# GigaScience

# De novo genome assembly of Indian Blue Peacock (Pavo cristatus) from Oxford Nanopore and Illumina sequencing reads
## --Manuscript Draft--

| | |
|---|---|
| Manuscript Number: | GIGA-D-18-00280 |
| Full Title: | De novo genome assembly of Indian Blue Peacock (Pavo cristatus) from Oxford Nanopore and Illumina sequencing reads |
| Article Type: | Data Note |
| Funding Information: | |

**Abstract:**

Background

Pavo cristatus the Indian blue peacock are geographically found distributed in natural habitats of South Asia. Peacock has been described among the bird species as one of the most elegant, majestic and beautiful bird (Fig. 1). Since prehistoric times they have been described or depicted in Indian culture and has been adopted as the national bird of India. Its length varies from 92-125 centimeter (without train), weighing about 4-8 Kilograms and lives up to 20 years in the wild. The avian species have been very important in the fields of phylogenetics, developmental studies, sexual reproduction and speciation. Avian genomics have contributed immensely towards understanding the vertebrate genome evolution. Here we present the first draft genome sequence of P. cristatus, yet another important bird species to further add values and gain insight into avian genomics.

Findings

For the first time in avian genomics, long reads using Oxford Nanopore technology have been used for the whole genome assembly. We sequenced different DNA insert size libraries from Illumina and long read Nanopore technologies from the peacock DNA. We performed de novo genome assembly by integrating the reads from Illumina short insert, long insert, multiple mate-pair reads along with Nanopore long reads using multiple genome improvement tools. A draft of the peacock genome of about 0.915 Gigabases (Gb) with a N50 of 0.23 Megabases (Mb) was assembled. Annotations with other avian species, protein families, KEGG were performed for functional understanding by in-silico approaches. Proteins were compared against Chicken, Turkey and Human to obtain evolutionary similarities and uniqueness of the Pavo species.

Conclusions

Our most important findings from the genome sequence of P. cristatus is to decipher the different gene families and to understand their role in body pattern development and other features that truly makes this bird unique. The genome sequence also gives insights on its genetic lineage and evolution with relation to other avian members. Several hypothesis and theories have been discussed with respect to sexual selection; now with the understanding of the genome sequence, some of these evolutionary theories will be better understood. The genome will also support future studies on population genetics and breeding for species conservation as well as in understanding its evolutionary ecology and sexual dimorphism. The comparative genomics with other avian species and specifically with Gallus gallus (Chicken) and Meleagris gallopavo (Turkey) have shown insights into the gene families and their conserved domains. Pavo proteins were also compared with human to understand the functional components that were conserved after the speciation split.

| | |
|---|---|
| Corresponding Author: | Subhradip Karmakar, PhD<br>All India Institute of Medical Sciences<br>New Delhi, Delhi INDIA |
| Corresponding Author Secondary Information: | |

| | |
|---|---|
| Corresponding Author's Institution: | All India Institute of Medical Sciences |
| Corresponding Author's Secondary Institution: | |
| First Author: | Subhradip Karmakar, PhD |
| First Author Secondary Information: | |
| Order of Authors: | Subhradip Karmakar, PhD |
| | Ruby Dhar |
| | Ashikh Seethy |
| | Karthikeyan Pethusamy |
| | Vishwajeet Rohil |
| | Sunil Singh |
| | Kakali Purkayastha |
| | Sandeep Goswami |
| | Rakesh Singh |
| | Indrani Mukherjee |
| | Ankita Raj |
| | Tryambak Srivastava |
| | Sovon Acharya |
| | Balaji Rajashekhar |
| Order of Authors Secondary Information: | |
| Additional Information: | |

| Question | Response |
|---|---|
| Are you submitting this manuscript to a special series or article collection? | No |
| **Experimental design and statistics**<br><br>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist. Information essential to interpreting the data presented should be made available in the figure legends.<br><br>Have you included all the information requested in your manuscript? | Yes |
| **Resources**<br><br>A description of all resources used, including antibodies, cell lines, animals | Yes |

| | |
|---|---|
| and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.<br><br>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist? | |
| **Availability of data and materials**<br><br>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in publicly available repositories (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the "Availability of Data and Materials" section of your manuscript.<br><br>Have you have met the above requirement as detailed in our Minimum Standards Reporting Checklist? | Yes |

1 *De novo* **genome assembly of Indian Blue Peacock (*Pavo***

2 *cristatus*)**, from Oxford Nanopore and Illumina sequencing**

3 **reads**

4

5

6 **Authors:** Ruby Dhar[1], Ashikh Seethy[1], Karthikeyan Pethusamy[1], Vishwajeet Rohil[2], Sunil

7 Singh[1], Kakali Purkayastha[2], Sandeep Goswami[1], Rakesh Singh[3], Indrani Mukherjee[1], Ankita

8 Raj[1], Tryambak Srivastava[1], Sovon Acharya[1], Balaji Rajashekhar[4,*] and Subhradip

9 Karmakar[1,*]

10

11 **Affiliation**: [1]Department of Biochemistry, AIIMS, New Delhi, India. [2]Vallabhbhai Patel Chest

12 Institute (VPCI), New Delhi, India. [3]Kanpur Zoo, Kanpur, India.[4] Institute of Computer

13 Science, University of Tartu, 50409 Tartu, Estonia.

14

15 **\*Corresponding Authors email:** subhradip.k@aiims.edu, balaji@ut.ee

16

17

18

19 **Running Title:** De novo Genome Assembly of the Peacock Bird

20

21

22 **Key words:** Peacock, *Pavo cristatus*, Indian National Bird, Genome Assembly, Oxford

23 Nanopore.

24 **Abstract**

**Background**

*Pavo cristatus* the Indian blue peacock are geographically found distributed in natural habitats of South Asia. Peacock has been described among the bird species as one of the most elegant, majestic and beautiful bird (Fig. 1). Since prehistoric times they have been described or depicted in Indian culture and has been adopted as the national bird of India. Its length varies from 92-125 centimeter (without train), weighing about 4-8 Kilograms and lives up to 20 years in the wild. The avian species have been very important in the fields of phylogenetics, developmental studies, sexual reproduction and speciation. Avian genomics have contributed immensely towards understanding the vertebrate genome evolution. Here we present the first draft genome sequence of P. *cristatus*, yet another important bird species to further add values and gain insight into avian genomics.

**Findings**

For the first time in avian genomics, long reads using Oxford Nanopore technology have been used for the whole genome assembly. We sequenced different DNA insert size libraries from Illumina and long read Nanopore technologies from the peacock DNA. We performed *de novo* genome assembly by integrating the reads from Illumina short insert, long insert, multiple mate-pair reads along with Nanopore long reads using multiple genome improvement tools. A draft of the peacock genome of about 0.915 Gigabases (Gb) with a N50 of 0.23 Megabases (Mb) was assembled. Annotations with other avian species, protein families, KEGG were performed for functional understanding by insilico approaches. Proteins were compared against Chicken, Turkey and Human to obtain evolutionary similarities and uniqueness of the *Pavo* species.

**Conclusions**

2

50 Our most important findings from the genome sequence of *P. cristatus* is to decipher the

51 different gene families and to understand their role in body pattern development and other

52 features that truly makes this bird unique. The genome sequence also gives insights on its

53 genetic lineage and evolution with relation to other avian members. Several hypothesis and

54 theories have been discussed with respect to sexual selection; now with the understanding of

55 the genome sequence, some of these evolutionary theories will be better understood. The

56 genome will also support future studies on population genetics and breeding for species

57 conservation as well as in understanding its evolutionary ecology and sexual dimorphism.

58 The comparative genomics with other avian species and specifically with *Gallus gallus*

59 (Chicken) and *Meleagris gallopavo* (Turkey) have shown insights into the gene families and

60 their conserved domains. *Pavo* proteins were also compared with human to understand the

61 functional components that were conserved after the speciation split.

62

## Introduction

*Pavo cristatus* commonly known as the blue Indian peacock obtained the status of National Bird of India in 1963. Peacocks have been distributed widely in Asian countries. The bird is very popular as it symbolizes beauty, love, grace and pride (Gadagkar, R; Kushwaha, et al.). It has been referred in ancient literatures of India and has been found closely associated with the life and culture of the peoples from South East Asia and particularly India (Kadgaonkar, SB). Peacocks are native to South Asia but have been introduced into many other countries usually as exhibits in park, zoos and also large number of aviculturists raise and breed these species as pets (Brickle, N; Jackson, C).

International Chicken Genome Sequencing Consortium sequenced the *Gallus gallus* genome in 2004, and this laid the foundation for study of avian genomics. A decade later, the avian genome consortium assembled 48 genomes of wide variety of avian species to understand the evolutionary relationships of modern birds (Zhang, G., et al.). Raw sequencing data for each species were generated at from coverage of 6X for zebra finch to a maximum 160X for budgerigar. The genome size varied from 1.04 to 1.26 Gb (http://avian.genomics.cn/en/). The studies on different bird species have provided a new perspective on vertebrate genome evolution. These genomes have also helped in improving the annotation of mammalian genomes. There are several distinguishing as well as unique features between chicken and human genome including genome size which is of one third of humans, conserved synteny blocks complete absence of active short interspersed nucleotide elements (SINE) to mention a few (International Chicken Genome sequencing consortium).

Despite the wealth of information from avian genomes sequencing projects, it is very important to genome sequence other new species to add value into aves and vertebrate

4

88 genomics. For the first time in bird genomics, Oxford Nanopore technology has been used to

89 sequence a bird genome in this present study. The long read chemistry will help in better

90 genome assembly of the TEs and repeat rich. The peacock genome will aid in understanding

91 about the uniqueness of this birds genome in comparison to other bird species. Comparative

92 genomes will help in understanding the development of this species, sexual selection and its

93 evolutionary relationships with other birds. The characterization of the genes involved in sex

94 determination could provide relevant information for the selective breeding of the peafowl

95 populations. We have unraveled some of the genomic signatures and thus have reported

96 unique gene pools of this bird by performing comparative genomics. Further different data

97 types will improve the assembly and gene/genome characterization will help to address the

98 sexual selection theory and key answers relevant to the evolution of this bird.

99

100

**Materials and methods**

**Sample collection and extraction of DNA**

103 The whole blood of male peacock was collected from Kanpur zoo, India after obtaining the

104 necessary ethical and institutional approval. 20μl of Proteinase K (PK) solution was taken

105 into a 1.5ml micro centrifuge tube. 200μl of blood was added and briefly mixed. 200μl of cell

106 lysis buffer was added to the tube, mixed by vortexing for 10seconds; incubated at 56°C for

107 10minutes. ReliaPrep™ Binding Column was placed into an empty collection tube. 250μl of

108 Binding Buffer (BBA) was added, capped the tube, and mixed by vortexing for 10 seconds

109 with a vortex mixer. Contents of the tube were added to the ReliaPrep™ binding column,

110 capped and placed in a refrigerated micro centrifuge. These were then centrifuged for 1

111 minute at maximum speed and flow through was discarded. Binding column was placed into

112 a fresh collection tube. 500μl of column wash solution was added to the column and

centrifuged for 3 minutes at maximum speed; Flow through was again discarded. Column

washing is repeated thrice. Columns were then placed in a nuclease free clean 1.5ml micro

centrifuge tube. 100 µl of Nuclease-Free Water was then added to the column and

centrifuged for an additional 1 minute at maximum speed. Column was discarded and the

elute was saved. The concentration and purity of the extracted DNA was evaluated using

Nanodrop Spectrophotometer (Thermo Scientific) and Qubit flurometer and integrity was

checked on a 0.8% agarose gel. The DNA sample was aliquoted for library preparation on

two different platforms: Illumina HiSeq4000 and Oxford Nanopore Technologies (ONT).

**HiSeq Paired-End library preparation and sequencing**

Whole genome sequencing (WGS) libraries were prepared with Illumina-compatible

NEXTflex DNA sequencing kit (BIOO Scientific, Austin, Texas, U.S.A.). Briefly, approx. 1

µg of genomic DNA was sheared using Covaris S2 sonicator (Covaris, Woburn,

Massachusetts, USA) to generate approx. fragment size distribution from 300 to 600 basepair

(bp). The fragment size distribution was checked on Agilent 2200 Tape Station with D1000

DNA screen tapes and reagents (Agilent Technologies, Palo Alto, CA, USA) and

subsequently purified using HighPrep magnetic beads (Magbio Genomics Inc, USA). The

purified fragments were end-repaired, adenylated and ligated to Illumina multiplex barcode

adaptors as per NEXTflex DNA sequencing kit protocol (BIOO Scientific, Austin, Texas,

USA).

The adapter-ligated DNA was purified with HighPrep beads (MagBio Genomics, Inc,

Gaithersburg, Maryland, USA) and then size selected on 2% low melting agarose gel and

cleaned using MinElute column (QIAGEN). The resultant fragments were amplified for 10

cycles of PCR using Illumina-compatible primers provided in the NEXTFlex DNA

138 sequencing kit. The final PCR product (sequencing library) was purified with HighPrep

139 beads, followed by library quality control check. The Illumina-compatible sequencing library

140 was initially quantified by Qubit fluorometer (Thermo Fisher Scientific, MA, USA) and its

141 fragment size distribution was analyzed on Agilent TapeStation. Finally, the sequencing

142 library was accurately quantified by quantitative PCR using Kapa Library Quantification Kit

143 (Kapa Biosystems, Wilmington, MA, USA). The qPCR-quantified library was subjected to

144 sequencing on an Illumina sequencer for 150 bp paired-end chemistry.

145

146 The Illumina-compatible sequencing library for the samples has a fragment size range

147 between 275 to 425 bp for Paired-End Short Insert (PE-SI) and 350 bp to 650bp for Paired-

148 End Long Insert (PE-LI). As the combined adapter size is approximately 120bp, the effective

149 user-defined insert size is 155 to 305 bp and 230 to 530 bp for PE-SI and PE-LI respectively.

150 Libraries were sequencing in Illumina HiSeq platform with 150*2 chemistry. The short reads

151 of Paired-End Short Insert (PE-SI), Paired-End Long Insert (PE-LI) and Mate-Pair (MP) from

152 Illumina HiSeq platform.

153

154

155 **Mate-Pair library preparation and sequencing**

156 Mate Pair sequencing library was prepared with Illumina-compatible Nextera Mate Pair

157 Sample Preparation Kit (Illumina Inc., Austin, TX, U.S.A.). Briefly, approx. 4 ug of genomic

158 DNA was simultaneously fragmented and tagged with Mate Pair adapters in a Transposon

159 based Tagmentation step. Tagmented DNA was then purified using AMPure XP Magnetic

160 beads (Beckman Coulter Life Sciences, Indianapolis, IN, U.S.A.) followed by Strand

161 Displacement to fill gaps in the Tagmented DNA. Strand displaced DNA was further purified

162 with AMPure XP beads before size-selecting the 3-5 Kilobases (Kb), 5-7 Kb & 7-10 Kb

163 fragments on low melting agarose gel. The fragments were circularized in an overnight blunt-

164 end intra-molecular ligation step, which will result in circularization of DNA with the insert

165 flanked mate pair adapter junction.

166

167 The circularized DNA was sheared using Covaris S220 sonicator (Covaris, Woburn,

168 Massachusetts, USA) to generate approx. fragment size distribution from 300 bp to 1000 bp.

169 The sheared DNA was purified to collect the Mate pair junction positive fragments using

170 Dynabeads M-280 Streptavidin Magnetic beads (Thermo Fisher Scientific, Waltham, MA,

171 U.S.A.). The purified fragments were end-repaired, adenylated and ligated to Illumina

172 multiplex barcode adaptors as per Nextera Mate Pair Sample Preparation Kit protocol.

173

174 The adapter-ligated DNA was then amplified for 15 cycles of PCR using Illumina-compatible

175 primers. The final PCR product (sequencing library) was purified with AMPure XP beads,

176 followed by library quality control check. The Illumina compatible sequencing library was

177 initially quantified by Qubit fluorometer (Thermo Fisher Scientific, MA, USA) and its

178 fragment size distribution was analyzed on Agilent TapeStation. Finally, the sequencing

179 library was accurately quantified by quantitative PCR using Kapa Library Quantification Kit

180 (Kapa Biosystems, Wilmington, MA, USA). The qPCR quantified libraries were pooled in

181 equimolar amounts to create a final multiplexed library pool for sequencing on an Illumina

182 sequencer.

183

**Nanopore MinION library preparation and sequencing**

185 Genomic DNA (1.5µg) was end-repaired (NEBnext ultra II end repair kit, New England

186 Biolabs, MA, USA), cleaned up with 1x AmPure beads (Beckmann Coulter,USA). Adapter

187 ligation were performed for 20 minutes using NEB blunt/ TA ligase (New England Biolabs,

188 MA, USA). Library mix were cleaned up using 0.4X AmPure beads (Beckmann Coulter,

189 USA) and eluted in 25 μl of elution buffer. Eluted Library were used for sequencing. Whole

190 genome library were prepared by using ligation sequencing kit SQK-LSK108-Oxford

191 Nanopore Technology (ONT) from Oxford Nanopore Technology. Sequencing were

192 performed on MinION Mk1b (Oxford Nanopore Technologies, Oxford, UK) using SpotON

193 flow cell (FLO-MIN106) in a 48hr sequencing protocol on MinKNOW 1.1.20 from ONT.

194

**Illumina raw data QC and processing**

196 The Illumina reads were de-multiplexed using Illumina bcl2fastq. The Illumina generated

197 raw data for genomic libraries was quality checked using FastQC (Andrews, S., 2010). The

198 paired-end Illumina reads were processed for clipping the adapter and low-quality bases

199 using customized script which retains minimum 70% bases/reads with Phred score (Q≥30 in

200 each base position) with a read length of 50 bp. The MP libraries were trimmed for adapter

201 and low-quality base trimming from the 3'-end using PLATANUS internal trimmer (Kajitani,

202 R., et al.).

203

**Nanopore reads base calling and processing**

205 Base calling was performed using Metrichor V.2.43.1 is a cloud based analysis tool provided

206 by Oxford Nanopore Technology software suite. The Nanopore reads were processed using

207 Poretools (Loman, NJ., et al.) for converting fast5 files to fasta format. For further

208 quantification and analysis the 2D reads or 1D high quality reads were selected for further

209 assembly.

210

211 *De novo* **genome assembly and genome size estimation**

212 The quality checked Nanopore reads were error-corrected using Illumina PE reads. For error-

213 correction the Illumina PE-reads were aligned to the Nanopore reads by using BWA aligner

214 (Li, H., et al.). The paired-end reads were assembled using Abyss (Birol. I., et al.) followed

215 by contig extension using Nanopore reads using SSPACE-LongRead (Boetzer, M., et al.).

216 Super scaffolding of the assembled scaffold was performed using SSPACE (Boetzer, M., et

217 al.) and PLATANUS using the Nanopore and Matepair data. Final draft genome resulted

218 after gap closure by GAPCLOSER and PLATANUS gap_close tool using Illumina data. The

219 genome size was estimated using a k-mer distribution plot using JELLYFISH (Marcais, G., et

220 al.). The repetitive elements were identified in the final assembled draft genome using Repeat

221 Masker tool. The draft genome was hardmasked by using reference genomic repeats of *G.*

222 *gallus*. The assembly and annotation workflow overview has been represented as Figure 2.

223

**Simple sequence repeats prediction**

225 Final assembled scaffolds were analysed for Simple Sequence Repeats (SSR) identification.

226 SSRs like the di, tri, tetra, penta and hexa-nucleotide repeats in the genome were obtained

227 using MISA (Version 1.0.0).

228

**Genome prediction and annotation**

230 Gene models was predicted on a hard masked draft genome, where the repetitive elements in

231 the draft genome were masked using genomic repeats of *G. gallus* with Repeatmasker tool

232 and further genes were predicted using AUGUSTUS with *G. gallus* as a reference model.

233 The predicted proteins were annotated by using BLASTP (Altshul, S., et al.) against all

234 Avian sequences downloaded from UniProt Protein Database.

235

**Pathway Analysis of the draft genome**

237 The predicted proteins were searched against the KEGG-KAAS server (Moriya, Y., et al.) for

238 pathway analysis. *G. gallus* (chicken), *Meleagris gallopavo* (turkey), *Taeniopygia guttata*

239 (zebra finch), *Falco peregrinus* (peregrine falcon) were used as reference organism for

240 pathway identification. The EuKaryotic Orthologous Groups (KOGs) were predicted using

241 homology based approach.

242

**Mitochondrial genome assembly and annotation**

244 The generated scaffolds from the draft assembly were aligned against the *P. cristatus*

245 mitochondria genome and the mapped reads were filtered and stitched using ABACUS

246 software using the same reference (Zhou, TC., et al.). Further gap closure were performed

247 with 3-7kb MP reads to generate an complete assembled mitochondrial genome. MITOS

248 (Bernt, M., et al.) was used for gene annotation. Circular plot generated using GenomeVx

249 (http://wolfe.ucd.ie/GenomeVx/) representing the localization of the gene in the assembled

250 mitochondrial genome.

251

**Phylogenetic tree construction**

253 The assembled Peacock mitochondrial genome was searched against 695 avian mitochondrial

254 genomes downloaded from NCBI. Based on the Blast-N homology results (with query

255 coverage> 100, subject coverage > 95, % identity >85 and with 1% gaps allowed in the

256 sequences). 51 mitochondrial genome sequences along with our assembled mitochondrial

257 genome were filtered. Multiple sequence alignment with default parameters were performed

258 using MUSCLE global sequence aligner. Phylogenetic trees were constructed using IQ-

259 TREE version 1.5.6 (www.iqtree.org). The parameters used for phylogenetic tree

260 construction were ultrafast boostrap (UFBoot, using the −bb option of 1000 replicates), and a

261 standard substitution model (-m MFP) was given for tree generation. The generated trees

262 from IQ-TREE tool were visualized using Figtree (http://tree.bio.ed.ac.uk/software/figtree/)

263 and the Brach-support values were recorded from the output ".treefile". The trees were

264 modified for better visualization under Trees section increasing order nodes were applied.

265

**Protein domain analysis**

267 Predicted proteins from Peacock, Chicken and Turkey with sequence length greater than 100

268 amino acids were considered for protein domain analysis. All the protein sequences from

269 each organism were searched against Pfam-A database using Pfam scan for protein domain

270 identification.

271

**Avian protein families**

273 The protein sequences of 48 avian genomes was downloaded from the link

274 http://avian.genomics.cn/en/jsp/database.shtml apart from the predicted proteins of the draft

275 genome. Sequences greater than 100 amino acids from all the avian genomes were filtered

276 and concatenated to a single fasta file. These sequences were clustered using CD-HIT (Fu, L.,

277 et al.) with 90% alignment coverage for the shorter sequence with a length difference cutoff

278 of 0.9. The single copy ortholog gene family present across all organisms and genes unique to

279 peacock were filtered and annotated.

280

**Genome conservation analysis**

282 The assembled draft genome was aligned against the *G. gallus* genome using Chromosomer

283 tool. Draft chromosomes were constructed based on alignment between fragments (scaffolds)

284 using a reference genome. These reordered assembled genome was aligned against the

285 Chicken genome using LAST aligner with NEAR (finding short-and-strong (near-identical)

286 similarities.) parameter allowing for substitution and gap frequencies leading to the

12

287 identification of orthologs. These query mapped regions were filtered with greater than 1% of

288 the maximum length for visualization using Circos.

289

290

291 **Results**

292 **DNA Sequencing data**

293 Five libraries of 150 bp paired-end from Illumina HiSeq technology were generated. The

294 short-insert reads of 489,114,747 accounted to genome coverage of 146.7X and long-insert

295 reads of 302,884,819 sequences was about 90.9X coverage with a total coverage of 236X.

296 Sequencing of three mate-pairs of 3-5Kb, 5-7Kb of and 7-10Kb yielded 72,915,033,

297 47,440,144 and 36,464,628 reads respectively with an approximate coverage of 21.9, 14.2

298 and 10.9 respectively, with a total of 156 million reads of 47X coverage. Oxford Nanopore

299 technology was used to generate 366,323 long reads having of 2,398,560,283 bp with

300 coverage of 2.3X. The complete sequencing was generated to a depth of ~287X from

301 Illumina and Nanopore platforms. The coverage was based on assuming the peacock genome

302 size of about 1 Gb (Table S1).

303

304 **Genome assembly**

305 The assembly was performed on Illumina reads with Abyss *de novo* assembler that resulted

306 in ~932 Mb (mega base) of genome with an N50 of 1639 bp. The extension of the contigs

307 were performed with Nanopore reads which generated scaffolds with N50 of 14,748 bp.

308 Super scaffolding of the assembled scaffold was performed using SSPACE and PLATANUS

309 with MP libraries that generated ~916 Mb genome with the N50 value of 168,140bp. The

310 final gap closer was executed using GAPCLOSER program with MP and PE-LI libraries

311 which generated a draft genome of 1.02 GB (giga base). The draft genome assembly of *Pavo*

312 *cristatus* consists of 179,346bp scaffolds, with a N50 of 189,886bp with 37 scaffolds having

313 sequence length >=1Mbp. Contigs above 5000 bp have covered a genome of ~0.915 Mb with

314 N50 0.23 Mb. In the assembled genome there were ~0.4% of non ATGC characters (Table

315 1).

316

317 Complete mitochondrial genome of 16699 bp was obatained. Total of 22 tRNA, tow rRNA

318 and 13 protein coding genes were identified in the assembled genome (Fig. 3a). A 100%

319 similarity was observed with the preiously published Peacocks mitochondrial genome (Fig.

320 S1a and S1b).

321

322 **Transposable Elements (TE)**

323 In the bird genome a total of 75,315,566 bp (7.3% of the genome) was predicted to have

324 5.5% of Retroelements (with SINEs 0.08% and LINEs 4.71%), 6.25 % total interspersed

325 repeats with 0.84 % simple repeats and 0.21% low complexity regions. The DNA

326 transposons identified in the genome was 0.71% (Table S2).

327

328 **Protein coding gene annotation**

329 A total of 23,153 proteins were predicted in the assembled draft genome using AUGUSTUS.

330 Among them 95% predicted genes were annotated against the other Aves proteins. The

331 21,854 annotated proteins showed top similarity to species *Gallus gallus* (Chicken) with

332 11,398 proteins, *Meleagris gallopavo* (Common turkey) with 4059 proteins, *Amazona aestiva*

333 (Blue-fronted Amazon parrot) with 1352 proteins and *Anas platyrhynchos* (Mallard) (*Anas*

334 *boschas*) with 849 proteins. Thirteen species had about 100 to 400 annotated proteins. The

335 remaining proteins were in the range of 1 to <100 proteins in about 62 species. From the

14

336    annotations a total of 13,161 proteins showed similarity to uncharacterized protein annotation

337    and some of the over represented proteins were Tyrosine-protein kinase, Sulfotransferase,

338    Phosphoinositide phospholipase C, Tetraspanin, Phospholipid-transporting ATPase,

339    Olfactory receptor, Polypeptide N-acetylgalactosaminyltransferase, Transporter, Keratin,

340    Hexosyltransferase, Protein Wnt, Kinesin-like protein, Gap junction protein, Claudin, POU

341    domain protein, Sodium/hydrogen exchanger, Phospholipid-transporting ATPase, Histone-

342    lysine N-methyltransferase and others (Table S3). The gene ontology annotations showed to

343    have Gene Ontology (GO) descriptions for 18,295 proteins. Among them, 14,490 proteins

344    have Molecular Function; 11,679 have Biological Process and 13,736 proteins have Cellular

345    Component as functional categories.

346    About 17.7% of proteins were found to have pathway information against the KEGG

347    database (Table S4). Some of the overrepresented annotations were Kinases like MAPK

348    (mitogen-activated protein kinase); JNK (c-Jun N-terminal kinase); RAF (RAF proto-

349    oncogene serine/threonine-protein kinase); AKT (RAC serine/threonine-protein kinase);

350    protein kinases and different GTPases.

351    Proteins searched against the KOG annotations showed a total of 20,937 proteins having

352    annotations. Among them, the most abundant annotations include Zn-finger, transmembrane

353    receptor, ubiquitin ligase, Leucine rich repeat, Cadherin repeats, Serine/threonine protein

354    kinase, Collagens, Ankryin repeat, Fibrillins, Voltage-gated Ca2+ channels and Hormone

355    receptors (Table S5). The peacock proteins when searched against the human proteins

356    showed expansions in ontologies for cell morphogenesis, neuronal projection and

357    development and GTPases (Table S10 and Fig. S4).

358

359

360    **Simple sequence repeats**

15

361 A total of 399,493 SSRs were identified from the peacock genome assembly. The largest

362 fraction of SSRs identified were mononucleotide (60.04%), followed by tetra nucleotide

363 (26%), di nucleotide (8.51%), tri nucleotide (4.31%), penta nucleotide (1.03%) and finally

364 hexa nucleotide (0.13%). Among the SSRs identified, A (49.2%) and T (44.9%) accounted

365 for 94.1% of the mono-nucleotide repeats. AT (23.8%), TA (16.5%), TG(13.7%), AC(10.6%)

366 and CA (10.32%) accounted for 75% of the di-nucleotide repeats. while TTG (9.9%), AAT

367 (9.6%), AAC (9.4%), TTA (7.1%), ATT (4.5%), TAA (3.5%), CAA (3.1%) and GGA

368 (2.69%) accounted for 49.7% of the tri-nucleotide repeats (Table S6).

369

**Avian protein family analysis**

371 A total of 748,544 protein sequences from 49 avian species have 653,497 protein sequences

372 of length above 100 amino acids (Table S7A). A total of 240,853 gene clusters were

373 generated of which 41 gene clusters had single copy orthologs in all avian species (Table S7B

374 and Table S7C). With the above stringent cutoff we observed 15,913 gene clusters were

375 unique to peacock species.

376

**Phylogeny and Genome comparisons**

378 The phylogeny of 51 mitochondrial genome sequences along with peacock genome showed a

379 clade consisting of *Pavo* species and *Gallus* (red junglefowl, Sri Lankan junglefowl, grey

380 junglefowl, grey junglefowl), Bambusicola (Mountain Bamboo-partridge, Chinese bamboo

381 partridge) and Francolinus (Chinese francolin). It can be observed many species were

382 distributed in two different clades where 34 species were found in one clade and five species

383 in other clade. Some of the endemic or native bird species like *Arborophila ardens*, *Acryllium*

384 *vulturinum* and *Numida meleagris* were found as clear outgroup of species in this study (Fig.

385 3b).

16

386 Predicted proteins from Peacock, Chicken and Turkey were searched for protein domain

387 analysis. 81% of the Pfams were common among the three species. About 94%, 98.4% and

388 99.7% predicted Pfam domains were identified in Peacock, Chicken and Turkey respectively.

389 There were 255, 69 and 14 Pfam domains found to be unique in the species mentioned above

390 respectively (Fig. 4).

391 The assembled Peacock genome was reordered for pseudo chromosomes generation against

392 the simple repeat masked Chicken genome (1.21 GB, Warren, WC., et al.) using

393 Chromosomer which generated a overall reordered Peacock genome of about 597MB. The

394 right side of the image represents the reference genome and left side of the image represents

395 the Peacock genome (Fig. 5).

396

**Conclusions**

398 Third generation sequencing in avian genomics where long reads having the substantiality to

399 improve genome assembly will benefit understanding the organisms in the structurally

400 complex regions having repeat elements and isoforms in the genome (Goodwin, S., et al.).

401 Using a combination of short reads of different insert sizes as well as mate pair reads

402 generated from Illumina technology along with long reads from Oxford Nanopore, we

403 obtained an improved assembly and a draft genome of the Indian Blue Peacock (*Pavo*

404 *cristatus*). In comparison to other avian genomes (Zhang, G., et al.), the current 290X

405 sequencing depth obtained from our study is one of the highest. With a N50 of 0.23Mb, we

406 presented here a reasonably reliable draft genome for the peacock species. The inclusion of

407 Nanopore reads 366,323 for scaffolding followed by subsequent gap-closing using Illumina

408 data led to a 26.2% reduction in the number of scaffolds and a 50.65% and 115% increase in

409 the scaffold and contig N50, respectively. With only 2.3X of long reads, a significant

410 improvement in the assembly was observed. On the contrary, the assembly contained less

17

411 than 0.4% of unknown nucleotides, which is very low for a draft assembly. With the addition

412 of more long reads along with transcriptomic sequencing along with scaffolding and/or gap

413 closure tools, further improvement in the assembly can be achieved.

414

415 Peacock seems to defy the Darwinian laws of natural selection. These concern were raised by

416 no other than Darwin himself. Hence, he proposed the theory of the sexual selection where

417 the female can choose for a male with a certain phenotypic feature such as brilliant color or a

418 long tail (Burgess, S.). Peacock's brilliantly colored long tail feathers seems to evolve at the

419 cost of finding its female partner thereby contributing its beneficial genes, even at the cost of

420 making itself venerable to predators. A female peafowl in turn tends to choose the mate with

421 the largest and decorated plumage, which indirectly reflects its healthiness and capacity to

422 wade off potential competitors thereby selecting the most suitable male. Peacocks beautiful

423 feathers with it all its artistry surely provides it with an advantage to impress the females

424 (Dakin, R., et al.). Understanding the formation of beautiful feathers from the genomic

425 context will help in resolving several evolutionary theories on sexual selection that have been

426 discussed on this species.

427 Pigment particles are embedded into the newly grown peacock feathers during the molting

428 season which seems to absorb light of selective wavelength there by imparting to the color of

429 the plumage (Mercedes Foster, Rennee Riedler, et al.) Pigment morphogenesis in *Gallus* has

430 been reported by a process of melanoblasts migration and colonization into feather bud where

431 they differentiate to produce the pigment, melanin (Kelsh, RN., et al.). The molecular

432 mechanisms that control the pigment cell migration have been narrowed down to proteins

433 Kit1 and FGFs which maintain the melanoblasts migration to feathers. Kit and FGF proteins

434 have also been identified in our current study (Table S3). Understanding of these proteins in

435 the patterned formation will help decode the pigment pattern morphology in peacocks. These

436 pigmented patterns play a role in communication, choice of mate and in some species it can

437 help in camouflage (Burgess, S)

438 It has been observed that the variations in the genome size among bird species are very low

439 (Table S8). The genome complexities of a species are influenced by the Transposable

440 elements (TE) that are believed to play a crucial role (Kapsuta, A., et al.). The long read

441 sequences have significantly helped in resolving the TEs in genome quality and assembly.

442 Peacock genome comparisons with Turkey and Chicken have showed closeness to the

443 Chicken species. The mitochondrial phylogeny also revealed similar findings. Homology

444 searches have shown several important gene family expansions such as Kinases, Zn finger

445 proteins, GTPases and others. Their roles in biology, development and evolution of the

446 Peacocks have to be further explored.

447 To summarize, we have assembled the *Pavo* genome using Illumina and Oxford nanopore

448 technology. The genome information can be valued and explored by avian enthusiasts to

449 further understand about this bird. Though not critically endangered yet, in India, peafowl

450 population is surely at a decline in the wild due to massive deforestation and habitat loss.

451 Thus is further compounded by increased poaching for meat and feathers. Our genome

452 sequencing initiative of *Pavo cristatus* is not just only from a conservational viewpoint, but

453 also to preserve a heritage associated with this bird that runs through centuries and that bears

454 a strong attachment to the national psyche.

455

**Availability of supporting data**

Supplementary data contains, read statistics, annotation, repeats identification, orthology analysis, mitochondria assembly and annotation. Figures, Gene ontology. DNA and library preparation protocols.

**Raw Data in SRA**

Raw reads (Illumina and Nanopore) are available in the Sequence Read Archive (SRA), and the Whole Genome Shotgun project has been deposited at GenBank under SRA Submission ID: SUB3108024, Bioproject: PRJNA413288 and Biosamples SAMN07739105 : SKPea2016_SI, SAMN07739104 : SKPea2016_LI, SAMN07739101 : FPL_3_5KB, SAMN07739102 : FPL_5_7KB, SAMN07739103 : FPL_7_10KB and SAMN07739107 : FPL_Nano.

**Competing interests**

The author(s) declare that they have no competing interests.

**Authors contributions**

RD, AS, KP performed wet lab experiments; RD designed work plan, experiments and logistics; SS, VR, KP SG IM and AR assisted with the work; RS provided samples from bird; BR, SK performed data analysis and interpretation; SK, BR, RD drafted the manuscript and SK overseen the whole project.

20

489 **Tables**

490 Table 1. *De novo* assembly statistics of the Peacock genome.

| Description | Contigs | Nanopore Scaffold | Super Scaffolds | GapClosed | >1000 Kb | >5000 Kb |
|---|---|---|---|---|---|---|
| Contigs | 685241 | 281272 | 179346 | 179346 | 34181 | 15026 |
| Maximum Length | 49159 | 251510 | 2390121 | 2488982 | 2488982 | 2488982 |
| Minimum Length | 300 | 5 | 265 | 265 | 1000 | 5000 |
| Average Length | 1360 | 3250 | 5111 | 5729 | - | - |
| Total Length | 932162464 | 914363908 | 916720956 | 1027551907 | 954483822 | 915373606 |
| Length >= 100 bp | 685241 | 281271 | 179346 | 179346 | 34181 | 15026 |
| Length >= 200 bp | 685241 | 281271 | 179346 | 179346 | 34181 | 15026 |
| Length >= 500 bp | 616120 | 186433 | 93727 | 93727 | 34181 | 15026 |
| Length >= 1 Kbp | 363428 | 104479 | 34168 | 34181 | 34181 | 15026 |
| Length >= 10 Kbp | 1591 | 24748 | 9249 | 10311 | 10311 | 10311 |
| Length >= 1 Mbp | 0 | 0 | 27 | 37 | 37 | 37 |
| Non-ATGC # | 350325 | 42696911 | 49169831 | 4034372 | 4032567 | 3978757 |
| Non-ATGC % | 0.038 | 4.67 | 5.364 | 0.393 | 0.422 | 0.435 |
| N50 value | 1639 | 14748 | 168140 | 189886 | 218023 | 232312 |

491

492

**Figure legend**

**Figure 1.** The beautiful and charismatic photo of Indian blue peacock (*Pavo cristatus)* bird.

**Figure 2.** Detailed workflow for *de novo* genome assembly and annotation.

**Figure 3a.** Circular representation of Peacock mitochondrial genome with genes predicted.

**Figure 3b.** Phylogenetic tree generated from mitochondrial genome from 52 different avian species.

**Figure 4.** Circular image of the assembled peacock genome aligned against the *G. gallus* genome using Chromosomer tool. Draft chromosomes were generated by similarity between scaffolds which were arranged on the reference chicken genome. Circos was used for visualization.

**Figure 5.** Venn diagram showing common and unique Protein family domains (Pfam) between Peacock, Chicken and Turkey proteins.

**References:**

Gadagkar, R., 2003. Is the peacock merely beautiful or also honest?. Current Science, 85(7), pp.1012-1020.

Kushwaha, S., and Kumar, A. 2016. A Review on Indian Peafowl (*Pavo cristatus*) Linnaeus, 1758. Journal of Wildlife and Research, 4, 42-59.

Kadgaonkar, Shivendra B. 1993. The peacock in ancient Indian art and literature. Bulletin of the Deccan College Research Institute, vol. 53, pp. 95–115. JSTOR, www.jstor.org/stable/42936434.

Brickle, N. 2002. Habitat use, predicted distribution and conservation of green peafowl (*Pavo muticus*) in Dak Lak Province, Vietnam. Biological Conservation, 105: 189-197.

Jackson, C. 2006. Peacock. London: Reaktion Books Ltd.

Zhang, G., Jarvis, E. D., and Gilbert, M. T. P. 2014. A flock of genomes. Science 346, 1308–1309.

International Chicken Genome Sequencing Consortium. 2004. Sequence and comparative analysis of the chicken genome provide unique perspectives on vertebrate evolution. Nature, 432(7018), 695.

Kajitani, R., Toshimoto, K., Noguchi, H., Toyoda, A., Ogura, Y., Okuno, M., Yabana, M., Harada, M., Nagayasu, E., Maruyama, H. and Kohara, Y., 2014. Efficient de novo assembly of highly heterozygous genomes from whole-genome shotgun short reads. Genome research, 24(8), pp.1384-1395.

Loman, N. J. and Quinlan, A. R. 2014. Poretools: a toolkit for analyzing nanopore sequence data. Bioinformatics, 30(23), 3399-3401.

Li H. and Durbin R. 2009. Fast and accurate short read alignment with Burrows-Wheeler Transform. Bioinformatics, 25:1754-60.

24

540 Birol, I., Jackman, S.D., Nielsen, C.B., Qian, J.Q., Varhol, R., Stazyk, G., Morin, R.D., Zhao,
541 Y., Hirst, M., Schein, J.E. and Horsman, D.E., 2009. De novo transcriptome assembly with
542 ABySS. Bioinformatics, 25(21), pp.2872-2877.

543

544 Boetzer, Marten, and Walter Pirovano. 2014. SSPACE-LongRead: scaffolding bacterial draft
545 genomes using long read sequence information. BMC bioinformatics 15.1: 211

546

547 Boetzer, M., Henkel, C.V., Jansen, H.J., Butler, D. and Pirovano, W., 2010. Scaffolding pre-
548 assembled contigs using SSPACE. Bioinformatics, 27(4), pp.578-579.

549

550 Marcais, G and Kingsford, C. 2011. A fast, lock-free approach for efficient parallel counting
551 of occurrences of k-mers. Bioinformatics 27(6): 764-770.

552

553 Altschul, S.F., Gish, W., Miller, W., Myers, E.W. and Lipman, D.J., 1990. Basic local
554 alignment search tool. Journal of molecular biology, 215(3), pp.403-410.

555

556 Moriya, Y., Itoh, M., Okuda, S., Yoshizawa, A.C. and Kanehisa, M., 2007. KAAS: an
557 automatic genome annotation and pathway reconstruction server. Nucleic acids research,
558 35(suppl_2), pp.W182-W185.

559

560 Zhou, T.C., Sha, T., Irwin, D.M. and Zhang, Y.P., 2015. Complete mitochondrial genome of
561 the Indian peafowl (*Pavo cristatus*), with phylogenetic analysis in phasianidae. Mitochondrial
562 DNA, 26(6), pp.912-913.

563

564 Bernt, M., Donath, A., Jühling, F., Externbrink, F., Florentz, C., Fritzsch, G., Pütz, J.,
565 Middendorf, M. and Stadler, P.F., 2013. MITOS: improved de novo metazoan mitochondrial
566 genome annotation. Molecular phylogenetics and evolution, 69(2), pp.313-319.

567

568 Fu, L., Niu, B., Zhu, Z., Wu, S. and Li, W., 2012. CD-HIT: accelerated for clustering the
569 next-generation sequencing data. Bioinformatics, 28(23), pp.3150-3152.

570

571 Warren, W.C., Hillier, L.W., Tomlinson, C., Minx, P., Kremitzki, M., Graves, T., Markovic,
572 C., Bouk, N., Pruitt, K.D., Thibaud-Nissen, F. and Schneider, V., 2016. A new chicken

573 genome assembly provides insight into avian genome structure. G3: Genes, Genomes,
574 Genetics, pp.g3-116.

575

576 Goodwin, S., McPherson, J.D. and McCombie, W.R., 2016. Coming of age: ten years of
577 next-generation sequencing technologies. Nature Reviews Genetics, 17(6), p.333.

578

579 Zhang, G., Li, C., Li, Q., Li, B., Larkin, D.M., Lee, C., Storz, J.F., Antunes, A., Greenwold,
580 M.J., Meredith, R.W. and Ödeen, A., 2014. Comparative genomics reveals insights into avian
581 genome evolution and adaptation. Science, 346(6215), pp.1311-1320.

582

583 Burgess, S., 2001. The beauty of the peacock tail and the problems with the theory of sexual
584 selection. Journal of Creation, 15(2), pp.94-102..

585

586 Dakin, R. 2008. The role of the visual train ornament in the courtship of peafowl, *Pavo*
587 *cristatus*. Masters Abstracts International, 47/03: 97.

588

589 Mercedes S Foster 1975. The overlap of molting and breeding in some birds. The Condor
590 77:304-314.

591

592 Renee Riedler, Christel Pesme, James Druzik, Molly Gleeson, Ellen Pearlstein The Journal of
593 the American Institute of Conservation  2014

594

595

596 Kelsh, R.N., Harris, M.L., Colanesi, S. and Erickson, C.A., 2009. Stripes and belly-spots—a
597 review of pigment cell morphogenesis in vertebrates. Seminars in cell & developmental
598 biology. Vol. 20, No. 1, pp. 90-104.

599

600 Kapusta, A. and Suh, A., 2017. Evolution of bird genomes—a transposon's- eye
601 view. Annals of the New York Academy of Sciences, 1389(1), pp.164-185.

602

603

604 **Webservers:**

605 FastQC : http://www.bioinformatics.babraham.ac.uk/projects/fastqc/
606 AUGUSTUS : http://augustus.gobics.de/

607    PLATANUS : http://platanus.bio.titech.ac.jp/

608    ABACUS : http://abacas.sourceforge.net/Manual.html

609    MISA: http://pgrc.ipk-gatersleben.de/misa/

610    REPEATMASKER : http://www.repeatmasker.org/

611    KOGs : https://genome.jgi.doe.gov/portal/help/kogbrowser.jsf

612    KAAS : http://www.genome.jp/tools/kaas/

613    MITOS : http://mitos.bioinf.uni-leipzig.de/

614    MUSCLE : https://www.drive5.com/muscle/

615    PFAMSCAN : https://www.ebi.ac.uk/seqdb/confluence/display/THD/PfamScan

616    LAST : http://last.cbrc.jp/

617    CHROMOSOMER : https://github.com/gtamazian/chromosomer

618    CIRCOS : http://circos.ca/

# *De novo* genome assembly of Indian Blue Peacock (*Pavo cristatus*), from Oxford Nanopore and Illumina sequencing reads

**Authors:** Ruby Dhar[1,], Ashikh Seethy[1], Karthikeyan Pethusamy[1,] Vishwajeet Rohil[2], Sunil Singh[1], Kakali Purkayastha[2], Sandeep Goswami[1], Rakesh Singh[3], Indrani Mukherjee[1,] Ankita Raj[1], Tryambak Srivastava[1], Sovon Acharya[1], Balaji Rajashekhar[4,*] and Subhradip Karmakar[1,*]

**Affiliation**: [1]Department of Biochemistry, AIIMS, New Delhi, India. [2] Vallabhbhai Patel Chest Institute (VPCI), New Delhi, India. [3]Kanpur Zoo, Kanpur, India.[4] Institute of Computer Science, University of Tartu, 50409 Tartu, Estonia.

**\*Corresponding Authors email:** subhradip.k@aiims.edu, balaji@ut.ee

**Running Title:** De novo Genome Assembly of the Peacock Bird

**Key words:** Peacock, *Pavo cristatus*, Indian National Bird, Genome Assembly, Oxford Nanopore.

**Abstract**

1

## Background

*Pavo cristatus* the Indian blue peacock are geographically found distributed in natural habitats of South Asia. Peacock has been described among the bird species as one of the most elegant, majestic and beautiful bird (Fig. 1). Since prehistoric times they have been described in Indian culture and has been adopted as the national bird of India. Its length varies from 92-125 centimeter (without train), weighing about 4-8 Kilograms and lives up to 20 years in the wild. The avian species have been very important in the fields of phylogenetics, developmental studies, sexual reproduction and speciation. Avian genomics have contributed immensely towards understanding the vertebrate genome evolution. Here we present the first draft genome sequence of P. *cristatus*, yet another important bird species to further add values and gain insight into avian genomics.

## Findings

For the first time in avian genomics, long reads using Oxford Nanopore technology have been used for the whole genome assembly. We sequenced different DNA insert size libraries from Illumina and long read Nanopore technologies from the peacock DNA. We performed *de novo* genome assembly by integrating the reads from Illumina short insert, long insert, multiple mate-pair reads along with Nanopore long reads using multiple genome improvement tools. A draft of the peacock genome of about 0.915 Gigabases (Gb) with a N50 of 0.23 Megabases (Mb) was assembled. Annotations with other avian species, protein families, KEGG were performed for functional understanding by insilico approaches. Proteins were compared against Chicken, Turkey and Human to obtain evolutionary similarities and uniqueness of the *Pavo* species.

## Conclusions

Our most important findings from the genome sequence of *P. cristatus* is to decipher the different gene families and to understand their role in body pattern development and other

52 features that truly makes this bird unique. The genome sequence also gives insights on its

53 genetic lineage and evolution with relation to other avian members. Several hypothesis and

54 theories have been discussed with respect to sexual selection; now with the understanding of

55 the genome sequence, some of these evolutionary theories will be better understood. The

56 genome will also support future studies on population genetics and breeding for species

57 conservation as well as in understanding its evolutionary ecology and sexual dimorphism. The

58 comparative genomics with other avian species and specifically with *Gallus gallus* (Chicken)

59 and *Meleagris gallopavo* (Turkey) have shown insights into the gene families and their

60 conserved domains. *Pavo* proteins were also compared with human to understand the

61 functional components that were conserved after the speciation split.

62

**Introduction**

*Pavo cristatus* commonly known as the blue Indian peacock obtained the status of National Bird of India in 1963. Peacocks have been distributed widely in Asian countries. The bird is very popular as it symbolizes beauty, love, grace and pride (Gadagkar, R; Kushwaha, et al.). It has been referred in ancient literatures of India and has been found closely associated with the life and culture of the peoples from South East Asia and particularly India (Kadgaonkar, SB). Peacocks are native to South Asia but have been introduced into many other countries usually as exhibits in park, zoos and also large number of aviculturists raise and breed these species as pets (Brickle, N; Jackson, C).

International Chicken Genome Sequencing Consortium sequenced the *Gallus gallus* genome in 2004, and this laid the foundation for study of avian genomics. A decade later, the avian genome consortium assembled 48 genomes of wide variety of avian species to understand the evolutionary relationships of modern birds (Zhang, G., et al.). Raw sequencing data for each species were generated at from coverage of 6X for zebra finch to a maximum 160X for budgerigar. The genome size varied from 1.04 to 1.26 Gb (http://avian.genomics.cn/en/). The studies on different bird species have provided a new perspective on vertebrate genome evolution. These genomes have also helped in improving the annotation of mammalian genomes. There are several distinguishing as well as unique features between chicken and human genome including genome size which is of one third of humans, conserved synteny blocks complete absence of active short interspersed nucleotide elements (SINE) to mention a few (International Chicken Genome sequencing consortium).

Despite the wealth of information from avian genomes sequencing projects, it is very important to genome sequence other new species to add value into aves and vertebrate genomics. For the

4

88 first time in bird genomics, Oxford Nanopore technology has been used to sequence a bird

89 genome in this present study. The long read chemistry will help in better genome assembly of

90 the TEs and repeat rich. The peacock genome will aid in understanding about the uniqueness

91 of this birds genome in comparison to other bird species. Comparative genomes will help in

92 understanding the development of this species, sexual selection and its evolutionary

93 relationships with other birds. The characterization of the genes involved in sex determination

94 could provide relevant information for the selective breeding of the peafowl populations. We

95 have unraveled some of the genomic signatures and thus have reported unique gene pools of

96 this bird by performing comparative genomics. Further different data types will improve the

97 assembly and gene/genome characterization will help to address the sexual selection theory

98 and key answers relevant to the evolution of this bird.

99

100

101 **Materials and methods**

102 **Sample collection and extraction of DNA**

103 The whole blood of male peacock was collected from Kanpur zoo, India after obtaining the

104 necessary ethical and institutional approval. 20µl of Proteinase K (PK) solution was taken into

105 a 1.5ml micro centrifuge tube. 200µl of blood was added and briefly mixed. 200µl of cell lysis

106 buffer was added to the tube, mixed by vortexing for 10seconds; incubated at 56°C for

107 10minutes. ReliaPrep™ Binding Column was placed into an empty collection tube. 250µl of

108 Binding Buffer (BBA) was added, capped the tube, and mixed by vortexing for 10 seconds

109 with a vortex mixer. Contents of the tube were added to the ReliaPrep™ binding column,

110 capped and placed in a refrigerated micro centrifuge. These were then centrifuged for 1 minute

111 at maximum speed and flow through was discarded. Binding column was placed into a fresh

112 collection tube. 500µl of column wash solution was added to the column and centrifuged for 3

113 minutes at maximum speed; Flow through was again discarded. Column washing is repeated

114 thrice. Columns were then placed in a nuclease free clean 1.5ml micro centrifuge tube. 100 μl

115 of Nuclease-Free Water was then added to the column and centrifuged for an additional 1

116 minute at maximum speed. Column was discarded and the elute was saved. The concentration

117 and purity of the extracted DNA was evaluated using Nanodrop Spectrophotometer (Thermo

118 Scientific) and Qubit flurometer and integrity was checked on a 0.8% agarose gel. The DNA

119 sample was aliquoted for library preparation on two different platforms: Illumina HiSeq4000

120 and Oxford Nanopore Technologies (ONT).

121

122 **HiSeq Paired-End library preparation and sequencing**

123 Whole genome sequencing (WGS) libraries were prepared with Illumina-compatible

124 NEXTflex DNA sequencing kit (BIOO Scientific, Austin, Texas, U.S.A.). Briefly, approx. 1

125 μg of genomic DNA was sheared using Covaris S2 sonicator (Covaris, Woburn,

126 Massachusetts, USA) to generate approx. fragment size distribution from 300 to 600 basepair

127 (bp). The fragment size distribution was checked on Agilent 2200 Tape Station with D1000

128 DNA screen tapes and reagents (Agilent Technologies, Palo Alto, CA, USA) and subsequently

129 purified using HighPrep magnetic beads (Magbio Genomics Inc, USA). The purified fragments

130 were end-repaired, adenylated and ligated to Illumina multiplex barcode adaptors as per

131 NEXTflex DNA sequencing kit protocol (BIOO Scientific, Austin, Texas, USA).

132

133 The adapter-ligated DNA was purified with HighPrep beads (MagBio Genomics, Inc,

134 Gaithersburg, Maryland, USA) and then size selected on 2% low melting agarose gel and

135 cleaned using MinElute column (QIAGEN). The resultant fragments were amplified for 10

136 cycles of PCR using Illumina-compatible primers provided in the NEXTFlex DNA sequencing

137 kit. The final PCR product (sequencing library) was purified with HighPrep beads, followed

138 by library quality control check. The Illumina-compatible sequencing library was initially

139 quantified by Qubit fluorometer (Thermo Fisher Scientific, MA, USA) and its fragment size

140 distribution was analyzed on Agilent TapeStation. Finally, the sequencing library was

141 accurately quantified by quantitative PCR using Kapa Library Quantification Kit (Kapa

142 Biosystems, Wilmington, MA, USA). The qPCR-quantified library was subjected to

143 sequencing on an Illumina sequencer for 150 bp paired-end chemistry.

144

145 The Illumina-compatible sequencing library for the samples has a fragment size range between

146 275 to 425 bp for Paired-End Short Insert (PE-SI) and 350 bp to 650bp for Paired-End Long

147 Insert (PE-LI). As the combined adapter size is approximately 120bp, the effective user-defined

148 insert size is 155 to 305 bp and 230 to 530 bp for PE-SI and PE-LI respectively. Libraries were

149 sequencing in Illumina HiSeq platform with 150*2 chemistry. The short reads of Paired-End

150 Short Insert (PE-SI), Paired-End Long Insert (PE-LI) and Mate-Pair (MP) from Illumina HiSeq

151 platform.

152

153

154 **Mate-Pair library preparation and sequencing**

155 Mate Pair sequencing library was prepared with Illumina-compatible Nextera Mate Pair

156 Sample Preparation Kit (Illumina Inc., Austin, TX, U.S.A.). Briefly, approx. 4 ug of genomic

157 DNA was simultaneously fragmented and tagged with Mate Pair adapters in a Transposon

158 based Tagmentation step. Tagmented DNA was then purified using AMPure XP Magnetic

159 beads (Beckman Coulter Life Sciences, Indianapolis, IN, U.S.A.) followed by Strand

160 Displacement to fill gaps in the Tagmented DNA. Strand displaced DNA was further purified

161 with AMPure XP beads before size-selecting the 3-5 Kilobases (Kb), 5-7 Kb & 7-10 Kb

162 fragments on low melting agarose gel. The fragments were circularized in an overnight blunt-

163  end intra-molecular ligation step, which will result in circularization of DNA with the insert

164  flanked mate pair adapter junction.

165

166  The circularized DNA was sheared using Covaris S220 sonicator (Covaris, Woburn,

167  Massachusetts, USA) to generate approx. fragment size distribution from 300 bp to 1000 bp.

168  The sheared DNA was purified to collect the Mate pair junction positive fragments using

169  Dynabeads M-280 Streptavidin Magnetic beads (Thermo Fisher Scientific, Waltham, MA,

170  U.S.A.). The purified fragments were end-repaired, adenylated and ligated to Illumina

171  multiplex barcode adaptors as per Nextera Mate Pair Sample Preparation Kit protocol.

172

173  The adapter-ligated DNA was then amplified for 15 cycles of PCR using Illumina-compatible

174  primers. The final PCR product (sequencing library) was purified with AMPure XP beads,

175  followed by library quality control check. The Illumina compatible sequencing library was

176  initially quantified by Qubit fluorometer (Thermo Fisher Scientific, MA, USA) and its

177  fragment size distribution was analyzed on Agilent TapeStation. Finally, the sequencing library

178  was accurately quantified by quantitative PCR using Kapa Library Quantification Kit (Kapa

179  Biosystems, Wilmington, MA, USA). The qPCR quantified libraries were pooled in equimolar

180  amounts to create a final multiplexed library pool for sequencing on an Illumina sequencer.

181

182  **Nanopore MinION library preparation and sequencing**

183  Genomic DNA (1.5μg) was end-repaired (NEBnext ultra II end repair kit, New England

184  Biolabs, MA, USA), cleaned up with 1x AmPure beads (Beckmann Coulter,USA). Adapter

185  ligation were performed for 20 minutes using NEB blunt/ TA ligase (New England Biolabs,

186  MA, USA). Library mix were cleaned up using 0.4X AmPure beads (Beckmann Coulter, USA)

187  and eluted in 25 μl of elution buffer. Eluted Library were used for sequencing. Whole genome

188 library were prepared by using ligation sequencing kit SQK-LSK108-Oxford Nanopore

189 Technology (ONT) from Oxford Nanopore Technology. Sequencing were performed on

190 MinION Mk1b (Oxford Nanopore Technologies, Oxford, UK) using SpotON flow cell (FLO-

191 MIN106) in a 48hr sequencing protocol on MinKNOW 1.1.20 from ONT.

192

**Illumina raw data QC and processing**

194 The Illumina reads were de-multiplexed using Illumina bcl2fastq. The Illumina generated raw

195 data for genomic libraries was quality checked using FastQC (Andrews, S., 2010). The paired-

196 end Illumina reads were processed for clipping the adapter and low-quality bases using

197 customized script which retains minimum 70% bases/reads with Phred score (Q≥30 in each

198 base position) with a read length of 50 bp. The MP libraries were trimmed for adapter and low-

199 quality base trimming from the 3'-end using PLATANUS internal trimmer (Kajitani, R., et al.).

200

**Nanopore reads base calling and processing**

202 Base calling was performed using Metrichor V.2.43.1 is a cloud based analysis tool provided

203 by Oxford Nanopore Technology software suite. The Nanopore reads were processed using

204 Poretools (Loman, NJ., et al.) for converting fast5 files to fasta format. For further

205 quantification and analysis the 2D reads or 1D high quality reads were selected for further

206 assembly.

207

***De novo* genome assembly and genome size estimation**

209 The quality checked Nanopore reads were error-corrected using Illumina PE reads. For error-

210 correction the Illumina PE-reads were aligned to the Nanopore reads by using BWA aligner

211 (Li, H., et al.). The paired-end reads were assembled using Abyss (Birol. I., et al.) followed by

212 contig extension using Nanopore reads using SSPACE-LongRead (Boetzer, M., et al.). Super

213 scaffolding of the assembled scaffold was performed using SSPACE (Boetzer, M., et al.) and

214 PLATANUS using the Nanopore and Matepair data. Final draft genome resulted after gap

215 closure by GAPCLOSER and PLATANUS gap_close tool using Illumina data. The genome

216 size was estimated using a k-mer distribution plot using JELLYFISH (Marcais, G., et al.). The

217 repetitive elements were identified in the final assembled draft genome using Repeat Masker

218 tool. The draft genome was hardmasked by using reference genomic repeats of *G. gallus*. The

219 assembly and annotation workflow overview has been represented as Figure 2.

220

**Simple sequence repeats prediction**

222 Final assembled scaffolds were analysed for Simple Sequence Repeats (SSR) identification.

223 SSRs like the di, tri, tetra, penta and hexa-nucleotide repeats in the genome were obtained using

224 MISA (Version 1.0.0).

225

**Genome prediction and annotation**

227 Gene models was predicted on a hard masked draft genome, where the repetitive elements in

228 the draft genome were masked using genomic repeats of *G. gallus* with Repeatmasker tool and

229 further genes were predicted using AUGUSTUS with *G. gallus* as a reference model. The

230 predicted proteins were annotated by using BLASTP (Altshul, S., et al.) against all Avian

231 sequences downloaded from UniProt Protein Database.

232

**Pathway Analysis of the draft genome**

234 The predicted proteins were searched against the KEGG-KAAS server (Moriya, Y., et al.) for

235 pathway analysis. *G. gallus* (chicken), *Meleagris gallopavo* (turkey), *Taeniopygia guttata*

236 (zebra finch), *Falco peregrinus* (peregrine falcon) were used as reference organism for

10

237 pathway identification. The EuKaryotic Orthologous Groups (KOGs) were predicted using

238 homology based approach.

239

**Mitochondrial genome assembly and annotation**

241 The generated scaffolds from the draft assembly were aligned against the *P. cristatus*

242 mitochondria genome and the mapped reads were filtered and stitched using ABACUS

243 software using the same reference (Zhou, TC., et al.). Further gap closure were performed with

244 3-7kb MP reads to generate an complete assembled mitochondrial genome. MITOS (Bernt, M.,

245 et al.) was used for gene annotation. Circular plot generated using GenomeVx

246 (http://wolfe.ucd.ie/GenomeVx/) representing the localization of the gene in the assembled

247 mitochondrial genome.

248

**Phylogenetic tree construction**

250 The assembled Peacock mitochondrial genome was searched against 695 avian mitochondrial

251 genomes downloaded from NCBI. Based on the Blast-N homology results (with query

252 coverage> 100, subject coverage > 95, % identity >85 and with 1% gaps allowed in the

253 sequences). 51 mitochondrial genome sequences along with our assembled mitochondrial

254 genome were filtered. Multiple sequence alignment with default parameters were performed

255 using MUSCLE global sequence aligner. Phylogenetic trees were constructed using IQ-TREE

256 version 1.5.6 (www.iqtree.org). The parameters used for phylogenetic tree construction were

257 ultrafast boostrap (UFBoot, using the –bb option of 1000 replicates), and a standard

258 substitution model (-m MFP) was given for tree generation. The generated trees from IQ-TREE

259 tool were visualized using Figtree (http://tree.bio.ed.ac.uk/software/figtree/) and the Brach-

260 support values were recorded from the output ".treefile". The trees were modified for better

261 visualization under Trees section increasing order nodes were applied.

**Protein domain analysis**

264 Predicted proteins from Peacock, Chicken and Turkey with sequence length greater than 100

265 amino acids were considered for protein domain analysis. All the protein sequences from each

266 organism were searched against Pfam-A database using Pfam scan for protein domain

267 identification.

268

**Avian protein families**

270 The protein sequences of 48 avian genomes was downloaded from the link

271 http://avian.genomics.cn/en/jsp/database.shtml apart from the predicted proteins of the draft

272 genome. Sequences greater than 100 amino acids from all the avian genomes were filtered and

273 concatenated to a single fasta file. These sequences were clustered using CD-HIT (Fu, L., et

274 al.) with 90% alignment coverage for the shorter sequence with a length difference cutoff of

275 0.9. The single copy ortholog gene family present across all organisms and genes unique to

276 peacock were filtered and annotated.

277

**Genome conservation analysis**

279 The assembled draft genome was aligned against the *G. gallus* genome using Chromosomer

280 tool. Draft chromosomes were constructed based on alignment between fragments (scaffolds)

281 using a reference genome. These reordered assembled genome was aligned against the Chicken

282 genome using LAST aligner with NEAR (finding short-and-strong (near-identical)

283 similarities.) parameter allowing for substitution and gap frequencies leading to the

284 identification of orthologs. These query mapped regions were filtered with greater than 1% of

285 the maximum length for visualization using Circos.

286

287

**Results**

**DNA Sequencing data**

Five libraries of 150 bp paired-end from Illumina HiSeq technology were generated. The short-insert reads of 489,114,747 accounted to genome coverage of 146.7X and long-insert reads of 302,884,819 sequences was about 90.9X coverage with a total coverage of 236X. Sequencing of three mate-pairs of 3-5Kb, 5-7Kb of and 7-10Kb yielded 72,915,033, 47,440,144 and 36,464,628 reads respectively with an approximate coverage of 21.9, 14.2 and 10.9 respectively, with a total of 156 million reads of 47X coverage. Oxford Nanopore technology was used to generate 366,323 long reads having of 2,398,560,283 bp with coverage of 2.3X. The complete sequencing was generated to a depth of ~287X from Illumina and Nanopore platforms. The coverage was based on assuming the peacock genome size of about 1 Gb (Table S1).

**Genome assembly**

The assembly was performed on Illumina reads with Abyss *de novo* assembler that resulted in ~932 Mb (mega base) of genome with an N50 of 1639 bp. The extension of the contigs were performed with Nanopore reads which generated scaffolds with N50 of 14,748 bp. Super scaffolding of the assembled scaffold was performed using SSPACE and PLATANUS with MP libraries that generated ~916 Mb genome with the N50 value of 168,140bp. The final gap closer was executed using GAPCLOSER program with MP and PE-LI libraries which generated a draft genome of 1.02 GB (giga base). The draft genome assembly of *Pavo cristatus* consists of 179,346bp scaffolds, with a N50 of 189,886bp with 37 scaffolds having sequence length >=1Mbp. Contigs above 5000 bp have covered a genome of ~0.915 Mb with N50 0.23 Mb. In the assembled genome there were ~0.4% of non ATGC characters (Table 1).

13

313 Complete mitochondrial genome of 16699 bp was obatained. Total of 22 tRNA, tow rRNA and

314 13 protein coding genes were identified in the assembled genome (Fig. 3a). A 100% similarity

315 was observed with the preiously published Peacocks mitochondrial genome (Fig. S1a and S1b).

316

**Transposable Elements (TE)**

318 In the bird genome a total of 75,315,566 bp (7.3% of the genome) was predicted to have 5.5%

319 of Retroelements (with SINEs 0.08% and LINEs 4.71%), 6.25 % total interspersed repeats with

320 0.84 % simple repeats and 0.21% low complexity regions. The DNA transposons identified in

321 the genome was 0.71% (Table S2).

322

**Protein coding gene annotation**

324 A total of 23,153 proteins were predicted in the assembled draft genome using AUGUSTUS.

325 Among them 95% predicted genes were annotated against the other Aves proteins. The 21,854

326 annotated proteins showed top similarity to species *Gallus gallus* (Chicken) with 11,398

327 proteins, *Meleagris gallopavo* (Common turkey) with 4059 proteins, *Amazona aestiva* (Blue-

328 fronted Amazon parrot) with 1352 proteins and *Anas platyrhynchos* (Mallard) (*Anas boschas*)

329 with 849 proteins. Thirteen species had about 100 to 400 annotated proteins. The remaining

330 proteins were in the range of 1 to <100 proteins in about 62 species. From the annotations a

331 total of 13,161 proteins showed similarity to uncharacterized protein annotation and some of

332 the over represented proteins were Tyrosine-protein kinase, Sulfotransferase, Phosphoinositide

333 phospholipase C, Tetraspanin, Phospholipid-transporting ATPase, Olfactory receptor,

334 Polypeptide N-acetylgalactosaminyltransferase, Transporter, Keratin, Hexosyltransferase,

335 Protein Wnt, Kinesin-like protein, Gap junction protein, Claudin, POU domain protein,

336 Sodium/hydrogen exchanger, Phospholipid-transporting ATPase, Histone-lysine N-

337 methyltransferase and others (Table S3). The gene ontology annotations showed to have Gene

338 Ontology (GO) descriptions for 18,295 proteins. Among them, 14,490 proteins have Molecular

339 Function; 11,679 have Biological Process and 13,736 proteins have Cellular Component as

340 functional categories.

341 About 17.7% of proteins were found to have pathway information against the KEGG database

342 (Table S4). Some of the overrepresented annotations were Kinases like MAPK (mitogen-

343 activated protein kinase); JNK (c-Jun N-terminal kinase); RAF (RAF proto-oncogene

344 serine/threonine-protein kinase); AKT (RAC serine/threonine-protein kinase); protein kinases

345 and different GTPases.

346 Proteins searched against the KOG annotations showed a total of 20,937 proteins having

347 annotations. Among them, the most abundant annotations include Zn-finger, transmembrane

348 receptor, ubiquitin ligase, Leucine rich repeat, Cadherin repeats, Serine/threonine protein

349 kinase, Collagens, Ankryin repeat, Fibrillins, Voltage-gated Ca2+ channels and Hormone

350 receptors (Table S5). The peacock proteins when searched against the human proteins showed

351 expansions in ontologies for cell morphogenesis, neuronal projection and development and

352 GTPases (Table S10 and Fig. S4).

353

354

**Simple sequence repeats**

356 A total of 399,493 SSRs were identified from the peacock genome assembly. The largest

357 fraction of SSRs identified were mononucleotide (60.04%), followed by tetra nucleotide

358 (26%), di nucleotide (8.51%), tri nucleotide (4.31%), penta nucleotide (1.03%) and finally hexa

359 nucleotide (0.13%). Among the SSRs identified, A (49.2%) and T (44.9%) accounted for

360 94.1% of the mono-nucleotide repeats. AT (23.8%), TA (16.5%), TG(13.7%), AC(10.6%) and

361 CA (10.32%) accounted for 75% of the di-nucleotide repeats. while TTG (9.9%), AAT (9.6%),

362 AAC (9.4%), TTA (7.1%), ATT (4.5%), TAA (3.5%), CAA (3.1%) and GGA (2.69%)

363 accounted for 49.7% of the tri-nucleotide repeats (Table S6).

364

**Avian protein family analysis**

365

366 A total of 748,544 protein sequences from 49 avian species have 653,497 protein sequences of

367 length above 100 amino acids (Table S7A). A total of 240,853 gene clusters were generated of

368 which 41 gene clusters had single copy orthologs in all avian species (Table S7B and Table

369 S7C). With the above stringent cutoff we observed 15,913 gene clusters were unique to

370 peacock species.

371

**Phylogeny and Genome comparisons**

372

373 The phylogeny of 51 mitochondrial genome sequences along with peacock genome showed a

374 clade consisting of *Pavo* species and *Gallus* (red junglefowl, Sri Lankan junglefowl, grey

375 junglefowl, grey junglefowl), Bambusicola (Mountain Bamboo-partridge, Chinese bamboo

376 partridge) and Francolinus (Chinese francolin). It can be observed many species were

377 distributed in two different clades where 34 species were found in one clade and five species

378 in other clade. Some of the endemic or native bird species like *Arborophila ardens*, *Acryllium*

379 *vulturinum* and *Numida meleagris* were found as clear outgroup of species in this study (Fig.

380 3b).

381 Predicted proteins from Peacock, Chicken and Turkey were searched for protein domain

382 analysis. 81% of the Pfams were common among the three species. About 94%, 98.4% and

383 99.7% predicted Pfam domains were identified in Peacock, Chicken and Turkey respectively.

384 There were 255, 69 and 14 Pfam domains found to be unique in the species mentioned above

385 respectively (Fig. 4).

386 The assembled Peacock genome was reordered for pseudo chromosomes generation against

387 the simple repeat masked Chicken genome (1.21 GB, Warren, WC., et al.) using Chromosomer

388 which generated a overall reordered Peacock genome of about 597MB. The right side of the

389 image represents the reference genome and left side of the image represents the Peacock

390 genome (Fig. 5).

391

**Conclusions**

393 Third generation sequencing in avian genomics where long reads having the substantiality to

394 improve genome assembly will benefit understanding the organisms in the structurally

395 complex regions having repeat elements and isoforms in the genome (Goodwin, S., et al.).

396 Using a combination of short reads of different insert sizes as well as mate pair reads generated

397 from Illumina technology along with long reads from Oxford Nanopore, we obtained an

398 improved assembly and a draft genome of the Indian Blue Peacock (*Pavo cristatus*). In

399 comparison to other avian genomes (Zhang, G., et al.), the current 290X sequencing depth

400 obtained from our study is one of the highest. With a N50 of 0.23Mb, we presented here a

401 reasonably reliable draft genome for the peacock species. The inclusion of Nanopore reads

402 366,323 for scaffolding followed by subsequent gap-closing using Illumina data led to a 26.2%

403 reduction in the number of scaffolds and a 50.65% and 115% increase in the scaffold and contig

404 N50, respectively. With only 2.3X of long reads, a significant improvement in the assembly

405 was observed. On the contrary, the assembly contained less than 0.4% of unknown nucleotides,

406 which is very low for a draft assembly. With the addition of more long reads along with

407 transcriptomic sequencing along with scaffolding and/or gap closure tools, further

408 improvement in the assembly can be achieved.

409

410 Peacock seems to defy the Darwinian laws of natural selection. These concern were raised by

411 no other than Darwin himself. Hence, he proposed the theory of the sexual selection where the

412 female can choose for a male with a certain phenotypic feature such as brilliant color or a long

413 tail (Burgess, S.). Peacock's brilliantly colored long tail feathers seems to evolve at the cost

414 of finding its female partner thereby contributing its beneficial genes, even at the cost of

415 making itself venerable to predators. A female peafowl in turn tends to choose the mate with

416 the largest and decorated plumage, which indirectly reflects its healthiness and capacity to

417 wade off potential competitors thereby selecting the most suitable male. Peacocks beautiful

418 feathers with it all its artistry surely provides it with an advantage to impress the females

419 (Dakin, R., et al.). Understanding the formation of beautiful feathers from the genomic context

420 will help in resolving several evolutionary theories on sexual selection that have been discussed

421 on this species.

422 Pigment particles are embedded into the newly grown peacock feathers during the molting

423 season which seems to absorb light of selective wavelength there by imparting to the color of

424 the plumage (Mercedes Foster, Rennee Riedler, et al.) Pigment morphogenesis in *Gallus* has

425 been reported by a process of melanoblasts migration and colonization into feather bud where

426 they differentiate to produce the pigment, melanin (Kelsh, RN., et al.). The molecular

427 mechanisms that control the pigment cell migration have been narrowed down to proteins Kit1

428 and FGFs which maintain the melanoblasts migration to feathers. Kit and FGF proteins have

429 also been identified in our current study (Table S3). Understanding of these proteins in the

430 patterned formation will help decode the pigment pattern morphology in peacocks. These

431 pigmented patterns play a role in communication, choice of mate and in some species it can

432 help in camouflage (Burgess, S)

433 It has been observed that the variations in the genome size among bird species are very low

434 (Table S8). The genome complexities of a species are influenced by the Transposable elements

435 (TE) that are believed to play a crucial role (Kapsuta, A., et al.). The long read sequences have

436 significantly helped in resolving the TEs in genome quality and assembly. Peacock genome

437 comparisons with Turkey and Chicken have showed closeness to the Chicken species. The

18

438    mitochondrial phylogeny also revealed similar findings. Homology searches have shown

439    several important gene family expansions such as Kinases, Zn finger proteins, GTPases and

440    others. Their roles in biology, development and evolution of the Peacocks have to be further

441    explored.

442    To summarize, we have assembled the *Pavo* genome using Illumina and Oxford nanopore

443    technology. The genome information can be valued and explored by avian enthusiasts to further

444    understand about this bird. Though not critically endangered yet, in India, peafowl population

445    is surely at a decline in the wild due to massive deforestation and habitat loss. Thus is further

446    compounded by increased poaching for meat and feathers. Our genome sequencing initiative

447    of *Pavo cristatus* is not just only from a conservational viewpoint, but also to preserve a

448    heritage associated with this bird that runs through centuries and that bears a strong attachment

449    to the national psyche.

450

**Availability of supporting data**

Supplementary data contains, read statistics, annotation, repeats identification, orthology analysis, mitochondria assembly and annotation. Figures, Gene ontology. DNA and library preparation protocols.

**Raw Data in SRA**

Raw reads (Illumina and Nanopore) are available in the Sequence Read Archive (SRA), and the Whole Genome Shotgun project has been deposited at GenBank under SRA Submission ID: SUB3108024, Bioproject: PRJNA413288 and Biosamples SAMN07739105 : SKPea2016_SI, SAMN07739104 : SKPea2016_LI, SAMN07739101 : FPL_3_5KB, SAMN07739102 : FPL_5_7KB, SAMN07739103 : FPL_7_10KB and SAMN07739107 : FPL_Nano.

**Competing interests**

The author(s) declare that they have no competing interests.

**Authors contributions**

RD, AS, KP performed wet lab experiments; RD designed work plan, experiments and logistics; SS, VR, KP SG IM and AR assisted with the work; RS provided samples from bird; BR, SK performed data analysis and interpretation; SK, BR, RD drafted the manuscript and SK overseen the whole project.

20

**Tables**

484    Table 1. *De novo* assembly statistics of the Peacock genome.

| Description | Contigs | Nanopore Scaffold | Super Scaffolds | GapClosed | >1000 Kb | >5000 Kb |
|---|---|---|---|---|---|---|
| Contigs | 685241 | 281272 | 179346 | 179346 | 34181 | 15026 |
| Maximum Length | 49159 | 251510 | 2390121 | 2488982 | 2488982 | 2488982 |
| Minimum Length | 300 | 5 | 265 | 265 | 1000 | 5000 |
| Average Length | 1360 | 3250 | 5111 | 5729 | - | - |
| Total Length | 932162464 | 914363908 | 916720956 | 1027551907 | 954483822 | 915373606 |
| Length >= 100 bp | 685241 | 281271 | 179346 | 179346 | 34181 | 15026 |
| Length >= 200 bp | 685241 | 281271 | 179346 | 179346 | 34181 | 15026 |
| Length >= 500 bp | 616120 | 186433 | 93727 | 93727 | 34181 | 15026 |
| Length >= 1 Kbp | 363428 | 104479 | 34168 | 34181 | 34181 | 15026 |
| Length >= 10 Kbp | 1591 | 24748 | 9249 | 10311 | 10311 | 10311 |
| Length >= 1 Mbp | 0 | 0 | 27 | 37 | 37 | 37 |
| Non-ATGC # | 350325 | 42696911 | 49169831 | 4034372 | 4032567 | 3978757 |
| Non-ATGC % | 0.038 | 4.67 | 5.364 | 0.393 | 0.422 | 0.435 |
| N50 value | 1639 | 14748 | 168140 | 189886 | 218023 | 232312 |

485

486

487 **Figure legend**

488 **Figure 1.** The beautiful and charismatic photo of Indian blue peacock (*Pavo cristatus)* bird.

489 **Figure 2.** Detailed workflow for *de novo* genome assembly and annotation.

490 **Figure 3a.** Circular representation of Peacock mitochondrial genome with genes predicted.

491 **Figure 3b.** Phylogenetic tree generated from mitochondrial genome from 52 different avian

492 species.

493 **Figure 4.** Circular image of the assembled peacock genome aligned against the *G. gallus*

494 genome using Chromosomer tool. Draft chromosomes were generated by similarity between

495 scaffolds which were arranged on the reference chicken genome. Circos was used for

496 visualization.

497 **Figure 5.** Venn diagram showing common and unique Protein family domains (Pfam) between

498 Peacock, Chicken and Turkey proteins.

499

**References:**

Gadagkar, R., 2003. Is the peacock merely beautiful or also honest?. Current Science, 85(7), pp.1012-1020.

Kushwaha, S., and Kumar, A. 2016. A Review on Indian Peafowl (*Pavo cristatus*) Linnaeus, 1758. Journal of Wildlife and Research, 4, 42-59.

Kadgaonkar, Shivendra B. 1993. The peacock in ancient Indian art and literature. Bulletin of the Deccan College Research Institute, vol. 53, pp. 95–115. JSTOR, www.jstor.org/stable/42936434.

Brickle, N. 2002. Habitat use, predicted distribution and conservation of green peafowl (*Pavo muticus*) in Dak Lak Province, Vietnam. Biological Conservation, 105: 189-197.

Jackson, C. 2006. Peacock. London: Reaktion Books Ltd.

Zhang, G., Jarvis, E. D., and Gilbert, M. T. P. 2014. A flock of genomes. Science 346, 1308–1309.

International Chicken Genome Sequencing Consortium. 2004. Sequence and comparative analysis of the chicken genome provide unique perspectives on vertebrate evolution. Nature, 432(7018), 695.

Kajitani, R., Toshimoto, K., Noguchi, H., Toyoda, A., Ogura, Y., Okuno, M., Yabana, M., Harada, M., Nagayasu, E., Maruyama, H. and Kohara, Y., 2014. Efficient de novo assembly of highly heterozygous genomes from whole-genome shotgun short reads. Genome research, 24(8), pp.1384-1395.

Loman, N. J. and Quinlan, A. R. 2014. Poretools: a toolkit for analyzing nanopore sequence data. Bioinformatics, 30(23), 3399-3401.

Li H. and Durbin R. 2009. Fast and accurate short read alignment with Burrows-Wheeler Transform. Bioinformatics, 25:1754-60.

534  Birol, I., Jackman, S.D., Nielsen, C.B., Qian, J.Q., Varhol, R., Stazyk, G., Morin, R.D., Zhao,

535  Y., Hirst, M., Schein, J.E. and Horsman, D.E., 2009. De novo transcriptome assembly with

536  ABySS. Bioinformatics, 25(21), pp.2872-2877.

537

538  Boetzer, Marten, and Walter Pirovano. 2014. SSPACE-LongRead: scaffolding bacterial draft

539  genomes using long read sequence information. BMC bioinformatics 15.1: 211

540

541  Boetzer, M., Henkel, C.V., Jansen, H.J., Butler, D. and Pirovano, W., 2010. Scaffolding pre-

542  assembled contigs using SSPACE. Bioinformatics, 27(4), pp.578-579.

543

544  Marcais, G and Kingsford, C. 2011. A fast, lock-free approach for efficient parallel counting

545  of occurrences of k-mers. Bioinformatics 27(6): 764-770.

546

547  Altschul, S.F., Gish, W., Miller, W., Myers, E.W. and Lipman, D.J., 1990. Basic local

548  alignment search tool. Journal of molecular biology, 215(3), pp.403-410.

549

550  Moriya, Y., Itoh, M., Okuda, S., Yoshizawa, A.C. and Kanehisa, M., 2007. KAAS: an

551  automatic genome annotation and pathway reconstruction server. Nucleic acids research,

552  35(suppl_2), pp.W182-W185.

553

554  Zhou, T.C., Sha, T., Irwin, D.M. and Zhang, Y.P., 2015. Complete mitochondrial genome of

555  the Indian peafowl (*Pavo cristatus*), with phylogenetic analysis in phasianidae. Mitochondrial

556  DNA, 26(6), pp.912-913.

557

558  Bernt, M., Donath, A., Jühling, F., Externbrink, F., Florentz, C., Fritzsch, G., Pütz, J.,

559  Middendorf, M. and Stadler, P.F., 2013. MITOS: improved de novo metazoan mitochondrial

560  genome annotation. Molecular phylogenetics and evolution, 69(2), pp.313-319.

561

562  Fu, L., Niu, B., Zhu, Z., Wu, S. and Li, W., 2012. CD-HIT: accelerated for clustering the next-

563  generation sequencing data. Bioinformatics, 28(23), pp.3150-3152.

564

565  Warren, W.C., Hillier, L.W., Tomlinson, C., Minx, P., Kremitzki, M., Graves, T., Markovic,

566  C., Bouk, N., Pruitt, K.D., Thibaud-Nissen, F. and Schneider, V., 2016. A new chicken genome

567 assembly provides insight into avian genome structure. G3: Genes, Genomes, Genetics, pp.g3-
568 116.

569

570 Goodwin, S., McPherson, J.D. and McCombie, W.R., 2016. Coming of age: ten years of next-
571 generation sequencing technologies. Nature Reviews Genetics, 17(6), p.333.

572

573 Zhang, G., Li, C., Li, Q., Li, B., Larkin, D.M., Lee, C., Storz, J.F., Antunes, A., Greenwold,
574 M.J., Meredith, R.W. and Ödeen, A., 2014. Comparative genomics reveals insights into avian
575 genome evolution and adaptation. Science, 346(6215), pp.1311-1320.

576

577 Burgess, S., 2001. The beauty of the peacock tail and the problems with the theory of sexual
578 selection. Journal of Creation, 15(2), pp.94-102..

579

580 Dakin, R. 2008. The role of the visual train ornament in the courtship of peafowl, *Pavo*
581 *cristatus*. Masters Abstracts International, 47/03: 97.

582

583 Mercedes S Foster 1975. The overlap of molting and breeding in some birds. The Condor
584 77:304-314.

585

586 Renee Riedler, Christel Pesme, James Druzik, Molly Gleeson, Ellen Pearlstein The Journal of
587 the American Institute of Conservation  2014

588

589

590 Kelsh, R.N., Harris, M.L., Colanesi, S. and Erickson, C.A., 2009. Stripes and belly-spots—a
591 review of pigment cell morphogenesis in vertebrates. Seminars in cell & developmental
592 biology. Vol. 20, No. 1, pp. 90-104.

593

594 Kapusta, A. and Suh, A., 2017. Evolution of bird genomes—a transposon's- eye view. Annals
595 of the New York Academy of Sciences, 1389(1), pp.164-185.

596

597

598 **Webservers:**
599 FastQC : http://www.bioinformatics.babraham.ac.uk/projects/fastqc/
600 AUGUSTUS : http://augustus.gobics.de/

601    PLATANUS : http://platanus.bio.titech.ac.jp/

602    ABACUS : http://abacas.sourceforge.net/Manual.html

603    MISA: http://pgrc.ipk-gatersleben.de/misa/

604    REPEATMASKER : http://www.repeatmasker.org/

605    KOGs : https://genome.jgi.doe.gov/portal/help/kogbrowser.jsf

606    KAAS : http://www.genome.jp/tools/kaas/

607    MITOS : http://mitos.bioinf.uni-leipzig.de/

608    MUSCLE : https://www.drive5.com/muscle/

609    PFAMSCAN : https://www.ebi.ac.uk/seqdb/confluence/display/THD/PfamScan

610    LAST : http://last.cbrc.jp/

611    CHROMOSOMER : https://github.com/gtamazian/chromosomer

612    CIRCOS : http://circos.ca/

Figure 1

Figure 2

Figure 3b

0.03

Figure 4

Figure 5

Mitochondria notes

Click here to access/download
**Supplementary Material**
Supplementary_Notes_Mitochondria.docx

Description of tables and figures

Click here to access/download
**Supplementary Material**
Description of all the tables and figures.docx

Table S1 Read stats, S2 TEs

Table S3 Gene annotations

Click here to access/download
**Supplementary Material**
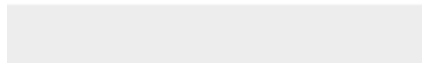Table_S3_Gene_annotation.xlsx

Table S4 KEGG annotations

Click here to access/download
**Supplementary Material**
Table_S4_KEGG_annotation.xlsx

Click here to access/download
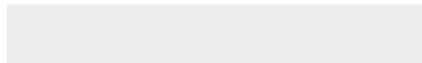**Supplementary Material**
Table_S5_KOG_annotation.xlsx

Table S6 Blast against Human proteins

Table S7 SSR analysis

Click here to access/download
**Supplementary Material**
Table_S7_SSR.xlsx

Click here to access/download
**Supplementary Material**
Table_S8_Orthology.xlsx

Table S9 Mitochondria species description

Click here to access/download
**Supplementary Material**
Table_S9_Species_for_mitochondriaxlsx.xlsx

Table S10 Pfam annotations

Click here to access/download
**Supplementary Material**
Table_S10_Pfam_Analysis.xlsx