1  *De Novo* **PacBio long-read and phased avian genome assemblies correct and**

2  **add to reference genes generated with intermediate and short reads**

3

4  Jonas Korlach[1*], Gregory Gedman[2], Sarah B. Kingan[1], Chen-Shan Chin[1], Jason T.

5  Howard[2], Jean-Nicolas Audet[2,3], Lindsey Cantin[2], and Erich D. Jarvis[2,4*]

6

7  [1]Pacific Biosciences, Menlo Park, CA 94025, USA; [2]Laboratory of Neurogenetics of Language,

8  The Rockefeller University, New York, NY 10065, USA; [3]Department of Biology, McGill

9  University, Montreal, Quebec H3A 1B1, Canada; [4]Howard Hughes Medical Institute, Chevy

10 Chase, MD 20815, USA

11

12 *Corresponding authors: jkorlach@pacb.com and ejarvis@rockefeller.edu

13

14 Jonas Korlach, Ph.D.

15 Chief Scientific Officer

16 Pacific Biosciences

17 1305 O'Brien Drive

18 Menlo Park, CA 94025

19 650-521-8006

20

21 Erich D. Jarvis, Ph.D.

22 Investigator, Howard Hughes Medical Institute

23 Professor, The Rockefeller University, Box 54

24 1230 York Avenue, New York, New York 10065

25 212 327-8806

26 ORCID: 0000-0001-8931-5049

27

28

29

## Abstract

**Background:** Reference quality genomes are expected to provide a resource for studying gene structure, function, and evolution. However, often genes of interest are not completely or accurately assembled, leading to unknown errors in analyses or additional cloning efforts for the correct sequences. A promising solution is long-read sequencing. Here we tested PacBio-based long-read sequencing and diploid assembly for potential improvements to the Sanger-based intermediate-read zebra finch reference and Illumina-based short-read Anna's hummingbird reference, two vocal learning avian species widely studied in neuroscience and genomics.

**Results:** With DNA of the same individuals used to generate the reference genomes, we generated diploid assemblies with the FALCON-Unzip assembler, resulting in contigs with no gaps in the megabase range, representing 150-fold and 200-fold improvements over the current zebra finch and hummingbird references, respectively. These long-read and phased assemblies corrected and resolved what we discovered to be numerous misassemblies in the references, including missing sequences in gaps, erroneous sequences flanking gaps, base call errors in difficult to sequence regions, complex repeat structure errors, and allelic differences between the two haplotypes. These improvements were validated by single long genome and transcriptome reads, and resulted for the first time in completely resolved protein-coding genes widely studied in neuroscience and specialized in vocal learning species.

**Conclusions:** These findings demonstrate the impact of long reads, sequencing of previously difficult-to-sequence regions, and phasing of haplotypes on generating high quality assemblies necessary for understanding gene structure, function, and evolution.

**Keywords:** De novo genome assembly, long reads, SMRT Sequencing, brain, language.

## Background

Having available genomes of species of interest provides a powerful resource to rapidly conduct investigations on genes of interest. For example, using the original Sanger method to sequence genomes of the two most commonly studied bird species, the chicken [1] and zebra finch [2], has impacted many studies. The zebra finch is a vocal learning songbird, with the rare ability to imitate sounds similar to as humans do for speech; comparative analyses of genes in its genome has allowed insights into the mechanisms and evolution of spoken-language in humans [2-4]. With the advent of more cost-effective next generation sequencing technologies using short reads, 10-fold more vertebrate genomes were sequenced [5], with one large successful project being the Avian Phylogenomics Consortium, which generated genomes of 45 new bird species across the family tree and several reptiles [6, 7]. The consortium was successful in conducting comparative genomics and phylogenomics with populations of genes [8-11]. However, when necessary to dig deeper into individual genes, it was discovered that many were incompletely

70 assembled or contained apparent misassemblies. For example, the *DRD4* dopamine receptor was
71 missing in half of the assemblies, in part due to sequence complexity [12]. The *EGR1* immediate
72 early gene transcription factor, a commonly studied gene in neuroscience and in vocal learning
73 species, was missing the promoter region in a GC-rich region in every bird genome we examined
74 (including the Sanger-based assemblies). Another immediate early gene, *DUSP1*, with
75 specialized vocalizing-driven gene expression in song nuclei of vocal learning species, has
76 microsatellite sequences in the promoters of vocal learning species that are missing or
77 misassembled, requiring single-molecule cloning and sequencing to resolve [13]. Such errors
78 create a great amount of effort to clone, sequence, and correct assemblies of individual genes of
79 interest.
80    High-throughput, single-molecule, long-read sequencing shows promise to alleviate these
81 problems [14-16]. As part of an effort to evaluate standards for the G10K vertebrate
82 (https://genome10k.soe.ucsc.edu) and the B10K bird (http://b10k.genomics.cn/index.html)
83 genome projects, here we applied PacBio single-molecule long-read (1,000-60,000 bp)
84 sequencing and diploid assembly on two vocal learning species, the zebra finch previously
85 assembled with Sanger-based intermediate reads (700-1,000 bp), and the Anna's hummingbird
86 previously assembled with Illumina-based short reads (100-150 bp). We found that the long-read
87 diploid assemblies resulted in major improvements in genome completeness and contiguity, and
88 completely resolved the problems in all of our genes of interest.
89
90

## Results

91

92

**The long-read assemblies result in 150-fold to 200-fold increases in contiguity**
93
94 To generate long-read assemblies, high molecular weight DNA was isolated from muscle tissue
95 of the same zebra finch male and Anna's hummingbird female used to create the current
96 reference genomes [2, 8]. The DNA was sheared, 35-40 kb libraries generated, size-selected for
97 inserts >17 kb (**Fig. S1**), and then SMRT sequencing performed on the PacBio RS II instrument
98 to obtain ~96X coverage for the zebra finch (19 kb N50 read length) and ~70X for the
99 hummingbird (22 kb N50 read length; **Fig. S2**). The long reads were originally assembled with
100 an early version of the FALCON assembler that only separates very divergent regions between
101 haplotypes, and merges the remaining sequence, which we and others found unintentionally
102 introduced indels in the merged regions for some nucleotides that differed between haplotypes
103 (tested on the hummingbird; data not shown) [17]. We then re-assembled using FALCON v0.4.0
104 followed by the FALCON-Unzip module [18] to prevent indel formation and generate longer-
105 range phased haplotypes. Thus, the new assemblies, unlike the current reference assemblies, are
106 phased diploids. This PacBio-based sequencing and assembly approach does not link contigs into
107 gapped scaffolds. Scaffolding requires additional approaches, which we will report on separately
108 in a study comparing scaffolding technologies with these assemblies. The results presented here
109 were found independent of scaffolding.

For the zebra finch, the long-read approach resulted in 1159 primary haplotype contigs with an estimated total genome size of 1.14 Gb (1.2 Gb expected; [19]) and contig N50 of 5.81 Mb, representing a 108-fold reduction in the number of contigs and a 150-fold improvement in contiguity compared to the current Sanger-based reference (**Table 1A**). The diploid assembly process produced 2188 associated, or secondary, haplotype contigs (i.e. haplotigs) with an estimated length of 0.84 Gb and contig N50 of 1.14 Mb (**Table 1A**), implying that about 75% of the genome contained sufficient heterozygosity to be phased into haplotypes by FALCON-Unzip. Since in FALCON-Unzip, the primary contigs are chosen as the longest path (i.e. longest contig) through the assembly string graph, whether it is from the maternal or paternal chromosome, since the latter information is not known; the secondary haplotigs are thus by definition shorter and more in number, resulting in lower contiguity for the haplotigs. Regions of the genome with very low heterozygosity remain as collapsed haplotypes in the primary contigs.

The PacBio long-read assembly for the hummingbird was of similar quality, with 1076 primary contigs generating a primary haploid genome size of 1.01 Gb (1.14 Gb expected; [19]), and a contig N50 of 5.36 Mb, representing a 116-fold reduction in the number of contigs and a 201-fold improvement in contiguity over the reference (**Table 1B**). The length of the assembled secondary haplotigs for the hummingbird was similar to that of the primary contig backbone (1.01 Gb) with a contig N50 of 1.01 Mb (**Table 1B**) indicating that there was sufficient heterozygosity to phase most of the diploid genome into the two haplotypes.

| Species | Reference assembly | PacBio-based primary haplotype | Improvement | PacBio-based secondary haplotype |
|---|---|---|---|---|
| **A. Zebra finch** | **Sanger-based** | | | |
| Number of contigs | 124,806 | 1,159 | **- 108 fold** | 2,188 |
| Contig N50 | 38,639 bp | 5,807,022 bp | **+ 150 fold** | 2,740,176 bp |
| Total size | 1,232,135,591 bp | 1,138,770,338 bp | | 843,915,757 bp |
| | | | | |
| **B. Hummingbird** | **Illumina-based** | | | |
| Number of contigs | 124,820 | 1,076 | **- 116 fold** | 4,895 |
| Contig N50 | 26,738 bp | 5,366,327 bp | **+ 201 fold** | 1,073,631 bp |
| Total size | 1,105,676,412 bp | 1,007,374,986 bp | | 1,013,746,550 bp |

**Table 1:** *De novo* genome assembly statistics comparing intermediate-read length and short-read length assemblies with the long-read assemblies. (A) Zebra finch intermediate-read length (Sanger-based, NCBI accession # GCF_000151805, version 3.2.4) compared to the long-read length PacBio-based assembly. (B) Anna's hummingbird short-read length (Illumina-based, accession # GCF_000699085) compared to the long-read length PacBio-based assembly. Improvement is calculated between the 2nd and 3rd columns for the primary PacBio-based haplotype. The higher number of contigs in the secondary haplotype (5th column) is a result of the arbitrary assignment of shorter haplotypes to the haplotig category ([18] and main text).

For comparison, using FALCON without the Unzip module [17] resulted in assemblies with high contiguity for the primary contigs (e.g. N50 5.9 Mb for the hummingbird), but much lower for

4

142 the associated contigs (N50 40 kb). Typical FALCON parameterization allows overlaps between
143 error-corrected reads that differ by ~5%, and therefore even somewhat divergent haplotypes are
144 collapsed (i.e. merged). Correspondingly, we observed smaller overall associated total assembly
145 sizes (204 Mb for the zebra finch, 187 Mb for the hummingbird, respectively) compared to the
146 more fully phased primary contig assembly sizes (1.11 Gb for the zebra finch, 1.05 Gb for the
147 hummingbird, respectively; **Table 1**). The FALCON-Unzip module generates larger haplotigs
148 through phasing of heterozygous SNPs, and also resolves smaller structural allelic variation. For
149 these reasons, all subsequent analyses were conducted on the more phased FALCON-Unzip
150 assemblies.
151

**The long-read assemblies have more complete conserved protein coding genes**

153 To assess gene completeness, we analyzed 248 highly conserved eukaryotic genes from the
154 CEGMA human set [20, 21] in each of the assemblies. Both the PacBio-based zebra finch and
155 hummingbird phased assemblies showed improved resolution of these gene sequences, with a
156 close to doubling (~71%) for the zebra finch and 26% increase for the hummingbird in the
157 number of complete or near-complete (>95%) CEGMA genes assembled, compared to the
158 references (**Fig. 1A**). Because updating the CEGMA gene sets was recently discontinued due to
159 lack of continued funding and ease of use (http://www.acgt.me/blog/2015/5/18/goodbye-cegma-
160 hello-busco), we also searched for a set of conserved, single-copy genes from the orthoDB9 [22]
161 gene set using the recommended replacement BUSCO pipeline [23]. When assessed using the
162 BUSCO v2.0 pipeline on a set of 303 single-copy conserved eukaryotic genes, we observed more
163 modest improvements (~10%) in the number of complete genes in the zebra finch (and no
164 change with the hummingbird; **Fig. 1B**), and barely any change (1-3%) when using a newly
165 generated BUSCO set of 4915 avian genes (**Fig. 1C**). However, we believe that the moderate
166 increase or no change is due to the fact that much of the BUSCO gene sets were generated from
167 incomplete genome assemblies with short- to intermediate-length reads; for example, the 4915
168 protein coding avian gene set is generated mostly from the 40+ avian species that the Avian
169 Phylogenomics Project sequenced with short reads [8], including the reference hummingbird
170 [24]. Supporting this view, we extracted the overlapping orthologous genes in the different
171 CEGMA and BUSCO datasets, and found that the CEGMA genes are on average significantly
172 longer than their BUSCO counterparts (**Fig. S3**). When we manually examined randomly chosen
173 genes, many of the BUSCO protein coding sequences were truncated relative to the
174 corresponding CEGMA gene and the PacBio-based assemblies (e.g. the ribosomal protein
175 RLP24 aves BUSCO gene is 117 a.a., whereas the CEGMA & PacBio assembly are 163 a.a.).
176 When compared to the CEGMA 303 eukaryotic set that includes several higher-quality genome
177 assemblies, the PacBio-based assemblies had very few fragmented genes compared to the
178 Sanger-based and Illumina-based assemblies (**Fig. 1B**). Thus, the new PacBio-based assemblies
179 have the potential to upgrade the BUSCO set with more complete and more accurately
180 assembled genes, a conclusion supported by analyses below.
181

**The long-read assemblies have greater and more accurate transcriptome and regulome representations**

To assess transcriptome gene completeness by an approach that does not depend on other species' genomes, we aligned zebra finch brain paired-end Illumina RNA-Seq reads to the zebra finch genome assemblies using TopHat2 [25]. We generated the RNA-Seq data from microdissected RA song nuclei, a region that has convergent gene expression specialization with the human laryngeal motor cortex (LMC) involved in speech production (**Fig. S4**; [4]). The PacBio-based assembly (primary haplotype) resulted in a ~7% increase in total transcript read mappings compared to the Sanger-based reference (**Fig. 2A**), suggesting more genic regions available for read alignments. This was explained by a decrease in unmapped reads and an increase in reads that mapped to the genome in multiple locations (2 or more) compared to the Sanger-based reference (**Fig. 2B**), supporting the idea that the long-read assemblies recovered more repetitive or closely related gene orthologs. The PacBio assembly also resulted in ~6% more concordant aligned paired-end reads (**Fig. 2A**), indicating a more structurally accurate assembly compared to the Sanger-based reference. RNA-Seq data from the other principle brain song nuclei (HVC, LMAN, and Area X) and adjacent brain regions containing multiple cell types (**Fig. S4A**; [26]) gave very similar results, with 7-11% increased mappings to the PacBio-based assembled genome (not shown).

Regulatory regions have been difficult to identify in the zebra finch genome, as they are often GC-rich and hard to sequence and assemble with short-read technologies. To assess the regulome, we aligned HK327ac ChIP-Seq reads generated from the RA song nucleus (see methods and [27]) to the zebra finch genome assemblies using Bowtie2 for single-end reads [28]. H3K27ac activity is generally high in active gene regulatory regions, such as promoters and enhancers [29]. Similar to the RNA-Seq transcriptome reads, there was an increase (~4%) of HK327ac Chip-Seq genomic reads that mapped to the PacBio-based assembly (primary haplotype) compared to the Sanger-based reference (**Fig. 2A**). However, unlike the RNA-Seq transcript reads, the ChIP-Seq genomic reads showed a significant 10% increase in unique mapped reads with a concomitant decrease in multiple mapped reads (**Fig. 2B**). We believe this difference is due to technical reasons. The RNA-Seq data was paired-end reads mapped to the genome, whereas the ChIP-Seq data was single-end reads; when just using the single-ends of the RNA-Seq data, the multiple-mapped increase to the Pacbio-based assembly was not detected ($p$=0.3, paired t-test, n=5), indicating that repetitive sequence in the paired end data influences read mapping. Overall, these findings are consistent with the PacBio-based assembly having a more complete and structurally accurate assembly for both coding and regulatory non-coding genomic regions.

**Completion and correction of genes important in vocal learning and neuroscience research**

The genome-wide analyses above demonstrate improvements to overall genome assembly quality using long reads, but they do not inform about real-life experiences with individual genes. We undertook a detailed analysis of four of our favorite genes that have been widely

222 studied in neuroscience and in vocal learning/language research in particular: *EGR1*, *DUSP1*,
223 *FOXP2*, and *SLIT1*.

224

225 **EGR1**. The early growth response gene 1 (*EGR1*) is an immediate early gene transcription factor
226 whose expression is regulated by activity in neurons, and is involved in learning and memory
227 [30]. It is up-regulated in song-learning nuclei when vocal learning birds produce song [31]; it
228 belongs to a large set of genes representing 10% of the transcribed genome that are up- or down-
229 regulated in response to activity in different cell types of the brain [27]. Studying the
230 mechanisms of regulation of *EGR1* and other immediate early genes has been an intensive area
231 of investigation [32, 33], but in all intermediate- and short-read bird genome assemblies we
232 examined thus far, part of the GC-rich promoter region is missing (**Fig. 3A, gap 1**).

233 In the zebra finch Sanger-based reference, *EGR1* is located on a 5.7 kb contig (on
234 chromosome 13), bounded by the gap in the GC-rich promoter region and 2 others downstream
235 of the gene; gaps between contigs in the published reference were given arbitrary 100 Ns [2]. We
236 found that the PacBio long-read assembly resolved all three gaps in the *EGR1* locus for both
237 alleles, resulting in complete protein coding and surrounding gene bodies in a 205.5 kb primary
238 contig and a 129.1 kb secondary haplotig (**Fig. 3B**; **Fig. S5A**). The promoter region gap was
239 resolved by PacBio-based 804 bp of 70.1% GC-rich sequence (**Fig. 3B, black**). In addition, to
240 the left and right of this gap there were 241 bp total of low quality sequence (<QV40; **Fig 3A,**
241 **blue; 3B, red**) that was not supported by the PacBio reads. For the second gap, located ~2.2 kb
242 downstream of the *EGR1* gene, there was an adjacent 210 bp low-similarity tandem repeat region
243 that also had low quality scores and was not supported by the PacBio-based reads (**Fig 3A,B,**
244 **gap 2**). The third 100 N gap, located ~3.5 kb downstream of the *EGR1* gene, was resolved by 18
245 bp of sequence in the PacBio assembly (**Fig. 3B, gap 3**). The PacBio-based differences in the
246 assembly were supported by numerous long-read (>10,000 bp) molecules that extended through
247 the entire gene, spanning all three gaps (**Fig. S6A**). The two haplotypes were >99.8% identical
248 over the region shown (**Fig. 3B**), with only one synonymous heterozygous SNP in the coding
249 sequence (G at position 169,283 in the primary contig 405; T at position 92,478 in secondary
250 haplotig 405_002; tick mark in **Fig. 3B**).

251 In the Illumina-based hummingbird reference, *EGR1* was represented by 3 contigs
252 separated by 2 large gaps of 544 Ns and 1987 Ns respectively (**Fig. 3C**), in a large 2.98 Mb
253 scaffold. In contrast, in the PacBio-based hummingbird assembly, *EGR1* was fully resolved in a
254 large 810 kb contig (**Fig. 3C**). Gene prediction (using Augustus [34]) yielded a protein of the
255 same length as the finch EGR1 protein (510 a.a.), and with high (93%) sequence identity (**Fig.**
256 **3D**). The PacBio-based assembly revealed that the larger gap in the Illumina-based assembly
257 harbors the beginning of the *EGR1* gene, including the entire first exon, two thirds of the first
258 intron, and the GC-rich promoter region (**Fig. 3C, black**). Due to this gap in the reference, the
259 corresponding NCBI gene prediction (accession XP_008493713.1) instead recruited a stretch of
260 sequence ~7 kb upstream of the gap, predicting a first exon with no sequence homology to *EGR1*
261 in the PacBio-based assembly or in other species (**Fig. 3C & D**). Upstream of this gap in the

262 Illumina-based assembly was also a 200 bp tandem repeat that was not supported by the PacBio
263 sequence reads and the assembly (**Fig. 3C, red; Fig. S5B**). The PacBio-based assembly was
264 further validated by single-molecule Iso-Seq mRNA long-reads of *EGR1* from a closely related
265 species (the Ruby-throated hummingbird; [35]) that fully contained both predicted exons (**Fig.
266 S6B**). The PacBio-based assembly did not generate a secondary haplotype for this region,
267 indicating that the two alleles are identical or nearly identical for the entire 810 kb contig in the
268 individual sequenced. Upstream and downstream of a high homology region that includes the
269 *EGR1* gene, there was little sequence homology between the hummingbird and zebra finch
270 assemblies (**Fig. S7**).

271        These findings indicate that relative to the intermediate- and short-read assemblies, the
272 PacBio-based long-read assembly can fill in missing gaps in a previously hard-to-sequence GC-
273 rich regulatory region, eliminate low quality erroneous sequences and base calls at the edges of
274 gaps in the Sanger-based assembly, and eliminate erroneous tandem duplications adjacent to
275 gaps, all preventing inaccurate gene predictions. In addition, using one species as a reference to
276 help assemble another may not work for such a gene, as the surrounding sequence to the gene
277 body in these two Neoaves species is highly divergent.

278

279 *DUSP1*. The dual specificity phosphatase 1 (*DUSP1*) is also an immediate early gene, but one
280 that regulates the cellular responses to stress [36]. In all species examined thus far it is mostly
281 up-regulated by activity in the highly active thalamic-recipient primary sensory neurons of the
282 cortex (i.e. mammal cortex layer 4 neurons and the comparable avian intercalated pallial
283 neurons), but within the motor pathways, it is only up-regulated to high levels by activity in the
284 vocal learning circuits of vocal learners [13, 37]. This specialized regulation in vocal learning
285 circuits has been proposed to be associated with convergent microsatellite sequences found in the
286 upstream promoter region of the gene mainly in vocal learning species [13]. This was determined
287 by PCR-cloning of single genomic molecules from multiple species, because the reference
288 assemblies did not have this region properly assembled [13].

289        In the zebra finch Sanger-based reference, *DUSP1* is located on the chromosome 13
290 scaffold, separated in 3 contigs, with 2 gaps, all surrounded by low quality sequences (**Fig. 4A**).
291 The NCBI gene prediction of this assembly resulted in 4 exons with 322 a.a. (XP_002192168.1),
292 which is ~13% shorter than *DUSP1* homologs of other species, e.g. chicken (369 a.a., Genbank
293 accession NP_001078828), rat (367 a.a., NP_446221), and human (367 a.a, NP_004408). The 2
294 gaps coincide with the end of the first predicted exon and the beginning of the third predicted
295 exon (**Fig. 4A**). An additional gap upstream of the coding sequence falls within the known
296 microsatellite repeat region (**Fig. 4A**). The PacBio-based assembly produced a completely
297 resolved region for both alleles, in an 8.4 Mb primary contig and an 8.0 Mb secondary haplotig
298 (**Fig. 4B, Fig. S8A**). The Augustus gene prediction resulted in a protein with 4 exons but now
299 larger, 369 a.a., that was homologous across its length to *DUSP1* of other vertebrate species
300 (e.g., 96% with chicken GGv5 assembly, also recently updated with long reads). Comparing the
301 two assemblies revealed that: 1) the first exon in the Sanger-based reference is truncated by 28

8

302 a.a. in the gap; 2) near the edge of that truncation are three a.a. that appear to be errors (**Fig. 4**; 303 residues 81, 89, and 98), as they are different from genomes of other songbird species using high 304 coverage Illumina reads (**Fig. S9A**), with strong support in the zebra finch PacBio reads (**Fig.** 305 **S9B**); 3) the second exon and adjacent intron is missing an 80.8% GC-rich 0.46 kb sequence in 306 the reference, and is instead replaced by a 1.7 kb contig of a partially repeated sequence from the 307 microsatellite region upstream of *DUSP1* (R2' in **Fig. 4B**), part of which was erroneously 308 recruited in the second exon of the NCBI reference gene prediction (**Fig. 4D**); and 4) the 309 microsatellite repeat itself is erroneously partially duplicated in the reference, flanking both sides 310 of gap 1 (R1'' and R2'' in **Fig. 4B**). The PacBio-based phased assembly revealed why both 311 instances of R' are not identical in the reference, because they in fact belong to the different 312 haplotypes: the 1.7 kb contig corresponds to the upstream region in the primary PacBio 313 haplotype (contig 32) whereas the actual upstream region in the reference corresponds to the 314 upstream region in the secondary PacBio haplotype (contig 32_022) (**Fig. 4B**). This main 315 microsatellite region is 76 bp longer (796 *vs.* 720 bp) in the primary haplotype, and the 316 neighboring smaller upstream microsatellite contains 3 additional 20-21 bp repeats (11 *vs.* 8) in 317 the primary haplotype (**Fig. S10A**). Within the protein coding sequence there were four 318 synonymous heterozygous SNPs between haplotypes (not shown). The assembled sequence of 319 the published Sanger-based single clone (AB574425.1) [13] is more consistent with the PacBio-320 based genome assembly than the Sanger-based reference genome assembly (**Fig. S11A**), and 321 does not support the erroneous tandem duplications and misplacements of repeat sequences in 322 the latter. Differences in the Sanger-based sequenced clone with the PacBio-based assembly are 323 that the main microsatellite region is smaller (~320 bp) and the upstream 20-21 bp microsatellite 324 has 10 repeats (instead of 11 or 8), which is consistent with the repeats differing in number 325 between haplotypes (this study) and also individuals [13]. We note that the *DUSP1* haplotypes in 326 the zebra finch PacBio-based FALCON-Unzip assembly are 4.8% divergent which was below 327 the 5% threshold for allelic segregation in the FALCON assembly without using the Unzip 328 module, but were successfully resolved when using FALCON-Unzip.

329 In the hummingbird Illumina-based assembly, the *DUSP1* region was represented by 2 330 contigs separated by a large 1005 N gap (**Fig. 4C**), on a 7 Mb scaffold. In the PacBio-based 331 assembly, the entire gene was fully resolved (**Fig. 4C; Fig. S8B)**, in a much larger gapless 12.8 332 Mb contig (the second allele is fully resolved in a 3.8 Mb contig). Comparing the two assemblies 333 revealed that the gap of the Illumina-based reference contains about half of the *DUSP1* gene, 334 including the first two exons and introns, and ~380 bp upstream of the start of the gene (**Fig.** 335 **4C**). As a result, the corresponding NCBI gene prediction (XP_008496991.1) recruited a 336 sequence ~44 kb upstream predicting 46 a.a. with no sequence homology to *DUSP1* of other 337 species, whereas the PacBio-based assembly yielded a 369 a.a. protein with 99% sequence 338 identity to the PacBio-based zebra finch and chicken *DUSP1* (**Fig. 4D**). A 200 bp tandem repeat 339 in the Illumina-based assembly downstream of the gap, erroneously in exon 3, is a misplaced 340 copy of the microsatellite region (**Fig. 4C; Fig. S8B**). This is the reason why two thirds of exon 341 3 is erroneously duplicated in the NCBI protein prediction (**Fig. 4D**). These differences in the

PacBio-based assembly were validated by single-molecule Iso-Seq mRNA long-reads (**Fig. S12A**) and a Sanger-based assembly of a single clone (AB574427.1; **Fig. S11B**) of *DUSP1*. The PacBio assemblies also revealed that the microsatellite region was significantly shorter in the hummingbird (~270 bp) than in the zebra finch genome (~1100 bp; **Fig. S10B**).

These findings in both species demonstrate that intermediate- and short-read assemblies not only have gaps with missing relevant repetitive microsatellite sequence, but that short-read misassemblies of these repetitive sequences lead to erroneous protein coding sequence predictions. Further, not only does the long-read assembly resolve them, but it helps generate a diploid assembly that resolves allelic differences and prevents erroneous assembly duplications and misplacement errors between haplotypes.

*FOXP2*. The forkhead box P2 (*FOXP2)* gene plays an important role in spoken-language acquisition [38]. In humans, a point mutation in the protein coding binding domain in the KE family [39] as well as deletions in the non-coding region of *FOXP2* [40] results in severe spoken language impairments in heterozygous individuals (homozygous is lethal). In songbirds, FOXP2 expression in the Area X song nucleus is differentially regulated by singing activity and during the song learning critical period, and is necessary to properly imitate song [41-43]. In mice, although vocalizations are mainly innate, animals with the KE mutation demonstrate a syntax apraxia-like deficit in syllable sequencing similar to that of humans [44, 45]. Thus, *FOXP2* has become the most studied gene for understanding the genetic mechanisms and evolution of spoken language [46], yet we find that the very large gene body of ~400 kb is incompletely assembled (**Fig. 5A**).

In the zebra finch Sanger-based reference, *FOXP2* is located on the chromosome 1A scaffold and separated into 10 contigs (1 to 231 kb in length) with nine 100 N gaps (**Fig. 5A**). These include 2 gaps immediately upstream of the first exon, making the beginning of the gene poorly resolved. The provisional RefSeq mRNA for *FOXP2* (NM_001048263.1) contains 19 exons and encodes a 711 a.a. protein (NP_001041728.1). In the PacBio-based assembly, the entire 400 kb gene is fully resolved for both haplotypes in 21.5 Mb and 7.6 Mb contigs, respectively (**Fig. S13A**). As observed in the previous examples, low quality sequences of various sizes surrounding all 9 gaps in the Sanger-based reference were unsupported by the PacBio higher quality data, resulting in a total of 2509 bp of corrected sequence in the PacBio-based primary haplotype (**Fig. 5B**). The two filled gaps in the upstream region and the next gap in the first intron were GC-rich (77.6%, 66.5%, and 67.8%, respectively; **Fig. 5A,C**), indicative of the likely cause of the poor quality Sanger-based read coverage (**Fig. 5D**). The DNA sequence between the two assembled PacBio haplotypes was >99% similar across the entire 400 kb *FOXP2* gene, and identical over the coding sequence, with differences occurring in the more complex non-coding gaps that were difficult to sequence and assemble by the Sanger method (**Fig. 5B** *61 nucleotide differences total). The predicted protein sequence from the PacBio-based assembly is identical to the predicted Sanger-based reference (NP_001041728.1), with the exception of a.a. residue 42 (threonine *vs.* serine; **Fig. S14A**). The PacBio nucleotide call also

382 exists in the mRNA sequence of another zebra finch animal in NCBI (NM_001048263.2) and in
383 other avian species we examined, and is thus likely a base call error in the Sanger-based zebra
384 finch reference.

385       In the hummingbird Illumina-based assembly, as expected with short-read assemblies
386 relative to the Sanger-based zebra finch reference, the *FOXP2* gene was even more fragmented,
387 in 23 contigs (ranging 0.025 to 2.28 kb in lengths) with 22 gaps (**Fig. S13B**). The two largest
388 gaps encompass the beginning of the gene and first (non-coding) exon, resulting in
389 corresponding low quality predicted mRNA (XM_008496149.1). The predicted protein
390 (XP_008494371.1) includes an introduced correction (a.a. 402; **Fig. S14A**, X nucleotide) to
391 account for a genomic stop codon, and an 88 N gap within exon 6 that artificially splits the exon
392 into two pieces (**Fig. S14B**). In the hummingbird PacBio-based assembly, the *FOXP2* gene is
393 fully resolved and phased into two haplotype contigs of 3.2 Mb each (**Fig. S13B**). The erroneous
394 stop codon is corrected (2170128C [ctg 110] and 2183088C [ctg 110_009], instead of 841788T
395 [Illumina assembly scaffold 125]), and exon 6 is accurately contiguous, removing the gap and an
396 additional 22 bp of erroneous tandem repeat sequence adjacent to the gap (**Fig. S14B & C**). The
397 PacBio-based assembly also corrects three other instances of erroneous tandem duplications over
398 the gene region in the Illumina-based assembly, as well as removes a 462 bp stretch of sequence
399 adjacent to a long homonucleotide A stretch in intron 1 of the Illumina-based assembly (position
400 972040; **Fig. S15A**). These PacBio-based differences in the assembly were validated by single-
401 molecule Iso-Seq mRNA long-reads of *FOXP2* (**Fig. S12B**). The two PacBio assembled
402 haplotypes are >99% similar, with one heterozygous SNP (2172601T (contig 110) *vs.* 2185560A
403 (contig 110_009)) in exon 6 that is silent, and a 708 bp deletion in the secondary haplotype
404 (contig 110_009 [at position 2128952] relative to contig 110; **Fig. S15B**). The Illumina-based
405 assembly has the deleted allele.

406       These findings replicate those of the previously discussed genes, and in addition show
407 that the PacBio-based assembly can fully resolve very large genes, resolve erroneous assembled
408 sequences in gaps due to repeats or homonucleotide stretches, and reveal large haplotype
409 differences. The phased diploid assembly also avoids the possibility of large missed sequences in
410 a haploid only assembly due to deletions in one allele.

411

412 *SLIT1*. Slit homolog 1 (*SLIT1*) is a repulsive axon guidance ligand for the *ROBO1* receptor, and
413 is involved in circuit formation in the developing brain [47]. Recently, *SLIT1* was shown to have
414 convergent specialized down-regulated expression compared to the surrounding brain region in
415 the RA song nucleus of all independently evolved vocal learning bird lineages and in the
416 analogous human LMC [4, 48] (**Fig. S4**), indicating a potential role of *SLIT1* in the evolution and
417 formation of vocal learning brain circuits. A fully resolved *SLIT1*, including regulatory regions,
418 is necessary to assess the mechanisms of its specialized regulation in vocal learning brain
419 regions.

420       In the zebra finch Sanger-based reference, *SLIT1* is located on chromosome 6, split
421 among 8 contigs with 7 gaps, and 7 additional contigs and gaps surrounding the ~40 kb gene

11

(**Fig. 6A**). The *SLIT1* gene is complex, with over 35 exons. We noted an incomplete predicted protein of the reference (XP_012430014.1) relative to some other species (chicken [NM_001277336.1], human [NM_003061.2], and mouse [NM_015748.3]); our *de novo* gene predictions of the reference also resulted in a truncated protein with two missing exons (**Fig. 6B**). The PacBio-based assembly fully resolved and phased the gene region, in two alleles on 15.7 Mb and 5.6 Mb contigs, respectively, and completely recovered all 35+ exons (**Fig. S16A**). Similar to above, reference sequences flanking the gaps were found to be erroneous and corrected, and an erroneous tandem duplication was also corrected (not shown). Filling in these gaps recovered the two missing exons: exon 1 within a 1 kb region of sequence in the PacBio-based assembly that is 75% GC-rich, replacing 390 bp of erroneous gap-flanking sequence; and exon 35 adjacent to a gap (**Fig. 6A,B**). A predicted exon upstream of exon 1 in a repeat region was not supported (**Fig. 6A,B**). The gene is heterozygous in the individual, with 3 codon differences between the two alleles (**Fig. 6B**, positions 90, 1006, and 1363, respectively), and an additional 24 silent heterozygous SNPs across the coding region.

In the hummingbird Illumina-based assembly, the *SLIT1* gene is separated on 9 contigs with 8 gaps ranging in length from 91 to 1018 bp, comprising 3320 bp of missing sequence, or 5.3% of the gene region (**Fig. S16B**). The PacBio-based assembly fully resolved and phased *SLIT1* into haplotypes on 9.9 Mb contigs (**Fig. S16B**). The resulting protein of 1538 a.a. has high sequence identity to the zebra finch PacBio-based *SLIT1* (95% a.a. identity; **Fig. 6B**) and the individual is homozygous for the SLIT1 protein. Comparisons revealed that as with the Sanger-based reference, the first exon (68 a.a.) is missing completely in the Illumina-based assembly (**Fig. 6B**), corresponding to a gap of 495 Ns, which the PacBio-based assembly replaced by a 567 bp 76% GC-rich sequence (**Fig. S16B**). In addition, there were two sequence errors in the Illumina-based assembly that were not found in the PacBio-based assembly or Sanger-based assemblies of other species, which resulted in erroneous amino acid predictions in the SLIT1 protein (**Fig. 6B**, positions 118 and 1381, respectively).

These findings demonstrate that long-read assemblies can fully resolve a complex multi-exon gene, as well as have a higher base-call accuracy than Sanger- or Illumina-based reads in difficult to sequence regions, including exons, leading to higher protein-coding sequence accuracy.

*Other genes.* We have manually compared several dozen other genes between the different assemblies, and found in all cases investigated errors in the Sanger-based and Illumina-based assemblies that were prevented in the PacBio-based long-read assemblies. These genes included other immediate early gene transcription factors, other genes in the *SLIT* and *ROBO* gene families, and the *SAP30* gene family. All had the same types of errors in the genes discussed above. In addition, we also found cases were genes were missing from the Sanger-based zebra finch or Illumina-based hummingbird assemblies entirely, and could have been interpreted as lost in these species. These included the DNA methyltransferase enzyme *DNMT3A* missing in the Sanger-based finch assembly and *DRD4* missing in the hummingbird assembly [12], with

462 both fully represented in the PacBio-based assemblies. We also noted cases where an assembled
463 gene was incorrectly localized on a scaffold in the Sanger-based assembly whose synteny was
464 corrected with the PacBio-based assembly, such as the vasopressin receptor AVPR1B, which
465 will be reported on in more detail separately. Data for these types of errors were not shown due
466 to space limitations, but they offer further examples of the important improvements of PacBio
467 long-read technology for generating more accurate genome assemblies.
468
469

## Discussion and Conclusions

471

472 Although the intermediate-read and short-read assemblies had correct sequences and assembled
473 regions in terms of total base pairs covered, the long-read assemblies revealed numerous errors
474 within and surrounding many genes. These errors are not simply in so-called "junk" intergenic
475 repetitive DNA known to be hard to assemble with short reads [49, 50], but within functional
476 regions of genes. **Table 2** summarizes 10 broad categories of errors we found, including gaps,
477 base call errors, gene prediction errors, missing genes, and assigns which of the three main
478 improvements prevented them in our *de novo* assemblies: long reads, SMRT sequencing reading
479 through normally difficult-to-sequence regions, and phasing of haplotypes. Some compounded
480 errors include the assemblers for the short reads sometimes erroneously inserting a repetitive
481 sequence in a non-repetitive region of a gene. These and other assembly and sequence errors, and
482 gaps in the sequences can all lead to gene and protein coding sequence prediction errors.
483     The long-read, phased assemblies prevented these problems and for the first time
484 resolved gene bodies of all the genes we examined into single, contiguous, gap-less sequences.
485 The phasing of haplotypes, although initially done to prevent a computationally introduced indel
486 error, reveal how important phasing is to prevent assembly and gene prediction errors. Thus far,
487 we have not seen an error (i.e. difference) in the genes we examined in the PacBio-based long-
488 read, phased assemblies relative to the other assemblies, with orthogonal support from both
489 PacBio-based datasets (single sequenced genomic DNA molecules, Iso-Seq mRNA molecules)
490 and other independent evidence (Illumina RNA-Seq and Sanger single clone data). With these
491 improvements, we now, for the first time, have complete and accurate assembled genes of
492 interest that can be pursued further without the need to individually and arduously clone,
493 sequence, and correct the assemblies one gene at a time.
494     Our study also highlights the value of maintaining frozen tissue or cells of the individuals
495 used to create previous reference genomes, as we could only discover some of the errors (e.g.
496 caused by haplotype differences) by long-read *de novo* genome assemblies of the same
497 individual used to create the reference. We are now using these PacBio-based assemblies with
498 several groups and companies as starting assemblies for scaffolding into phased, diploid,
499 chromosome-level zebra finch and hummingbird assemblies to upgrade the references, which
500 will be reported on separately. However, even without scaffolding, these more highly contiguous
501 assemblies will be helpful to researchers to extract more accurate assemblies of their genes of

502 interests, saving a great amount of time and energy, while adding new knowledge and biological
503 insights necessary for understanding gene structure, function, and evolution.

504

| Error type | Caused by | Sequence and/or assembly approach | Improved by |
|---|---|---|---|
| Gaps | Difficult to sequence, assembly algorithm errors with short reads | Sanger- & Illumina-based | Long reads |
| Low-quality sequences surrounding gaps | Low coverage of difficult to sequence GC-rich & other seq | Sanger- & Illumina-based | SMRT sequencing read through |
| Base call errors | Difficult to sequence GC-rich & other seq | Sanger- & Illumina-based | SMRT sequencing read through |
| Tandem, microsatellite, and other repeat errors | Difficult to assemble with short reads | Sanger- & Illumina-based | Long reads |
| Homonucleotide stretch assembly errors | Misassembly with short reads | Illumina-based | Long reads |
| InDel errors | Assembly algorithm | PacBio-based | Phasing |
| Misplaced/merged haplotype errors | Assembly algorithm | Sanger-, Illumina-, & PacBio-based unphased | Long reads and phasing |
| Gene prediction errors | All errors above, and haplotype merging errors | Sanger-, Illumina-, and PacBio-based unphased | Long reads, phasing, and coverage |
| Missing gene errors | Short- & intermediate-raw reads not able to be assembled | Sanger-, & Illumina-based | Long reads |
| Misplaced gene synteny | Insufficient sequence data around paralogous genes | Sanger-, & Illumina-based approach | Long reads |

505
506 **Table 2:** Summary of error types found in the different sequencing/assembly approaches, and the three
507 main factors that improved them in the *de novo* assemblies and gene predictions presented in this study.
508

509

## Materials & Methods

511

### DNA isolation

513 For both the zebra finch and hummingbird, frozen muscle tissue from the same animals used to
514 create the Sanger-based [2] and Illumina-based [8] references, respectively, was processed for
515 DNA isolation using the KingFisher Cell and Tissue DNA Kit (97030196). Tissue was
516 homogenized in 1 ml of lysis buffer in M tubes (Miltenyi Biotec) using the gentleMACS™
517 Dissociator at the Brain 2.01 setting for 1 minute. The cell lysate was treated with 40 ul of
518 protease K (20mg/ml) and incubated overnight. DNA was purified using the KingFisher Duo
519 system (5400100) using the built in KFDuoC_T24 DW program.

520

521 **Library preparation and sequencing**

For the zebra finch, two samples were used for library construction. Each DNA sample was mechanically sheared to 60 kb using the Megaruptor system (Diagenode). Then >30 kb libraries were created using the SMRTbell Template Prep Kit 1.0 (Pacific Biosciences), which includes a DNA Damage Repair step after size selection. Size selection was made for 15 kb for the first sample and 20 kb for the second sample, using a Blue Pippin instrument (Sage Science) according to the protocol "Procedure & Checklist – 20 kb Template Preparation Using BluePippin Size-Selection System". For the hummingbird, 70 ug of input DNA was mechanically sheared to 35 and 40 kb using the Megaruptor system, a SMRTbell library constructed, and size selected to > 17 kb with the BluePippin. Library quality and quantity were assessed using the Pippin Pulse field inversion gel electrophoresis system (Sage Science), as well as with the dsDNA Broad Range Assay kit and Qubit Fluorometer (Thermo Fisher).

SMRT sequencing was performed on the Pacific Biosciences RS II instrument at Pacific Biosciences using an on plate concentration of 125 pM, P6-C4 sequencing chemistry, with magnetic bead loading, and 360 minute movies. A total of 124 SMRT Cells were run for the zebra finch and 63 SMRT Cells for the hummingbird. Sequence coverage for the zebra finch was ~96 fold, with half of the 114 Gb of data contained in reads longer than 19 kb. For the hummingbird, coverage was ~70 fold, with half of the 40.4 Gb of data contained in reads longer than 22 kb (**Fig. S2**).

**Assembly**

Assemblies were carried out using FALCON v0.4.0 followed by the FALCON-Unzip module [18]. FALCON is based on a hierarchical genome assembly process [51]. It constructs a string graph from error-corrected PacBio reads that contains 'haplotype-fused' genomic regions as well as "bubbles" that capture divergent haplotypes from homologous genomic regions. The FALCON-Unzip module then assigns reads to haplotypes using heterozygous SNP variants identified in the FALCON assembly to generate phased contigs corresponding to the two alleles. The diploid nature of the genome is thereby captured in the assembly by a set of primary contigs with divergent haplotypes represented by a set of additional contigs called haplotigs. Genomic regions with low heterozygosity are represented as collapsed haplotypes in the primary contigs. Genome assemblies were run on an SGE-managed cluster using up to 30 nodes, where each node has 512 Gb of RAM distributed over 64 slots. The same configuration files were used for both species (**Additional file 1**). Three rounds of contig polishing were performed. For the first round, as part of the FALCON-Unzip pipeline, primary contigs and secondary haplotigs were polished using haplotype-phased reads and the Quiver consensus caller. For the second and third rounds of polishing, using the "resequencing" pipeline in SMRTLink v3.1, primary contigs and haplotigs were concatenated into a single reference and BLASR (version 3.1.0) was used to map all raw reads back to the assembly, followed by consensus calling with Arrow.

**Genome completeness**

To assess quality and completeness of the assemblies, we used a set of 248 highly conserved eukaryotic genes from the CEGMA human set (CEGMA, RRID:SCR_015055) [21] and located them in each of the assemblies compared in this study. We used the human gene set because it is the phylogenetically closest set to birds available, since all other CEGMA gene sets are from non-vertebrates. Briefly, the CEGMA human peptides were aligned to each genome using genblastA [52] (command: genblast_v138_linux_x86_64 -p genblasta -t ${genome} -q ${CEGMA_genes} -c 0.3 -e 0.00001 -gff -pid -r 1, where ${genome} is the assembly and ${CEGMA_genes} is the CEGMA file; the output file contains the alignment percentage for each gene). The regions showing homology were then used to build gene models with exonerate [53] which were then assessed for frameshifts (command: exonerate -m protein2genome --percent 30 --bestn 1 --showtargetgff --ryo ">%qi\n%tcs\n%m\n" -q CEGMA_prot.fa -t contig.fa, where CEGMA_prot.fa is a CEGMA peptide and contig.fa is the corresponding contig in the assembly). In addition, we queried each genome for a set of 303 eukaryotic conserved single-copy genes as well as from 4915 conserved single-copy genes from 40 different avian species using the BUSCOv2.0 pipeline (BUSCO , RRID:SCR_015008)[23].

To compare protein amino acid sequence size between the CEGMA and BUSCO datasets, we performed blastp of each CEGMA sequence against the ancestral proteins of the target BUSCO dataset. We took the single best hit with an e-value cut off of 0.001 and extracted the CEGMA and BUSCO protein length values. We then ran a one-sided paired Wilcoxon signed-rank test of the two lengths for each protein (using the "wilcox.test" function with "paired = T." in R).

**Gene prediction**

Gene predictions for the zebra finch PacBio-based assembly were conducted by running Augustus gene prediction software (Augustus: Gene Prediction , RRID:SCR_008417)(v3.2.2, [34]) on the contigs, and incorporating the Illumina short read RNA-Seq brain data aligned with Tophat2 (TopHat , RRID:SCR_013035)(v2.0.14, [25]) as hints for possible gene structures. The data consisted of 146,126,838 paired-end reads with an average base quality score of 36. Augustus produces a distribution of possible gene models for a given locus and models that are supported by our RNA-Seq data are given a "bonus" while the gene models not supported by RNA-Seq data are given a "penalty". This results in the gene model most informed by biological data being selected as the most likely gene model for that locus.

We did not have Illumina transcriptome data for Anna's hummingbird, so standard Augustus gene prediction (v3.2.2) was used with both chicken and human training background to determine the sequence predictions of the genes examined. The human-based predictions captured more of the divergent 5' ends of the longer genes (*SLIT1* and *FOXP2*) then the chicken-based predictions, so a combination of both were used to produce the final sequences in this manuscript.

**RNA-Seq**

601 RNA sequencing was centered around vocal learning brain regions in the zebra finch and will be
602 described in more detail in a later publication. We utilized our data here for population analyses
603 of assembly quality and for initial annotations. In brief, following modifications of a previously
604 described protocol [27], nine adult male zebra finches were kept isolated in soundproof chambers
605 for 12 hours in the dark to obtain brain tissue from silent animals. Then brains were dissected
606 from the skull and sectioned to 400 microns using a Stoelting tissue slicer (51415). The sections
607 were moved to a petri dish containing cold PBS with proteinase inhibitor cocktail
608 (11697498001). Under a dissecting microscope (Olympus MVX10), the four principle song
609 nuclei (Area X, LMAN, HVC, and RA) as well as their immediate adjacent brain regions were
610 microdissected using 2mm fine scissors and placed in microcentrifuge tubes.  The samples were
611 stored at -80 °C.  Then RNA was isolated and quantified, and samples of two birds were then
612 pooled for each replicate, resulting in 5 replicates (one single animal in one). RNA was
613 converted to cDNA and library preparation was performed using the NEXTflex™ Directional
614 RNA-Seq Kit (Illumina) and paired-end reads were sequenced on an Illumina HiSeq 2500
615 system. Adapters and poor quality bases (<30) were trimmed using fastq-mcf from the ea-
616 utilities package, and reads were aligned to assemblies using Tophat2 (v2.0.14).
617

618 **Chip-Seq**
619 Three adult male zebra finches were treated as above, the brains dissected, and the RA and
620 surrounding arcopallium of each bird was then processed individually using the native ChIP
621 protocol described in [54] with an H3K27ac antibody (Ab#4729). The DNA libraries were
622 prepared using the MicroPlex Library Preparation Kit v2 (C05010012). 50 bp single-end
623 sequencing was done on the Illumina HiSeq 4000 system. The reads were aligned to the
624 assemblies using Bowtie2 (Bowtie , RRID:SCR_005476)(v2.2.9, [28]). More detail will be
625 provided in a later publication focusing on vocal learning brain regions.
626

627 **Comparative analyses between assemblies for individual genes**
628 The Sanger-based reference zebra finch assembly in the UCSC browser and the Ilumina-based
629 reference Anna's Hummingbird in Avianbase (http://avianbase.narf.ac.uk/index.html), and both
630 in NCBI where used for comparing with the Pacbio assembly. In the UCSC browser, there are
631 two annotations, one from 2008 (http://genome.ucsc.edu/cgi-bin/hgGateway?db=taeGut1) and
632 the other from 2013 (http://genome.ucsc.edu/cgi-bin/hgGateway?db=taeGut2), with some
633 differences between them. Our findings were similar, although not always identical, with both
634 annotations, with errors being present in both annotations based on the Pacbio assembly. The
635 nucleotide quality score tract was only available in the 2008 browser.
636     Multiple species sequence alignments were done with BioEdit v7.2.5
637 (http://www.mbio.ncsu.edu/bioedit/bioedit.html) [55]; Dotplots of alignments were generated
638 with Gepard v1.4 (http://cube.univie.ac.at/gepard) [56]; Alignments of raw SMRT genome reads
639 to the assembled genomes were done with Blasr, which is part of SMRTLink software from

640 Pacbio; Iso-Seq reads were aligned with GMAP (http://research-pub.gene.com/gmap/, version
641 2016-08-16) [57].

## Availability of data

This Whole Genome Shotgun project has been deposited at DDBJ/ENA/GenBank under BioProject PRJNA368994. The zebra finch accession number is MUGN00000000 and SRA for raw reads is SRS1954332. The Anna's Hummingbird accession number is MUGM00000000 and SRA is SRP061272. The NCBI accessions also contain translation tables of the PacBio contig designations and their corresponding NCBI accession labels, for all contigs (ftp://ftp.ncbi.nlm.nih.gov/genomes/all/GCA/002/008/985/GCA_002008985.2_Tgut_diploid_1.0 /GCA_002008985.2_Tgut_diploid_1.0_assembly_report.txt and ftp://ftp.ncbi.nlm.nih.gov/genomes/all/GCA/002/021/895/GCA_002021895.1_Canna_diploid_1. 0/GCA_002021895.1_Canna_diploid_1.0_assembly_report.txt, respectively; the first column contains the PacBio assembly contig ID, the 5th column designates the corresponding NCBI contig accession number). We have also included these tables here as additional files (Tables S1 & S2). Supporting assemblies, BUSCO & CEGMA output files, and RNASeq and ChipSeq data are also available from the *GigaScience* GigaDB repository[58].

## Competing interest

Jonas Korlach, Sarah Kingan, Chen-Shan Chin are full-time employees at Pacific Biosciences, a company developing single-molecule sequencing technologies.

## Author contributions

J.K. and E.D.J. designed the project and wrote the manuscript; C.S.C. and S.K. carried out genome assemblies; J.K., G.G. and S.K. conducted analyses on single genes as well as CEGMA and BUSCO analyses; G.G. and J-N.A. conducted RNA-Seq experiments, L.C. conducted Chip-Seq experiments; J.H. processed samples; and all authors contributed to writing and editing the manuscript.

## Additional Files

Supporting data is included in supplementary figures S1-S15, and supplementary tables S1-S2.

## Figure legends

**Figure 1.** Gene completeness within assemblies. *(A)* Comparison to a 248 highly conserved core CEGMA eukaryote gene set using human genes [21], between the Sanger-based zebra finch and Illumina-based Anna's hummingbird references and their respective PacBio-based assemblies. We used a more stringent cut-off ($> 95\%$) for completeness than usually done ($> 90\%$), because we felt 90% was too permissive, as it could allow entire missing exons and still call a gene as complete. Gene count is the percentage of genes in each of the assemblies that met this criterion. *(B)* Comparison to a 303 single-copy conserved eukaryotic BUSCO gene set [23]. Complete is $\geq$ 95% complete; fragmented is $< 95\%$ complete; missing is not found. *(C)* Comparison to 4915 single-copy conserved genes from the avian BUSCO gene [23].

**Figure 2.** Transcriptome and regulome representation within assemblies. *(A)* Percentage of RNA-Seq and H3K27Ac ChIP-Seq reads from the zebra finch RA song nucleus mapped back to the zebra finch Sanger-based and PacBio-based genome assemblies. *(B)* Pie charts of the distributions of the RNA-Seq reads mapped to the zebra finch genome assemblies. *(C)* Pie charts of the distribution of ChIP-Seq reads mapped to the zebra finch genome assemblies. * $p < 0.05$; ** $p < 0.002$; *** $p < 0.0001$; paired t-test within animals between assemblies; n = 5 RNA-Seq and n = 3 ChIP-Seq independent replicates from different animals.

**Figure 3.** Comparison of *EGR1* assemblies. *(A)* UCSC Genome browser view of the Sanger-based zebra finch *EGR1* assembly, highlighting (from top to bottom) four contigs (light and dark brown) with three gaps, GC percent, nucleotide quality score (blue), RefSeq gene prediction (purple), and areas of repeat sequences. *(B)* Summary comparison of the Sanger-based and PacBio-based zebra finch assemblies, showing in the latter filling the gaps (black) and correcting erroneous reference sequences surrounding the gaps (red). Tick mark is a synonymous heterozygous SNP in the coding region between the primary (1) and secondary (2) haplotypes. Panels *A* and *B* are of the same scale. *(C)* Comparison of the hummingbird Illumina- and PacBio-based assemblies, showing similar corrections that further lead to a correction in the protein coding sequence prediction (blue). *(D)* Multiple sequence alignment of the EGR1 protein for the four assemblies (two zebra finch and two hummingbird) in *B* and *C*, showing corrections to the Illumina-based hummingbird protein prediction by the PacBio-based assembly.

**Figure 4.** Comparison of *DUSP1* assemblies. *(A)* UCSC Genome browser view of the Sanger-based zebra finch *DUSP1* assembly, highlighting four contigs with three gaps, GC percent, nucleotide quality score, Blat alignment of the NCBI gene prediction (XP_002193168.1, blue), and repeat sequences. *(B)* Resolution of the region by the PacBio-based zebra finch assembly, filling the gaps (black) and correcting erroneous reference sequences in repeat regions (red) and gene predictions (blue). Panels *A* and *B* are of the same scale. *(C)* Resolution and correction to the hummingbird Illumina-based assembly with the PacBio-based assembly (same color scheme

as in *B*). *(D)* Multiple sequence alignment of the DUSP1 protein for the four assemblies in *B* and *C*, showing numerous corrections to the Sanger-based and Ilumina-based protein predictions by both PacBio-based assemblies.

**Figure 5.** Comparison of *FOXP2* assemblies. *(A)* UCSC Genome browser view of the Sanger-based zebra finch *FOXP2* assembly, highlighting 10 contigs with 9 gaps, GC percent, nucleotide quality score, RefSeq gene prediction, and repeat sequences. *(B)* Table showing the number of resolved and corrected erroneous base pairs in the gaps by the PacBio-based primary and secondary haplotype assemblies; * indicates differences between haplotypes. *(C)* Dot plot of the Sanger-based reference (x-axis) and the PacBio-based primary assembly (y-axis) corresponding to the three GC-rich region gaps immediately upstream and surrounding the first exon of the *FOXP2* gene. *(D)* Schematic summary of corrections to the three gaps shown in *C*, in the two haplotypes of the PacBio-based assembly. The protein coding sequence alignments are in Figure S13A.

**Figure 6.** Comparison of *SLIT1* assemblies. *(A)* UCSC Genome browser view of the Sanger-based zebra finch *SLIT1* assembly, highlighting 15 contigs with 14 gaps, GC percent, nucleotide quality score, NCBI *SLIT1* gene prediction (XP_012430014.1, blue), and repeat sequences. Red circles, gaps that correspond to the missing exon 1 and part of the missing exon 35, respectively. *(B)* Multiple sequence alignment comparison of the SLIT1 protein for the four assemblies compared, including the two different haplotypes from the PacBio-based zebra finch assembly (rows 2 and 3).

**Supplementary Figure S1.** DNA isolation, library construction, and size selection. *(A)* Pulsed-field gel showing original size of starting genomic DNA (lane 3), the sheared DNA (1), and the size selected library (2). *(B)* Bioanalyzer trace before (blue) and after (red) library size selection for fragments > 17 kb.

**Supplementary Figure S2.** Read and insert length distributions. *(A, B)* Sequence read length distributions from SMRT cell sequencing for both species. *(C, D)* Sequenced DNA insert length distributions from SMRT cell sequencing for both species.

**Supplementary Figure S3.** Box plots comparing protein coding sequence lengths of orthologous proteins between the CEGMA and BUSCO eukaryotic and avian datasets. ** $p < 0.001$; *** $p < 0.0001$, one-sided paired Wilcoxon signed-rank test, prediction of the proteins being longer in CEGMA datasets.

**Supplementary Figure S4.** Vocal learning and adjacent brain regions in songbirds used for RNA-Seq and ChIP-Seq analyses, and comparison with humans. *(A)* Drawing of a zebra finch male brain section showing specialized vocal learning pathway and associated profiled song

771 nuclei RA, HVC, LMAN, and Area X. *(B)* Drawing of a human brain section showing spoken-
772 language pathway and analogous brain regions. Black arrows, posterior vocal motor pathway;
773 White arrows, anterior vocal learning pathway; Dashed arrows, connections between the two
774 pathways; Red arrow, specialized direct projection from forebrain to brainstem vocal motor
775 neurons in vocal learners. Italicized letters adjacent to the song and speech regions indicates
776 regions (in songbirds) that show mainly show motor *(m)*, auditory *(a)*, equally both motor and
777 auditory *(m/a)* neural activity or activity-dependent gene expression. Figure from [59] and [4].

779 *Abbreviations*: A1-L4, primary auditory cortex – layer 4; Am, nucleus ambiguous; Area X, a
780 vocal nucleus in the striatum; aSt, anterior striatum vocal region; aT, anterior thalamus speech
781 area; Av, avalanche; aDLM, anterior dorsolateral nucleus of the thalamus; DM, dorsal medial
782 nucleus of the midbrain; HVC, a vocal nucleus (no abbreviation); L2, auditory area similar to
783 human cortex layer 4; LSC, laryngeal somatosensory cortex; LMC, laryngeal motor cortex;
784 MAN, magnocellular nucleus of the anterior nidopallium; MO, oval nucleus of the anterior
785 mesopallium; NIf, interfacial nucleus of the nidopallium; PAG, peri-aqueductal gray; RA, robust
786 nucleus of the arcopallium; v, ventricle space

788 **Supplementary Figure S5.** Dot plot of sequence comparisons for genome assemblies of the
789 *EGR1* region. *(A)* Comparison of zebra finch PacBio-based versus Sanger-based assemblies for
790 the region containing *EGR1*, showing the GC-rich promoter region and closing and corrections
791 of gaps for the PacBio-based assembly. *(B)* Comparison of hummingbird Illumina-based versus
792 PacBio-based assemblies for the region containing *EGR1*, showing an erroneous tandem
793 duplication in the Ilumina-based assembly and closing of gaps for the PacBio-based assembly.

795 **Supplementary Figure S6.** Single SMRT genomic reads and Iso-Seq mRNA reads supporting
796 Pacbio *EGR1* assembly. *(A)* Zebra finch PacBio SMRT reads (rows) mapped against the zebra
797 finch PacBio assembly (contig 405, entire *EGR1* region, same as Fig. 3A). Reads are shaded by
798 length (>10 kb reads = black). *(B)* Example of a single Ruby-throated hummingbird Iso-Seq read
799 mapped against Illumina-based (top) and PacBio-based (bottom) Anna's hummingbird genome
800 assemblies using GMAP. Note the first exon (blue) which is present in the Iso-Seq read is
801 missing in the Illumina-based assembly, but present in the PacBio-based assembly.

803 **Supplementary Figure S7.** Dot plot of sequence comparison for the PacBio-based hummingbird
804 and zebra finch *EGR1* region assemblies. Note regions of high species conservation and
805 divergence surrounding *EGR1*. Blue box, location of the *EGR1* exons and intron.

807 **Supplementary Figure S8.** Dot plot comparisons for *DUSP1* region assemblies. *(A)*
808 Comparison of the Sanger-based and PacBio-based zebra finch *DUSP1* region assemblies,
809 showing problems in the Sanger-based assembly with microsatellite repeats. *(B)* Comparison of
810 the Illumina-based and PacBio-based hummingbird *DUSP1* region assemblies, showing a large

811 gap including the microsatellite region and the beginning of the gene, and an erroneous tandem
812 duplication in the Illumina-based assembly.

813

814 **Supplementary Figure S9.** Pacbio correction of base call errors found in Sanger reference *(A)*
815 Confirmation of the PacBio sequence in the three locations different from the zebra finch Sanger
816 reference by alignments to DUSP1 sequences of other songbirds. *(B)* PacBio reads (rows)
817 corresponding to the genomic region in DUSP1 that differs in the three locations from the zebra
818 finch Sanger reference, resulting in a.a. changes. The codons in question are highlighted.

819

820 **Supplementary Figure S10.** Dot plot comparison of assemblies for the *DUSP1* microsatellite
821 region. *(A)* Differences in the microsatellite region upstream of the *DUSP1* protein coding
822 sequence between the primary and the secondary haplotypes in the fully assembled zebra finch
823 PacBio-based assembly. *(B)* Differences in microsatellites region upstream of *DUSP1* between
824 the zebra finch and hummingbird in the fully assembled PacBio-based assemblies.

825

826 **Supplementary Figure S11.** Dot plot comparisons for PacBio-based *DUSP1* region assemblies
827 with orthogonal validation. Comparison of the PacBio-based genome assembly and Sanger-
828 based single clone of the *(A)* zebra finch and *(B)* hummingbird *DUSP1* upstream region
829 assemblies showing more consistency between the two (than in **Fig S8A**). Not visible in this
830 high-level alignment view is an 11-bp deletion and several SNPs in this allele of the PacBio
831 contig relative to the other allele; the single clone of the individual is more consistent with the
832 alternate allele without the 11-bp deletion.

833

834 **Supplementary Figure S12.** Single Iso-Seq mRNA reads supporting Pacbio assemblies. *(A)*
835 Full-length PacBio mRNA sequence Iso-Seq ruby throated hummingbird reads for DUSP1
836 aligned against the exons of the corresponding primary contigs from Anna's hummingbird
837 Illumina (top panel) and PacBio (bottom panel) assemblies. *(B)* Similar alignments for FOXP2
838 IsoSeq reads.

839

840 **Supplementary Figure S13.** Dot plot comparison of assemblies for the *FOXP2* region. *(A)*
841 zebra finch, *(B)* hummingbird.

842

843 **Supplementary Figure S14.** *(A)* Multiple sequence alignment of the FOXP2 protein for the four
844 assemblies (two zebra finch and two hummingbird) compared in this study, showing correction
845 of a nucleotide error in the Sanger-based zebra finch assembly, and correction of an erroneous
846 stop codon (x) in the Illumina-based hummingbird assembly. Note an extra 18 a.a. stretch in the
847 hummingbird sequence validated by gene prediction of both assemblies, that was not present in
848 the zebra finch. *(B)* Missing 88bp of sequence in exon 6 of Illumina-based assembly. *(C)*
849 Resolution of exon 6 in Pacbio-based assembly, also revealing a SNP.

850

851 **Supplementary Figure S15.** Large regional correction made by the PacBio diploid assembly.
852 *(A)* Correction of an erroneous stretch of 462 bp in the first intron of *FOXP2* in the hummingbird
853 Illumina assembly by the PacBio assembly. *(B)* Dot plot of haplotype variation in the *FOXP2*
854 gene revealed by the PacBio diploid assembly: a 708 bp deletion in the secondary haplotype
855 contig relative to the primary contig.
856

857 **Supplementary Figure S16.** Dot plot comparison of assemblies for the *SLIT1* region. *(A)* zebra
858 finch, *(B)* hummingbird.

859

## References

860

1. Hillier LW, Miller W, Birney E, Warren W, Hardison RC, Ponting CP, Bork P, Burt DW, Groenen MA, Delany ME, et al: **Sequence and comparative analysis of the chicken genome provide unique perspectives on vertebrate evolution.** *Nature* 2004, **432:**695-716.

2. Warren WC, Clayton DF, Ellegren H, Arnold AP, Hillier LW, Kunstner A, Searle S, White S, Vilella AJ, Fairley S, et al: **The genome of a songbird.** *Nature* 2010, **464:**757-762.

3. Shi Z, Luo G, Fu L, Fang Z, Wang X, Li X: **miR-9 and miR-140-5p target FoxP2 and are regulated as a function of the social context of singing behavior in zebra finches.** *J Neurosci* 2013, **33:**16510-16521.

4. Pfenning AR, Hara E, Whitney O, Rivas MV, Wang R, Roulhac PL, Howard JT, Wirthlin M, Lovell PV, Ganapathy G, et al: **Convergent transcriptional specializations in the brains of humans and song-learning birds.** *Science* 2014, **346:**1256846.

5. Koepfli K-P, Paten B, O'Brien SJ, Scientists GKC, Lewin H, Roberts R: **The Genome 10K Project: A Way Forward.** In *Annual Review of Animal Biosciences, Vol 3. Volume 3*; 2015: 57-111

6. Zhang GJ, Jarvis ED, Gilbert MTP: **A flock of Genomes.** *Science* 2014, **346:**1308-1309.

7. Green RE, Braun EL, Armstrong J, Earl D, Nguyen N, Hickey G, Vandewege MW, St John JA, Capella-Gutiérrez S, Castoe TA, et al: **Three crocodilian genomes reveal ancestral patterns of evolution among archosaurs.** *Science* 2014, **346:**1254449.

8. Zhang GJ, Li C, Li QY, Li B, Larkin DM, Lee C, Storz JF, Antunes A, Greenwold MJ, Meredith RW, et al: **Comparative genomics reveals insights into avian genome evolution and adaptation.** *Science* 2014, **346:**1311-1320.

9. Jarvis ED, Mirarab S, Aberer AJ, Li B, Houde P, Li C, Ho SY, Faircloth BC, Nabholz B, Howard JT, et al: **Whole-genome analyses resolve early branches in the tree of life of modern birds.** *Science* 2014, **346:**1320-1331.

10. Joseph L, Buchanan KL: **A quantum leap in avian biology.** *Emu* 2015, **115:**1-5.

11. Kraus RHS, Wink M: **Avian genomics: fledging into the wild!** *Journal of Ornithology* 2015, **156:**851-865.

12. Haug-Baltzell A, Jarvis ED, McCarthy FM, Lyons E: **Identification of dopamine receptors across the extant avian family tree and analysis with other clades uncovers a polyploid expansion among vertebrates.** *Frontiers in Neuroscience* 2015, **9**.

13. Horita H, Kobayashi M, Liu WC, Oka K, Jarvis ED, Wada K: **Specialized Motor-Driven dusp1 Expression in the Song Systems of Multiple Lineages of Vocal Learning Birds.** *PLoS ONE* 2012, **7:**e42173.

14. Roberts RJ, Carneiro MO, Schatz MC: **The advantages of SMRT sequencing.** *Genome Biol* 2013, **14:**405.

15. Bradnam KR, Fass JN, Alexandrov A, Baranay P, Bechner M, Birol I, Boisvert S, Chapman JA, Chapuis G, Chikhi R, et al: **Assemblathon 2: evaluating de novo methods of genome assembly in three vertebrate species.** *Gigascience* 2013, **2:**10.

16. Gordon D, Huddleston J, Chaisson MJ, Hill CM, Kronenberg ZN, Munson KM, Malig M, Raja A, Fiddes I, Hillier LW, et al: **Long-read sequence assembly of the gorilla genome.** *Science* 2016, **352:**aae0344.

17. **FALCON assembler** [https://github.com/PacificBiosciences/FALCON/commit/a1180264c3c7d2de1c5eb55b3663dce093354dd7]

18. Chin CS, Peluso P, Sedlazeck FJ, Nattestad M, Concepcion GT, Clum A, Dunn C, O'Malley R, Figueroa-Balderas R, Morales-Cruz A, et al: **Phased diploid genome**

25

910 **assembly with single-molecule real-time sequencing.** *Nat Methods* 2016, **13:**1050-
911 1054.

912 19. Gregory TR: **Animal Genome Size Database.** 2017.

913 20. Parra G, Bradnam K, Korf I: **CEGMA: a pipeline to accurately annotate core genes in**
914 **eukaryotic genomes.** *Bioinformatics* 2007, **23:**1061-1067.

915 21. Parra G, Bradnam K, Ning Z, Keane T, Korf I: **Assessing the gene space in draft**
916 **genomes.** *Nucleic Acids Res* 2009, **37:**289-297.

917 22. Zdobnov EM, Tegenfeldt F, Kuznetsov D, Waterhouse RM, Simão FA, Ioannidis P,
918 Seppey M, Loetscher A, Kriventseva EV: **OrthoDB v9.1: cataloging evolutionary and**
919 **functional annotations for animal, fungal, plant, archaeal, bacterial and viral**
920 **orthologs.** *Nucleic Acids Res* 2017, **45:**D744-D749.

921 23. Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM: **BUSCO:**
922 **assessing genome assembly and annotation completeness with single-copy**
923 **orthologs.** *Bioinformatics* 2015, **31:**3210-3212.

924 24. Zhang G, Li B, Li C, Gilbert MTP, Mello CV, Jarvis ED, Wang J, The Avian Genome C:
925 **Genomic data of the Anna's Hummingbird (Calypte anna).** *GigaDB* 2014.

926 25. Kim D, Pertea G, Trapnell C, Pimentel H, Kelley R, Salzberg SL: **TopHat2: accurate**
927 **alignment of transcriptomes in the presence of insertions, deletions and gene**
928 **fusions.** *Genome Biol* 2013, **14:**R36.

929 26. Jarvis ED, Yu J, Rivas MV, Horita H, Feenders G, Whitney O, Jarvis SC, Jarvis ER,
930 Kubikova L, Puck AEP, et al: **Global View of the Functional Molecular Organization**
931 **of the Avian Cerebrum: Mirror Images and Functional Columns.** *Journal of*
932 *Comparative Neurology* 2013, **521:**3614-3665.

933 27. Whitney O, Pfenning AR, Howard JT, Blatti CA, Liu F, Ward JM, Wang R, Audet JN,
934 Kellis M, Mukherjee S, et al: **Core and region-enriched networks of behaviorally**
935 **regulated genes and the singing genome.** *Science* 2014, **346:**1256780.

936 28. Langmead B, Salzberg SL: **Fast gapped-read alignment with Bowtie 2.** *Nat Methods*
937 2012, **9:**357-359.

938 29. Shlyueva D, Stampfel G, Stark A: **Transcriptional enhancers: from properties to**
939 **genome-wide predictions.** *Nat Rev Genet* 2014, **15:**272-286.

940 30. Veyrac A, Besnard A, Caboche J, Davis S, Laroche S: **The transcription factor**
941 **Zif268/Egr1, brain plasticity, and memory.** *Prog Mol Biol Transl Sci* 2014, **122:**89-129.

942 31. Jarvis ED, Nottebohm F: **Motor-driven gene expression.** *Proc Natl Acad Sci U S A*
943 1997, **94:**4097-4102.

944 32. Flavell SW, Greenberg ME: **Signaling mechanisms linking neuronal activity to gene**
945 **expression and plasticity of the nervous system.** *Annu Rev Neurosci* 2008, **31:**563-
946 590.

947 33. Cortés-Mendoza J, Díaz de León-Guerrero S, Pedraza-Alva G, Pérez-Martínez L:
948 **Shaping synaptic plasticity: the role of activity-mediated epigenetic regulation on**
949 **gene transcription.** *Int J Dev Neurosci* 2013, **31:**359-369.

950 34. Stanke M, Diekhans M, Baertsch R, Haussler D: **Using native and syntenically**
951 **mapped cDNA alignments to improve de novo gene finding.** *Bioinformatics* 2008,
952 **24:**637-644.

953 35. Workman RE, Myrka AM, Tseng E, Wong GW, Welch KC, Timp W: **Single molecule,**
954 **full-length transcript sequencing provides insight into the extreme metabolism of**
955 **ruby-throated hummingbird Archilochus colubris.** *bioRxiv* 2017.
956 https://doi.org/10.1101/117218

957 36. Liu YX, Wang J, Guo J, Wu J, Lieberman HB, Yin Y: **DUSP1 is controlled by p53**
958 **during the cellular response to oxidative stress.** *Mol Cancer Res* 2008, **6:**624-633.

959 37. Horita H, Wada K, Rivas MV, Hara E, Jarvis ED: **The dusp1 immediate early gene is regulated by natural stimuli predominantly in sensory input neurons.** *J Comp Neurol* 2010, **518:**2873-2901.

962 38. Fisher SE, Scharff C: **FOXP2 as a molecular window into speech and language.** *Trends Genet* 2009, **25:**166-177.

964 39. Lai CS, Fisher SE, Hurst JA, Vargha-Khadem F, Monaco AP: **A forkhead-domain gene is mutated in a severe speech and language disorder.** *Nature* 2001, **413:**519-523.

966 40. Turner SJ, Hildebrand MS, Block S, Damiano J, Fahey M, Reilly S, Bahlo M, Scheffer IE, Morgan AT: **Small intragenic deletion in FOXP2 associated with childhood apraxia of speech and dysarthria.** *Am J Med Genet A* 2013, **161A:**2321-2326.

969 41. Haesler S, Wada K, Nshdejan A, Morrisey EE, Lints T, Jarvis ED, Scharff C: **FoxP2 expression in avian vocal learners and non-learners.** *J Neurosci* 2004, **24:**3164-3175.

972 42. Teramitsu I, White SA: **FoxP2 regulation during undirected singing in adult songbirds.** *J Neurosci* 2006, **26:**7390-7394.

974 43. Haesler S, Rochefort C, Georgi B, Licznerski P, Osten P, Scharff C: **Incomplete and inaccurate vocal imitation after knockdown of FoxP2 in songbird basal ganglia nucleus Area X.** *PLoS Biol* 2007, **5:**e321.

977 44. Castellucci GA, McGinley MJ, McCormick DA: **Knockout of Foxp2 disrupts vocal development in mice.** *Sci Rep* 2016, **6:**23305.

979 45. Chabout J, Sarkar A, Patel SR, Radden T, Dunson DB, Fisher SE, Jarvis ED: **A Foxp2 Mutation Implicated in Human Speech Deficits Alters Sequencing of Ultrasonic Vocalizations in Adult Male Mice.** *Front Behav Neurosci* 2016, **10:**197.

982 46. Condro MC, White SA: **Recent Advances in the Genetics of Vocal Learning.** *Comp Cogn Behav Rev* 2014, **9:**75-98.

984 47. Blockus H, Chédotal A: **The multifaceted roles of Slits and Robos in cortical circuits: from proliferation to axon guidance and neurological diseases.** *Curr Opin Neurobiol* 2014, **27:**82-88.

987 48. Wang R, Chen CC, Hara E, Rivas MV, Roulhac PL, Howard JT, Chakraborty M, Audet JN, Jarvis ED: **Convergent differential regulation of SLIT-ROBO axon guidance genes in the brains of vocal learners.** *J Comp Neurol* 2014.

990 49. Treangen TJ, Salzberg SL: **Repetitive DNA and next-generation sequencing: computational challenges and solutions.** *Nat Rev Genet* 2011, **13:**36-46.

992 50. Palazzo AF, Gregory TR: **The case for junk DNA.** *PLoS Genet* 2014, **10:**e1004351.

993 51. Chin CS, Alexander DH, Marks P, Klammer AA, Drake J, Heiner C, Clum A, Copeland A, Huddleston J, Eichler EE, et al: **Nonhybrid, finished microbial genome assemblies from long-read SMRT sequencing data.** *Nat Methods* 2013, **10:**563-569.

996 52. She R, Chu JS, Wang K, Pei J, Chen N: **GenBlastA: enabling BLAST to identify homologous gene sequences.** *Genome Res* 2009, **19:**143-149.

998 53. Slater GS, Birney E: **Automated generation of heuristics for biological sequence comparison.** *BMC Bioinformatics* 2005, **6:**31.

1000 54. Brind'Amour J, Liu S, Hudson M, Chen C, Karimi MM, Lorincz MC: **An ultra-low-input native ChIP-seq protocol for genome-wide profiling of rare cell populations.** *Nat Commun* 2015, **6:**6033.

1003 55. Hall TA: **BioEdit: a user-friendly biological sequence alignment editor and analysis program for Windows 95/98/NT.** *Nucl Acids Symp Ser* 1999, **41:**95-98.

1005 56. Krumsiek J, Arnold R, Rattei T: **Gepard: a rapid and sensitive tool for creating dotplots on genome scale.** *Bioinformatics* 2007, **23:**1026-1028.

1007 57. Wu TD, Watanabe CK: **GMAP: a genomic mapping and alignment program for mRNA and EST sequences.** *Bioinformatics* 2005, **21:**1859-1875.
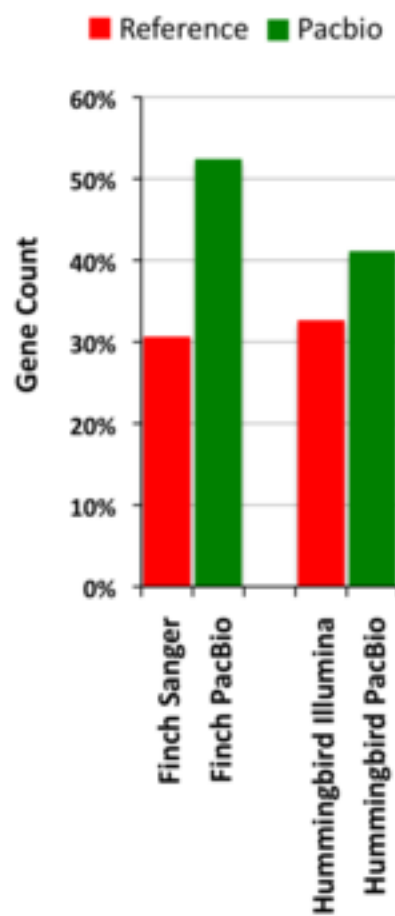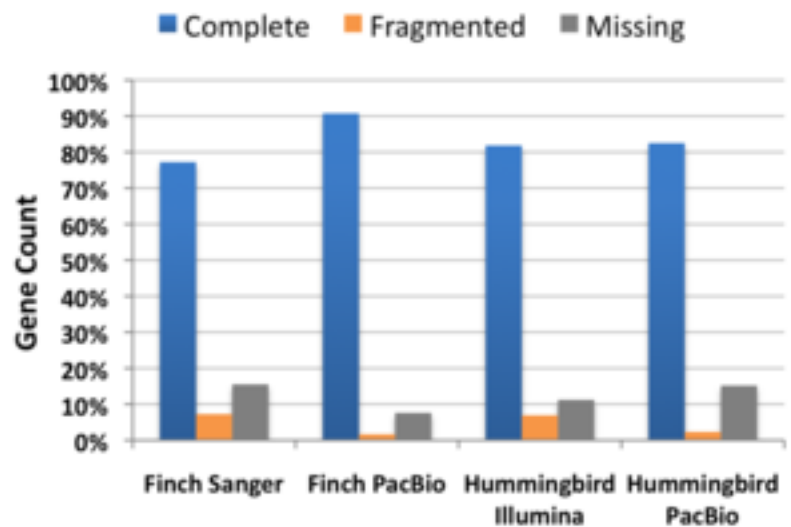
58. Korlach, J; Gedman, G; Kingan, S, B; Chin, C; Howard, J, T; Audet, J; Cantin, L; Jarvis, E, D (2017): **Supporting data for "De Novo PacBio long-read and phased avian genome assemblies correct and add to reference genes generated with intermediate and short reads"** GigaScience Database. http://dx.doi.org/10.5524/100311

59. Chakraborty M, Jarvis ED: **Brain evolution by brain pathway duplication.** *Philosophical Transactions of the Royal Society B-Biological Sciences* 2015, **370:**50056-50056.

Figure 1

Figure 1

Figure 2

Click here to download Figure Korlach Fig. 2 v9.png ⬇



## A. Aligned RNA-Seq reads

## B. Aligned RNA-Seq read distribution

## C. Aligned Chip-Seq read distribution
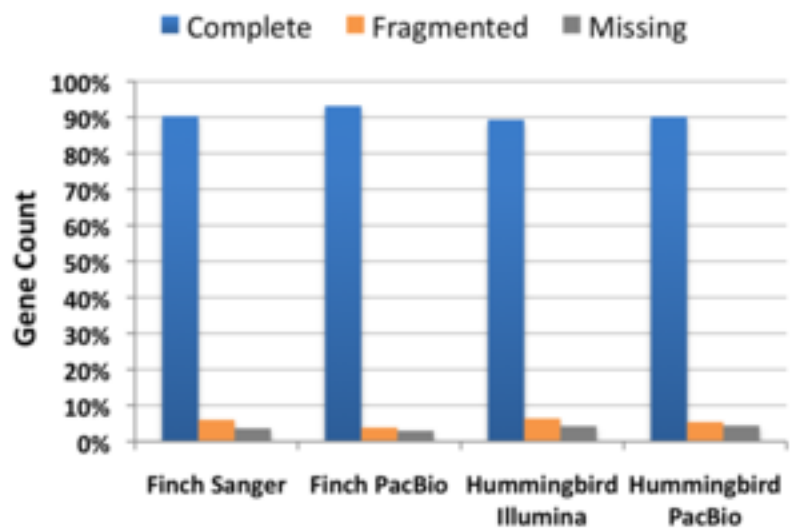
**Figure 2**

Figure 3

**Figure 3**

Figure 4

Figure 4

Figure 5                                        Click here to download Figure Korlach Fig. 5 v9.png  ⬇

**A.**



**B.**

| Gap # Sanger assembly | Gap and discordant length in Sanger assembly | Corrected length | | |
|---|---|---|---|---|
| | | Primary PacBio haplotype assembly | Secondary PacBio haplotype assembly | |
| 1 | 100N + 269 bp | 237 bp | 237 bp | |
| 2 | 100N + 541 bp | 320 bp* | 262 bp* | |
| 3 | 100N | 281 bp | 281 bp | |
| 4 | 100N | 345 bp* | 354 bp* | |
| 5 | 100N + 3 bp | 55 bp | 55 bp | |
| 6 | 100N + 835 bp | 241 bp | 241 bp | |
| 7 | 100N + 332 bp | 38 bp* | 36 bp* | |
| 8 | 100N + 492 bp | 426 bp* | 416 bp* | |
| 9 | 100N + 279 bp | 566 bp | 566 bp | |
| Total | 900N + 2751 bp | 2509 bp | 2448 bp | |

\* Differences between haplotypes

**C.**



Gaps (100 Ns) in Sanger reference, GC-rich

**D.**

Zebra finch Sanger-based reference



Zebra finch PacBio-based assembly (contigs 5 and 5_002)

# Figure 5

Figure 6

**Figure 6**

Supplementary figures

Click here to access/download
**Supplementary Material**
Korlach et al supplementary figures v15.pdf

**The Rockefeller University**
SCIENCE FOR THE BENEFIT OF HUMANITY

May 16th, 2017

To: Scott Edmunds, Editor, GigaScience

Dear Scott

Please find uploaded our revised paper "***De Novo* PacBio long-read and phased avian genome assemblies correct and add to reference genes generated with intermediate and short reads"**. We are glad that you and the reviewers were positive about the manuscript, with recommend minor revisions. The reviewer's comments were very useful, and we were able to address all of them. This meant mostly revisions in manuscript writing style, additional links to sources of the original data, and one additional supplementary figure (Figure S11) with additional orthogonal data requested by reviewer #2 for validating the improvements in the Pacbio-based long read assemblies over the Sanger-based intermediate-read and Illumina-based short read assemblies. We think the manuscript is improved as a result. We submitted the text with tracked changes as a supplementary file to help show where the changes have been made.

We hope that this manuscript is suitable for publication in GigaScience.

Sincerely,

Erich D. Jarvis, Ph.D.
Investigator, Howard Hughes Medical Institute
Professor, The Rockefeller University, Box 54
1230 York Avenue, New York, New York 10065
http://www.jarvislab.net/
http://www.jarvislab.net/Publications.html