



THE UNIVERSITY *of* EDINBURGH

Edinburgh Research Explorer

## Recovering the 'missing' avian genes using multi-omics data

**Citation for published version:**

Zhong-Tao , Y, Smith, J & Zhuo-Cheng , H 2023, 'Recovering the 'missing' avian genes using multi-omics data', *Cytogenetic and Genome Research*. <https://doi.org/10.1159/000529376>

**Digital Object Identifier (DOI):**

[10.1159/000529376](https://doi.org/10.1159/000529376)

**Link:**

[Link to publication record in Edinburgh Research Explorer](#)

**Document Version:**

Peer reviewed version

**Published In:**

Cytogenetic and Genome Research

**General rights**

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

**Take down policy**

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact [openaccess@ed.ac.uk](mailto:openaccess@ed.ac.uk) providing details, and we will remove access to the work immediately and investigate your claim.



1

2

## Recovering the ‘missing’ avian genes using multi-omics data

3

Zhong-Tao Yin<sup>1</sup>, Jacqueline Smith<sup>2</sup>, Zhuo-Cheng Hou<sup>1,\*</sup>

4

5

### **Affiliation**

6

1. National Engineering Laboratory for Animal Breeding and Key Laboratory of  
7 Animal Genetics, Breeding and Reproduction, MARA; College of Animal Science  
8 and Technology, China Agricultural University, No. 2 Yuanmingyuan West Rd,  
9 Beijing 100193, China

10

2. The Roslin Institute & R(D)SVS, University of Edinburgh, Easter Bush, Midlothian,  
11 EH25 9RG, UK

12

13

### **Email addresses**

14

Co-authors: [yinzhngtao@cau.edu.cn](mailto:yinzhngtao@cau.edu.cn) (Z.T.Y.)

15

Co-authors: [jacqueline.smith@roslin.ed.ac.uk](mailto:jacqueline.smith@roslin.ed.ac.uk) (JS)

16

Corresponding author: [zhou@cau.edu.cn](mailto:zhou@cau.edu.cn) (Z.C.H.)

17

18

19 Gene gain and loss are common events in the evolution of species, especially for birds,  
20 which have evolved many unique characteristics such as feathers, wings, and flight  
21 capabilities, strong and lightweight skeletons, toothless beaks, high metabolic rates, and  
22 heat absorption sex, and unique respiratory and excretory systems [Kennedy and Ververs,  
23 1976, Blomme et al., 2006]. The release of the first chicken genome provided the basis  
24 for systematic analysis of the similarities and differences between vertebrate and avian  
25 genomes [International Chicken Genome Sequencing et al., 2004]. In comparison with  
26 other amniotes, bird genomes are more compact, and this difference may be related to  
27 the overall smaller cell size [Hughes and Hughes, 1995, Hughes and Friedman, 2008].  
28 The reductions in genome size may be the result of the loss of noncoding DNA  
29 sequences, with bird genomes having less repetitive DNA, fewer pseudogenes, and  
30 shorter introns than mammalian genomes [Hillier et al., 2004, Hughes and Piontkivska,  
31 2005]. Importantly, the evolution of avian genomes also appears to involve the loss of  
32 protein-coding genes, as the total number of uniquely identified avian-coding genes is  
33 much smaller than in other tetrapods (i.e., 23,294 in humans, GRCh38.p14; 19,404 in  
34 lizards, AnoCar2.0; 17,007 in chickens, GRCg7b). Paralog analysis revealed a higher  
35 overall incidence of gene families with fewer members in birds compared to other  
36 vertebrates [Hughes and Friedman, 2008]. Likewise, birds have a high rate of  
37 chromosomal rearrangements compared to other organisms, all of which may result in  
38 the deletion of protein-coding genes [Backstrom et al., 2010]. In recent years, the  
39 genomes of a large number of birds and lizards have been assembled and annotated,  
40 including zebra finches [Warren et al., 2010], chickens [Hillier et al., 2004], turkeys  
41 [Dalloul et al., 2010] and duck [Zhu et al., 2021]. Moreover, large-scale bird genome  
42 projects [Jarvis et al., 2014, Zhang et al., 2014], and chicken pan-genomes [Wang et al.,  
43 2021, Li et al., 2022] have also generated considerable genomic data. These large  
44 comparative genomic datasets identified hundreds of lost genomic-blocks in the bird  
45 genomes, and also suggested that hundreds of genes are missing in birds [Lovell et al.,  
46 2014, Zhang et al., 2014].

47 The missing genes seem to be directly related to the unique physiological  
48 phenomena of birds. Several functionally important genes in mammals are supposed

49 'missing' in chickens and have caused long-debated questions in bird biology. Spurious  
50 discovery of the missing/hidden genes in the bird genome has continued for decades.  
51 Previously, *BGN* [Blaschke et al., 1996], *COROIA* [Xavier et al., 2008], *MAPK3*  
52 [Lemoine et al., 2009], *MMP14* [Simsa et al., 2007], *TBX6* [Lardelli et al., 2003, Ahn  
53 et al., 2012], *TSSK4* [Shang et al., 2013] and five adipokine genes [Dakovic et al., 2014]  
54 were reported to be missing in birds, however, several long-debated genes including  
55 *TNF-alpha*, and *leptin* have been cloned in birds [Prokop et al., 2014, Seroussi et al.,  
56 2016, Rohde et al., 2018]. This hide-and-seek game still continues, and does not appear  
57 to be ending anytime soon [Elleder and Kaspers, 2019]. Here we summarize recent  
58 efforts using multi-omics data to probe those genes missing/hidden in avian genomes.

59

### 60 **Reconstruction of missing genes in the chicken genome**

61 While the hypothesis of missing genes in birds has been proposed for decades,  
62 researchers have found that some of the missing genes were, in fact, present in chickens  
63 or other birds. In the presence of large gaps and imperfect gene annotation in the  
64 genome, the *de novo* assembly of gene sequences using RNA-seq is considered to be  
65 an efficient way to identify unannotated genes in the genome. Attempts that only used  
66 a few tissues/organs have identified many missing genes in birds [Hron et al., 2015,  
67 Bornelov et al., 2017, Botero-Castro et al., 2017]. Recently, we used the raw data from  
68 26 chicken tissues downloaded from the GenBank database to assemble and obtain  
69 2,048,631 transcripts and identified 589 missing genes in birds [Yin et al., 2019].

70 At the same time, the continuity and integrity of chicken genome assemblies have  
71 been rapidly improving. The chicken genome released in 2017 was assembled by third-  
72 generation sequencing technology, and the number of annotated genes increased  
73 significantly (2,768 noncoding and 1,911 protein-coding genes) [Warren et al., 2017].  
74 In the *Gallus\_gallus-5.0* genome, 442 (77.41%, from a total of 571) genes thought to  
75 be missing in chickens (in Lovell et al., 2014 see Table S1 and Table S6, plus select  
76 entries in Table S4 and Table S18) were annotated, indicating that there is no systematic  
77 deletion of genes in birds. With the development of sequencing and hybrid assembly  
78 technology, the genomes of different chicken breeds continue to be assembled and

79 another 136 missing genes were further annotated in our recently assembled *Silkie*  
80 genome (unpublished). To date, it has now been shown that 528 (92.47%) genes that  
81 were thought to be missing, actually exist in chickens. This has been made possible by  
82 exploiting a large amount of multi-omics data available in chicken and has led to the  
83 revelation of genes with important functions such as *TNF- $\alpha$*  and *Leptin* [Seroussi et al.,  
84 2016, Rohde et al., 2018]. Recent large-scale chicken pan-genome data have also  
85 identified thousands of genes that are not presented in the current chicken reference  
86 genome [Li et al., 2022].

87

### 88 **Reconstruction of missing genes from other birds**

89 In addition to chicken, researchers have reconstructed many genes thought to be  
90 missing from other birds. We collected data from various important tissues from duck  
91 (24), pigeon (11), goose (8), and zebra finch (22) [Yin et al., 2019], and an avian  
92 transcriptomic database containing a total of 9,296,247 transcripts was constructed by  
93 *de novo* transcriptome assembly. From this, we identified several genes in duck (583),  
94 pigeon (558), goose (537), and zebra finch (543) from 806 genes that were thought to  
95 be missing in birds (in Lovell et al., 2014 see Table S1 and Zhang et al., 2014 see Table  
96 S10). Only 135 genes were not found in this bird transcriptome database. The number  
97 of missing genes reconstructed in different birds by *de novo* assembly of large  
98 transcriptome data is similar, indicating that these genes thought to be missing exist  
99 across different bird species.

100 In recent years, duck functional genomics has developed rapidly. We have  
101 assembled the Mallard, Pekin duck, and Shaoxing laying duck genomes using a  
102 combination of third-generation sequencing, Bionano, and Hi-C sequencing  
103 technologies. These have proved to be a rich source of genetic information [Zhu et al.,  
104 2021]. In the Mallard duck the CAU\_wild 1.0 genome has 1,872 more protein-coding  
105 genes annotated than the previous CAU 1.0 genome, including 89 genes previously  
106 thought to be missing in birds. Among these 89 genes, 5 genes have become  
107 pseudogenes, losing part of their gene function, 3 genes have been annotated as  
108 lncRNAs, and the remaining 81 genes remain as protein-coding genes. In addition, 240

109 genes were annotated as paralogous genes and 108 genes had similar segments in the  
110 genome. Mining large multi-omics data assemblies and annotations now reveals that  
111 only 10 genes (from a total of 806 missing genes), to date, have not been reconstructed  
112 in birds, with the rest of the genes thought to be missing in birds having been shown to  
113 actually exist. The recovered gene list is shown in Supplementary Table 1.

114

### 115 **Development of new methods to identify more missing genes**

116 Summarizing the characteristics of these reconstructed missing genes in birds and the  
117 reasons why they are thought to be missing can provide insights and methods for us to  
118 identify more missing genes. First, these reconstructed gene sequences have high GC  
119 content and length in many birds. The GC content of most of these ‘missing’ genes is  
120 more than 60%, and few genes even have over 80% (the median GC content of the  
121 chicken genome is 42.22% and the median GC content of the duck genome is 41.99%)  
122 [Hron et al., 2015, Bornelov et al., 2017, Botero-Castro et al., 2017, Yin et al., 2019].  
123 At the same time, the multi-tissue transcriptome expression profiles of birds showed  
124 that most of the reconstructed genes usually have strong tissue-specific expression.  
125 These genes are generally expressed predominantly in one tissue and are rarely  
126 expressed in the other tissues [Yin et al., 2019]. High-throughput transcriptome-based  
127 assembly approaches have limitations for fully recovering missing genes due to  
128 technical factors such as the PCR amplification bias against GC-rich fragments  
129 [Beauclair et al., 2019]. Expression patterns, i.e. tissue-specific expression patterns, and  
130 low expression, also limit the ability for full transcriptome assembly. Now, the third-  
131 generation sequencing technologies, which have less GC bias, such as single-molecule  
132 real-time (SMRT) and nanopore sequencing technologies, can obtain full-length  
133 transcripts directly, without assembly [Yin et al., 2019; Kuo et al. 2020]. The missing  
134 genes will continue to be discovered with the accumulation of full-length transcriptome  
135 data from more avian tissues from different physiological conditions.

136 Furthermore, the missing genes annotated in the chicken and duck genomes are  
137 mainly distributed on the micro-chromosomes, the ends of the chromosomes, and  
138 within regions showing a high content of tandem repeats clustering with non-canonical

139 DNA structures. [Zhu et al., 2021, Li et al., 2022]. Long repetitive regions [Treangen  
140 and Salzberg, 2011], regions of high GC content [Chen et al., 2013], telomeric regions,  
141 fragmented micro-chromosomes [O'Connor et al., 2019], and adaptive assembly  
142 strategies have always proved problematic for enabling complete bird genome assembly.  
143 To fully resolve the whole chicken gene sets, a Telomere-to-Telomere (T2T) genome is  
144 necessary. The recently completed human T2T genome has now paved the way for the  
145 finished bird genome assembly [Miga et al., 2020, Hoyt et al., 2022, Mao and Zhang,  
146 2022, Nurk et al., 2022]. Ultra-long ONT sequencing, high-precision HiFi sequencing  
147 data, multi-type auxiliary assembly data, and hybrid assembly using multiple strategies  
148 will greatly promote the quality of bird genome assembly [Sohn and Nam, 2018]. For  
149 large presence/absence variations within species, we can enrich genomic information  
150 by constructing high-quality multi-breed pan-genomes [Vernikos et al., 2015]. The Bird  
151 10,000 Genomes (B10K) Project [Zhang et al., 2015] has generated insightful results  
152 and the future bird T2T genome and pan-genome will undoubtedly reveal more genes.  
153 This complete gene map of birds will be critical for the further understanding of the  
154 biology and evolution of birds.

155 Finally, precise genome annotation will also provide the necessary sequence and  
156 structural information for mining more genes in birds. Annotation errors are  
157 unavoidable in genome annotation using automated processes, especially for some  
158 protein-coding genes that cannot be annotated in complex and high GC regions  
159 [Salzberg et al., 2019]. While applying full-length transcriptomic data for genome  
160 annotation [Nudelman et al., 2018, Wang et al., 2019; Kuo et al., 2020], the use of novel  
161 annotation methods developed based on machine learning can further improve the  
162 accuracy of annotation [Mahood et al., 2020, Stiehler et al., 2020]. More accurate  
163 manual annotation of important genome regions is also necessary for novel gene  
164 identification [Dunn et al., 2019]. It can be seen that, with the continuous development  
165 of omics technology and analysis methods, the genome information will be more  
166 complete, the annotation will be more accurate, and the genes that were previously  
167 thought to be missing in birds will continue to be discovered.

168

169

170



171 **Statement of Ethics**

172 All experiments with birds were performed under the guidance of ethical regulations  
173 from the Animal Care and Use Committee of China Agricultural University, Beijing,  
174 China.

175

176 **Conflict of Interest Statement**

177 The authors declare no competing interests.

178

179 **Funding Sources**

180 The work was supported by the National Waterfowl-Industry Technology Research  
181 System (CARS-42), the National Nature Science Foundation of China (31972525,  
182 31572388), Beijing Joint Research Program for Germplasm Innovation and New  
183 Variety Breeding (G20220628007), Beijing Municipal Science & Technology  
184 Commission (Z211100004621007).

185

186 **Author Contributions**

187 Z.C.H designed the study. Z.T.Y collected the samples and performed the analyses of  
188 the identification of the missing genes in birds. Z.C.H, J.S. and Z.T.Y wrote and revised  
189 the paper.

190

191 **Data Availability Statement**

192 Data have been submitted to the public databases under the following accession  
193 numbers: The raw data for assembling the transcriptome database of chicken, duck,  
194 goose, pigeon, and zebra finch were deposited in Sequence Read Archive (SRA)  
195 database under the accession number SRP141084. The Mallard genome is stored in  
196 NCBI under accession number PRJNA554956. The raw data for the Silkie genome  
197 assembly can be found in the SRA database under the accession number PRJNA805080  
198 (Unpublished data).

199

200

201

202 **References**

203 Ahn D, You KH, Kim CH. Evolution of the *tbx6/16* subfamily genes in vertebrates: insights from  
204 zebrafish. *Mol Biol Evol.* 2012;29(12):3959-83.

205 Backstrom N, Forstmeier W, Schielzeth H, Mellenius H, Nam K, Bolund E, et al. The recombination  
206 landscape of the zebra finch *Taeniopygia guttata* genome. *Genome Res.* 2010;20(4):485-95.

207 Beauclair L, Rame C, Arensburger P, Piegu B, Guillou F, Dupont J, et al. Sequence properties of  
208 certain GC rich avian genes, their origins and absence from genome assemblies: case studies. *BMC*  
209 *Genomics.* 2019;20(1):734.

210 Blaschke UK, Hedbom E, Bruckner P. Distinct isoforms of chicken decorin contain either one or  
211 two dermatan sulfate chains. *J Biol Chem.* 1996;271(48):30347-53.

212 Blomme T, Vandepoele K, De Bodt S, Simillion C, Maere S, Van de Peer Y. The gain and loss of  
213 genes during 600 million years of vertebrate evolution. *Genome Biol.* 2006;7(5):R43.

214 Bornelov S, Seroussi E, Yosefi S, Pendavis K, Burgess SC, Grabherr M, et al. Correspondence on  
215 Lovell et al.: identification of chicken genes previously assumed to be evolutionarily lost. *Genome Biol.*  
216 2017;18(1):112.

217 Botero-Castro F, Figuet E, Tilak MK, Nabholz B, Galtier N. Avian Genomes Revisited: Hidden  
218 Genes Uncovered and the Rates versus Traits Paradox in Birds. *Mol Biol Evol.* 2017;34(12):3123-31.

219 Chen YC, Liu T, Yu CH, Chiang TY, Hwang CC. Effects of GC bias in next-generation-sequencing  
220 data on de novo genome assembly. *PLoS One.* 2013;8(4):e62856.

221 Dakovic N, Terezol M, Pitel F, Maillard V, Elis S, Leroux S, et al. The loss of adipokine genes in  
222 the chicken genome and implications for insulin metabolism. *Mol Biol Evol.* 2014;31(10):2637-46.

223 Dalloul RA, Long JA, Zimin AV, Aslam L, Beal K, Blomberg Le A, et al. Multi-platform next-  
224 generation sequencing of the domestic turkey (*Meleagris gallopavo*): genome assembly and analysis.  
225 *PLoS Biol.* 2010;8(9).

226 Dunn NA, Unni DR, Diesh C, Munoz-Torres M, Harris NL, Yao E, et al. Apollo: Democratizing  
227 genome annotation. *PLoS Comput Biol.* 2019;15(2):e1006790.

228 Elleder D, Kaspers B. After TNF-alpha, still playing hide-and-seek with chicken genes. *Poult Sci.*  
229 2019;98(10):4373-4.

230 Hillier LW MW, Birney E, Warren WC, Consortium. *atmotCG*. Sequence and comparative analysis  
231 of the chicken genome provide unique perspectives on vertebrate evolution. *Nature.*  
232 2004;432(7018):695-716.

233 Hoyt SJ, Storer JM, Hartley GA, Grady PGS, Gershman A, de Lima LG, et al. From telomere to  
234 telomere: The transcriptional and epigenetic state of human repeat elements. *Science.*  
235 2022;376(6588):eabk3112.

236 Hron T, Pajer P, Paces J, Bartunek P, Elleder D. Hidden genes in birds. *Genome Biol.* 2015;16:164.

237 Hughes AL, Friedman R. Genome size reduction in the chicken has involved massive loss of  
238 ancestral protein-coding genes. *Mol Biol Evol.* 2008;25(12):2681-8.

239 Hughes AL, Hughes MK. Small genomes for better flyers. *Nature.* 1995;377(6548):391.

240 Hughes AL, Piontkivska H. DNA repeat arrays in chicken and human genomes and the adaptive  
241 evolution of avian genome size. *BMC Evol Biol.* 2005;5:12.

242 International Chicken Genome Sequencing C. Sequence and comparative analysis of the chicken  
243 genome provide unique perspectives on vertebrate evolution. *Nature.* 2004;432(7018):695-716.

244 Jarvis ED, Mirarab S, Aberer AJ, Li B, Houde P, Li C, et al. Whole-genome analyses resolve early  
245 branches in the tree of life of modern birds. *Science*. 2014;346(6215):1320-31.

246 Kennedy GY, Vevers HG. A survey of avian eggshell pigments. *Comp Biochem Physiol B*.  
247 1976;55(1):117-23.

248 Kuo RI, Cheng Y, Zhang R, Brown JWS, Smith J, Archibald AL, Burt DW. Illuminating the dark  
249 side of the human transcriptome with long read transcript sequencing. *BMC Genomics*. 2020 Oct  
250 30;21(1):751.

251 Lardelli M. The evolutionary relationships of zebrafish genes *tbx6*, *tbx16/spadetail* and *mga*. *Dev*  
252 *Genes Evol*. 2003;213(10):519-22.

253 Lemoine M, Dupont J, Guillory V, Tesseraud S, Blesbois E. Potential involvement of several  
254 signaling pathways in initiation of the chicken acrosome reaction. *Biol Reprod*. 2009;81(4):657-65.

255 Li M, Sun C, Xu N, Bian P, Tian X, Wang X, et al. De Novo Assembly of 20 Chicken Genomes  
256 Reveals the Undetectable Phenomenon for Thousands of Core Genes on Microchromosomes and  
257 Subtelomeric Regions. *Mol Biol Evol*. 2022;39(4).

258 Lovell PV, Wirthlin M, Wilhelm L, Minx P, Lazar NH, Carbone L, et al. Conserved syntenic clusters  
259 of protein coding genes are missing in birds. *Genome Biol*. 2014;15(12):565.

260 Mahood EH, Kruse LH, Moghe GD. Machine learning: A powerful tool for gene function prediction  
261 in plants. *Appl Plant Sci*. 2020;8(7):e11376.

262 Mao Y, Zhang G. A complete, telomere-to-telomere human genome sequence presents new  
263 opportunities for evolutionary genomics. *Nat Methods*. 2022;19(6):635-8.

264 Miga KH, Koren S, Rhie A, Vollger MR, Gershman A, Bzikadze A, et al. Telomere-to-telomere  
265 assembly of a complete human X chromosome. *Nature*. 2020;585(7823):79-84.

266 Nudelman G, Frasca A, Kent B, Sadler KC, Sealson SC, Walsh MJ, et al. High resolution annotation  
267 of zebrafish transcriptome using long-read sequencing. *Genome Res*. 2018;28(9):1415-25.

268 Nurk S, Koren S, Rhie A, Rautiainen M, Bzikadze AV, Mikheenko A, et al. The complete sequence  
269 of a human genome. *Science*. 2022;376(6588):44-53.

270 O'Connor RE, Kiazim L, Skinner B, Fonseka G, Joseph S, Jennings R, et al. Patterns of  
271 microchromosome organization remain highly conserved throughout avian evolution. *Chromosoma*.  
272 2019;128(1):21-9.

273 Prokop JW, Schmidt C, Gasper D, Duff RJ, Milsted A, Ohkubo T, et al. Discovery of the elusive  
274 leptin in birds: identification of several 'missing links' in the evolution of leptin and its receptor. *PLoS*  
275 *One*. 2014;9(3):e92751.

276 Rohde F, Schusser B, Hron T, Farkasova H, Plachy J, Hartle S, et al. Characterization of Chicken  
277 Tumor Necrosis Factor-alpha, a Long Missed Cytokine in Birds. *Front Immunol*. 2018;9:605.

278 Salzberg SL. Next-generation genome annotation: we still struggle to get it right. *Genome Biol*.  
279 2019;20(1):92.

280 Seroussi E, Cinnamon Y, Yosefi S, Genin O, Smith JG, Rafati N, et al. Identification of the Long-  
281 Sought Leptin in Chicken and Duck: Expression Pattern of the Highly GC-Rich Avian leptin Fits an  
282 Autocrine/Paracrine Rather Than Endocrine Function. *Endocrinology*. 2016;157(2):737-51.

283 Shang P, Hoogerbrugge J, Baarends WM, Grootegoed JA. Evolution of testis-specific kinases  
284 TSSK1B and TSSK2 in primates. *Andrology*. 2013;1(1):160-8.

285 Simsa S, Genina O, Ornan EM. Matrix metalloproteinase expression and localization in turkey  
286 (*Meleagris gallopavo*) during the endochondral ossification process. *J Anim Sci*. 2007;85(6):1393-401.

287 Sohn JI, Nam JW. The present and future of de novo whole-genome assembly. *Brief Bioinform*.

288 2018;19(1):23-40.

289 Stiehler F, Steinborn M, Scholz S, Dey D, Weber APM, Denton AK. Helixer: Cross-species gene  
290 annotation of large eukaryotic genomes using deep learning. *Bioinformatics*. 2020.

291 Treangen TJ, Salzberg SL. Repetitive DNA and next-generation sequencing: computational  
292 challenges and solutions. *Nat Rev Genet*. 2011;13(1):36-46.

293 Vernikos G, Medini D, Riley DR, Tettelin H. Ten years of pan-genome analyses. *Curr Opin  
294 Microbiol*. 2015;23:148-54.

295 Wang K, Hu H, Tian Y, Li J, Scheben A, Zhang C, et al. The Chicken Pan-Genome Reveals Gene  
296 Content Variation and a Promoter Region Deletion in IGF2BP1 Affecting Body Size. *Mol Biol Evol*.  
297 2021;38(11):5066-81.

298 Wang X, You X, Langer JD, Hou J, Rupprecht F, Vlatkovic I, et al. Full-length transcriptome  
299 reconstruction reveals a large diversity of RNA and protein isoforms in rat hippocampus. *Nat Commun*.  
300 2019;10(1):5009.

301 Warren WC, Clayton DF, Ellegren H, Arnold AP, Hillier LW, Kunstner A, et al. The genome of a  
302 songbird. *Nature*. 2010;464(7289):757-62.

303 Warren WC, Hillier LW, Tomlinson C, Minx P, Kremitzki M, Graves T, et al. A New Chicken  
304 Genome Assembly Provides Insight into Avian Genome Structure. *G3 (Bethesda)*. 2017;7(1):109-17.

305 Xavier CP, Eichinger L, Fernandez MP, Morgan RO, Clemen CS. Evolutionary and functional  
306 diversity of coronin proteins. *Subcell Biochem*. 2008;48:98-109.

307 Yin Z, Zhang F, Smith J, Kuo R, Hou ZC. Full-length transcriptome sequencing from multiple  
308 tissues of duck, *Anas platyrhynchos*. *Sci Data*. 2019;6(1):275.

309 Yin ZT, Zhu F, Lin FB, Jia T, Wang Z, Sun DT, et al. Revisiting avian 'missing' genes from de novo  
310 assembled transcripts. *BMC Genomics*. 2019;20(1):4.

311 Zhang G, Li C, Li Q, Li B, Larkin DM, Lee C, et al. Comparative genomics reveals insights into  
312 avian genome evolution and adaptation. *Science*. 2014;346(6215):1311-20.

313 Zhang G, Rahbek C, Graves GR, Lei F, Jarvis ED, Gilbert MT. Genomics: Bird sequencing project  
314 takes off. *Nature*. 2015;522(7554):34.

315 Zhu F, Yin ZT, Wang Z, Smith J, Zhang F, Martin F, et al. Three chromosome-level duck genome  
316 assemblies provide insights into genomic variation during domestication. *Nat Commun*.  
317 2021;12(1):5932.

318

319

**Supplementary Table 1.** The recovered ‘missing genes’ from birds

Previously reported gene list	No. of recovered missing genes	Gene Symbols
Evidence supporting the missing genes in five birds [Yin et al., 2019]	446	<p><i>ABCD1 ATL3 ADCY4 AIF1 ANKRD23 AP2A1 APLP1 ARAF ARHGEF25 ATP6AP1 AVIL AVPR2 BBS1 BCAT2 BGN BRSKI CA11 CALM3 CD37 CD97 CDC42EP5 CDH24 CDKN2D CEBPE CORO1A CPTIC CYP2F1 DNAJC4 DOC2A DOCK6 DUSP9 EGLN2 EPN1 ERF FAM120C FBXL19 FERMT3 FGD1 FOXA3 FOXH1 GDII GLS2 GMPR2 GPR173 GSK3A HAS1 HIF3A HIGD1C HOMEZ HOOK2 HSPB6 IGLON5 IRAK1 IRF2BP1 IRF9 ITPKC KANK2 KCNAB3 KCNK4 KCNK6 KCNN4 KDELRI KIFC2 KIRREL2 KLC2 KMT5C KREMEN2 LOC100555519 MACROD1 MAMSTR MAP3K10 MAP4K1 MAPK3 MARCH9 METTL21B MMP14* MRPL52* MYH14 NFATC4 NR1H2 NTN5 ORAI3 PACS1 PALM3 PARP2 PIM2 PLXNA3 PLXNB3 PNCK PODNLI POU5F1 PPPIA3 PPP1R3E PPP6R1 PRDX2 PRKACA PRRG2 PRSS8* PRX PSPN PTGIR PTPRH RAB1B RAB3D RBCK1 RCN3* RGL3 RNF181 RPS6KA4 RRAS SIPR5 SERPINI2 SLC22A17 SLC6A16 SLC6A8 SLC7A7 SNX15 SPRED3 STX4 SYNI SYNGR4 SYT5 TBX6 TEAD2 TEPI TIMM29 TIMP1 TMC4 TMEM150A TMEM91 TRMT1* TRPT1 TSPAN31* TSSK4 TUBB4A TULP2 USP11 WRAP53 YIPF2 YPEL3 ZNF784 ZSWIM4 ASF1B ATP2B3 BCL7C CCDC88B CHMP4A CNTD2 DENND1C DNAH2 EFS ELK1 EMP3 EPS8L1 ESRR A FAM98C FBRS FGF21 FKBPL FUZ GAPDHS GNG8 GRAMD1A HSD17B14 IRF3 JPH4 JUNB KCND1 KCNJ14 KMT2B LRP10 LRRC8E MMP25 MYL6B NPHS1* NTF4 NUCB1 PDZD4 PELI3 PLCB3* PLEKHA4 PLPPR2* PRDX5 PRR14 PRRC2A PRSS53 PSMB11 PSME1 RAB2B RASGRP4 REC8 REM2 RPGRIP1 SLC25A23SLC26A10 SLC44A4 SLC7A8 SPHK2 SPIB STX10 STXPB2 SYP TBC1D10B TCF19 TRIB3 TRPM4 TTC5 TTC9C TTYHI UBL4A YIF1A ZNF385A NFE2 OR52N4 CGREF1 OPA3 SCAND1 RCE1 FIS1 PIP4P1 SOAT2 SIX5 MPDU1 EPSTI1 PLOD3 OR52N2 ACP4 GRWD1 PRODH2 CIB4 GPR137 ABHD16B SMG9 RCOR2 IDO1 KDM6B ITGB7 PAX8 MAP3K11 TMEM147 SCYL1 PSRC1 FAM131C CADM4 B4GALNT2 BAZ2A PCMTD2 ZNF583 PHF23 ANKRD45 GATA1 IGHV1-2 IGHV4-34 ABCB5 DDX39A PELP1 PRKCSH GTF2IRD2 MED25 TMEM30B S100A5 IKBKG ABHD16A PPP1R12C GRIPAPI BRMS1 AQP6 RBM42 HCFC1 WDR6 ARHGAP11B MBD1 OR7C1 DBP CHD3 AKRIC3 KCTD13 COX4I2 ARG1 CABP5 NELFE OR1C1 CPA3 THAP8 RASIP1 ZNF653 NANOS2 ECT2L TIMM17B LMTK3 ALAS2 EHBPI1 SLC52A1 B3GAT3 EPOR PXDNL INO80E EFEMP2 SLC25A18 RNF31 OR52B4 PIWIL4 KHNYN BTBD18 OLIG1 PRLH STK19 SLC43A1 ADH4 KRII CD2BP2 LYL1 GPR108 GNB2 PRMT5 AXL CYB561D1 MEGF8 TSSK1B CALML6 SPNS1 ARR3 STAC2 REL A SOWAHB PGAM2 OXAIL CHRNB1 TP53 THOC6 SH2D6 GPR45 CIART ACY3 AIF1L PHF1 MUS81 OR11H6 NOSIP ARHGAP27 SRRM2 LSR TNS2 CXXC1 CUTA SRSF9 CDIPT ARL4D CHRMI PSMA8 SART1 CFL1 PLSCR3 OR10V1</i></p>

		<p><i>TBC1D10C ALKBH7 DNAJC30 RNF183 DCAF11 ACAP1 NOVA2 RMND5B TNFAIP8L2 NCKAP5L BEST2 ENTHD1 ABT1 SLC12A6 KLC3 CLIC1 ABHD4 GDPD3 DVL2 FAM89B ZDHHC11 COX7A1 PLD2 METTL21A TAOK2 MORC1 UXT SHANK1 OR51A7 NDRG2 CNOT3 XRCCI HUWE1 HIGD1B RARRES3 RNASEH2A TBC1D17 DGKA DNAH11 CLN3 NXPH3 NUMBL HOXD1 OR2B3 TNNT1 ZNF500 PSMB8 AKR1C1 WAS CNFN CCDC22 ARHGEF1 POU2F2 WDR74 CIC PPP1R13L FXR2 ACINI EMC9 RFX1 ZFPL1 B4GALNT1 SLC35E3 PHKG2 GNL1 OR52A5 EML3 TSC22D4 DUSP2 LIN37 BAG6 BCL6B ANXA9 MZT2A ARRB2 IER2 TOX4 CDK20 FAM151A SLC1A5 OR10G6 RASGRP2 MMP12 SLC25A45 OR5L1 TFPT C10ORF131 C11ORF58 ENSG00000238163 ENSG00000250246 ENSG00000255613 U82695.1</i></p>
<p>Newly annotated missing genes in Mallard genome [Zhu et al., 2021]</p>	89	<p><i>ZBTB39 GPR182 OLFM2 HIPK4 HSPA12B WNT1 ADCY6 HOXC6 GPAA1 KLHL33 GABBR1 SYT3 NR4A1 SOX12 SLC44A2 USP39 GUCY2D FMNL3 PRPH SHMT2 CCDC65 CACNB3 DCC KANSL2 MPZ OPLAH SGCA ILF3 CYP27B1 OSBPL7 FARSA TBX21 ERBB3 TBKBP1 ALDOA PPP5C MYL6 POLI TARBP2 COL5A3 RAB27B STAT6 SLC6A8 RAB5B PDLIM2 ATG4D PACS1 JOSD2 BLVRB RCE1 YIF1B ANKRD39 SAE1 ESYT1 FLOT1 KRII LETMD1 METTL1 TSFM RPS26 STAC3 PPOX TMEM150A ARHGAP9 UBA1 TINF2 TECR B9D2 DDIT3 FKBP11 TMEM88 RBCK1 DAZAP2 DCTN2 UBL5 TMEM147 TMEM205 STAP2 STARD6 RLN3 IRF3 RABGGTA CCDC68 PDZD4 C12ORF44 C18ORF54 C1ORF192 C2ORF68 CCDC88B</i></p>
<p>Newly annotated missing genes in chicken genome (Gallus_gallus-5.0) [Warren et al., 2017]</p>	240	<p><i>ASNA1 ATP5B AVIL B4GALT3 CALR CDK2 CFAP126 (aka C1orf192) COPZ1 ECSIT ERBB3 ESYT1 EXOSC5 FARSA FBXL12 FBXW9 GPR182 HNRNPUL1 HOXC6 IKZF4 KCNH2 KLHL33 MIP NOS3 PRIM1 PRPH RASAL3 RDH5 RPS26 S100A10 SCNM1 SDHC SEMA4C SMARCC2 SOX12 STAP2 TARBP2 TBX21 TGFB1 TNPO2 TTC9C YIF1B ZBTB39 ADCY6 ADGRL1 aka LPHN1 AKT2 ANKRD39 APEX1 ARF3 ATAT1 ATG4D BBS1 BCAP31 C12orf44 aka ATG101 C19orf52 CACNA1A CACNB3 CACNG7 CACNG8 CAMSAP3 CCDC120 CCDC130 CCDC65 CCDC97 CCNT1 CHD8 CLASRP CLEC17A CLPP CSAD CSRNP2 CYTH2 DAZAP2 DCTN2 DDIT3 DNAJB1 DNM2 DPF1 ETFB EXOSC4 FKBP11 FLOT1 FMNL3 FUS FUZ GABBR1 GATA1 GEMIN7 GNG3 GNL3L GPAA1 GPKOW GRIK5 GTF2F1 HCFC1 HDAC6 HIF3A HOOK2 HSPBP1 ILF3 IPO4 JOSD2 JUNB KANSL2 KCNA7 KEAP1 KHSRP KMT5C KRII L1CAM LDLR LENG8 LETMD1 LIN7B LRRC4B LSMD1 aka NAA38 MAP2K7 MAP4K1 MARK2 MARS MBOAT7 METTL21B METTL3 MMP14 MPZ MRPL52 MYBPC2 NDUFB7 NKPD1 NOSIP NR4A1 NRXN2 OS9 OTUD5 PIH1D1 PLD3 POU2F2 POU6F1 PPOX PPP1R10 PPP1R12C PPP1R9B PPP4C PPP5C PRKCG PRMT5 PRPF31 PRX PSMD8 QPRT RAB4B RABGGTA REC8 RELB RENBP RING1 RNASEH2A RNF31 RUVBL2 SAE1 SCAF1 SCN1B SETD1A SHANK1 SHMT2 SIPAIL3 SLC11A2 SLC17A7 SLC35A2 SLC39A7 SLC44A2 SMC1A SMG9 SNRNP70 SPINT2 SPRYD3 SPTBN4 SRPK3 SSR4 STAC3 STAT6 STIP1 STRN4 TECR TFCP2 TFE3 TFPT TINF2 TNNT1 TRMT112 TRPT1 TSFM TSPYL2 UBA1 UBL4A UBL5 USP39</i></p>

		<p><i>WNT1 ZNF385A ZNF653 ZNF668 ZNF865 AGAP2 ARL2 ASPDH AXL BLVRB C19orf53 C6orf136 CATSPERB CDK16 CYP27B1 DCAF11 FAM50A FKBP2 FLRT1 GPT GRM6 GUCY2D KDMB KIF5A LENG1 LMTK3 LTBP4 NOP9 OPLAH OSGEP PHF8 PIP4K2C PPP1R18 PQBP1 PRKCSH PSENE1 RCE1 SRCAP STX1B THOC6 TRAPPC1 UXTWDR45XAB2</i></p>
<p>Newly annotated missing genes in chicken genome (Silkie)</p>	136	<p><i>ZBTB12 ZNF865 GABBRI LENG9 PPP1R10 BAG6 FLOT1 ABCF1 C6orf136 LSM2 MRPS18B VARS ATATI FKBPL ERBB3 DNMT2 OSBPL7 MARK2 PACS1 RASAL3 ILF3 RUVBL2 ESYT1 FARSA ZBTB39 RFX1 POLD1 CSAD UBA1 OPLAH KEAP1 PRPF31 KANSL2 ALDOA XAB2 U2AF2 SPRYD3 WNT1 MAP2K7 PRIM1 POU6F1 B4GALNT1 KLHL33 KHSRP PPP1R37 RASIP1 STX1B ADCK5 PORCN STAT6 OSGEP GPR182 TARBP2 WDR45 PPP5C RDH5 MAP4K1 TBKBP1 CCDC65 RBM23 PPP1R12C CACNG7 TRMT1 ARF3 TECR USP39 RNF31 SMG9 PSMB5 DCTN2 SLC35A2 CLPP PRKCSH OS9 SYMPK RPL13A SNRNP70 CCDC22 EXOSC5 NPHS1 METTL1 DTX3 RNASEH2A MPZ PABPN1 NOSIP RPL18 CARM1 ASPDH LETMD1 TGFB1 STAC3 SGCA SOX12 SYP GTF2F1 BCAP31 TTC5 GPAA1 CCDC130 MED25 KRI1 EMC4 RENBP TMEM147 ANKRD39 ATP6A1 QPRT RCE1 PIH1D1 LIN37 BAX WRAP53 CYB5D1 PSMD8 BSCL2 RBCK1 TMEM150A TSPAN31 CCDC97 MBOAT7 DAZAP2 ALKBH7 THOC6 FUZ PPP1R3E XHI CCDC106 TFPT UXT BBS1 RLN3 TSFM DDIT3 CD2BP2 NOP9 TNF-<math>\alpha</math> LEPTIN</i></p>

321 \*: The gene has also been recovered from chicken RNA-Seq data in recent studies [Hron et al., 2015,

322 Bornelov et al., 2017, Botero-Castro et al., 2017].

323