# A comprehensive prediction of transcript isoforms in 19 chicken tissues by Oxford Nanopore long-read sequencing

**Dailu Guan[1], Michelle M. Halstead[1], Alma D. Islas-Trejo[1], Daniel E. Goszczynski[1], Hans H. Cheng[2], Pablo Ross[1,*], Huaijun Zhou[1,*]**

[1]Department of Animal Science, University of California Davis, Davis, CA, 95616 USA

[2] USDA, ARS, USNPRC, Avian Disease and Oncology Laboratory, East Lansing, MI 48823, USA

**\* Correspondence:**
Corresponding Author
Pablo Ross pross@ucdavis.edu

Huaijun Zhou hzhou@ucdavis.edu

## Abstract

To comprehensively identify and annotate transcript isoforms in the chicken genome, we generated Nanopore long-read sequencing data from a diverse set of 19 chicken tissues comprising 68 samples collected from experimental line $6 \times$ line 7 $F_1$ adult males and females. More than 23.8 million reads with mean read length of 790 bases and average quality of 18.2 were generated. The annotation and subsequent filtering resulted in identification of 55,382 transcripts with mean length of 1,700 bases at 40,547 loci, representing ~1.4 transcripts per locus. Among them, we predicted 30,967 potential coding transcripts at 19,461 loci and 16,495 potential lncRNA transcripts at 15,512 loci. Compared to reference annotations, we found 52% of annotated transcripts could partially to fully match while 47% were novel and potentially transcribed from lncRNA loci. Based on our annotation, we quantified transcript expression across tissues and found brain tissues (i.e. cerebellum, cortex) expressed highest number of transcripts and loci. The further tissue specificity revealed that ~22% of the transcripts displaying tissue specificity. Of them, the reproductive tissues (i.e. testis, ovary) contained the most tissue-specific transcripts. Despite sequencing 68 transcriptomes derived from 19 tissues, still ~20% of Ensembl reference loci were not detected. This suggests that including additional samples from different cell types, developmental and physiological conditions, is needed to fully annotate the chicken genome. The application of Nanopore sequencing transcriptomes in this study demonstrated the usefulness of long-read data in discovering additional novel loci (e.g., lncRNA loci) and resolving complex transcripts (e.g., the longest transcript for the *TTN* locus).

## 1    Introduction

36    Chicken (*Gallus gallus domesticus*) is one of the most widespread and common domesticated farm
37    animals for egg and meat production, with a total population of 37.2 billion in stocks for the year
38    2020 (http://www.fao.org/). As the most popular studied bird species, moreover, its importance to the
39    study of evolution, development, immunology, etc. is self-evident. In 2004, the first draft whole
40    chicken genome was assembled with an estimated set of 20-23,000 protein-coding genes (PCGs)
41    (Hillier et al., 2004). This effort offered a genome-wide view for understanding the configuration of
42    the chicken genome, and the evolution of coding and noncoding vertebrate genomes. Since then,
43    continuous efforts have been made to improve the completeness of chicken genome. For instance,
44    Warren et al. (2017) added an additional 183 Mb sequences and assembled chromosomes 30-33 for
45    the chicken reference genome. To fill the gaps of the chicken reference genome, recently two
46    pangenomes were built that reported additional sequences absent from the GRCg6a reference
47    genome (Wang et al., 2021a; Li et al., 2022).

48    The functional annotation of the chicken genome is also being produced in parallel. The two most
49    commonly used databases, i.e. Ensembl (https://uswest.ensembl.org) and National Center for
50    Biotechnology Information (NCBI, https://www.ncbi.nlm.nih.gov/) regularly update the chicken
51    genome annotation. For instance, the Ensembl release (V102) included 16,779 PCGs and 39,288
52    transcripts, representing 2.34 transcripts per gene. Compared to the human ~10 transcripts per gene,
53    this estimate is quite low. The high estimate in human is partly attributed to the global efforts, such as
54    GENCODE, which is part of the ENCODE (ENCyclopedia Of DNA Elements) consortium which
55    aims to identify and classify all gene features in the human and mouse genomes. In farm animals,
56    likewise, the consortium of the Functional Annotation of ANimal Genome (FAANG) was formed in
57    order to improve the annotation of livestock genomes (Giuffra et al., 2019; Clark et al., 2020). In
58    prior work, Kern et al. (2021) annotated noncoding genomes of three important livestock including
59    chicken, and predicted 29,526 regulatory elements-gene interactions in chickens. In addition, Kern et
60    al. (2018) also identified a total of 9,393 lncRNAs (including 5,288 novel lncRNAs) by utilizing
61    short-read transcriptomes from eight chicken tissues.

62    The transcribed genomic region, though it only accounts for ~3% of the genome, is very complex due
63    to the alternative usage of transcription start, splicing, and polyadenylation sites. Alternative splicing
64    has been shown to play important roles in evolution, phenotypic diversity, and organ development
65    (Keren et al., 2010; Baralle and Giudice, 2017; Wright et al., 2022). For example, Yu et al. ( 2019)
66    identified five alternative splicing variants of the *TYR* gene that were associated with skin
67    melanogenesis in chickens. To annotate these features, transcriptome profiling provides important
68    and useful resources (Yandell and Ence, 2012). For example, Jehl et al. ( 2020) annotated additional
69    1,199 PCGs and 13,009 long non-coding RNA genes (Compared to Ensembl V94) using 364 short-
70    read transcriptomes derived from 25 chicken tissues. In human, a comprehensive annotation using
71    transcriptomes of 41 tissues generated by Genotype-Tissue Expression (GTEx) Consortium improved
72    transcript prediction for 13,429 genes, including 1,831 (63%) Online Mendelian Inheritance in Man
73    (OMIM) genes and 317 neurodegeneration-associated genes (Zhang et al., 2020). This analysis
74    demonstrated that a detailed annotation is better for understanding the phenome-to-genome
75    connections. Although the short-read sequencing is widely used for annotating human and animal
76    genomes, it is difficult to accurately resolve the complex structure of transcript isoforms. Chen et al.
77    (2021), for instance, demonstrated that Nanopore long-read transcriptome sequencing classified
78    individual isoforms better than Illumina short-reads despite they generated comparable gene
79    expression estimates.

80    The contiguity of the long-read sequencing technology can sequence full-length transcript, thus is
81    better suitable for dissecting the complexity of transcript structure compared to short-read sequencing

82   (Muret et al., 2017). The Iso-Seq by the Pacific Biosciences is long-read sequencing technology that
83   is widely used in profiling full-length transcriptome in human (Kuo et al., 2020), pig (Beiki et al.,
84   2019), rabbit (Chen et al., 2017). In chickens, Thomas et al. (2014) used Iso-Seq long-read
85   sequencing and identified 9,221 new transcript isoforms in embryonic chicken heart tissue. Later on,
86   Kuo et al. (2017) annotated 64,277 additional distinct transcripts (55,315 in brain and 9,206 in
87   embryo) using Iso-Seq plus 5′ cap selection in chicken brain and embryo tissues. However, the few
88   tissues studied (only brain and embryo included in Kuo et al. (2017)) make it difficult to capture the
89   diversity of chicken transcript variations.

90   The Oxford Nanopore Technologies has provided an alternative long-read sequencing approach
91   (Amarasinghe et al., 2020), which has been applied in cattle (Halstead et al., 2021), duck (Lin et al.,
92   2021) and many other species, but not yet in chickens. The Nanopore long-read sequencing allows
93   for accurate identification and quantification of transcript isoforms and for resolving complex
94   isoforms (Byrne et al., 2017; Soneson et al., 2019; Chen et al., 2021). In this study, we aimed to
95   identify and characterize transcripts in a diverse set of chicken tissues, including cerebellum,
96   hypothalamus, cortex, duodenum, jejunum, ileum, cecum, colon, testis, ovary, adipose, gizzard,
97   heart, kidney, liver, lung, muscle, spleen and thymus, using Oxford Nanopore long-read sequencing
98   technology. The data generated from this study will be a valuable source to improve our
99   understanding of the complexity of the chicken transcriptome, and also aid in dissecting the
100  connection of gene expression and phenotypic traits.

101

## 2   Methods and Materials

### 2.1   Sample collection

104  All animals and samples used in this study were obtained in concordance with the Protocol for
105  Animal Care and Use no. 18464 (approved by Institutional Animal Care and Use Committee at the
106  University of California at Davis). All tissues were from one of two FAANG pilot projects
107  (FarmENCODE) (Tixier-Boichard et al., 2021). In brief, ADOL experimental White Leghorn lines $6_3$
108  and $7_2$ were intermated to produce $F_1$ progeny, and 4 male and 2 female individuals were euthanized
109  at 20 weeks of age. Tissues were collected within 1-2 hours and stored at -80 °C until further use.

110

### 2.2   RNA extraction and library preparation

112  RNA extraction and library preparation were done by following the protocols reported in (Halstead et
113  al., 2021). Briefly, frozen tissues were mashed using a pestle in a mortar filled with liquid nitrogen.
114  Then, Trizol reagent (Invitrogen, Carlsbad, CA, United States) was added to extract total RNA using
115  the Direct-zol RNA Mini Prep Plus kit (Zymo Research, Irvine, CA, United States). The integrity and
116  quality of extracted RNA were checked using an Experion electrophoresis system (Bio-Rad,
117  Hercules, CA, United States) and those passing quality control were used for library preparation.
118  First, 50 ng of total RNA in a volume of 9 μl was mixed with 1 μl 10 μM VNP primer, 1 μl 10 mM
119  dNTPs for incubation 5 min at 65 °C. The products were used for strand-switching and reverse
120  transcription reactions. Then, barcodes were ligated to the cDNA products generated from the last
121  step using the Oxford Nanopore PCR barcoding expansion 1-96 kit (cat. no. EXP-PBC096), which
122  were further ligated with adapters from the SQK-DCS109 kit following the manufacturer's
123  guidelines. Products were loaded onto a PromethION flow cell (vR9.4.1) for sequencing.

124

## 2.3 Base calling, quality control and preprocessing

126  After base calling and de-multiplexing with the ont-guppy-for-minknow (v3.0.5) tool
127  (https://nanoporetech.com/), we run quality control using NanoPlot (v1.0.0) software in order to
128  summarize read length, average quality, among others. Then, the Pychopper software
129  (https://github.com/nanoporetech/pychopper) was employed to identify and orient full-length reads,
130  which were mapped against the reference genomes (GRCg6a, Ensembl V102) with options of "-ax
131  splice -uf -k14 -G 1000000" using the minimap2 software (Li, 2018). We discarded reads with a
132  minimum quality score of 10 using SAMtools (v1.9) (Li et al., 2009) and counted the gene
133  expression using the HTSeq 0.13.5 software (Anders et al., 2015). Read counts were then normalized
134  into variance stabilizing transformation (VST), which was used for sample clustering analysis with
135  the function of "plotPCA" implemented in the DEseq2 software (Love et al., 2014).

136

## 2.4 Reference-guided prediction of transcript isoforms

138  To predict transcripts, we used a computational pipeline supported by the Oxford Nanopore
139  Technology community (https://github.com/nanoporetech/pipeline-nanopore-ref-isoforms). Briefly,
140  the oriented full-length reads with fastq format were pooled together and then mapped against the
141  Ensembl annotation (GRCg6a, V102) using minimap2 (Li, 2018) in order to carry out a reference-
142  guided transcriptome assembly. The mapped reads were then used to annotate transcripts using the
143  StringTie2 software (Kovaka et al., 2019) in the long-read mode (with the option of "-L").
144  Transcripts on unplaced scaffolds, as well as those with exon coverage < 100% and read depth < 2
145  were excluded. Then we only keep single-exon transcripts with expression TPM > 1 in > 2 samples
146  of a tissue, and multi-exon transcripts with expression TPM >0.1 in >2 samples of a tissue. After that,
147  we further excluded transcripts categorized as potential artifacts (see **Comparing predicted**
148  **transcripts to previous annotations** section).

149

## 2.5 Prediction of coding and non-coding transcripts and loci

151  To predict the coding potential of predicted transcripts, we employed the TransDecoder
152  (https://github.com/TransDecoder/TransDecoder) and CPP2 software (Kang et al., 2017). The final
153  predictions are composed of either TransDecoder or CPP2 ones. After prediction of coding potential,
154  we obtained the list of non-coding transcripts, which were further used for predicting whether they
155  are lncRNA loci using the FEElnc software (Wucher et al., 2017).

156

## 2.6 Comparing predicted transcripts to previous annotations

158  The predicted transcripts were compared to the Ensembl (V102) and NCBI reference (V105)
159  annotations using GffCompare tool (version 0.11) (Pertea and Pertea, 2020) and classified into 14
160  classes. According to Halstead et al. (2021), the predicted transcripts classified into four categories:
161  exact match (class code "="), which means the intron chains of our annotated transcripts can exactly
162  match to reference annotations; novel isoform (class codes 'c,' 'k,' 'j,' 'm,' 'n,' or 'o'), which means
163  predicted transcript can not match a reference transcript but can match a reference gene; novel loci

164 (class codes 'i,' 'u,' 'y,' or 'x'), which means predicted transcript can not match either a reference
165 transcript or a reference locus; and potential artifacts (class codes 'e,' 's,' or 'p'), which are possibly
166 due to mapping error, e.g. pre-mRNA fragments, polymerase run-on, etc. To compare our prediction
167 with novel transcripts reported by Thomas et al. (2014), we first converted positions of their
168 transcripts from galGal4 to GRCg6a using the liftover software (Kuhn et al., 2013). Then the
169 GffCompare tool was used for comparing our annotation to their transcripts.

170

### 171 2.7 Quantification of predicted transcripts

172 We extracted sequences of predicted transcripts using GffRead v0.12 (Pertea and Pertea, 2020),
173 which constituted a reference transcriptome in the FASTA format. Then, we mapped the full-length
174 reads generated by Pychopper (https://github.com/nanoporetech/pychopper) to the predicted
175 transcriptome using minimap2 (v2.1) (Li, 2018). The transcript expression was quantified using
176 Nanocount (v0.2.4) (Leger, 2020). Based on the metric of the transcripts per million (TPM), we
177 categorized transcripts as highly (average TPM $\geq$ 10), moderately ($1 \leq$ average TPM $< 10$), and lowly
178 expressed (average TPM $< 1$) (Halstead et al., 2021).

179

### 180 2.8 Tissue-specificity analysis

181 The tissue specificity of transcripts expression across tissues were evaluated by using a tissue
182 specificity index (*TSI*) (Julien et al., 2012; Halstead et al., 2021):

$$TSI = \frac{\max_{1 \leq i \leq n}(x_i)}{\sum_{i=1}^{n} x_i}$$

183 Where $x_i$ is an average of transcript expression (TPM) in a given tissue, n is the number of tissues.
184 Transcripts were then categorized as tissue-specific (TSI $\geq$ 0.8), broadly expressed (TSI $< 0.5$), or
185 biased towards a group of tissues ($0.5 \leq$ TSI $< 0.8$). To reveal functional biology of tissue specific
186 transcripts, we extracted tissue-specific transcript sequences and blast them against the SwissProt
187 (protein sequence database, V5) using the Diamond blastx tool (v2.0.11.149) (Buchfink et al., 2015).
188 We then carried out functional enrichment (only considering Gene Ontology Biological Process
189 terms) using the matched UniProt identifiers using the PANTHER tool (Mi et al., 2013). The false
190 discovery rate (FDR) approach (Benjamini and Hochberg, 1995) was used for multiple testing
191 corrections and FDR value less than 0.05 was set as the significance threshold.

192

### 193 2.9 Differential alternative splicing analysis

194 To detect differential alternative splicing (DAS) events, we employed the LIQA software (Hu et al.,
195 2021). Based on our annotation, we quantified isoform expression using the "quantify" function.
196 Then the DAS events between tissues were detected using the "diff" within the LIQA tool (Hu et al.,
197 2021). After multiple testing correction using FDR approach (Benjamini and Hochberg, 1995), the
198 threshold of significance was set as FDR $< 0.05$.

199

## 3    Results

To comprehensively annotate transcripts of the chicken genome, we sequenced 68 samples comprising of 19 tissues collected from six individuals (two females: CC and CD; and four males: CA, CB, M1, M2) (**Supplementary Table 1**). The tissues collected were cerebellum, hypothalamus, cortex, duodenum, jejunum, ileum, cecum, colon, testis, ovary, adipose, gizzard, heart, kidney, liver, lung, muscle, spleen and thymus. Sequencing generated a total of 23.8 million reads, with an average of 344,650 reads per tissue and an average length of 790 bp (**Figure 1a**, **Supplementary Table 2**).

Principal component analysis (PCA) and hierarchical clustering of mapped sequencing reads to the Ensembl annotation (GRCg6a, version 102) revealed that samples generally clustered according to the origin of tissues or organs as expected (**Figure 1b** and **Supplementary Figure 1**). Moreover, we found samples from the same biological system tend to cluster together, such as brain cortex, cerebellum and hypothalamus from the central neural system; cecum, colon, duodenum, ileum, and jejunum from the intestinal system (**Figure 1b**). However, an outlier (i.e., Cecum_CA) in the PCA plot and hierarchical clustering (indicated by a red arrow in **Figure 1b** and **Supplementary Figure 2**) was observed, even separated from the same cecum tissue. The summary statistic indicates that the unexpected clustering is possibly due to the insufficient sequencing depth (number of reads = 1,279), which was lower than the rest (average number of reads for remaining samples = 355,008, **Supplementary Figure 2**). However, 895 out of 1,279 reads from Cecum_CA reads aligned to the GRCg6a genome, corresponding to a mapping rate of 70%. In the light of these analyses, we included the Cecum_CA sample in transcript prediction, but not in the transcript expression analysis, e.g., tissue specificity of transcript expression, differential alternative splicing (DAS) analysis.

To assemble potential transcripts, we identified, oriented and trimmed full-length reads using the Pychopper v2 software. Then, the StringTie tool with the long read mode was used for predicting transcripts (https://github.com/nanoporetech/pipeline-nanopore-ref-isoforms). As a result, 79,757 transcripts in 54,551 loci in total were identified. After filtering out transcripts on unplaced scaffolds, as well as those with exon coverage < 100% and read depth < 2, we obtained 74,665 transcripts in 50,569 loci, of which there were 45,132 multi-exon and 29,533 single-exon transcripts. Moreover, we required multi-exon transcripts with expression TPM (Transcripts Per Million) > 0.1 and single-exon transcripts with expression TPM > 1 in at least 2 sample of a tissue. By doing so, there were 61,556 transcripts in 45,284 loci remained. To further exclude potential artifacts, we compared assembled transcripts with NCBI (V105) and Ensembl (V102) reference annotations. The result is shown in **Figure 2a** and **Table 1** (see **Methods**). Overall, we found ~14% of predicted transcripts can exactly match the reference annotations (**Figure 2a**). With the Ensembl annotation, 77% of them were considered as novel transcripts, either novel isoforms (35%) or novel loci (42%). In addition, ~8% were potential artifacts, which possibly caused by pre-mRNA fragment, polymerase run-on or mapping error (**Figure 2a**). After excluding these potential artifacts, we finally kept 55,382 transcripts in 40,547 loci, representing 1.4 transcripts per locus (**Supplementary Data 1**).

The length of predicted transcripts ranged from 49 to 34,500 bp, with a mean length of 1,767 bp (**Figure 2b**). The longest transcript, for instance, is located on chromosome 7 (15,343,033-15,384,347), which highly matched to the *TTN* gene encoding the giant protein titin (NCBI reference sequence XM_046921719.1, E-value = 0.0, percent of identity = 99.99%) (**Figure 2d**). This protein plays important roles in the movement of skeletal muscle, but its gene locus has not been annotated in both NCBI (V105) and Ensembl (V102) GRCg6a reference annotations (**Figure 2d**). Moreover, we found the annotated 55,382 transcripts are supported by 171,651 unique exons, with an average estimate of 4.34 exons per transcript (**Figure 2c**).

6

245 To predict the coding potential of predicted transcripts, we employed the CPC2 and TransDecoder
246 software. The former predicted 21,984 transcripts at 12,999 loci with coding potential, and the latter
247 one predicted open reading frames for 30,727 transcripts corresponding to 19,306 loci. In total, we
248 predicted 30,967 uniquely potential coding transcripts at 19,461 loci, representing 1.6 transcripts per
249 locus (**Supplementary Table 3**). Furthermore, we surveyed whether the remaining 24,415 transcripts
250 were long non-coding RNAs (lncRNAs). To do so, we employed the FEELnc software and found
251 16,495 potential lncRNA transcripts at 15,512 loci (**Supplementary Table 3**).

252 We compared our prediction to two reference annotations and found the number of transcripts per
253 locus of our annotation (~1.4) is lower than reference annotations (Ensembl v102: ~1.8 transcripts
254 per locus; NCBI v105: ~3.3 transcripts per locus), but we predicted ~20K more loci, of which a
255 substantial proportion is lncRNA loci (**Figures 3a** and **3c**). At the transcript level, we classified
256 transcripts into three categories (see **Methods**): 1) exact match: predicted transcripts completely
257 match to reference annotations; 2) novel isoform: predicted transcripts do not match reference
258 transcripts but match reference loci; 3) novel loci: predicted transcripts do not match any reference
259 loci and transcripts (**Figure 3b**). Concordantly, we found our prediction identified high proportion of
260 "novel loci" transcripts (47%), followed by "novel isoforms" (37%) when comparing to Ensembl
261 annotation (V102) (**Figure 3b**). A similar pattern was observed when comparing to NCBI annotation
262 (**Supplementary Figure 3**). By further comparing lncRNA loci predicted in this study with those
263 predicted by Jehl et al. (2020), we found ~ 83% of our predicted lncRNA transcripts can match their
264 annotations (**Supplementary Figure 4**). Thomas et al. (2014) also reported 9K novel transcripts
265 from long-read sequenced embryonic chicken heart transcriptomes. By comparing these available
266 novel transcripts to our annotation, we found 89% of them can completely or partially match to our
267 annotation, while there were still 1,000 transcripts categorized as "novel loci" (**Supplementary
268 Figure 5**). Moreover, we found the transcripts grouped into the "novel isoform" and "novel loci"
269 categories tend to be lowly expressed, while the expressions of transcripts in "exact match" group are
270 higher (**Figure 3d**).

271 Considering the largest set of tissues used, we then sought to identify tissue-specifically expressed
272 transcripts. By quantifying transcript expressions, we found the number of expressed transcripts and
273 loci ranged from 14,841 (liver) to 28,648 (cerebellum), and from 10,285 (liver) to 21,662
274 (cerebellum), respectively (**Supplementary Figure 6**). The tissue specificity index (TSI) indicated
275 that the set of "exact match" transcripts tend to be lowly tissue-specific, while "novel isoform" and
276 "novel loci" transcripts are highly tissue-specific (**Figure 4a**). We observed that the set of transcripts
277 with low expression tended to have high tissue-specificity, while in contrast, highly expressed
278 transcripts are commonly found across tissues (**Figure 4b**). Moreover, we identified tissue-specific
279 transcripts and found the reproductive tissues (i.e., testis and ovary) have high proportion of tissue-
280 specific transcripts, followed by brain-related tissues (i.e., cerebellum and cortex) (**Figure 4c**). For
281 instance, we identified a novel transcript located on chromosome 4 (52,482,563-52,492,561), which
282 is specifically expressed in testis samples (**Figures 4d** and **4e**). This transcript was predicted as a
283 sense intergenic lncRNA by using the FEELnc software (Wucher et al., 2017) (**Supplementary
284 Tables 3** and **4**). To reveal the function of tissue-specific transcripts, we aligned sequences of tissue-
285 specific transcripts to SwissProt (V5) database with the blastx function implemented in the Diamond
286 tool (v2.0.11.149) (Buchfink et al., 2015). Then, the matched UniProt identifiers were used for
287 carrying out functional enrichment analysis with the PANTHER tool (Mi et al., 2013). This analysis
288 revealed that tissue-specific transcripts recapitulated the tissue biology (**Figure 5a, Supplementary
289 Table 5**), such as muscle contraction, muscle cell differentiation enriched in muscle and heart tissues,
290 trans-synaptic signaling and nervous system development in cerebellum and brain cortex, and B cell

7

291    receptor signaling pathway in spleen (**Figure 5a**, **Supplementary Table 5**), a finding concordant
292    with previous results (Yang et al., 2018; Fang et al., 2020).

293    The utilization of large scale of tissues allows us to investigate which tissue is better to capture more
294    transcripts and to annotate chicken genome. Herein we tried to detect the number of unique
295    transcripts expressed as a function of more tissues added. By doing so, we found brain-related tissues
296    (i.e., cerebellum and cortex) could detect the higher number of transcripts as expected (**Figure 5b**,
297    **Supplementary Table 6**). In addition, our design including a diverse set of 19 chicken tissues offers
298    the opportunity to analyze DAS events between chicken tissues. To do so, we quantified isoform
299    expression and identified differential alternative splicing events using the LIQA software (Hu et al.,
300    2021). The results are shown in **Supplementary Figure 7** and **Supplementary Table 7**. In total, we
301    found a list of 4,211 loci showing DAS events between tissues (FDR < 0.05). For instance, the top
302    significant locus is the *CYB561A3* gene showing DAS between heart and testis (FDR = 9.12E-16,
303    **Figure 5c**). This gene encodes cytochrome B561 family member A3 whose functions are related
304    cellular iron ion homeostasis and mitochondrial respiration (Wang et al., 2021b).

305

## 4    Discussion

307    A well-annotated chicken genome is essential in associating genomic variation to phenotypic
308    variation, and there are a number of ongoing efforts through the Functional Annotation of Animal
309    Genomes (FAANG) consortium (Andersson et al., 2015), primarily focus on non-coding functional
310    elements in farm animals including chicken (Kern et al., 2021). In this study, using Oxford Nanopore
311    long-read sequencing in 19 chicken tissues, we preliminary annotated 79,757 transcripts in 54,551
312    loci, while the subsequent filtering resulted in the exclusion of ~2K transcripts. Finally, our
313    prediction resulted in the identification of 55,382 clean transcripts derived from 40,547 loci,
314    representing ~1.4 transcripts per locus, an estimate lower than the Ensembl (~1.8 transcripts per
315    locus), and the NCBI annotations (~3.3 transcripts per locus). The lower estimate in our study might
316    be due to the higher number of annotated loci (N = 40,547), i.e. around 2.6-fold higher than both
317    reference annotations.

318    The number of transcripts of loci predicted in this study is substantially higher than two reference
319    annotations (Esembl V102: 27,955 transcripts in 15,305 loci; NCBI V105: 51,222 in 15,706 loci),
320    while our prediction is lower than Kuo et al. who annotated 60K transcripts and 29K genes using the
321    Iso-Seq approach (Kuo et al., 2017). Unfortunately, the unavailability of their annotation hinders us
322    to exclusively make a comparison. Specifically, we predicted higher proportion of lncRNA loci,
323    indicating that reference annotations are not annotated lncRNAs well. Indeed, Jehl et al. (2020)
324    annotated additional 13,009 lncRNA genes (compared to Ensembl V94) using 364 chicken short-read
325    transcriptomes derived from 25 tissues. Indeed, when we compared our lncRNA transcripts to Jehl et
326    al. ( 2020), we found over 80% of them completely or partially match to their lncRNA loci. Still, our
327    annotation contains 4,953 additional novel lncRNA transcripts despite we used the sample lncRNA
328    prediction tool (FEELnc, Wucher et al., 2017), which may be due to the increased sensitivity of long-
329    read sequencing (Lagarde et al., 2017). Moreover, we found > 89% of novel transcripts reported by
330    Thomas et al. (2014) could match our prediction. These evidences collectively indicate our
331    annotation is reliable.

332    Comparing to the reference annotations, we observed a higher percentage of novel loci (~47%) than
333    that of cattle (6% predicted transcripts did not match to a reference gene), whilst the exact matched

8

334 transcripts predicted in this study were also lower (16% in our study vs. 21% in cattle), though the
335 cattle study included more tissues (Halstead et al., 2021). Potential reasons are low number of
336 samples, possible degradations of RNA samples or low sequencing depth. We also cannot rule out
337 the possibility that the annotation of the bovine reference genome is better than the chicken one in the
338 database. It should be noted that a substantial proportion of novel loci predicted by us are lncRNA
339 loci which can to some extend match to a previous study (Jehl et al., 2020). These results suggest
340 more efforts for annotating the chicken genome is needed in the future. The human genome is
341 considered to be better-annotated than farm animals', while 36.4% of full-length transcripts identified
342 by long-read Iso-Seq methods are classified as "novel" in human cortex tissue (Leung et al., 2021).
343 Using the same approach, another study also reported 17 to 55% of novel isoforms in human breast
344 cancer samples (Veiga et al., 2022). These studies, together with ours, indicate long-read sequencing
345 is better approach for discovering novel isoforms and being able to better annotate animal genomes.

346 The number of transcripts reported by this study, reference genome annotations, as well as by Kuo et
347 al. (2017) varies widely, ranging from 27,955 to 74,665. Although sequencing depth could be one of
348 reasons, another possible interpretation is that the number of detectable transcripts is tissue-
349 dependent. Indeed, our study with similar sequencing depth also detected variable number of
350 expressed transcripts across tissues, ranging from 14,841 (liver) to 28,648 (cerebellum). These
351 observations suggest that including as diverse and many tissues as possible can detect tissue-specific
352 transcripts and better annotate the chicken genome. It is reported that brain tissues have a higher level
353 of alternative splicing, such as skipped exons, alternative 3' splice site exons or 5' splice site exons,
354 (Yeo et al., 2004; Melé et al., 2015). Our analysis supported this notion, suggesting brain-related
355 tissues are better for annotating an animal genome if available tissues are limited. The consistent
356 pattern of the higher number of transcripts observed in brain possibly reflects the complexity of
357 tissue biology (Naumova et al., 2013; Fang et al., 2020). Moreover, the whole embryo was also
358 expected to include as many transcripts as possible since it contains all organs. Unfortunately, our
359 study design did not include the whole embryo, but in Kuo et al. study they identified 55,932
360 transcripts in brain while only 9,368 transcripts in embryo (Kuo et al., 2017).

361 Although our study has annotated a substantial proportion of novel transcripts, there were still some
362 limitations, e.g., our study only includes a single developmental stage (adult). Previous reports
363 indicate that detecting gene expression using long-read sequencing approaches requires lower
364 number of reads, such as Nanopore sequencing needs ~ 40-fold less reads or ~ 8-fold less bases than
365 Illumina technology, which required over 36 million reads for accurately quantifying highly
366 expressed genes (FPKM > 10), and over 80 million reads for lowly expressed genes (FPKM < 10)
367 (Sims et al., 2014; Su et al., 2014; Oikonomopoulos et al., 2020). Based on that estimate, at least 7.5
368 million long-reads are likely to be required per tissue, while this will be the cost-prohibitive given the
369 cost of Nanopore long-read sequencing we did in 2019 with so many samples. This indicates our
370 study, very possibly, missed a proportion of transcripts due to the low sequencing depth, though we
371 reported a higher number of transcripts and loci than reference annotations. This is also reflected by
372 the ratio that each gene can only produce 1.4 transcripts per locus based on our data, while each
373 human gene can produce ~10 splicing transcripts (Mathur et al., 2019). With continued DNA
374 technologies development in cost-effective, tissues from more developmental stages and
375 physiological status, and more in-depth sequencing on full-length transcriptome are warranted to
376 improve the annotation of transcript isoforms in the chicken genome.

377 **1    Conflict of Interest**

378 None.

379

9

## 2 Author Contributions

HZ and PR conceived and designed the experiments. ADI, DEG, HC collected samples and carried out nanopore sequencing experiments. DG and MM developed the computational pipeline and analyzed all data. DG and HJ wrote the paper. All authors read, edited and approved the final manuscript.

## 5 References

Amarasinghe, S. L., Su, S., Dong, X., Zappia, L., Ritchie, M. E., and Gouil, Q. (2020). Opportunities and challenges in long-read sequencing data analysis. *Genome Biology* 21, 30. doi: 10.1186/s13059-020-1935-5.

Anders, S., Pyl, P. T., and Huber, W. (2015). HTSeq—a Python framework to work with high-throughput sequencing data. *Bioinformatics* 31, 166–169. doi: 10.1093/bioinformatics/btu638.

Andersson, L., Archibald, A. L., Bottema, C. D., Brauning, R., Burgess, S. C., Burt, D. W., et al. (2015). Coordinated international action to accelerate genome-to-phenome with FAANG, the Functional Annotation of Animal Genomes project. *Genome Biology* 16, 57. doi: 10.1186/s13059-015-0622-4.

Baralle, F. E., and Giudice, J. (2017). Alternative splicing as a regulator of development and tissue identity. *Nat Rev Mol Cell Biol* 18, 437–451. doi: 10.1038/nrm.2017.27.

Beiki, H., Liu, H., Huang, J., Manchanda, N., Nonneman, D., Smith, T. P. L., et al. (2019). Improved annotation of the domestic pig genome through integration of Iso-Seq and RNA-seq data. *BMC Genomics* 20, 344. doi: 10.1186/s12864-019-5709-y.

Benjamini, Y., and Hochberg, Y. (1995). Controlling the False Discovery Rate: A Practical and Powerful Approach to Multiple Testing. *Journal of the Royal Statistical Society: Series B (Methodological)* 57, 289–300. doi: 10.1111/j.2517-6161.1995.tb02031.x.

417   Buchfink, B., Xie, C., and Huson, D. H. (2015). Fast and sensitive protein alignment using
418         DIAMOND. *Nat Methods* 12, 59–60. doi: 10.1038/nmeth.3176.

419   Byrne, A., Beaudin, A. E., Olsen, H. E., Jain, M., Cole, C., Palmer, T., et al. (2017). Nanopore long-
420         read RNAseq reveals widespread transcriptional variation among the surface receptors of
421         individual B cells. *Nat Commun* 8, 16027. doi: 10.1038/ncomms16027.

422   Chen, S.-Y., Deng, F., Jia, X., Li, C., and Lai, S.-J. (2017). A transcriptome atlas of rabbit revealed
423         by PacBio single-molecule long-read sequencing. *Sci Rep* 7, 7648. doi: 10.1038/s41598-017-
424         08138-z.

425   Chen, Y., Davidson, N. M., Wan, Y. K., Patel, H., Yao, F., Low, H. M., et al. (2021). A systematic
426         benchmark of Nanopore long read RNA sequencing for transcript level analysis in human cell
427         lines. *bioRxiv*, 2021.04.21.440736. doi: 10.1101/2021.04.21.440736.

428   Clark, E. L., Archibald, A. L., Daetwyler, H. D., Groenen, M. A. M., Harrison, P. W., Houston, R.
429         D., et al. (2020). From FAANG to fork: application of highly annotated genomes to improve
430         farmed animal production. *Genome Biology* 21, 285. doi: 10.1186/s13059-020-02197-8.

431   De Coster, W., D'Hert, S., Schultz, S. T., Cruts, M., and Broeckhoven, C. V. (2018). NanoPack:
432         visualizing and processing long-read sequencing data. *Bioinformatics* 34, 26666–2669. doi:
433         10.1093/bioinformatics/bty149.

434   Fang, L., Cai, W., Liu, S., Canela-Xandri, O., Gao, Y., Jiang, J., et al. (2020). Comprehensive
435         analyses of 723 transcriptomes enhance genetic and biological interpretations for complex
436         traits in cattle. *Genome Res.* doi: 10.1101/gr.250704.119.

437   Giuffra, E., Tuggle, C. K., and FAANG Consortium (2019). Functional Annotation of Animal
438         Genomes (FAANG): Current Achievements and Roadmap. *Annu Rev Anim Biosci* 7, 65–88.
439         doi: 10.1146/annurev-animal-020518-114913.

440   Halstead, M. M., Islas-Trejo, A., Goszczynski, D. E., Medrano, J. F., Zhou, H., and Ross, P. J.
441         (2021). Large-Scale Multiplexing Permits Full-Length Transcriptome Annotation of 32
442         Bovine Tissues From a Single Nanopore Flow Cell. *Front. Genet.* 0. doi:
443         10.3389/fgene.2021.664260.

444   Hillier, L. W., Miller, W., Birney, E., Warren, W., Hardison, R. C., Ponting, C. P., et al. (2004).
445         Sequence and comparative analysis of the chicken genome provide unique perspectives on
446         vertebrate evolution. *Nature* 432, 695–716. doi: 10.1038/nature03154.

447   Hu, Y., Fang, L., Chen, X., Zhong, J. F., Li, M., and Wang, K. (2021). LIQA: long-read isoform
448         quantification and analysis. *Genome Biology* 22, 182. doi: 10.1186/s13059-021-02399-8.

449   Jehl, F., Muret, K., Bernard, M., Boutin, M., Lagoutte, L., Désert, C., et al. (2020). An integrative
450         atlas of chicken long non-coding genes and their annotations across 25 tissues. *Sci Rep* 10,
451         20457. doi: 10.1038/s41598-020-77586-x.

452   Julien, P., Brawand, D., Soumillon, M., Necsulea, A., Liechti, A., Schütz, F., et al. (2012).
453         Mechanisms and Evolutionary Patterns of Mammalian and Avian Dosage Compensation.
454         *PLOS Biology* 10, e1001328. doi: 10.1371/journal.pbio.1001328.

455    Kang, Y.-J., Yang, D.-C., Kong, L., Hou, M., Meng, Y.-Q., Wei, L., et al. (2017). CPC2: a fast and
456          accurate coding potential calculator based on sequence intrinsic features. *Nucleic Acids*
457          *Research* 45, W12–W16. doi: 10.1093/nar/gkx428.

458    Keren, H., Lev-Maor, G., and Ast, G. (2010). Alternative splicing and evolution: diversification,
459          exon definition and function. *Nat Rev Genet* 11, 345–355. doi: 10.1038/nrg2776.

460    Kern, C., Wang, Y., Xu, X., Pan, Z., Halstead, M., Chanthavixay, G., et al. (2021). Functional
461          annotations of three domestic animal genomes provide vital resources for comparative and
462          agricultural research. *Nat Commun* 12, 1821. doi: 10.1038/s41467-021-22100-8.

463    Kovaka, S., Zimin, A. V., Pertea, G. M., Razaghi, R., Salzberg, S. L., and Pertea, M. (2019).
464          Transcriptome assembly from long-read RNA-seq alignments with StringTie2. *Genome*
465          *Biology* 20, 278. doi: 10.1186/s13059-019-1910-1.

466    Kuhn, R. M., Haussler, D., and Kent, W. J. (2013). The UCSC genome browser and associated tools.
467          *Briefings in Bioinformatics* 14, 144–161. doi: 10.1093/bib/bbs038.

468    Kuo, R. I., Cheng, Y., Zhang, R., Brown, J. W. S., Smith, J., Archibald, A. L., et al. (2020).
469          Illuminating the dark side of the human transcriptome with long read transcript sequencing.
470          *BMC Genomics* 21, 751. doi: 10.1186/s12864-020-07123-7.

471    Kuo, R. I., Tseng, E., Eory, L., Paton, I. R., Archibald, A. L., and Burt, D. W. (2017). Normalized
472          long read RNA sequencing in chicken reveals transcriptome complexity similar to human.
473          *BMC Genomics* 18, 323. doi: 10.1186/s12864-017-3691-9.

474    Lagarde, J., Uszczynska-Ratajczak, B., Carbonell, S., Pérez-Lluch, S., Abad, A., Davis, C., et al.
475          (2017). High-throughput annotation of full-length long noncoding RNAs with capture long-
476          read sequencing. *Nat Genet* 49, 1731–1740. doi: 10.1038/ng.3988.

477    Leger, A. (2020). a-slide/NanoCount. Available at: https://zenodo.org/badge/latestdoi/142873004.

478    Leung, S. K., Jeffries, A. R., Castanho, I., Jordan, B. T., Moore, K., Davies, J. P., et al. (2021). Full-
479          length transcript sequencing of human and mouse cerebral cortex identifies widespread
480          isoform diversity and alternative splicing. *Cell Reports* 37, 110022. doi:
481          10.1016/j.celrep.2021.110022.

482    Li, H. (2018). Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* 34, 3094–
483          3100.

484    Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., et al. (2009). The Sequence
485          Alignment/Map format and SAMtools. *Bioinformatics* 25, 2078–2079. doi:
486          10.1093/bioinformatics/btp352.

487    Li, M., Sun, C., Xu, N., Bian, P., Tian, X., Wang, X., et al. (2022). De novo assembly of 20 chicken
488          genomes reveals the undetectable phenomenon for thousands of core genes on micro-
489          chromosomes and sub-telomeric regions. *Molecular Biology and Evolution*, msac066. doi:
490          10.1093/molbev/msac066.

491    Lin, J., Guan, L., Ge, L., Liu, G., Bai, Y., and Liu, X. (2021). Nanopore-based full-length
492        transcriptome sequencing of Muscovy duck (Cairina moschata) ovary. *Poultry Science* 100,
493        101246. doi: 10.1016/j.psj.2021.101246.

494    Love, M. I., Huber, W., and Anders, S. (2014). Moderated estimation of fold change and dispersion
495        for RNA-seq data with DESeq2. *Genome Biology* 15, 550. doi: 10.1186/s13059-014-0550-8.

496    Mathur, M., Kim, C. M., Munro, S. A., Rudina, S. S., Sawyer, E. M., and Smolke, C. D. (2019).
497        Programmable mutually exclusive alternative splicing for generating RNA and protein
498        diversity. *Nat Commun* 10, 2673. doi: 10.1038/s41467-019-10403-w.

499    Melé, M., Ferreira, P. G., Reverter, F., DeLuca, D. S., Monlong, J., Sammeth, M., et al. (2015). The
500        human transcriptome across tissues and individuals. *Science* 348, 660–665. doi:
501        10.1126/science.aaa0355.

502    Mi, H., Muruganujan, A., Casagrande, J. T., and Thomas, P. D. (2013). Large-scale gene function
503        analysis with the PANTHER classification system. *Nat Protoc* 8, 1551–1566. doi:
504        10.1038/nprot.2013.092.

505    Naumova, O. Yu., Lee, M., Rychkov, S. Yu., Vlasova, N. V., and Grigorenko, E. L. (2013). Gene
506        Expression in the Human Brain: The Current State of the Study of Specificity and Spatio-
507        temporal Dynamics. *Child Dev* 84, 76–88. doi: 10.1111/cdev.12014.

508    Oikonomopoulos, S., Bayega, A., Fahiminiya, S., Djambazian, H., Berube, P., and Ragoussis, J.
509        (2020). Methodologies for Transcript Profiling Using Long-Read Technologies. *Frontiers in
510        Genetics* 11, 606. doi: 10.3389/fgene.2020.00606.

511    Pertea, G., and Pertea, M. (2020). GFF Utilities: GffRead and GffCompare. *F1000Research*. doi:
512        10.12688/f1000research.23297.2.

513    Sims, D., Sudbery, I., Ilott, N. E., Heger, A., and Ponting, C. P. (2014). Sequencing depth and
514        coverage: key considerations in genomic analyses. *Nat Rev Genet* 15, 121–132. doi:
515        10.1038/nrg3642.

516    Soneson, C., Yao, Y., Bratus-Neuenschwander, A., Patrignani, A., Robinson, M. D., and Hussain, S.
517        (2019). A comprehensive examination of Nanopore native RNA sequencing for
518        characterization of complex transcriptomes. *Nat Commun* 10, 3359. doi: 10.1038/s41467-
519        019-11272-z.

520    Su, Z., Łabaj, P. P., Li, S., Thierry-Mieg, J., Thierry-Mieg, D., Shi, W., et al. (2014). A
521        comprehensive assessment of RNA-seq accuracy, reproducibility and information content by
522        the Sequencing Quality Control Consortium. *Nat Biotechnol* 32, 903–914. doi:
523        10.1038/nbt.2957.

524    Thomas, S., Underwood, J. G., Tseng, E., Holloway, A. K., and Subcommittee,  on behalf of the B.
525        T. B. C. I. (2014). Long-Read Sequencing of Chicken Transcripts and Identification of New
526        Transcript Isoforms. *PLOS ONE* 9, e94650. doi: 10.1371/journal.pone.0094650.

527    Tixier-Boichard, M., Fabre, S., Dhorne-Pollet, S., Goubil, A., Acloque, H., Vincent-Naulleau, S., et
528         al. (2021). Tissue Resources for the Functional Annotation of Animal Genomes. *Frontiers in*
529         *Genetics* 12, 847. doi: 10.3389/fgene.2021.666265.

530    Veiga, D. F. T., Nesta, A., Zhao, Y., Mays, A. D., Huynh, R., Rossi, R., et al. (2022). A
531         comprehensive long-read isoform analysis platform and sequencing resource for breast
532         cancer. *Science Advances*. doi: 10.1126/sciadv.abg6711.

533    Wang, K., Hu, H., Tian, Y., Li, J., Scheben, A., Zhang, C., et al. (2021a). The chicken pan-genome
534         reveals gene content variation and a promoter region deletion in IGF2BP1 affecting body
535         size. *Molecular Biology and Evolution*. doi: 10.1093/molbev/msab231.

536    Wang, Z., Guo, R., Trudeau, S. J., Wolinsky, E., Ast, T., Liang, J. H., et al. (2021b). CYB561A3 is
537         the key lysosomal iron reductase required for Burkitt B-cell growth and survival. *Blood* 138,
538         2216–2230. doi: 10.1182/blood.2021011079.

539    Warren, W. C., Hillier, L. W., Tomlinson, C., Minx, P., Kremitzki, M., Graves, T., et al. (2017). A
540         New Chicken Genome Assembly Provides Insight into Avian Genome Structure. *G3*
541         *Genes|Genomes|Genetics* 7, 109–117. doi: 10.1534/g3.116.035923.

542    Wright, C. J., Smith, C. W. J., and Jiggins, C. D. (2022). Alternative splicing as a source of
543         phenotypic diversity. *Nat Rev Genet*, 1–14. doi: 10.1038/s41576-022-00514-4.

544    Wucher, V., Legeai, F., Hédan, B., Rizk, G., Lagoutte, L., Leeb, T., et al. (2017). FEELnc: a tool for
545         long non-coding RNA annotation and its application to the dog transcriptome. *Nucleic Acids*
546         *Research* 45, e57. doi: 10.1093/nar/gkw1306.

547    Yandell, M., and Ence, D. (2012). A beginner's guide to eukaryotic genome annotation. *Nat Rev*
548         *Genet* 13, 329–342. doi: 10.1038/nrg3174.

549    Yang, R. Y., Quan, J., Sodaei, R., Aguet, F., Segrè, A. V., Allen, J. A., et al. (2018). A systematic
550         survey of human tissue-specific gene expression and splicing reveals new opportunities for
551         therapeutic target identification and evaluation. *bioRxiv*, 311563. doi: 10.1101/311563.

552    Yeo, G., Holste, D., Kreiman, G., and Burge, C. B. (2004). Variation in alternative splicing across
553         human tissues. *Genome Biology* 5, R74. doi: 10.1186/gb-2004-5-10-r74.

554    Yu, S., Wang, G., Liao, J., and Tang, M. (2019). Five alternative splicing variants of the TYR gene
555         and their different roles in melanogenesis in the Muchuan black-boned chicken. *British*
556         *Poultry Science* 60, 8–14. doi: 10.1080/00071668.2018.1533633.

557    Zhang, D., Guelfi, S., Garcia-Ruiz, S., Costa, B., Reynolds, R. H., D'Sa, K., et al. (2020). Incomplete
558         annotation has a disproportionate impact on our understanding of Mendelian and complex
559         neurogenetic disorders. *Science Advances* 6, eaay8299. doi: 10.1126/sciadv.aay8299.

560

# 6    Supplementary Material

562    **Supplementary data 1** Predicted transcripts in the General Feature Format (GTF) format

563 **Supplementary Table 1** Information about tissue sampling used in this study

564 **Supplementary Table 2** Summary statistics of sequencing samples

565 **Supplementary Table 3** Predicted transcript types (including protein-coding, lncRNA and other
566 non-coding)

567 **Supplementary Table 4** A list of tissue-specific transcripts

568 **Supplementary Table 5** Functional enrichment of tissue-specific transcripts (only Biological
569 Process of Gene Ontology terms)

570 **Supplementary Table 6** Number of unique transcripts detected when adding more tissues

571 **Supplementary Table 7** A list of loci showing differential alternative splicing (DAS) events
572 between tissues

573 **Supplementary Figure 1** Hierarchical clustering of samples used in this study. The dendrogram is
574 built based on gene expressions quantified with Transcripts Per Million (TPM > 0.1). The distance
575 between individuals is indicated by 1-r, where r is the Pearson correlation coefficient

576 **Supplementary Figure 2** Dotplot depicting the number of sequencing reads (x-axis) and the number
577 of expressed genes (y-axis). A given gene was considered as expressed when Transcripts Per Million
578 (TPM) > 0.1. The red text indicates the outlier sample in principal component analysis (PCA) plot
579 (Figures 1b) and hierarchical clustering (Supplementary Figure 1).

580 **Supplementary Figure 3** GffCompare types when comparing our predicted transcripts to NCBI
581 annotation (V105)

582 **Supplementary Figure 4** GffCompare types when comparing protein-coding (a) and lncRNA loci
583 (b) predicted in this study with those predicted in Jehl et al., (2020).

584 **Supplementary Figure 5** GffCompare types when comparing novel transcripts reported by Thomas
585 et al. (2014) to our annotation.

586 **Supplementary Figure 6** Number of expression loci and transcripts (TPM > 0.1) across tissues

587 **Supplementary Figure 7** Number of loci showing differential alternative splicing (DAS) between
588 tissues

589

590 **7    Data Availability Statement**

591 The Nanopore sequencing data are accessible in the Sequence Read Archive (SRA) database of the
592 National Center for Biotechnology Information with the identifier PRJNA671673
593 (https://www.ncbi.nlm.nih.gov/bioproject/PRJNA671673). The code used for annotating full-length
594 transcripts can be accessed by the link: https://github.com/guandailu/nanopore_annotation.

595

596 **Table 1** Comparison of reference and predicted transcripts using GffCompare tool

| Level | Predicted vs Ensembl | | Predicted vs NCBI | | NCBI vs Ensembl | |
|---|---|---|---|---|---|---|
| | Sensitivity | Precision | Sensitivity | Precision | Sensitivity | Precision |
| Base | 70.9 | 30.6 | 54.1 | 41.3 | 86.6 | 43.6 |
| Exon | 62.6 | 55.3 | 55.1 | 55.6 | 78.5 | 64.5 |
| Intron | 66.3 | 74.2 | 58.8 | 77.5 | 88.6 | 72.2 |

15

| Transcript | 38.7 | 14.5 | 21.1 | 14.5 | 41.6 | 21.1 |
|---|---|---|---|---|---|---|
| Locus | 57.8 | 17.5 | 54.3 | 17.0 | 59.7 | 47.3 |
| Missed exons | 44,538/179919 (24.8%) | | 60,304/211468 (28.5%) | | 10,378/202,369 (5.1%) | |
| Novel exons | 63,322/201,393 (31.4%) | | 54,465/201,393 (27.0%) | | 50,528/252,210 (20.0%) | |
| Missed introns | 41,164/157,463 (26.1%) | | 53,133/185,508 (28.6%) | | 6,790/175,889 (3.9%) | |
| Novel introns | 22,985/140,865 (16.3%) | | 19,416/140,865 (13.8%) | | 30,813/215,950 (14.3%) | |
| Novel loci | 32,725/50,569 (64.7%) | | 29,332/50,569 (58.0%) | | 5,656/23,336 (24.2%) | |

597    The annotation versions of NCBI and Ensembl are V105 and V102, respectively.

598

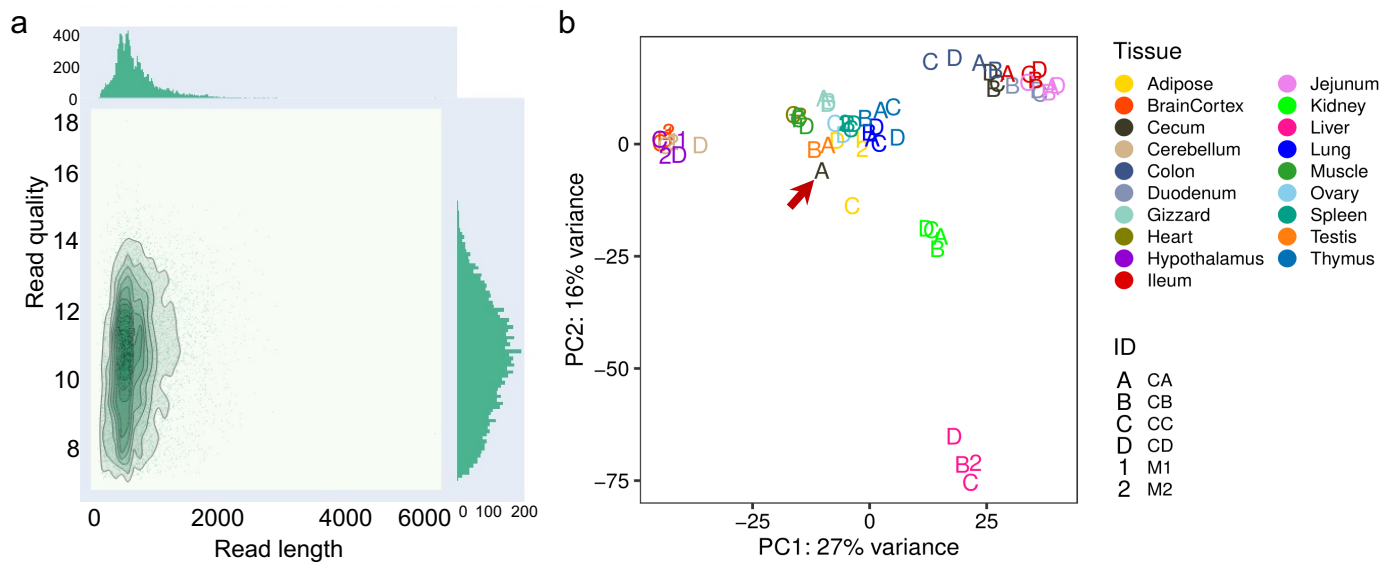599    **Figure Legends**

600    **Figure 1** (a) Bivariate plot (De Coster et al., 2018) depicting read length (x-axis) and quality (y-axis)
601    of Nanopore long-read sequencing in 68 samples (b) Principal component analysis of 68 chicken
602    Nanopore long-read transcriptomes. The red arrow indicated the sample, CA_Cecum, which was not
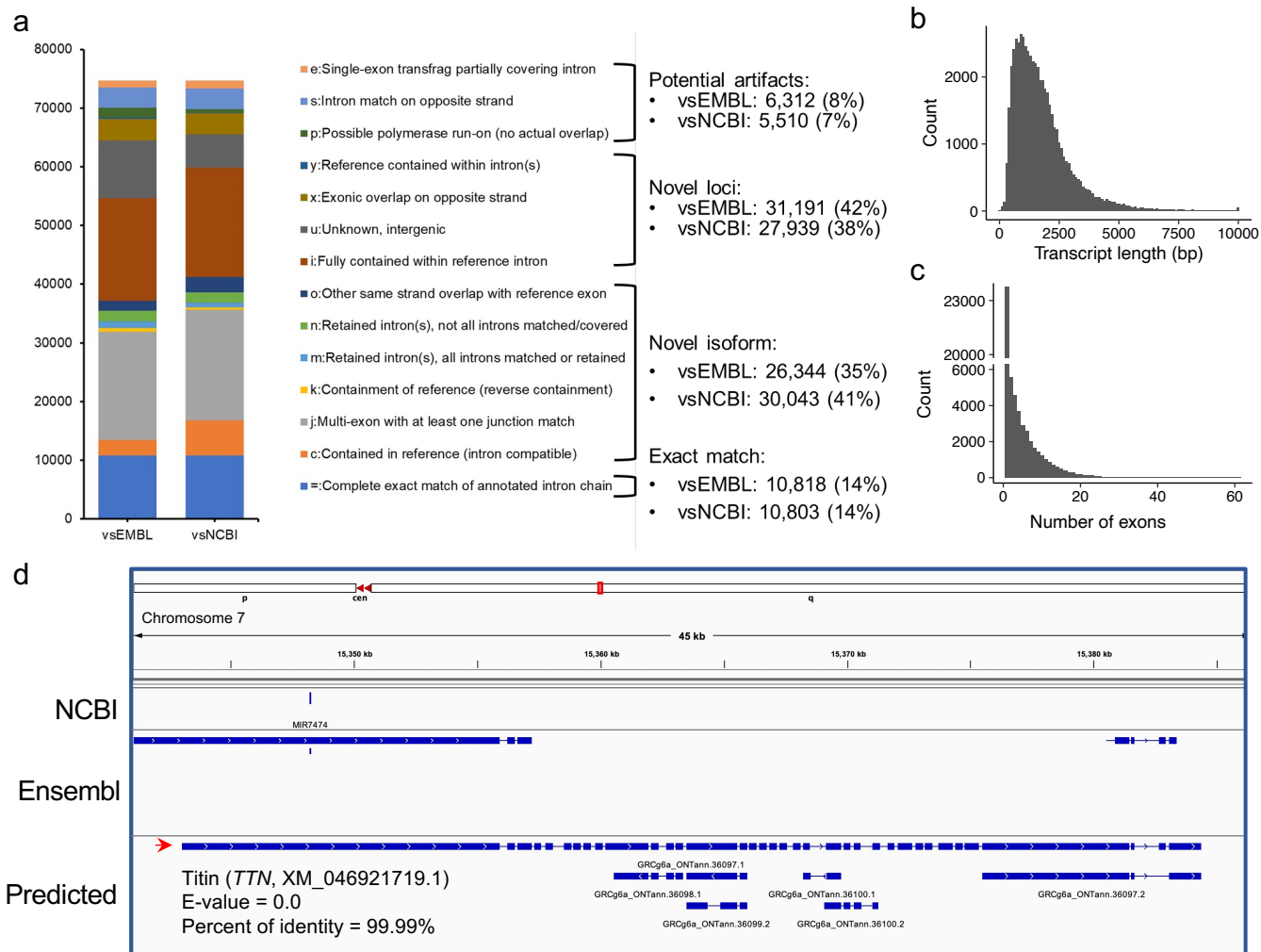603    clustered with other samples from the cecum tissue.

604

605    **Figure 2** (a) Comparisons of predicted transcripts against Ensembl (V102, vsEMBL) and NCBI
606    annotation (V105, vsNCBI). The transcripts were classified according to the GffCompare software
607    (Pertea and Pertea, 2020). The panels (b) to (c) depict the distributions of predicted transcript length
608    and exon numbers, respectively. (d) A screenshot showing the predicted longest transcript, which is
609    located on chromosome 7 (15,343,033-15,384,347). Blast analysis indicated the transcript matched to
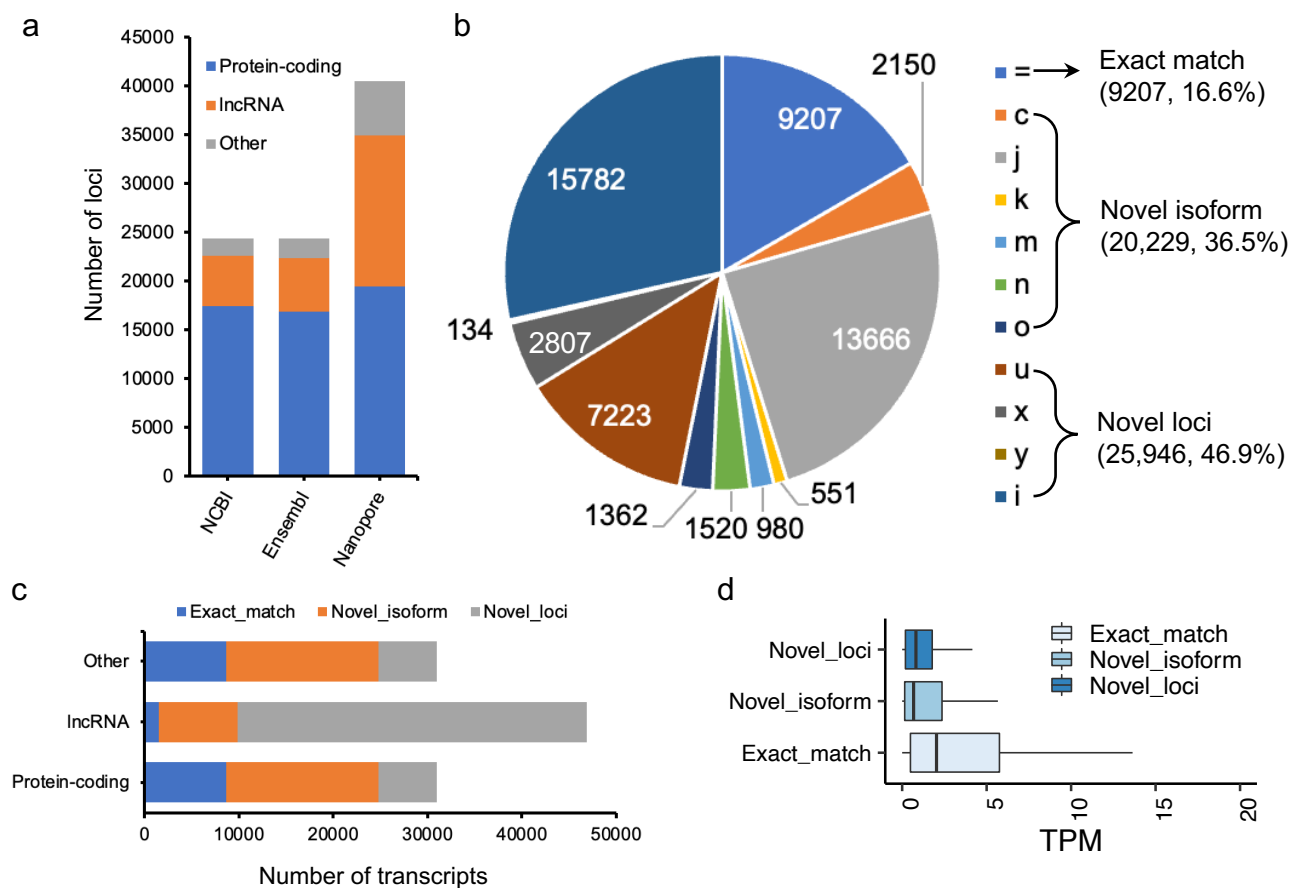610    the *TTN* gene locus encoding the titin protein.

611    **Figure 3** (a) Number of loci in NCBI (V105), Ensembl (V102) and our annotations. (b) Pie chart
612    depicting GffCompare types to Ensembl annotation (V102). See **Methods** for explanation of the type
613    codes. (c) Number of transcripts as a function of protein-coding, lncRNA, and other non-coding loci.
614    (d) Transcript expression measured as transcript per million (TPM) as a function of different types of
615    transcripts classified by GffCompare (See **Methods**).

616

617    **Figure 4** (a) Tissue specificity index (TSI) as a function of different types of transcripts classified by
618    GffCompare (See **Methods**). (b) Transcript expression measured as transcript per million (TPM) as a
619    function of tissue specificity index (TSI). We grouped transcripts according to their expressions (see
620    **Methods**). (c) Number of tissue-specific transcripts in each tissue. (d) A screenshot showing a novel
621    transcript only predicted by our data, which is located on chromosome 4 (52,482,563-52,492,561).
622    The transcript is highly expressed in testis samples, but not any other tissue samples. The FEELnc
623    predicted it as a sense intergenic lncRNA.

624

625    **Figure 5** (a) Heatmap depicting the negative $\log_{10}$ FDR (false discovery rate) values for the top 10
626    Gene Ontology (GO) Biological Process terms. At the right side, we show several examples of GO
627    terms, as well as their FDR values. (b) Number of unique transcripts detected as a function of tissues
628    added. Transcripts are categories into three types (see **Methods**). (c). Sashimi plots of *CYB561A3*
629    gene which showed DAS between heart (red) and testis (blue).

a



Potential artifacts:
- vsEMBL: 6,312 (8%)
- vsNCBI: 5,510 (7%)

Novel loci:
- vsEMBL: 31,191 (42%)
- vsNCBI: 27,939 (38%)

Novel isoform:
- vsEMBL: 26,344 (35%)
- vsNCBI: 30,043 (41%)

Exact match:
- vsEMBL: 10,818 (14%)
- vsNCBI: 10,803 (14%)

e:Single-exon transfrag partially covering intron
s:Intron match on opposite strand
p:Possible polymerase run-on (no actual overlap)
y:Reference contained within intron(s)
x:Exonic overlap on opposite strand
u:Unknown, intergenic
i:Fully contained within reference intron
o:Other same strand overlap with reference exon
n:Retained intron(s), not all introns matched/covered
m:Retained intron(s), all introns matched or retained
k:Containment of reference (reverse containment)
j:Multi-exon with at least one junction match
c:Contained in reference (intron compatible)
=:Complete exact match of annotated intron chain

b



c



d

a

b

c

d

a



Actin cytoskeleton organization (FDR = 5.08E-23)
Cytoskeleton organization (FDR = 1.14E-21)
Muscle cell differentiation (FDR = 1.95E-21)
Myofibril assembly (FDR = 3.39E-21)

Muscle contraction (FDR = 6.94E-40)
Muscle structure development (FDR = 4.06E-34)
Striated muscle cell differentiation (FDR = 1.84E-22)

Synaptic signaling (FDR = 2.1E-24)
Trans-synaptic signaling (FDR = 2.16E-24)
Anterograde trans-synaptic signaling (FDR = 3.27E-23)
Chemical synaptic transmission (FDR = 4.37E-23)

B cell receptor signaling pathway (FDR = 9.19E-18)
Immune response-activating signal transduction
(FDR = 7.05E-17)

Negative regulation of endopeptidase activity (FDR = 2.83E-35)
Negative regulation of proteolysis (FDR = 5.06E-35)
Regulation of peptidase activity (FDR = 4.3E-26)
Regulation of hydrolase activity (FDR = 8.89E-19)

$-\log_{10}(\text{FDR})$
0  10  20  30  40

b



c