.org/about/repbase.html)

955     together with the consensus sequences from (Prost et al. 2019), then ran RepeatModeler. The new

956     consensus sequences generated by RepeatModeler were aligned back to the reference genome; the 20

957     best BLASTN 2.7.1+ results were collected, extended by 2 kb on both sides and aligned to one

958     another with MAFFT 7.4.07. The alignments were manually curated applying the majority rule and

959     the superfamily of repeat assessed following the (Wicker et al. 2007) classification.

960

961     All    the    new    consensus    sequences    were    masked    in    CENSOR

962     (http://www.girinst.org/censor/index.php) and named according to homology to known repeats in the

963     Repbase database. Sequences with high similarity to known repeats for their entire lengths were given

964     the name of the known repeat + suffix "_lycPyr"; repeats with partial homology have been named

965     with the suffix "-L_lycPyr" where "L" stands for "like" (Suh et al. 2018). Repeats with no homology

966     with known ones have been considered as new families and named with the prefix "lycPyr" followed

967     by the name of their superfamilies.

968

969     The final repeat library also contains the manually curated version of the consensus sequences

970     previously generated on other two birds-of-paradise *Astrapia rothschildi* "astRot", *Ptiloris*

971     *paradiseus* "ptiPar" (Prost et al. 2019), the ones from *Corvus cornix* (Weissensteiner et al. 2019),

972     *Uraeginthus cyanocephalus* (Boman et al. 2019), *Ficedula albicollis* and all the avian repeats

973     available on Repbase (mostly from chicken and zebra finch).

974

975     **G4 motif identification**

976     The *de-novo* assemblies and the final version have been scanned for G-quadruplex (G4) motifs with

977     the software Quadron (Sahakyan et al. (2017); https://github.com/aleksahak/Quadron). Only non-

978    overlapping hits with a score greater than 19 were used for subsequent analysis as suggested in

979    (Sahakyan et al. 2017). The density of such motifs per chromosome model was calculated using

980    bedtools coverage (BEDTools 2.27.1; Quinlan (2014)).

981

982    **MHC class IIB analysis**

983    To infer how highly duplicated genes are assembled with different input data and assembly

984    strategies, we investigated the distribution of major histocompatibility class IIB (MHCIIB)

985    sequence hits in seven assemblies: lycPyrIL, lycPyrPB, lycPyrSN1, lycPyrSN1, lycPyrILPB,

986    lycPyr2 and lycPyr6 (the intermediate assemblies like lycPyr3 are not shown here because the MHC

987    content did not change from lycPyr2 to lycPyr5). We performed BLAST (Altschul et al. 1990)

988    searches both with sequences of the highly variable exon 2 that encodes the peptide binding region,

989    and with the much more conserved exon 3 (Hughes and Yeager 1998), as the disparate levels of

990    polymorphism within these regions may provide insights into different aspects of challenges with

991    genome assembly. We conducted tBLASTn (BLAST 2.7.1+) searches using alignments available

992    from Goebel et al. (2017) that include sequences from across the entire avian phylogeny. We chose

993    this strategy to ensure the identification MHCIIB sequences, as with sequences of only a single-

994    species BLAST search might miss highly divergent sequences as they are often present in the MHC,

995    where within-species diversity of MHC genes often equals between-species divergence. From the

996    available alignments, we exclusively retained sequences spanning the entire 270 bp of exon 2 and

997    sequences covering 220 bp of exon 3. This left query alignments including 233 sequences from 22

998    bird orders/families for exon 2, and 314 sequences from 26 bird orders/families for exon 3.

999    Overlapping blast hit intervals were merged. To ensure that these intervals contained sequences

1000    corresponding to MHCIIB, we first BLAST searched them back against the GenBank database

1001    using BLASTn queries, and retained only intervals producing hits with MHCIIB. We then aligned

1002    the remaining sequences using the MAFFT alignment server with the --add option and default

1003    settings, and manually screened the alignments to identify non-MHCIIB sequences. Finally, we

41

1004 determined the alignment lengths of BLAST hit intervals after removing insertions relative to the

1005 query alignment. We report only hits longer than 240 bp for exon 2 and longer than 195 bp for exon

1006 3, corresponding to approximately 90% of the respective query alignment lengths.

1007

1008 **Gap analysis**

1009 For each assembly produced, we estimated the number of gaps caused by repeats by intersecting the

1010 gap and repeat coordinates using bedtools window (Quinlan 2014) with a window size of 100 bp

1011 (**Figure 4a**). Only gaps longer than 10 bp were taken into consideration. This filter is particularly

1012 important for lycPyrIL since there are many small gaps of 1-5 Ns that are probably caused by

1013 sequencing or base-calling errors.

1014

1015 We estimated what is missing in the draft assemblies with respect to the final multiplatform assembly

1016 lycPyr6 by aligning the flanking regions to the gaps onto the final version. We then assessed the

1017 presence of annotated repeats on lycPyr6 between the aligned flanking regions to the draft assembly

1018 gaps. To do these pairwise alignments, we extracted 500 bp of flanking regions from the intra-scaffold

1019 gaps of lycPyrIL, lycPyrSN1, lycPyrSN2, lycPyrPB and lycPyrILPB and BLASTn searched the

1020 sequences to lycPyr6 with BLAST 2.7.1+. The alignments were filtered to retain only unambiguously

1021 orthologous positions on lycPyr6, namely there was only one alignment (98% identity, 90% coverage)

1022 of both flanks on the same lycPyr6 scaffold. The coordinates of the draft genome gaps projected onto

1023 lycPyr6 were then intersected with the RepeatMasker annotation using bedtools intersect. Draft

1024 genome gaps containing only one type of repeat on lycPyr6 were classified according to the type of

1025 repeat. In case the draft genome gaps corresponded to a region containing more than one type of

1026 repeat, the gaps were classified as 'complex'. Finally, in case that the draft genome gaps could not be

1027 mapped unambiguously (e.g., no homology, only one flank aligned or the two flanking regions

1028 mapped to different scaffolds) or mapped to gaps on lycPyr6, they were classified as 'not scorable

1029 gaps' (**Figure 4a**)

42

1030

1031    We also compared how many repeats were assembled in the draft assemblies compared to lycPyr6

1032    (**Figure 4b**) by calculating the proportion of repeat base pairs present in the draft assemblies relative

1033    to the total bp in lycPyr6. This was done for each major repeat group using the RepeatMasker table

1034    (.tbl) files; more details in **Supplementary Table S10**.

1035

## Data Access

1036

1037

## Acknowledgements

1069
1070

## Disclosure declaration

1072   Shawn Sullivan and Ivan Liachko are employed at Phase Genomics.

1073

## References

1075   Adey A, Kitzman JO, Burton JN, Daza R, Kumar A, Christiansen L, Ronaghi M, Amini S, L.
1076        Gunderson K, Steemers FJ et al. 2014. In vitro, long-range sequence information for de novo
1077        genome assembly via transposase contiguity. *Genome Research* **24**: 2041-2049.

1078   Aird D, Ross MG, Chen WS, Danielsson M, Fennell T, Russ C, Jaffe DB, Nusbaum C, Gnirke A.
1079        2011. Analyzing and minimizing PCR amplification bias in Illumina sequencing libraries.
1080        *Genome Biol* **12**.

1081   Alkan C, Sajjadian S, Eichler EE. 2010. Limitations of next-generation genome sequence assembly.
1082        *Nature Methods* **8**: 61.

1083   Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. 1990. Basic local alignment search tool.
1084        *Journal of Molecular Biology* **215**: 403-410.

1085   Bellott DW, Skaletsky H, Cho T-J, Brown L, Locke D, Chen N, Galkina S, Pyntikova T, Koutseva
1086        N, Graves T et al. 2017. Avian W and mammalian Y chromosomes convergently retained
1087        dosage-sensitive regulators. *Nature Genetics* **49**: 387.

1088   Belser C, Istace B, Denis E, Dubarry M, Baurens F-C, Falentin C, Genete M, Berrabah W, Chèvre
1089        A-M, Delourme R et al. 2018. Chromosome-scale assemblies of plant genomes using
1090        nanopore long reads and optical maps. *Nature Plants* **4**: 879-887.

1091   Bickhart DM, Rosen BD, Koren S, Sayre BL, Hastie AR, Chan S, Lee J, Lam ET, Liachko I, Sullivan
1092        ST et al. 2017. Single-molecule sequencing and chromatin conformation capture enable de
1093        novo reference assembly of the domestic goat genome. *Nature Genetics* **49**: 643.

1094   Biffi G, Tannahill D, Balasubramanian S. 2012. An Intramolecular G-Quadruplex Structure Is
1095        Required for Binding of Telomeric Repeat-Containing RNA to the Telomeric Protein TRF2.
1096        *Journal of the American Chemical Society* **134**: 11974-11976.

Boman J, Frankl-Vilches C, da Silva dos Santos M, de Oliveira EHC, Gahr M, Suh A. 2019. The Genome of Blue-Capped Cordon-Bleu Uncovers Hidden Diversity of LTR Retrotransposons in Zebra Finch. *Genes* **10**: 301.

Botero-Castro F, Figuet E, Tilak M-K, Nabholz B, Galtier N. 2017. Avian Genomes Revisited: Hidden Genes Uncovered and the Rates versus Traits Paradox in Birds. *Molecular Biology and Evolution* **34**: 3123-3131.

Bourque G, Burns KH, Gehring M, Gorbunova V, Seluanov A, Hammell M, Imbeault M, Izsvák Z, Levin HL, Macfarlan TS et al. 2018. Ten things you should know about transposable elements. *Genome Biology* **19**: 199.

Burt DW. 2002. Origin and evolution of avian microchromosomes. *Cytogenetic and Genome Research* **96**: 97-112.

Burton JN, Adey A, Patwardhan RP, Qiu R, Kitzman JO, Shendure J. 2013. Chromosome-scale scaffolding of de novo genome assemblies based on chromatin interactions. *Nature Biotechnology* **31**: 1119.

Butler J, MacCallum I, Kleber M, Shlyakhter IA, Belmonte MK, Lander ES, Nusbaum C, Jaffe DB. 2008. ALLPATHS: De novo assembly of whole-genome shotgun microreads. *Genome Research* **18**: 810-820.

Chaisson MJP, Huddleston J, Dennis MY, Sudmant PH, Malig M, Hormozdiari F, Antonacci F, Surti U, Sandstrom R, Boitano M et al. 2014. Resolving the complexity of the human genome using single-molecule sequencing. *Nature* **517**: 608.

Chaisson MJP, Wilson RK, Eichler EE. 2015. Genetic variation and the de novo assembly of human genomes. *Nature Reviews Genetics* **16**: 627.

Chalopin D, Volff J-N, Galiana D, Anderson JL, Schartl M. 2015. Transposable elements and early evolution of sex chromosomes in fish. *Chromosome Research* **23**: 545-560.

Chang C-H, Chavan A, Palladino J, Wei X, Martins NMC, Santinello B, Chen C-C, Erceg J, Beliveau BJ, Wu C-T et al. 2019. Islands of retroelements are major components of Drosophila centromeres. *PLOS Biology* **17**: e3000241.

Charlesworth B, Harvey PH, Charlesworth B, Charlesworth D. 2000. The degeneration of Y chromosomes. *Philosophical Transactions of the Royal Society of London Series B: Biological Sciences* **355**: 1563-1572.

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. 2016. Phased diploid genome assembly with single-molecule real-time sequencing. *Nature Methods* **13**: 1050.

Chuong EB, Elde NC, Feschotte C. 2016. Regulatory activities of transposable elements: from conflicts to benefits. *Nature Reviews Genetics* **18**: 71.

Coombe L, Zhang J, Vandervalk BP, Chu J, Jackman SD, Birol I, Warren RL. 2018. ARKS: chromosome-scale scaffolding of human genome drafts with linked read kmers. *BMC Bioinformatics* **19**: 234.

Cowley M, Oakey RJ. 2013. Transposable Elements Re-Wire and Fine-Tune the Transcriptome. *PLOS Genetics* **9**: e1003234.

Deschamps S, Zhang Y, Llaca V, Ye L, Sanyal A, King M, May G, Lin H. 2018. A chromosome-scale assembly of the sorghum genome using nanopore sequencing and optical mapping. *Nature Communications* **9**: 4844.

Dohm JC, Lottaz C, Borodina T, Himmelbauer H. 2008. Substantial biases in ultra-short read data sets from high-throughput DNA sequencing. *Nucleic Acids Res* **36**.

Domanska D, Kanduri C, Simovski B, Sandve GK. 2018. Mind the gaps: overlooking inaccessible regions confounds statistical testing in genome analysis. *BMC Bioinformatics* **19**: 481.

Dréau A, Venu V, Avdievich E, Gaspar L, Jones FC. 2019. Genome-wide recombination map construction from single individuals using linked-read sequencing. *Nature Communications* **10**: 4309.

1147  Du Z, Zhao Y, Li N. 2008. Genome-wide analysis reveals regulatory role of G4 DNA in gene
1148         transcription. *Genome Research* **18**: 233-241.
1149  Du Z, Zhao Y, Li N. 2009. Genome-wide colonization of gene regulatory elements by G4 DNA
1150         motifs. *Nucleic Acids Research* **37**: 6784-6798.
1151  Dudchenko O, Batra SS, Omer AD, Nyquist SK, Hoeger M, Durand NC, Shamim MS, Machol I,
1152         Lander ES, Aiden AP et al. 2017. De novo assembly of the *Aedes aegypti* genome using Hi-
1153         C yields chromosome-length scaffolds. *Science* **356**: 92-95.
1154  Dudchenko O, Shamim MS, Batra SS, Durand NC, Musial NT, Mostofa R, Pham M, Glenn St Hilaire
1155         B, Yao W, Stamenova E et al. 2018. The Juicebox Assembly Tools module facilitates *de novo*
1156         assembly of mammalian genomes with chromosome-length scaffolds for under $1000.
1157         *bioRxiv* doi:10.1101/254797: 254797.
1158  Durand NC, Shamim MS, Machol I, Rao SSP, Huntley MH, Lander ES, Aiden EL. 2016. Juicer
1159         Provides a One-Click System for Analyzing Loop-Resolution Hi-C Experiments. *Cell
1160         Systems* **3**: 95-98.
1161  Eid J, Fehr A, Gray J, Luong K, Lyle J, Otto G, Peluso P, Rank D, Baybayan P, Bettman B et al.
1162         2009. Real-Time DNA Sequencing from Single Polymerase Molecules. *Science* **323**: 133-
1163         138.
1164  Emera D, Wagner GP. 2012. Transposable element recruitments in the mammalian placenta: impacts
1165         and mechanisms. *Briefings in Functional Genomics* **11**: 267-276.
1166  English AC, Richards S, Han Y, Wang M, Vee V, Qu J, Qin X, Muzny DM, Reid JG, Worley KC et
1167         al. 2012. Mind the Gap: Upgrading Genomes with Pacific Biosciences RS Long-Read
1168         Sequencing Technology. *PLOS ONE* **7**: e47768.
1169  Etherington GJ, Heavens D, Baker D, Lister A, McNelly R, Garcia G, Clavijo B, Macaulay I, Haerty
1170         W, Di Palma F. 2019. Sequencing smart: *De novo* sequencing and assembly approaches for
1171         non-model mammals. *bioRxiv* doi:10.1101/723890: 723890.
1172  Faino L, Seidl MF, Datema E, van den Berg GCM, Janssen A, Wittenberg AHJ, Thomma BPHJ.
1173         2015. Single-Molecule Real-Time Sequencing Combined with Optical Mapping Yields
1174         Completely Finished Fungal Genome. *mBio* **6**: e00936-00915.
1175  Goebel J, Promerová M, Bonadonna F, McCoy KD, Serbielle C, Strandh M, Yannic G, Burri R,
1176         Fumagalli L. 2017. 100 million years of multigene family evolution: origin and evolution of
1177         the avian MHC class IIB. *BMC Genomics* **18**: 460.
1178  Goodwin S, McPherson JD, McCombie WR. 2016. Coming of age: ten years of next-generation
1179         sequencing technologies. *Nature Reviews Genetics* **17**: 333.
1180  Gordon D, Huddleston J, Chaisson MJP, Hill CM, Kronenberg ZN, Munson KM, Malig M, Raja A,
1181         Fiddes I, Hillier LW et al. 2016. Long-read sequence assembly of the gorilla genome. *Science*
1182         **352**: aae0344.
1183  Grabherr MG, Russell P, Meyer M, Mauceli E, Alföldi J, Di Palma F, Lindblad-Toh K. 2010.
1184         Genome-wide synteny through highly sensitive sequence alignment: Satsuma. *Bioinformatics
1185         (Oxford, England)* **26**: 1145-1151.
1186  Gregory TR. 2019. Animal Genome Size Database, http://www.genomesize.com.
1187  Griffin D, Burt DW. 2014. All chromosomes great and small: 10 years on. *Chromosome Research*
1188         **22**: 1-6.
1189  Guiblet WM, Cremona MA, Cechova M, Harris RS, Kejnovská I, Kejnovsky E, Eckert K,
1190         Chiaromonte F, Makova KD. 2018. Long-read sequencing technology indicates genome-wide
1191         effects of non-B DNA on polymerization speed and error rate. *Genome Research* **28**: 1767-
1192         1778.
1193  Harris RS. 2007. Improved pairwise alignment of genomic DNA. *PhD Thesis, The Pennsylvania
1194         State University*.
1195  Hartley G, O'Neill RJ. 2019. Centromere Repeats: Hidden Gems of the Genome. *Genes* **10**: 223.
1196  Hobza R, Cegan R, Jesionek W, Kejnovsky E, Vyskot B, Kubat Z. 2017. Impact of Repetitive
1197         Elements on the Y Chromosome Formation in Plants. *Genes (Basel)* **8**.

46

Hron T, Pajer P, Pačes J, Bartůněk P, Elleder D. 2015. Hidden genes in birds. *Genome Biology* **16**: 164.

Hughes AL, Yeager M. 1998. NATURAL SELECTION AT MAJOR HISTOCOMPATIBILITY COMPLEX LOCI OF VERTEBRATES. *Annual Review of Genetics* **32**: 415-435.

Irestedt M, Jønsson KA, Fjeldså J, Christidis L, Ericson PGP. 2009. An unexpectedly long history of sexual selection in birds-of-paradise. *BMC Evolutionary Biology* **9**: 235.

Jackman SD, Coombe L, Chu J, Warren RL, Vandervalk BP, Yeo S, Xue Z, Mohamadi H, Bohlmann J, Jones SJM et al. 2018. Tigmint: correcting assembly errors using linked reads from large molecules. *BMC Bioinformatics* **19**: 393.

Jain M, Olsen HE, Turner DJ, Stoddart D, Bulazel KV, Paten B, Haussler D, Willard HF, Akeson M, Miga KH. 2018. Linear assembly of a human centromere on the Y chromosome. *Nature Biotechnology* **36**: 321.

Johnson JM, Edwards S, Shoemaker D, Schadt EE. 2005. Dark matter in the genome: evidence of widespread transcription detected by microarray tiling experiments. *Trends in Genetics* **21**: 93-102.

Kapitonov VV, Koonin EV. 2015. Evolution of the RAG1-RAG2 locus: both proteins came from the same transposon. *Biology Direct* **10**: 20.

Kapusta A, Suh A. 2017. Evolution of bird genomes-a transposon's-eye view. *Ann N Y Acad Sci* **1389**: 164-185.

Kozarewa I, Ning Z, Quail MA, Sanders MJ, Berriman M, Turner DJ. 2009. Amplification-free Illumina sequencing-library preparation facilitates improved mapping and assembly of (G+C)-biased genomes. *Nat Methods* **6**.

Lajoie BR, Dekker J, Kaplan N. 2015. The Hitchhiker's guide to Hi-C analysis: practical guidelines. *Methods* **72**: 65-75.

Law JA, Jacobsen SE. 2010. Establishing, maintaining and modifying DNA methylation patterns in plants and animals. *Nature Reviews Genetics* **11**: 204-220.

Lerat E, Casacuberta J, Chaparro C, Vieira C. 2019. On the Importance to Acknowledge Transposable Elements in Epigenomic Analyses. *Genes* **10**: 258.

Levin HL, Moran JV. 2011. Dynamic interactions between transposable elements and their hosts. *Nature Reviews Genetics* **12**: 615-627.

Levis RW, Ganesan R, Houtchens K, Tolar LA, Sheen F-m. 1993. Transposons in place of telomeric repeats at a Drosophila telomere. *Cell* **75**: 1083-1093.

Li H, Durbin R. 2010. Fast and accurate long-read alignment with Burrows–Wheeler transform. *Bioinformatics* **26**: 589-595.

Li H, Genome Project Data Processing S, Wysoker A, Handsaker B, Marth G, Abecasis G, Ruan J, Homer N, Durbin R, Fennell T. 2009. The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**: 2078-2079.

Li Q, Li HB, Huang W, Xu YC, Zhou Q, Wang SH, Ruan J, Huang SW, Zhang Z. 2019. A chromosome-scale genome assembly of cucumber (Cucumis sativus L.). *Gigascience* **8**: 10.

Lieberman-Aiden E, van Berkum NL, Williams L, Imakaev M, Ragoczy T, Telling A, Amit I, Lajoie BR, Sabo PJ, Dorschner MO et al. 2009. Comprehensive Mapping of Long-Range Interactions Reveals Folding Principles of the Human Genome. *Science* **326**: 289-293.

Ligon RA, Diaz CD, Morano JL, Troscianko J, Stevens M, Moskeland A, Laman TG, Scholes E, III. 2018. Evolution of correlated complexity in the radically different courtship signals of birds-of-paradise. *PLOS Biology* **16**: e2006962.

Loomis EW, Eid JS, Peluso P, Yin J, Hickey L, Rank D, McCalmon S, Hagerman RJ, Tassone F, Hagerman PJ. 2013. Sequencing the unsequenceable: Expanded CGG-repeat alleles of the fragile X gene. *Genome Research* **23**: 121-128.

Lowe TM, Eddy SR. 1997. tRNAscan-SE: A Program for Improved Detection of Transfer RNA Genes in Genomic Sequence. *Nucleic Acids Research* **25**: 955-964.

1248  Marks P, Garcia S, Barrio AM, Belhocine K, Bernate J, Bharadwaj R, Bjornson K, Catalanotti C,
1249      Delaney J, Fehr A et al. 2019. Resolving the full spectrum of human genome variation using
1250      Linked-Reads. *Genome Research* **29**: 635-645.
1251  McGurk MP, Dion-Côté A-M, Barbash DA. 2019. Rapid evolution at the telomere: transposable
1252      element dynamics at an intrinsically unstable locus. *bioRxiv* doi:10.1101/782904: 782904.
1253  Michael TP, Jupe F, Bemm F, Motley ST, Sandoval JP, Lanz C, Loudet O, Weigel D, Ecker JR. 2018.
1254      High contiguity Arabidopsis thaliana genome assembly with a single nanopore flow cell.
1255      *Nature Communications* **9**: 541.
1256  Miller MM, Taylor RL, Jr. 2016. Brief review of the chicken Major Histocompatibility Complex: the
1257      genes, their distribution on chromosome 16, and their contributions to disease resistance.
1258      *Poultry Science* **95**: 375-392.
1259  Montoliu-Nerin M, Sánchez-García M, Bergin C, Grabherr M, Ellis B, Kutschera VE, Kierczak M,
1260      Johannesson H, Rosling A. 2019. From single nuclei to whole genome assemblies. *bioRxiv*
1261      doi:10.1101/625814: 625814.
1262  Nishimura O, Hara Y, Kuraku S. 2017. gVolante for standardizing completeness assessment of
1263      genome and transcriptome assemblies. *Bioinformatics* **33**: 3635-3637.
1264  O'Leary NA, Wright MW, Brister JR, Ciufo S, Haddad D, McVeigh R, Rajput B, Robbertse B, Smith-
1265      White B, Ako-Adjei D et al. 2015. Reference sequence (RefSeq) database at NCBI: current
1266      status, taxonomic expansion, and functional annotation. *Nucleic Acids Research* **44**: D733-
1267      D745.
1268  O'Connor EA, Westerdahl H, Burri R, Edwards SV. 2019. Avian MHC Evolution in the Era of
1269      Genomics: Phase 1.0. *Cells* **8**: 1152.
1270  Ou S, Chen J, Jiang N. 2018. Assessing genome assembly quality using the LTR Assembly Index
1271      (LAI). *Nucleic Acids Research* **46**: e126-e126.
1272  Oyola SO, Otto TD, Gu Y, Maslen G, Manske M, Campino S, Turner DJ, MacInnis B, Kwiatkowski
1273      DP, Swerdlow HP et al. 2012. Optimizing illumina next-generation sequencing library
1274      preparation for extremely at-biased genomes. *BMC Genomics* **13**: 1.
1275  Paajanen P, Kettleborough G, López-Girona E, Giolai M, Heavens D, Baker D, Lister A, Cugliandolo
1276      F, Wilde G, Hein I et al. 2019. A critical comparison of technologies for a plant genome
1277      sequencing project. *GigaScience* **8**.
1278  Peñalba JV, Deng Y, Fang Q, Joseph L, Moritz C, Cockburn A. 2019. Genome of an iconic Australian
1279      bird: High-quality assembly and linkage map of the superb fairy-wren (*Malurus cyaneus*).
1280      *Molecular Ecology Resources* doi:10.1111/1755-0998.13124.
1281  Peona V, Weissensteiner MH, Suh A. 2018. How complete are "complete" genome assemblies?-An
1282      avian perspective. *Mol Ecol Resour* doi:10.1111/1755-0998.12933.
1283  Pettersson ME, Rochus CM, Han F, Chen J, Hill J, Wallerman O, Fan G, Hong X, Xu Q, Zhang H et
1284      al. 2019. A chromosome-level assembly of the Atlantic herring genome—detection of a
1285      supergene and other signals of selection. *Genome Research* doi:10.1101/gr.253435.119.
1286  Platt RN, II, Blanco-Berdugo L, Ray DA. 2016. Accurate Transposable Element Annotation Is Vital
1287      When Analyzing New Genome Assemblies. *Genome Biology and Evolution* **8**: 403-410.
1288  Prost S, Armstrong EE, Nylander J, Thomas GWC, Suh A, Petersen B, Dalen L, Benz BW, Blom
1289      MPK, Palkopoulou E et al. 2019. Comparative analyses identify genomic features potentially
1290      involved in the evolution of birds-of-paradise. *GigaScience* **8**: 1-12.
1291  Putnam NH, O'Connell BL, Stites JC, Rice BJ, Blanchette M, Calef R, Troll CJ, Fields A, Hartley
1292      PD, Sugnet CW et al. 2016. Chromosome-scale shotgun assembly using an in vitro method
1293      for long-range linkage. *Genome Research* **26**: 342-350.
1294  Quinlan AR. 2014. BEDTools: The Swiss-Army Tool for Genome Feature Analysis. *Current
1295      Protocols in Bioinformatics* **47**: 11.12.11-11.12.34.
1296  Raiber E-A, Kranaster R, Lam E, Nikan M, Balasubramanian S. 2011. A non-canonical DNA
1297      structure is a binding motif for the transcription factor SP1 in vitro. *Nucleic Acids Research*
1298      **40**: 1499-1508.

1299 Rhoads A, Au KF. 2015. PacBio Sequencing and Its Applications. *Genomics, Proteomics &*
1300          *Bioinformatics* **13**: 278-289.
1301 Ribeiro FJ, Przybylski D, Yin S, Sharpe T, Gnerre S, Abouelleil A, Berlin AM, Montmayeur A, Shea
1302          TP, Walker BJ et al. 2012. Finished bacterial genomes from shotgun sequence data. *Genome*
1303          *Research* **22**: 2270-2277.
1304 Sahakyan AB, Chambers VS, Marsico G, Santner T, Di Antonio M, Balasubramanian S. 2017.
1305          Machine learning model for sequence-driven DNA G-quadruplex formation. *Scientific*
1306          *Reports* **7**: 14535.
1307 Schadt EE, Turner S, Kasarskis A. 2010. A window into third-generation sequencing. *Human*
1308          *Molecular Genetics* **19**: R227-R240.
1309 Schiavone D, Guilbaud G, Murat P, Papadopoulou C, Sarkies P, Prioleau M-N, Balasubramanian S,
1310          Sale JE. 2014. Determinants of G quadruplex-induced epigenetic instability in REV1-
1311          deficient cells. *The EMBO Journal* **33**: 2507-2520.
1312 Sedlazeck FJ, Lee H, Darby CA, Schatz MC. 2018. Piercing the dark matter: bioinformatics of long-
1313          range sequencing and mapping. *Nature Reviews Genetics* **19**: 329-346.
1314 Seo J-S, Rhie A, Kim J, Lee S, Sohn M-H, Kim C-U, Hastie A, Cao H, Yun J-Y, Kim J et al. 2016.
1315          De novo assembly and phasing of a Korean human genome. *Nature* **538**: 243.
1316 Shedlock AM, Takahashi K, Okada N. 2004. SINEs of speciation: tracking lineages with retroposons.
1317          *Trends Ecol Evol* **19**: 545-553.
1318 Shiina T, Hosomichi K, Inoko H, Kulski JK. 2009. The HLA genomic loci map: expression,
1319          interaction, diversity and disease. *Journal Of Human Genetics* **54**: 15.
1320 Slotkin RK. 2018. The case for not masking away repetitive DNA. *Mobile DNA* **9**: 15.
1321 Smeds L, Kawakami T, Burri R, Bolivar P, Husby A, Qvarnström A, Uebbing S, Ellegren H. 2014.
1322          Genomic identification and characterization of the pseudoautosomal region in highly
1323          differentiated avian sex chromosomes. *Nature Communications* **5**: 5448.
1324 Smeds L, Warmuth V, Bolivar P, Uebbing S, Burri R, Suh A, Nater A, Bureš S, Garamszegi LZ,
1325          Hogner S et al. 2015. Evolutionary analysis of the female-specific avian W chromosome.
1326          *Nature Communications* **6**: 7330.
1327 Smith JJ, Timoshevskaya N, Timoshevskiy VA, Keinath MC, Hardy D, Voss SR. 2019. A
1328          chromosome-scale assembly of the axolotl genome. *Genome Research* **29**: 317-324.
1329 Su X-z, Wu Y, Sifri CD, Wellems TE. 1996. Reduced Extension Temperatures Required for PCR
1330          Amplification of Extremely A+T-rich DNA. *Nucleic Acids Research* **24**: 1574-1575.
1331 Suh A, Smeds L, Ellegren H. 2018. Abundant recent activity of retrovirus-like retrotransposons
1332          within and among flycatcher species implies a rich source of structural variation in songbird
1333          genomes. *Molecular Ecology* **27**: 99-111.
1334 Sun H, Rowan BA, Flood PJ, Brandt R, Fuss J, Hancock AM, Michelmore RW, Huettel B,
1335          Schneeberger K. 2019. Linked-read sequencing of gametes allows efficient genome-wide
1336          analysis of meiotic recombination. *Nature Communications* **10**: 4310.
1337 Tanaka Y, Asano T, Kanemitsu Y, Goto T, Yoshida Y, Yasuba K, Misawa Y, Nakatani S, Kobata K.
1338          2019. Positional differences of intronic transposons in pAMT affect the pungency level in
1339          chili pepper through altered splicing efficiency. *The Plant Journal* **0**.
1340 Thomma BPHJ, Seidl MF, Shi-Kunne X, Cook DE, Bolton MD, van Kan JAL, Faino L. 2016. Mind
1341          the gap; seven reasons to close fragmented genome assemblies. *Fungal Genetics and Biology*
1342          **90**: 24-30.
1343 Tilak M-K, Botero-Castro F, Galtier N, Nabholz B. 2018. Illumina Library Preparation for
1344          Sequencing the GC-Rich Fraction of Heterogeneous Genomic DNA. *Genome Biology and*
1345          *Evolution* **10**: 616-622.
1346 Vertebrate Genome Project V. https://vertebrategenomesproject.org.
1347 Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, Cuomo CA, Zeng Q, Wortman
1348          J, Young SK et al. 2014. Pilon: An Integrated Tool for Comprehensive Microbial Variant
1349          Detection and Genome Assembly Improvement. *PLOS ONE* **9**: e112963.

Wallberg A, Bunikis I, Pettersson OV, Mosbech M-B, Childers AK, Evans JD, Mikheyev AS, Robertson HM, Robinson GE, Webster MT. 2019. A hybrid de novo genome assembly of the honeybee, Apis mellifera, with chromosome-length scaffolds. *BMC Genomics* **20**: 275.

Warren RL, Yang C, Vandervalk BP, Behsaz B, Lagman A, Jones SJM, Birol I. 2015. LINKS: Scalable, alignment-free scaffolding of draft genomes with long reads. *GigaScience* **4**: 35.

Warren WC, Hillier LW, Tomlinson C, Minx P, Kremitzki M, Graves T, Markovic C, Bouk N, Pruitt KD, Thibaud-Nissen F et al. 2017. A New Chicken Genome Assembly Provides Insight into Avian Genome Structure. *G3: Genes|Genomes|Genetics* **7**: 109-117.

Waterhouse RM, Seppey M, Simão FA, Manni M, Ioannidis P, Klioutchnikov G, Kriventseva EV, Zdobnov EM. 2017. BUSCO Applications from Quality Assessments to Gene Prediction and Phylogenomics. *Molecular Biology and Evolution* **35**: 543-548.

Watson M, Warr A. 2019. Errors in long-read assemblies can critically affect protein prediction. *Nature Biotechnology* **37**: 124-126.

Weisenfeld NI, Kumar V, Shah P, Church DM, Jaffe DB. 2017. Direct determination of diploid genome sequences. *Genome Res* **27**: 757-767.

Weissensteiner MH, Bunikis I, Catalan A, Francoijs K-J, Knief U, Heim W, Peona V, Pophaly S, Sedlazeck F, Suh A et al. 2019. The population genomics of structural variation in a songbird genus. *bioRxiv* doi:10.1101/830356: 830356.

Weissensteiner MH, Pang AWC, Bunikis I, Hoijer I, Vinnere-Petterson O, Suh A, Wolf JBW. 2017. Combination of short-read, long-read, and optical mapping assemblies reveals large-scale tandem repeat arrays with population genetic implications. *Genome Res* **27**: 697-708.

Weissensteiner MH, Suh A. 2019. Repetitive DNA: The Dark Matter of Avian Genomics. In *Avian Genomics in Ecology and Evolution: From the Lab into the Wild*, doi:10.1007/978-3-030-16477-5_5 (ed. RHS Kraus), pp. 93-150. Springer International Publishing, Cham.

Westbrook CJ, Karl JA, Wiseman RW, Mate S, Koroleva G, Garcia K, Sanchez-Lockhart M, O'Connor DH, Palacios G. 2015. No assembly required: Full-length MHC class I allele discovery by PacBio circular consensus sequencing. *Human Immunology* **76**: 891-896.

Wicker T, Sabot F, Hua-Van A, Bennetzen JL, Capy P, Chalhoub B, Flavell A, Leroy P, Morgante M, Panaud O et al. 2007. A unified classification system for eukaryotic transposable elements. *Nature Reviews Genetics* **8**: 973-982.

Willard HF, Waye JS. 1987. Hierarchical order in chromosome-specific human alpha satellite DNA. *Trends in Genetics* **3**: 192-198.

Wyman SK, Jansen RK, Boore JL. 2004. Automatic annotation of organellar genomes with DOGMA. *Bioinformatics* **20**: 3252-3255.

Xu L, Auer G, Peona V, Suh A, Deng Y, Feng S, Zhang G, Blom MPK, Christidis L, Prost S et al. 2019. Dynamic evolutionary history and gene content of sex chromosomes across diverse songbirds. *Nature Ecology & Evolution* **3**: 834-844.

Yazdi HP, Ellegren H. 2018. A Genetic Map of Ostrich Z Chromosome and the Role of Inversions in Avian Sex Chromosome Evolution. *Genome Biology and Evolution* **10**: 2049-2060.

Yeo S, Coombe L, Warren RL, Chu J, Birol I. 2017. ARCS: scaffolding genome drafts with linked reads. *Bioinformatics* **34**: 725-731.

Yoshimura J, Ichikawa K, Shoura MJ, Artiles KL, Gabdank I, Wahba L, Smith CL, Edgley ML, Rougvie AE, Fire AZ et al. 2019. Recompleting the Caenorhabditis elegans genome. *Genome Research* **29**: 1009-1022.

Zhang G Li C Li Q Li B Larkin DM Lee C Storz JF Antunes A Greenwold MJ Meredith RW et al. 2014. Comparative genomics reveals insights into avian genome evolution and adaptation. *Science* **346**: 1311-1320.

Zhang Y, Cheng TC, Huang G, Lu Q, Surleac MD, Mandell JD, Pontarotti P, Petrescu AJ, Xu A, Xiong Y et al. 2019. Transposon molecular domestication and the evolution of the RAG recombinase. *Nature* **569**: 79-84.

1400　Zheng GXY, Lau BT, Schnall-Levin M, Jarosz M, Bell JM, Hindson CM, Kyriazopoulou-
1401　　　　Panagiotopoulou S, Masquelier DA, Merrill L, Terry JM et al. 2016. Haplotyping germline
1402　　　　and cancer genomes with high-throughput linked-read sequencing. *Nature Biotechnology* **34**:
1403　　　　303.
1404　Zhou Q, Zhang J, Bachtrog D, An N, Huang Q, Jarvis ED, Gilbert MTP, Zhang G. 2014. Complex
1405　　　　evolutionary trajectories of sex chromosomes across bird taxa. *Science* **346**: 1246338.
1406