## *De Novo* PacBio long-read and phased avian genome assemblies correct and add to reference genes generated with intermediate and short reads

Jonas Korlach[1*], Gregory Gedman[2], Sarah B. Kingan[1], Chen-Shan Chin[1], Jason T. Howard[2], Jean-Nicolas Audet[2,3], Lindsey Cantin[2], and Erich D. Jarvis[2,4*]

[1]Pacific Biosciences, Menlo Park, CA 94025, USA; [2]Laboratory of Neurogenetics of Language, The Rockefeller University, New York, NY 10065, USA; [3]Department of Biology, McGill University, Montreal, Quebec H3A 1B1, Canada; [4]Howard Hughes Medical Institute, Chevy Chase, MD 20815, USA

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

Jonas Korlach, Ph.D.
Chief Scientific Officer
Pacific Biosciences
1305 O'Brien Drive
Menlo Park, CA 94025
650-521-8006

Erich D. Jarvis, Ph.D.
Investigator, Howard Hughes Medical Institute
Professor, The Rockefeller University, Box 54
1230 York Avenue, New York, New York 10065
212 327-8806

## 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 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 corrections 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 and phasing 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 costly Sanger method to sequence genomes of the two most commonly studied bird species, the chicken [1] and zebra finch [2], have impacted many studies. The zebra finch is a vocal learning songbird, with the rare ability to imitate sounds 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 genomes were sequenced, 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 [5]. The consortium was successful in conducting comparative genomics and phylogenetics with populations of genes [6-9]. However, when it was necessary to dig deeper into individual genes, it was discovered that many were incompletely assembled or

2

contained apparent misassemblies. For example, the *DRD4* dopamine receptor was missing in half of the assemblies, in part due to sequence complexity [10]. The *EGR1* immediate early gene transcription factor, a commonly studied gene in neuroscience and in vocal learning species, was missing the promoter region in an GC-rich region in every bird genome we examined. Another immediate early gene, *DUSP1*, with specialized vocalizing-driven gene expression in song nuclei of vocal learning species, has microsatellite sequences in the promoters of vocal learning species that are missing or misassembled, requiring single-molecule cloning and sequencing to resolve [11]. Such errors create a great amount of effort to clone, sequence, and correct assemblies of individual genes of interest.

High-throughput, single-molecule, long-read sequencing shows promise to alleviate these problems [12-14]. Here, we applied PacBio single-molecule long-read (1,000-60,000 bp) sequencing and diploid assembly on two vocal learning species, the zebra finch previously assembled with Sanger-based intermediate reads (700-1,000 bp), and the Anna's hummingbird previously assembled with Illumina-based short reads (100-150 bp). We found that the long-read diploid assemblies resulted in major improvements in genome completeness and contiguity, and completely resolved the problems in all of our genes of interest. This study is part of an effort to help evaluate standards for the G10K vertebrate (https://genome10k.soe.ucsc.edu) and the B10K bird (http://b10k.genomics.cn/index.html) genome consortiums.


## Results

### The long-read assemblies result in 150-fold to 200-fold increases in contiguity

To generate long-read assemblies, high molecular weight DNA was isolated from muscle tissue of the same zebra finch male and Anna's hummingbird female used to create the current reference genomes [2, 6]. The DNA was sheared, 35-40 kb libraries generated, size-selected for inserts >17 kb (**Fig. S1**), and then SMRT sequencing performed on the PacBio RS II instrument to obtain ~96X coverage for the zebra finch (19 kb N50 read length) and ~70X for the hummingbird (22 kb N50 read length; **Fig. S2**). The long reads were originally assembled into a merged haplotype with an early version of the FALCON assembler [15], which we found unintentionally introduced indels for some nucleotides that differed between haplotypes (tested on the hummingbird; data not shown). We then re-assembled using FALCON v0.4.0 followed by the FALCON-Unzip module [16] to prevent indel formation and generate long-range phased haplotypes. Thus, the new assemblies, unlike the current reference assemblies, are phased diploids. This PacBio-based sequencing and assembly approach does not link contigs into gapped scaffolds. Scaffolding requires additional approaches, which we will report on separately in a study comparing scaffolding technologies with these assemblies. The results presented here were found independent of scaffolding.

For the zebra finch, our long-read approach resulted in 1159 primary haplotype contigs with an estimated total genome size of 1.14 Gb (1.2 Gb expected; [17]) and contig N50 of 5.81

109 Mb, representing a 108-fold reduction in the number of contigs and a 150-fold improvement in
110 contiguity compared to the current Sanger-based reference (**Table 1A**). The diploid assembly
111 process produced 2188 associated, or secondary, haplotype contigs (i.e. haplotigs) with an
112 estimated length of 0.84 Gb (**Table 1A**), implying that about 75% of the genome contained
113 sufficient heterozygosity to be phased into haplotypes by FALCON-Unzip. Since in FALCON-
114 Unzip, the primary contigs are the longest path through the assembly string graph, the secondary
115 haplotigs are by definition shorter and can be more numerous, resulting in lower contiguity for
116 the haplotigs. Regions of the genome with very low heterozygosity remain as collapsed
117 haplotypes in the primary contigs.

118     The PacBio long-read assembly for the hummingbird was of similar quality, with 1076
119 primary contigs generating a primary haploid genome size of 1.01 Gb (1.14 Gb expected; [17]),
120 and a contig N50 of 5.36 Mb, representing a 116-fold reduction in the number of contigs and a
121 201-fold improvement in contiguity over the reference (**Table 1B**). The length of the assembled
122 secondary haplotigs for the hummingbird was similar to that of the primary contig backbone
123 (1.01 Gb; **Table 1B**) indicating that there was sufficient heterozygosity to phase most of the
124 diploid genome into the two haplotypes.

125

| 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 |

126
127 **Table 1:** *De novo* genome assembly statistics comparing intermediate-read length and short-read length
128 assemblies with the long-read assemblies. (A) Zebra finch intermediate-read length (Sanger-based, NCBI
129 accession # GCF_000151805, version 3.2.4) compared to the long-read length PacBio-based assembly.
130 (B) Anna's hummingbird short-read length (Illumina-based, accession # GCF_000699085) compared to
131 the long-read length PacBio-based assembly. Improvement is calculated between the 2nd and 3rd columns
132 for the primary PacBio-based haplotype. The higher number of contigs in the secondary haplotype (5th
133 column) are a result of the arbitrary assignment of shorter haplotypes to the haplotig category.

134

135 **The long-read assemblies have more complete conserved protein coding genes**
136 To assess gene completeness, we analyzed 248 highly conserved eukaryotic genes from the
137 CEGMA human set [18, 19] in each of the assemblies. Both the PacBio-based zebra finch and
138 hummingbird assemblies showed improved resolution of these gene sequences, with a close to
139 doubling (~71%) for the zebra finch and 26% increase for the hummingbird in the number of
140 complete or near-complete (>95%) CEGMA genes assembled, compared to the references (**Fig.**

4

**1A**). Because updating the CEGMA gene sets was recently discontinued due to lack of continued funding and ease of use (http://www.acgt.me/blog/2015/5/18/goodbye-cegma-hello-busco), we also searched for a set of conserved, single-copy genes from the orthoDB9 [20] gene set using the recommended replacement BUSCO pipeline [21]. We observed more modest improvements (~10%) in the number of complete genes in the zebra finch (and no change with the hummingbird) when assessed using the BUSCO v2.0 pipeline on a set of 303 single-copy conserved eukaryotic genes (**Fig. 1B**), and barely any change (1-3%) when using a newly generated BUSCO set of 4915 avian genes (**Fig. 1C**). However, we believe that the moderate increase or no change is due to the fact that much of the BUSCO gene sets were generated from incomplete genome assemblies with short- to intermediate-length reads; for example, the 4915 protein coding avian gene set is generated mostly from the 40+ avian species that the Avian Phylogenomics Project sequenced with short reads [6], including the reference hummingbird [22]. Supporting this view, we extracted the overlapping orthologous genes in the different CEGMA and BUSCO datasets, and found that the CEGMA genes are on average significantly longer than their BUSCO counterparts (**Fig. S3**). When we manually examined genes randomly, many of the BUSCO protein coding sequences were truncated relative to the corresponding CEGMA gene and the PacBio-based assemblies (e.g. the ribosomal protein RLP24 aves BUSCO gene is 117 a.a., whereas the CEGMA & PacBio assembly are 163 a.a.). When compared to the CEGMA 303 eukaryotic set that includes several higher-quality genome assemblies, the PacBio-based assemblies had very few fragmented genes compared to the Sanger-based and Illumina-based assemblies (**Fig. 1B**). Thus, our new assemblies have the potential to upgrade the BUSCO set to more complete and more accurately assembled genes, a conclusion supported by our analyses below.

**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 [23]. We generated the RNA-Seq data from microdissected RA song nuclei, a region that has convergent gene regulation with the human laryngeal motor cortex (LMC) involved in speech production (**Fig. S4**; [4]). The PacBio-based assembly 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 increase in reads that mapped to the genome more than once (multiple) 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

5

180  brain regions containing multiple cell types (**Fig. S4A**; [24]) gave very similar results, with 7-
181  11% increased mappings to the PacBio assembled genome (not shown).

182  Regulatory regions have been difficult to identify in the zebra finch genome, as they are
183  often GC-rich and hard to sequence and assemble with short-read technologies. To assess the
184  regulome, we aligned HK327ac ChIP-Seq reads generated from the RA song nucleus (see
185  methods and [25]) to the zebra finch genome assemblies using Bowtie2 for single-end reads [26].
186  H3K27ac activity is generally high in active gene regulatory regions, such as promoters and
187  enhancers [27]. Similar to the transcriptome, there was an increase (~4%) of HK327ac Chip-Seq
188  genomic reads that mapped to the PacBio-based assembly compared to the Sanger-based
189  reference (**Fig. 2A**). Unlike the RNA-Seq transcript reads, the ChIP-Seq genomic reads showed a
190  significant 10% increase in unique mapped reads with a concomitant decrease in multiple
191  mapped reads (**Fig. 2B**). We believe this difference is due to technical reasons in using paired-
192  end transcript (RNA-Seq) versus single-end genomic (ChIP-Seq) read data, as a multiple-
193  mapped increase with the RNA-Seq transcript data was not detected when using only one read of
194  each pair-end ($p$=0.3, paired t-test, n=5). Overall, these findings are consistent with the PacBio-
195  based assembly having a more complete and structurally accurate assembly for both coding and
196  regulatory non-coding genomic regions.
197

**Completion and correction of genes important in vocal learning and neuroscience research**
199  The genome-wide analyses above demonstrate improvements to overall genome assembly
200  quality using long reads, but they do not inform about real-life experiences with individual genes
201  where there have been challenges with assemblies. We undertook a detailed analysis of four of
202  our favorite genes that have been widely studied in neuroscience and in vocal learning/language
203  research in particular: *EGR1*, *DUSP1*, *FOXP2*, and *SLIT1*.
204

205  *EGR1*. The early growth response gene 1 (*EGR1*) is an immediate early gene transcription factor
206  whose expression is regulated by activity in neurons, and is involved in learning and memory
207  [28]. It is up-regulated in song-learning nuclei when vocal learning birds produce song [29]; it
208  belongs to a large set of genes representing 10% of the transcribed genome that are up- or down-
209  regulated in response to activity in different cell types of the brain [25]. Studying the
210  mechanisms of regulation of *EGR1* and other immediate early genes has been an intensive area
211  of investigation [30, 31], but in all intermediate- and short-read bird genome assemblies we
212  examined thus far, part of the GC-rich promoter region is missing (**Fig. 3A, gap 1**).

213  In the zebra finch Sanger-based reference, *EGR1* is located on a 5.7 kb contig (on
214  chromosome 13), bounded by the gap in the GC-rich promoter region and 2 others downstream
215  of the gene; gaps between contigs in the published reference were given arbitrary 100 Ns [2]. We
216  found that the PacBio long-read assembly completely resolved all three gaps in the zebra finch
217  *EGR1* locus for both alleles, resulting in complete protein coding and surrounding gene bodies in
218  a 205.5 kb primary contig and a 129.1 kb secondary haplotig (**Fig. 3B**; **Fig. S5A**). The promoter
219  region gap, located 572 bp upstream of the start of the first exon, was resolved by an 804 bp

220 70.1% GC-rich PacBio-based sequence (**Fig. 3B, black**). In addition to the 100 Ns in the
221 reference, there were 241 bp to the left and right of this gap of low quality sequence (<QV40;
222 **Fig 3A, blue; 3B, red**) that was not supported by the PacBio data. For the second gap located
223 ~2.2 kb downstream of the *EGR1* gene, there was an adjacent 210 bp low-similarity tandem
224 repeat region that was also not supported by the PacBio data and also had low quality scores (**Fig**
225 **3A,B, gap 2**). The third 100 N gap, located ~3.5 kb downstream of the *EGR1* gene, was resolved
226 by 18 bp of sequence in the PacBio assembly (**Fig. 3B, gap 3**). The PacBio-based differences in
227 the assembly were supported by numerous long-read (>10,000 bp) molecules that extended
228 through the entire gene, spanning all three gaps (**Fig. S6A**). The two haplotypes were >99.8%
229 identical over the region shown (**Fig. 3B**), with only one synonymous heterozygous SNP in the
230 coding sequence (G at position 169,283 in the primary contig 405; T at position 92,478 in
231 secondary haplotig 405_002; tick mark in **Fig. 3B**).

232 In the Illumina-based hummingbird reference, *EGR1* was represented by 3 contigs
233 separated by 2 large gaps of 544 Ns and 1987 Ns respectively (**Fig. 3C**), in a large 2.98 Mb
234 scaffold. In contrast, in the PacBio-based hummingbird assembly, *EGR1* was fully resolved in a
235 large 810 kb contig (**Fig. 3C**). Gene prediction (using Augustus [32]) yielded a protein of the
236 same length as the finch EGR1 protein (510 a.a.), and with high (93%) sequence homology (**Fig.**
237 **3D**). The PacBio-based assembly revealed that the larger gap in the Illumina-based assembly
238 harbors the beginning of the *EGR1* gene, including the entire first exon, two thirds of the first
239 intron, and the GC-rich promoter region (**Fig. 3C, black**). Due to this gap in the reference, the
240 corresponding NCBI gene prediction (accession XP_008493713.1) instead recruited a stretch of
241 sequence ~7 kb upstream of the gap, predicting a first exon that has no sequence homology with
242 *EGR1* in the PacBio-based assembly or to sequences of other species (**Fig. 3C & D**). Upstream
243 of this gap in the Illumina-based assembly was also a 200 bp tandem repeat that was not
244 supported by the PacBio sequence reads and the assembly (**Fig. 3C, red; Fig. S5B**). These
245 PacBio-based differences in the assembly were further validated by single-molecule Iso-Seq
246 mRNA long-reads of *EGR1* from a closely related species (the Ruby-throated hummingbird;
247 kindly provided by R. Workman & W. Timp) that fully contained both predicted exons (**Fig.**
248 **S6B**). The PacBio-based assembly did not generate a secondary haplotype for this region,
249 indicating that the two alleles are identical or nearly identical for the entire 810 kb contig in the
250 individual sequenced. Upstream and downstream of a high homology region that includes the
251 *EGR1* exons, intron, and GC-rich promoter, there was little sequence homology between the
252 PacBio-based hummingbird and zebra finch assemblies (**Fig. S7**).

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

260

261   *DUSP1*. The dual specificity phosphatase 1 (*DUSP1*) is also an immediate early gene, but one
262   that regulates the cellular responses to stress [33]. In all species examined thus far it is mostly
263   up-regulated by activity in the highly active thalamic-recipient primary sensory neurons of the
264   cortex (i.e. mammal cortex layer 4 cells and the comparable avian intercalated pallial cells), but
265   within the motor pathways, it is only up-regulated to high levels by activity in the vocal learning
266   circuits of vocal learners [11, 34]. This specialized regulation in vocal learning circuits has been
267   proposed to be associated with convergent microsatellite sequences found in the upstream
268   promoter region of the gene mainly in vocal learning species [11]. This was determined by PCR-
269   cloning of single genomic molecules from multiple species, because the reference assemblies did
270   not have this region properly assembled [11].

271       In the zebra finch Sanger-based reference, *DUSP1* is located on the chromosome 13
272   scaffold, separated in 3 contigs, with 2 gaps, all surrounded by low quality sequences (**Fig. 4A**).
273   The NCBI gene prediction of this assembly resulted in 4 exons generating a 322 a.a.
274   (XP_002192168.1), which is ~13% shorter than the *DUSP1* homologs of other species, e.g.
275   chicken (369 a.a., Genbank accession NP_001078828), rat (367 a.a., NP_446221), and human
276   (367 a.a, NP_004408). The 2 gaps coincide with the end of the first predicted exon and the
277   beginning of the third predicted exon (**Fig. 4A**). An additional gap upstream of the coding
278   sequence falls within the known microsatellite repeat region (**Fig. 4A**). The PacBio-based
279   assembly completely resolved the entire region for both alleles, in an 8.4 Mb primary contig and
280   an 8.0 Mb secondary haplotig (**Fig. 4B, Fig. S8A**). The Augustus gene prediction resulted in a
281   protein with 4 exons but now with a total length of 369 a.a. that was homologous across its
282   length to *DUSP1* of other vertebrate species (e.g., 96% with chicken GGv5 assembly, also
283   recently updated with long reads). Comparing the two assemblies revealed that: 1) the first exon
284   in the Sanger-based reference is truncated by 28 a.a. in the gap; 2) near the edge of that
285   truncation are three a.a. that appear to be errors (**Fig. 4**; residues 81, 89, and 98), as they are
286   different from genomes of other songbird species using high coverage Illumina reads (**Fig. S9A**),
287   with strong support in the zebra finch PacBio reads (**Fig. S9B**); 3) the second exon and adjacent
288   intron is missing a 80.8% GC-rich 0.46 kb sequence in the reference, and is instead replaced by a
289   1.7 kb contig of a partially repeated sequence from the microsatellite region upstream of *DUSP1*
290   (R' in **Fig. 4B**), part of which was erroneously recruited in the second exon of the NCBI
291   reference gene prediction (**Fig. 4D**); and 4) the microsatellite repeat itself is erroneously partially
292   duplicated in the reference, flanking both sides of gap 1 (R'' in **Fig. 4B**). Our PacBio phased
293   assembly revealed why both instances of R' are not identical in the reference, because they in
294   fact belong to the different haplotypes: the 1.7 kb contig corresponds to the upstream region in
295   the primary PacBio haplotype (contig 32) whereas the actual upstream region in the reference
296   corresponds to the upstream region in the secondary PacBio haplotype (contig 32_022) (**Fig.
297   4B**). This main microsatellite region is 76 bp longer (796 *vs.* 720 bp) in the primary haplotype,
298   and the neighboring smaller upstream microsatellite contains 3 additional 20-21 bp repeats (11

299    *vs.* 8) in the primary haplotype (**Fig. 10A**). Within the protein coding sequence there were four
300    synonymous heterozygous SNPs between haplotypes (not shown).

301        In the hummingbird Illumina-based assembly, the *DUSP1* region was represented by 2
302    contigs separated by a large 1005 N gap (**Fig. 4C**), on a 7 Mb scaffold. In the PacBio-based
303    assembly, the entire gene was fully resolved (**Fig. 4C; Fig. S8B**), in a much larger gapless 12.8
304    Mb contig (the second allele is fully resolved in a 3.8 Mb contig). Comparing the two assemblies
305    revealed that because of the gap in the Illumina-based reference, it lacks about half of the
306    *DUSP1* gene, including the first two exons and introns, and ~380 bp upstream of the start of the
307    gene (**Fig. 4C**). As a result, the corresponding NCBI gene prediction (XP_008496991.1)
308    recruited a sequence ~44 kb upstream predicting 46 a.a. with no sequence homology to *DUSP1*
309    of other species, whereas the PacBio-based assembly yielded a 369 a.a. protein with 99%
310    sequence homology to the PacBio-based zebra finch and chicken *DUSP1* (**Fig. 4D**). A 200 bp
311    tandem repeat in the Illumina-based assembly downstream of the gap, erroneously in exon 3, is a
312    misplaced copy of the microsatellite region (**Fig. 4C; Fig. S8B**). This is the reason why two
313    thirds of exon 3 is erroneously duplicated in the NCBI protein prediction (**Fig. 4D**). These
314    PacBio-based differences in the assembly were validated by single-molecule Iso-Seq mRNA
315    long-reads of *DUSP1* (**Fig. S11A**). The PacBio assemblies also revealed that the microsatellite
316    region was significantly shorter in the hummingbird (~270 bp) than the zebra finch genome
317    (~1100 bp; **Fig. S10B**).

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

324

325    *FOXP2*. The forkhead box P2 (*FOXP2)* gene plays an important role in spoken-language
326    acquisition [35]. In humans, a point mutation in the protein coding binding domain in the KE
327    family [36] as well as deletions in the non-coding region of *FOXP2* [37] results in severe spoken
328    language impairments in heterozygous individuals (homozygous is lethal). In songbirds, FOXP2
329    expression in the Area X song nucleus is differentially regulated by singing activity and during
330    the song learning critical period, and is necessary to properly imitate song [38-40]. In mice,
331    although vocalizations are mainly innate, animals with the KE mutation demonstrate a syntax
332    apraxia-like deficit in syllable sequencing similar to that of humans [41, 42]. Thus, *FOXP2* has
333    become the most studied gene for understanding the genetic mechanisms and evolution of
334    spoken language [43], yet we find that the very large gene body of ~400 kb is incompletely
335    assembled, including in vocal learning species (**Fig. 5A**).

336        In the zebra finch Sanger-based reference, *FOXP2* is located on the chromosome 1A
337    scaffold and separated into 10 contigs (1 to 231 kb in length) with nine 100 N gaps each (**Fig.
338    5A**). These include 2 gaps immediately upstream of the first exon, making the beginning of the

339 gene poorly resolved. The provisional RefSeq mRNA for *FOXP2* (NM_001048263.1) contains
340 19 exons and encodes a 711 a.a. protein (NP_001041728.1). In the PacBio-based assembly, the
341 entire 400 kb gene is fully resolved for both haplotypes in 21.5 Mb and 7.6 Mb contigs,
342 respectively (**Fig. S12A**). As observed in the previous examples, sequences of various sizes
343 surrounding all 9 gaps in the Sanger-based reference were unsupported by the PacBio data,
344 resulting in a total of 2509 bp of corrected sequence in the PacBio-based primary haplotype (**Fig.**
345 **5B**). The two filled gaps in the upstream region and the next gap in the first intron were GC-rich
346 (77.6%, 66.5%, and 67.8%, respectively; **Fig. 5A,C**), indicative of the likely cause of the poor
347 quality Sanger-based reads (**Fig. 5D**). The DNA sequence between the two assembled PacBio
348 haplotypes was >99% similar across the entire 400 kb *FOXP2* gene, and identical over the
349 coding sequence, with differences occurring in the more complex non-coding gaps that were
350 difficult to sequence and assemble by the Sanger method (**Fig. 5B \***61 nucleotide differences
351 total). The predicted protein sequence from the PacBio-based assembly is identical to the
352 predicted Sanger-based reference (NP_001041728.1), with the exception of a.a. residue 42
353 (threonine *vs.* serine) (**Fig. S13A**). The PacBio nucleotide call also exists in the mRNA sequence
354 of another zebra finch animal in NCBI (NM_001048263.2) and in other avian species we
355 examined, and is thus likely a base call error in the Sanger-based zebra finch reference.

356       In the hummingbird Illumina-based assembly, as expected with short-read assemblies
357 relative to the Sanger-based zebra finch reference, the *FOXP2* gene was even more fragmented,
358 in 23 contigs (ranging 0.025 to 2.28 kb in lengths) with 22 gaps (**Fig. S12B**). The two largest
359 gaps encompass the beginning of the gene and first (non-coding) exon, resulting in
360 corresponding low quality predicted mRNA (XM_008496149.1). The predicted protein
361 (XP_008494371.1) includes an introduced correction (a.a. 402; **Fig. S13A**, X nucleotide) to
362 account for a genomic stop codon, and an 88 N gap within exon 6 that artificially splits the exon
363 into two pieces (**Fig. S13B**). In the hummingbird PacBio-based assembly, the *FOXP2* gene is
364 fully resolved and phased into two haplotype contigs of 3.2 Mb each (**Fig. S12B**). The erroneous
365 stop codon is corrected (2170128C [ctg 110] and 2183088C [ctg 110_009], instead of 841788T
366 [Illumina assembly scaffold 125]), and exon 6 is accurately contiguous, removing the gap and an
367 additional 22 bp of erroneous tandem repeat sequence adjacent to the gap (**Fig. S13B & C**). The
368 PacBio-based assembly also corrects three other instances of erroneous tandem duplications over
369 the gene region in the Illumina-based assembly, as well as removes a 462 bp stretch of sequence
370 adjacent to a long homonucleotide A stretch in intron 1 of the Illumina-based assembly (position
371 972040; **Fig. S14A**). These PacBio-based differences in the assembly were validated by single-
372 molecule Iso-Seq mRNA long-reads of *FOXP2* (**Fig. S11B**). The two PacBio assembled
373 haplotypes are >99% similar, with one heterozygous SNP (2172601T (contig 110) *vs.* 2185560A
374 (contig 110_009)) in exon 6 that is silent, and a 708 bp deletion in the secondary haplotype
375 (contig 110_009 [at position 2128952] relative to contig 110; **Fig. S14B**). The Illumina-based
376 assembly has the deleted allele.

377       These findings replicate those of the previously discussed genes, and in addition show
378 that the PacBio-based assembly can fully resolve very large genes, resolve erroneous assembled

379 sequences in gaps due to repeats or homonucleotide stretches, and reveal large haplotype
380 differences. The phased diploid assembly also avoids the possibility of large missed sequences in
381 a haploid only assembly due to deletions in one allele.
382

383 *SLIT1*. Slit homolog 1 (*SLIT1*) is a repulsive axon guidance ligand for the *ROBO1* receptor, and
384 is involved in circuit formation in the developing brain [44]. Recently, *SLIT1* was shown to have
385 convergent specialized down-regulated expression compared to the surrounding brain region in
386 the RA song nucleus of all independently evolved vocal learning bird lineages and in the
387 analogous human LMC [4, 45] (**Fig. S4**), indicating a potential role of *SLIT1* in the evolution and
388 formation of vocal learning brain circuits. A fully resolved *SLIT1*, including regulatory regions,
389 is necessary to assess the mechanisms of its specialized regulation in vocal learning brain
390 regions.
391     In the zebra finch Sanger-based reference, *SLIT1* is located on chromosome 6, split
392 among 8 contigs with 7 gaps, and 7 additional contigs and gaps surrounding the ~40 kb gene
393 (**Fig. 6A**). The SLIT1 gene is complex, with over 35 exons. We noted an incomplete predicted
394 protein of the reference (XP_012430014.1) relative to some other species (chicken
395 [NM_001277336.1], human [NM_003061.2], and mouse [NM_015748.3]), and our *de novo* gene
396 predictions from the reference also resulted in a truncated protein with two missing exons (**Fig.
397 6B**). The PacBio-based assembly fully resolved the gene region, in two alleles on 15.7 Mb and
398 5.6 Mb contigs, respectively, and completely recovered all 35+ exons (**Fig. S15A**). Similar to
399 above, reference sequences flanking the gaps were found to be erroneous and corrected, and an
400 erroneous tandem duplication was also corrected (not shown). Filling in these gaps recovered the
401 two missing exons: exon 1 within a 1 kb region of sequence in the PacBio-based assembly that is
402 75% GC-rich, replacing 390 bp of erroneous gap-flanking sequence; and exon 35 adjacent to a
403 gap (**Fig. 6A,B**). A predicted exon upstream of exon 1 in a repeat region was not supported (**Fig.
404 6A,B**). The PacBio-based assembly thereby generates a complete *SLIT1* gene prediction of 1538
405 a.a. (**Fig. 6B**). The gene is heterozygous in the individual, with 3 codon differences between the
406 two alleles (**Fig. 6B**, positions 90, 1006, and 1363, respectively), and an additional 24 silent
407 heterozygous SNPs across the coding region. The two alleles were phased along the entire length
408 of the gene.
409     In the hummingbird Illumina-based assembly, the *SLIT1* gene is separated on 9 contigs
410 with 8 gaps ranging in length from 91 to 1018 bp, comprising 3320 bp of missing sequence, or
411 5.3% of the gene region (**Fig. S15B**). The PacBio-based assembly fully resolved and phased
412 *SLIT1* into haplotypes on 9.9 Mb contigs (**Fig. S15B**). The resulting protein of 1538 a.a. has high
413 homology to the zebra finch PacBio-based *SLIT1* (95% a.a. identity; **Fig. 6B**) and the individual
414 is homozygous for the SLIT1 protein. Comparisons revealed that as with the Sanger-based
415 reference, the first exon (68 a.a.) is missing completely in the Illumina-based assembly (**Fig. 6B**),
416 corresponding to a gap of 495 Ns, which the PacBio-based assembly replaced by a 567 bp 76%
417 GC-rich sequence (**Fig. S15B**). In addition, there were two sequence errors in the Illumina-based

11

assembly, 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 that errors in the Sanger-based and Illumina-based assemblies were corrected 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, which 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 [10], with both fully represented in the PacBio-based assemblies. We also noted cases where an assembled gene was incorrectly localized on a scaffold in the Sanger-based assembly whose synteny was corrected with the PacBio-based assembly, such as the vasopressin receptor AVPR1B, which will be reported on in more detail separately. Data for these types of errors were not shown due to space limitations, but they offer further examples of the important improvements of PacBio long-read technology for generating more accurate genome assemblies.

## Discussion and Conclusions

Although the intermediate-read and short-read assemblies had correct sequences and assembled regions in terms of total base pairs covered, the long-read assemblies revealed numerous errors within and surrounding many genes. These errors are not simply in so-called "junk" intergenic repetitive DNA known to be hard to assemble with short reads [46, 47], but within functional regions of genes. The assemblers for the short reads sometimes take a repetitive sequence, some in functional repetitive regulatory regions, and insert them in a non-repetitive region of a gene, resulting in an error. Some of these assembly errors and gaps in the sequences lead to gene and protein coding sequence prediction errors, sometimes recruiting completely wrong sequence in the protein.

The PacBio-based long-read assemblies corrected these problems, and for the first time resolved gene bodies of all the genes we examined into single, contiguous, gap-less sequences. The phasing of haplotypes, although initially done to prevent a computationally introduced indel error, reveal how important phasing is to prevent assembly and gene prediction errors. Thus far, we have not seen an error (i.e. difference) in the genes we examined in the PacBio-based long-

read assembly relative to the other assemblies that was supported by single sequenced genomic DNA molecules, RNA-Seq and Iso-Seq mRNA molecules, or other independent evidence. With these improvements, we now, for the first time, have complete and accurate assembled genes of interest that we now can pursue further without the need to individually and arduously clone, sequence, and correct the assemblies one gene at a time.

Our study highlights the value of maintaining frozen tissue or cells of the individuals used to create previous reference genomes, as we could only discover some of the errors (e.g. caused by haplotype differences) by long-read *de novo* genome assemblies of the same individual used to create the reference. We are now using these PacBio-based assemblies with several groups and companies as starting assemblies for scaffolding into phased, diploid, chromosome-level zebra finch and hummingbird assemblies to upgrade the references, which will be reported on separately. However, even without scaffolding, these more highly contiguous assemblies will be helpful to researchers to extract more accurate assemblies of their genes of interests, saving a great amount of time and energy, while adding new knowledge and biological insights necessary for understanding gene structure, function, and evolution.

## Materials & Methods

### DNA isolation
For both the zebra finch and hummingbird, frozen muscle tissue from the same animals used to create the Sanger-based [2] and Illumina-based [6] references, respectively, was processed for DNA isolation using the KingFisher Cell and Tissue DNA Kit (97030196). Tissue was homogenized in 1 ml of lysis buffer in M tubes (Miltenyi Biotec) using the gentleMACS™ Dissociator at the Brain 2.01 setting for 1 minute. The cell lysate was treated with 40 ul of protease K (20mg/ml) and incubated overnight. DNA was purified using the KingFisher Duo system (5400100) using the built in KFDuoC_T24 DW program.

### 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).

13

498      SMRT sequencing was performed on the Pacific Biosciences RS II instrument at Pacific
499 Biosciences using an on plate concentration of 125 pM, P6-C4 sequencing chemistry, with
500 magnetic bead loading, and 360 minute movies. A total of 124 SMRT Cells were run for the
501 zebra finch and 63 SMRT Cells for the hummingbird. Sequence coverage for the zebra finch was
502 ~96 fold, with half of the 114 Gb of data contained in reads longer than 19 kb. For the
503 hummingbird, coverage was ~70 fold, with half of the 40.4 Gb of data contained in reads longer
504 than 22 kb (**Fig. S2**).
505

### Assembly

507 Assemblies were carried out using FALCON v0.4.0 followed by the FALCON-Unzip module
508 [16]. FALCON is based on a hierarchical genome assembly process [48]. It constructs a string
509 graph from error-corrected PacBio reads that contains 'haplotype-fused' genomic regions as well
510 as "bubbles" that capture divergent haplotypes from homologous genomic regions. The
511 FALCON-Unzip module then assigns reads to haplotypes using heterozygous SNP variants
512 identified in the FALCON assembly to generate phased contigs corresponding to the two alleles.
513 The diploid nature of the genome is thereby captured in the assembly by a set of primary contigs
514 with divergent haplotypes represented by a set of additional contigs called haplotigs. Genomic
515 regions with low heterozygosity are represented as collaped haplotypes in the primary contigs.
516 Genome assemblies were run on an SGE-managed cluster using up to 30 nodes, where each node
517 has 512 Gb of RAM distributed over 64 slots. The same configuration files were used for both
518 species (**Additional file 1**). Three rounds of contig polishing were performed. For the first round,
519 as part of the FALCON-Unzip pipeline, primary contigs and secondary haplotigs were polished
520 using haplotype-phased reads and the Quiver consensus caller. For the second and third rounds
521 of polishing, using the "resequencing" pipeline in SMRTlink v3.1, primary contigs and haplotigs
522 were concatenated into a single reference and BLASR was used to map all raw reads back to the
523 assembly, followed by consensus calling with Arrow.
524

### Genome completeness

526 To assess quality and completeness of the assemblies, we used a set of 248 highly conserved
527 eukaryotic genes from the CEGMA human set [19] and located them in each of the assemblies
528 compared in this study. Briefly, the CEGMA human peptides were aligned to each genome using
529 genblastA [49]. The regions showing homology were then used to build gene models with
530 exonerate [50] which were then assessed for frameshifts using custom shell scripts. In addition,
531 we queried each genome for a set of 303 eukaryotic conserved single-copy genes as well as from
532 4915 conserved single-copy genes from 40 different avian species using the BUSCOv2.0
533 pipeline [21].
534      To compare protein amino acid sequence size between the CEGMA and BUSCO
535 datasets, we performed blastp of each CEGMA sequence against the ancestral proteins of the
536 target BUSCO dataset. We took the single best hit with an e-value cut off of 0.001 and extracted
537 the CEGMA and BUSCO protein length values. We then ran a one-sided paired Wilcoxon

14

538    signed-rank test of the two lengths for each protein.

539

**Gene prediction**

541    Gene predictions for the zebra finch PacBio-based assembly were conducted by running
542    Augustus gene prediction software (v3.2.2, [32]) on the contigs, and incorporating the Illumina
543    short read RNA-Seq brain data aligned with Tophat2 (v2.0.14, [23]) as hints for possible gene
544    structures. The data consisted of 146,126,838 paired-end reads with an average base quality
545    score of 36. Augustus produces a distribution of possible gene models for a given locus and
546    models that are supported by our RNA-Seq data are given a "bonus" while the gene models not
547    supported by RNA-Seq data are given a "penalty". This results in the gene model most informed
548    by biological data being selected as the most likely gene model for that locus.

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

555

**RNA-Seq**

557    RNA sequencing was centered around vocal learning brain regions in the zebra finch and will be
558    described in more detail in a later publication. We utilized our data here for population analyses
559    of assembly quality and for initial annotations. In brief, following modifications of a previously
560    described protocol [25], nine adult male zebra finches were isolated in soundproof chambers for
561    12 hours in the dark to obtain brain tissue from silent animals. Then brains were dissected from
562    the skull and sectioned to 400 microns using a Stoelting tissue slicer (51415). The sections were
563    moved to a petri dish containing cold PBS with proteinase inhibitor cocktail (11697498001).
564    Under a dissecting microscope (Olympus MVX10), the four principle song nuclei (Area X,
565    LMAN, HVC, and RA) as well as their immediate adjacent brain regions were microdissected
566    using 2mm fine scissors and placed in microcentrifuge tubes. The samples were stored at -80
567    °C. Then RNA was isolated and quantified, and samples of two birds were then pooled for each
568    replicate, resulting in 5 replicates (one single animal in one). RNA was converted to cDNA and
569    library preparation was performed using the NEXTflex™ Directional RNA-Seq Kit (Illumina)
570    and paired-end reads were sequenced on an Illumina HiSeq 2500 system. Adapters and poor
571    quality bases (<30) were trimmed using fastq-mcf from the ea-utilities package, and reads were
572    aligned to assemblies using Tophat2 (v2.0.14).

573

**Chip-Seq**

575    Three adult male zebra finches were treated as above, the brains dissected, and the RA and
576    surrounding arcopallium of each bird was then processed individually using the native ChIP
577    protocol described in [51] with an H3K27ac antibody (Ab#4729). The DNA libraries were

15

578 prepared using the MicroPlex Library Preparation Kit v2 (C05010012). 50 bp single-end
579 sequencing was done on the Illumina HiSeq 4000 system. The reads were aligned to the
580 assemblies using Bowtie2 (v2.2.9, [26]). More detail will be provided in a later publication
581 focusing on vocal learning brain regions.
582

**Comparative analyses between assemblies for individual genes**

583
584 The Sanger-based reference zebra finch assembly in the UCSC browser and the Ilumina-based
585 reference Anna's Hummingbird in Avianbase (http://avianbase.narf.ac.uk/index.html), and both
586 in NCBI where used for comparing with the Pacbio assembly. In the UCSC browser, there are
587 two annotations, one from 2008 (http://genome.ucsc.edu/cgi-bin/hgGateway?db=taeGut1) and
588 the other from 2013 (http://genome.ucsc.edu/cgi-bin/hgGateway?db=taeGut2), with some
589 differences between them. Our findings were similar, although not always identical, with both
590 annotations, with errors being present in both annotations based on the Pacbio assembly. The
591 nucleotide quality score tract was only available in the 2008 browser.
592 Multiple species sequence alignments were done with BioEdit v7.2.5
593 (http://www.mbio.ncsu.edu/bioedit/bioedit.html) [52]; Dotplots of alignments were generated
594 with Gepard v1.4 (http://cube.univie.ac.at/gepard) [53]; Alignments of raw SMRT genome reads
595 to the assembled genomes were done with Blasr, which is part of SMRTLink software from
596 Pacbio; Iso-Seq reads were aligned with GMAP (http://research-pub.gene.com/gmap/) [54].
597
598

## Availability of data

599
600 This Whole Genome Shotgun project has been deposited at DDBJ/ENA/GenBank under
601 BioProject PRJNA368994. The zebra finch accession number is MUGN00000000 and SRA for
602 raw reads is SRS1954332. The Anna's Hummingbird accession number is MUGM00000000 and
603 SRA is SRP061272.
604

## Competing interest

605
606 Jonas Korlach, Sarah Kingan, Chen-Shan Chin are full-time employees at Pacific Biosciences, a
607 company developing single-molecule sequencing technologies.
608

## Funding

609
611

## Author contributions

612
613 J.K. and E.D.J. designed the project and wrote the manuscript; C.S.C. and S.K. carried out
614 genome assemblies; J.K., G.G. and S.K. conducted analyses on single genes as well as CEGMA
615 and BUSCO analyses; G.G. and J-N.A. conducted RNA-Seq experiments, L.C. conducted Chip-
616 Seq experiments; J.H. processed samples; and all authors contributed to writing and editing the
617 manuscript.

618

## Acknowledgements

631

632

633

## Additional Files

634 Supporting data is included in supplementary figures S1-S15.

635

17

## Figure legends

**Figure 1.** Gene completeness within assemblies. *(A)* Comparison to a 248 highly conserved core CEGMA eukaryote gene set using human genes [19], 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%). Gent 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 [21]. Complete is ≥ 95% complete; fragmented is < 95% complete; missing is not found. *(C)* Comparison to 4915 single-copy conserved genes from the avian BUSCO gene [21].

**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

18

677 *C*, showing numerous corrections to the Sanger-based and Ilumina-based protein predictions by
678 both PacBio-based assemblies.

679

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

690

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

698

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

703

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

707

708 **Supplementary Figure S3.** Box plots comparing protein coding sequence lengths of
709 orthologous proteins between the CEGMA and BUSCO eukaryotic and avian datasets. ** p <
710 0.001; *** p < 0.0001, one-sided paired Wilcoxon signed-rank test, prediction of the proteins
711 being longer in CEGMA datasets.

712

713 **Supplementary Figure S4.** Vocal learning and adjacent brain regions in songbirds used for
714 RNA-Seq and ChIP-Seq analyses, and comparison with humans. *(A)* Drawing of a zebra finch
715 male brain section showing specialized vocal learning pathway and associated profiled song
716 nuclei RA, HVC, LMAN, and Area X. *(B)* Drawing of a human brain section showing spoken-

19

717 language pathway and analogous brain regions. Black arrows, posterior vocal motor pathway;
718 White arrows, anterior vocal learning pathway; Dashed arrows, connections between the two
719 pathways; Red arrow, specialized direct projection from forebrain to brainstem vocal motor
720 neurons in vocal learners. Italicized letters adjacent to the song and speech regions indicates
721 regions (in songbirds) that show mainly show motor *(m)*, auditory *(a)*, equally both motor and
722 auditory *(m/a)* neural activity or activity-dependent gene expression. Figure from [55] and [4].
723

724 *Abbreviations*: A1-L4, primary auditory cortex – layer 4; Am, nucleus ambiguous; Area X, a
725 vocal nucleus in the striatum; aSt, anterior striatum vocal region; aT, anterior thalamus speech
726 area; Av, avalanche; aDLM, anterior dorsolateral nucleus of the thalamus; DM, dorsal medial
727 nucleus of the midbrain; HVC, a vocal nucleus (no abbreviation); L2, auditory area similar to
728 human cortex layer 4; LSC, laryngeal somatosensory cortex; LMC, laryngeal motor cortex;
729 MAN, magnocellular nucleus of the anterior nidopallium; MO, oval nucleus of the anterior
730 mesopallium; NIf, interfacial nucleus of the nidopallium; PAG, peri-aqueductal gray; RA, robust
731 nucleus of the arcopallium; v, ventricle space
732

733 **Supplementary Figure S5.** Dot plot of sequence comparisons for genome assemblies of the
734 *EGR1* region. *(A)* Comparison of zebra finch PacBio-based versus Sanger-based assemblies for
735 the region containing *EGR1*, showing the GC-rich promoter region and closing and corrections
736 of gaps for the PacBio-based assembly. *(B)* Comparison of hummingbird Illumina-based versus
737 PacBio-based assemblies for the region containing *EGR1*, showing an erroneous tandem
738 duplication in the Ilumina-based assembly and closing of gaps for the PacBio-based assembly.
739

740 **Supplementary Figure S6.** Single SMRT genomic reads and Iso-Seq mRNA reads supporting
741 Pacbio *EGR1* assembly. *(A)* Zebra finch PacBio SMRT reads (rows) mapped against the zebra
742 finch PacBio assembly (contig 405, entire *EGR1* region, same as Fig. 3A). Reads are shaded by
743 length (>10 kb reads = black). *(B)* Example of a single Ruby-throated hummingbird Iso-Seq read
744 mapped against Illumina-based (top) and PacBio-based (bottom) Anna's hummingbird genome
745 assemblies using GMAP. Note the first exon (blue) which is present in the Iso-Seq read is
746 missing in the Illumina-based assembly, but present in the PacBio-based assembly.
747

748 **Supplementary Figure S7.** Dot plot of sequence comparison for the PacBio-based hummingbird
749 and zebra finch *EGR1* region assemblies. Note regions of high species conservation and
750 divergence surrounding *EGR1*. Blue box, location of the *EGR1* exons and intron.
751

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

20

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

758

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

764

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

770

771     **Supplementary Figure S11.** Single Iso-Seq mRNA reads supporting Pacbio assemblies. *(A)*
772     Full-length PacBio mRNA sequence Iso-Seq ruby throated hummingbird reads for DUSP1
773     aligned against the exons of the corresponding primary contigs from Anna's hummingbird
774     Illumina (top panel) and PacBio (bottom panel) assemblies. *(B)* Similar alignments for FOXP2
775     IsoSeq reads.

776

777     **Supplementary Figure S12.** Dot plot comparison of assemblies for the *FOXP2* region. *(A)*
778     zebra finch, *(B)* hummingbird.

779

780     **Supplementary Figure S13.** *(A)* Multiple sequence alignment of the FOXP2 protein for the four
781     assemblies (two zebra finch and two hummingbird) compared in this study, showing correction
782     of a nucleotide error in the Sanger-based zebra finch assembly, and correction of an erroneous
783     stop codon (x) in the Illumina-based hummingbird assembly. Note an extra 18 a.a. stretch in the
784     hummingbird sequence validated by gene prediction of both assemblies, that was not present in
785     the zebra finch. *(B)* Missing 88bp of sequence in exon 6 of Illumina-based assembly. *(C)*
786     Resolution of exon 6 in Pacbio-based assembly, also revealing a SNP.

787

788     **Supplementary Figure S14.** Large regional correction made by the PacBio diploid assembly.
789     *(A)* Correction of an erroneous stretch of 462 bp in the first intron of *FOXP2* in the hummingbird
790     Illumina assembly by the PacBio assembly. *(B)* Dot plot of haplotype variation in the *FOXP2*
791     gene revealed by the PacBio diploid assembly: a 708 bp deletion in the secondary haplotype
792     contig relative to the primary contig.

793

794     **Supplementary Figure S15.** Dot plot comparison of assemblies for the *SLIT1* region. *(A)* zebra
795     finch, *(B)* hummingbird.

796

# References

798 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.

802 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.

804 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.

807 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.

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

812 6. 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.

815 7. 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.

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

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

821 10. 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**.

825 11. 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.

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

830 13. 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.

833 14. 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.

836 15. **FALCON assembler** [https://github.com/PacificBiosciences/FALCON/commit/a1180264c3c7d2de1c5eb55b36 63dce093354dd7]

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

843 17. Gregory TR: **Animal Genome Size Database.** 2017.

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

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

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

852 21. Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM: **BUSCO:**
853 **assessing genome assembly and annotation completeness with single-copy**
854 **orthologs.** *Bioinformatics* 2015, **31:**3210-3212.

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

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

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

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

867 26. Langmead B, Salzberg SL: **Fast gapped-read alignment with Bowtie 2.** *Nat Methods*
868 2012, **9:**357-359.

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

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

873 29. Jarvis ED, Nottebohm F: **Motor-driven gene expression.** *Proc Natl Acad Sci U S A*
874 1997, **94:**4097-4102.

875 30. Flavell SW, Greenberg ME: **Signaling mechanisms linking neuronal activity to gene**
876 **expression and plasticity of the nervous system.** *Annu Rev Neurosci* 2008, **31:**563-
877 590.

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

881 32. Stanke M, Diekhans M, Baertsch R, Haussler D: **Using native and syntenically**
882 **mapped cDNA alignments to improve de novo gene finding.** *Bioinformatics* 2008,
883 **24:**637-644.

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

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

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

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

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

896    38.    Haesler S, Wada K, Nshdejan A, Morrisey EE, Lints T, Jarvis ED, Scharff C: **FoxP2**
897           **expression in avian vocal learners and non-learners.** *J Neurosci* 2004, **24:**3164-
898           3175.
899    39.    Teramitsu I, White SA: **FoxP2 regulation during undirected singing in adult**
900           **songbirds.** *J Neurosci* 2006, **26:**7390-7394.
901    40.    Haesler S, Rochefort C, Georgi B, Licznerski P, Osten P, Scharff C: **Incomplete and**
902           **inaccurate vocal imitation after knockdown of FoxP2 in songbird basal ganglia**
903           **nucleus Area X.** *PLoS Biol* 2007, **5:**e321.
904    41.    Castellucci GA, McGinley MJ, McCormick DA: **Knockout of Foxp2 disrupts vocal**
905           **development in mice.** *Sci Rep* 2016, **6:**23305.
906    42.    Chabout J, Sarkar A, Patel SR, Radden T, Dunson DB, Fisher SE, Jarvis ED: **A Foxp2**
907           **Mutation Implicated in Human Speech Deficits Alters Sequencing of Ultrasonic**
908           **Vocalizations in Adult Male Mice.** *Front Behav Neurosci* 2016, **10:**197.
909    43.    Condro MC, White SA: **Recent Advances in the Genetics of Vocal Learning.** *Comp*
910           *Cogn Behav Rev* 2014, **9:**75-98.
911    44.    Blockus H, Chédotal A: **The multifaceted roles of Slits and Robos in cortical**
912           **circuits: from proliferation to axon guidance and neurological diseases.** *Curr Opin*
913           *Neurobiol* 2014, **27:**82-88.
914    45.    Wang R, Chen CC, Hara E, Rivas MV, Roulhac PL, Howard JT, Chakraborty M, Audet
915           JN, Jarvis ED: **Convergent differential regulation of SLIT-ROBO axon guidance**
916           **genes in the brains of vocal learners.** *J Comp Neurol* 2014.
917    46.    Treangen TJ, Salzberg SL: **Repetitive DNA and next-generation sequencing:**
918           **computational challenges and solutions.** *Nat Rev Genet* 2011, **13:**36-46.
919    47.    Palazzo AF, Gregory TR: **The case for junk DNA.** *PLoS Genet* 2014, **10:**e1004351.
920    48.    Chin CS, Alexander DH, Marks P, Klammer AA, Drake J, Heiner C, Clum A, Copeland
921           A, Huddleston J, Eichler EE, et al: **Nonhybrid, finished microbial genome assemblies**
922           **from long-read SMRT sequencing data.** *Nat Methods* 2013, **10:**563-569.
923    49.    She R, Chu JS, Wang K, Pei J, Chen N: **GenBlastA: enabling BLAST to identify**
924           **homologous gene sequences.** *Genome Res* 2009, **19:**143-149.
925    50.    Slater GS, Birney E: **Automated generation of heuristics for biological sequence**
926           **comparison.** *BMC Bioinformatics* 2005, **6:**31.
927    51.    Brind'Amour J, Liu S, Hudson M, Chen C, Karimi MM, Lorincz MC: **An ultra-low-input**
928           **native ChIP-seq protocol for genome-wide profiling of rare cell populations.** *Nat*
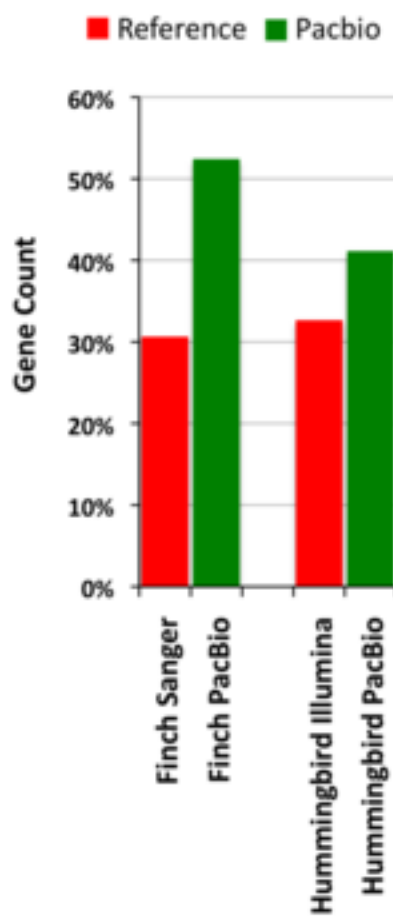929           *Commun* 2015, **6:**6033.
930    52.    Hall TA: **BioEdit: a user-friendly biological sequence alignment editor and analysis**
931           **program for Windows 95/98/NT.** *Nucl Acids Symp Ser* 1999, **41:**95-98.
932    53.    Krumsiek J, Arnold R, Rattei T: **Gepard: a rapid and sensitive tool for creating**
933           **dotplots on genome scale.** *Bioinformatics* 2007, **23:**1026-1028.
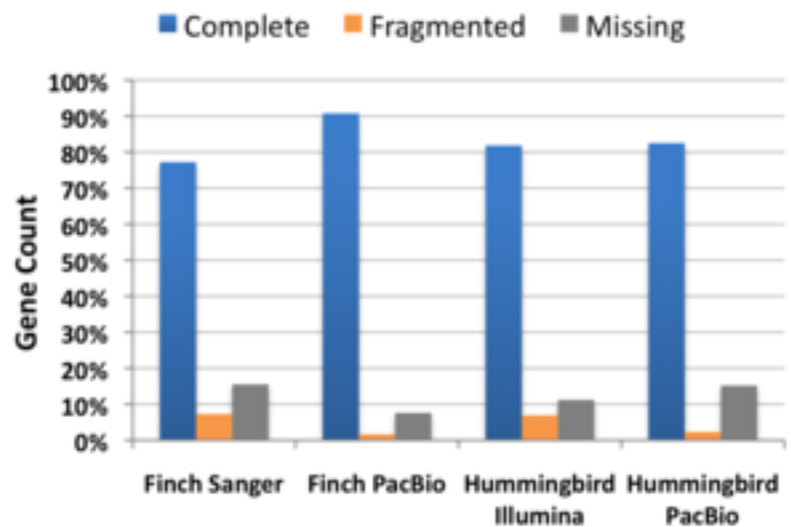934    54.    Wu TD, Watanabe CK: **GMAP: a genomic mapping and alignment program for**
935           **mRNA and EST sequences.** *Bioinformatics* 2005, **21:**1859-1875.
936    55.    Chakraborty M, Jarvis ED: **Brain evolution by brain pathway duplication.**
937           *Philosophical Transactions of the Royal Society B-Biological Sciences* 2015, **370:**50056-
938           50056.
939

Figure 1

Figure 2

Figure 2

Figure 3

Click here to download Figure Korlach Fig. 3 v9.png ⬇



Figure 3

Figure 4

Click here to download Figure Korlach Fig. 4 v9.png ±



Figure 4

Figure 5

**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
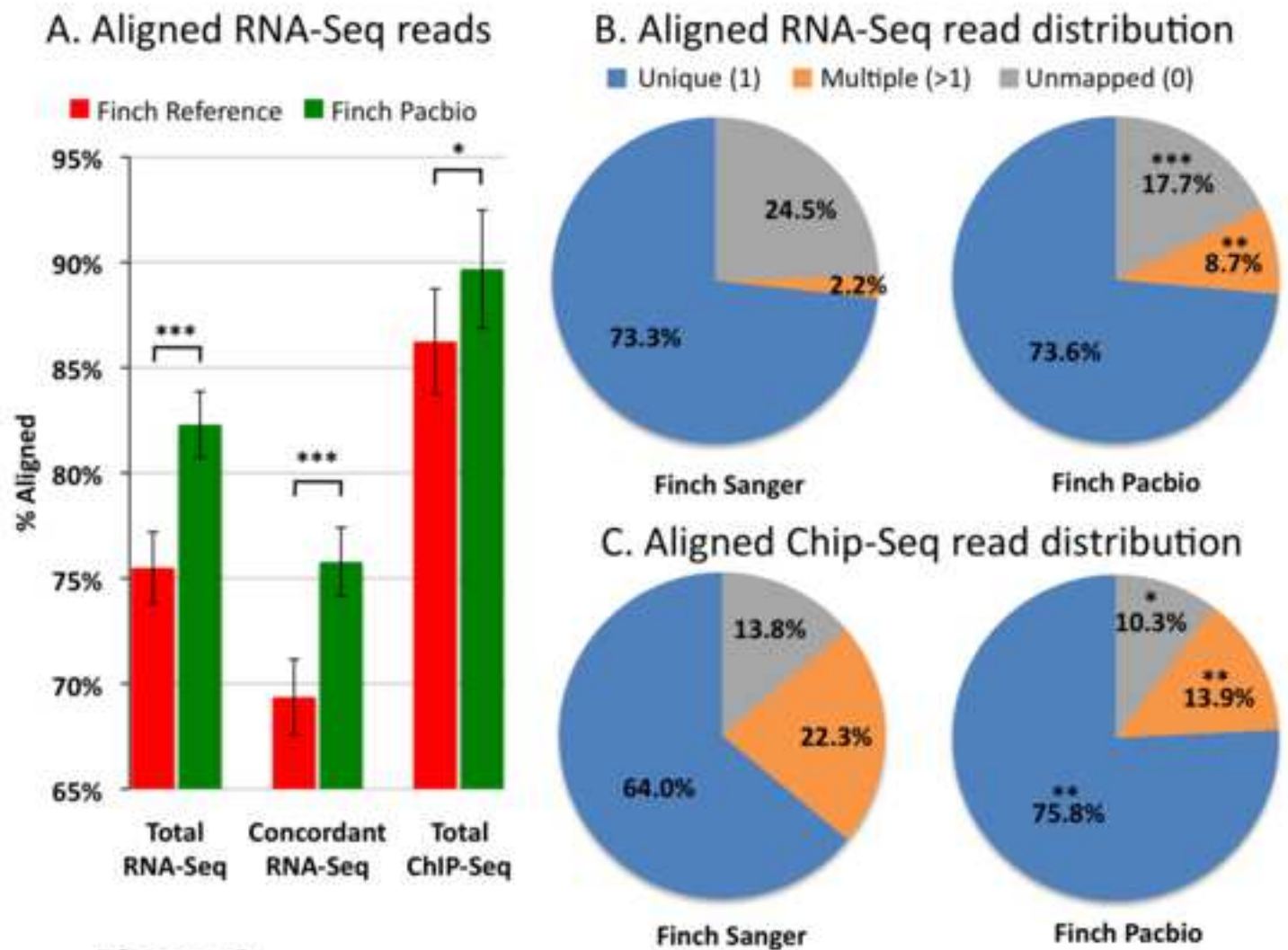


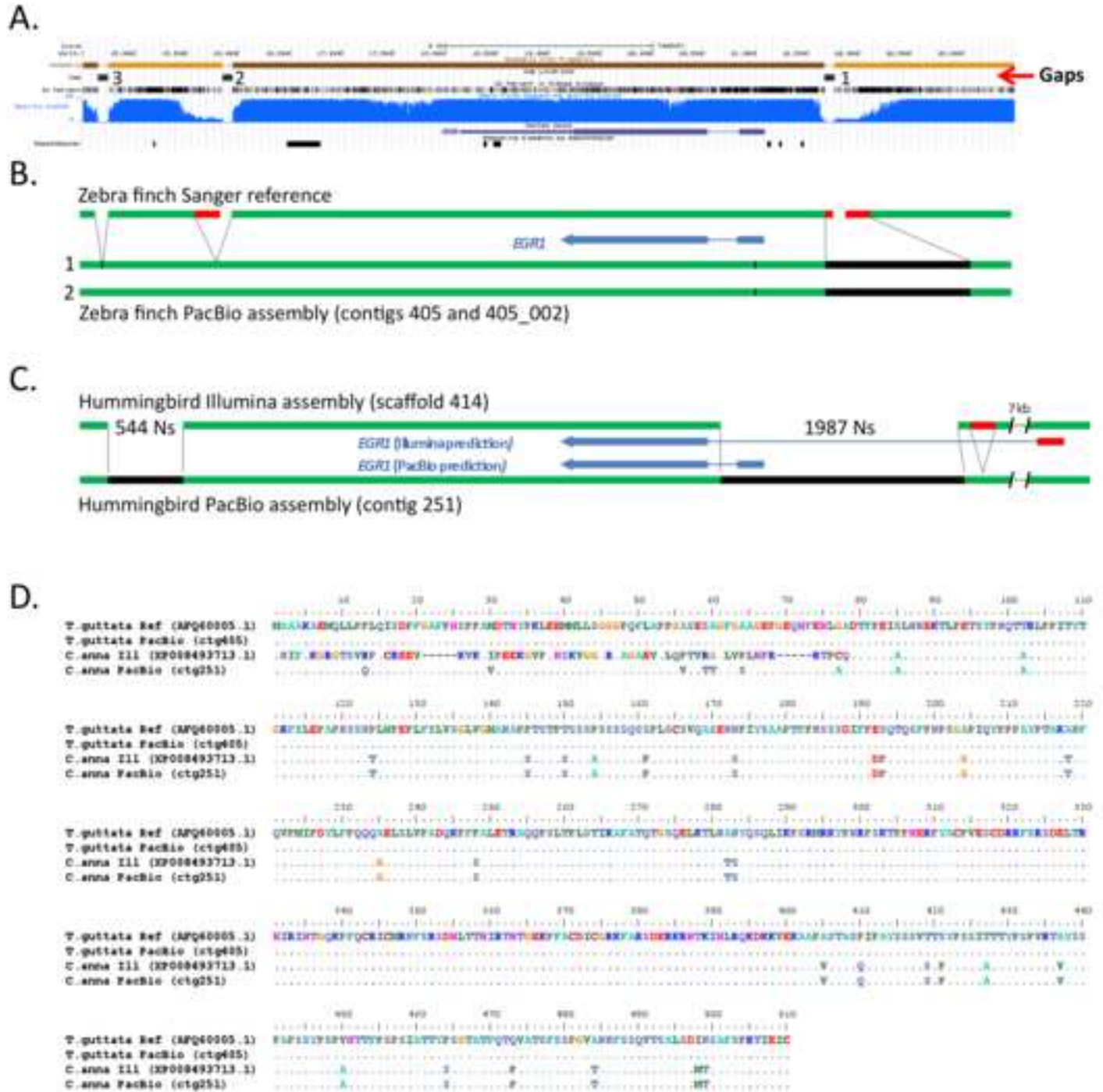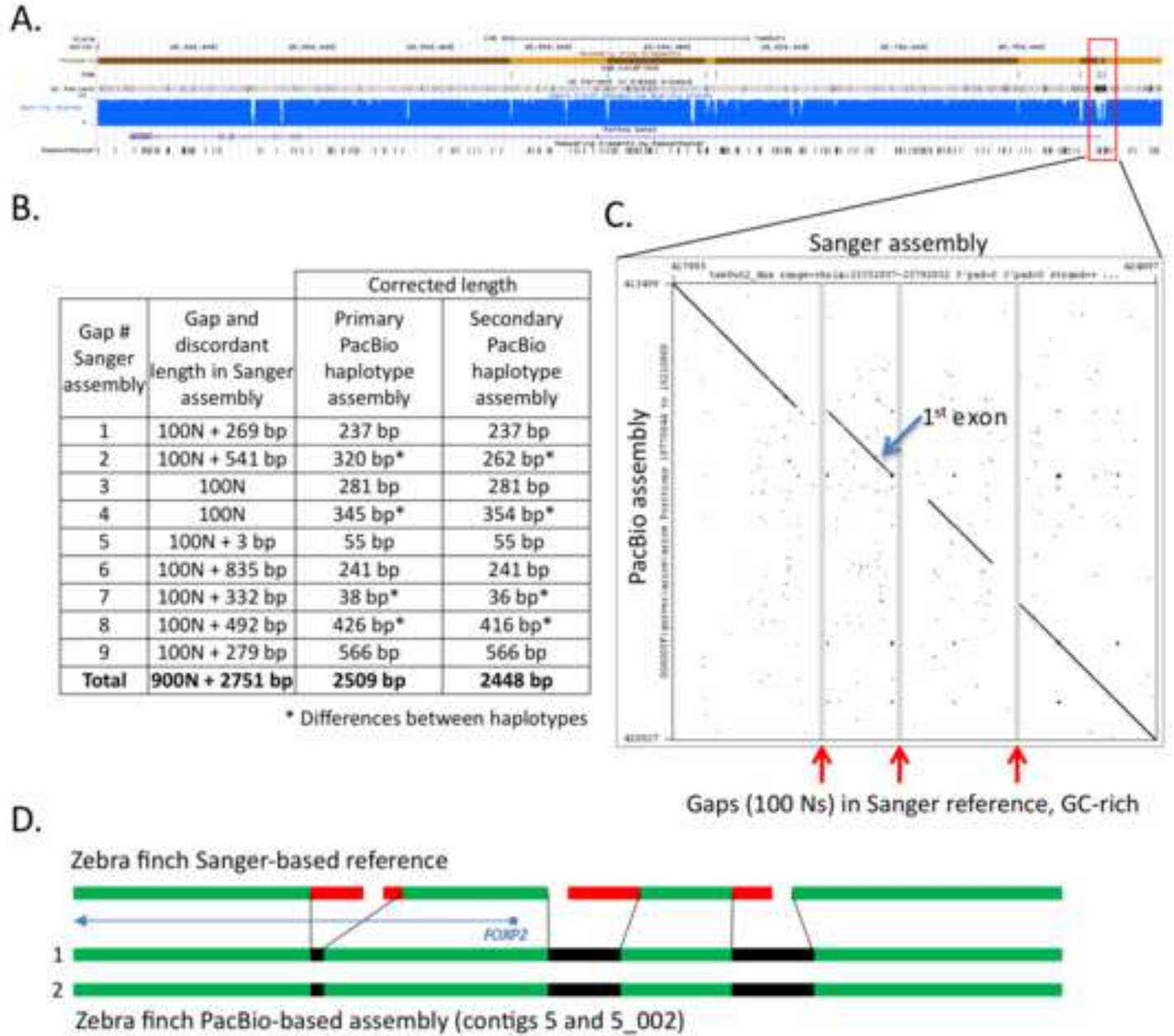Zebra finch PacBio-based assembly (contigs 5 and 5_002)

# Figure 5

Figure 6

# Figure 6

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