# Brain transcriptome sequencing and assembly of three songbird model systems for the study of social behavior

Emberizid sparrows (emberizidae) have played a prominent role in the study of avian vocal communication and social behavior. We present here brain transcriptomes for three emberizid model systems, song sparrow *Melospiza melodia*, white-throated sparrow *Zonotrichia albicollis*, and Gambel's white-crowned sparrow *Zonotrichia leucophrys gambelii*. Each of the assemblies covered fully or in part, over 89% of the previously annotated protein coding genes in the zebra finch *Taeniopygia guttata*, with 16,846, 15,805, and 16,646 unique BLAST hits in song, white-throated and white-crowned sparrows, respectively. As in previous studies, we find tissue of origin (auditory forebrain versus hypothalamus and whole brain) as an important determinant of overall expression profile. We also demonstrate the successful isolation of RNA and RNA-sequencing from *post-mortem* samples from building strikes and suggest that such an approach could be useful when traditional sampling opportunities are limited. These transcriptomes will be an important resource for the study of social behavior in birds and for data driven annotation of forthcoming whole genome sequences for these and other bird species.

1  **Brain transcriptome sequencing and assembly of three songbird model systems for the**

2  **study of social behavior**

3  Christopher N. Balakrishnan[1, *], Motoko Mukai[2,3], Rusty A. Gonser,[4] John C. Wingfield,[3] Sarah E.

4  London[5], Elaina M. Tuttle[4], and David F. Clayton[6]

5  [1]Department of Biology, East Carolina University, Greenville, North Carolina, USA

6  [2] Department of Food Science, College of Agriculture and Life Sciences, Cornell University,

7  Ithaca, New York, USA

8  [3]Department of Neurobiology, Physiology and Behavior, University of California, Davis,

9  California, USA

10  [4]Department of Biology and The Center for Genomic Advocacy (TCGA), Indiana State

11  University, Terre Haute, Indiana, USA

12  [5]Department of Psychology, University of Chicago, Chicago, Illinois, USA

13  [6] Division of Biological & Experimental Psychology, School of Biological and Chemical

14  Sciences, Queen Mary University of London, London, UK

15  *Author for correspondence:

16  Christopher N. Balakrishnan

17  East Carolina University

18  Howell Science Complex

19  Greenville, NC 27858

20  balakrishnanc@ecu.edu

21  252 328 2910

## Introduction

22

23       The comparative method, broadly speaking, is a powerful approach for understanding

24   adaptations including behavior and central control of physiological responses to environmental

25   change. Natural variation in behavior among species has been used in various taxonomic groups

26   to begin to unravel the molecular underpinnings of animal social behavior. Among these

27   comparative studies of behavior, different strategies and technologies have been deployed in

28   order to gain an understanding of the proximate mechanisms at play. For example, experimental

29   hormonal manipulations and gene sequence comparisons in different species of *Microtus* voles

30   led to insights into the mechanisms of parental care (Young et al. 1999). Similarly, quantitative

31   trait locus (QTL) mapping studies have recently revealed the genetic architecture of burrowing

32   behavior in *Peromyscus* mice (Weber et al. 2013). Phylogenetic analyses of rates of molecular

33   evolution based on transcriptomes in eusocial and solitary bees has also led to insights into

34   potential underpinnings of social behavior variation (Woodard et al. 2011).

35       Songbirds, or oscine passerines, comprise roughly half of avian diversity and also serve as

36   important models for the study of social behavior. Arguably the most prominent of the songbird

37   species for behavioral research is the zebra finch *Taeniopygia guttata*, which now boasts a full

38   suite of genomic and molecular tools including a complete genome sequence (Warren et al.

39   2010), RNA-seq based mRNA (Warren et al. 2010; Balakrishnan et al. 2012) and microRNA data

40   (Gunaratne et al. 2011; Luo et al. 2012), transgenics (Agate et al. 2009) and cell lines  (Itoh &

41   Arnold 2011; Balakrishnan et al. 2012). A key strength of songbirds as a model system, however,

42   has always been the behavioral complexity and diversity of songbirds as a group  (Beecher &

43   Brenowitz 2005; Brenowitz & Beecher 2005; Clayton et al. 2009).

44       Among songbirds, many comparative neurobiological studies have focused on three species

45   of new world sparrows (emberizidae). Before the zebra finch assumed its role as a model system

46   for vocal learning, Peter Marler and colleagues had demonstrated age-limited song learning and

47 cultural transmission of song dialects in the white-crowned sparrow, *Zonotrichia leucophrys*

48 (Marler & Tamura 1964). There is also a striking behavioral polymorphism in which some

49 subspecies, such as Gambel's white-crowned sparrow *Z. l. gambelii*, are migratory, living in large

50 non-territorial flocks during non-breeding seasons, whereas other subspecies are non-migratory

51 and are territorial throughout the year (DeWolfe et al. 1989). White-throated sparrows

52 *Zonotrichia albicollis* also show polymorphism in behavior but in this case, the polymorphism is

53 known to be caused by a large chromosomal rearrangement on chromosome 2 (Thorneycroft

54 1966; Thorneycroft 1975). Tan morph individuals are homozygotic for the metacentric form of

55 the chromosome whereas white morphs are almost always heterozygous. In addition to

56 coloration, the two morphs differ in a suite of behaviors including increased aggression and

57 promiscuity and decreased parental care in birds of the white morph (Knapton and Falls 1983,

58 Collins & Houtman 1999; Tuttle 2003). Male song sparrows *Melospiza melodia* are distinctive in

59 that they are territorial during both the breeding season (summer) and much of the non- breeding

60 season (autumn and winter) (Wingfield & Hahn 1994; Mukai et al. 2009). Different hormonal

61 mechanisms, however, appear to underlie this similar behavioral phenotype with increased

62 plasma testosterone levels driving intensity and persistence of aggression during breeding, but not

63 at other times of year (Wingfield 1994; Wingfield & Soma 2002). With this comparative

64 perspective in mind, we have generated brain transcriptomes for these three historically important

65 emberizid songbird models for the study of social behavior: white-throated sparrow,  Gambel's

66 white-crowned sparrow, and song sparrow.


67 **Methods**

68 *Sample Collection*

69        Samples for each of the three species were collected for diverse research purposes of the

70 laboratories involved, so sampling strategy for each species was unique. Animal procedures were

71    approved by the Institutional Animal Care and Use Committees of the University of California,

72    Davis (protocol 07-13208) and the University of Illinois (protocol 11062) and were conducted in

73    accordance with the NIH Guide for the Principles of Animal Care.

74        *White-throated Sparrow*: During migration, white-throated sparrows and other birds are

75    often killed in collisions with buildings. We took advantage of this unfortunate fact by sampling

76    white-throated sparrows that had been opportunistically collected following night migration and

77    collision into McCormick Place, Chicago, IL. Birds that had been killed overnight were collected

78    first thing in the morning beginning at dawn by David Willard, Collection Manager - Birds, Field

79    Museum of Natural History, Chicago, IL. Specimens used in this study were collected during the

80    spring migration in 2010.  Each specimen was immediately vouchered at the Field Museum

81    where they were dissected to determine sex.  Whole brain tissue was stored in RNA-later (Life

82    Technologies, Carlsbad, CA).  Prior to analysis we determined the morph of each bird sampled

83    using a modification of Michopoulos et al. (2007), which is based on the identification of a

84    morph-specific SNP present in the vasoactive intestinal peptide (VIP) gene.  We modified the

85    protocol by using labelled PCR primers, so that the amplification products could be analyzed on

86    an ABI PRISM Genetic Analyzer (Life Technologies). For RNA sequencing we used the brains

87    from six males, three white and three tan.

88        *Gambel's white-crowned sparrow:* We captured Gambel's white-crowned sparrows within

89    the University of California, Davis campus in February 2008, using Potter traps baited with seed,

90    and determined their sex using published PCR methods (Griffiths et al. 1998).  After two weeks

91    of acclimation in captivity we anesthetized 12 male birds with with isoflurane, decapitated them

92    and collected the whole hypothalamus from each bird. After dissection we immediately froze the

93    samples in liquid nitrogen. Fieldwork in California was covered by the US Fish and Wildlife

94    permit (MB713321-0) and State of California permit (SC-004400).

95        *Song sparrow*: Between July and August 2011 we captured seven male song sparrows using

96   song playbacks from behind a mist net. We conducted fieldwork at two locations in central

97   Illinois: "Phillips Tract" (40 07' 54.74" N 88 08' 39.66" W) and Vermillion River Observatory (40

98   03' 50.79" N 87 33' 30.30" W). We euthanized the birds immediately following capture in the net,

99   and then dissected auditory forebrain tissue (auditory lobule, or AL).  AL is a composite brain

100  area including the caudomedial nidopallium (NCM), caudomedial mesopallium (CMM) and Field

101  L and can be readily dissected following bisection of the brain along the midline (Chen and

102  Clayton 2004). We immediately froze the specimens on dry ice. Flat skins of collected song

103  sparrows have been accessioned in the Illinois Natural History Survey, Urbana, Illinois. We

104  conducted fieldwork in Illinois under US Fish and Wildlife Service Permit SCCL-41077A.


105  *RNA Extraction, Library Preparation and Sequencing*

106      *White-throated Sparrow and Song Sparrow*: In order to broadly describe the brain-

107  expressed transcriptome of the white-throated sparrow, we extracted RNA from whole brain. We

108  homogenized the entire brain in Tri-Reagent (Molecular Research Center, Cincinnati, OH) for

109  RNA purification and extracted total RNA following the Tri-Reagent protocol. We then DNase

110  treated (Qiagen, Valencia CA) the total RNA to remove any genomic DNA contamination, and

111  further purified the resulting RNA using Qiagen RNeasy columns. We assessed the purified total

112  RNA for quality using an Agilent Bioanalyzer (Agilent Technologies, Wilmington, DE). Library

113  preparation and sequencing were done at the University of Illinois Roy J. Carver Biotechnology

114  Center. The RNAseq libraries were constructed with the Illumina TruSeq RNA Sample Prep Kit

115  (Illumina, San Diego, CA). Briefly, polyA+ messenger RNA was selected from 1ug of total RNA

116  and chemically fragmented. First-strand cDNA was synthesized with a random hexamer and

117  SuperScript II (ThermoFisher, Waltham, MA). After second-strand synthesis, the double-stranded

118  DNA was blunt-ended, 3'-end A tailed, ligated to barcoded adaptors and amplified with 15 cycles

119  of PCR using Kapa HiFi polymerase (Kapa Biosystems, Woburn, MA). The six barcoded

120   libraries were quantitated with Qubit (ThermoFisher) and the average size was determined on a

121   Bioanalyzer DNA7500 DNA chip (Agilent). The libraries were pooled in equimolar

122   concentration and the pool was quantitated by qPCR on an ABI 7900HT (ThermoFisher).

123   Sequencing was done in a single lane of an Illumina HiSeq2000 using a TruSeq SBS sequencing

124   kit version 3. Fastq files were demultiplexed and generated with the software Casava 1.8.2

125   (Illumina). The same basic procedure was used to sequence the song sparrow except for the fact

126   that we extracted RNA from the dissected AL (rather than whole brain) tissue, and that samples

127   from seven individuals were run in a single lane of paired end (rather than single end)

128   sequencing.

129       *Gambel's White-crowned Sparrow*: We extracted total RNA from each hypothalamus using

130   TRIzol reagent (Life Technologies) followed by RNA cleanup using Qiagen RNeasy Mini Kits.

131   We then pooled RNA samples, quantified them using a Nanodrop (ThermoFisher) and ran them

132   on a Bioanalyzer for quality control (RIN = 8.5). We used this pooled RNA sample to generate a

133   mRNA-seq library of 400 bp size with a mRNA-Seq 8 sample prep kit (Illumina) following

134   manufacturer's protocol with slight modifications. We began by isolating mRNA using oligo(dT)

135   and then fragmented it using divalent cations under elevated temperature. We then reverse

136   transcribed the RNA into cDNA using random primers, modified and ligated with GEN PE

137   adapters. We ran the resulting cDNA on an agarose gel, excised a 400 bp band and enriched the

138   library with 15 cycles of PCR. We validated the final library using a Bioanalyzer and confirmed a

139   distinct band at approximately 400 bp. Pair-end sequencing (100 bp x 2) was performed by the

140   Genome Center DNA Technologies Core at the University of California, Davis, using an Illumina

141   HiSeq 2000 and TruSeq SBS kit version 2.

142       *Zebra Finch:* To provide a benchmark for comparison, we compared our newly collected

143   data with previously published data from zebra finches *Taeniopygia guttata* (Balakrishnan et al.

144   2012, GenBank Accession: SRX493920- SRX493922). These data were derived from RNA

145 extracted from the AL of female zebra finches. The three libraries were derived from pools of 10

146 female finches each, and sequenced on an Illumina Genome Analyser and processed with

147 Illumina pipeline 1.6.


148 *Transcriptome Assembly, Annotation and Assessment*

149      We checked overall sequence quality using FastQC

150 (http://www.bioinformatics.babraham.ac.uk/projects/fastqc/) and trimmed reads using

151 ConDeTriV2.2 (Smeds & Kunstner 2011). We used default settings for trimming except for the

152 high quality (hq) threshold which was set to 20 and lfrac, the maximum fraction of reads with

153 quality < 10, which was set to 0.2. The lfrac parameter allows for trimming, rather than complete

154 removal, of reads with low quality ends.

155      We used the Trinity (version r20131110) assembler (Grabherr et al. 2011) to generate *de*

156 *novo* assemblies for each species. For white-throated sparrow we assembled the reads for the two

157 color morphs both separately and combined. Assembling the reads separately was reasonable

158 given evidence of sequence divergence within the inversion (Thomas et al. 2008) and assembling

159 the reads together was reasonable to improve coverage outside such areas. We used default

160 settings in Trinity besides those specific to our computing system (we generally used 24 CPUs

161 and allowed for 100G of memory). We used TransDecoder (included in the Trinity package) to

162 identify open reading frames (ORFs) in our predicted transcripts.

163      We assessed the quality of our assembly by estimating N50 and average transcript length.

164 The shortcomings of such metrics for transcriptome assessment have been described (O'Neil and

165 Emrich 2013) and we use them here primarily to facilitate comparison with previously published

166 studies. To provide further insight into assembly quality, we also assessed 5' to 3' gene model

167 coverage relative to annotated zebra finch genes (see details below) and quantified the number of

168 transcripts containing both start and stop codons using the annotation information provided by

169    TransDecoder ("type:complete" in the fastq header).

170    We used BLAST (Altschul et al. 1990) searches against a database of Ensembl (release 74)

171    zebra finch transcripts to annotate our ORF-containing transcripts. Functional description of

172    annotated transcripts was conducted using Gene Ontology, and statistical over and under

173    representation was tested using CORNA software (Wu & Watson 2009) and Fisher's Exact Tests

174    with *p* values adjusted for multiple testing (Benjamini & Hochberg 1995). For each assembly we

175    tested our identified set of putative zebra finch orthologs relative to the full population of

176    Ensembl transcripts.


177    *Gene Expression and Read-Mapping Profiling*

178    In order to compare read mapping and gene expression profiles across libraries, we mapped

179    RNA-seq reads to the zebra finch whole genome assembly (2.3.4) using Stampy, a read mapper

180    tailored for divergent reads relative to the reference genome (Lunter & Goodson 2011). We

181    mapped reads for all six individual white-throated sparrows, three of the seven song sparrows,

182    and the pooled white-crowned sparrow using default settings but with the substitution rate set to

183    0.05 to accommodate sequence divergence. In addition, we mapped reads from previously

184    published zebra finch auditory forebrain reads (Balakrishnan et al. 2012) using substitution rate =

185    0.01.

186    To quantify gene expression, we used htseq-count (Anders et al. 2014) and tallied reads

187    relative to Ensembl gene models and normalized them read counts using the regularized log

188    transformation in DE-Seq2 (Anders & Huber 2010). Expression profiles were then visualized by

189    Euclidean distance based clustering and principal components analysis (PCA) using heatmap.2 in

190    the gplots R package, and the plotPCA function in DE-Seq2. We then also used the geneBody.py

191    script within the RseqC package (Wang et al. 2012) to describe read coverage across gene models

192    and to test specifically for a 3' bias in transcript coverage in *post-mortem* samples.

193 **Results & Discussion**

194 *RNA extraction and sequencing*

195       Despite collecting tissues for the white-throated sparrow opportunistically from building

196 strikes, we were able to extract reasonably high quality RNA from all samples (Fig. 1). This

197 finding suggests that *post-mortem* collected birds can be used as a viable source of RNA for

198 transcriptome sequencing. From a total of twelve samples, we selected a set of six (three per

199 morph) with Bioanalyzer RNA integrity numbers (RIN) above 7 (10-083 (7.2), 10-092 (7.2), 10-

200 093 (7.7) and 10-118 (8.5), 10-124 (8.0) and 10-308 (7.9)). Samples for sequencing were also

201 chosen such that tan and white morphs were collected at the same time of year (spring migration

202 2010). By chance, our tan samples had higher average RINs than the white morph samples did

203 (Fig. 1). RNA from the other two species were of good quality and met Illumina's standard QC

204 benchmark of RIN > 8.  All of our sequencing runs yielded high quality sequence data. After

205 fairly stringent quality trimming, we retained over 89% of the initial nucleotides sequenced

206 (Table 1). Raw RNA seq reads have been deposited to the GenBank Short Read Archive under

207 accession numbers SRX342288-SRX342293, SRX493875- SRX493882, and SRX514971.


208 *Transcriptome Assembly and Annotation*

209       We reconstructed a large number of transcripts (> 95,000) and open reading frame (ORF)

210 containing transcripts (>54,000) in all of our assemblies, exceeding the likely number of coding

211 genes (Table 2). These transcripts reflect a combination of partial transcripts, alternative

212 isoforms, allelic variants, and noncoding transcripts. We were able to generate high quality

213 transcriptomes based on N50 and average transcript length (Table 2). N50s for the assemblies

214 were 1,942 for the white morphed white-throated sparrow, 2,557 for the tan morph, 3,415 for

215 Gambel's white-crowned sparrow and 4,072 for the song sparrow (Table 2). For the song

216   sparrow, this is an improvement over a recent 454-based transcriptome (N50=482; Srivastava *et*

217   *al.* 2012). As expected, N50 in general improved with increased sequencing depth (with paired

218   end data sets benefitting from both the reads being paired and having more reads). One exception

219   to this rule was in the white-throated sparrow, where combining reads from the two morphs

220   actually generated a worse assembly in terms of N50 relative to the "tan morph only" assembly

221   (combined N50=2,284, tan only N50 = 2,557). Tan morph individuals are homozygous for a large

222   structural polymorphism spanning much of chromosome 2 whereas white morph individuals are

223   heterozygous. Recombination within the inversion is suppressed, allowing genetic divergence in

224   this region (Thomas et al. 2008), and potentially explaining the drop in N50 in the combined

225   assembly. For the purposes of annotation of the white-throated sparrow we therefore used the two

226   morph-specific assemblies, merging them after the assembly process.

227        Although N50s were generally high, the white-throated sparrow assemblies, which were

228   based on smaller, single-end datasets and *post-mortem* samples, had the lowest scores. This effect

229   was even more dramatic when assemblies were assessed in terms of the number of complete

230   transcripts possessing both a start and stop codon. Gambel's white-crowned, song, and white-

231   throated sparrow transcriptomes contained 115,515, 79,451, and 24,388 complete transcripts,

232   respectively (Table 2).  Because the white-throated sparrow samples were collected *post-mortem*

233   and had the fewest reads, we cannot determine whether *post-mortem* sampling itself influenced

234   assembly quality. Given the relatively high quality (RINs) of the white-throated sparrow RNA,

235   however, it is more likely that the reduced quality of the assembly is a result of it being generated

236   from a smaller dataset.

237        For white-throated sparrow we were able to find predicted transcripts with significant

238   BLAST hits to 15,805 zebra finch genes (89% of Ensembl annotated zebra finch genes), whereas

239   for song sparrow we found 16,846 (94%) and Gambel's white-crowned sparrow 16,646 (93%).

240   Therefore, in terms of unique BLAST hits, the song sparrow and Gambel's white-crowned

241 assemblies were also better than that of the white-throated sparrows. All three assemblies,

242 however, cover a large proportion of known genes and represent an improvement of over recent

243 454-based bird transcriptomes (e.g., violet-eared waxbill, 11,084 genes, Balakrishnan et al.

244 2013).

245      We evaluated and compared the general composition of genes present in each of the new

246 assemblies by performing a Gene Ontology (GO) analysis, using the GO annotation of the

247 complete zebra finch genome as the point of reference for the statistical tests of enrichment

248 (Table 3). All three datasets shared a number of similarities, including significant enrichment for

249 eight GO categories ("cytoplasm", "intracellular', "mitochondrion", "nucleic acid binding",

250 "nucleolus", "protein binding", "protein phosphorylation" and "transferase activity") and under-

251 representation of six categories ("cytokine activity", "DNA integration", "extracellular region",

252 "hormone activity", "immune response" and MCH Class I protein complex"). The under-

253 represented categories may in part reflect the well-described pattern of limited immune activity,

254 or "immune privilege" in the brain (Galea et al. 2007). As in previous studies of avian brain gene

255 expression, however, we did see some evidence of expression of the MHC Class I gene itself

256 (Ekblom et al. 2010; Balakrishnan et al. 2013).

257      Interestingly, genes annotated with the GO term "olfactory receptor activity" are well

258 represented in all three assemblies (where observed/expected were 165/150 in white-throated

259 sparrows, 165/156 in song sparrow, and 165/158 in Gambel's white-crowned sparrow, out of a

260 total of 168 annotated genes). This was notable as a previous 454-based whole brain

261 transcriptome of another songbird did not detect any olfactory receptor genes at all (Balakrishnan

262 *et al.* 2013). The detection of such genes here suggests that the increased sequencing depth

263 provided by the Illumina platform has aided in this regard. Despite the generally tissue-restricted

264 distribution of olfactory receptor expression, we were able to pick up these genes in all of our

265 tissue samples irrespective of the brain region targeted. High depth RNA-sequencing data

266    including those presented here will therefore be useful for annotating these diverse olfactory

267    receptor transcripts.

268         Thirteen other GO terms were significantly under-represented only in the white-throated

269    sparrow assembly (Table 4).  These categories were relatively well-represented in the other two

270    sparrow assemblies (Table 4) and included "visual function", "G-protein coupled receptor

271    activity", and "neurotransmitter transport". The white-throated sparrow assembly differs from the

272    others in several factors that could contribute to this difference in gene composition, including

273    tissue of origin (whole brain, versus auditory forebrain or hypothalamus), physiological condition

274    (spring migration, versus breeding season or captive/wintering) and *post-mortem* tissue

275    collection.


276    *Transcriptome Coverage of Zebra Finch Gene Models*

277         We performed further analysis of read distribution and the relative abundance of different

278    transcripts in each of the source tissues, by mapping RNAseq reads back to the zebra finch

279    genome reference.  For comparison we also included previously published RNAseq read data

280    from the zebra finch auditory forebrain (Balakrishnan et al. 2013). White-throated sparrow reads

281    mapped at a lower rate (average = 83% of reads mapped) than reads from Gambel's white-

282    crowned sparrow (90%), song sparrow (94%) and zebra finch (93%). Among the reads that did

283    map to the genome, however, all of the species were similar in showing a large proportion of

284    reads (53.2 +/- 3.6%) mapping outside of currently defined zebra finch genes, suggesting

285    extensive transcription outside of known genes.

286         Based on this read mapping we were able to assess coverage of annotated genes. This was

287    important given our *post-mortem* sampling of white-throated sparrows.  In highly degraded

288    samples we would expect to see a strong 3' bias in gene coverage. RNA quality as measured by

289    RIN was only slightly lower in white-throated sparrow samples and thus, we found that 3' bias

290   was similar across all of our samples (Fig. 2). This finding further suggests that RNA degradation

291   may not be the primary factor associated with the lower assembly quality in the white-throated

292   sparrow assembly.

293        Cheviron et al. (2011) documented the time course of RNA degradation *post-mortem*, and

294   also suggest that such samples can provide a useful source of RNA, even though such specimens

295   are often overlooked.  Similarly, a recent RNA-sequencing study of pinnipeds successfully used

296   *post-mortem* samples (Hoffman et al. 2013). Although clearly not an ideal strategy for studies

297   aimed at quantifying gene expression, the use of recently killed samples is viable strategy for

298   initial transcriptome description, and in our study gave access to a large portion of the

299   transcriptome. This approach could be particularly useful for rare species where collection of

300   fresh specimens is impossible.


301   *Impacts of Ancestry, Tissue of Origin, and Library Preparation on Expression Profile*

302        We used cluster analysis to compare the broad structure of gene expression in the different

303   samples, recognizing that the samples differed in multiple dimensions (i.e., species, sex, brain

304   region, physiological condition, collection method, sequencing method). If species or sex were

305   the dominant factors driving the differences in gene expression patterns, one would expect to see

306   a clustering pattern with zebra finch as the most divergent profile (Fig. 3a). Similarly, if the

307   sequencing facility and platform were dominant technical factors one would expect to see either

308   the zebra finch or the white crowned-sparrow as most divergent (Fig 3b). However, the zebra

309   finch samples clustered closely with the song sparrow samples taken from the same brain region

310   (auditory forebrain), with the white throated sparrow samples from the whole brain clustering

311   together as most divergent, and the Gambel's white-crowned sparrow samples from

312   hypothalamus in between (Fig 3c, Fig 4).  This echoes previous findings that brain region is a

313   major determinant of gene expression pattern in songbirds (Replogle et al. 2008; Drnevich et al.

314   2012). Both euclidean distance-based clustering and PCA also highlight the fact that zebra

315   finches, which were sacrificed in captivity and sequenced in pools of ten, had much reduced

316   variance in expression profile relative to our non-pooled, field-collected white-throated sparrow

317   and song sparrow samples (Fig. 4).


318   **Conclusion**

319   Transcriptome assemblies are a valuable resource, particularly for species without reference

320   genomes, providing access to a large proportion of the coding and noncoding expressed genome.

321   For taxa with genomes, or with genomes in progress, transcriptome data provides empirical (as

322   opposed to model based) information on transcript structures including alternative isoforms that

323   are not well-annotated in most species. We have presented here neuro-transcriptomic data for

324   three important model species for the study of social behavior and neurobiology building on a

325   growing body of such data (e.g., Balakrishnan et al. 2013, Ekblom et al. 2014; MacManes &

326   Lacey 2012; Moghadam *et al.* 2013).

332    **References**

333    Agate RJ, Scott BB, Haripal B, Lois C, Nottebohm F (2009) Transgenic songbirds offer an

334        opportunity to develop a genetic model for vocal learning. *Proceedings of the National*

335        *Academy of Sciences of the United States of America*, **106**, 17963–17967.

336    Altschul SF, Gish W, Miller W, Myers EW, LIPMAN DJ (1990) Basic local alignment search

337        tool. *Journal of Molecular Biology*, **215**, 403–410.

338    Anders S, Huber W (2010) Differential expression analysis for sequence count data. *Genome*

339        *Biology*, **11**, R106.

340    Anders S, Pyl PT, Huber W (2014) *HTSeq — A Python framework to work with high-throughput*

341        *sequencing data*. *bioRxiv preprint*.

342    Balakrishnan CN, Lin Y-C, London SE, Clayton DF (2012) RNA-seq transcriptome analysis of

343        male and female zebra finch cell lines. *Genomics,* **100**, 363–369.

344    Balakrishnan C, N., Chapus C, Brewer M, S., Clayton D, F. (2013) Brain transcriptome of the

345        violet-eared waxbill *Uraeginthus granatina* and recent evolution in the songbird genome.

346        *Open Biology*, **3**, 130063.

347    Beecher MD, Brenowitz EA (2005) Functional aspects of song learning in songbirds. *Trends in*

348        *Ecology & Evolution*, **20**, 143–149.

349    Benjamini Y, Hochberg Y (1995) Controlling the false discovery rate - a practical and powerful

350        approach to multiple testing. *Journal of the Royal Statistical Society Series B-Methodological*,

351        **57**, 289–300.

352    Brenowitz EA, Beecher MD (2005) Song learning in birds: diversity and plasticity, opportunities

353        and challenges. *Trends in Neurosciences* **28**, 127–132.

354    Cheng HY, Clayton DF (2004) Activation and habituation of extracellular signal-regulated kinase

355     phosphorylation in zebra finch auditory forebrain during song presentation. *Journal of*

356     *Neuroscience* **24**(34):7503–7513.

357     Cheviron ZA, Carling MD, Brumfield RT (2011) Effects of postmortem interval and preservation

358     method on rna isolated from field-preserved avian tissues. *Condor*, **113**, 483–489.

359     Clayton D, F., Balakrishnan C, N., London S, E. (2009) Integrating Genomes, Brain and

360     Behavior in the Study of Songbirds. *Current Biology*, **19**, R865–R873.

361     Collins CE, Houtman AM (1999) Tan and white color morphs of White-throated Sparrows differ

362     in their non-song vocal responses to territorial intrusion. *Condor*, **101**, 842–845.

363     DeWolfe BB, Baptista LF, Petrinovich L (1989) Song development and territory establishment in

364     Nuttals White-Crowned Sparrows. *Condor*, **91**, 397–407.

365     Drnevich J, Replogle KL, Lovell P et al. (2012) Impact of experience-dependent and

366     -independent factors on gene expression in songbird brain. *Proceedings of the National*

367     *Academy of Sciences of the United States of America*, **109**, 17245–17252.

368     Ekblom R, Balakrishnan CN, Burke T, Slate J (2010) Digital gene expression analysis of the

369     zebra finch genome. *BMC Genomics*, **11**, 219.

370     Ekblom, R, Wennekes, P, Horsburgh, GJ, Burke T. 2014 Characterization of the house sparrow

371     (*Passer domesticus*) transcriptome: a resource for molecular ecology and immunogenetics.

372     *Molecular Ecology Resources,* 14(3):636p646.

373     Galea I, Bechmann I, Perry VH (2007) What is immune privilege (not)? *Trends in Immunology*,

374     **28**, 12–18.

375     Goodson JL, Kelly AM, Kingsbury MA, Thompson RR (2012) An aggression-specific cell type

376     in the anterior hypothalamus of finches. *Proceedings of the National Academy of Sciences of*

377    *the United States of America*, **109**, 13847–13852.

378    Goodson JL, Wang YW (2006) Valence-sensitive neurons exhibit divergent functional profiles in

379    gregarious and asocial species. *Proceedings of the National Academy of Sciences of the United*

380    *States of America*, **103**, 17013–17017.

381    Grabherr MG, Haas BJ, Yassour M et al. (2011) Full-length transcriptome assembly from RNA-

382    Seq data without a reference genome. *Nature Biotechnology*, **29**, 644–U130.

383    Griffiths R, Double MC, Orr K, Dawson RJG (1998) A DNA test to sex most birds. *Molecular*

384    *Ecology*, **7**, 1071–1075.

385    Gunaratne PH, Lin YC, Benham AL et al. (2011) Song exposure regulates known and novel

386    microRNAs in the zebra finch auditory forebrain. *BMC Genomics*, **12**, 277.

387    Hoffman JI, Thorne MAS, Trathan PN, Forcada J (2013) Transcriptome of the dead:

388    characterisation of immune genes and marker development from necropsy samples in a free-

389    ranging marine mammal. *BMC Genomics*, **14**, 52.

390    Itoh Y, Arnold AP (2011) Zebra finch cell lines from naturally occurring tumors. *In Vitro Cellular*

391    *& Developmental Biology-Animal*, **47**, 280–282.

392    Knapton, R.W. & Falls, J.B. 1983. Differences in parental contribution among pair types in the

393    polymorphic white-throated sparrow. *Canadian Journal of Zoology*. 61: 1288-1292.

394    Lunter G, Goodson M (2011) Stampy: A statistical algorithm for sensitive and fast mapping of

395    Illumina sequence reads. *Genome Research*, **21**, 936–939.

396    Luo GZ, Hafner M, Shi ZM et al. (2012) Genome-wide annotation and analysis of zebra finch

397    microRNA repertoire reveal sex-biased expression. *BMC Genomics*, **13**, 727.

398    MacManes MD, Lacey EA (2012) The Social Brain: Transcriptome Assembly and

399    Characterization of the Hippocampus from a Social Subterranean Rodent, the Colonial Tuco-

400    Tuco (Ctenomys sociabilis). *PLoS One*, **7**, e45524.

401    Marler P, Tamura M (1964) Culturally transmitted patterns of vocal behavior in sparrows.

402    *Science*, **146**, 1483–148.

403    Michopoulos, V. Maney, D.L., Morehouse, C.B. & Thomas, J.W. 2007. A genotyping assay to

404    determine plumage morph in the White-throated Sparrow (*Zonotrichia albicollis*). *The Auk*

405    124 No. 4 1330-1335.

406    Moghadam HK, Harrison PW, Zachar G, Szekely T, Mank JE (2013) The plover

407    neurotranscriptome assembly: transcriptomic analysis in an ecological model species without a

408    reference genome. *Molecular Ecology Resources*, **13**, 696–705.

409    Mukai M, Replogle K, Drnevich J et al. (2009) Seasonal Differences of Gene Expression Profiles

410    in Song Sparrow (*Melospiza melodia*) Hypothalamus in Relation to Territorial Aggression.

411    *PLoS One*, **4**, e8182.

412    O'Neil ST, Emrich SJ (2013) Assesing *De Novo* transcriptome assembly metrics for consistency

413    and utility. *BMC Genomics,* **14**, 465.

414    Replogle K, Arnold AP, Ball GF et al. (2008) The Songbird Neurogenomics (SoNG) Initiative:

415    Community-based tools and strategies for study of brain gene function and evolution. *BMC*

416    *Genomics*, **9**, 131.

417    Smeds L, Kunstner A (2011) CONDETRI - A Content Dependent Read Trimmer for Illumina

418    Data. *PLoS One*, **6**, e26314.

419    Srivastava A, Winker K, Shaw TI, Jones KL, Glenn TC (2012) Transcriptome Analysis of a North

420    American Songbird, *Melospiza melodia*. *DNA Research*, **19**, 325–333.

421    Thomas J, W., Caceres M, Lowman J, J. et al. (2008) The chromosomal polymorphism linked to

422    variation in social behavior in the white-throated sparrow (*Zonotrichia albicollis*) is a complex

423     rearrangement and suppressor of recombination. *GENETICS*, **179**, 1455–1468.

424     Thorneycroft HB (1966) Chromosomal polymorphism in white-throated sparrow *Zonotrichia*

425     *albicollis* (Gmelin). *Science*, **154**, 1571–157.

426     Thorneycroft HB (1975) Cytogenetic study of white-throated sparrow, *Zonotrichia albicollis*

427     (Gmelin). *Evolution*, **29**, 611–621.

428     Tuttle EM (2003) Alternative reproductive strategies in the white-throated sparrow: behavioral

429     and genetic evidence. *Behavioral Ecology*, **14**, 425–432.

430     Wang L, Wang S, Li W (2012) RSeQC: quality control of RNA-seq experiments. *Bioinformatics*,

431     **28**, 2184–2185.

432     Warren W, C., Clayton D, F., Ellegren H et al. (2010) The genome of a songbird. *Nature*, **464**,

433     757–762.

434     Weber JN, Peterson BK, Hoekstra HE (2013) Discrete genetic modules are responsible for

435     complex burrow evolution in *Peromyscus* mice. *Nature*, **493**, 402–U145.

436     Wingfield JC (1994) Regulation of territorial behavior in the sedentary song sparrow, *Melospiza*

437     *melodia morphna*. *Hormones and Behavior*, **28**, 1–15.

438     Wingfield JC, Hahn TP (1994) Testosterone and territorial behavior in sedentary and migratory

439     sparrows. *Animal Behaviour*, **47**, 77–89.

440     Wingfield JC, Soma KK (2002) Spring and autumn territoriality in song sparrows: Same

441     behavior, different mechanisms? *Integrative and Comparative Biology*, **42**, 11–20.

442     Woodard SH, Fischman BJ, Venkat A et al. (2011) Genes involved in convergent evolution of

443     eusociality in bees. *Proceedings of the National Academy of Sciences of the United States of*

444     *America*, **108**, 7472–7477.

445     Wu X, Watson M (2009) CORNA: testing gene lists for regulation by microRNAs.

446    *Bioinformatics*, **25**, 832–833.

447    Young LJ, Nilsen R, Waymire KG, MacGregor GR, Insel TR (1999) Increased affiliative

448    response to vasopressin in mice expressing the V-1a receptor from a monogamous vole.

449    *Nature,* **400**, 766–768.

# Figure 1

RNA quality from post-mortem sampled sparrows

Bioanalyzer gel image showing RNA extracted from 12 white-throated sparrows sampled *post-mortem*. RNA integrity numbers (RIN) are given at the bottom and ranged from 6.4 to 8.5. Samples chosen for sequencing are indicated by tan and white circles, representing tan and white morph sparrows, respectively.
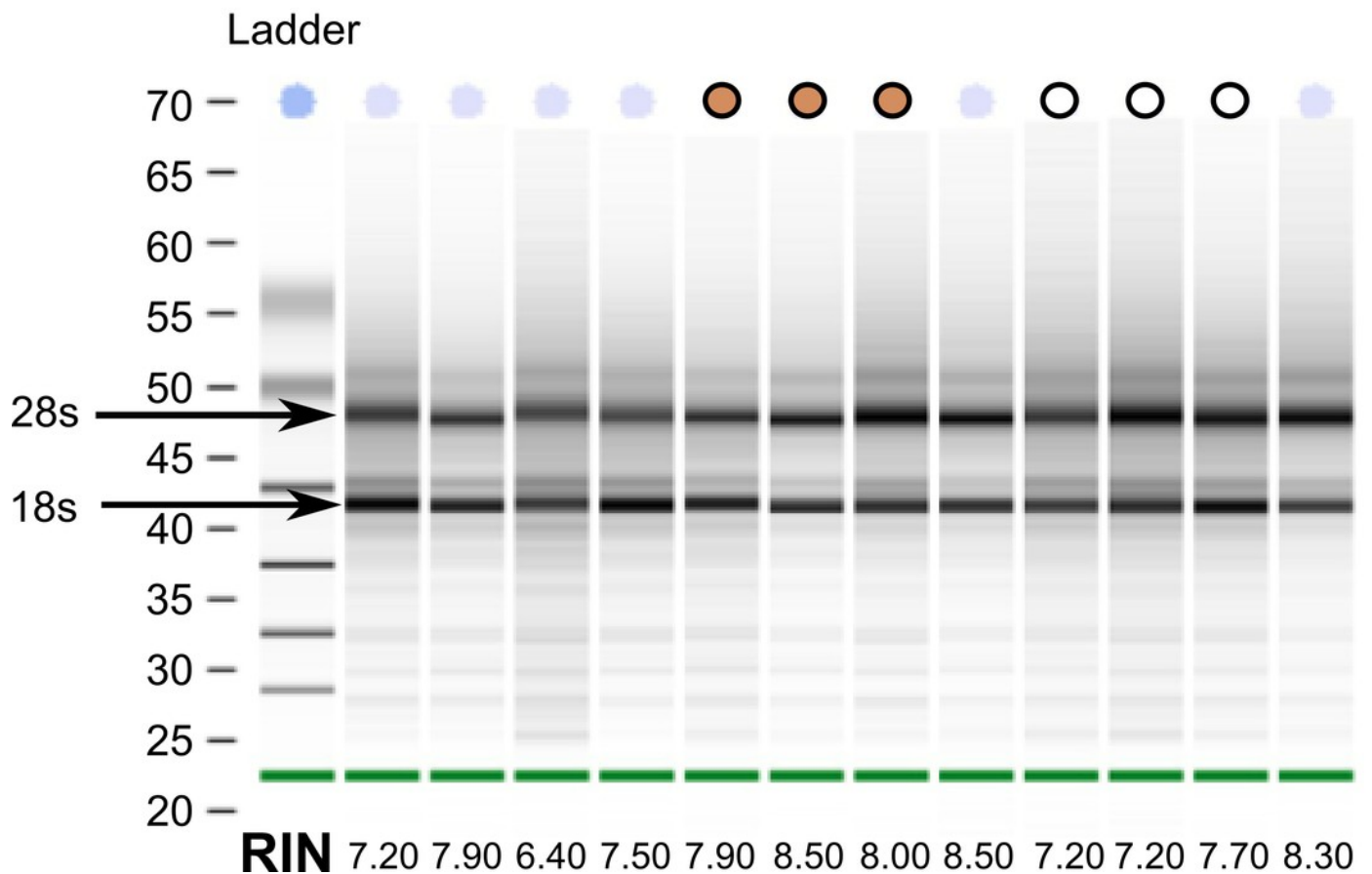
# Figure 2

Coverage of zebra finch gene models by RNA-seq reads

Gene model coverage across all genes based on mapping of reads to the zebra finch genome. Samples collected *post-mortem* from white-throated sparrow show a similar gene coverage profile to freshly collected samples. Zebra finch data included fewer total reads, explaining the lower depth across genes.
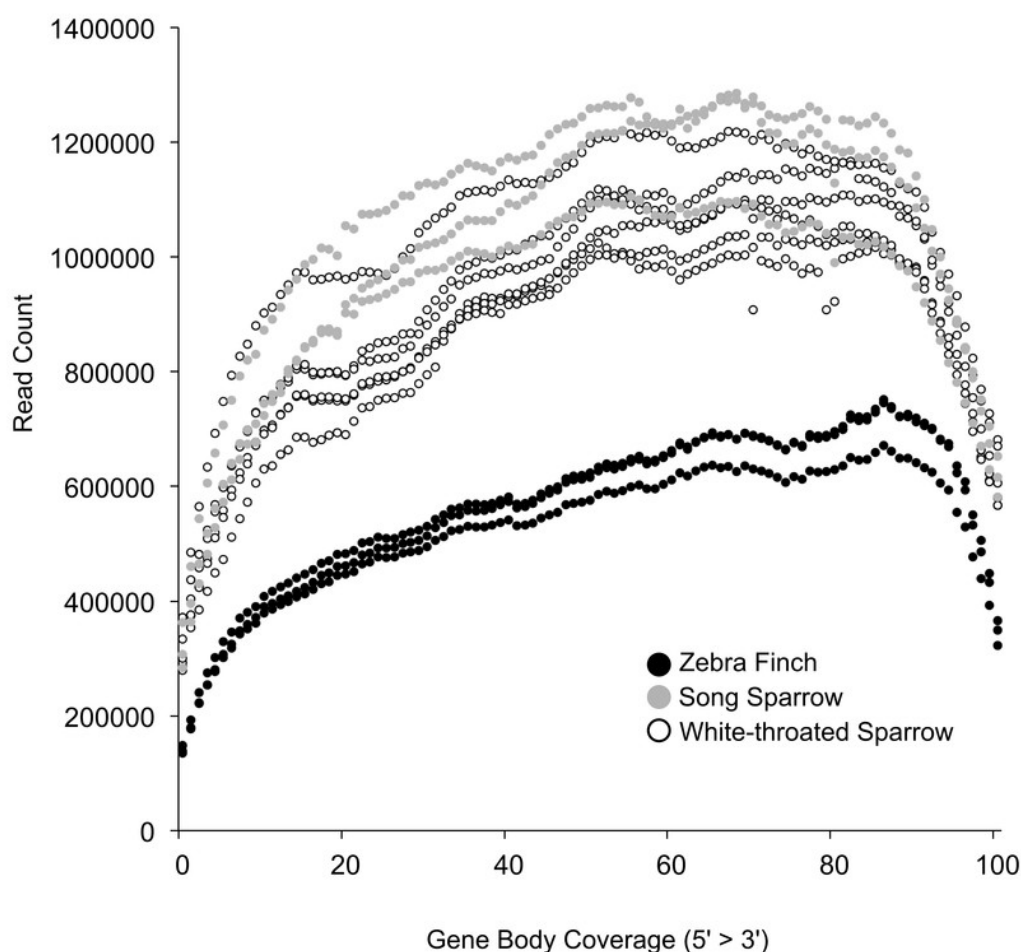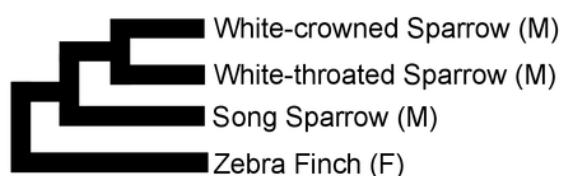
# Figure 3

Alternative expectations for expression profile clustering

Alternative expectations if A) phylogeny or sex B) sequencing platform or library preparation protocols or C) tissue of origin, were the dominant factor underlying expression clustering. Only tissue of origin unites zebra finch and song sparrow samples together as observed in the clustering analysis (Fig. 4).



A) Phylogeny, Sex
- White-crowned Sparrow (M)
- White-throated Sparrow (M)
- Song Sparrow (M)
- Zebra Finch (F)

B) Sequencing platform, strategy, and facility
- Song Sparrow (Hi-Seq, paired-end, Illinois)
- White-throated Sparrow (Hi-Seq, single-end, Illinois)
- White-crowned Sparrow (Hi-Seq, paired-end, Davis)
- Zebra Finch (Genome Analyzer, single-end, Illinois)

C) Tissue
- Zebra Finch (Auditory Lobule)
- Song Sparrow (Auditory Lobule)
- White-crowned Sparrow (Hypothalamus)
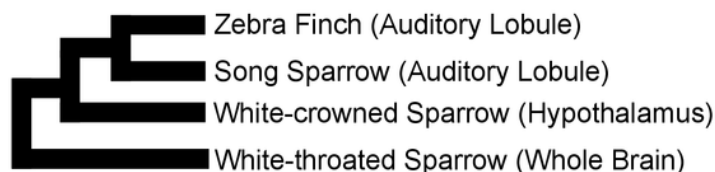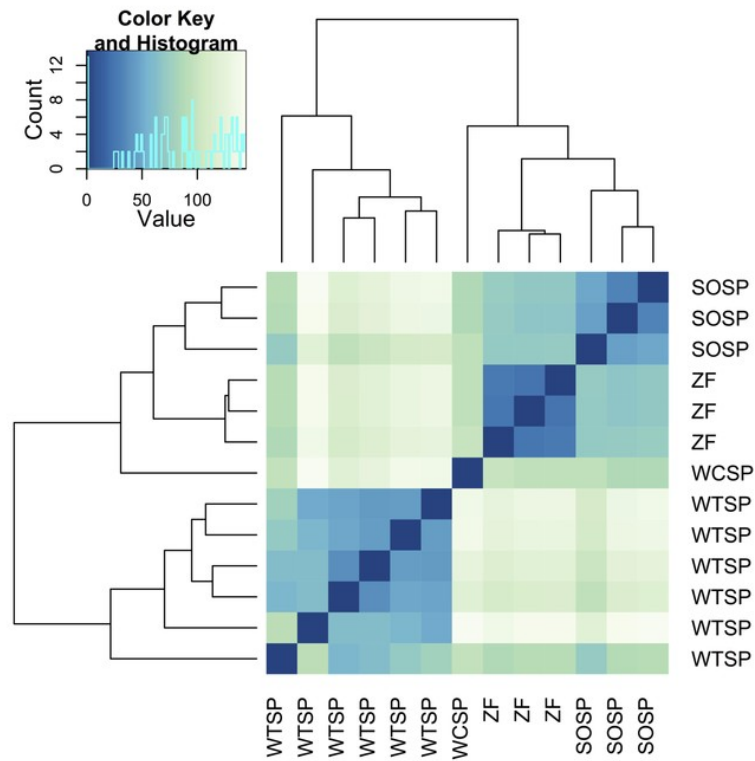- White-throated Sparrow (Whole Brain)

# Figure 4

Clustering of expression profiles from four songbird species

A) Hierarchical clustering and B) Principal components analysis of expression profiles for six white-throated sparrow (WTSP), three song sparrow (SOSP), three zebra finch (ZF) and one white-crowned sparrow libraries. Libraries derived from auditory lobule (AL) tissue cluster (SOSP and ZF) to the exclusion of the others. White-throated sparrow samples, taken from whole brain (rather than forebrain as the other samples are) show divergent and variable profiles. Zebra Finch (ZF) samples collected in captivity and generated from pools of 10 individuals, show much reduced sample variability.
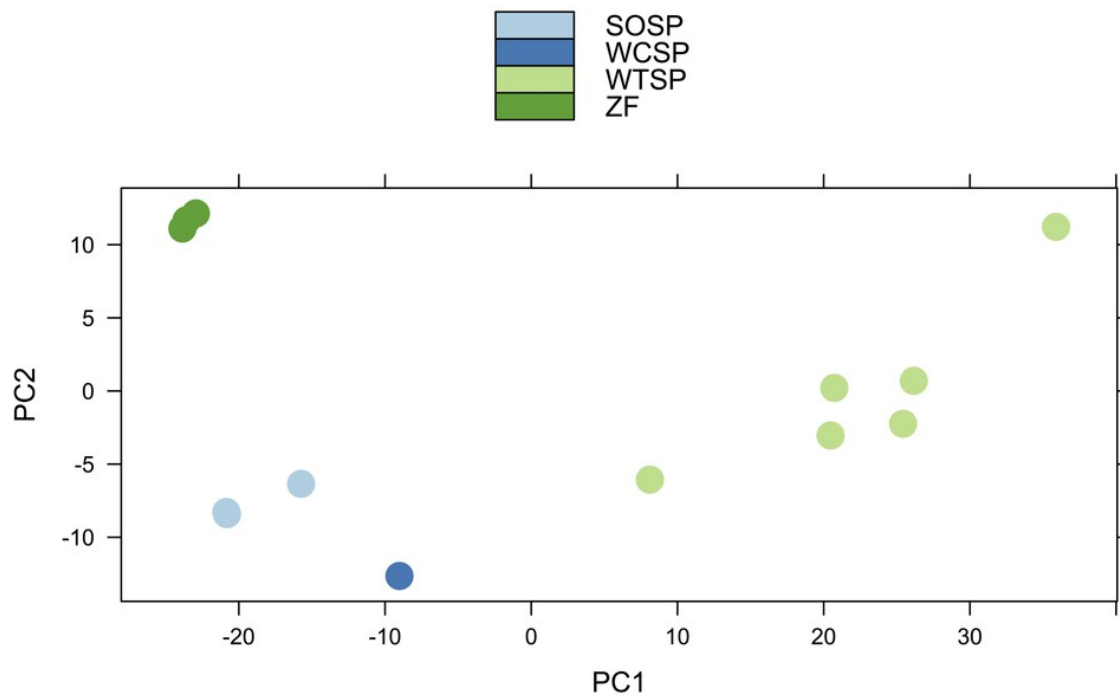
A.



B.

# Table 1(on next page)

RNA-seq dataset

Raw number of reads and bases before and after trimming with ConDeTri.

| Species | Reads Before | Bases Before | Paired Reads After | Paired Read Bases After | Single Reads After | Single Read Bases After |
|---|---|---|---|---|---|---|
| WTSP-Tan | 99,374,744 | 9,937,474,400 | NA | NA | 97,162,587 | 9,014,814,467 |
| WTSP-White | 97,605,312 | 9,760,531,200 | NA | NA | 95,347,015 | 8,779,352,471 |
| SOSP-Paired | 271,249,550 | 27,124,855,000 | 245,289,038 | 23,613,455,033 | 11,228,223 | 992,474,010 |
| WCSP-Paired | 160,229,712 | 16,022,971,200 | 153,636,836 | 14,171,465,431 | 2,871,235 | 213,815,184 |

**Table 2**(on next page)

Transcriptome assembly description

Tissue of origin, pool size, assembly statistics (N50, average transcript length, number of transcripts) and annotation description (number of zebra finch genes with significant BLAST hit) for whole assembly and open reading frame (ORF) containing transcripts. "Complete Transcripts" are those containing both a start and stop codon. We used the individual tan and white morph assemblies in the subsequent BLAST search and annotation which yielded 15,805 genes.

| Species | Tissue | pool size | N50 | Mean Length | # Transcripts | # ORF | Complete Transcripts | ZF genes |
|---|---|---|---|---|---|---|---|---|
| WTSP-Tan | Whole Brain | 3 | 2,557 | 1,119 | 116,894 | 54,868 | 22,799 | - |
| WTSP-White | Whole Brain | 3 | 1,942 | 960 | 95,129 | 37,910 | 11,855 | - |
| WTSP-Both | Whole Brain | 6 | 2,284 | 982 | 149,184 | 58,284 | 24,388 | 15,805 |
| SOSP | Auditory Forebrain | 7 | 4,072 | 1,416 | 276,670 | 133,740 | 79,451 | 16,864 |
| WCSP | Hypothalamus | 12 | 3,415 | 1,591 | 307,617 | 206,926 | 115,515 | 16,646 |

# Table 3 (on next page)

Functional description of transcriptome assemblies

Gene Ontology categories significantly A) over- and B) under-represented in song (SOSP), white-crowned (WCSP) and white-throated (WTSP) sparrows (observed/expected, FDR adjusted Fisher's Exact Test, $p < 0.05$).

A.

| GO Category | SOSP | WCSP | WTSP |
|---|---|---|---|
| cytoplasm | 1810/1739 | 1793/1718 | 1751/1650 |
| intracellular | 1629/1575 | 1632/1555 | 1577/1494 |
| mitochondrion | 790/753 | 788/744 | 781/715 |
| nucleic acid binding | 935/903 | 935/892 | 900/857 |
| nucleolus | 244/231 | 243/229 | 241/220 |
| protein binding | 5298/5218 | 5258/5154 | 5037/4951 |
| protein phosphorylation | 558/539 | 558/532 | 542/511 |
| transferase activity, transferring phosphorous containing groups | 538/519 | 538/513 | 522/493 |

B.

| GO Category | SOSP | WCSP | WTSP |
|---|---|---|---|
| cytokine activity | 43/58 | 40/58 | 37/55 |
| DNA integration | 8/13 | 7/13 | 4/12 |
| extracellular region | 263/320 | 264/316 | 238/303 |
| hormone activity | 31/43 | 32/43 | 26/41 |
| immune response | 68/88 | 61/87 | 57/84 |
| MHC Class I protein complex | 3/8 | 2/7 | 2/7 |

# Table 4 <span style="font-size:smaller">(on next page)</span>

Functional differences between post-mortem and fresh tissues

GO terms underrepresented in post-mortem white-throated sparrow samples (observed/expected, adjusted p < 0.01), but not in song sparrow and white-crowned sparrow (adjusted p > 0.05).

| GO Category | WTSP | WCSP | SOSP |
|---|---|---|---|
| **photoreceptor activity** | 3/12 | 10/13 | 9/13 |
| **protein-chromophore linkage** | 3/12 | 10/13 | 9/13 |
| **visual perception** | 7/18 | 16/19 | 15/19 |
| **response to stimulus** | 7/17 | 14/18 | 13/18 |
| **G-protein coupled receptor activity** | 345/381 | 391/397 | 389/402 |
| **G-protein coupled purinergic nucleotide receptor activity** | 11/21 | 18/22 | 18/23 |
| **G-protein coupled purinergic nucleotide receptor signaling pathway** | 11/21 | 18/22 | 18/23 |
| **transporter activity** | 136/157 | 153/164 | 157/166 |
| **receptor activity** | 497/532 | 552/554 | 551/561 |
| **G-protein coupled receptor signaling pathway** | 463/496 | 513/517 | 514/523 |
| **integral to membrane** | 1564/1617 | 1683/1687 | 1692/1704 |
| **neurotransmitter transport** | 16/24 | 23/25 | 21/25 |