# GigaScience

## Facilitating Functional genomics of cattle through integration of multi-omics data
### --Manuscript Draft--

| Manuscript Number: | GIGA-D-23-00037 | |
|---|---|---|
| Full Title: | Facilitating Functional genomics of cattle through integration of multi-omics data | |
| Article Type: | Research | |
| Funding Information: | National Institute of Food and Agriculture (2018-67015-27500) | Dr Huaijun Zhou |
| | National Institute of Food and Agriculture (2015-67015-22940) | Dr Huaijun Zhou |

| Abstract: | Background |
|---|---|
| | The accurate identification of the functional elements in the bovine genome is a fundamental requirement for high quality analysis of data informing both genome biology and genomic selection. Functional annotation of the bovine genome was performed to identify a more complete catalogue of transcript isoforms across bovine tissues. |
| | Results |
| | A total number of 171,985 unique transcripts (50% protein-coding) representing 35,150 unique genes (64% protein-coding) were identified across tissues. Among them, 118,563 transcripts (70% of the total) were structurally validated by independent datasets (PacBio Iso-seq data, ONT-seq data, de novo assembled transcripts from RNA-seq data) and comparison with Ensembl and NCBI gene sets. In addition, all transcripts were supported by extensive data from different technologies such as WTTS-seq, RAMPAGE, ChIP-seq, and ATAC-seq. A large proportion of identified transcripts (69%) were un-annotated, of which 87% were produced by annotated genes and 13% by un-annotated genes. A median of two 5' untranslated regions were expressed per gene. Around 50% of protein-coding genes in each tissue were bifunctional and transcribed both coding and noncoding isoforms. Furthermore, we identified 3,744 genes that functioned as non-coding genes in fetal tissues, but as protein coding genes in adult tissues. Our new bovine genome annotation extended more than 11,000 annotated gene borders compared to Ensembl or NCBI annotations. The resulting bovine transcriptome was integrated with publicly available QTL data to study tissue-tissue interconnection involved in different traits and construct the first bovine trait similarity network. |
| | Conclusions |
| | These validated results show significant improvement over current bovine genome annotations. |

| Corresponding Author: | James Reecy<br>Iowa State University<br>Ames, IA UNITED STATES |
|---|---|
| Corresponding Author Secondary Information: | |
| Corresponding Author's Institution: | Iowa State University |
| Corresponding Author's Secondary Institution: | |
| First Author: | Hamid Beiki |
| First Author Secondary Information: | |
| Order of Authors: | Hamid Beiki |
| | Brenda M. Murdoch |

| | |
|---|---|
| | Carissa A. Park |
| | Chandlar Kern |
| | Denise Kontechy |
| | Gabrielle Becker |
| | Gonzalo Rincon |
| | Honglin Jiang |
| | Huaijun Zhou |
| | Jacob Thorne |
| | James E. Koltes |
| | Jennifer J. Michal |
| | Kimberly Davenport |
| | Monique Rijnkels |
| | Pablo J. Ross |
| | Rui Hu |
| | Sarah Corum |
| | Stephanie McKay |
| | Timothy P.L. Smith |
| | Wansheng Liu |
| | Wenzhi Ma |
| | Xiaohui Zhang |
| | Xiaoqing Xu |
| | Xuelei Han |
| | Zhihua Jiang |
| | Zhi-Liang Hu |
| | James Reecy |
| **Order of Authors Secondary Information:** | |
| **Additional Information:** | |
| **Question** | **Response** |
| Are you submitting this manuscript to a special series or article collection? | No |
| **Experimental design and statistics**<br><br>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our [Minimum Standards Reporting Checklist](#). Information essential to interpreting the data presented should be made available in the figure legends. | Yes |

| | |
|---|---|
| Have you included all the information requested in your manuscript? | |
| **Resources**<br><br>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.<br><br>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist? | Yes |
| **Availability of data and materials**<br><br>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in publicly available repositories (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the "Availability of Data and Materials" section of your manuscript.<br><br>Have you have met the above requirement as detailed in our Minimum Standards Reporting Checklist? | Yes |

# Facilitating Functional genomics of cattle through integration of multi-omics data

Hamid Beiki[1], Brenda M. Murdoch[2], Carissa A. Park[1], Chandlar Kern[3], Denise Kontechy[2], Gabrielle Becker[2], Gonzalo Rincon[4], Honglin Jiang[5], Huaijun Zhou[6], Jacob Thorne[2], James E. Koltes[1], Jennifer J. Michal[7], Kimberly Davenport[2], Monique Rijnkels[8], Pablo J. Ross[6], Rui Hu[5], Sarah Corum[4], Stephanie McKay[9], Timothy P.L. Smith[10], Wansheng Liu[3], Wenzhi Ma[3], Xiaohui Zhang[7], Xiaoqing Xu[6], Xuelei Han[7], Zhihua Jiang[7], Zhi-Liang Hu[1], James M. Reecy[1]


[1]Department of Animal Science, Iowa State University; [2]Department of Animal and Veterinary and Food Science, University of Idaho; [3]Department of Animal Science, Pennsylvania State University; [4]Zoetis; [5]Department of Animal and Poultry Sciences, Virginia Tech; [6]Department of Animal Science, University of California, Davis; [7]Department of Animal Science, Washington State University; [8]Department of Veterinary Integrative Biosciences, Texas A&M University; [9]University of Vermont; [10]USDA, ARS, USMARC.


**Corresponding author:**

James M. Reecy

19    Professor of Animal Breeding and Genetics, Department of Animal Science, Ames, IA, USA

20    jreecy@iastate.edu

21

## Abstract

### Background

24    The accurate identification of the functional elements in the bovine genome is a fundamental

25    requirement for high quality analysis of data informing both genome biology and genomic

26    selection. Functional annotation of the bovine genome was performed to identify a more

27    complete catalogue of transcript isoforms across bovine tissues.

### Results

29    A total number of 171,985 unique transcripts (50% protein-coding) representing 35,150 unique

30    genes (64% protein-coding) were identified across tissues. Among them, 118,563 transcripts

31    (70% of the total) were structurally validated by independent datasets (PacBio Iso-seq data,

32    ONT-seq data, *de novo* assembled transcripts from RNA-seq data) and comparison with

33    Ensembl and NCBI gene sets. In addition, all transcripts were supported by extensive data from

34    different technologies such as WTTS-seq, RAMPAGE, ChIP-seq, and ATAC-seq. A large

35    proportion of identified transcripts (69%) were un-annotated, of which 87% were produced by

36    annotated genes and 13% by un-annotated genes. A median of two 5' untranslated regions

37    were expressed per gene. Around 50% of protein-coding genes in each tissue were bifunctional

38    and transcribed both coding and noncoding isoforms. Furthermore, we identified 3,744 genes

39    that functioned as non-coding genes in fetal tissues, but as protein coding genes in adult

40    tissues. Our new bovine genome annotation extended more than 11,000 annotated gene

41    borders compared to Ensembl or NCBI annotations. The resulting bovine transcriptome was

42    integrated with publicly available QTL data to study tissue-tissue interconnection involved in

43    different traits and construct the first bovine trait similarity network.

44    **Conclusions**

45    These validated results show significant improvement over current bovine genome

46    annotations.

47    # Introduction

48    Domestic bovine (*Bos taurus*) provides a valuable source of nutrition and an important disease

49    model for humans [1]. Furthermore, cattle have the greatest number of genotype associations

50    and genetic correlations of the domesticated livestock species, which means they provide an

51    excellent model to close the genotype-to-phenotype gap. Therefore, the accurate identification

52    of the functional elements in the bovine genome is a fundamental requirement for high quality

53    analysis of data informing both genome biology and genomic selection.

54    Current annotations of farm animal genomes largely focus on the protein-coding regions and

55    fall short of explaining the biology of many important traits that are controlled at the

56    transcriptional level [2]. In humans, 88% of trait-associated single nucleotide polymorphisms

57    (SNP) identified by genome-wide association studies (GWAS) are found in non-coding regions

58    [3]. Therefore, elucidating non-coding functional elements of the genome is essential for

59    understanding the mechanisms that control complex biological processes.

60    Untranslated regions play critical roles in the regulation of mRNA stability, translation, and

61    localization [4], but these regions have been poorly annotated in farm animals [2, 5]. A recent

62    study of the pig transcriptome using single-molecule long-read isoform sequencing technology

63    resulted in the extension of more than 6000 annotated gene borders compared to Ensembl or

64    National Center for Biotechnology Information (NCBI) annotations [2].

65    Small non-coding RNAs, such as microRNAs (miRNA), are known to be involved in gene

66    regulation through post-transcriptional regulation of expression via silencing, degradation, or

67    sequestering to inhibit translation [6-8]. The number of annotated miRNAs in the current

68    bovine genome annotation (Ensembl release 2018-11; 951 miRNAs) is much lower than the

69    number reported in the highly annotated human genome (Ensembl release 2021-03; 1,877

70    miRNAs).

71    This study applied a comprehensive set of transcriptome and chromatin state data from 47

72    cattle tissues and cell types to identify previously unannotated genes and improve the

73    annotation of thousands of protein-coding and non-coding genes. Predicted un-annotated

74    genes and transcripts were highly supported by independent Pacific Biosciences single-

75    molecule long-read isoform sequencing (PacBio Iso-Seq), Oxford Nanopore Technologies

76    sequencing (ONT-seq), Illumina high-throughput RNA sequencing (RNA-seq), Whole

77    Transcriptome Termini Site Sequencing (WTTS-seq), RNA Annotation and Mapping of

78    Promoters for the Analysis of Gene Expression (RAMPAGE), chromatin immunoprecipitation

79    sequencing (ChIP-seq), and Assay for Transposase-Accessible Chromatin using sequencing

80 (ATAC-seq) data. The transcriptome data was integrated with publicly available Quantitative

81 Trait Loci (QTL) and gene association data to construct the first bovine trait similarity network

82 that recapitulates published genetic correlations. Thus, it may be possible to begin to examine

83 the genetic mechanisms underlying genetic correlations.

84 **Results**

85 The diversity of RNA and miRNA transcript diversity among 47 different bovine tissues and cell

86 types was assessed using miRNA-seq and polyadenylation (poly(A)) selected RNA-seq data.

87 Most of the tissues studied were from Hereford cattle closely related to L1 Dominette 01449,

88 the individual from which the bovine reference genome (ARS-UCD1.2) was sequenced. The 47

89 tissues and cell samples included follicular cells, myoblasts, five mammary gland samples from

90 various stages of mammary gland development and lactation, eight fetal tissues (78-days of

91 gestation), eight tissues from adult digestive tract, and 16 other adult organs. A total of

92 approximately 4.1 trillion RNA-seq reads and 1.2 billion miRNA-seq reads were collected, with a

93 minimum of 27.5 million RNA-seq and 9.3 million miRNA-seq reads from each tissue/cell type

94 (average $87.8 \pm 49.7$ million and $27.6 \pm 12.9$ million, respectively) (Supplemental file 1: Fig. S1

95 and Supplemental file 2).

96 **Transcript level analyses**

97 A total of 171,985 unique transcripts (76% spliced) were identified (Table 1) with a median of

98 51,231 transcripts per tissue. There was a median of 9.1 exons per spliced transcript, and all of

99 the predicted acceptor and donor splice sites conformed to the canonical consensus sequences.

100 All of the predicted splice junctions across tissues were supported by RNA-seq reads that

101    spanned the splice junction, substantiating the accuracy of the transcript definition from RNA-

102    seq reads.

103    A total of 31,476 transcripts appeared tissue-specific by virtue of being assembled from RNA-

104    seq reads in just a single tissue, but 20,100 of those transcripts (64%) were actually expressed in

105    multiple tissues. Thus, reliance solely on assembled transcripts in a given tissue to predict a

106    tissue transcript atlas may overestimate tissue specificity due to a high false-negative rate for

107    transcript detection. To solve this problem of over-prediction of tissue specificity, we marked a

108    transcript as "expressed" in a given tissue only if (1) it had been assembled from RNA-seq data

109    in that tissue; or (2) its expression and all of its splice junctions has been quantified using RNA-

110    seq reads in the tissue of interest with an expression level more than 1 reads per kilobase of

111    transcript per Million reads mapped (RPKM) (see Methods section). This resulted in 11,375

112    apparently tissue-specific transcripts (7%) and 156,423 transcripts (91%) expressed in more

113    than one tissue (Fig. 1), among which 9,125 transcripts (5%) were found in all 47 tissues

114    examined.

115    The unique transcripts identified were equally distributed between 85,658 (50%) protein-

116    coding transcripts and 86,327 (50%) non-coding transcripts (ncRNAs) (Fig. 2). Non-coding

117    transcripts were further classified as long non-coding RNAs (lncRNAs) (56%), nonsense-

118    mediated decay (NMD) transcripts (38%), non-stop decay (NSD) transcripts (5%), and small non-

119    coding RNAs (sncRNAs) (1%). While the majority of expressed transcripts in each tissue were

120    protein coding (median of 62% of tissue transcripts), NMD transcripts (median of 14.58% of

121    tissue transcripts) and antisense lncRNAs (median of 12% of tissue transcripts) each made up

122    more than 10% of the transcripts (Supplemental file 1: Fig. S2A and B, Supplemental file 3 and

123    4). Fetal muscle and fetal gonad tissues showed the highest proportion of antisense lncRNAs

124    compared to that observed in other tissues (Supplemental file 1: Fig. S2B) and around 60% of

125    antisense lncRNAs (17,982 transcripts) were expressed from these two tissues. Compared to

126    non-coding transcripts, protein-coding transcripts were more likely to have spliced exons (p-

127    value < 2.2e-16) and were expressed in a higher number of tissues (median of 11 tissues for

128    protein-coding transcripts versus six tissues for non-coding transcripts; p-value < 2.2e-16)

129    (Additional file1: Fig. S2C). The lncRNAs had a significantly lower splice rate (36%) compared to

130    other non-coding transcripts (p-value < 2.2e-16). Splice rate was highest (70%) in sncRNAs (p-

131    value < 2.2e-16; NMD transcripts were not included in this analysis, as they were all spliced

132    transcripts by definition).

133    There were no significant correlations between the number of RNA-seq reads for a given tissue

134    and the number of unique transcripts identified, except for a modest correlation for the

135    antisense lncRNA class (Supplemental file 1: Fig. S3A). There was a significant positive

136    correlation (p-value 1.3e-04) between the number of unique NMD transcripts in a tissue and

137    the number of protein-coding transcripts, and the NMD transcript class showed the lowest

138    median expression level across tissues, followed by antisense-lncRNAs and sense intronic-

139    lncRNAs (Supplemental file 1: Fig. S2D and Fig. S3B). In addition, there was a significant positive

140    correlation (p-value 3.4e-03) between the number of NMD transcripts and the number of

141    protein-coding transcripts across tissues (Supplemental file 1: Fig. S3A). The expression levels of

142    sncRNAs and protein-coding transcripts were higher (p-values: 1.1e-02 and 2.6e-06,

143    respectively) than that observed for other transcript biotypes (Supplemental file 1: Fig. S2D and

144    Fig. S3B).

145 **Transcript similarity to other species**

146 Protein/peptide homology analysis of transcripts with an open reading frame (protein-coding

147 transcripts, lncRNAs, and sncRNAs) revealed a higher conservation of protein-coding transcripts

148 (86%) compared to lncRNA and sncRNA transcripts (8%; p-value < 2.2e-16) (Table 2). Bovine

149 non-coding transcripts had significantly (p-value < 2.2e-16) less similarity to other species than

150 protein-coding transcripts (Table 2 and Table 3). Within non-coding transcripts, NSD transcripts

151 showed the lowest conservation rate (35%), followed by sncRNAs (37%), lncRNAs (49%), and

152 NMD transcripts (55%), while sense intronic lncRNAs had the highest conservation rate (60%)

153 compared to other non-coding transcripts (Table 4).

154 **Transcript expression diversity across tissues**

155 A median of 70% of protein-coding transcripts were shared between pairs of tissues

156 (Supplemental file 1: Fig. S4A), significantly higher than that was observed for non-coding

157 transcripts (53%; p-value < 2.2e-16; Supplemental file 1: Fig. S5). Clustering of tissues based on

158 protein-coding transcripts was different than that observed based on non-coding transcripts

159 (Supplemental file 1: Fig. S4B and Fig. S5B, Fig. S35F). The fetal tissues clustered together and

160 were generally more similar to one another than to the corresponding adult tissue in both

161 dendrograms, but thymus was closely related to fetal tissues for protein-coding transcript

162 content, while it appeared more similar to lymph nodes, myoblasts, and pregnant/lactating

163 mammary tissue using non-coding transcript profiles. The digestive tract tissues clustered

164 together in the non-coding dendrogram with ileum as a slight outlier, while both jejunum and

165 ileum were distant from the other digestive tissues in the protein-coding transcript profile. The

166  "adult mammary gland" (78 day pregnant) and "virgin mammary gland" samples did not cluster

167  with the three other pregnant/lactating mammary samples nor with each other in either

168  dendrogram. This is mostly likely because: 1) these are from different physiological stages, 2)

169  these were whole tissue samples while the other three pregnant/lactating samples are enriched

170  for mammary gland epithelial cells, 3) the virgin and 78-day pregnant samples are from

171  Hereford background while other pregnant/lactating samples are from Holstein-Frisian breed.

172  Fetal tissues had significantly higher proportions than adult tissues of unique non-coding

173  transcripts (specifically NSDs, antisense lncRNAs, and intragenic lncRNAs) compared to protein-

174  coding transcripts (p-value < 2.2e-16; Supplemental file 5).

175  **Transcript validation**

176  Prediction of transcripts and isoforms from RNA-seq data may produce erroneous predicted

177  isoforms. The validity of transcripts was therefore examined by comparison to a library of

178  isoforms taken from Ensembl (release 2021-03) and NCBI gene sets (Release 106), as well as

179  isoforms identified through complete isoform sequencing with Pacific Biosciences, a de novo

180  assembly produced from its matched RNA-seq reads, and isoforms identified from Oxford

181  Nanopore platforms (see Methods section). A total of 118,563 transcripts (70% of predicted

182  transcripts) were structurally validated by independent datasets (PacBio Iso-seq data, ONT-seq

183  data, *de novo* assembled transcripts from RNA-seq data) and comparison with Ensembl and

184  NCBI gene sets. A total of 160,610 transcripts were expressed in multiple tissues (93% of

185  predicted transcripts), providing further support for their validity (Fig. 3). All transcripts were

186  also extensively supported by data from different technologies such as WTTS-seq, RAMPAGE,

187  histone modification (H3K4me3, H3K4me1, H3K27ac), CTCF-DNA binding, and ATAC-seq (Fig. 3).

188    Comparison of predicted transcript structures with annotated transcripts in the current bovine

189    genome annotations (Ensembl release 2021-03 and NCBI Release 106) resulted in a total of

190    52,645 annotated transcripts that exactly matched previously annotated transcripts (31% of all

191    transcripts), including 47,054 annotated NCBI transcripts, 31,740 annotated Ensembl

192    transcripts, and 26,149 transcripts that were common to both annotated gene sets (Fig. 3). The

193    median expression level of annotated transcripts in their expressed tissues (1.8 RPKM) was

194    similar to that observed for un-annotated transcripts (1.4 RPKM) (Supplemental file 1: Fig. S6).

195    Annotated transcripts were expressed in a median of 17 tissues, which was higher (p-value

196    7.4e-03) than that observed for un-annotated transcripts (median of seven tissues)

197    (Supplemental file 1: Fig. S6). In addition, compared to un-annotated transcripts, annotated

198    transcripts were enriched with protein-coding (p-value 1.37e-02) and spliced transcripts (p-

199    value 3.76e-02).

200    The median length of coding sequence (CDS) of annotated transcripts was 1,014 nt, significantly

201    longer than that observed in un-annotated transcripts (510 nt; p-value 0.0) (Additional file1: Fig.

202    S7A). In addition, un-annotated transcripts had longer 5' untranslated regions (UTR) (400 bp)

203    compared to that was observed in annotated transcripts (300 nt, p-value 2.631E-06; Additional

204    file1: Fig. S7A). Un-annotated transcripts encoding proteins with homology to proteins

205    annotated in other species had longer CDS (687 bp) compared to transcripts without such

206    homology (192 nt; p-value 0.0). Annotated protein-coding transcripts showed a higher GC

207    content in their 5' UTRs (61%) than un-annotated transcripts (53%; p-value 5.562E-18), but both

208    classes of transcripts showed similar GC content within their CDS (Supplemental file 1: Fig. S7B).

209 **Gene level analyses**

210 The transcripts correspond to a total of 35,150 genes, which were classified into protein coding

211 (21,193), non-coding (10,928), and pseudogenes (3,029) (Supplemental file 3 and 4, and Fig. 4).

212 Genes transcribed at least a single "expressed" transcript (see Transcript level analysis section)

213 in a given tissue, were marked as "expressed gene" in that tissue. Most genes expressed in each

214 tissue were protein coding (median of 83% of tissue genes), followed by non-coding (median of

215 14% of tissue genes) and pseudogenes (median of 3% of tissue genes) (Supplemental file 1: Fig.

216 S8). Testis showed the highest number of expressed genes with observed transcripts compared

217 to other tissues (Supplemental file 1: Fig. S8). Fetal brain and fetal muscle tissues showed the

218 highest number and percentage of non-coding genes compared to that observed in other

219 tissues (Supplemental file 1: Fig. S8). In addition, more than 40% of transcripts corresponded to

220 non-coding genes (1,271 genes) in fetal brain and fetal muscle. The proportion (6%) and

221 number (1,271) of transcript-producing pseudogenes was higher in testis than in other tissues.

222 There was no significant correlation between the number of input reads and the number of

223 expressed genes across tissues, but the numbers of genes from different coding potential

224 classes were significantly correlated across tissues (Supplemental file 1: Fig. S9).

225 Transcripts corresponding to the predicted genes that had at least one exon overlapping an

226 Ensembl- or NCBI-annotated gene were considered to belong to an annotated gene. This

227 supported an intersection analysis of predicted and previously annotated genes that indicated

228 22,452 (64%) of our predicted genes correspond to previously annotated genes. Approximately

229 87% of un-annotated transcripts (103,387) were associated with this set of annotated genes.

230 The remaining 12,698 genes (36% of predicted genes) represent un-annotated genes, i.e., genes

231  not found on Ensembl (release 2021-03) or NCBI (release 106), with which 15% of un-annotated

232  transcripts (22,364 transcripts) were associated. The median number of unique transcripts per

233  annotated gene (tpg) was four, which was higher than that observed in either the Ensembl (1.5

234  tpg) or NCBI (2.3 tpg) annotated gene sets, while the median number of transcripts per un-

235  annotated gene was one, with an average of 1.31 and standard deviation of 1.36. Most of the

236  transcripts identified were transcribed from annotated genes, including 96% of protein-coding

237  transcripts (82,060), 79% of lncRNA transcripts (38,662), 78% of sncRNA transcripts (413), and

238  more than 95% of NMD transcripts (31,422). Annotated genes were enriched with protein-

239  coding genes (p-value < 2.2e-16). The median transcript abundance from annotated genes in

240  their expressed tissues (6.59 RPKM) was significantly higher than that observed for un-

241  annotated genes (median of 1.68 RPKM; p-value < 2.2e-16; Supplemental file 1: Fig. S10A). The

242  median number of tissues in which annotated genes were expressed (42 tissues) was also

243  significantly higher than that observed for un-annotated genes (median of four tissues; p-value

244  < 2.2e-16; Supplemental file 1: Fig. S10B).

245  More than a third (37%) of genes with at least one predicted protein-coding transcript

246  displayed either multiple 5' UTRs or multiple 3' UTRs (median of three 5' UTRs and three 3'

247  UTRs per gene) among associated transcript isoforms (Fig. 5). The 496 genes with the highest

248  number of UTRs (the top 5% in this metric) were highly enriched (q-value 1.7E-7) for the

249  "response to protozoan" Biological Process (BP) Gene Ontology (GO) term (Supplemental file 1:

250  Fig. S11 and Supplemental file 6).

251  A median of 51% of the expressed protein-coding genes in each tissue transcribed both protein-

252  coding and non-coding transcripts and were denoted as bifunctional genes. These genes were

253 mostly previously annotated (95%) and had both coding and non-coding transcripts in a median

254 of 21 tissues, representing 57% of their expressed tissues (Fig. 6A and B). Protein-coding

255 transcripts and NMD transcripts covered more than 90% of the exonic length in bifunctional

256 genes (Fig. 6C). This percentage was significantly lower for other types of non-coding transcripts

257 transcribed from bifunctional genes (77%, 81%, and 62% for NSD transcripts, sncRNAs, and

258 intragenic lncRNAs, respectively) (Fig. 6C). Although transcript terminal sites (TTS) of transcripts

259 encoded by bifunctional genes were centralized around these genes' 3' ends, transcript start

260 sites (TSS) varied greatly among transcript biotypes (Fig. 6C). The TTSs of NSD transcripts,

261 sncRNAs, and intragenic lncRNAs were shifted from their protein-coding genes' start sites (Fig.

262 6C). Genes that transcribed both protein-coding and non-coding transcripts in all of their

263 expressed tissues (1,661 genes) were highly enriched for "mRNA processing" (q-value 6.08E-16)

264 and "RNA splicing" (q-value 1.35E-14) BP GO terms that were mostly (65%) related to different

265 aspects of transcription and translation (Fig. 6D and Supplemental file 7).

266 A total of 3,744 protein-coding genes (17% of all predicted protein-coding genes) only

267 transcribed non-coding transcripts in a median of two tissues (equivalent to 15% of their

268 expressed tissues). Detailed investigation of these genes in tissues from both adult and fetal

269 samples (brain, kidney, muscle, and spleen) revealed the total of 106 non-coding genes (90%

270 annotated) in fetal tissues that were switched to protein-coding genes with only protein-coding

271 transcripts in their matched adult tissues (Supplemental file 1: Fig. S12). Functional enrichment

272 analysis of these genes resulted in the identification of enriched BP GO terms related to

273 "humoral immune response", "sphingolipid biosynthetic process", "negative regulation of

274  wound healing", "cellular senescence", "symporter activity", "regulation of lipid biosynthetic

275  process", and "filopodium assembly" (Supplemental file 1: Fig. S12, Supplemental file 8).

276  A median of 32% of protein-coding genes in each tissue expressed at least a single potentially

277  aberrant transcript (PAT), i.e., NMDs and NSDs. In this group of genes, the number of PATs was

278  strongly correlated with the total number of transcripts (median correlation of 0.61 across all

279  tissues). The median expression level of these genes in their expressed tissues (11.52 RPKM)

280  was significantly higher (p-value < 2.2e-16) than for protein-coding genes with no PATs (4.48

281  RPKM). In each tissue, protein-coding genes with PATs showed a significantly higher number of

282  introns (p-value < 2.2e-16; median of 65 introns per gene) than that observed in the remainder

283  of protein-coding genes (median of 15 introns per gene). In addition, genes from this group

284  were expressed in a median of 47 tissues, significantly higher (p-value < 2.2e-16) than that

285  observed for the other coding genes (median of 24 tissues), non-coding genes (median of five

286  tissues), and pseudogenes (median of four tissues) (Supplemental file 1: Fig. S13A and B). These

287  genes transcribed a median of two PATs in half (median 54%) of their expressed tissues,

288  equivalent to a median of 22% of all their transcripts in each tissue. Protein-coding genes that

289  transcribed PATs as their main transcripts (PATs comprised >50% of their transcripts) in all of

290  their expressed tissues were highly enriched with RNA splicing–related BP GO terms

291  (Supplemental file 9).

292  **Gene similarity to other species**

293  Eighty-five percent of protein-coding genes (18,087) encoded either homologous proteins

294  (17,150 genes or 80% of protein-coding genes) or homologous ncRNAs (7,347 genes or 35% of

295    protein-coding genes) (Supplemental file 1: Fig. S14A). Nineteen percent of protein-coding

296    genes (4,043) encoded cattle-specific proteins (Supplemental file 1: Fig. S14A). Most of these

297    genes (2,750 or 68%) were either annotated genes or genes with homology to another cattle

298    gene(s) that has established homology to genes in other species (Supplemental file 1: Fig.

299    S14C). The remaining 32% of cattle-specific, protein-coding genes (1,293 genes or six percent of

300    protein-coding genes) were denoted as protein-coding orphan genes (Supplemental file 1: Fig.

301    S14C). A median of 70 protein-coding orphan genes were expressed in each tissue. The

302    expression level of these genes was significantly lower than other types of protein-coding genes

303    (Additional file1: Fig. S15A and B). The median number of expressed tissues for protein-coding

304    orphan genes (one tissue) was lower than for other types of protein-coding genes (46 tissues)

305    (Supplemental file 1: Fig. S15C). In addition, protein-coding orphan genes only transcribed

306    protein-coding transcripts in their expressed tissue(s).

307    Fifty percent of non-coding genes (5,559) encoded either homologous short peptides (9-43

308    amino acids; 5.8% of non-coding genes) or homologous ncRNAs (49% of non-coding genes)

309    (Supplemental file 1: Fig. S14B). There were 5,546 non-coding genes (51% of non-coding genes)

310    that encoded cattle-specific ncRNAs (Supplemental file 1: Fig. S14B). Ninety-nine percent of

311    these genes (5,537 genes) were either annotated genes or genes with homology to another

312    cattle gene(s) that has established homology to genes in other species (Supplemental file 1: Fig.

313    S14C). The remaining 1% (nine non-coding genes) were denoted as non-coding orphan genes

314    (Supplemental file 1: Fig. S14C). The median number of expressed tissues for non-coding

315    orphan genes was 17 tissues, which was higher (p-value < 2.2e-16) than for homologous non-

15

316    coding genes (six tissues) and protein-coding orphan genes (one tissue) (Supplemental file 1:

317    Fig. S15C).

318    A total of 3,029 pseudogenes were expressed. The median expression level of these genes in

319    their expressed tissues was 2.15 RPKM, which was lower than that observed for protein-coding

320    genes (7.08 RPKM) and similar to that observed for non-coding genes (1.7 RPKM)

321    (Supplemental file 1: Fig. S16A). Pseudogenes were expressed in a median of four tissues

322    (Supplemental file 1: Fig. S16B). The median number of expressed tissues for protein-coding

323    and non-coding genes was 44 tissues and five tissues, respectively (Supplemental file 1: Fig.

324    S16B). In addition, a total of 1,038 pseudogene-derived lncRNAs were expressed. The median

325    expression of pseudogene-derived lncRNAs was 1.8 RPKM, similar to that observed for other

326    lncRNAs (1.6 RPKM) (Supplemental file 1: Fig. S17A). In addition, pseudogene-derived lncRNAs

327    were expressed in a median of four different tissues, which was lower than observed for other

328    lncRNAs (seven tissues) (Supplemental file 1: Fig. S17B).

329    Testis had the highest number of expressed pseudogene-derived lncRNAs (427), followed by

330    fetal brain (315) (Supplemental file 1: Fig. S8A and B). The correlation between the number of

331    input reads and the number of pseudogene-derived lncRNAs was not significant (0.25, p-value

332    0.09).

333    **Gene expression diversity across tissues**

334    Tissue similarities increased dramatically from transcript level to gene level (Supplemental file

335    1: Fig. S4A, Fig. S5A, Fig. S18A, Fig. S19A). The median percentage of shared genes between

336    pairs of tissues was significantly higher in protein-coding genes compared to non-coding genes

16

337    (90% and 57%, respectively; p-value < 2.2e-16; Supplemental file 1: Fig. S18A, Fig. S19A).

338    Clustering of tissues based on protein-coding genes was similar to that observed based on

339    protein-coding transcripts (Supplemental file 1: Fig. S18B, Fig. S19B). The same result was

340    observed in non-coding genes and transcripts. In addition, clustering of tissues based on

341    protein-coding genes was different than that of non-coding genes (Supplemental file 1: Fig. S4B,

342    Fig. S5B, Fig. S18B, Fig. S19B, Fig. S35F).

343    Tissues with both fetal and adult samples (brain, kidney, muscle, and spleen) were used to

344    investigate gene biotype differences between these developmental stages. Similar to what was

345    observed at transcript level, fetal tissues were significantly enriched for non-coding genes and

346    pseudogenes and were depleted for protein-coding genes (p-value < 2.2e-16; Supplemental file

347    10). These results were consistent across all tissues with both adult and fetal samples

348    (Supplemental file 10).

349    **Gene validation**

350    A total of 32,460 genes (92% of predicted genes) were structurally validated by independent

351    datasets (PacBio Iso-seq data, ONT-seq data, *de novo* assembled transcripts from RNA-seq data)

352    and comparison with Ensembl and NCBI gene sets (see Method section). In addition, a total of

353    31,635 genes (90% of predicted genes) were expressed in multiple tissues (31,635 genes or

354    90%) (Fig. 7). All genes were extensively supported by data from different technologies such as

355    WTTS-seq, RAMPAGE, histone modification (H3K4me3, H3K4me1, H3K27ac) and CTCF-DNA

356    binding, and ATAC-seq data generated from the samples (Fig. 7).

**357**     **Identification and validation of annotated gene border extensions**

**358**     This new bovine gene set annotation extended (5' end extension, 3' end extension, or both)

**359**     more than 11,000 annotated Ensembl or NCBI gene borders. Extensions were longer on the 3'

**360**     side, but the median increase was 104 nt for the 5' end (Table 5). To validate gene border

**361**     extensions, independent WTTS-seq (24 tissues) and RAMPAGE datasets (30 tissues) were

**362**     utilized. More than 80% of annotated gene border extensions were validated by independent

**363**     data (Fig. 8). The extension of annotated gene borders on both ends resulted in an approximate

**364**     nine-fold expression increase of these genes in the new bovine gene set annotation compared

**365**     to their matched Ensembl and NCBI genes (Table 6). This effect was smaller in annotated genes

**366**     extended only on 5' or 3' ends (Table 6).

**367**     **Alternative splicing events**

**368**     Alternative splicing (AS) events (Supplemental file 1: Fig. S20A) are commonly distinguished in

**369**     terms of whether RNA transcripts differ by inclusion or exclusion of an exon, in which case the

**370**     exon involved is referred to as a "skipped exon" (SE) or "cassette exon", "alternative first exon",

**371**     or "alternative last exon". Alternatively, spliced transcripts may also differ in the usage of a 5'

**372**     splice site or 3' splice site, giving rise to alternative 5' splice site exons (A5Es) or alternative 3'

**373**     splice site exons (A3Es), respectively. A sixth type of alternative splicing is referred to as

**374**     "mutually exclusive exons" (MXEs), in which one of two exons is retained in RNA but not both.

**375**     However, these types are not necessarily mutually exclusive; for example, an exon can have

**376**     both an alternative 5' splice site and an alternative 3' splice site, or have an alternative 5' splice

**377**     site or 3' splice site, but be skipped in other transcripts. A seventh type of alternative splicing is

378     "intron retention", in which two transcripts differ by the presence of an unspliced intron in one

379     transcript that is absent in the other. An eighth type of alternative splicing is "unique splice site

380     exons" (USEs), in which two exons overlap with no shared splice junction. A total of 102,502

381     transcripts (85% of spliced transcripts) were involved in different types of AS events, a large

382     increase over Ensembl (63% of spliced transcripts) and NCBI (75% of spliced transcripts)

383     annotations (Additional file1: FigureS20B). Skipped exons were observed in a greater number of

384     transcripts compared to other types of AS events (Supplemental file 1: Fig. S21).

385     A median of 60% of tissue transcripts showed at least one type of AS event (Supplemental file

386     1: Fig. S22A). There was no significant correlation between the number of input reads and the

387     number of AS event transcripts across tissues (Supplemental file 1: Fig. S22B).

388     The median expression level of AS transcripts (111,366 transcripts or 65% of transcripts) was

389     1.38 RPKM, which was similar to that observed for other types of transcripts (1.58RPKM)

390     (Supplemental file 1: Fig. S23A). In addition, AS transcripts were expressed in a median of 10

391     tissues (Supplemental file 1: Fig. S23B), which was higher than for the other transcript types

392     (median of nine tissues). Alternatively spliced transcripts were enriched with protein-coding

393     transcripts (p-value < 2.2e-16). Meanwhile, transcripts that did not show AS events, i.e.,

394     unspliced transcripts and spliced transcripts from single transcript genes, were enriched for

395     non-coding transcripts (p-value < 2.2e-16). A median of 67% of protein-coding genes showed at

396     least one type of AS event. In contrast, this was only 3% in non-coding genes. In most cases, AS

397     events did not change transcript biotypes (Supplemental file 1: Fig. S24). In addition, a switch

398     from protein-coding to ncRNAs was the main biotype change resulting from AS events

399     (Supplemental file 1: Fig. S24).

400　A median of four AS events were expressed in alternatively spliced genes (14,260 genes or 40%

401　of genes) (Supplemental file 1: Fig. S25). The top five percent of genes with the highest number

402　of AS events (2,734 genes, Fig. 35A) were highly enriched for several BP GO terms related to

403　different aspects of RNA splicing (Supplemental file 1: Fig. S26B, Supplemental file 11).

404　Comparison of tissues with both fetal and adult samples (brain, kidney, Longissimus Dorsi (LD)

405　muscle, and spleen) revealed a significantly higher rate of AS events in fetal tissues (only genes

406　expressed in both fetal and adult samples were included in this analysis) (Supplemental file 1:

407　Fig. S27).

408　**Tissue specificity**

409　Nine percent of all genes (3,174) and transcripts (15,562) were only expressed in a single tissue

410　and were denoted as tissue-specific (Supplemental file 1: Fig. S28A). Most tissue-specific genes

411　(75%) and transcripts (84%) were un-annotated. Forty-nine percent of tissue-specific transcripts

412　(11,748) were produced by annotated genes. Most tissue-specific genes (61%) and transcripts

413　(57%) were protein-coding (Supplemental file 1: Fig. S28A and B). In addition, more than 70% of

414　tissue-specific transcripts (11,222) were transcribed from non-tissue-specific genes. Compared

415　to other tissues, testis and thymus had the highest number of tissue-specific genes and

416　transcripts (Supplemental file 1: Fig. S28C, Supplemental file 12). The expression level of tissue-

417　specific genes and transcripts was significantly lower than that of their non-tissue-specific

418　counterparts (p-value < 2.2e-16; Supplemental file 1: Fig. S28D). A median of 71% of tissue-

419　specific transcripts showed any type of AS event in their expressed tissues (Supplemental file 1:

420　Fig. S29). This was only 3.9% for tissue-specific genes (Supplemental file 1: Fig. S29). Testis,

421    myoblasts, mammary gland, and thymus had the highest proportion of tissue-specific genes

422    displaying any type of AS event (Supplemental file 1: Fig. S29).

423    A total of 16,806 multi-tissue expressed genes (53% of all multi-tissue expressed genes) and

424    74,487 multi-tissue expressed transcripts (51% of all multi-tissue expressed transcripts) showed

425    Tissue Specificity Index (TSI) scores (Supplemental file 13) greater than 0.9 and were expressed

426    in a tissue-specific manner. These genes and transcripts were expressed in a median of six

427    tissues and four tissues, respectively (Supplemental file 1: Fig. S30A and B). Functional

428    enrichment analysis of the top five percent of genes with the highest TSI score (3,171 genes)

429    resulted in the identification of "sexual reproduction" (p-value 3.06e-24) and "fertilization" (p-

430    value 1.04e-8) as their top enriched BP GO terms (Supplemental file 1: Fig. S30C-E,

431    Supplemental file 14).

432    **Tying genes to phenotypes**

433    There were 9,800 predicted genes identified as the closest expressed gene to an existing QTL

434    (QTL-associated genes) in their expressed tissues (Supplemental file 15). These genes had either

435    QTLs located inside (6,511 genes) or outside (5,306 genes) their genomic borders (either from

436    their 5' end or 3' end) with a median distance of 51.9 kilobases (KB) and a maximum distance of

437    2.6 million bases (MB) (Supplemental file 1: Fig. S31). Most QTL-associated genes were

438    annotated genes (8,130 genes or 83%). In addition, the median number of AS events in these

439    genes (eight) was significantly higher than that observed in other genes (median of seven AS

440    events; p-value 5.69e-09).

441 **Potential testis-pituitary axis**

442 Testis tissue was not clustered with any other tissues and had the highest number of tissue-

443 specific genes (1,195 genes) compared to the rest of the tissues (Supplemental file 1: Fig. S4,

444 Fig. S5, Fig. S18, and Fig. S19). Testis-specific genes were highly enriched with different traits

445 related to fertility (e.g., percentage of normal sperm and scrotal circumference), body weight

446 (e.g., body weight gain and carcass weight), and feed efficiency (e.g., residual feed intake)

447 (Supplemental file 16). The extent of testis-pituitary axis involvement in the "percentage of

448 normal sperm" was investigated using animals with both testis and pituitary samples (three

449 samples per tissue). The *SPACA5* gene was the only testis-specific gene encoded protein with a

450 signal peptide (SP) that was close to the "percentage of normal sperm" QTLs. The expression of

451 this gene in testis samples showed significant positive correlation with 70 pituitary expressed

452 genes that were closest to the "percentage of normal sperm" QTLs (Supplemental file 1: Fig.

453 S32, Supplemental file 17). These pituitary genes were enriched with the "signal transduction in

454 response to DNA damage" BP GO term (Supplemental file 1: Fig. S32). In addition, the

455 expression of testis genes that encoded protein with a signal peptide that were close to the

456 "percentage of normal sperm" QTLs was significantly correlated with expression of pituitary

457 genes close to this trait (Fig. 9, Supplemental file 18). The same result was observed for the

458 pituitary-testis tissue axis (Supplemental file 1: Fig. S33, Supplemental file 19).

459 **Trait similarity network**

460 The extent of genetic similarity between different bovine traits was investigated using their

461 associated QTLs. A total of 1,857 significantly similar trait pairs (184 different traits) were

462    identified and used to create a bovine trait similarity network

463    (https://www.animalgenome.org/host/reecylab/a; Supplemental file 20).

464    **miRNAs**

465    A total of 2,007 miRNAs (at least ten mapped reads in each tissue) comprised of 973 annotated

466    and 1,034 un-annotated miRNAs were expressed (Supplemental file 21). In each tissue, a

467    median of 704 annotated miRNAs and 549 un-annotated miRNAs were expressed (Fig. 10A).

468    The median expression of un-annotated miRNAs was 0.10 Reads Per Million (RPM), which was

469    significantly lower than that observed for annotated miRNAs (0.41 RPM; p-value 3.25e-25; Fig.

470    10B). In addition, un-annotated miRNAs were expressed in a median of 23 tissues, significantly

471    lower than for annotated miRNAs (43 tissues; p-value 1.00e-45; Fig. 10C). A median of 84.53%

472    of miRNAs were shared between pairs of tissues (Supplemental file 1: Fig. S34). Clustering of

473    tissues based on miRNAs was similar to what was observed based on non-coding genes

474    (Supplemental file 1: Fig. S35).

475    A total of 113 miRNAs (5.6%) were expressed in a single tissue and were denoted as tissue-

476    specific (Supplemental file 1: Fig. S36A). The proportion of tissue-specific miRNAs was higher for

477    un-annotated miRNAs, such that 75% of the tissue-specific miRNAs (85) were un-annotated.

478    The number of un-annotated miRNAs was higher in pre-adipocytes compared to other tissues,

479    followed by fetal gonad and testis (Supplemental file 1: Fig. S36B). Un-annotated miRNAs

480    showed a significantly lower expression level compared to annotated miRNAs (p-value 1.4e-19;

481    Supplemental file 1: FigureS36 C). In addition, a total of 1,047 multi-tissue expressed miRNAs

482    (55% of all multi-tissue expressed miRNAs) had a TSI score greater than 0.9 and were expressed

483    in a tissue-specific manner (Additional file1: Fig. S36D). These miRNAs were expressed in a

484    median of 19 tissues (Supplemental file 1: Fig. S36E).

485    Chromatin features across 500-base pair (bp) windows surrounding upstream of miRNA

486    precursors' start sites or downstream of miRNA precursors' terminal sites from independent

487    cattle experiments were used to investigate the relationship between miRNAs and chromatin

488    accessibility. More than 99% of un-annotated miRNAs (1,027) and 94% of annotated miRNAs

489    (923) were supported by at least one of the H3K4me3, H3K4me1, H3K27ac, CTCF-DNA binding,

490    or ATAC-seq peaks (Fig. 11).

491    **Summary of** expressed **transcripts, genes, and miRNAs**

492    The numbers of expressed transcripts, genes, and miRNAs in different tissues are summarized

493    in Supplemental file 1: Fig. S37. In addition, the number of annotated and un-annotated genes,

494    transcripts, and miRNAs in different tissues are summarized in Supplemental file 1: Fig. S38.

495    **Discussion**

496    Despite many improvements in the current bovine genome annotation ARS-UCD1.2 assembly

497    (Ensembl release 2021-03 and NCBI release 106) compared to the previous genome assembly

498    (UMD3.1), these annotations are still far from complete [9, 10]. In this study, using RNA-seq and

499    miRNA-seq data from 47 different bovine tissues/cell types, 12,698 un-annotated genes and

500    1,034 un-annotated miRNAs were identified that have not been reported in current bovine

501    genome annotations (Ensembl release 2021-03, NCBI release 106 and miRbase [11]). In

502    addition, we identified protein-coding transcripts with a median ORF length of 270 nt for 822

503     annotated bovine genes that have been annotated as non-coding in current bovine genome

504     annotations (Supplemental file 1: Fig. S14C). The high frequency of validation of these un-

505     annotated genes and un-annotated miRNAs using multiple independent datasets from different

506     technologies verifies the improvement in terms of the number of genes and miRNAs using our

507     methods.

508     Five prime and 3'untranslated region length plays a critical role in regulation of mRNA stability,

509     translation, and localization [4]. However, only a single 5' UTR and 3' UTR per gene is annotated

510     in current bovine genome annotations (Ensembl release 2021-03 and NCBI release 106), and

511     variations in UTR length are not available. In this study, 7,909 genes (22% of predicted genes)

512     with multiple UTRs were identified. Genes with multiple 5' UTRs are common, primarily due to

513     the presence of multiple promoters [12] or alternative splicing mechanisms within 5' UTRs [12].

514     Fifty-four percent of human genes have multiple transcription start sites [12]. In addition, the

515     length of 3' UTRs often varies within a given gene, due to the use of different poly(A) sites [4,

516     13].

517     In this study, around 50% of expressed protein-coding genes in each tissue transcribed both

518     coding and non-coding transcript isoforms. Several studies have shown evidence of the

519     existence of bifunctional genes with coding and non-coding potential using RNA-seq and

520     ribosome footprinting followed by sequencing (Ribo-seq) [14-16]. More than 20% of human

521     protein-coding genes have been reported to transcribe non-coding isoforms, often generated

522     by alternative splicing [17] and recurrently expressed across tissues and cell lines [16]. A

523     considerable number of non-coding isoform variants of protein-coding genes appear to be

524     sufficiently stable to have functional roles in cells [18]. It has been shown that the proportion of

525 non-coding isoforms from protein-coding genes dramatically increases during myogenic

526 differentiation of primary human satellite cells and decreases in myotonic dystrophy muscles

527 [19]. In this study, 106 non-coding genes were identified in fetal tissues that switched to

528 protein-coding genes in their matched adult tissues. Taken together this supports the notion

529 that protein-coding/non-coding transcript switching plays an important role in tissue

530 development in cattle as well.

531 Nonsense-mediated RNA decay is an evolutionarily conserved process involved in RNA quality

532 control and gene regulatory mechanisms [20]. For instance, the RNA-binding protein

533 polypyrimidine tract binding protein 1 (*PTBP1*) can promote the transcription of NMD

534 transcripts via alternative splicing, which negatively regulates its own expression [21]. In this

535 study, NMD transcripts comprised 19% of bovine transcripts that were transcribed from 30% of

536 bovine genes (10,498). In humans, NMD-mediated degradation can affect up to 25% of

537 transcripts [22] and 53% of genes [23]. As expected, in this study, most genes that transcribed

538 NMD transcripts were protein coding (83% or 8,687 genes), while a considerable portion (17%)

539 were pseudogenes. Many pseudogenes are annotated to give rise to NMD transcripts [24, 25].

540 Bioinformatic study of the human transcriptome revealed that 78% of NMD transcript–

541 producing genes were protein coding, followed by pseudogenes (nine percent), long intergenic

542 noncoding RNAs (six percent), and antisense transcripts (four percent) [25].

543 Despite the important regulatory function of lncRNAs and miRNAs, very low numbers of these

544 elements have been annotated in the current bovine genome annotations (Table 7). In this

545 study, a total of 10,789 lncRNA genes and 2,007 miRNA genes were expressed in the bovine

546 transcriptome, which is similar to what has been reported for the human transcriptome (Table

26

547    7). While, a total of 3,770 human miRNAs and 1,203 cattle miRNAs have been reported in

548    miRbase [11].

549    In this study, 1,038 pseudogene-derived lncRNAs were identified that were recurrently

550    expressed across tissues and cell types. Ever-increasing evidence from different studies

551    suggests pseudogene derived RNAs are key components of lncRNAs [26-28]. lncRNAs expressed

552    from pseudogenes have been shown to regulate genes with which they have sequence

553    homology [26, 27] or to coordinate development and disease in metazoan systems [26].

554    Correct annotation of gene borders has an important role in defining promoter and regulatory

555    regions. Our novel transcriptome analysis extended (5'-end extension, 3'-end extension, or

556    both) more than 11,000 annotated Ensembl or NCBI gene borders. Extensions were longer on

557    the 3' side, which was relatively similar to that we observed in the pig transcriptome using

558    PacBio Iso-Seq data [2].

559    A growing body of evidence indicates that a considerably large portion of lncRNAs encode

560    microproteins that are less conserved than canonical open reading frames [29-33]. In this study,

561    a vast majority (98%) of predicted lncRNAs had short ORFs (<44 amino acids) that were less

562    conserved than canonical ORFs (Table 2).

563    Alternative splicing is the key mechanism to increase the diversity of the mRNA expressed from

564    the genome and is therefore essential for response to diverse environments. In this study,

565    skipped exons and retained introns were the most prevalent AS events identified in the bovine

566    transcriptome, similar to what has been observed in other vertebrates and invertebrates [34]. A

567    higher rate of AS events was observed in fetal tissues compared to their adult tissue

568    counterparts. The same result has been observed in a recently published study in humans [35].

569    We hypothesized that the integration of the gene/transcript data with previously published

570    QTL/gene association data would allow for the identification of potential molecular

571    mechanisms responsible for a) tissue-tissue communication as well as b) genetic correlations

572    between traits. To test the first hypothesis, we developed a novel approach to study the

573    involvement of tissue-tissue interconnection in different traits based on the integration of the

574    transcriptome with publicly available QTL data. In particular, the interconnection between

575    testis and pituitary tissues with respect to the "percentage of normal sperm" trait was

576    investigated in more detail. This resulted in the identification of the regulation of ubiquitin-

577    dependent protein catabolic process, the regulation of Nuclear factor-κB (NF-κB) transcription

578    factor activity, and Rab protein signal transduction as key components of this tissue-tissue

579    interaction (Supplemental file 18 and 19). Interestingly, expressed genes that were closest to

580    "percentage of normal sperm" QTLs, and also encoded protein with a signal peptide (short

581    peptide present at the N-terminus of proteins that are destined toward the secretory

582    pathway[36])  in both testis and pituitary tissues, were highly enriched for the BP GO term

583    "regulation of ubiquitin-dependent protein catabolic process" (Supplemental file 18 and 19).

584    The expression of these genes in testis tissue was significantly correlated with expression levels

585    of pituitary expressed genes closest to "percentage of normal sperm" QTLs that were highly

586    enriched for the "positive regulation of NF-kappaB transcription factor activity" BP GO term

587    (Supplemental file 1: Fig. S32 and Supplemental file 18). Activation of NF-κB requires

588    ubiquitination, and this modification is highly conserved across different species [37]. NF-κB

28

589    induces secretion of adrenocorticotropic hormone from the pituitary [38], which directly

590    stimulates testosterone production by the testis [39]. In addition, ubiquitinated proteins in

591    testis cells are required for the progression of mature spermatozoa [40]. The expression levels

592    of pituitary expressed genes closest to "percentage of normal sperm" QTLs that also encoded

593    signal peptides were significantly correlated with expression levels of testis expressed genes

594    closest to "percentage of normal sperm" QTLs (Supplemental file 1: Fig. S33). These testis genes

595    were highly enriched for the "Rab protein signal transduction" BP GO term (Supplemental file

596    19). Rab proteins have been reported to be involved in male germ cell development [41]. Thus,

597    it appears that integration of gene data with QTL/association data can be used to identify

598    putative molecular pathways underlying tissue-tissue communication mechanisms.

599    To test the second hypothesis, we also developed a novel approach to study trait similarities

600    based on the integration of the transcriptome with publicly available QTL data. Using this

601    approach, we could identify significant similarity between 184 different bovine traits. For

602    example, clinical mastitis showed significant similarity with 23 different cattle traits that were

603    greatly supported by published studies, such as milk yield [42], milk composition traits [43],

604    somatic cell score [44], foot traits [45], udder traits [46], daughter pregnancy rate [47], length

605    of productive life [48] and net merit [49]. Similar results were observed for residual feed intake,

606    which showed significant similarity with 14 different traits such as average daily feed intake

607    [50], average daily gain [51], carcass weight [52], feed conversion ratio [53], metabolic body

608    weight [54], subcutaneous fat [55], and dry matter intake [56].

609    Taken together, these results identify a list of candidate genes that might harbor genetic

610    variation responsible for the genetic mechanisms underlying genetic correlations

611 (Supplemental file 18 and 1. If this is the case, in the future, these novel methods should be

612 able to predict the impact of a given set of genetic variants that are associated with a trait of

613 interest on other traits that were not measured in a given study. This might then lead to the

614 optimization of variants used (or not used) in genomic selection to minimize any non-beneficial

615 effect of selection on selected traits. However, it is important to acknowledge that the nearest

616 neighbor gene to a genotype association is not necessarily the causal gene.  None the less,

617 these results are intriguing in that meaningful genetic correlation can be recapitulated.

618 **Conclusions**

619 In-depth analysis of multi-omics data from 47 different bovine tissues/cell types provided

620 evidence to improve the annotation of thousands of protein-coding, lncRNA, and miRNA genes.

621 These validated results increase the complexity of the bovine transcriptome (number of

622 transcripts per gene, number of UTRs per gene, lncRNA transcripts, AS events, and miRNAs),

623 comparable to that reported for the highly annotated human genome. We provided direct

624 evidence that the predicted un-annotated transcripts extend existing annotated gene models,

625 by verifying such extensions using independent WTTS-seq and RAMPAGE data. We utilized a

626 novel approach to integrate the transcriptome with publicly available QTL data and showed its

627 application in a study of tissue axis involvement in different traits and genetic similarity

628 between different traits. This approach is particularly important in the selection of indicator

629 traits for breeding purposes, study of artificial selection side effects in livestock species, and

630 functional annotation of poorly annotated livestock genomes.

631 **Methods**

632 **Tissue and cell collection, total RNA extraction and construction of RNA-seq, miRNA-seq,**

633 **WTTS-seq, ATAC-seq, and ChIP-seq libraries**

634 **Cell sample collections.** Skeletal muscle and subcutaneous fat samples were collected from

635 Angus-crossbred steers slaughtered at the Virginia Tech Meat Center. Satellite cells were

636 isolated from skeletal muscle by pronase digestion as described previously (Leng et al. 2019).

637 The isolated satellite cells were activated to proliferate as myoblasts by culturing in growth

638 medium composed of Dulbecco's Modified Eagle Medium (DMEM), 10% fetal bovine serum

639 (FBS), and 1% antibiotics-antimycotics. To induce myoblasts to differentiate into myocytes,

640 myoblasts cultured in growth medium were switched to differentiation medium composed of

641 DMEM and 2% horse serum for 2 days. Preadipocytes from subcutaneous fat were isolated by

642 collagenase digestion as previously described (Hausman et al. 2008). To induce preadipocytes

643 to differentiate into adipocytes, preadipocytes were initially cultured in growth medium

644 (DMEM/F12, 10% FBS, 1% antibiotics-antimycotics) to reach confluency, then in induction

645 medium (DMEM/F12, 10% FBS, 1% antibiotics-antimycotics, 10 μg/mL insulin, 1 μM

646 dexamethasone, 0.5 mM isobutyl methylxanthine, and 200 μM indomethacin) for 2 days, and

647 lastly in maintenance medium (DMEM/F12, 10% FBS, 1% antibiotics-antimycotics, 1 μg/mL

648 insulin) for 10 days.

649 **Adult tissue collections.** Procedures for tissue collection followed the Animal Care and Use

650 protocol (#18464) approved by the Institutional Animal Care and Use Committee (IACUC),

651 University of California, Davis (UCD). Four cattle (2 males and 2 females) were slaughtered at

652 UCD using captive bolt under USDA inspection at 14 months old and were intact male and

653 female Line 1 Herefords that had the same sire, provided by Fort Keogh Livestock and Range

654 Research Lab [57]. Tissue samples were flash frozen in liquid nitrogen then stored at –80 °C

655 until further assay processing.

656 **Fetal tissue collections.** Fetal sample collection and tissue collection were approved by IACU),

657 University of Idaho (2017-67). Four pregnant females at day 78 of gestation Line 1 Herefords

658 were slaughtered at UI meats lab using captive bolt under USDA inspection. Animals were

659 provided by Fort Keogh Livestock and Range Research Lab (Tixier-Boichard et al. 2021). Tissue

660 samples were flash frozen in liquid nitrogen then stored at –80 °C until further assay processing.

661 **RNA-seq library construction.** Tissue samples (Supplemental file 22) were collected from

662 storage at -80 °C and ground to a powder using a mortar and pestle and liquid nitrogen. The

663 tissue was next homogenized in QIAzol Lysis Reagent (Qiagen Catalog No. 79306) using a

664 QIAshredder spin column (Qiagen Catalog No. 79656). After centrifugation, the lysate was

665 mixed with chloroform, shaken vigorously for 15 sec, incubated for 2 – 3 min at room

666 temperature, and centrifuged for 15 min at 12,000 x g at 4°C. The upper, aqueous phase was

667 transferred to a new collection tube and 1.5 vol of 100% ethanol was added and mixed

668 thoroughly by pipetting up and down several times. Total RNA was then isolated from the

669 sample using the RNeasy Mini Kit (Qiagen Catalog No. 74106) according to the manufacturer's

670    instructions. Contaminating DNA was removed by treating total RNA with DNase (AM1906,

671    Ambion). Total RNA quantity was measured with the Quant-It RiboGreen RNA Assay Kit (Life

672    Technologies Corp., Carlsbad, CA) and quality assessed by fragment analysis (Advance Analytical

673    Technologies, Inc., Ankeny IA).

674    **Mammary gland tissue collection and RNA-seq library construction.** The 16 animals used in

675    this study were Holstein-Friesian heifers from a single herd managed at the AgResearch

676    Research Station in Ruakura, NZ. All experimental protocols were approved by the AgResearch,

677    NZ, ethics committee, and carried out according to their guidelines. Samples were collected

678    from the same animals at 5 time points: virgin state before pregnancy between 13 and 15

679    months of age (virgin), mid-pregnant at day 100 of pregnancy, late pregnant ~2 weeks pre-

680    calving, early lactation ~2 weeks post-calving, and at peak lactation, 34-38 days post-calving.

681    Tissue samples were obtained by mammary biopsy using the Farr method [58]. Lactating cows

682    were milked before biopsy and sampled within 5 hours of milking. Biopsy sites were clipped and

683    given aseptic skin preparation (povidone iodine base scrub and iodine tincture) and

684    subcutaneous local anesthetic (4 ml per biopsy site). Core biopsies were taken using a powered

685    sampling cannula (4.5 mm internal diameter) inserted into a 2 cm incision. The resulting

686    samples of mammary gland parenchyma measured 70 mm in length, with a 4 mm diameter.

687    Small slices from each sample were preserved for histology before mammary epithelial

688    organoids were separated from surrounding adipose and connective tissue to allow for

689    secretory-specific signals in the RNA-seq analysis. In preparation for isolating organoids, tissue

690    samples were digested in a freshly prepared collagenase solution containing 0.2% collagenase A

691    (Roche), 0.05% trypsin (1:250 powder, 100U/ml Gibco), hyaluronidase (Sigma), 5% fetal calf

692  serum (Hyclone), Pen/Strep/Fungizone solution (Hyclone) or 5 µg/ml Gentamycin (Sigma) in

693  DMEM/F12 (Gibco) with 10 ng/ml insulin. Samples were minced to a fine slurry and incubated

694  in this freshly prepared collagenase solution (10 ml solution/g tissue) for 3.5 hours at 37°C in a

695  50 ml conical tube with slow shaking (120 rpm). Digested tissue was centrifuged at 453 x g for

696  10 min at 4°C, after which the supernatant and fat layers were discarded, and the pellet was

697  gently resuspended in 5 ml DMEM/F12 without serum. A further 5 ml DMEM/F12 without

698  serum was added, and the sample was centrifuged at 453 x g for 10 min at 4°C. The media was

699  discarded, and the pellet was gently resuspended in 10 ml DMEM/F12 and centrifuged for

700  another 10 min at 453 x g and 4°C. The media was discarded, and pellet resuspended in 10 ml

701  DMEM/F12 for a third time, and the sample centrifuged in a series of brief spins achieved by

702  allowing the centrifuge to reach 453 x g for two seconds before applying the brake. These brief

703  pulse spins were repeated at least 4 times, or until examining the sample under a microscope

704  revealed primarily epithelial organoid clusters and very few single cells. At this point, the

705  organoid pellet was resuspended in 1 ml TRIzol and stored at -80°C until RNA isolation. High-

706  quality total RNA (RIN > 7) was extracted from frozen mammary epithelial organoid pellets

707  using NucleoSpin® miRNA isolation kit (MACHEREY-NAGEL) according to the manufacturer's

708  protocol, isolating large and small (<200 bp) fractions separately. The "large" RNA fraction was

709  used to prepare strand-specific poly(A)+ RNA-seq libraries for sequencing. The "small" RNA

710  fraction was used to make miRNA-seq libraries using NEXTflex™ Small RNA-Seq Kit v3.

711 **miRNA-seq library construction.** Tissue samples (Supplemental file 22) were collected similarly

712 to the method described in the previous section. QIAseq miRNA Library Kit (Qiagen, cat no.

713 331505) and QIAseq miRNA NGS 96 Index IL Kit (Qiagen, cat no. 331565) were used to isolate

714 miRNAs from all tissues except mammary gland. miRNAs from mammary gland were isolated

715 using NEXTflex™ Small RNA-Seq Kit v3 (Illumina) according to the manufacturer's instructions.

716 The isolated miRNA was subjected to 3' ligation to ligate a pre-adenylated DNA adaptor to the

717 3' ends of all miRNAs. An RNA adaptor was then ligated to the 5' end of the mature miRNA to

718 complete 5' ligation. cDNA synthesis was completed using a reverse transcriptase (RT) primer

719 containing integrated unique molecular identifiers (UMI). The RT primer bound to the 3'

720 adaptor region and facilitated conversion of the 3'/5' ligated miRNAs into cDNA while a UMI

721 was assigned to every miRNA molecule. After reverse transcription, a clean-up of the cDNA was

722 performed using a streamlined magnetic bead-based method. Library amplification was

723 accomplished by a universal forward primer from a plate being paired with 1 of 96 dried

724 reverse primers in the same plate (Qiagen, cat no. 331565) to assign each sample a unique

725 custom index. Following library amplification, a clean-up of the miRNA library was performed

726 using a streamlined magnetic bead-based method. Libraries were then evaluated for quantity

727 and quality measures before being normalized and pooled for Ilumina sequencing (1×50bp).

728 **WTTS-seq library construction.** Construction of the WTTS-seq libraries from tissue samples

729 (Supplemental file 22) involved fragmentation, poly(A)+ RNA enrichment, first-strand cDNA

730 synthesis by reverse transcription and second-strand cDNA synthesis by PCR as described

731 previously [59]. The starting material was 2.5 µg of total RNA per library, which was fragmented

732 with 1 µl of 10X RNA fragmentation buffer (Ambion, AM8740), followed by enrichment of

733 poly(A)+ RNA using Dynabeads (Ambion 61002). The poly(A)+ RNA molecules were then used

734 for the first-strand cDNA synthesis with both 5' adaptor (switching primer, 100 µM) and 3'

735 adaptor (containing oligo (dT10), 100 µM) catalyzed by the SuperScript III reverse transcriptase

736 (200 U/µl) (Invitrogen, 18080). The first-strand cDNA molecules were chemically enriched with

737 RNases I and H and used to synthesize the second-strand cDNA using PCR. Base PCR conditions

738 were as follow: initial denaturation at 98 °C for 30 s, PCR cycles of 98 °C for 10 s, 50°C for 30 s,

739 and 72°C for 30 s, and final extension at 72°C for 10 min. The size-selected cDNA (200 – 500 bp)

740 was purified with SPRI beads (Agencourt AMPure XP beads, Beckman Coulter, Brea, CA) and

741 sequenced using an Ion PGM™ Sequencer at Washington State University.

742 **ATAC-seq library construction.** Frozen tissue samples (Supplemental file 22) were pulverized

743 under liquid nitrogen using mortar and pestle. Permeabilized nuclei were obtained by

744 resuspending pulverized tissue (5-15 mg) in 250 µL Nuclear Permeabilization Buffer (0.2%

745 IGEPAL-CA630 [I8896, Sigma], 1 mM DTT [D9779, Sigma], Protease inhibitor [05056489001,

746 Roche], and 5% BSA [A7906, Sigma] in PBS [10010-23, Thermo Fisher Scientific]), and incubating

747 for 10 min on a rotator at 4°C. Nuclei were then pelleted by centrifugation for 5 min at 500 x g

748 at 4°C. The pellet was resuspended in 25 µL ice-cold Tagmentation Buffer (33 mM Tris-acetate

749 [pH = 7.8; BP-152, Thermo Fisher Scientific], 66 mM K-acetate [P5708, Sigma], 11 mM Mg-

750 acetate [M2545, Sigma], 16% DMF [DX1730, EMD Millipore] in molecular biology grade water

751 [46000-CM, Corning]). An aliquot was then taken and counted by hemocytometer to determine

752 nuclei concentration. Approximately 50,000 nuclei were resuspended in 20 µL ice-cold

753 Tagmentation Buffer and incubated with 1 µL Tagmentation enzyme (FC-121-1030, Illumina) at

754 37 °C for 30 min with shaking at 500 rpm. The tagmentated DNA was purified using MinElute

755   PCR purification kit (28004, Qiagen). The libraries were amplified using NEBNext High-Fidelity

756   2X PCR Master Mix (M0541, NEB) with primer extension at 72°C for 5 min, denaturation at 98°C

757   for 30 s, followed by 8 cycles of denaturation at 98°C for 10 s, annealing at 63°C for 30 s and

758   extension at 72°C for 60 s. Amplified libraries were then purified using MinElute PCR

759   purification kit (28004, Qiagen), and two size selection steps were performed using SPRIselect

760   bead (B23317, Beckman Coulter) at 0.55X and 1.5X bead-to-sample volume ratios, respectively.

761   ATAC-seq libraries were sequenced on an Illumina Nextseq 500 platform using Nextra V2

762   sequencing chemistry to generate 2 × 75 paired-end reads.

763   **Sequencing the transcriptomes of seven bovine tissues by using the PacBio Iso-Seq and**

764   **Illumina RNA-Seq technologies**

765   Publicly available PacBio Iso-seq reads and matched RNA-seq reads (PRJNA386670) were used

766   in this experiment. Sequence reads were generated using the following procedure. Frozen

767   tissue samples (Supplemental file 22) were pulverized by grinding with disposable mortar and

768   pestle in liquid nitrogen. RNA was extracted using TRIzol reagent as directed by the

769   manufacturer (Invitrogen) with integrity examined using a BioAnalyzer (Agilent). Only samples

770   with RIN values >8 were used for cDNA synthesis. Libraries for RNA-seq short-read sequencing

771   were prepared using the TruSeq RNA Kit following the "TruSeq RNA Sample Preparation v2

772   Guide" as recommended by the manufacturer (Illumina). RNA-seq libraries were sequenced on

773   a NextSeq500 instrument. IsoSeq libraries for long-read sequencing were prepared using the

774   SMRTbell Template Prep Kit 1.0. First strand cDNA synthesis was performed with approximately

775   1 μg of extracted RNA from each tissue using the Clontech SMARTer PCR cDNA Synthesis Kit

776   (Clontech) as directed by the manufacturer. cDNA was then converted to SMRTbell template

777 library following the "Iso-Seq using Clontech cDNA Synthesis and BluePippin Size Selection"

778 protocol as directed by the manufacturer (Pacific Biosciences). Three size fraction pools for

779 each tissue were prepared using the BluePippin instrument (Sage Science), representing insert

780 sizes of 1-2 kb, 2-3 kb, and 3-6 kb. The two smaller fractions were sequenced in three to five

781 SMRT cells on an RSII instrument (Pacific Biosciences), and the largest fraction sequenced in five

782 or six cells, using P6/C4 chemistry. The sequences were processed into HQ isoforms using SMRT

783 Analysis v6.0 for each tissue independently but with all size fractions within tissue included in

784 the analysis.

785 **RNA-seq data analysis and transcriptome assembly**

786 Single-end Illumina RNA-Seq reads (75 bp) from each tissue sample were trimmed to remove

787 the adaptor sequences and low-quality bases using Trim Galore (version 0.6.4) [60] with --

788 quality 20 and --length 20 option settings. The resulting reads were aligned against ARS-UCD1.2

789 bovine genome using STAR (version 020201) [61] with a cut-off of 95% identity and 90%

790 coverage. FeatureCounts (version 2.0.2) [62] was used to quantify genes reported in the NCBI

791 gene build (version 1.21) with -Q 255 -s 2 --ignoreDup --minOverlap 5 option settings. The

792 resulting gene counts were adjusted for library size and converted to Counts Per Million (CPM)

793 values using SVA R package (version 3.30.0) [63]. In each tissue, sample similarities were

794 checked using hierarchical clustering and regression analysis of gene expression values (log2

795 based CPM), and outlier samples were expressed and removed from downstream analysis.

796 Samples from each tissue were combined to get the most comprehensive set of data in each

797 tissue. To reduce the processing time due to huge sequencing depth, the trimmed reads were

798 in silico normalized using insilico_read_normalization.pl from Trinity package (version 2.6.6)

799    [64] with --JM 350G and --max_cov 50 option settings. Normalized RNA-seq reads were aligned

800    against ARS-UCD1.2 bovine genome using STAR (version 020201) [61] with a cut-off of 95%

801    identity and 90% coverage. The normalized reads were assembled using *de novo* Trinity

802    software (version 2.6.6) [64] combined with massively parallelized computing using

803    HPCgridRunner (v1.0.1) [65] and GNU parallel software [66]. The resulting transcript reads were

804    collapsed and grouped into putative gene models (clustering transcripts that had at least a one-

805    nucleotide overlap) by the pbtranscript-ToFU from SMRT Analysis software (v2.3.0) [67]  with

806    min-identity = 95%, min-coverage = 90% and max_fuzzy_junction = 15 nt, whereas the 5'-end

807    and 3'-end difference were not considered when collapsing the reads. Base coverage of the

808    resulting transcripts was calculated using mosdepth (version 0.2.5) [68]. Predicted transcripts

809    were required to have a minimum of three times base coverage in their assembled tissues. The

810    predicted acceptor and donor splice sites were required to be canonical and supported by

811    Illumina-seq reads that spanned the splice junction with 5-nt overhang. Spliced transcripts with

812    the exact same splice junctions as their reference transcripts but that contained retained

813    introns were removed from analysis, as they were likely pre-RNA sequences. Unspliced

814    transcripts with a stretch of at least 20 A's (allowing one mismatch) in a genomic window

815    covering 30 bp downstream of their putative terminal site were removed from analysis, as they

816    were likely genomic-DNA contaminations. To decrease the false positive rate, unspliced

817    transcripts that were only expressed in a single tissue were removed from downstream

818    analysis. In addition, single-exon genes without histone mark (H3K4me3, H3K4me1, H3K27ac)

819    or ATAC-seq peaks mapped to their promoter (see Relating transcripts and genes to epigenetic

820    data section) were removed from downstream analysis as they were likely transcriptional noise.

39

821   The resulting transcripts from each tissue were re-grouped into gene models using an in-house

822   Python script. Structurally similar transcripts from the different tissues (see Comparison of

823   transcript structures across datasets/tissues section) were collapsed using an in-house Python

824   script to create the RNA-seq based bovine transcriptome.

825   The resulting transcripts and genes were quantified using align_and_estimate_abundance.pl

826   from the Trinity package (version 2.6.6) [64] with --aln_method bowtie --est_method RSEM --

827   SS_lib_type R option settings.

828   "Isoform" and "transcript" terms are used interchangeably throughout the manuscript.

829   **PacBio Iso-Seq data analysis**

830   PacBio Iso-seq data has been processed as described for the pig transcriptome [2] with the

831   following exceptions. Errors in the full-length, non-chimeric (FLNC) cDNA reads were corrected

832   with the preprocessed RNA-Seq reads from the same tissue samples using the combination of

833   proovread (v2.12) [69] and FMLRC (v1.0.0) [70] software packages. Error rates were computed

834   as the sum of the numbers of bases of insertions, deletions, and substitutions in the aligned

835   FLCN error-corrected reads divided by the length of aligned regions for each read (Table 8).

836   The RNA-seq-based transcriptome was assembled as described in the previous section.

837 **Oxford Nanopore data analysis**

838 Assembled isoforms from a previously published Oxford Nanopore experiment were used in

839 this study (Halstead et al. 2021).

840 **Comparison of transcript structures across datasets/tissues**

841 When comparing transcripts across datasets/tissues, transcripts whose 5' and 3' borders were

842 supported by RAMPAGE and/or WTTS data (see Transcript and gene border validation section)

843 and whose splice junctions were identical (maximum fuzzy junction was set to 15 bp) were

844 considered "structurally equivalent transcripts". The maximum of 100 nt fuzzy 5' and 3'

845 transcript borders were applied when comparing transcripts were not supported by RAMPAGE

846 and/or WTTS data. Other transcripts that did not met these criteria were considered

847 "structurally different transcripts".

848 A pair of genes was considered as structurally equivalent across datasets if they transcribed at

849 least single "structurally equivalent transcript".

850

851 **Prediction of transcript and gene biotypes**

852 Transcripts' open reading frames (ORFs) were predicted using the stand-alone version of

853 ORFfinder [71] with "ATG and alternative initiation codons" as ORF start codon. The longest

854 three ORFs were matched to the NCBI non-redundant vertebrate database and Uniprot

855 vertebrate database using Blastp [71] with E-value cutoff of $10^{-6}$, min coverage 60%, and min

856 identity 95%. The ORFs with the lowest E-value to a protein were used as the representative, or

857     if no matches were found, the longest ORF was used. Putative transcripts that had

858     representative ORFs longer than 44 amino acids were labelled as protein-coding transcripts. If

859     the representative ORF had a stop codon that was more than 50 bp upstream of the final splice

860     junction, it was labelled as a nonsense-mediated decay transcript [72]. Transcripts with start

861     codon but no stop codon before their poly(A) site were labelled non-stop decay RNAs. Putative

862     non-coding transcripts (ORFs shorter than 44 amino acids and lack of coding potential predicted

863     by CPC2 [73]) with lengths less than 200 bp that did not overlap with annotated or un-

864     annotated miRNA precursors (see miRNA-seq data analysis section) were labelled as small non-

865     coding RNAs [72]. Putative non-coding transcripts with lengths greater than 200 bp were

866     labelled as long non-coding RNAs [72]. Long non-coding RNAs overlapping one or more coding

867     loci on the opposite strand were labelled as antisense lncRNAs. Long non-coding RNAs located

868     in introns of coding genes on the same strand were labelled as sense-intronic lncRNAs. Long

869     non-coding RNAs that had an exon(s) that overlapped with a protein-coding gene were labeled

870     as Intragenic lncRNAs. Long non-coding RNAs located in intergenic regions of the genome were

871     labeled as Intergenic lncRNAs.

872     Putative genes that transcribed at least a single protein-coding transcript were labelled as

873     protein-coding genes. Putative genes with homology to existing vertebrate protein-coding

874     genes (Blastx [71], E-value cut-off $10^{-6}$, min coverage 90%, and min identity 95%) but containing

875     a disrupted coding sequence, i.e., transcribe only nonsense-mediated decay or non-stop decay

876     transcripts in all of their expressed tissues, were labelled as pseudogenes. The rest of the

877     putative genes were labeled as non-coding.

878 **ncRNAs homology analysis**

879 Putative non-coding transcripts were matched to NCBI and Ensembl vertebrate ncRNA

880 databases using Blastn [71] with E-value cutoff of $10^{-6}$, min coverage 90%, and min identity

881 95%. Transcripts with at least one hit were considered as homologous ncRNAs.

882 **Transcriptome termini site sequencing data analysis**

883 T-rich stretches located at the 5′ end of each WTTS-seq raw read were removed using an in-

884 house Perl script, as described previously [59]. T-trimmed reads were error-corrected using

885 Coral (version 1.4.1) [74] with -v -Y -u -a 3 option settings. The resulting reads were quality

886 trimmed using FASTX Toolkit (version 0.0.14) [75] with -q 20 and -p 50 option settings. High-

887 quality, error-corrected WTTS-seq reads were aligned against the ARS-UCD1.2 bovine genome

888 using STAR (version 020201) [61] with a cut-of of 95% identity and 90% coverage.

889 **ChIP-seq data analysis**

890 Regions of signal enrichment ("peaks") from a previously published ChIP-seq experiment were

891 used in this study [76].

892 **ATAC-seq data analysis**

893 The UC Davis FAANG Functional Annotation Pipeline was applied to process the ATAC-seq data,

894 as previously described [76]. Briefly, the ARS-UCD1.2 genome assembly and Ensembl genome

895 annotation (v100) were used as references for cattle. Sequencing reads were trimmed with

896 Trim Galore! (Krueger et al. 2015) (v.0.6.5) and aligned with either STAR (Dobin et al. 2012)

897 (v.2.5.4a) or BWA (Li et al. 2013) (v0.7.17) to the respective genome assemblies. Alignments

898    with MAPQ scores <30 were filtered using Samtools (Li et la. 2009) (v.1.9). Duplicate reads were

899    marked and removed using Picard (v.2.18.7). Regions of signal enrichment were called by

900    MACS2 (Zhang et al. 2008) (v.2.1.1).

901    **Relating transcripts and genes to epigenetic data**

902    The promoter was defined as the genomic region that spans from 500 bp 5' to 100 bp 3' of the

903    gene/transcript start site. Histone mark (H3K4me3, H3K4me1, H3K27ac), CTCF-DNA binding or

904    ATAC-seq peaks mapped to the promoter of a given gene/transcript were related to that

905    gene/transcript.

906    **Transcript and gene border validation**

907    RAMPAGE peaks from a previously published experiment [10] were used to validate

908    gene/transcript start site. Peaks within the genomic region that spans from 30 bp 5' to 10 bp 3'

909    of a gene/transcript start site were assigned to that gene/transcript. WTTS-seq reads (median

910    length of 161 bp) within the genomic region that spans from 10 bp 5' to 165 bp 3' of a

911    gene/transcript terminal site were assigned to that gene/transcript.

912    **Functional enrichment analysis**

913    The potential mechanism of action of a group of genes was deciphered using ClueGO [77]. The

914    latest update (May 2021) of the Gene Ontology Annotation database (GOA)  [78] was used in

915    the analysis. The list of genes with at least one transcript expressed in a given tissue was used

916    as background for that tissue. The GO tree interval ranged from 3 to 20, with the minimum

917    number of genes per cluster set to three. Term enrichment was tested with a right-sided hyper-

918 geometric test that was corrected for multiple testing using the Benjamini-Hochberg procedure

919 [79]. The adjusted p-value threshold of 0.05 was used to filter enriched GO terms.

920 **Alternative splicing analysis**

921 Alternative splicing events, except Unique Splice Site Exons, were detected using

922 generateEvents from SUPPA (version 2.3) [80] with default settings. Unique Splice Site Exons

923 were detected using an in-house Python script.

924 **miRNA-seq data analysis**

925 Single-end Qiagen miRNA-seq reads (50 bp) from each tissue sample were trimmed to remove

926 the adaptor sequences and low-quality bases using Trim Galore (version 0.6.4) [60] with --

927 quality 20, --length 16, --max_length 30 -a AACTGTAGGCACCATCAAT option settings. miRNA

928 reads were aligned against the ARS-UCD1.2 bovine genome using mapper.pl from mirDeep2

929 (version 0.1.3) [81] with -e -h -q -j -l 16 -o 40 -r 1 -m -v -n option settings. miRNA mature

930 sequences along with their hairpin sequences for Bos taurus species were downloaded from

931 miRbase [11]. These sequences, along with the aligned miRNA reads, were used to quantify

932 annotated miRNAs in each sample using miRDeep2.pl from mirDeep2 (version 0.1.3) [81] with -t

933 bta -c -v 2 setting options. miRNA normalized Reads Per Million (RPM) were used to check

934 sample similarities using hierarchical clustering and regression analysis of gene expression

935 values (log2 based CPM), and outlier samples were detected and removed from downstream

936 analysis. In order to predict the most comprehensive set of un-annotated miRNAs, samples

937 from different tissues were concatenated into a single file that were aligned against the ARS-

938 UCD1.2 bovine genome using mapper.pl from mirDeep2 (version 0.1.3) [81] with the

939    aforementioned settings. Aligned reads from the previous step were used, along with

940    annotated miRNAs' mature sequences and their hairpins, to predict un-annotated miRNAs

941    using miRDeep2.pl from mirDeep2 (version 0.1.3) [81] with the aforementioned settings.

942    Samples from each tissue were combined to get the most comprehensive set of data for that

943    tissue. Mature miRNA sequences and their hairpins for both annotated and predicted un-

944    annotated miRNAs' sequences along with the aligned miRNA reads from each tissue were used

945    to quantify annotated and un-annotated miRNAs in each tissue using mirDeep2 (version 0.1.3)

946    [81] with the aforementioned settings.

947    **Tissue-specificity index**

948    Tissue Specificity Index (TSI) calculations were utilized to present more comprehensive

949    information on transcript/gene/miRNA expression patterns across tissues. This index has a

950    range of zero to one with a score of zero corresponding to ubiquitously expressed

951    transcripts/genes/miRNAs (i.e., "housekeepers") and a score of one for

952    transcripts/genes/miRNAs that are expressed in a single tissue (i.e., "tissue-specific") [82]. The

953    TSI for a transcript/gene/miRNA j was calculated as [82]:

954

955
$$TSI_j = \frac{\sum_{i=1}^{N}(1 - x_{j,i})}{N - 1}$$

956

957　　where $N$ corresponds to the total number of tissues measured, and $x_{j,i}$ is the expression

958　　intensity of tissue $i$ normalized by the maximal expression of any tissue for

959　　transcript/gene/miRNA $j$.

960　　**QTL enrichment analysis**

961　　Publicly available bovine QTLs were retrieved from Animal QTLdb [83]. Closest expressed gene

962　　to a given trait's QTLs were denoted as QTL-associated genes for that trait. The median distance

963　　of QTLs located outside gene borders to the closest expressed gene was 51.9 kilobases and the

964　　maximum distance was 2.6 million bases. QTL enrichment was tested with a right-sided Fisher

965　　Exact test using an in-house Python script. The resulting p-values were corrected for multiple

966　　testing by the Benjamini-Hochberg procedure [79]. The adjusted p-value threshold of 0.05 was

967　　used to filter QTLs.

968　　**Trait similarity network**

969　　For a given pair of traits, trait A was denoted as "similar" to trait B if a significant portion of trait

970　　A's QTL-associated genes were also the closest expressed genes to trait B QTLs based on 1000

971　　permutation tests. The resulting p-values were corrected for multiple testing using the

972　　Benjamini-Hochberg procedure [79]. The same procedure was used to test trait B's similarity to

973　　trait A. The adjusted p-value threshold of 0.05 was used to filter significant trait similarities. A

974　　graphical presentation of the method used to construct the tissue similarity network is

975　　presented in Supplemental file 1: Fig. S39. The resulting network was visualized using

976　　Cystoscape software [84].

977

978 **Testis-pituitary axis correlation significance test**

979 The presence of signal peptides on representative ORFs of protein-coding transcripts was

980 predicted using SignalP-5.0 [85]. Spearman correlation coefficients were used to study

981 expression similarity between testis genes encoding signal peptides that were closest to the

982 "percentage of normal sperm" QTLs (62 genes) and pituitary expressed genes closest to the

983 "percentage of normal sperm" QTLs (246 genes). To test the statistical difference between

984 these correlation coefficients (reference correlations) and random chance, 1000 random sets of

985 246 pituitary genes were selected, and their correlation coefficients with 62 previously

986 described testis genes were calculated (random correlations). The reference correlations were

987 compared with 1000 sets of random correlations using a right-sided t-test. The resulting p-

988 values were corrected for multiple testing by the Benjamini-Hochberg procedure [79]. The

989 distribution-adjusted p-values were used to determine the significance level of expression

990 similarities for genes involved in the testis-pituitary axis related to "percentage of normal

991 sperm". The same analysis was conducted to determine the significance of pituitary-testis axis

992 involvement in this trait.

993 **Tissue dendrogram comparison across different transcript and gene biotypes**

994 Tissues were clustered based on the percentage of their transcripts/genes that were shared

995 between tissue pairs using the hclust function in R. Cophenetic distances for tissue

996 dendrograms were calculated using the cophenetic R function. The degree of similarity

997 between dendrograms constructed based on different gene/transcript biotypes was obtained

998 using the Spearman correlation coefficient between the dendrograms' Cophenetic distances.

999 **Figure legends**

1000 **Figure 1.** Distribution of the number of expressed transcripts (A) and genes (B) across tissues.

1001 **Figure 2.** Classification of the predicted transcripts into different biotypes.

1002 **Figure 3.** Support of predicted transcripts using data from different technologies and datasets.

1003 **Figure 4.** Classification of the predicted genes into different biotypes.

1004 **Figure 5.** Distribution of the number of 5' UTRs and 3' UTRs per gene in genes with multiple

1005 UTRs.

1006 **Figure 6.** (A) Classification of protein-coding genes based on their novelty and types of encoded

1007 transcripts. (B) Number of expressed tissues for bifunctional genes. Dots have been color coded

1008 based on their density. (C) Location of different transcript biotypes on bifunctional genes. (D)

1009 Functional enrichment analysis of genes that remained bifunctional in all of their expressed

1010 tissues.

1011 **Figure 7.** Support of predicted genes using data from different technologies and datasets

1012 **Figure 8.** Functional enrichment analysis of non-coding genes in fetal tissues that were switched

1013 to protein coding with only coding transcripts in their matched adult tissue.

1014 **Figure 9**- (A) Correlation between testis genes encoded protein with a signal peptide that were

1015 close to the "percentage of normal sperm" QTL and pituitary expressed genes closest to this

1016 trait (reference correlations). (B) Distribution of p-values resulting from a right-sided t-test

49

1017   between reference correlation coefficients and correlation coefficients derived from random

1018   chance (see methods for details).

1019   **Figure 10-** (A) Distribution of the number of expressed annotated and un-annotated miRNAs

1020   across tissues. (B) Expression of annotated and un-annotated miRNAs across their expressed

1021   tissues. (C) Number of expressed tissues for annotated and un-annotated miRNAs.

1022   **Figure 11-** Support of annotated (A) and un-annotated (B) miRNAs using different histone marks

1023   and CTCF-DNA binding data.

1024

1025 **Tables**

**Table 1.** Summary of expressed transcripts/genes

| Feature | Annotation[1] | | |
|---|---|---|---|
| | Current project | Ensembl | NCBI |
| | | (Release 2021-03) | (Release 106) |
| Number of genes | 35,150 (21,193) | 27,607 (21,880) | 35,143 (21,355) |
| Number of transcripts | 171,985 (85,658) | 43,984 (37,538) | 83,195 (47,280) |
| Number of spliced transcripts | 130,531 | 37,299 | 73,423 |
| Number of transcripts per gene | 4.9 | 1.5 | 2.3 |
| Median number of 5' UTRs per gene | 2 | 1 | 1 |
| Median number of 3' UTRs per gene | 1 | 1 | 1 |

[1]Numbers in parentheses indicate the number of protein-coding genes/transcripts.

1026

1027

1028

**Table 2.** Protein/peptide homology of transcripts with coding potential

| Transcript biotype | Number of transcripts | Transcripts with protein/peptide homology to other species[1] |
|---|---|---|
| Protein-coding transcripts | 85,658 | 73,268 (86%) |
| sncRNAs and lncRNAs that encode short peptides[2] | 48,425 | 4,054 (8%) |

[1]Number in parentheses indicates the percentage of each transcript biotype.

[2]Open reading frame of 9 to 43 amino acids

1029

1030

1031

**Table 3.** Sequence homology of non-coding transcripts

| Transcript biotype | Number of transcripts | Transcripts with sequence homology to ncRNAs in other species[1] |
|---|---|---|
| Long non-coding RNAs | 48,661 | 23,707 (49%) |
| Small non-coding RNAs | 526 | 194 (37%) |
| Non-stop decay RNAs | 4,359 | 1,551 (35%) |
| Nonsense-mediated decay RNAs | 32,781 | 18,195 (55%) |

[1]Number in parentheses indicates the percentage of each transcript biotype.

1032

1033

1034

**Table 4.** Sequence homology of different types of lncRNAs

| lncRNA biotype | Number of transcripts | Transcripts with sequence homology to ncRNAs in other species[1] |
|---|---|---|
| antisense lncRNAs | 29,987 | 13,793 (46%) |
| sense-intronic lncRNAs | 1,694 | 1,029 (60%) |
| intragenic lncRNAs | 5,569 | 2,314 (41%) |
| intergenic lncRNAs | 11,841 | 5,820 (49%) |

[1]Number in parentheses indicates the percentage of each transcript biotype.

1035

1036

1037

**Table 5.** Gene border extensions in current ARS-UCD1.2 genome annotations by *de novo*

assembled transcriptome from short-read RNA-seq data

| Annotation | Type of gene extension | Number of genes | Median extension (nucleotides) |
|---|---|---|---|
| Ensembl | 5' extension only | 1,848 | 128 |
| (Release 2021-03) | 3' extension only | 5,701 | 422 |
| | Both ends extended | 4,874 | 122, 5' |
| | | | 439, 3' |
| NCBI | 5' extension only | 2,214 | 80 |
| (Release 106) | 3' extension only | 5,496 | 126 |
| | Both ends extended | 3,613 | 66, 5' |
| | | | 210, 3' |

1038

1039

1040

1041

1042

**Table 6.** Median number of reads mapped to the extended region of annotated genes[1]

| Annotation | 5' end extension | 3' end extension | Both ends extension |
|---|---|---|---|
| Ensembl (release 2021-03) | 92 (1.10) | 220 (1.24) | 1,766 (8.90) |
| NCBI (release 106) | 72 (1.05) | 95 (1.10) | 2,009 (9.05) |

[1]Numbers in parentheses indicate the median fold change in expression level resulting from gene extensions.

1043

1044

1045

**Table 7.** Comparison of different gene builds based on gene biotypes

| Species | Gene build | Protein-coding genes | lncRNA genes | miRNA genes | Other types of small non-coding genes[1] | Pseudo-genes |
|---|---|---|---|---|---|---|
| Bovine (ARS-UCD1.2) | Ensembl (Release 2021-03) | 21,880 | 1,480 | 951 | 2,209 | 492 |
| | NCBI (Release 106) | 21,039 | 5,179 | 797 | 3,249 | 4,569 |
| | Current project[2] | 21,193 (18,096) | 10,789 (2,847) | 2,007 (973) | 139 (0) | 3,029 (1,509) |
| Human (GRCh38.104) | Ensembl (release 2021-03) | 20,442 | 16,876 | 1,877 | 2,930 | 15,266 |

[1]Small nucleolar RNAs, small non-coding RNAs, small Cajal body specific RNAs, small conditional RNAs, and tRNAs

[2]Numbers in parentheses indicate the number of un-annotated RNAs in each biotype.

1046

1047

Table 8. Summary of error-corrected, FLNC Iso-Seq reads and their matched RNA-seq

reads

| Tissue | Error-corrected FLNC Iso-Seq reads[1] | Median error rate in error-corrected FLNC Iso-Seq reads | Normalized RNA-seq reads used for error correction[2] |
|---|---|---|---|
| Thalamus | 664,900 (90%) | 0.21% | 32,452,612 |
| Testes | 711,821 (86%) | 1.43% | 31,939,024 |
| Liver | 1,064,146 (84%) | 1.84% | 13,657,156 |
| Medulla | 380,531 (86%) | 0.43% | 48,256,918 |
| Subcutaneous fat | 215,759 (93%) | 0.45% | 42,043,313 |
| Cerebral cortex | 440,797 (87%) | 1.01% | 21,285,864 |
| Jejunum | 604,436 (90%) | 2.331% | 34,457,447 |

[1] Number in parentheses indicates mapping rate (90% coverage and 95% identity).

[2] In silico normalized using insilico_read_normalization.pl from Trinity (version 2.6.6) with the

following settings: --max_cov 50 --max_pct_stdev 100 --single

1048

1049

## Supplemental files

1050

**Supplemental file 1: Fig. S1** Distribution of the number of RNA-seq reads across tissues. **Fig. S2** 1051

1052 (A) Comparison of tissues based on number of transcript biotypes and (B) percentage of

1053 transcript biotypes. (C) Comparison of transcript biotypes based on their number of expressed

1054 tissues and (D) their expression level across expressed tissues. **Fig. S3** (A) Relation between the

1055 number of input reads and the number of transcript biotypes (B) Comparison of expression

1056 level between different transcript biotypes. **Fig. S4** Tissue similarities (A) and clustering (B)

1057 based on the percentage of protein-coding transcripts shared between pairs of tissues. **Fig. S5**

1058 Tissue similarities (A) and clustering (B) based on the percentage of non-coding transcripts

1059 shared between pairs of tissues. **Fig. S6** Comparison of annotated and un-annotated transcripts

1060 based on their expression (A) and number of expressed tissues (B). **Fig. S7** Comparison of

1061 annotated and un-annotated protein-coding transcripts based on the length (A) and GC content

1062 (B) of their 5' UTR, CDS, and 3' UTR. **Fig. S8** (A) Comparison of tissues based on number of gene

1063 biotypes and (B) percentage of gene biotypes. **Fig. S9** Relation between the number of input

1064 reads and the number of gene biotypes. **Fig. S10** Comparison of annotated and un-annotated

1065 genes based on their expression (A) and number of expressed tissues (B). **Fig. S11** Functional

1066 enrichment analysis of the top five percent of genes with the highest number of UTRs. **Fig. S12**

1067 Similarity of tissues based on the number of non-coding genes in their fetal samples that

1068 switched to protein-coding genes with only coding transcripts in their adult samples. **Fig. S13**

1069 (A) Distribution of genes that transcribed PATs, based on their number of expressed tissues,

1070 percentage of genes' transcripts that are PATs and percentage of genes' expressed tissues in

1071 which PATs were transcribed. (B) Comparison of genes that transcribed PATs with other gene

1072      biotypes. **Fig. S14** (A) Homology analysis of protein-coding genes. (B) Homology analysis of non-

1073      coding genes. (C) Detection of orphan genes based on homology classification of cattle-specific

1074      protein-coding genes and non-coding genes. **Fig. S15** Comparison of the expression level of

1075      homologous and orphan genes across (A) and within (B) their expressed tissues. (C)

1076      Comparison of homologous and orphan genes based on the number of expressed tissues. **Fig.**

1077      **S16** Comparison of different gene biotypes based on the expression (A) and the number of

1078      expressed tissues (B). **Fig. S17** Comparison of different pseudogene-derived lncRNAs and non-

1079      pseudogene derived lncRNAs based on the expression level (A) and the number of expressed

1080      tissues (B)**. Fig. S18** Tissue similarities (A) and clustering (B) based on the percentage of protein-

1081      coding genes shared between pairs of tissues. **Fig. S19** Tissue similarities (A) and clustering (B)

1082      based on the percentage of non-coding genes shared between pairs of tissues. **Fig. S20** (A)

1083      Different types of alternative splicing events. (B) Comparison of bovine genome builds based on

1084      the number of transcripts that showed any type of alternative splicing (AS) events**. Fig. S21**

1085      Comparison of tissues based on the number (A) and the percentage (B) of transcripts that

1086      showed different types of alternative splicing events. Comparison of tissues based on the

1087      number (C) and the percentage (D) of alternative splicing events**. Fig. S22** (A) Comparison of

1088      tissues based on the percentage of transcripts that showed any type of alternative splicing

1089      events, spliced transcripts from single-transcript genes, and unspliced transcripts and (B) the

1090      relation between the number of input reads and the number of these transcripts across tissues.

1091      **Fig. S23** Comparison of transcripts that showed different types of alternative splicing events

1092      based on (A) the expression level in the expressed tissues and (B) the number of expressed

1093      tissues. **Fig. S24** Transcript biotype switching due to alternative splicing events**. Fig. S25**

1094     Comparison of tissues based on the number of alternative splicing events per alternatively

1095     spliced gene. **Fig. S26** (A) Distribution of the number of alternative splicing events per

1096     alternatively spliced gene. The 5% quantile is shown using a dashed red line. (B) Functional

1097     enrichment analysis of the top five percent of genes with the highest number of alternative

1098     splicing events. **Fig. S27** Comparison of the alternative splicing rate between adult and fetal

1099     tissues. **Fig. S28** (A) Distribution of gene's number of expressed tissues. Tissue-specific gene

1100     biotypes are shown in the pie chart. (B) Distribution of transcript's number of expressed tissues.

1101     Tissue-specific transcript biotypes are shown in the pie chart. (C) Comparison of tissues based

1102     on the number of tissue-specific genes and transcripts. (D) Comparison of the expression level

1103     of tissue-specific genes and transcripts versus their non-tissue-specific counterparts. **Fig. S29**

1104     Relationship between tissue specificity and alternative splicing events. **Fig. S30** Relationship

1105     between tissue specificity index and the number of multi-tissue expressed genes (A) and

1106     transcripts (B). Distribution of tissue specificity indexes in multi-tissue expressed genes (C) and

1107     transcripts (D). The 5% quantile is shown using dashed red lines. (E) Functional enrichment

1108     analysis of the top five percent of multi-tissue expressed genes with the highest tissue

1109     specificity indexes. **Fig. S31** Distribution of QTLs located outside gene borders in relation to the

1110     closest expressed gene. **Fig. S32** (A) Distribution of correlation coefficients between *SPACA5*

1111     gene expression and pituitary expressed genes closest to "percentage of normal sperm" QTLs.

1112     Dashed lines show the minimum significant positive and negative correlation (p-value <0.05).

1113     (B) Expression atlas of *SPACA5* gene in human tissues from The Human Protein Atlas [86]. **Fig.**

1114     **S33** (A) Correlation between pituitary genes with signal peptides that were close to the

1115     "percentage of normal sperm" QTL and testis expressed genes closest to this trait's QTL

1116 (reference correlations). (B) Distribution of p-values resulting from right-sided t-test between

1117 reference correlation coefficients and correlation coefficients derived from random chance (see

1118 methods for details)**. Fig. S34** Tissue similarities (A) and clustering (B) based on the percentage

1119 of miRNAs shared between pairs of tissues. **Fig. S35** Clustering of tissues based on protein-

1120 coding genes (A), protein-coding transcripts (B), non-coding genes (C), non-coding transcripts

1121 (D), and miRNAs (E). (F) Comparison of tissue dendrograms based on the correlation between

1122 their Cophenetic distances. **Fig. S36** (A) Distribution of the number of expressed tissues for

1123 annotated and un-annotated miRNAs. Classification of miRNAs as annotated, or un-annotated

1124 is presented in the pie chart. (B) Comparison of tissues based on their number of tissue-specific

1125 miRNAs. (C) Expression of annotated and un-annotated miRNAs in their expressed tissues. (D)

1126 Distribution of multi-tissue expressed miRNAs' tissue specificity indexes. (E) Relationship

1127 between tissue specificity index and number of expressed tissues in multi-tissue expressed

1128 miRNAs. Dots have been color coded based on their density. **Fig. S37** Distribution of the

1129 number of expressed genes (A), transcripts (B), and miRNAs (C) across tissues. **Fig. S38**

1130 Distribution of the number of annotated and un-annotated genes (A), transcripts (B), and

1131 miRNAs (C) across tissues. **Fig. S39** Graphical representation of the method used to construct

1132 the tissue similarity network.

1133 **Supplemental file 2:** Summary of RNA-seq and miRNA-seq reads

1134 **Supplemental file 3:** Detailed description of the number of transcripts, genes, and miRNAs

1135 expressed in each tissue

1136    **Supplemental file 4:** List of transcripts and genes expressed in each tissue and their expression

1137    values (RPKM)

1138

1139    **Supplemental file 5:** Transcript biotype enrichment analysis in adult and fetal tissues

1140    **Supplemental file 6:** Functional enrichment analysis of the top five percent of genes with the

1141    highest number of UTRs

1142    **Additional file7:** Functional enrichment analysis of genes that remained bifunctional in all their

1143    expressed tissues

1144    **Additional file8:** Functional enrichment analysis of non-coding genes in fetal tissues that were

1145    switched to protein coding with only coding transcripts in their matched adult tissue

1146    **Additional file9:** Functional enrichment analysis of protein-coding genes that transcribed PATs

1147    as their main transcripts (PATs comprised >50% of their transcripts) in all their expressed

1148    tissues

1149    **Supplemental file 10:** Gene biotype enrichment analysis in adult and fetal tissues

1150    **Supplemental file 11:** Functional enrichment analysis of the top five percent of genes with the

1151    highest number of alternative splicing events

1152    **Additional file12:** List of tissue-specific genes and transcripts

1153    **Additional file13:** Genes and transcripts tissue specificity indexes

1154    **Additional file14:** Functional enrichment analysis of the top five percent of multi-tissue

1155    expressed genes with the highest tissue specificity indexes

1156    **Additional file15:** List of QTL's closest expressed genes in each tissue

1157    **Additional file16:** Trait enrichment analysis of testis-specific genes

1158    **Additional file16:** Pituitary expressed genes closest to "percentage of normal sperm" QTLs that

1159    showed positive significant correlation with SPACA5 gene in testis

1160    **Additional file18:** List of expressed genes closest to "percentage of normal sperm" QTLs that

1161    were involved in testis-pituitary tissue axis    and their functional enrichment analysis results

1162    **Additional file19:** List of genes expressed closest to "percentage of normal sperm" QTLs that

1163    were involved in pituitary-testis tissue axis and their functional enrichment analysis results

1164    **Additional file20:** Similarity of traits based on the integration of the assembled bovine

1165    transcriptome with publicly available QTLs

1166    **Additional file21:** List of miRNAs expressed in each tissue and their expression values

1167    **Additional file22:** List of tissues related to different omics datasets used in the experiment

1168    **Abbreviations**

1169    A3Es: Alternative 3' splice site Exons; A5Es: Alternative 5' splice site Exons; AFEs: Alternative

1170    First Exon; ALEs: Alternative Last Exon; AS: Alternative Splicing; ATAC-seq: Assay for

1171    Transposase-Accessible Chromatin using sequencing; bp: base pair; BP: Biological Process; CDS:

1172    coding sequence; ChIP-seq: Chromatin Immunoprecipitation Sequencing; CPM: Counts Per

1173    Million; CTCF: CCCTC-binding factor; DMEM: Dulbecco's Modified Eagle Medium; FLNC: Full-

1174    Length, Non-Chimeric; GO:  Gene Ontology; GOA: Gene Ontology Annotation database; GWAS:

1175    Genome-Wide Association Studies; H3K27ac: N-terminal acetylation of lysine 27 on histone H3;

1176    H3K4me1: tri-methylation of lysine 4 on histone H1; H3K4me3: tri-methylation of lysine 4 on

1177    histone H3; IACUC: Institutional Animal Care and Use Committee; LD:  Longissimus Dorsi;

1178    lncRNAs: long non-coding RNAs; miRNA: microRNAs; MXEs: Mutually Exclusive Exons; NCBI:

1179    National Center for Biotechnology Information; ncRNAs: non-coding RNAs; NMD: Nonsense-

1180    Mediated Decay; NSD: Non-Stop Decay; ONT-seq: Oxford Nanopore Technologies sequencing;

1181    ORFs:  Open Reading Frames; PacBio Iso-Seq: Pacific Biosciences single-molecule long-read

1182    isoform sequencing; PAT: Potentially Aberrant Transcript; poly(A): Polyadenylation; PTBP1:

1183    polypyrimidine tract binding protein 1; QTL: Quantitative Trait Loci; RAMPAGE: RNA Annotation

1184    and Mapping of Promoters for the Analysis of Gene Expression; Ribo-seq: Ribosome

1185    footprinting followed by Sequencing; RIEs: Retained Intron Exons; RNA-seq: Illumina high-

1186    throughput RNA sequencing; RPKM: Reads Per Kilobase of Transcript per Million reads mapped;

1187    RPM: Reads Per Million; SEs: Skipped Exons; sncRNAs: small non-coding RNAs; SNP: Single

1188    Nucleotide Polymorphism; tpg: transcripts per annotated gene; TSI: Tissue Specificity Index;

1189    TSS: Transcript Start Sites; TTS: Transcript Terminal Sites; UCD: University of California, Davis;

1190    USEs: Unique Splice Site Exons; UTR: untranslated region; WTTS-seq: Whole Transcriptome

1191    Termini Site Sequencing.

1192    **Data availability**

1193    RNA-seq and miRNA-seq, ATAC-seq, and WTTS-seq datasets generated in this study are

1194    submitted to the ArrayExpress database (https://www.ebi.ac.uk/biostudies/arrayexpress)

1195    under accession numbers E-MTAB-11699, E-MTAB-11815, and E-MTAB-12052, respectively. The

1196    constructed bovine trait similarity network is publicly available through the Animal Genome

1197    database (https://www.animalgenome.org/host/reecylab/a). The constructed cattle

1198    transcriptome and related sequences are publicly available in the Open Science Framework

1199    database (https://osf.io/jze72/?view_only=d2dd1badf37e4bafae1e12731a0cc40d). Custom

1200    code used is available at https://github.com/hamidbeiki/Cattle-Genome.

1201    **Ethics approval and consent to participate**

1202    Procedures for tissue collection followed the Animal Care and Use protocol (#18464) approved

1203    by the Institutional Animal Care and Use Committee (IACUC), University of California, Davis

1204    (UCD).

1205    **Consent for publication**

1206    Not applicable

1207    **Competing interests**

1208    The authors declare no competing interests.

1209    **Funding**

## 1213 Acknowledgments

## 1216 Authors' contributions

1217 H.B., B.M.M., H.J., H.Z., M.R., P.J.R., S.M., T.P.L.S., W.L., Z.J., and J.M.R. conceived and designed

1218 the project; C.K., W.M., and W.L. generated RNA-seq and miRNA-seq data; D.K., G.B., J.T., and

1219 K.D. participated in tissue collection; R.H and H.J prepared cells; J.J.M., X.Z., X.H., and Z.J.

1220 generated W.T.T.S-seq data, X.X., P.J.R. and H.J generated ChIP-seq data; M.R.J. generated

1221 ATAC-seq data; T.P.L.S. generated PacBio Iso-seq data; G.R. and S.C. conducted sequencing of

1222 RNA-seg, miRNA-seq, ChIP-seq, and ATAC-seq data;  H.B. conducted bioinformatics data

1223 analysis and drafted the manuscript, which was edited by C.A.P., B.M.M., H.J., H.Z., J.E.K., M.R.,

1224 P.J.R., S.M., T.P.L.S., W.L., Z.J. and J.M.R.; Z.H. created the web-based database for the trait

1225 similarity network; all authors read and approved the final manuscript.

## 1226 Endnotes

1227 Mention of trade names or commercial products in this publication is solely for the purpose of

1228 providing specific information and does not imply recommendation or endorsement by the U.S.

1229 Department of Agriculture. USDA is an equal opportunity provider and employer.

1232

## 1233    References

1234    1.    Roth JA and Tuggle CK. Livestock models in translational medicine. ILAR J. 2015;56 1:1-6.
1235          doi:10.1093/ilar/ilv011.
1236    2.    Beiki H, Liu H, Huang J, Manchanda N, Nonneman D, Smith TPL, et al. Improved
1237          annotation of the domestic pig genome through integration of Iso-Seq and RNA-seq
1238          data. BMC Genomics. 2019;20 1:344. doi:10.1186/s12864-019-5709-y.
1239    3.    Hindorff LA, Sethupathy P, Junkins HA, Ramos EM, Mehta JP, Collins FS, et al. Potential
1240          etiologic and functional implications of genome-wide association loci for human
1241          diseases and traits. Proc Natl Acad Sci U S A. 2009;106 23:9362-7.
1242          doi:10.1073/pnas.0903103106.
1243    4.    Jereb S, Hwang HW, Van Otterloo E, Govek EE, Fak JJ, Yuan Y, et al. Differential 3'
1244          Processing of Specific Transcripts Expands Regulatory and Protein Diversity Across
1245          Neuronal Cell Types. Elife. 2018;7  doi:10.7554/eLife.34042.
1246    5.    Schurch NJ, Cole C, Sherstnev A, Song J, Duc C, Storey KG, et al. Improved annotation of
1247          3' untranslated regions and complex loci by combination of strand-specific direct RNA
1248          sequencing, RNA-Seq and ESTs. PLoS One. 2014;9 4:e94270.
1249          doi:10.1371/journal.pone.0094270.
1250    6.    Ambros V. The functions of animal microRNAs. Nature. 2004;431 7006:350-5.
1251          doi:10.1038/nature02871.
1252    7.    Bartel DP. MicroRNAs: genomics, biogenesis, mechanism, and function. Cell. 2004;116
1253          2:281-97. doi:10.1016/s0092-8674(04)00045-5.
1254    8.    Yates LA, Norbury CJ and Gilbert RJ. The long and short of microRNA. Cell. 2013;153
1255          3:516-9. doi:10.1016/j.cell.2013.04.003.
1256    9.    Halstead MM, Islas-Trejo A, Goszczynski DE, Medrano JF, Zhou H and Ross PJ. Large-
1257          Scale Multiplexing Permits Full-Length Transcriptome Annotation of 32 Bovine Tissues
1258          From a Single Nanopore Flow Cell. Front Genet. 2021;12:664260.
1259          doi:10.3389/fgene.2021.664260.
1260    10.   Goszczynski DE, Halstead MM, Islas-Trejo AD, Zhou H and Ross PJ. Transcription
1261          initiation mapping in 31 bovine tissues reveals complex promoter activity, pervasive
1262          transcription, and tissue-specific promoter usage. Genome Res. 2021;31 4:732-44.
1263          doi:10.1101/gr.267336.120.
1264    11.   Kozomara A, Birgaoanu M and Griffiths-Jones S. miRBase: from microRNA sequences to
1265          function. Nucleic Acids Res. 2019;47 D1:D155-D62. doi:10.1093/nar/gky1141.

68

1266    12.    Araujo PR, Yoon K, Ko D, Smith AD, Qiao M, Suresh U, et al. Before It Gets Started:
1267            Regulating Translation at the 5' UTR. Comp Funct Genomics. 2012;2012:475731.
1268            doi:10.1155/2012/475731.
1269    13.    Gerber S, Schratt G and Germain PL. Streamlining differential exon and 3' UTR usage
1270            with diffUTR. BMC Bioinformatics. 2021;22 1:189. doi:10.1186/s12859-021-04114-7.
1271    14.    Andrews SJ and Rothnagel JA. Emerging evidence for functional peptides encoded by
1272            short open reading frames. Nat Rev Genet. 2014;15 3:193-204. doi:10.1038/nrg3520.
1273    15.    Kumari P and Sampath K. cncRNAs: Bi-functional RNAs with protein coding and non-
1274            coding functions. Semin Cell Dev Biol. 2015;47-48:40-51.
1275            doi:10.1016/j.semcdb.2015.10.024.
1276    16.    Nam JW, Choi SW and You BH. Incredible RNA: Dual Functions of Coding and Noncoding.
1277            Mol Cells. 2016;39 5:367-74. doi:10.14348/molcells.2016.0039.
1278    17.    Gonzàlez-Porta M, Frankish A, Rung J, Harrow J and Brazma A. Transcriptome analysis of
1279            human tissues and cell lines reveals one dominant transcript per gene. Genome Biol.
1280            2013;14 7:R70. doi:10.1186/gb-2013-14-7-r70.
1281    18.    Mayba O, Gilbert HN, Liu J, Haverty PM, Jhunjhunwala S, Jiang Z, et al. MBASED: allele-
1282            specific expression detection in cancer tissues and cell lines. Genome Biol. 2014;15
1283            8:405. doi:10.1186/s13059-014-0405-3.
1284    19.    Hubé F, Velasco G, Rollin J, Furling D and Francastel C. Steroid receptor RNA activator
1285            protein binds to and counteracts SRA RNA-mediated activation of MyoD and muscle
1286            differentiation. Nucleic Acids Res. 2011;39 2:513-25. doi:10.1093/nar/gkq833.
1287    20.    Kurosaki T, Popp MW and Maquat LE. Quality and quantity control of gene expression
1288            by nonsense-mediated mRNA decay. Nat Rev Mol Cell Biol. 2019;20 7:406-20.
1289            doi:10.1038/s41580-019-0126-2.
1290    21.    Wollerton MC, Gooding C, Wagner EJ, Garcia-Blanco MA and Smith CW. Autoregulation
1291            of polypyrimidine tract binding protein by alternative splicing leading to nonsense-
1292            mediated decay. Mol Cell. 2004;13 1:91-100. doi:10.1016/s1097-2765(03)00502-1.
1293    22.    Nickless A, Bailis JM and You Z. Control of gene expression through the nonsense-
1294            mediated RNA decay pathway. Cell Biosci. 2017;7:26. doi:10.1186/s13578-017-0153-7.
1295    23.    Supek F, Lehner B and Lindeboom RGH. To NMD or Not To NMD: Nonsense-Mediated
1296            mRNA Decay in Cancer and Other Genetic Diseases. Trends Genet. 2021;37 7:657-68.
1297            doi:10.1016/j.tig.2020.11.002.
1298    24.    Mitrovich QM and Anderson P. mRNA surveillance of expressed pseudogenes in C.
1299            elegans. Curr Biol. 2005;15 10:963-7. doi:10.1016/j.cub.2005.04.055.
1300    25.    Colombo M, Karousis ED, Bourquin J, Bruggmann R and Mühlemann O. Transcriptome-
1301            wide identification of NMD-targeted human mRNAs reveals extensive redundancy
1302            between SMG6- and SMG7-mediated degradation pathways. RNA. 2017;23 2:189-201.
1303            doi:10.1261/rna.059055.116.
1304    26.    Milligan MJ and Lipovich L. Pseudogene-derived lncRNAs: emerging regulators of gene
1305            expression. Front Genet. 2014;5:476. doi:10.3389/fgene.2014.00476.
1306    27.    Stewart GL, Enfield KSS, Sage AP, Martinez VD, Minatel BC, Pewarchuk ME, et al.
1307            Aberrant Expression of Pseudogene-Derived lncRNAs as an Alternative Mechanism of
1308            Cancer Gene Regulation in Lung Adenocarcinoma. Front Genet. 2019;10:138.
1309            doi:10.3389/fgene.2019.00138.

1310 28. Lou W, Ding B and Fu P. Pseudogene-Derived lncRNAs and Their miRNA Sponging
1311 Mechanism in Human Cancer. Front Cell Dev Biol. 2020;8:85.
1312 doi:10.3389/fcell.2020.00085.
1313 29. Anderson DM, Anderson KM, Chang CL, Makarewich CA, Nelson BR, McAnally JR, et al. A
1314 micropeptide encoded by a putative long noncoding RNA regulates muscle
1315 performance. Cell. 2015;160 4:595-606. doi:10.1016/j.cell.2015.01.009.
1316 30. Mackowiak SD, Zauber H, Bielow C, Thiel D, Kutz K, Calviello L, et al. Extensive
1317 identification and analysis of conserved small ORFs in animals. Genome Biol.
1318 2015;16:179. doi:10.1186/s13059-015-0742-x.
1319 31. Olexiouk V, Crappé J, Verbruggen S, Verhegen K, Martens L and Menschaert G.
1320 sORFs.org: a repository of small ORFs identified by ribosome profiling. Nucleic Acids Res.
1321 2016;44 D1:D324-9. doi:10.1093/nar/gkv1175.
1322 32. Li J and Liu C. Coding or Noncoding, the Converging Concepts of RNAs. Front Genet.
1323 2019;10:496. doi:10.3389/fgene.2019.00496.
1324 33. Wei L-H and Guo JU. Coding functions of "noncoding" RNAs. Science. 2020;367
1325 6482:1074-5. doi:10.1126/science.aba6117.
1326 34. Sammeth M, Foissac S and Guigó R. A general definition and nomenclature for
1327 alternative splicing events. PLoS Comput Biol. 2008;4 8:e1000147.
1328 doi:10.1371/journal.pcbi.1000147.
1329 35. Mazin PV, Khaitovich P, Cardoso-Moreira M and Kaessmann H. Alternative splicing
1330 during mammalian organ development. Nature Genetics. 2021;53 6:925-34.
1331 doi:10.1038/s41588-021-00851-w.
1332 36. Wu Z, Yang KK, Liszka MJ, Lee A, Batzilla A, Wernick D, et al. Signal Peptides Generated
1333 by Attention-Based Neural Networks. ACS Synth Biol. 2020;9 8:2154-61.
1334 doi:10.1021/acssynbio.0c00219.
1335 37. Chen J and Chen ZJ. Regulation of NF-κB by ubiquitination. Curr Opin Immunol. 2013;25
1336 1:4-12. doi:10.1016/j.coi.2012.12.005.
1337 38. Karalis KP, Venihaki M, Zhao J, van Vlerken LE and Chandras C. NF-kappaB participates in
1338 the corticotropin-releasing, hormone-induced regulation of the pituitary
1339 proopiomelanocortin gene. J Biol Chem. 2004;279 12:10837-40.
1340 doi:10.1074/jbc.M313063200.
1341 39. O'Shaughnessy PJ, Fleming LM, Jackson G, Hochgeschwender U, Reed P and Baker PJ.
1342 Adrenocorticotropic hormone directly stimulates testosterone production by the fetal
1343 and neonatal mouse testis. Endocrinology. 2003;144 8:3279-84. doi:10.1210/en.2003-
1344 0277.
1345 40. Richburg JH, Myers JL and Bratton SB. The role of E3 ligases in the ubiquitin-dependent
1346 regulation of spermatogenesis. Semin Cell Dev Biol. 2014;30:27-35.
1347 doi:10.1016/j.semcdb.2014.03.001.
1348 41. Kumar S, Lee HJ, Park HS and Lee K. Testis-Specific GTPase (TSG): An oligomeric protein.
1349 BMC Genomics. 2016;17 1:792. doi:10.1186/s12864-016-3145-9.
1350 42. Rajala-Schultz PJ, Gröhn YT, McCulloch CE and Guard CL. Effects of clinical mastitis on
1351 milk yield in dairy cows. J Dairy Sci. 1999;82 6:1213-20. doi:10.3168/jds.S0022-
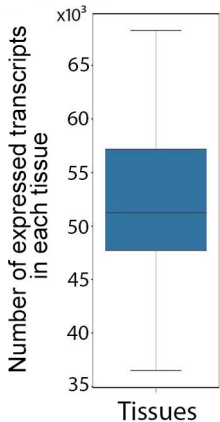1352 0302(99)75344-0.

1353 43. Martí De Olives A, Díaz JR, Molina MP and Peris C. Quantification of milk yield and
1354      composition changes as affected by subclinical mastitis during the current lactation in
1355      sheep. J Dairy Sci. 2013;96 12:7698-708. doi:10.3168/jds.2013-6998.
1356 44. Halasa T and Kirkeby C. Differential Somatic Cell Count: Value for Udder Health
1357      Management. Front Vet Sci. 2020;7:609055. doi:10.3389/fvets.2020.609055.
1358 45. Remnant J, Green MJ, Huxley J, Hirst-Beecham J, Jones R, Roberts G, et al. Association of
1359      lameness and mastitis with return-to-service oestrus detection in the dairy cow. Vet
1360      Rec. 2019;185 14:442. doi:10.1136/vr.105535.
1361 46. Miles AM, McArt JAA, Leal Yepes FA, Stambuk CR, Virkler PD and Huson HJ. Udder and
1362      teat conformational risk factors for elevated somatic cell count and clinical mastitis in
1363      New York Holsteins. Prev Vet Med. 2019;163:7-13.
1364      doi:10.1016/j.prevetmed.2018.12.010.
1365 47. Lima FS, Silvestre FT, Peñagaricano F and Thatcher WW. Early genomic prediction of
1366      daughter pregnancy rate is associated with improved reproductive performance in
1367      Holstein dairy cows. J Dairy Sci. 2020;103 4:3312-24. doi:10.3168/jds.2019-17488.
1368 48. Hertl JA, Schukken YH, Tauer LW, Welcome FL and Gröhn YT. Does clinical mastitis in the
1369      first 100 days of lactation 1 predict increased mastitis occurrence and shorter herd life in
1370      dairy cows? J Dairy Sci. 2018;101 3:2309-23. doi:10.3168/jds.2017-12615.
1371 49. Kaniyamattam K, De Vries A, Tauer LW and Gröhn YT. Economics of reducing antibiotic
1372      usage for clinical mastitis and metritis through genomic selection. J Dairy Sci. 2020;103
1373      1:473-91. doi:10.3168/jds.2018-15817.
1374 50. Green TC, Jago JG, Macdonald KA and Waghorn GC. Relationships between residual feed
1375      intake, average daily gain, and feeding behavior in growing dairy heifers. J Dairy Sci.
1376      2013;96 5:3098-107. doi:10.3168/jds.2012-6087.
1377 51. Elolimy AA, Abdelmegeid MK, McCann JC, Shike DW and Loor JJ. Residual feed intake in
1378      beef cattle and its association with carcass traits, ruminal solid-fraction bacteria, and
1379      epithelium gene expression. J Anim Sci Biotechnol. 2018;9:67. doi:10.1186/s40104-018-
1380      0283-8.
1381 52. Weber C, Hametner C, Tuchscherer A, Losand B, Kanitz E, Otten W, et al. Variation in fat
1382      mobilization during early lactation differently affects feed intake, body condition, and
1383      lipid and glucose metabolism in high-yielding dairy cows. J Dairy Sci. 2013;96 1:165-80.
1384      doi:10.3168/jds.2012-5574.
1385 53. Yi Z, Li X, Luo W, Xu Z, Ji C, Zhang Y, et al. Feed conversion ratio, residual feed intake and
1386      cholecystokinin type A receptor gene polymorphisms are associated with feed intake
1387      and average daily gain in a Chinese local chicken population. J Anim Sci Biotechnol.
1388      2018;9:50. doi:10.1186/s40104-018-0261-1.
1389 54. Liu E and VandeHaar MJ. Relationship of residual feed intake and protein efficiency in
1390      lactating cows fed high- or low-protein diets. J Dairy Sci. 2020;103 4:3177-90.
1391      doi:10.3168/jds.2019-17567.
1392 55. Clare M, Richard P, Kate K, Sinead W, Mark M and David K. Residual feed intake
1393      phenotype and gender affect the expression of key genes of the lipogenesis pathway in
1394      subcutaneous adipose tissue of beef cattle. J Anim Sci Biotechnol. 2018;9:68.
1395      doi:10.1186/s40104-018-0282-9.

1396    56.    Houlahan K, Schenkel FS, Hailemariam D, Lassen J, Kargo M, Cole JB, et al. Effects of
1397            Incorporating Dry Matter Intake and Residual Feed Intake into a Selection Index for
1398            Dairy Cattle Using Deterministic Modeling. Animals (Basel). 2021;11 4
1399            doi:10.3390/ani11041157.

1400    57.    Tixier-Boichard M, Fabre S, Dhorne-Pollet S, Goubil A, Acloque H, Vincent-Naulleau S, et
1401            al. Tissue Resources for the Functional Annotation of Animal Genomes. Front Genet.
1402            2021;12:666265. doi:10.3389/fgene.2021.666265.

1403    58.    Farr VC, Stelwagen K, Cate LR, Molenaar AJ, McFadden TB and Davis SR. An improved
1404            method for the routine biopsy of bovine mammary tissue. J Dairy Sci. 1996;79 4:543-9.
1405            doi:10.3168/jds.S0022-0302(96)76398-1.

1406    59.    Zhou X, Li R, Michal JJ, Wu XL, Liu Z, Zhao H, et al. Accurate Profiling of Gene Expression
1407            and Alternative Polyadenylation with Whole Transcriptome Termini Site Sequencing
1408            (WTTS-Seq). Genetics. 2016;203 2:683-97. doi:10.1534/genetics.116.188508.

1409    60.    Krueger F: https://www.bioinformatics.babraham.ac.uk/projects/trim_galore/.  (2019).

1410    61.    Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, et al. STAR: ultrafast
1411            universal RNA-seq aligner. Bioinformatics. 2013;29 1:15-21.
1412            doi:10.1093/bioinformatics/bts635.

1413    62.    Liao Y, Smyth GK and Shi W. featureCounts: an efficient general purpose program for
1414            assigning sequence reads to genomic features. Bioinformatics. 2014;30 7:923-30.
1415            doi:10.1093/bioinformatics/btt656.

1416    63.    Leek J, Johnson W, Parker HS, Fertig EJ, Jaffe AE, Zhang Y, et al. *sva: Surrogate Variable
1417            Analysis* . R package version 3.30.0. 2021.

1418    64.    Grabherr MG, Haas BJ, Yassour M, Levin JZ, Thompson DA, Amit I, et al. Full-length
1419            transcriptome assembly from RNA-Seq data without a reference genome. Nat
1420            Biotechnol. 2011;29 7:644-52. doi:10.1038/nbt.1883.

1421    65.    Hass B: https://hpcgridrunner.github.io/.  (2015).

1422    66.    Tange O: GNU Parallel. https://doi.org/10.5281/zenodo.1146014.  (2018).

1423    67.    PacificBiosciences: https://www.pacb.com/products-and-services/analytical-
1424            software/smrt-analysis/.  (2018).

1425    68.    Pedersen BS and Quinlan AR. Mosdepth: quick coverage calculation for genomes and
1426            exomes. Bioinformatics. 2018;34 5:867-8. doi:10.1093/bioinformatics/btx699.

1427    69.    Hackl T, Hedrich R, Schultz J and Förster F. proovread: large-scale high-accuracy PacBio
1428            correction through iterative short read consensus. Bioinformatics. 2014;30 21:3004-11.
1429            doi:10.1093/bioinformatics/btu392.

1430    70.    Wang JR, Holt J, McMillan L and Jones CD. FMLRC: Hybrid long read error correction
1431            using an FM-index. BMC Bioinformatics. 2018;19 1:50. doi:10.1186/s12859-018-2051-3.

1432    71.    Wheeler DL, Church DM, Federhen S, Lash AE, Madden TL, Pontius JU, et al. Database
1433            resources of the National Center for Biotechnology. Nucleic Acids Res. 2003;31 1:28-33.
1434            doi:10.1093/nar/gkg033.

1435    72.    Aken BL, Ayling S, Barrell D, Clarke L, Curwen V, Fairley S, et al. The Ensembl gene
1436            annotation system. Database (Oxford). 2016;2016  doi:10.1093/database/baw093.

1437    73.    Kang YJ, Yang DC, Kong L, Hou M, Meng YQ, Wei L, et al. CPC2: a fast and accurate
1438            coding potential calculator based on sequence intrinsic features. Nucleic Acids Res.
1439            2017;45 W1:W12-W6. doi:10.1093/nar/gkx428.

1440    74.    Salmela L and Schröder J. Correcting errors in short reads by multiple alignments.
1441           Bioinformatics. 2011;27 11:1455-61. doi:10.1093/bioinformatics/btr170.
1442    75.    Hannon GJ: FASTX-Toolkit.   http://hannonlab.cshl.edu/fastx_toolkit.  (2010).
1443    76.    Kern C, Wang Y, Xu X, Pan Z, Halstead M, Chanthavixay G, et al. Functional annotations
1444           of three domestic animal genomes provide vital resources for comparative and
1445           agricultural research. Nat Commun. 2021;12 1:1821. doi:10.1038/s41467-021-22100-8.
1446    77.    Bindea G, Mlecnik B, Hackl H, Charoentong P, Tosolini M, Kirilovsky A, et al. ClueGO: a
1447           Cytoscape plug-in to decipher functionally grouped gene ontology and pathway
1448           annotation networks. Bioinformatics. 2009;25 8:1091-3.
1449           doi:10.1093/bioinformatics/btp101.
1450    78.    Huntley RP, Sawford T, Mutowo-Meullenet P, Shypitsyna A, Bonilla C, Martin MJ, et al.
1451           The GOA database: gene Ontology annotation updates for 2015. Nucleic Acids Res.
1452           2015;43 Database issue:D1057-63. doi:10.1093/nar/gku1113.
1453    79.    Kim KI and van de Wiel MA. Effects of dependence in high-dimensional multiple testing
1454           problems. BMC Bioinformatics. 2008;9 1:114. doi:10.1186/1471-2105-9-114.
1455    80.    Trincado JL, Entizne JC, Hysenaj G, Singh B, Skalic M, Elliott DJ, et al. SUPPA2: fast,
1456           accurate, and uncertainty-aware differential splicing analysis across multiple conditions.
1457           Genome Biol. 2018;19 1:40. doi:10.1186/s13059-018-1417-1.
1458    81.    Friedländer MR, Mackowiak SD, Li N, Chen W and Rajewsky N. miRDeep2 accurately
1459           identifies known and hundreds of novel microRNA genes in seven animal clades. Nucleic
1460           Acids Res. 2012;40 1:37-52. doi:10.1093/nar/gkr688.
1461    82.    Ludwig N, Leidinger P, Becker K, Backes C, Fehlmann T, Pallasch C, et al. Distribution of
1462           miRNA expression across human tissues. Nucleic Acids Res. 2016;44 8:3865-77.
1463           doi:10.1093/nar/gkw116.
1464    83.    Hu ZL, Park CA and Reecy JM. Building a livestock genetic and genomic information
1465           knowledgebase through integrative developments of Animal QTLdb and CorrDB. Nucleic
1466           Acids Res. 2019;47 D1:D701-D10. doi:10.1093/nar/gky1084.
1467    84.    Shannon P, Markiel A, Ozier O, Baliga NS, Wang JT, Ramage D, et al. Cytoscape: a
1468           software environment for integrated models of biomolecular interaction networks.
1469           Genome Res. 2003;13 11:2498-504. doi:10.1101/gr.1239303.
1470    85.    Almagro Armenteros JJ, Tsirigos KD, Sønderby CK, Petersen TN, Winther O, Brunak S, et
1471           al. SignalP 5.0 improves signal peptide predictions using deep neural networks. Nature
1472           Biotechnology. 2019;37 4:420-3. doi:10.1038/s41587-019-0036-z.
1473    86.    Uhlén M, Fagerberg L, Hallström BM, Lindskog C, Oksvold P, Mardinoglu A, et al.
1474           Proteomics. Tissue-based map of the human proteome. Science. 2015;347
1475           6220:1260419. doi:10.1126/science.1260419.

1476

Figure 1

**A**

**B**

Figure 2

Figure 3

Figure 3

Legend:
- Transcripts reported in Ensembl gene-build (Release 2021-03)
- Transcripts reported in NCBI gene-build (Release 106)
- Transcripts structurally validated by an independent Oxford Nanopore experiment (32 pooled tissues)
- Transcripts supported by RAMPAGE data from an independent experiment (30 tissues)
- Transcripts structurally validated by an independent PacBio Iso-seq experiment (7 tissues)
- Transcripts structurally validated by de novo-assembled transcripts from an independent RNA-seq experiment (7 tissues)
- Transcripts supported by WTTS data (24 tissues)
- Transcripts expressed in multiple tissues
- Transcripts supported by ATAC-seq data (47 tissues) and H3K4me3, H3K4me1, H3K27ac, and CTCF-DNA binding data from independent experiments (8 tissues)
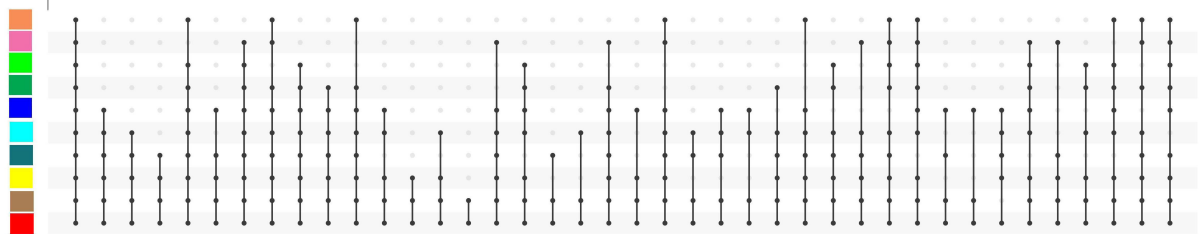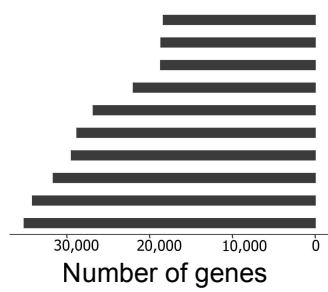- All transcripts (47 tissues)

Figure 4

Figure 5

# Figure 6

Figure 7

Figure 7

Figure 8

Figure 9

**A**



Pituitary genes that are close to "percentage of normal sperm" QTLs
(246 genes)

**B**

Figure 10

**A**

Number of miRNAs found in different combination of groups

annotated mirnas supported by H3K4me3 peaks from an independent experiment (8 tissues)

annotated mirnas supported by H3K4me1 peaks from an independent experiment (8 tissues)

annotated mirnas supported by H3K27ac peaks from an independent experiment (8 tissues)

annotated mirnas supported by H3K27ac peaks from an independent experiment (8 tissues)

annotated mirnas supported by CTFC peaks from an independent experiment (8 tissues)

annotated mirnas supported by ATAC-seq peaks (8 tissues)

Number of miRNAs

**B**

Number of miRNAs found in different combination of groups

un-annotated mirnas supported by H3K4me3 peaks from an independent experiment (8 tissues)

un-annotated mirnas supported by H3K4me1 peaks from an independent experiment (8 tissues)

un-annotated mirnas supported by H3K27ac peaks from an independent experiment (8 tissues)

un-annotated mirnas supported by H3K27ac peaks from an independent experiment (8 tissues)

un-annotated mirnas supported by CTFC peaks from an independent experiment (8 tissues)
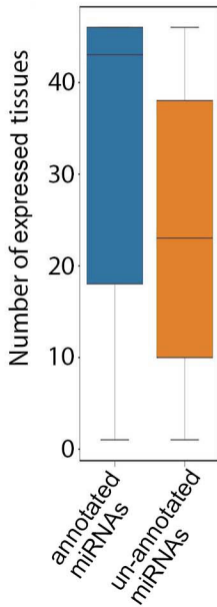
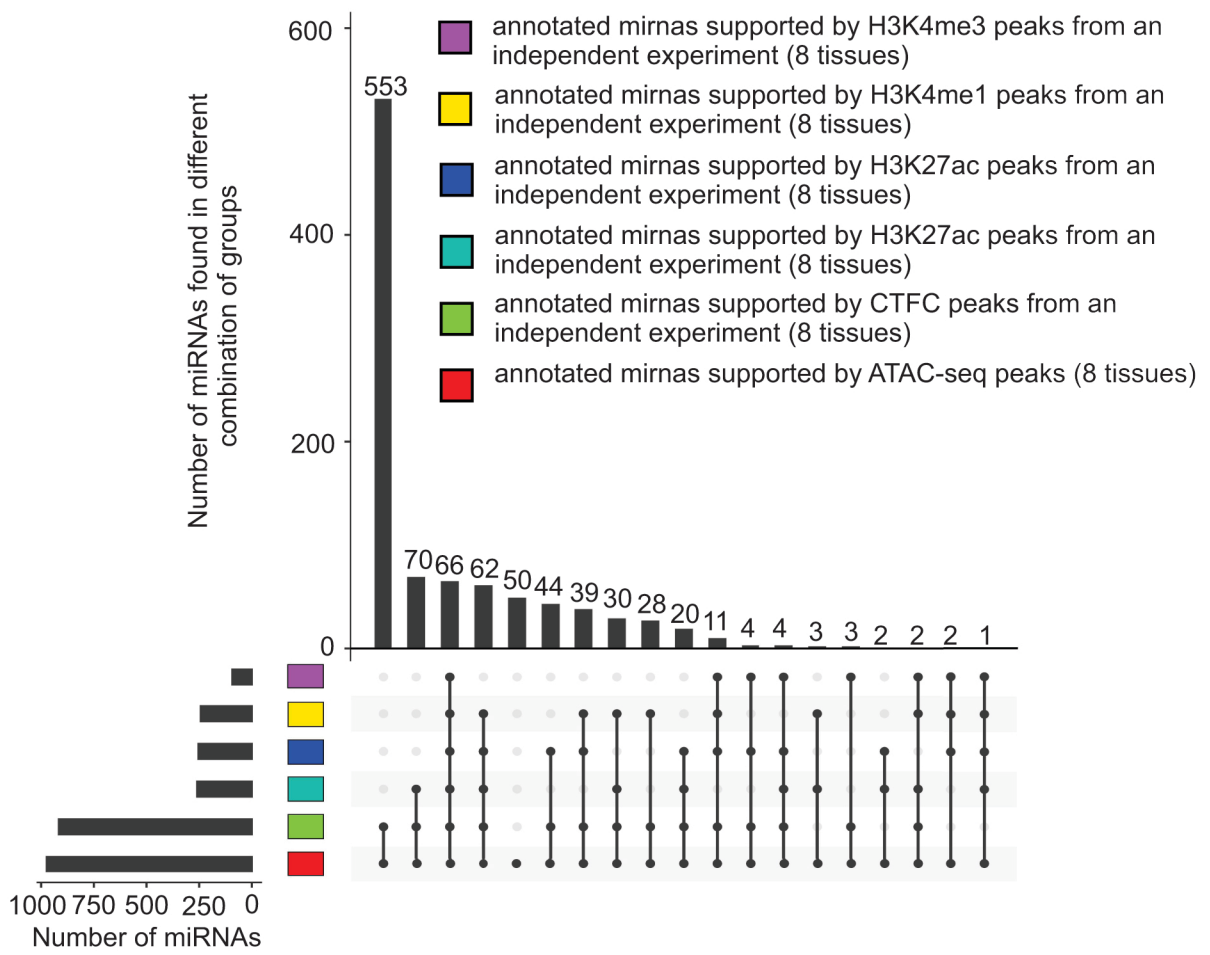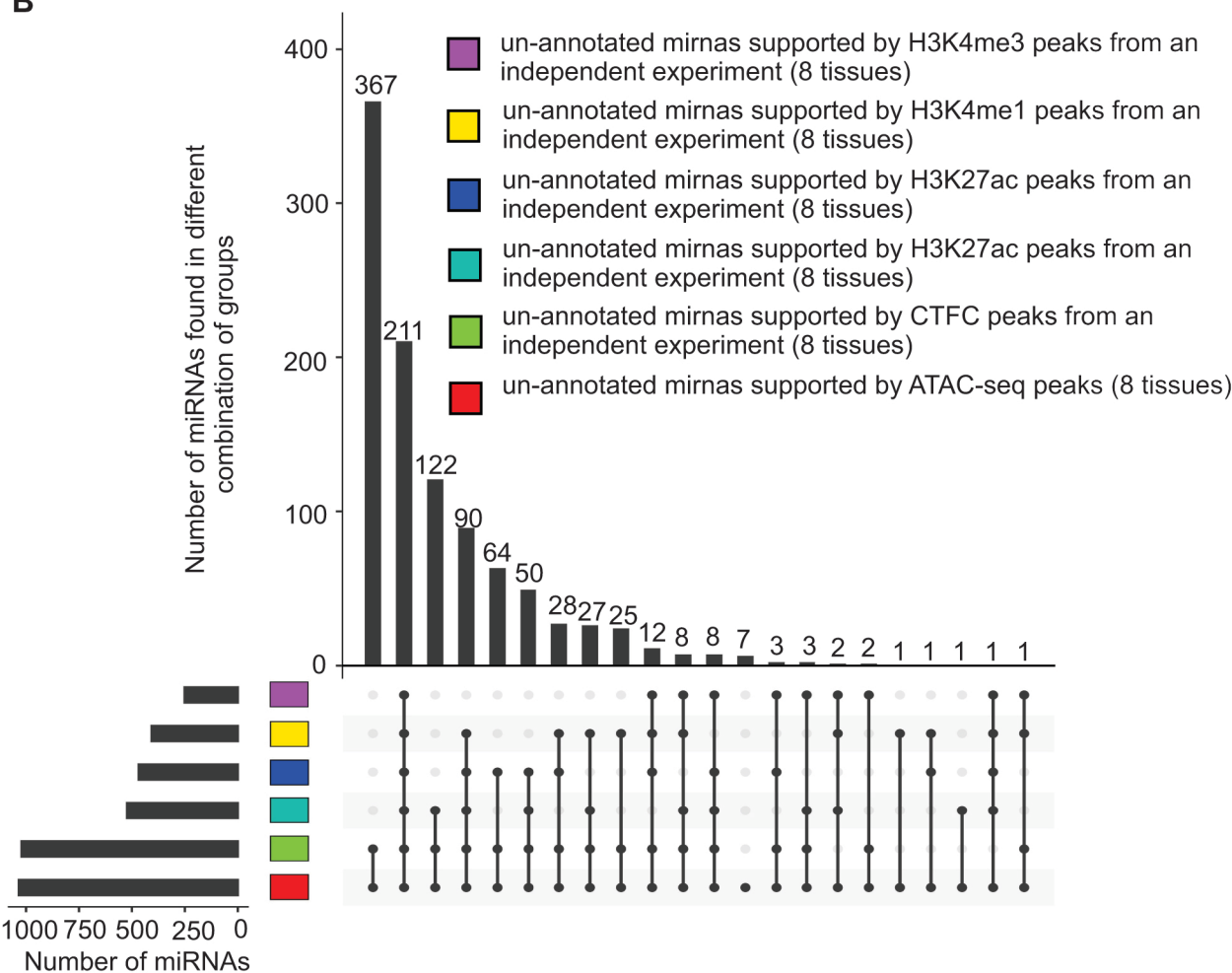un-annotated mirnas supported by ATAC-seq peaks (8 tissues)

Number of miRNAs

Click here to access/download
**Supplementary Material**
Supplemental_file1.docx

Click here to access/download

**Supplementary Material**

Supplemental_file2.xlsx

Supplemental File 3

Click here to access/download
**Supplementary Material**
Supplemental_file3.xlsx

Click here to access/download
**Supplementary Material**
Supplemental_file4.xlsx

Click here to access/download
**Supplementary Material**
Supplemental_file5.xlsx

Click here to access/download
**Supplementary Material**
Supplemental_file6.xlsx

Click here to access/download
**Supplementary Material**
Supplemental_file7.xlsx

Click here to access/download
**Supplementary Material**
Supplemental_file8.xlsx

Click here to access/download
**Supplementary Material**
Supplemental_file9.xlsx

Click here to access/download

**Supplementary Material**

Supplemental_file10.xlsx

Click here to access/download
**Supplementary Material**
Supplemental_file11.xlsx

Click here to access/download
**Supplementary Material**
Supplemental_file12.xlsx

Click here to access/download

**Supplementary Material**

Supplemental_file13.xlsx

Click here to access/download
**Supplementary Material**
Supplemental_file14.xlsx

Click here to access/download
**Supplementary Material**
Supplemental_file15.xlsx

Click here to access/download
**Supplementary Material**
Supplemental_file16.xlsx

Click here to access/download
**Supplementary Material**
Supplemental_file17.xlsx

Click here to access/download

**Supplementary Material**

Supplemental_file18.xlsx

Click here to access/download

**Supplementary Material**

Supplemental_file19.xlsx

Click here to access/download
**Supplementary Material**
Supplemental_file20.xlsx

Click here to access/download
**Supplementary Material**
Supplemental_file21.xlsx

Click here to access/download
**Supplementary Material**
Supplemental_file22.xlsx

**IOWA STATE UNIVERSITY**
OF SCIENCE AND TECHNOLOGY

Office of the Vice President for Research

2610 Beardshear Hall

Ames, Iowa 50011-2036

515 294-6344

FAX 515 294-6100

February 11, 2023

Dear GigaScience Editor,

I am very excited to share our new manuscript, entitled "Facilitating Functional genomics of cattle through integration of multi-omics data" for your consideration. Given the impact and novelty of our findings, we are confident that this work is ideally suited for publication in *Genome Research*. The comprehensive analyses presented herein crystallized into a decisive advance for the cattle transcriptomics field and beyond. This study applied a comprehensive set of transcriptome and chromatin state data from 47 cattle tissues and cell types to identify previously unannotated genes and improve the annotation of thousands of protein-coding and non-coding genes.

Predicted novel genes and transcripts were highly supported by independent Pacific Biosciences single-molecule long-read isoform sequencing (PacBio Iso-Seq), Oxford Nanopore Technologies sequencing (ONT-seq), Illumina high-throughput RNA sequencing (RNA-seq), Whole Transcriptome Termini Site Sequencing (WTTS-seq), RNA Annotation and Mapping of Promoters for the Analysis of Gene Expression (RAMPAGE), chromatin immunoprecipitation sequencing (ChIP-seq), and Assay for Transposase-Accessible Chromatin using sequencing (ATAC-seq) data.

Our key findings show that around half of protein-coding genes in each tissue are bifunctional and transcribe both in coding and noncoding isoforms. Critically, we identified 3,744 genes that functioned as non-coding genes in fetal tissues, but as protein coding genes in adult tissues. Most interestingly when the transcriptome was integrated with publicly available cattle QTL/association data using a novel bioinformatics approach, we were able to study tissue-tissue interconnection involved in different traits and construct the first bovine trait similarity network. These independent findings agree with published trait correlation data and move us closer to being able to identify the gene networks that underlie genetic correlations between traits.

Given these results, we strongly maintain that our work will be of general interest to the broad readership of *Genome Research*. All datasets generated in this study have been submitted to public databases, and all gene/protein names and symbols used in the paper adhere to approved nomenclature guidelines for specific species. We also confirm that this work is original and has not been published elsewhere, nor is it currently under consideration for publication elsewhere.

Thank you for your consideration and we look forward to hearing from you.

Sincerely,

James Reee

James Reecy
Associate Vice President for Research