

GigaScience

De novo assembly of the Indian Blue Peacock (*Pavo cristatus*) genome using Oxford Nanopore Technology and Illumina sequencing --Manuscript Draft--

Manuscript Number:	GIGA-D-18-00280R4
Full Title:	De novo assembly of the Indian Blue Peacock (<i>Pavo cristatus</i>) genome using Oxford Nanopore Technology and Illumina sequencing
Article Type:	Data Note
Funding Information:	
Abstract:	<p>Background: The Indian peafowl (<i>Pavo cristatus</i>) is native to South Asia and is the national bird of India. Here we present a draft genome sequence of the male blue peacock using Illumina and Oxford Nanopore Technology (ONT).</p> <p>Results: ONT sequencing gave approximately 2.3-fold sequencing coverage, whereas Illumina generated 150-bp paired-end sequence data at 284.6-fold coverage from five libraries. Subsequently, we generated a 0.915-Gb de novo assembly of the peacock genome with a scaffold N50 of 0.23 Mb. We predict that the peacock genome contains 23,153 protein-coding genes and 75.3 Mb (7.33%) of repetitive sequences.</p> <p>Conclusions: We report a high-quality assembly of the peacock genome using a hybrid approach of sequences generated by both Illumina and ONT. The long-read chemistry generated by ONT was useful for addressing challenges related to de novo assembly, particularly at regions containing repetitive sequences spanning longer than the read length, and which could not be resolved with only short-read-based assembly. Contig assembly of Illumina short reads gave an N50 of 1,639 bases, whereas with ONT, the N50 increased by more than nine-fold to 14,749 bases. The initial contig assembly based on Illumina sequencing reads alone gave 685,241 contigs. Further scaffolding on assembled contigs using both Illumina and ONT sequencing reads resulted in a final assembly of 15,025 super-scaffolds, with an N50 of about 0.23 Mb. Ninety-five per cent of proteins predicted by homology matched with those in a public repository, verifying the completeness of our assembly. Like other phylogenetic studies of avian conserved genes, we found <i>P. cristatus</i> to be most closely related to <i>Gallus gallus</i>, followed by <i>Meleagris gallopavo</i> and <i>Anas platyrhynchos</i>. Compared with the recently published peacock genome assembly, the current, superior, hybrid assembly has greater sequencing depth, fewer non-ATGC sequences, and fewer scaffolds.</p>
Corresponding Author:	Subhradip Karmakar, PhD All India Institute of Medical Sciences 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:	Ruby Dhar
First Author Secondary Information:	
Order of Authors:	Ruby Dhar Ashikh Seethy Karthikeyan Pethusamy Vishwajeet Rohil Sunil Singh Kakali Purkayastha Indrani Mukherjee

	Sandeep Goswami
	Rakesh Singh
	Ankita Raj
	Tryambak Srivastava
	Sovon Acharya
	Balaji Rajashekhar, Ph.D
	Subhradip Karmakar, PhD
Order of Authors Secondary Information:	
Response to Reviewers:	<p>Dear Dr. Scott, Kindly find our response to the minor corrections in the manuscript. All the comments have been replied and corrected in the manuscript, figure 2 and supplementary table and figure description.</p> <p>Thank you for all the help and review provided. Regards,</p>
Additional Information:	
Question	Response
Are you submitting this manuscript to a special series or article collection?	No
<p>Experimental design and statistics</p> <p>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.</p> <p>Have you included all the information requested in your manuscript?</p>	Yes
<p>Resources</p> <p>A description of all resources used, including antibodies, cell lines, animals 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.</p> <p>Have you included the information</p>	Yes

<p>requested as detailed in our Minimum Standards Reporting Checklist?</p>	
<p>Availability of data and materials</p> <p>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.</p> <p>Have you have met the above requirement as detailed in our Minimum Standards Reporting Checklist?</p>	<p>Yes</p>

[Click here to view linked References](#)

Dhar et al.

De novo genome assembly of the peacock

1 **De novo assembly of the Indian Blue Peacock (*Pavo***
2 ***cristatus*) genome using Oxford Nanopore Technology and**
3 **Illumina sequencing**

4
5 Ruby Dhar¹, Ashikh Seethy¹, Karthikeyan Pethusamy¹, Sunil Singh¹, Vishwajeet Rohil²,
6 Kakali Purkayastha², Indrani Mukherjee¹, Sandeep Goswami¹, Rakesh Singh³, Ankita Raj¹,
7 Tryambak Srivastava¹, Sovon Acharya¹, Balaji Rajashekar^{4,5*}, Subhradip Karmakar^{1*}

8
9 ¹Department of Biochemistry, Room No 3020, AIIMS, New Delhi, India

10 ²Vallabhbhai Patel Chest Institute (VPCI), Delhi University, New Delhi, India

11 ³Kanpur Zoo, Hastings Ave, Azad Nagar, Nawabganj, Kanpur, Uttar Pradesh 208002, India
12 India

13 ⁴Institute of Computer Science, University of Tartu, 50409 Tartu, Estonia

14 ⁵Celixa, Bangalore, 560020, India

15
16 *Correspondence address: Balaji Rajashekar, Institute of Computer Science, University of
17 Tartu, 50409 Tartu, Estonia; Tel: +91-9844677993; Email: balaji@ut.ee

18 Subhradip Karmakar, Department of Biochemistry, Room No 3020, AIIMS, New Delhi,
19 India; Tel: +91-9999612564; Email: subhradip.k@aiims.edu

20
21 **ORCID IDs:** Ruby Dhar: 0000-0003-3600-6554; Ashikh Seethy: 0000-0001-6825-5753;

22 Tryambak Pratap Srivastava: 0000-0002-7903-5876; Balaji Rajashekar: 0000-0002-1665-
23 5584; Subhradip Karmakar: 0000-0002-4757-8729

24

1 **Abstract**

2 **Background:** The Indian peafowl (*Pavo cristanus*) is native to South Asia and is the national
3 bird of India. Here we present a draft genome sequence of the male blue peacock using
4 Illumina and Oxford Nanopore Technology (ONT).

5 **Results:** ONT sequencing gave approximately 2.3-fold sequencing coverage, whereas
6 Illumina generated 150-bp paired-end sequence data at 284.6-fold coverage from five
7 libraries. Subsequently, we generated a 0.915-Gb de novo assembly of the peacock genome
8 with a scaffold N50 of 0.23 Mb. We predict that the peacock genome contains 23,153
9 protein-coding genes and 75.3 Mb (7.33%) of repetitive sequences.

10 **Conclusions:** We report a high-quality assembly of the peacock genome using a hybrid
11 approach of sequences generated by both Illumina and ONT. The long-read chemistry
12 generated by ONT was useful for addressing challenges related to de novo assembly,
13 particularly at regions containing repetitive sequences spanning longer than the read length,
14 and which could not be resolved with only short-read-based assembly. Contig assembly of
15 Illumina short reads gave an N50 of 1,639 bases, whereas with ONT, the N50 increased by
16 more than nine-fold to 14,749 bases. The initial contig assembly based on Illumina
17 sequencing reads alone gave 685,241 contigs. Further scaffolding on assembled contigs using
18 both Illumina and ONT sequencing reads resulted in a final assembly of 15,025 super-
19 scaffolds, with an N50 of about 0.23 Mb. Ninety-five per cent of proteins predicted by
20 homology matched with those in a public repository, verifying the completeness of our
21 assembly. Like other phylogenetic studies of avian conserved genes, we found *P. cristatus* to
22 be most closely related to *Gallus gallus*, followed by *Meleagris gallopavo* and *Anas*
23 *platyrhynchos*. Compared with the recently published peacock genome assembly, the current,
24 superior, hybrid assembly has greater sequencing depth, fewer non-ATGC sequences, and
25 fewer scaffolds.

1 *Keywords:* Peacock, *Pavo cristatus*, Indian National Bird, Genome Assembly, Oxford
2 Nanopore.

4 **Data description**

5 **Background**

6 *Pavo cristatus*, commonly known as the Indian Blue Peafowl, is native to South Asian
7 countries. Apart from the wild, they are usually found as park and zoo exhibits, or are raised
8 for breeding and conservation purposes [1, 2] (Fig. 1). Peafowl have been widely referred to
9 in ancient Indian literature [3] and are closely associated with the life and culture of
10 Southeast Asian, symbolizing beauty, love, grace and pride [4, 5]. For these reasons, the
11 peafowl – specifically the peacock – was chosen to be the national bird of India in 1963.

12 Genome sequencing of the avian model organism *Gallus gallus* (the red junglefowl, or
13 chicken) [6] and other avian species [7] has provided novel perspectives on vertebrate
14 genome evolution, such as a better understanding of genome structure and annotating the
15 mammalian genome. Genome studies of *G. gallus* have revealed high conservation within
16 orthologous regions of the human genome [8], thus showing promise as a good candidate for
17 studies on developmental biology, immunology and vertebrate genome architecture [9, 10].

18 Despite a wealth of information from existing avian genome sequencing projects, it remains
19 important to sequence the genomes of other species to add value to avian and vertebrate
20 genomics. Here, we use Oxford Nanopore Technology (ONT) to sequence a bird genome for
21 the first time. The long reads generated from this sequencing technology were helpful during
22 the de novo assembly of this genome, especially in the GC-rich repeat regions, which
23 invariably pose serious challenges. By comparing this genome with those of other birds, we
24 will understand more about the uniqueness of the peacock genome; the development of this
25 species, its complex plumage pigmentation and sexual dimorphism; and its evolutionary

1 relationships with other birds. Characterization of genes and their specific functions will
2 facilitate better understanding of the peafowl species. By comparing proteins between the
3 peacock, chicken and *Meleagris gallopavo* (domestic turkey), conserved domains and
4 functional annotations may be revealed.

6 **Methods**

7 **Sample collection and extraction of DNA**

8 Blood was collected from an Indian male peacock (Figure 1) at Kanpur Zoo, India, after
9 obtaining the necessary ethical and institutional approvals.

10 DNA from blood was prepared for sequencing as follows: firstly, 200 µl of blood was added
11 to a 1.5-ml microcentrifuge tube containing approximately 20 µl of proteinase K (PK)
12 solution, and briefly mixed. Cell lysis buffer (200 µl) was added to the tube, which was
13 mixed by vortexing for 10 seconds, then incubated at 56°C for 10 minutes. Then, 250 µl of
14 binding buffer (BBA) was added to the tube, which was mixed by vortexing again for
15 10 seconds. The contents of the tube were added to a ReliaPrep™ (Promega, Madison, USA)
16 binding column, which had been placed into an empty collection tube, then capped and
17 placed in a refrigerated microcentrifuge. The binding column and tube were then centrifuged
18 for 1 minute at 12000 rpm and flow-through was discarded. The binding column was placed
19 into a fresh collection tube, 500 µl of column wash solution was added, and then centrifuged
20 for 3 minutes at maximum speed, again discarding flow-through. Column washing was
21 repeated three times. Columns were then placed in a clean, nuclease-free 1.5-ml
22 microcentrifuge tube. Nuclease-free water (100 µl) was then added to the column and
23 centrifuged for 1 minute more at maximum speed before discarding the column and saving
24 the elute.

1 The concentration and purity of the extracted DNA was evaluated using a Nanodrop 2000
2 spectrophotometer (Thermo Fisher Scientific, MA, USA) and a Qubit fluorometer (Thermo
3 Fisher Scientific, MA, USA), and integrity was checked on 0.8% agarose gel. The DNA
4 sample was aliquoted for library preparation on two different platforms: Illumina HiSeq 2000
5 (Illumina, CA, USA) and Oxford Nanopore Technology (ONT) (Oxford, UK) MinION
6 sequencing platform. The genome sequencing was performed by Genotypic Technology,
7 Bengaluru, India in accordance to standard protocols.

9 **Library preparation and sequencing**

10 *Paired-end library preparation and sequencing*

11 Whole genome sequencing (WGS) libraries were prepared with an Illumina-compatible
12 NEXTflex DNA sequencing kit (BIOO Scientific, Austin, TX, USA). Approximately 1 µg of
13 genomic DNA was sheared using a Covaris S2 sonicator (Covaris, Woburn, MA, USA) to
14 generate fragment sizes of approximately 300–600 bp. The fragment size distribution was
15 checked using an Agilent 2200 TapeStation system with D1000 DNA screen tapes and
16 reagents (Agilent Technologies, Palo Alto, CA, USA), and subsequently purified using
17 HighPrep magnetic beads (Magbio Genomics Inc, USA). The purified fragments were end-
18 repaired, adenylated and ligated to Illumina multiplex barcode adaptors, as per the NEXTflex
19 DNA sequencing kit protocol (BIOO Scientific, Austin, TX, USA).

20 The adapter-ligated DNA was purified with HighPrep beads (MagBio Genomics, Inc,
21 Gaithersburg, MD, USA), then size selected on 2% low melting agarose gel, and cleaned
22 using a MinElute column (QIAGEN). The resulting fragments were amplified for 10 cycles
23 of polymerase chain reaction (PCR) using the Illumina-compatible primers provided in the
24 NEXTFlex DNA sequencing kit. The final PCR product (sequencing library) was purified
25 with HighPrep beads, followed by a library quality control check. The Illumina-compatible

1 sequencing library was initially quantified using a Qubit fluorometer (Thermo Fisher
2 Scientific, MA, USA), and fragment size distribution was analyzed on an Agilent
3 TapeStation. Finally, the sequencing library was quantified by quantitative PCR (qPCR)
4 using the Kapa Library Quantification Kit (Kapa Biosystems, Wilmington, MA, USA). The
5 qPCR-quantified library was sequenced on an Illumina sequencer for 150-bp paired-end
6 chemistry.

7 For each sample, the Illumina-compatible sequencing library had a fragment size range of
8 275–425 bp for paired-end short inserts (PE-SI), and 350–650 bp for paired-end long inserts
9 (PE-LI). As the combined adapter size was approximately 120 bp, the effective user-defined
10 insert size was 155–305 bp and 230–530 bp for PE-SI and PE-LI, respectively. Libraries
11 were sequenced using the Illumina HiSeq platform [11] with 150 PE chemistry.

12 *Mate-pair library preparation and sequencing*

13 The mate-pair sequencing library was prepared using the Illumina-compatible NextEra Mate
14 Pair Sample Preparation Kit (Illumina Inc., Austin, TX, USA). Approximately 4 µg of
15 genomic DNA was simultaneously fragmented and tagged with mate-pair adapters in a
16 transposon-based tagmentation step. Tagmented DNA was then purified using AMPure XP
17 magnetic beads (Beckman Coulter Life Sciences, Indianapolis, IN, USA), followed by strand
18 displacement to fill gaps in the tagmented DNA. Strand-displaced DNA was further purified
19 with AMPure XP beads before size-selecting fragments of 3–5 Kb, 5–7 Kb and 7–10 Kb on
20 low melting agarose gel. The fragments were circularized in an overnight blunt-end intra-
21 molecular ligation step, which resulted in circularization of DNA with the insert mate-pair
22 adapter junction. Circularized DNA was sheared using a Covaris S220 sonicator (Covaris,
23 Woburn, MA, USA) to generate approximate fragment sizes of 300–1000 bp. The sheared
24 DNA was purified to collect the mate-pair junction-positive fragments using Dynabeads M-

1 280 streptavidin magnetic beads (Thermo Fisher Scientific, Waltham, MA, USA). The
2 purified fragments were end-repaired, adenylated and ligated to Illumina multiplex barcode
3 adaptors, as per the NextEra Mate Pair Sample Preparation Kit protocol.

4 The adapter-ligated DNA was then amplified for 15 cycles of PCR using Illumina-compatible
5 primers. The final PCR product (sequencing library) was purified with AMPure XP beads,
6 followed by a library quality control check. The Illumina compatible sequencing library was
7 initially quantified using a Qubit fluorometer (Thermo Fisher Scientific, MA, USA), and its
8 fragment size distribution was analyzed with an Agilent TapeStation. Finally, the sequencing
9 library was accurately quantified by qPCR using the Kapa Library Quantification Kit (Kapa
10 Biosystems, Wilmington, MA, USA). The qPCR-quantified libraries were pooled in
11 equimolar amounts to create a final multiplexed library pool for sequencing on an Illumina
12 sequencer.

14 *ONT MinION library preparation and sequencing*

15 Genomic DNA (1.5 µg) was end-repaired using the NEBnext Ultra II End Repair kit (New
16 England Biolabs, MA, USA), and cleaned up with 1x AmPure beads (Beckmann Coulter,
17 USA). Adapter ligations were performed for 20 minutes using NEB blunt/TA ligase (New
18 England Biolabs, MA, USA). The library mixtures were cleaned up using 0.4X AmPure
19 beads (Beckmann Coulter, USA), and eluted in 25 µl of elution buffer. The eluted library was
20 used for sequencing. Whole genome libraries were prepared using the ligation sequencing
21 SQK-LSK108 Oxford Nanopore sequencing kit (ONT, Oxford, UK). Sequencing was
22 performed on a MinION Mk1b (ONT, Oxford, UK) using SpotON flow cell (FLO-MIN106)
23 in a 48-hour sequencing protocol on MinKNOW (version 1.1.20, ONT, Oxford, UK).

25 **Raw data quality control and processing**

1 *Illumina raw data: quality control and processing*

2 Illumina reads were de-multiplexed using bcl2fastq (Illumina). Raw genomic library data
3 generated by Illumina was quality-checked using FastQC (FastQC, RRID:SCR_014583)
4 [12]. Paired-end Illumina reads were processed for clipping adapter and low quality bases
5 using a customized script that retains a minimum of 70% bases/reads with Phred score ($Q \geq 30$
6 in each base position) with a read length of 50 bp. Mate-pair libraries were trimmed for
7 adapter sequences and low-quality bases, trimming from the 3-end using the PLATANUS
8 internal trimmer (Platanus version 1.2.4, RRID:SCR_015531)[13].

10 *ONT reads: base calling and processing*

11 Raw data were base-called with the cloud-based Metrichor workflow 2D Basecalling plus
12 Barcoding (Metrichor version 2.43.1, ONT, Oxford, UK [14]. ONT reads were processed
13 using Poretools [15] to convert fast5 files to fasta format. The 2D reads or 1D high quality
14 reads were selected for further assembly.

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

17 Quality-checked ONT reads were error-corrected using Illumina PE reads. For error-
18 correction, the Illumina PE reads were aligned to the ONT reads using BWA aligner (BWA ,
19 RRID:SCR_010910) [16]. Paired-end reads were assembled using Abyss (ABYSS,
20 RRID:SCR_010709) [17], followed by contig extension using ONT reads using SSPACE-
21 LongRead [18]. Super-scaffolding of the assembled scaffold was performed using SSPACE
22 (SSPACE, RRID:SCR_005056) [19] and PLATANUS on the ONT and mate-pair data. A
23 final draft genome resulted after gap closure using GAPCLOSER (GapCloser,
24 RRID:SCR_015026) [20] and the PLATANUS gap_close tool, with Illumina data. The

1 genome size was estimated with a k-mer distribution plot using JELLYFISH (Jellyfish,
2 RRID:SCR_005491) [21]. The assembly and annotation workflow is shown in Figure 2.

4 **Identification of repetitive elements and simple sequence repeat (SSR) markers**

5 Repetitive elements, retrotransposons and DNA transposons were identified in the draft
6 genome, and hard-masked by using reference genomic repeats of *G. gallus* using
7 Repeatmasker (RRID:SCR_012954) [22]. Final assembled scaffolds were analyzed to
8 identify simple sequence repeats (SSRs). SSRs, such as di-, tri-, tetra-, penta- and hexa-
9 nucleotide repeats in the genome, were identified using MISA (version 1.0.0) [23].

11 **Annotation of the draft genome**

12 Gene models were predicted on a hard-masked draft genome using AUGUSTUS
13 (RRID:SCR_008417) [24], with *G. gallus* as a reference model. Predicted proteins were
14 annotated using BLASTP (RRID:SCR_001010) [25] against the National Center for
15 Biotechnology Information (NCBI)'s NR (non-redundant) database, with default parameters
16 at an E-value cutoff of 1E-5.

17 Predicted proteins were searched against the Kyoto Encyclopedia of Genes and Genomes'
18 Automatic Annotation Server (KEGG-KAAS) for pathway analysis [26]. *G. gallus*, *M.*
19 *gallopavo*, *Taeniopygia guttata* (zebra finch), and *Falco peregrinus* (peregrine falcon) were
20 used as reference organisms for pathway identification. EuKaryotic Orthologous Groups
21 (KOGs) [27] were predicted using a homology-based approach.

23 **Prediction of protein domains**

1 Predicted proteins from peacock, chicken and turkey, with sequence lengths greater than 100
2 amino acids, were considered for protein domain analysis. All protein-coding sequences from
3 each organism were searched against the Pfam-A database using Pfam scan [28].

5 **Identification of avian protein families**

6 A total of 748,544 protein sequences from 49 avian species (including peacock proteins from
7 this study) and others were downloaded from the Avian Phylogenomics Project [29, 30].
8 Sequences with lengths greater than 100 amino acids from all the avian genomes were
9 selected and concatenated to a single fasta file. These sequences were clustered using CD-
10 HIT [31], with 70% alignment coverage for the shorter sequences, with a length difference
11 cutoff of 0.7. Single-copy gene family orthologs present across all avian species, and not
12 clustered peacock proteins, were annotated.

14 **Phylogenetic tree construction**

15 To construct a phylogenetic tree, we considered single-copy gene clusters present as single
16 copies in all the avian species analyzed. These protein sequences from each species were
17 concatenated and further aligned using the multiple sequence alignment tool Clustalw [32].
18 Poorly aligned positions and divergent regions were removed using Gblock [33]. Sequences
19 in fasta format were converted to phylip format using Phylip [34]. Phylogenetic trees were
20 constructed using IQ-TREE (version 1.5.6) [35]. The parameters used to construct the
21 phylogenetic tree were ultrafast bootstrap (UFBoot, using the `-bb` option of 1000 replicates),
22 and a standard substitution model (`-st AA -m TEST`), and `alrt 1000 -nt AUTO` was given to
23 generate the tree. Trees generated from IQ-TREE were visualized using FigTree [36], and the
24 branch-support values were recorded from the output `'treefile'`. For better visualization, trees
25 were modified under the `'Trees'` section, and increasing order nodes were applied.

1

2 **Genome conservation analysis**

3 Draft chromosome visualizations were constructed by aligning the assembled peacock
4 genome against that for *G. gallus* using the Chromosomer tool [37]. The reordered,
5 assembled genome was aligned to the chicken genome using LAST aligner [38], with NEAR
6 (finding short-and-strong [near-identical] similarities) parameters to allow for substitution
7 and gap frequencies, leading to the identification of orthologs. For visualization, these query-
8 mapped regions were filtered for >1% of the maximum length using Circos [39].

9

10 **Results**

11 **Genome sequencing assessment**

12 Five libraries were generated from 150-bp paired-end Illumina sequences. Short-insert reads
13 (489,114,747) represented genome coverage of 146.7x, and 302,884,819 long-insert reads
14 represented about 90.9x coverage, with a total coverage of 237.6x. Sequencing of three mate-
15 pairs of 3–5 Kb, 5–7 Kb and 7–10 Kb yielded 72,915,033, 47,440,144 and 36,464,628 reads,
16 respectively, with an approximate coverage of 21.9x, 14.2x and 10.9x, respectively, and a
17 grand total of 156 million mate-pair reads representing 47x coverage.

18 ONT was used to generate 366,323 long reads, having 2,398,560,283 bp and coverage of
19 2.3x. The complete genome was sequenced to a depth of ~287x, using both Illumina and
20 ONT platforms (Table 1). Coverage was based on the assumption that the peacock genome is
21 1 Gb in size.

22

23 **Genome assembly**

24 The first assembly was based on Illumina reads only, using the Abyss de novo assembler,
25 which resulted in a genome size of ~932 Mb and an N50 of 1639 bp. Contig extension was

1 performed using ONT-generated reads, which gave scaffolds with an N50 of 14,748 bp.
2 SSPACE AND PLATANUS were used to super-scaffold the assembled scaffold with mate-
3 pair libraries, which generated a genome size of ~916 Mb and an N50 of 168,140 bp. Finally,
4 gaps were closed using GAPCLOSER with mate-pair and PE-LI libraries, which generated a
5 draft genome size of 1.02 Gb.

6 The draft genome assembly of *P. cristatus* comprises 179,346-bp scaffolds, with an N50 of
7 189,886 bp with 37 scaffolds, having a sequence length ≥ 1 Mb. Contigs greater than 5000 bp
8 in length covered a genome of ~0.915 Mb, with an N50 of 0.23 Mb. In the assembled
9 genome, there were ~0.4% non-ATGC characters (Table 2).

11 Repetitive genome elements and SSR markers

12 It was estimated that 75 Mb (7.33%) of the peacock genome consisted of repeat sequences
13 (Table S1). About 56 Mb (5.5%) of class I retrotransposons were identified (long interspersed
14 nuclear elements [LINEs], 4.7%; short interspersed nuclear elements [SINEs], 0.08%; and
15 total LTR elements, 0.72%). Subsequently, 7,277,390 bp (0.71%) class II DNA transposons
16 and 467,719 (0.05%) unclassified elements were identified (Table S1). The median
17 percentages of LINEs, SINEs, LTR, DNA, unknown and total masked bases of other avian
18 birds were 3.94, 0.11, 1.31, 0.22, 0.85 and 6.93, respectively (Table S2). A total of 399,493
19 SSRs were obtained from the peacock genome assembly. The largest fraction of SSRs
20 identified were mononucleotides (60.04%), followed by tetranucleotides (26%), dinucleotides
21 (8.51%), trinucleotides (4.31%), pentanucleotides (1.03%), and hexanucleotides (0.13%).
22 Among these SSRs, A (49.2%) and T (44.9%) accounted for 94.1% of the mononucleotide
23 repeats. AT (23.8%), TA (16.5%), TG (13.7%), AC (10.6%) and CA (10.32%) accounted for
24 75% of the dinucleotide repeats, whereas TTG (9.9%), AAT (9.6%), AAC (9.4%), TTA

1 (7.1%), ATT (4.5%), TAA (3.5%), CAA (3.1%) and GGA (2.69%) accounted for 49.7% of
2 the trinucleotide repeats (Table S3).

4 **Gene prediction and annotation**

5 A total of 23,153 proteins were predicted from the assembled draft peacock genome using
6 AUGUSTUS. Of these, 21,854 (94.4%) predicted proteins showed homology to other
7 sequences from the NCBI NR database (Fig. 3). The top four organisms with which peacock
8 proteins showed homology were *G. gallus* (11,398 proteins), *M. gallopavo* (4,059 proteins),
9 *Amazona aestiva* (blue-fronted Amazon parrot; 1352 proteins), and *Anas platyrhynchos*
10 (mallard duck; 849 proteins). Detailed annotations of all proteins are available in Table S4.

11 Gene Ontology (GO) descriptions were assigned for 18,294 (79%) peacock proteins. Of
12 these, 14,489 proteins were identified as having molecular function; 11,678 as biological
13 processes, and 13,735 proteins as cellular components (Table S4).

14 A total of 4,091 (17.7%) peacock proteins had pathway information from the KEGG database
15 (Table S5), whereas 20,937 (88.1%) peacock proteins were similar to KOG annotations
16 (Table S6). When peacock proteins were searched against human proteins, gene family
17 expansions were found in cell morphogenesis, neuronal projection and development and
18 GTPases (Table S7 and Fig. S3).

20 **Analysis of avian protein families**

21 From a total of 748,544 protein sequences from 49 avian species, 653,497 protein sequences
22 were found to have a length of 100 amino acids or greater (Table S8A). Based on their level
23 of identity, CD-HIT clustered the proteins into 114,121 gene clusters. Of these, 68 highly
24 homologous gene clusters were present as single copies in all the 49 avian species (Table

1 S8B and Table S8C). We also observed 13,860 peacock protein clusters that were not
2 clustered with other avian species (Table S8D).

4 **Phylogenetic analysis**

5 Phylogenetic analysis of 48 avian species and peacock proteins showed *P. cristatus* to be
6 clustered in a clade with *G. gallus*, *M. gallopavo*, *A. platyrhynchos*, *Tinamus guttatus* (white-
7 throated tinamou), and *Struthio camelus* (ostrich). This is the largest clade, with six species,
8 having bootstrap support of 100. All species within this clade, except the mallard duck, are
9 flightless or low flying birds. Bootstrap support between *P. cristatus* and *G. gallus* was 96,
10 followed by *M. gallopavo*, with bootstrap support of 100 (Fig. 4).

12 **Comparison with other species and databases**

13 When searching Pfam for conserved protein domains between the predicted proteins from
14 peacock, chicken and turkey, it was revealed that about 81% of domains were common to
15 these three species (Fig. 5, Table S9). Compared with the total number of Pfam domains from
16 these three species, 94%, 98.4% and 99.7% Pfam domains were present in peacock, chicken
17 and turkey, respectively, but 255, 69 and 14 Pfam domains were absent between the species
18 comparisons, respectively (Table S9H).

19 There were 15,470 (78%), 12,794 (85%) and 11,745 (85%) of the peacock, chicken and
20 turkey proteins were found to match with Pfam domains, respectively (Table S9). Domain
21 comparisons between these species showed gene family expansions such as kinases, zinc
22 finger proteins, GTPases, and others, in either one of the species (Fig. 6).

23 A total of 9,974 peacock proteins were annotated in all four databases (NCBI NR, KOG,
24 Pfam and GO) (Fig. 7). When reordered for the generation of pseudo-chromosomes, 597 Mb

1 of the assembled peacock genome was reordered peacock genome compared with the 1.21-
2 Gb masked chicken genome [40] (Fig. 8).

3 Around 60 different avian species have been sequenced using various sequencing
4 technologies (Table S10). The depth of these sequences varies, from as low as 6x to as high
5 as 390x coverage. These results, which were obtained using different bioinformatics methods
6 to assemble the sequencing data, are measured as scaffold N50; i.e., from 30 Kb to 14 Mb.

8 **Discussion and conclusions**

9 In recent years, there has been a rapid surge in the de novo genome sequence assembly of
10 diverse species [41]. This surge has largely been driven by a more affordable cost per base
11 sequencing, and the development of smarter algorithms that have been refined and equipped
12 to handle large datasets [42–44]. The challenge for newer genome analysis pipeline is to
13 generate assemblies with lower contig numbers and longer contigs per genome. To achieve
14 this, technologies that generate longer reads or greater read depths are very helpful [45]; but
15 the use of combinations of different sequencing technologies also plays a significant role in
16 improving genome assemblies [46] (Table S10). Libraries generated using more than one
17 type of chemistry have been found to generate superior assemblies [47], and have been
18 shown to reduce the number of scaffolds – even with very low coverage. Thus, we need to
19 consider combinations of sequencing technologies, along with the use of different
20 bioinformatics software programs, to obtain assemblies with fewer numbers of scaffolds, or
21 which are closer to chromosome-level sequencing [48].

22 Compared with other avian genomes [49], the sequencing depth of 290x that we achieved for
23 the peacock is one of the highest. The final draft peacock genome assembly resulted in an
24 N50 of 0.23 MB. Including 2.3x of reads generated by ONT helped to improve the assembly

1 by reducing the number of scaffolds by 26.2% and increasing the scaffold and contig N50s
2 by about 50.7% and 115%, respectively.

3 The draft assembly contained fewer than 0.4% unknown nucleotides, which is very low for a
4 draft assembly. Our hybrid peacock assembly outperforms the currently available draft
5 peacock assembly (Table S11) by sequencing six different libraries, including long reads
6 from ONT, and 2.1-fold increased sequencing data generation. Greater sequencing depth and
7 the use of multiple libraries enabled us to obtain a better assembly with 6.6-fold fewer
8 scaffolds and an improvement in N50 length by 9.1-fold. The longest scaffold in our
9 assembly is 8.7-fold longer than in the previously published draft assembly, and has a 5-fold
10 lower percentage of non-ATGC. Thus, for the first time in avian genomics, we have
11 demonstrated how low-cost, third-generation sequencing data generated by ONT can help to
12 improve draft genome assembly. Assemblies with longer scaffolds will further help us to
13 understand more about organisms with structurally complex genomic regions, repeat
14 elements and isoforms [39].

15 Our confidence in the peacock proteins predicted from our assembly was strengthened when
16 we discovered that about 95% of them showed significant homology to various genomic
17 features from different databases (Fig. 7). Based on proteins conserved across the avian
18 species, our phylogenetic analysis revealed that the peacock is most closely related to the
19 chicken, followed by turkey and duck. This concurs with previous data based on
20 mitochondrial phylogeny [50]. Thus, our genome sequence provides further insight into the
21 peacock's genetic lineage and evolution with respect to other avian species. The estimated
22 median divergence time of *P. cristatus* from *G. gallus* is about 35 million years ago (MYA),
23 whereas the divergence time estimated between *G. gallus* and *M. gallopavo* is about 37 MYA
24 [51]. The large gap between peacock and other avians may be attributed to the non-

1 availability of other avian genome sequences. The gap may be closed by sequencing other
2 avian species.

3 Among the vertebrates, it has been observed that variations in transposable elements (TEs)
4 between avians are very low [52] (Table S8). The genome complexities of a species are
5 influenced by the TEs that are believed to play a crucial role [53]. In this peacock genome
6 assembly, inclusion of ONT long-read sequences has significantly improved the assembly,
7 thus helping to resolve repetitive regions across the genome. The roles of TEs in the
8 development and evolution of the peacock should be further explored.

9 Information about the peacock genome will be valued, and may be explored, by avian
10 enthusiasts to further understand the avian world. Though not yet critically endangered in
11 India, the wild peafowl population is declining because of massive deforestation, habitat loss
12 [54], and increased poaching for their meat and feathers. Our *P. cristatus* genome sequencing
13 initiative is not only valuable from a conservational viewpoint, but also to preserve the
14 history and heritage that is associated with this bird, which bears a strong attachment to the
15 national psyche.

17 **Availability of supporting data**

18 The data sets supporting the results of this article are available on the study website [55] and
19 the *GigaScience* GigaDB repository [58].

20 Raw reads (Illumina and ONT) are available in the Sequence Read Archive (SRA) database,
21 and the whole genome shotgun project has been deposited at GenBank under SRA
22 submission ID SUB3108024, Bioproject PRJNA413288, and biosamples
23 SUB3108018/SAMN07739105:SKPea2016_SI,
24 SUB3108017/SAMN07739104:SKPea2016_LI,
25 SUB3107930/SAMN07739101:FPL_3_5KB, SUB3108015/SAMN07739102:FPL_5_7KB,

1 SUB3108016/SAMN07739103:FPL_7_10KB, and

2 SUB3108020/SAMN07739107:FPL_Nano (Table 1).

3 The de novo genome assembly can be accessed under SUB4504869/SAMN07739105.

4

5 **Declarations**

6 **List of abbreviations**

7 bp, base pair; Gb, Gigabase pairs; Kb, Kilobase pairs; LINE, Long interspersed nuclear
8 elements; LTR, Long terminal repeat; Mb, Megabase pairs; ONT, Oxford Nanopore
9 Technology; Pfam, Protein families; SINE, Short interspersed nuclear elements; SSR, Simple
10 sequence repeat; WGS, whole genome sequencing.

11

12 **Ethics approval and consent to participate**

13 Not applicable.

14

15 **Consent for publication**

16 Not applicable

17

18 **Competing interests**

19 The authors declare that they have no competing interests.

20

21 **Funding**

22 R.D. provided partial funding. This work was also supported by Estonian Research Council

23 ETAG grants IUT 34–4 for B.R.

24

25 **Authors' contributions**

1 R.D., A.S. and K.P. performed the wet lab experiments; R.D. designed the work plan,
 2 experiments and logistics; S.S., V.R., K.P., S.G., I.M. and A.R. assisted with the work; R.S.
 3 provided bird samples; B.R. and S.K. analyzed and interpreted the data, and drafted the
 4 manuscript. S.K. oversaw the whole project. All authors read and approved the final version
 5 of the manuscript.

7 Acknowledgements

8 We acknowledge the Department of Biochemistry, AIIMS, New Delhi, India, for providing
 9 space and infrastructure to carry on the work. We would like to thank Genotypic Technology
 10 for providing sequencing services.

12 Tables

13 Table 1. Raw data statistics of peacock genome reads generated by Illumina HiSeq and
 14 Oxford Nanopore Technology

Sample	Platform	Library and chemistry	Number of reads	Coverage	SRA ID
SO_6221_SKPea2016_SI	HiSeq	PE – SI (150 * 2)	489114747	146.73	SUB3108018, SAMN07739105
SO_6221_SKPea2016_LI	HiSeq	PE – LI (150 * 2)	302884819	90.87	SUB3108017, SAMN07739104
SO_6221_FPL_3_5KB	HiSeq	MP (150 * 2)	72915033	21.87	SUB3107930, SAMN07739101
SO_6221_FPL_5_7KB	HiSeq	MP (150 * 2)	47440144	14.23	SUB3108015, SAMN07739102
SO_6221_FPL_7_10KB	HiSeq	MP (150 * 2)	36464628	10.94	SUB3108016, SAMN07739103
SO_6221_NP	ONT	5 - 341124	366323	2.3	SUB3108020, SAMN07739107

1 Abbreviations: KB, kilobases; LI, long insert; MP, mate-pair; ONT, Oxford Nanopore Technology; PE, paired-
2 end; SI, short insert; SRA, Sequence Read Archive

3 Table 2. De novo assembly statistics of the peacock genome

Description	Contigs	ONT scaffolds	Super-scaffolds	GapClosed	>1000 Kb	>5000 Kb
Contigs	685,241	281,272	179,346	179,332	34,178	15,025
Maximum length	49,159	251,510	2,390,121	2,488,982	2,488,982	2,488,982
Minimum length	300	5	265	265	1,000	5,000
Average length	1,360	3,250	5,111	5,729	-	-
Total length	932,162,464	914,363,908	916,720,956	1,027,510,962	954,449,349	915,342,012
Length \geq 100 bp	685,241	281,271	179,346	179,332	34,178	15,025
Length \geq 200 bp	685,241	281,271	179,346	179,332	34,178	15,025
Length \geq 500 bp	616,120	186,433	93,727	93,718	34,178	15,025
Length \geq 1 Kb	363,428	104,479	34,168	34,178	34,178	15,025
Length \geq 10 Kb	1,591	24,748	9,249	10,310	10,310	10,310
Length \geq 1 Mb	0	0	27	37	37	37
Non-ATGC #	350,325	42,696,911	49,169,831	4,043,129	4,040,790	3,986,487
Non-ATGC %	0.038	4.67	5.36	0.393	0.423	0.436
N50 value	1,639	14,748	168,140	190,304	218,023	232,312

4

5 Figure legends

6 **Figure 1.** Photograph of the Indian blue peacock (*Pavo cristatus*).

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

8 **Figure 3.** Peacock proteins showing homology. Pie chart showing significant similarity
9 scores of peacock proteins against the National Center for Biotechnology Information Non-
10 Redundant (NCBI NR) database. The pie chart colors are grouped based on the E-value
11 scores from most significant E-value of 0.0 (red) going clockwise to least significant of about
12 1E-5 (blue).

1 **Figure 4.** Phylogenetic tree generated from homologous proteins from 49 different avian
2 species.

3 **Figure 5.** Venn diagram showing common and absent protein family domains (Pfam)
4 between peacock, chicken and turkey proteins.

5 **Figure 6.** Heatmap showing protein families (Pfam) distributed in peacock, chicken or turkey
6 species. The number represents the Pfam domain count predicted from the protein sequences.
7 Pfam domains of 50 and above identified in any one of the species are compared in the
8 heatmap.

9 **Figure 7.** Venn diagram showing peacock proteins with significant homology to the NCBI
10 NR database, the EuKaryotic Orthologous Groups (KOG) database, and protein family
11 (Pfam) and Gene Ontology (GO) ontologies.

12 **Figure 8.** Circular image of the assembled peacock genome, aligned against the *Gallus gallus*
13 genome. The right side of the image represents the reference chicken genome; left side
14 represents the peacock genome.

16 **References**

- 17 1. Brickle NW. Habitat use, predicted distribution and conservation of green peafowl
18 (*Pavo muticus*) in Dak Lak Province, Vietnam. *Biological Conservation*.
19 2002;105:189–97.
- 20 2. Jackson CE. Peacock. London: *Reaktion*, 2006.
- 21 3. Kadgaonkar, Shivendra B. The peacock in ancient Indian art and literature. *Bulletin of*
22 *the Deccan College Research Institute*. 1993;53:95–115.
- 23 4. Gadagkar R. Is the peacock merely beautiful or also honest? *Current*
24 *Science*. 2003;85:1012–20.

- 1 5. Kushwaha S, Kumar A. A Review on Indian Peafowl (*Pavo cristatus*) Linnaeus,
2 1758. *Journal of Wildlife Research*. 2016;4:42–59.
3
4
- 5 6. Hillier LW, Miller W, Birney E et al. Sequence and comparative analysis of the
6 chicken genome provide unique perspectives on vertebrate evolution. International
7 Chicken Genome Sequencing Consortium (ed.). *Nature*. 2004;432:695–716.
8
9
- 10 7. Zhang G, Jarvis ED, Gilbert MTP. A flock of genomes. *Science*. 2014;346:1308–9.
11
12
- 13 8. Bejerano G, Pheasant M, Makunin I et al. Ultraconserved elements in the human
14 genome. *Science*. 2004;304:1321–5.
15
16
- 17 9. Burt DW. Emergence of the Chicken as a Model Organism: Implications for
18 Agriculture and Biology. *Poultry Science*. 2007;86:1460–71.
19
20
- 21 10. Furlong RF. Insights into vertebrate evolution from the chicken genome
22 sequence. *Genome Biology*. 2005;6:207.
23
24
- 25 11. Edmunds S. Hiseq 4000 Sequencing protocol v1
26 (protocols.io.q58dy9w). *protocols.io* 2018,
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
- 50 12. Andrews S. FastQC A Quality Control tool for High Throughput Sequence
51 Data. 2015. <http://www.bioinformatics.babraham.ac.uk/projects/fastqc>
52
53
- 54 13. Kajitani R, Toshimoto K, Noguchi H et al. Efficient de novo assembly of highly
55 heterozygous genomes from whole-genome shotgun short reads. *Genome*
56
57
58
59
- 60 14. Metrichor. <https://nanoporetech.com/products/metrichor> Accessed 1st Feb 2019
61
62
63
64
65

- 1 15. Loman NJ, Quinlan AR. Poretools: a toolkit for analyzing nanopore sequence
2 data. *Bioinformatics*. 2014;30:3399–401.
3
4
- 5 16. Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler
6 transform. *Bioinformatics*. 2009;25:1754–60.
7
8
- 9 17. Birol I, Jackman SD, Nielsen CB et al. De novo transcriptome assembly with
10 ABySS. *Bioinformatics*. 2009;25:2872–7.
11
12
- 13 18. Boetzer M, Pirovano W. SSPACE-LongRead: scaffolding bacterial draft genomes
14 using long read sequence information. *BMC Bioinformatics*. 2014;15:211.
15
16
- 17 19. Boetzer M, Henkel CV, Jansen HJ et al. Scaffolding pre-assembled contigs using
18 SSPACE. *Bioinformatics* 2010;27:578–9.
19
20
- 21 20. Luo R, Liu B, Xie Y et al. SOAPdenovo2: an empirically improved memory-efficient
22 short-read de novo assembler. *GigaScience* 2012;1(1):18
23
24
- 25 21. Marçais G, Kingsford C. A fast, lock-free approach for efficient parallel counting of
26 occurrences of k-mers. *Bioinformatics*. 2011;27:764–70.
27
28
- 29 22. Smit A, Hubley R, Green P. RepeatMasker Open-4.0. 2013-
30 2015. www.repeatmasker.org. Accessed 1st Feb 2019
31
32
- 33 23. Thiel T. MISA - MIcroSAteLLite identification tool. 2012. [http://pgrc.ipk-
34 gatersleben.de/misa/](http://pgrc.ipk-gatersleben.de/misa/) Accessed 1st Feb 2019
35
36
- 37 24. Stanke M, Diekhans M, Baertsch R et al. Using native and syntenically mapped
38 cDNA alignments to improve de novo gene finding. *Bioinformatics*. 2008;24:637–44.
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

- 1 25. Altschul SF, Gish W, Miller W et al. Basic local alignment search tool. *Journal of*
2 *molecular biology*. 1990;215:403–10.
3
4
5
6 3 26. Moriya Y, Itoh M, Okuda S et al. KAAS: an automatic genome annotation and
7
8 4 pathway reconstruction server. *Nucleic Acids Research*. 2007;35:W182–5.
9
10
11 5 27. Nordberg H, Cantor M, Dusheyko S et al. The genome portal of the Department of
12
13 6 Energy Joint Genome Institute: 2014 updates. *Nucleic Acids Research*. 2013;42:D26–
14
15 7 31.
16
17
18
19
20 8 28. El-Gebali S, Mistry J, Bateman A et al. The Pfam protein families database in
21
22 9 2019. *Nucleic Acids Research*. 2018;47:D427–32.
23
24
25
26 10 29. Zhang G, Li B, Li C et al. Comparative genomic data of the Avian Phylogenomics
27
28 11 Project. *GigaScience*. 2014;3:26.
29
30
31
32 12 30. Jarvis ED, Mirarab S, Aberer AJ et al. Phylogenomic analyses data of the avian
33
34 13 phylogenomics project. *GigaScience*. 2015;4:4.
35
36
37
38 14 31. Fu L, Niu B, Zhu Z et al. CD-HIT: accelerated for clustering the next-generation
39
40 15 sequencing data. *Bioinformatics*. 2012;28:3150–2.
41
42
43
44 16 32. Larkin MA, Blackshields G, Brown NP et al. Clustal W and Clustal X version
45
46 17 2.0. *Bioinformatics*. 2007;23:2947–8.
47
48
49
50 18 33. Talavera G, Castresana J. Improvement of Phylogenies after Removing Divergent and
51
52 19 Ambiguously Aligned Blocks from Protein Sequence Alignments. *Systematic*
53
54 20 *Biology*. 2007;56:564–77.
55
56
57
58
59
60
61
62
63
64
65

- 1 34. Felsenstein J. PHYLIP (Phylogeny Inference Package) Version 3.57c. 1993
2 <http://evolution.genetics.washington.edu/phylip.html> Accessed 1st Feb 2019
3
4
5
6 35. Nguyen L-T, Schmidt HA, von Haeseler A et al. IQ-TREE: A Fast and Effective
7
8 Stochastic Algorithm for Estimating Maximum-Likelihood Phylogenies. *Molecular*
9
10 *Biology and Evolution*. 2014;32:268–74.
11
12
13
14 36. FigTree, a graphical viewer of phylogenetic trees.
15
16 <http://tree.bio.ed.ac.uk/software/figtree/> Accessed 1st Feb 2019
17
18
19
20 37. Tamazian G, Dobrynin P, Krasheninnikova K et al. Chromosomer: a reference-based
21
22 genome arrangement tool for producing draft chromosome sequences. *GigaScience*
23
24 2016;5(1):38.
25
26
27
28 38. Frith MC, Kawaguchi R. Split-alignment of genomes finds orthologies more
29
30 accurately. *Genome Biology*. 2015;16:106.
31
32
33
34 39. Krzywinski M, Schein J, Birol I et al. Circos: an information aesthetic for
35
36 comparative genomics. *Genome research*. 2009;19:1639–45.
37
38
39
40 40. Warren WC, Hillier LW, Tomlinson C et al. A New Chicken Genome Assembly
41
42 Provides Insight into Avian Genome Structure. *G3: Genes, Genomes,*
43
44 *Genetics*. 2016;7:109–17.
45
46
47
48 41. Peona V, Weissensteiner MH, Suh A. How complete are “complete” genome
49
50 assemblies?-An avian perspective. *Molecular Ecology Resources*. 2018;18:1188–95.
51
52
53
54 42. Muir P, Li S, Lou S et al. The real cost of sequencing: scaling computation to keep
55
56 pace with data generation. *Genome Biology*. 2016;17:53.
57
58
59
60
61
62
63
64
65

- 1 43. Goodwin S, McPherson JD, McCombie WR. Coming of age: ten years of next-
2 generation sequencing technologies. *Nature Reviews Genetics*. 2016;17:333–51.
3
4
5
6 3 44. Levy SE, Myers RM. Advancements in Next-Generation Sequencing. *Annual Review*
7
8 4 *of Genomics and Human Genetics*. 2016;17:95–115.
9
10
11 5 45. Lu H, Giordano F, Ning Z. Oxford Nanopore MinION Sequencing and Genome
12
13 Assembly. *Genomics, Proteomics & Bioinformatics*. 2016;14:265–79.
14
15
16
17 7 46. Rice ES, Green RE. New Approaches for Genome Assembly and Scaffolding. *Annual*
18
19 *Review of Animal Biosciences*. 2018;7.
20
21
22
23 9 47. Weissensteiner MH, Pang AWC, Bunikis I et al. Combination of short-read, long-
24
25 read, and optical mapping assemblies reveals large-scale tandem repeat arrays with
26 10
27 population genetic implications. *Genome Research*. 2017;27:697–708.
28
29
30
31
32 12 48. Sohn J, Nam J-W. The present and future of de novowhole-genome assembly.
33
34 13 *Briefings in Bioinformatics*. 2016:bbw096.
35
36
37
38 14 49. Zhang G, Li C, Li Q et al. Comparative genomics reveals insights into avian genome
39
40 15 evolution and adaptation. *Science*. 2014;346:1311–20.
41
42
43
44 16 50. Dalloul RA, Long JA, Zimin AV et al. Multi-Platform Next-Generation Sequencing
45
46 17 of the Domestic Turkey (*Meleagris gallopavo*): Genome Assembly and Analysis.
47
48 Roberts RJ (ed.). *PLoS Biology*. 2010;8:e1000475.
49
50
51
52 19 51. Kumar S, Stecher G, Suleski M et al. TimeTree: A Resource for Timelines,
53
54 20 Timetrees, and Divergence Times. *Molecular Biology and Evolution*. 2017;34:1812–
55
56 21 9.
57
58
59
60
61
62
63
64
65

- 1 52. Sotero-Caio CG, Platt RN, Suh A et al. Evolution and Diversity of Transposable
2 Elements in Vertebrate Genomes. *Genome Biology and Evolution*. 2017;9:161–77.
3
4
5
6 53. Kapusta A, Suh A. Evolution of bird genomes—a transposon’s-eye view. *Annals of the*
7
8 *New York Academy of Sciences*. 2016;1389:164–85.
9
10
11 54. Ramesh K, McGowan P. On the current status of Indian Peafowl *Pavo cristatus*
12 (Aves: Galliformes: Phasianidae): keeping the common species common. *Journal of*
13 *Threatened Taxa*. 2009;1:106–8.
14
15
16
17
18
19
20 55. Additional data for De novo genome assembly of the peacock.
21
22 <https://biit.cs.ut.ee/supplementary/peacock/> Accessed 1st Feb 2019
23
24
25
26 56. Reimand J, Arak T, Adler P et al. g:Profiler—a web server for functional
27 interpretation of gene lists (2016 update). *Nucleic Acids Research*. 2016;44:W83–9.
28
29
30
31
32 57. Kolde R, Vilo J. GOsummaries: an R Package for Visual Functional Annotation of
33 Experimental Data. *F1000Research*. 2015, DOI: 10.12688/f1000research.6925.1.
34
35
36
37
38 58. Dhar R, Seethy A, Pethusamy K, Singh S, Rohil V, Purkayastha K, Mukherjee I,
39 Goswami S, Singh R, Raj A, Srivastava T, Acharya S, Rajashekhar B and Karmakar S
40 (2019): Supporting data for "De novo genome assembly of the Indian Blue Peacock
41 (Pavo cristatus), from Oxford Nanopore and Illumina sequencing" *GigaScience*
42 *Database*. <http://dx.doi.org/10.5524/100559>
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65



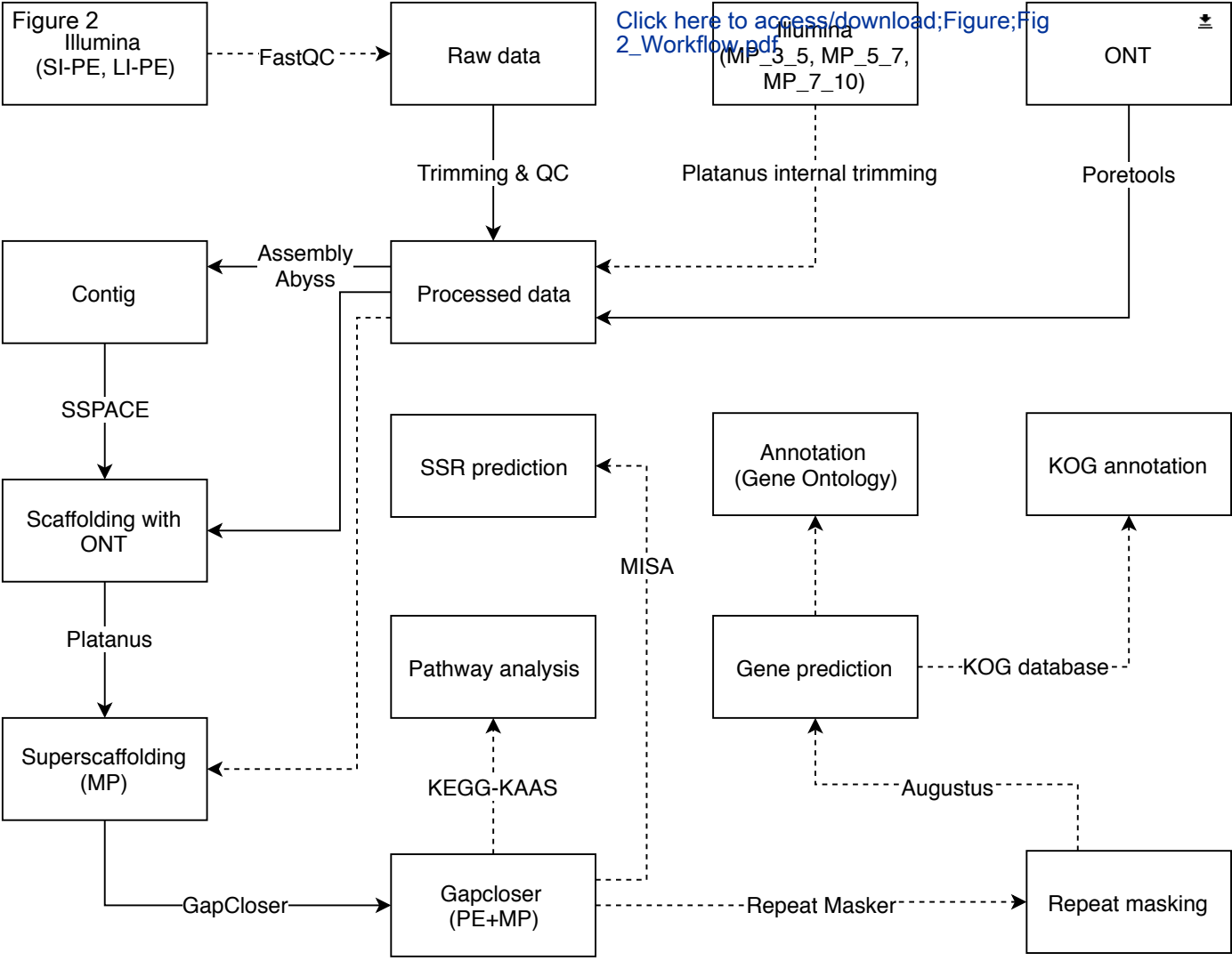
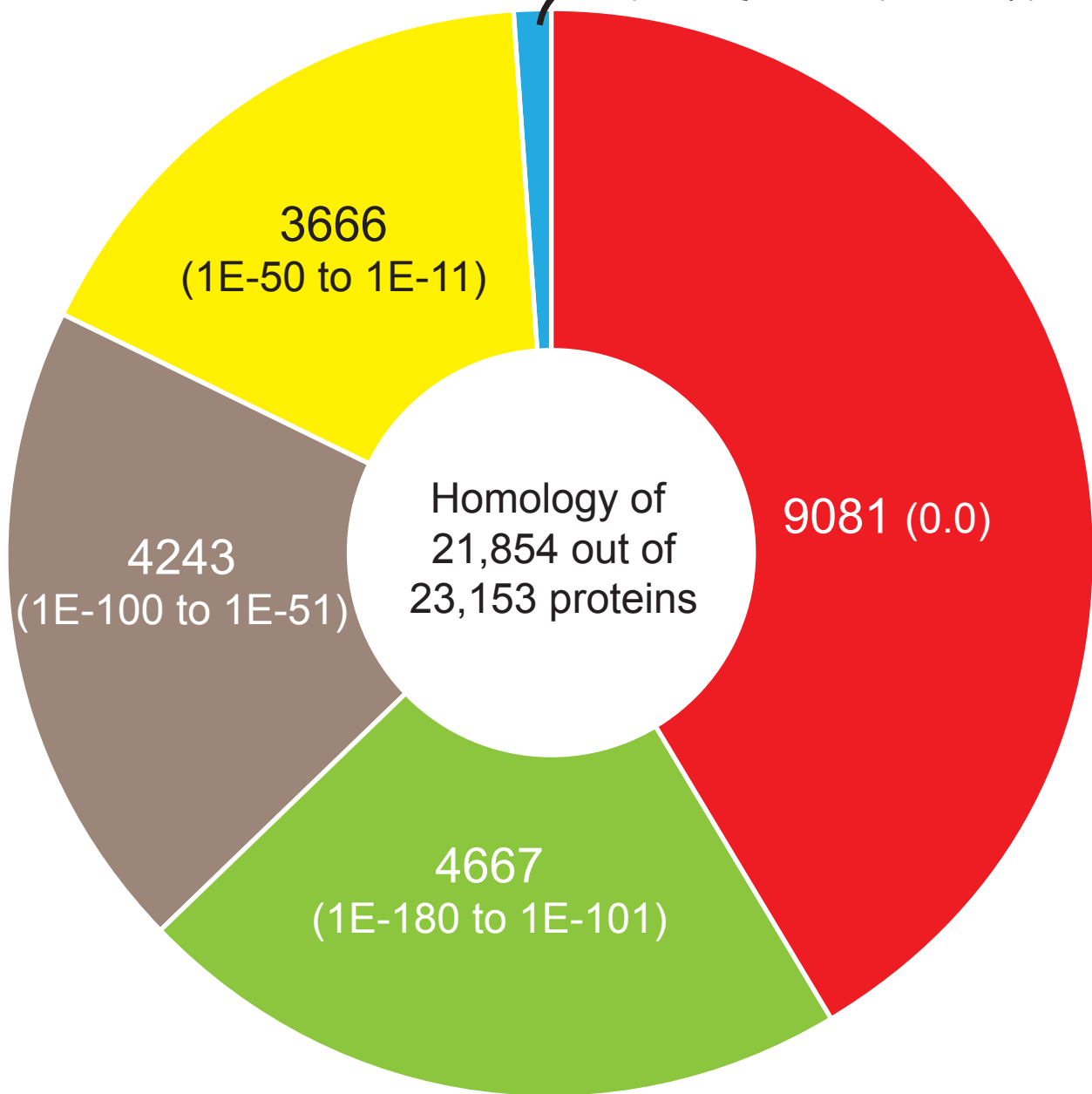


Figure 3

[Click here to access/download;Figure;Fig 3; Similarity scores against the Uniprot database.pdf](#)



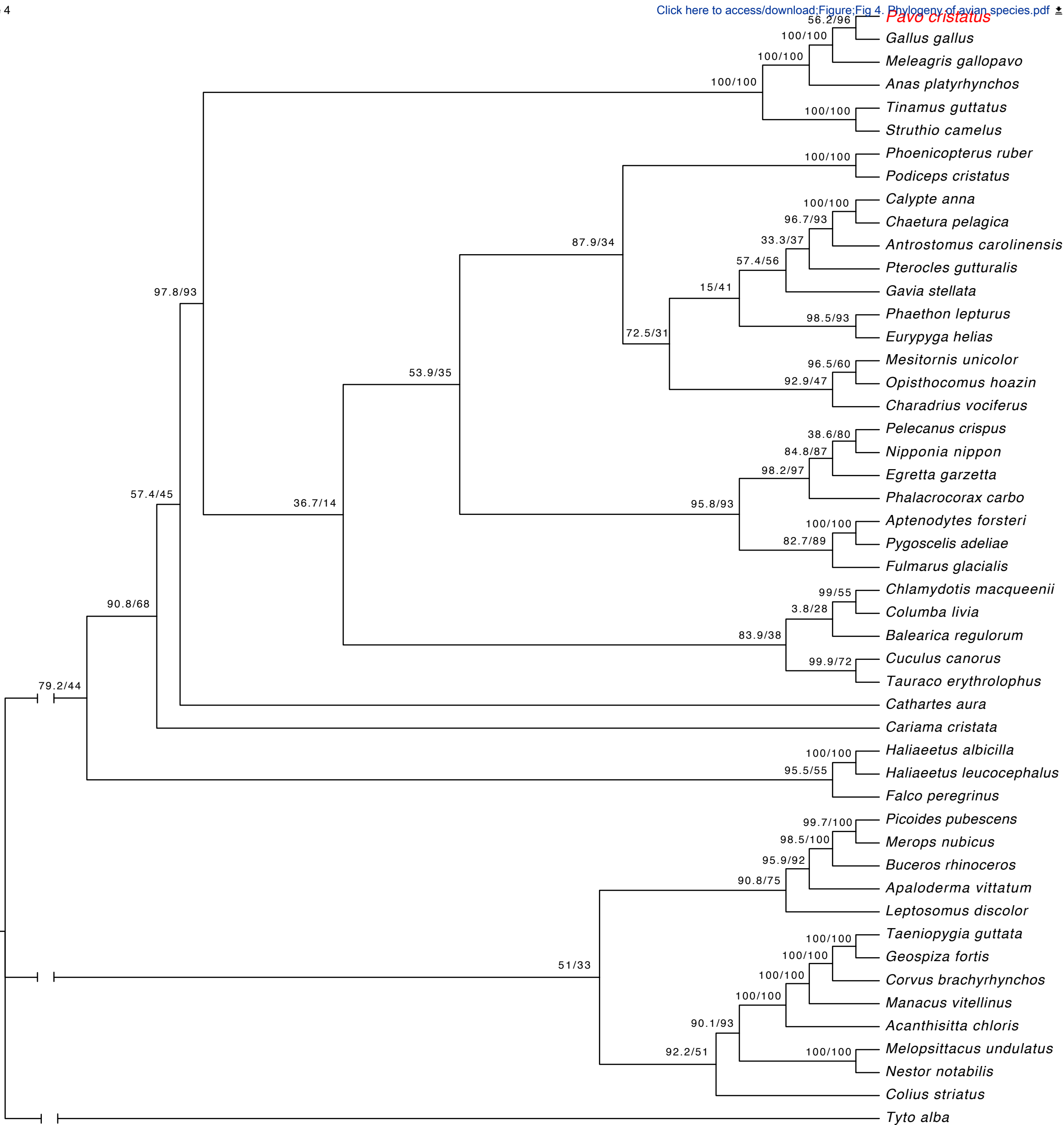
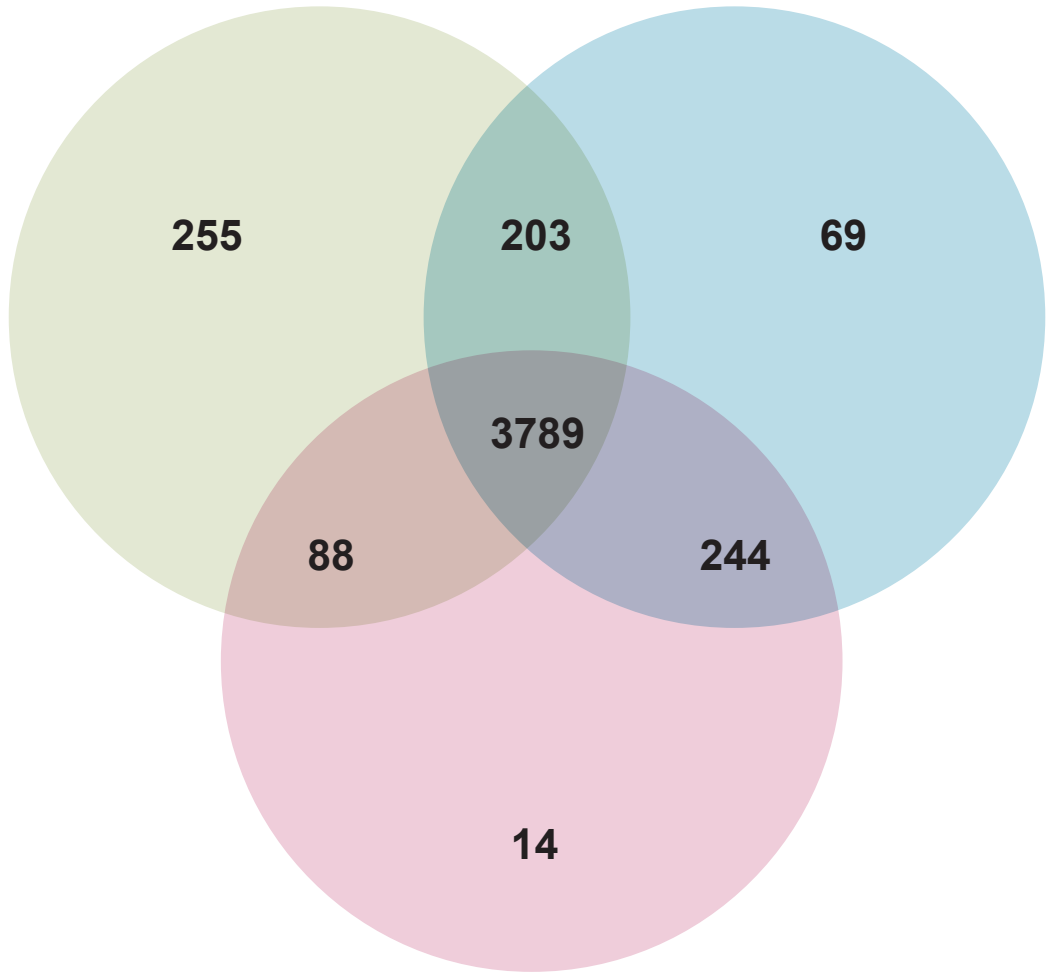
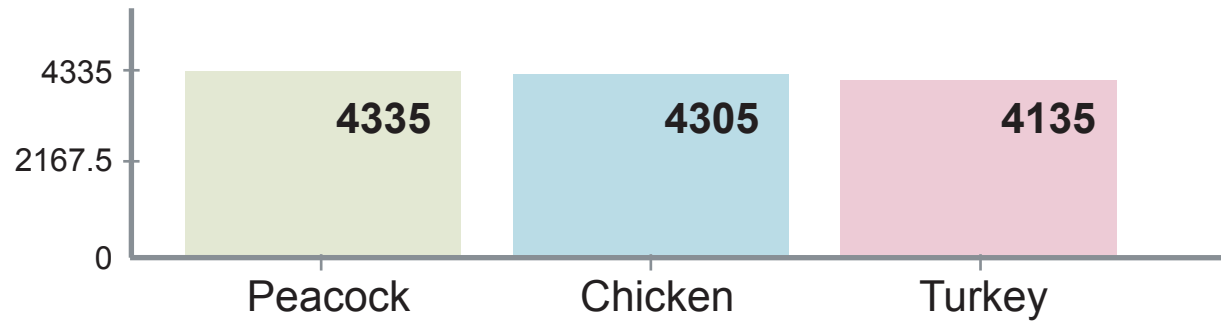


Figure 5

[Click here to access/download;Figure;Fig 5.](#) Unique Pfam domains common with chicken



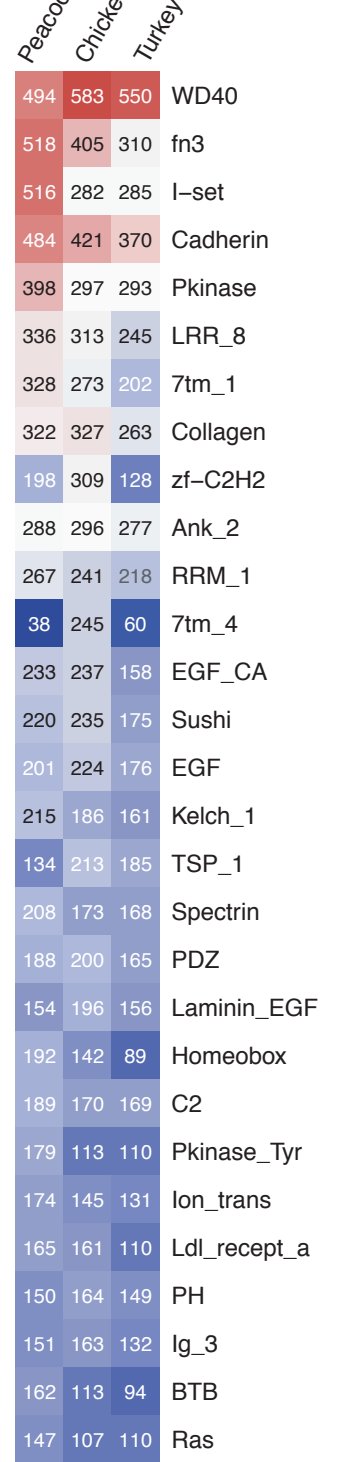
Turkey



Number of Pfam domains unique to 1 species or shared between 2 or all 3



Figure 6



Click here to access/download/figure;Figure6. Heatmap

Figure 6

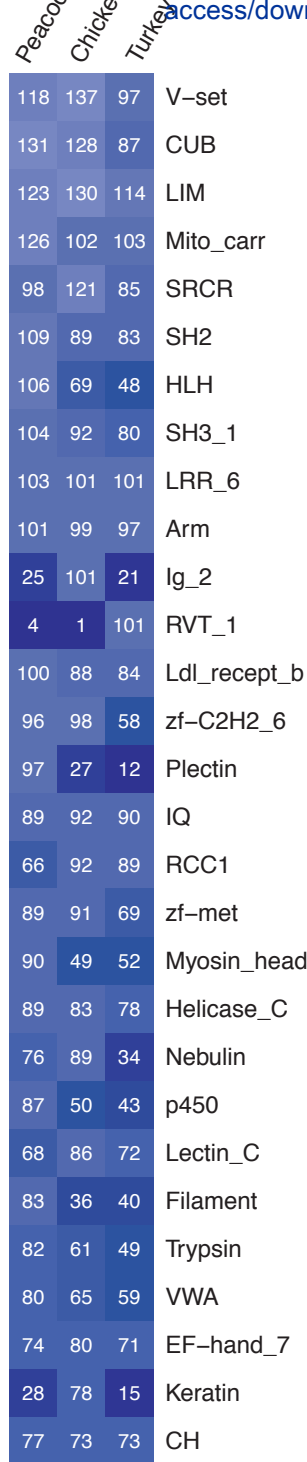
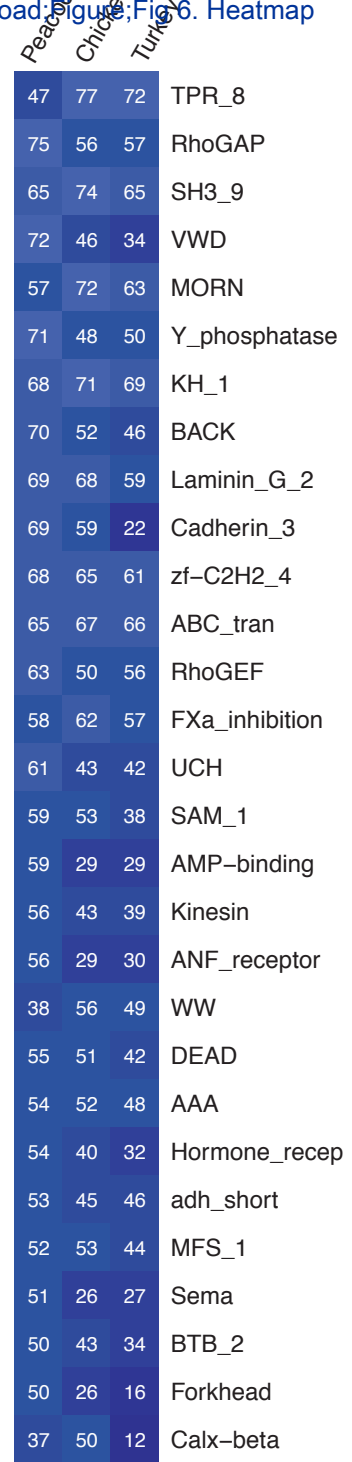
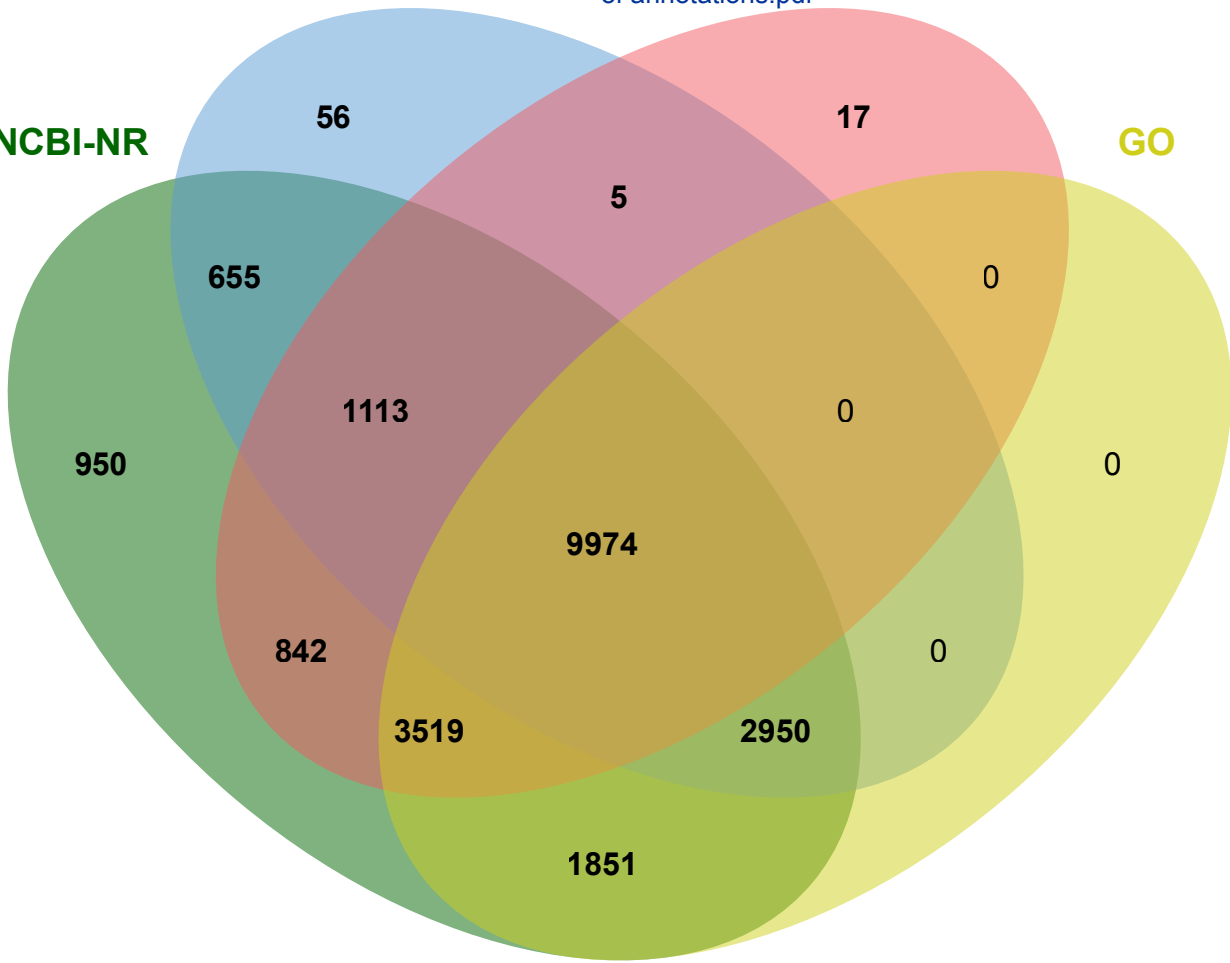


Figure 6

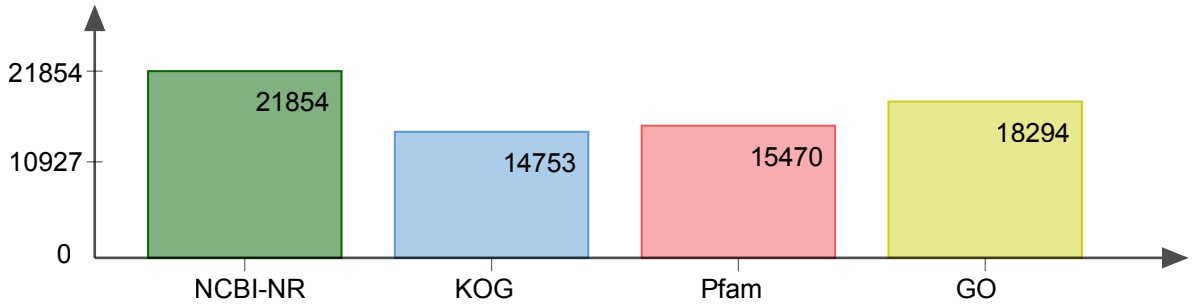


NCBI-NR

GO



Proteins annotated from different sources



Number of common proteins: specific to 1 or shared by 2, 3, or 4 annotations

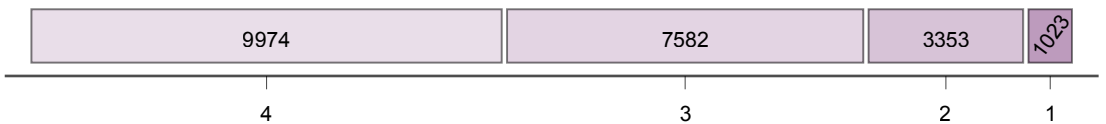
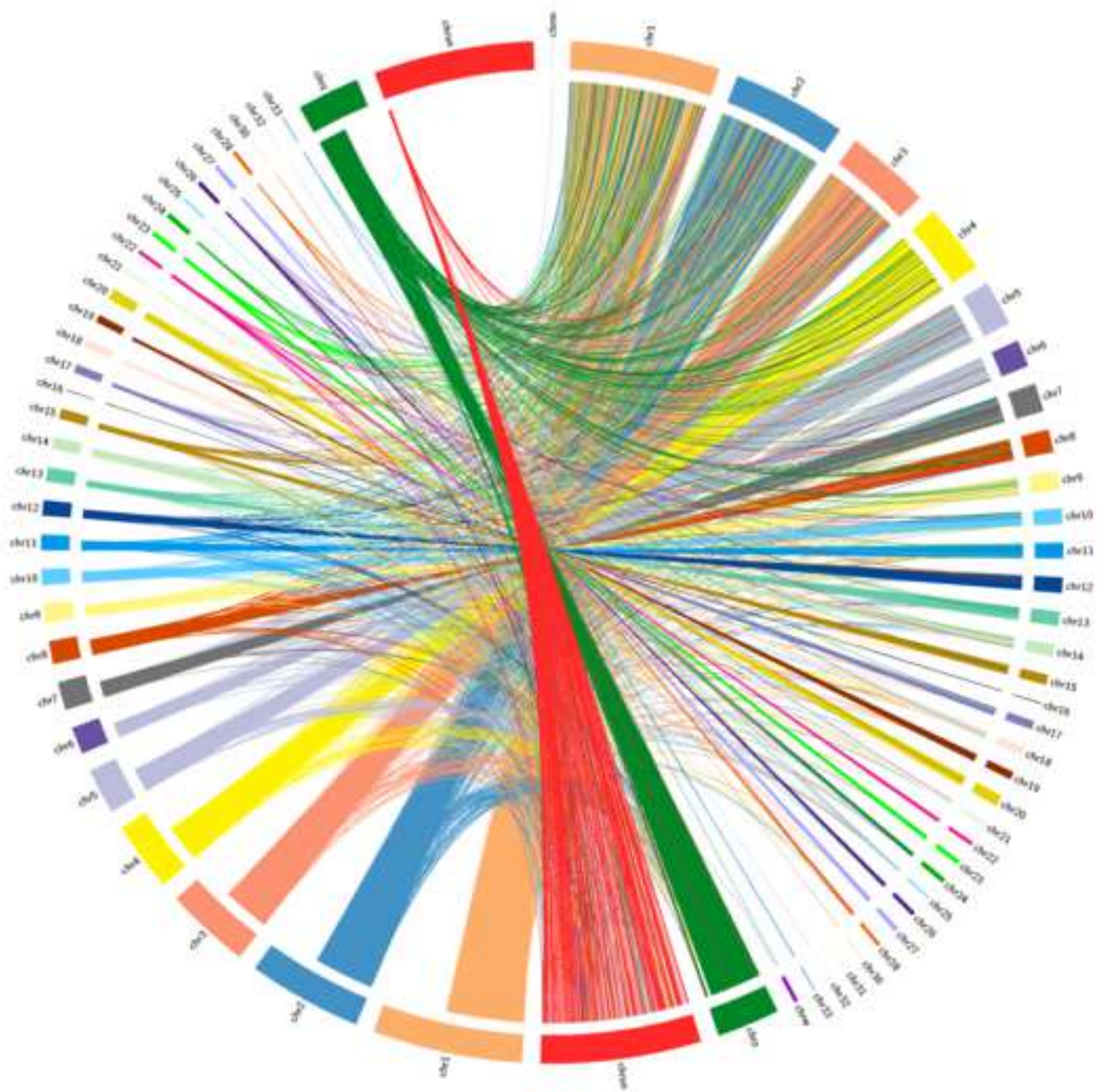


Figure 8

[Click here to access/download;Figure;Fig 8. Peacock scaffolds against Gallus circular synteny.png](#)





[Click here to access/download](#)

Supplementary Material

[Table_S1_ReadStats_Table_S2_TEs.xlsx](#)





Click here to access/download
Supplementary Material
Table_S3_Repeats.xlsx





[Click here to access/download](#)

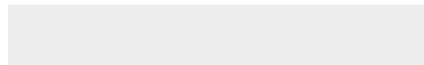
Supplementary Material

[Table_S4_Gene_annotations_of_peacock_proteins.xlsx](#)





Click here to access/download
Supplementary Material
Table_S5_KEGG_annotation.xlsx





Click here to access/download
Supplementary Material
Table_S6_KOG_annotation.xlsx





[Click here to access/download](#)

Supplementary Material

[Table_S7_BlastVsHumanProteins.xlsx](#)





Click here to access/download
Supplementary Material
Table_S8_Orthologous_proteins

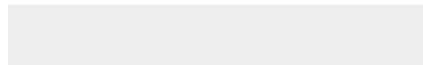




Click here to access/download

Supplementary Material

Table_S9_Pfam_Analysis.xlsx





[Click here to access/download](#)

Supplementary Material

[Table_S10_Bird_Species_with_counts.xlsx](#)

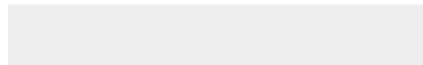




[Click here to access/download](#)

Supplementary Material

Fig S1. Proteins showing similarity to Pfam domains.pdf

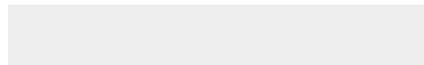




[Click here to access/download](#)

Supplementary Material

Fig S2. Gene Ontology of top 10 WGS.png





Click here to access/download
Supplementary Material
Fig S3.Peacock vs Human_GO.pdf

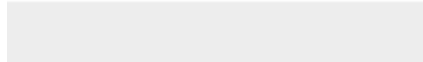




[Click here to access/download](#)

Supplementary Material

[Table_S11_Peacock_assembly_comparison.xlsx](#)

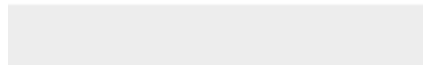




[Click here to access/download](#)

Supplementary Material

[Table_S10_Bird_Species_with_counts.xlsx](#)





[Click here to access/download](#)

Supplementary Material

[Supplementary_Description of tables and figures.docx](#)

