1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47

**Tissue Resolved, Gene Structure Refined Equine Transcriptome**

Mansour,T.A.*[1,2], Scott, E.Y.*[3], Finno, C.J.[1], Bellone, R.R.[4], Mienaltowski, M.J.[3], Penedo, M.C.[4], Ross, P.J.[3], Valberg, S.J.[5], Murray, J.D.[1,3], Brown, C.T[1].

[1] Department of Population Health and Reproduction, University of California, Davis, [2] Department of Clinical Pathology, College of Medicine, Mansoura University, Egypt [3] Department of Animal Science, University of California, Davis, [4] Veterinary Genetics Laboratory, University of California, Davis, [5] Large Animal Clinical Sciences, Michigan State University, College of Veterinary Medicine

*Both authors contributed equally to this manuscript

**To be submitted to BMC Genomics**

48 **Abstract**

49 *Background:* Transcriptome interpretation relies on a good-quality reference transcriptome for accurate

50 quantification of gene expression as well as functional analysis of genetic variants. The current

51 annotation of the horse genome lacks the specificity and sensitivity necessary to assess gene expression

52 especially at the isoform level, and suffers from insufficient annotation of untranslated regions (UTR).

53 We built an annotation pipeline for horse and used it to integrate 1.9 billion reads from multiple RNA-

54 seq data sets into a new refined transcriptome.

55 *Results:* This equine transcriptome integrates eight different tissues from 59 individuals and improves

56 gene structure and isoform resolution while providing considerable tissue-specific information. We

57 utilized four levels of transcript filtration in our pipeline, aimed at producing several transcriptome

58 versions that are suitable for different downstream analyses. Our most refined transcriptome includes

59 36,876 genes and 76,125 isoforms, with 6474 candidate transcriptional loci novel to the equine

60 transcriptome.

61 *Conclusions:* We have employed a variety of descriptive statistics and figures that demonstrate the

62 quality and content of the transcriptome. The equine transcriptomes that are provided by this pipeline

63 show the best tissue-specific resolution of any equine transcriptome to date and can serve several types

64 of downstream analyses.

65

66 **Keywords:** Equine transcriptome, tissue-specificity, RNA-seq
67

68 **Introduction**

69     Transcriptomics is rapidly evolving from a focus on novel gene identification to resolving structural

70 gene details. The transcriptomes of better-studied organisms, such as Drosophila, mouse and human

71 have been updated to accommodate for this transition [1-3]. However, for less well characterized

72 animals, such as the horse, there is often only annotation of a single variant of a gene with insufficient

73 annotation of multiple splice variants, UTR extensions and non-protein coding RNA. This lack of

74 information can challenge subsequent differential gene expression analyses and functional studies. There

75 have been several attempts to improve the equine transcriptome with single tissue transcriptomes from

76 lamellar tissue [4] or peripheral blood mononuclear cells [5] and from pooled composites of various

77 tissues [6, 7], however a broader effort defining and integrating many tissue-specific transcriptomes and

78 obtaining the library depth and strand information required to capture gene complexity is still needed.

79 ENSEMBL and NCBI provide publically available annotations for several vertebrate genomes

80 including horses [8]. Both underlying annotation pipelines integrate homology search and *ab initio*

81 prediction however accurate UTR prediction and isoform recognition require species-specific

82 transcriptional evidence [9, 10]. For this equine transcriptome, the transcriptional evidence provided by

83 total RNA sequencing (RNA-seq) was the basis of our gene annotation. This approach permits more

84 reliable discovery of novel genes and isoforms, extension of UTRs and the flexibility necessary to

85 establish a balance between sensitivity and specificity of gene detection for downstream applications.

86 Our annotation integrates the benefits of increased depth in reads and strand-specificity, for some

87 tissues, as well as using a range of tissues from many horses, which allows tissue-specific

88 transcriptomes to be extracted. We have incorporated RNA-seq from a diverse set of 8 tissues ranging

89 from the central nervous system (CNS), skin and skeletal muscle tissues in adults to the inner cell mass

90 (ICM) and trophectoderm (TE) in embryonic tissues (Table 1). The diversity in age, sex and tissue of the

91 samples included in our assembly supply the equine transcriptome with its best spatiotemporal

92 resolution and most complete gene UTR definition to date.

93 We recognize that availability of annotation criteria and integration of transcriptome data is

94 paramount for systematically improving the equine transcriptome. Our goal is to encourage equine

95   researchers to incorporate their transcriptomic data using our pipeline as the common annotation

96   pipeline and our initial transcriptomes as a reference framework. We intend to continue improving

97   equine gene annotation through better UTR definition, isoform splicing characterization and novel gene

98   identification. The annotation presented in this paper will improve the gene structure definition in

99   current databases and the accuracy of downstream analyses, including both differential gene expression

100   analysis and genetic variant annotation in the horse.

101

102

103   **Results**

104

105   *Overall Mapping Statistics and Gene Counts After Filtration*

106

107         RNA-seq of 59 samples in 12 libraries from 8 different horse tissues provided 1.9 billion

108   fragments and 364 Gb of sequence bases. A summary of the library preparation, number of horses per

109   library and total number of fragments and bases provided by each tissue library can be found in Table 1.

110   The overall average mapping rate for Tophat2 was ~83% with concordance rates ranging from 29% to

111   89%  (average 75%) for paired end libraries. Concordance rates seem to be affected by the type of

112   library preparation, where polyA selected and strand-specific libraries have the best rates.  Library

113   specific mapping rates can be found in Supplementary Table 1.  The initial Cufflinks assembly identified

114   117,019 genes/211,562 transcripts. After this initial analysis we applied four steps of filtration (Figure

115   1). Primary filtration of transcripts removed the likely pre-mRNA fragments by eliminating single exon

116   transcripts that were present within introns or overlapping with exons of other multi-exon transcripts.

117   After primary filtration there were 75,102 genes/162,261 transcripts. The second filter was implemented

118    to remove isoforms likely to be experimental artifacts by excluding low abundant transcripts with less

119    than 5% of total expression for their locus. The remaining 114,830 transcripts represented 75,375 genes.

120    In the third filter, non-coding transcripts that lack any supporting evidence from NCBI or ENSEMBL

121    annotations, non-horse gene models ("Other RefSeq" and "TransMap RefGene" UCSC tracks) or *ab*

122    *initio* predictions ("Augustus", "Geneid", "Genscan" and "N-SCAN" UCSC gene prediction tracks)

123    were excluded. This third filtered version of the transcriptome has 76,323 transcripts in 37,062 genes.

124    The last filter was for removing likely erroneous transcripts. The mtDNA in mammals is known for gene

125    overlapping and polycistronic expression [11], permitting inaccurate prediction of mitochondrial

126    transcripts by Cufflinks; we therefore excluded the mitochondrial contigs from our filtered assembly.

127    Also, short transcripts less than 201bp (192 transcripts in 184 genes) were removed because they are

128    more likely to represent repetitive sequences or incomplete gene fragments. Once erroneous transcripts

129    were removed, our final refined version of the transcriptome contained 36,876 genes (76,125 transcripts)

130    including 15,343 single exon transcripts, 8,808 two-exon transcripts, and 51,974 transcripts with three or

131    more exons.  A version of our refined transcriptome that is merged with the NCBI and ENSEMBL

132    annotations, with redundant transcripts removed, is also available. This is the most comprehensive

133    product of our pipeline and is valuable for differential gene expression analysis in tissues other than

134    those provided in our assembly.  Summary statistics including N50, number of genes and Mb and

135    average length of fragment for all six versions of the transcriptome can be found in Supplementary

136    Table 2.

137    *Comparison between Our Transcriptome and Currently Available Equine Transcriptomes*

138

139       We performed a comparison between our transcriptome and gene models from NCBI,

140    ENSEMBL and two published equine transcriptomes, that we refer to as Hestand [7] and ISME [5]

141 (Table 2 and Supplementary Table 3). In our comparisons, transcripts sharing one or more splice

142 junctions are considered similar but only those with identical intron chains are matching. The

143 comparison shows that the matching transcripts between our refined transcriptome and NCBI annotation

144 are greater than 2.5-fold in number those matching the ENSEMBL annotation. However the highest

145 number of matching transcripts occurred with the ISME transcriptome with 12,849 transcripts (Figure

146 2A). About 50% of the refined transcripts have a similar match in all the public transcriptomes.

147 Evidence of improvements to the annotation of genes with a similar match to other assemblies can be

148 found in genes such as *MUTYH*, where the three major isoforms annotated in humans [12] are now

149 distinguishable in the horse (Figure 2B). The gene *CYP7A1* is another example where a novel first exon

150 has been annotated and extended in our version of the transcriptome [13] (Figure 2C). About 20% and

151 28% of the refined transcripts are novel when compared to NCBI and ENSEMBL annotations

152 respectively. Combined, there are 22,641 transcripts in candidate novel loci. Our approach of applying

153 four successive steps of filtration strictly qualifies our novel isoforms as transcripts with ORFs or exonic

154 overlap with candidate gene models. Mainly, novel transcripts contained within introns of other genes

155 were excluded to avoid artifacts from retained intronic reads, common in rRNA depleted libraries. Using

156 NCBI as a reference for comparison, our novel transcripts from the refined transcriptome have no bias

157 towards any particular chromosome after accounting for chromosome size (Supplementary Figure 1). In

158 order to calculate the gene and isoform detectability of our transcriptome compared to current

159 annotation, we calculated sensitivity and specificity [14] between our transcriptome and a reference and

160 found that, using NCBI as the reference, our transcriptome had a 78.8% sensitivity and 23.8%

161 specificity at the base level and a 32% sensitivity and 21.1% specificity at the locus level. Detailed

162 pairwise assessment for all equine annotations can be found in Supplementary Table 4. We developed a

163 statistic to assess the conflict between different assemblies, termed "complex loci", which refer to the

164    loci that represent one gene locus in one transcriptome and two or more gene loci in another. Our

165    transcriptome has 1355 and 997 transcripts that were considered complex loci between our

166    transcriptome and NCBI and ENSEMBL, respectively. The Hestand transcriptome, however, has fewer

167    with 660 and 798 complex loci when compared to NCBI and ENSEMBL, respectively. The ISME

168    transcriptome has substantially more, with 1546 and 1226 complex loci when compared to NCBI and

169    ENSEMBL, respectively.

170

171    *UTR extension*

172

173    To test the effect of the new assembly on the UTRs of known genes, we identified the protein

174    coding isoforms sharing the exact intron chain with NCBI isoforms, which yielded 9736 isoforms from

175    7419 genes.   The difference in the total length of each transcript was then calculated and we found that

176    we extended the length of 8899 isoforms (6817 genes) by 29.7 Mb in total. 831 isoforms (718 genes)

177    lost 0.3 Mb in total with an average of 0.4 kb per isoform, while 6 isoforms did not change.

178

179    *Gene and Isoform Distinctions between Tissue-Specific Transcriptomes*

180

181    We selected genes with high expression (a sum of TPMs across all tissues above 200) and

182    substantial expression differences across tissues (a standard deviation above 200). Unsupervised

183    hierarchical clustering grouped genes that may be co-expressed as well as illustrating the relationship

184    between the tissue-specific transcriptomes.  As expected, the transcriptomes from the three central

185    nervous system (CNS) tissues clustered together, as did the two embryonic tissues, with the skin and

186    skeletal muscle furthest from these clusters (Figure 3A). Blocks of genes showing uniquely high

187   expression in a given tissue were further annotated with NCBI gene names and then summarized with

188   Panther biological processes annotations.  The top two Panther pathways (lowest p-values) for each of

189   these gene blocks are reported in the text below, with the corresponding p-values (p) and fold-

190   enrichment vales (FE).  The full Panther annotation tables are detailed in Supplementary Table 5.  The

191   CNS cluster contained overrepresented processes regarding brain function and development: nervous

192   system development (p=2.10E-06, FE=7.12) and neurogenesis (p=9.36E-05, FE=7.73).  The retina

193   contained processes consisting of photoreception and visual perception: phototransduction (p=3.52E-08,

194   FE=80.75) and visual perception (p=3.69E-18, FE=37.64).  The skeletal muscle encompassed genes

195   pertaining to muscle physiology and assembly: muscle contraction (p=1.48E-27, FE =57.58) and

196   myofibril assembly (p=6.05E-11, FE=72.8). The embryonic tissues have the most general processes

197   assigned to their distinct clusters: translation (p=1.15E-11, FE=16.35) and peptide biosynthetic

198   processes (p=1.95E-11, FE=15.8).  And finally, the skin consisted of processes concerning epithelial

199   organization and production: intermediate filament organization (p=1.69E-07, FE > 100) and skin

200   development (p=1.89E-09, FE=22.49).

201         When attention is given to the isoforms showing unique presence or sole absence in a tissue, the

202   retina and cerebellum possess the most isoforms that are uniquely present, with the retina also

203   containing the largest amount of solely absent isoforms (Figure 3B). The uniquely present transcripts in

204   the retina have Panther annotations of visual perception (p=2.96E-23, FE=7.47) and photoreceptor cell

205   differentiation (p=1.59E-09, FE=7.56) and in the cerebellum nervous system development (p=3.08E-09,

206   FE=1.85) and generation of neurons (p=1.72E-08, FE=2.02). Full Panther annotation tables of

207   transcripts uniquely present in tissues can be found in Supplementary Table 6. The transcripts solely

208   absent from the retina pertain mainly to positive regulation of DNA replication (p=2.49E-03, FE=3.53)

209   and anatomical structure development (p=2.72E-19, FE=1.48).  Full Panther annotation tables of

210   transcripts solely absent in tissues can be found in Supplementary Table 7. Utility of these isoforms, in

211   terms of expression, is strongest in the skin, retina, skeletal muscle and to a small extent the cerebellum

212   (Figure 3B). Despite these differences in unique isoforms, multi-exons transcripts and multi-transcript

213   loci, the splicing rate across tissues, as calculated by Cuffcompare [15], ranges from 1.7 to 1.9 (Table 3).

214          Nuclear coding versus mitochondrial encoded genes were parsed out per tissue to determine how

215   much of the sequencing resources are allocated to genes of the mitochondria (Figure 3C), with the

216   conclusion that the brainstem, spinal cord, embryonic tissues and skeletal muscle exhibit the largest

217   proportions of transcriptional output devoted to mitochondrial genes.

218

219   *Classification and Annotation of Novel Genes*

220

221          In total there were 22,640 novel transcripts, with varying levels of support from current equine

222   annotations (Figure 4A). Classification of these candidate novel genes was necessary to better describe

223   our novel gene identification and aid in interpreting the degree of support available to each category.

224   Three categories of novel genes based upon the supportive evidence within and across species were

225   made with each successive category being less supported by equine or orthologous gene models. Our

226   first category of novel genes contains those missing from NCBI and/or ENSEMBL annotations, but

227   supported by either NCBI, ENSEMBLE, Hestand or ISME annotations (Category I). The second group

228   of novel genes were novel to all public equine annotations, but conserved by means of orthologous gene

229   similarity or supported by possible gene prediction (Category II). The third category of novel genes

230   were unsupported by any candidate gene models, but had an ORF (Category III). Category I has a total

231   of 8459 transcripts, with 2/3 of these transcripts novel to the ENSEMBL annotation (Supplementary

232   Figure 2). Another 1849 transcripts in this category are absent in both NCBI and ENSEMBL

233     annotations, yet supported by Hestand or ISME annotations. Homology with the SWISS-prot database

234     identified at least one significant (p<1E-10) hit for almost half the transcripts in this Category I

235     (Supplementary Table 8).  The second category has 7494 transcripts that – unless on the opposite strand

236     - do not overlap with known gene models in public annotations. Annotation of these transcripts was

237     performed partially by testing overlap with non-horse gene models and also by homology search. Only

238     16% of these transcripts have significant hits against the SWISS-prot database (Supplementary Table 8).

239     The third category of novel genes includes 6687 transcripts with an ORF as the only functional support

240     for these transcripts. The first category of novel genes shows the most diverse distribution of exon

241     numbers comprising the genes (ranging from 1 to 28), whereas the unsupported genes contained mainly

242     single exon genes (Figure 4B).  The expression analysis of the three novel gene categories shows a clear

243     reduction of cumulative expression from category I to III.  There was also an obvious tissue-specific

244     pattern in the expression of novel genes.  Supported novel genes (Category I) had the highest expression

245     in the cerebellum.  However, when looking at only the second category of novel genes, the embryo

246     contributed the highest expression of novel transcripts. Category III novel transcripts mainly consisted

247     of single exon transcripts and showed similarly low expression across all tissues (Figure 4C).

248

249     **Discussion**

250         Using RNA-seq from fifty-nine horses across eight tissues has allowed us to capture

251     transcriptome complexity and provide spatial resolution in terms of tissue-specificity in manner that

252     exceeds any current equine annotations.  Our descriptive statistics and accessible pipeline make this

253     project open to modifications and further integration of transcriptomes.

254         RNA-seq based transcriptomes are prone to false inflation of gene numbers for several reasons.

255     Technical limitations such as limited sequencing read length and amplification errors, false splicing

256    events, and assembler deficiencies are among several reasons of misassembly. Pervasive transcription is

257    another predominant source of such inflation [16-18].  Some types of sequencing libraries increase the

258    problem as well; for example rRNA depletion inflates the assembly with primary transcripts and false

259    isoforms exhibiting intronic retention [19]. Our pipeline takes these factors into account and runs

260    unguided by previous transcriptome annotations with several transcript filtration steps that reduce

261    inclusion of inaccurate transcripts, while retaining the sensitivity for novel transcript detection. The

262    effect of this procedure can be seen by comparing the gene numbers between our initial unfiltered and

263    final filtered transcriptomes, where gene inflation was reduced by 68% (Table 2) and our final refined

264    transcriptome contained 36,876 genes and 76,125 isoforms.

265         Although not indicative of transcriptome quality, we calculated specificity, as a measure of

266    difference between our transcriptome and other annotations, and sensitivity, which indicates how our

267    transcriptome covers another annotation.  These parameters demonstrate that our aggressive filtering

268    does sacrifice sensitivity at the locus level only by a margin of approximately 5%, and increases our

269    specificity often by more than 10%, relative to NCBI and ENSEMBL (Supplementary Table 4).  We

270    have a comparable sensitivity to the Hestand transcriptome, which could be explained by adopting strict

271    filtering approaches in both pipelines. However, the numbers of unstranded and multi-exon transcripts in

272    the Hestand transcriptome relative to our refined version serve as the more discriminating statistics. We

273    have approximately six fold less unstranded transcripts and more than double the multi-exon transcripts

274    (Table 2).  Regarding how our transcriptome compares to the other recent equine ISME assembly [5],

275    which is ENSEMBL annotation guided, we have three times more matching transcripts to the ISME

276    assembly than to the ENSEMBL annotation itself (Figure 2A), suggesting significant improvement

277    made by the ISME annotation.  However their improvements are impaired by false inflation in the

278    number of genes identified due to presenting most of the transcripts in two copies representing the

279    forward and reverse strands.  This inflation of the ISME annotation can explain why it has different

280    statistics from Hestand as well as our new transcriptome.  Hestand et al (2015) also observed a bias

281    towards single exon genes, which represented approximately 55% of their whole transcriptome [7]. Our

282    native assembly identified similar percentage of single exon transcripts, however those numbers went

283    down to 20% after filtration of single exon pre-mRNA.   Our statistic, complex loci, also highlights a

284    level of sensitivity as well as an area for further investigation in our transcriptome.  We have more than

285    two times more complex loci, using NCBI as a reference, than Hestand. The inflated ISME complex loci

286    numbers could be attributed to the double reporting of their transcripts.  Awareness of these complex

287    loci allows for refinement of transcriptome-wide gene structure, while a pipeline to appropriately

288    process these loci has yet to be established. The evaluation of these alternative descriptive statistics

289    permits comparisons with transcriptomes that have pipeline-specific limitations.

290         Accurate identification of UTRs is often difficult for *ab initio* programs and requires sufficient

291    support of transcription evidence. Our integrative analysis of several tissues using different library preps

292    enabled us to achieve unpreceded extension of equine transcripts' UTRs by an average of 3.3 kb per

293    transcript. Indeed high coverage of CNS tissues in our analysis was an important factor as reported

294    previously with several transcriptomes [20], [21] and [22]. Further improvements to this transcriptome

295    would include providing tissue-specific UTR lengths and allowing for a more clear depiction of

296    differences in gene structure between tissues.  The improved UTR structure provided by our

297    transcriptome has already shown its utility in the horse community by defining isoform and gene

298    boundaries of *MUTYH* and *TOE1* [23] as well as providing an alternative start exon for *CYP7A1* [13]

299         The distinct RNA-seq libraries from eight tissues,allow us to extract tissue-specific features

300    regarding gene expression, isoform usage and mitochondrial gene expression.  Tissue specificity, in

301    terms of gene expression, was demonstrated on two levels, first as biases in exon usage, especially one

302    and two exon gene expression in the embryo (Supplementary Table 3). Second, gene clustering and

303    Panther annotation revealed inherent functions associated with each tissue (Figure 3A, Supplementary

304    Table 5). The retina displayed the most uniquely present and absent isoforms, in agreement with other

305    studies [25], and had correlated isoform annotations unique to the retina: visual perception and

306    photoreceptor cell differentiation. The skeletal muscle was a tissue with a relatively low amount of

307    unique tissue-specific isoforms, however, it shows utility of these isoforms with relatively high

308    cumulative expression values (Figure 3B) as well as comparable splicing rates (Table 3). Three tissues:

309    retina, skin, and embryo, had shorter read lengths and were not prepared as stranded libraries and thus

310    these data may be artificially understated in terms of transcriptome complexity. In addition to the

311    nuclear gene expression, the amount of transcription occurring from the mitochondrial chromosome can

312    show how much of the sequencing resources are being allocated to genes of mitochondrial origin.

313    Across the eight tissues, one would expect the tissue with the largest numbers of mitochondria to have

314    the largest proportion of mitochondrial allocated transcriptional output [2]. Our data demonstrates these

315    trends in the brainstem, skeletal muscle and spinal cord, however the cerebellum and retina do show an

316    unexpectedly low mitochondrial gene load (Figure 3C). Further research establishing the relationship

317    between the amount of mitochondria processed in a sample for RNA-seq and the resulting mitochondrial

318    expression loads would be beneficial to understanding how much of the transcriptional output is

319    dominated by an individual mitochondrion.

320       We identified 7494 candidate novel transcripts. These novel transcripts are selected based on

321    having no overlap with genes in current equine annotation and authenticated by their protein-coding

322    ability and/or overlap with aligned non-horse genes or *ab initio* gene predictions. Our novel transcripts

323    have a diversity of coding exons in Category I and a particular expression bias in the embryonic tissues

324    of Category II, in which a majority of these novel transcripts contain two exons (Figure 4C). The

325    Category I novel transcripts highlight the deficient equine ENSEMBL annotation, the need to pool the

326    databases to get the most transcriptome coverage and the ability of our transcriptome to capture the

327    potentially rare novel gene models (Figure 4A). Despite the ORF requirement for Category III novel

328    transcripts, there is an obvious enrichment for single exon transcripts and a marked reduction of total

329    transcription level (Figure 4C), which is indicative of non-coding RNA. Our novel gene analysis also

330    produced a category of novel transcripts that were removed due to not having ORFs and were presumed

331    to represent noisy transcription relating to primary transcripts, repetitive elements, sequencing errors and

332    genome-based errors. The collection of Category III and these excluded novel transcripts may represent

333    a repository of non-annotated non-coding RNA, which is an area that needs further annotation in the

334    horse genome.

335        This transcriptome assembly pipeline not only produces flexible incorporation of additional

336    transcriptomes, it also provides several products regarding levels of transcript filtering and

337    appropriateness for downstream analysis. The different extractable transcriptome versions include

338    transcriptomes after each individual filter, with the final refined transcriptome containing only genes

339    with complete ORFs and genes aligning with other non-horse genes or *ab initio* gene predictions. A

340    version of our transcriptome merged with NCBI and ENSEMBL annotations achieves breadth not

341    covered by our tissues. These transcriptomes as well as the pipeline to make each of these

342    transcriptomes are publically available on our GitHub repository. By making the workflow public and

343    easy to execute and manipulate, we aim to expand the spectrum of tissues embodying this transcriptome

344    and eliminate biases in annotated genes and thus downstream differential gene analysis. Increasing and

345    providing the option for tissue-specific transcriptomes, allows for more targeted and refined usage of a

346    certain tissue's transcriptome for differential gene analysis, resulting in downstream analysis of more

347    significant differentially expressed genes. As stated in our overall goals of this project, we have

348    provided a framework for further improving the equine transcriptome and produced an equine

349    transcriptome that expands on current equine annotations in the manner of UTR extension, isoform

350    detection and novel gene identification.

351

352    **Methods**

353

354    *RNA-seq library preparation*

355

356        A total of twelve RNA-seq libraries in 8 tissues from 59 individuals (20 female, 27 males and 12

357    embryos) were used to prepare our transcriptome. The brainstem, spinal cord and cerebellum were

358    strand-specific 100 bp paired-end (PE) libraries.  The skeletal muscle tissues were strand-specific PE125

359    bp libraries.  A subset of the embryo ICM (3 samples) and TE (3 samples) were unstranded PE100 bp

360    libraries, the other subset (3 ICM and 3 TE) consisted of single end (SE) 100 bp reads.  The retina RNA-

361    seq libraries were unstranded SE80 bp libraries. And the skin libraries were all unstranded and consisted

362    of PE80 bp, SE80bp and SE95 bp reads.  The brainstem, spinal cord and cerebellum RNA libraries were

363    all rRNA-depleted, the skin, retina and skeletal muscle libraries were poly-A captured.  The embryonic

364    libraries were neither poly-A selected nor rRNA depleted, they were prepared with the Ovation® RNA-

365    seq System V2 (NuGEN, San Carlos, CA, USA), which aims to amplify mRNA as well as non-

366    polyadenylated transcripts. Table 1 summarizes the tissue-specific RNA-seq library parameters.

367

368

369    *Trimming and Mapping of Reads*

370

371     The Illumina adaptors as well as the reads were trimmed with the sliding window quality

372     trimmer Trimmomatic [26] with a window size of 4 and a softer quality threshold of 2 [27]. Mapping of

373     the trimmed reads was done with Tophat2 [28] to EquCab2.0, 2007

374     (ftp://hgdownload.cse.ucsc.edu/goldenPath/equCab2). Cufflinks [15] was used to assemble transcripts

375     from the aligned RNA-Seq reads. Two cerebellar samples failed assembly due to computation

376     limitations (8 CPUs, 250 Gb RAM and 7 days) and required digital normalization [29] to 200X coverage

377     before mapping with Tophat2.

378

379

380     *Filtering Transcripts*

381

382     Four categories of filters were used to remove likely pre-mRNA and artifactual transfrags, as

383     summarized in Figure 1, resulting in six versions of the transcriptome. Primary transcript filtration was

384     done using Cuffcompare [15] between our assembly and a version of our assembly containing only

385     multi-exon transcripts and removing transcripts overlapping with intronic regions and with class codes

386     "i", "e" and "o". The input trimmed RNA-seq reads were then back-mapped to the pre-mRNA free

387     transcriptome using the quasi-mapping based software package Salmon [30]. While back-mapping, a

388     second filtration step was implemented: low abundance transcripts in every locus were excluded with

389     the lower threshold of a TPM (normalized read count standing for transcripts per million) less than 5%

390     of the total TPM per locus. For the third filter, Transdecoder [31] was used to predict the ORFs and

391     Cuffcompare [15] to determine any exonic overlap with any candidate gene locus using the class codes

392     "j", "o", "x" and "c". In the Transdecoder analysis, the longest open reading frames were extracted as

393     well as any sequences having significant homology to the Pfam and Swissprot protein databases.

394  Finally, the removal of likely erroneous mitochondrial and short transcripts was done by a homemade

395  script.

396

397  *Transcriptome Comparisons*

398

399       Comparisons of our refined transcriptome to the four public horse transcriptomes were done

400  using Cuffcompare [15]. In any pairwise comparison, two transcripts are considered matching if they

401  have the exact intron chain, despite differing terminal exons (class code "="). If the transcripts are not

402  matching but sharing one or more splice junctions (class code "j"), these would be considered similar

403  transcripts. A transcript is considered novel if it does not overlap with any gene model in the 2$^{nd}$

404  reference assembly (class code "u"). All other class codes including any kind of overlap with a reference

405  annotation on the opposite strand were considered as "other". For more detailed descriptions of the class

406  codes provided by Cuffcompare, please see their manual [15]. Complex loci were flagged if a gene

407  model of one assembly overlapped with 2 different gene models in the other assembly. Sensitivity and

408  specificity relative to a given reference transcriptome were calculated per base, intron and locus for each

409  transcriptome and reference combination as described by Burset and Guigó [14].

410

411  *Novel Gene Prediction*

412

413       Any transcript in our final *refined* transcriptome is defined as novel if it does not overlap with a

414  gene model in at least one of the two public equine assemblies, NCBI and ENSEMBL (Cuffcompare

415  class code "u"). Transcripts considered novel were divided into three groups according to the degree of

416  supportive evidence. Transcripts novel to either the NCBI or ENSEMBL assemblies with transcriptional

417   supportive evidence from the other or any other public assembly [5, 7] were in the first category of

418   novel transcripts. Supportive evidence is defined as any overlapping with exon sequence (Cuffcompare

419   class code "=","j","o","x" or "c"). The second and third categories of novel transcripts required that the

420   transcript be absent in all current equine transcriptomes. Transcripts in the second category have

421   supportive evidences in non-horse alignment gene models or *ab initio* gene prediction tracks from the

422   UCSC genome browser.  The third category of novel transcripts included transcripts that lack such

423   evidence but have ORFs.

424

425   *Tissue-Specific Characterization of The Transcriptome*

426

427        Tissue-specific transcriptomes were generated by back-mapping the input trimmed RNA-seq

428   reads with Salmon [30] to the refined version of the transcriptome to obtain expression information on a

429   tissue-specific level. A transcript is considered expressed in a given tissue if it has a TPM more than 5%

430   of the total TPM per locus calculated from the tissue specific libraries only.  Biological processes

431   identified within the tissue-specific gene blocks were annotated with Panther [32] and reported if the p-

432   values were below the Bonferroni-corrected threshold (5% experiment-wide).

433

434   *UCSC Track hubs*

435

436        Gene Annotation Format (GTF) files were converted into the binary bigbed files [33] using

437   UCSC kintUtils (https://github.com/ENCODE-DCC/kentUtils). The track hub directory structure was

438   designed as recommend by UCSC genome browser [34]. Tracks were constructed using "bigBed 12"

439   format and multiple libraries of the same tissue were organized in composite tracks. The hub files are

440    hosted on a github server as a part of the horse_trans repository (https://github.com/dib-lab/horse_trans).

441

442    **List of Abbreviations:**

443    Single Nucleotide polymorphism (SNP)

444    Untranslated region (UTR)

445    Central nervous system (CNS)

446    Inner cell mass (ICM)

447    Trophectoderm (TE)

448    RNA sequencing (RNA-seq)

449

450

451    **Declarations**

452    *Ethics approval and consent to participate*

453    Not applicable

454    *Consent for publication*

455    All consent for publication is available

456    *Availability of data and materials*

457    The data including the scripts for the pipeline as well as the GTF files for the transcriptome can be found

458    at: https://github.com/dib-lab/horse_trans.  (This repository is archived by Zenodo at

459    10.5281/zenodo.56934).  All sequencing reads used in this study have been submitted to NCBI

460    Sequence Read Archive; SRA SRP082284 for muscle samples, SRP073514 for brainstem, SRP073514

461    and SRP082291 for spinal cord, SRP082342 for cerebellum, ERP001525 for retina, SRP031504 and

462    SRP082454 for embryonic tissues and ERP001524, ERP001525 and ERP005568 for skin.

463

464 *Competing Interests*

465      The authors declare that there are no competing interests.

466 *Funding*

470 *Authors' Contributions*

471      All authors contributed to study conception and design. TAM produced the annotation pipeline

472 and did the data analysis with oversight on writing of the manuscript as well as figure formation.  EYS

473 wrote the manuscript and made the figures.  CTB, CJF and JDM aided with experimental design and

474 supervised the whole project.   RRB and MJM provided retina and skin data.  SJV provided muscle data.

475 PJR provided both embryonic tissue data.  CJF provided spinal cord and brainstem data.  EYS, JDM and

476 MCP provided the cerebellar RNA-seq data.  All authors reviewed and approved the final  version of the

477 manuscript.

478 *Acknowledgements*

486
487     **References**

488
489    1.     Brown JB, Boley N, Eisman R, May GE, Stoiber MH, Duff MO, Booth BW, Wen J, Park S,
490        Suzuki AM *et al*: **Diversity and dynamics of the Drosophila transcriptome**. *Nature* 2014,
491        **512**(7515):393-399.
492    2.     Mele M, Ferreira PG, Reverter F, DeLuca DS, Monlong J, Sammeth M, Young TR, Goldmann
493        JM, Pervouchine DD, Sullivan TJ *et al*: **Human genomics. The human transcriptome**
494        **across tissues and individuals**. *Science* 2015, **348**(6235):660-665.
495    3.     Okazaki Y, Furuno M, Kasukawa T, Adachi J, Bono H, Kondo S, Nikaido I, Osato N, Saito R,
496        Suzuki H *et al*: **Analysis of the mouse transcriptome based on functional annotation of**
497        **60,770 full-length cDNAs**. *Nature* 2002, **420**(6915):563-573.
498    4.     Holl HM, Gao S, Fei Z, Andrews C, Brooks SA: **Generation of a de novo transcriptome**
499        **from equine lamellar tissue**. *BMC Genomics* 2015, **16**(1):739.
500    5.     Pacholewska A, Drogemuller M, Klukowska-Rotzler J, Lanz S, Hamza E, Dermitzakis ET,
501        Marti E, Gerber V, Leeb T, Jagannathan V: **The transcriptome of equine peripheral blood**
502        **mononuclear cells**. *PLoS One* 2015, **10**(3):e0122011.
503    6.     Coleman SJ, Zeng Z, Wang K, Luo S, Khrebtukova I, Mienaltowski MJ, Schroth GP, Liu J,
504        MacLeod JN: **Structural annotation of equine protein-coding genes determined by**
505        **mRNA sequencing**. *Anim Genet* 2010, **41 Suppl 2**:121-130.
506    7.     Hestand MS, Kalbfleisch TS, Coleman SJ, Zeng Z, Liu J, Orlando L, MacLeod JN: **Annotation**
507        **of the Protein Coding Regions of the Equine Genome**. *PLoS One* 2015, **10**(6):e0124375.
508    8.     Yandell M, Ence D: **A beginner's guide to eukaryotic genome annotation**. *Nat Rev Genet*
509        2012, **13**(5):329-342.
510    9.     Curwen V, Eyras E, Andrews TD, Clarke L, Mongin E, Searle SM, Clamp M: **The Ensembl**
511        **automatic gene annotation system**. *Genome Res* 2004, **14**(5):942-950.
512   10.     Kitts P: **The NCBI Handbook: Chapter 14. Genome Assembly and Annotation Process**.
513        In: *National Center for Biotechnology*. 2003 edn: McEntyre, J. & Ostell, J.; 2003.
514   11.     Taanman JW: **The mitochondrial genome: structure, transcription, translation and**
515        **replication**. *Biochim Biophys Acta* 1999, **1410**(2):103-123.
516   12.     Plotz G, Casper M, Raedle J, Hinrichsen I, Heckel V, Brieger A, Trojan J, Zeuzem S: **MUTYH**
517        **gene expression and alternative splicing in controls and polyposis patients**. *Hum*
518        *Mutat* 2012, **33**(7):1067-1074.
519   13.     Finno CJ, Bordbari M, Monsour T, Bannasch DL, Mickelson DB, Valberg SJ: **Spinal cord**
520        **transcriptome profiling of equine vitamin E deficient neuroaxonal dystrophy**
521        **identifies dysregulation of cholesterol homeostasis with upregulation of liver X**
522        **receptor target genes.** *In Review, Free Rad Biol Med* 2016.
523   14.     Burset M, Guigo R: **Evaluation of gene structure prediction programs**. *Genomics* 1996,
524        **34**(3):353-367.

525 15.   Trapnell C, Hendrickson DG, Sauvageau M, Goff L, Rinn JL, Pachter L: **Differential analysis**
526       **of gene regulation at transcript resolution with RNA-seq**. *Nat Biotechnol* 2013,
527       **31**(1):46-53.
528 16.   Bertone P, Stolc V, Royce TE, Rozowsky JS, Urban AE, Zhu X, Rinn JL, Tongprasit W, Samanta
529       M, Weissman S *et al*: **Global identification of human transcribed sequences with**
530       **genome tiling arrays**. *Science* 2004, **306**(5705):2242-2246.
531 17.   Khaitovich P, Kelso J, Franz H, Visagie J, Giger T, Joerchel S, Petzold E, Green RE, Lachmann
532       M, Paabo S: **Functionality of intergenic transcription: an evolutionary comparison**.
533       *PLoS Genet* 2006, **2**(10):e171.
534 18.   Schadt EE, Edwards SW, GuhaThakurta D, Holder D, Ying L, Svetnik V, Leonardson A, Hart
535       KW, Russell A, Li G *et al*: **A comprehensive transcript index of the human genome**
536       **generated using microarrays and computational approaches**. *Genome Biol* 2004,
537       **5**(10):R73.
538 19.   Sultan M, Amstislavskiy V, Risch T, Schuette M, Dokel S, Ralser M, Balzereit D, Lehrach H,
539       Yaspo ML: **Influence of RNA extraction methods and library selection schemes on**
540       **RNA-seq data**. *BMC Genomics* 2014, **15**:675.
541 20.   Smibert P, Miura P, Westholm JO, Shenker S, May G, Duff MO, Zhang D, Eads BD, Carlson J,
542       Brown JB *et al*: **Global patterns of tissue-specific alternative polyadenylation in**
543       **Drosophila**. *Cell Rep* 2012, **1**(3):277-289.
544 21.   Mangone M, Manoharan AP, Thierry-Mieg D, Thierry-Mieg J, Han T, Mackowiak SD, Mis E,
545       Zegar C, Gutwein MR, Khivansara V *et al*: **The landscape of C. elegans 3'UTRs**. *Science*
546       2010, **329**(5990):432-435.
547 22.   Ulitsky I, Shkumatava A, Jan CH, Subtelny AO, Koppstein D, Bell GW, Sive H, Bartel DP:
548       **Extensive alternative polyadenylation during zebrafish development**. *Genome Res*
549       2012, **22**(10):2054-2066.
550 23.   Scott EY, Penedo MC, Murray JD, Finno CJ: **Defining Trends in Global Gene Expression in**
551       **Arabian Horses with Cerebellar Abiotrophy**. *In Review, Cerebellum* 2016.
552 24.   Wu JQ, Habegger L, Noisa P, Szekely A, Qiu C, Hutchison S, Raha D, Egholm M, Lin H,
553       Weissman S *et al*: **Dynamic transcriptomes during neural differentiation of human**
554       **embryonic stem cells revealed by short, long, and paired-end sequencing**. *Proc Natl*
555       *Acad Sci U S A* 2010, **107**(11):5254-5259.
556 25.   Farkas MH, Grant GR, White JA, Sousa ME, Consugar MB, Pierce EA: **Transcriptome**
557       **analyses of the human retina identify unprecedented transcript diversity and 3.5 Mb**
558       **of novel transcribed sequence via significant alternative splicing and novel genes**.
559       *BMC Genomics* 2013, **14**:486.
560 26.   Bolger AM, Lohse M, Usadel B: **Trimmomatic: a flexible trimmer for Illumina sequence**
561       **data**. *Bioinformatics* 2014, **30**(15):2114-2120.
562 27.   Macmanes MD: **On the optimal trimming of high-throughput mRNA sequence data**.
563       *Front Genet* 2014, **5**:13.
564 28.   Kim D, Pertea G, Trapnell C, Pimentel H, Kelley R, Salzberg SL: **TopHat2: accurate**
565       **alignment of transcriptomes in the presence of insertions, deletions and gene**
566       **fusions**. *Genome Biol* 2013, **14**(4):R36.
567 29.   Brown C, Howe, A., Zhang, Q., Pyrkosz, A.B. & Brom, T.H.: **A Reference-Free Algorithm for**
568       **Computational Normalization of Shotgun Sequencing Data**. 2013.
569 30.   Rob Patro GD, Carl Kingsford: **Accurate, fast, and model-aware transcript expression**
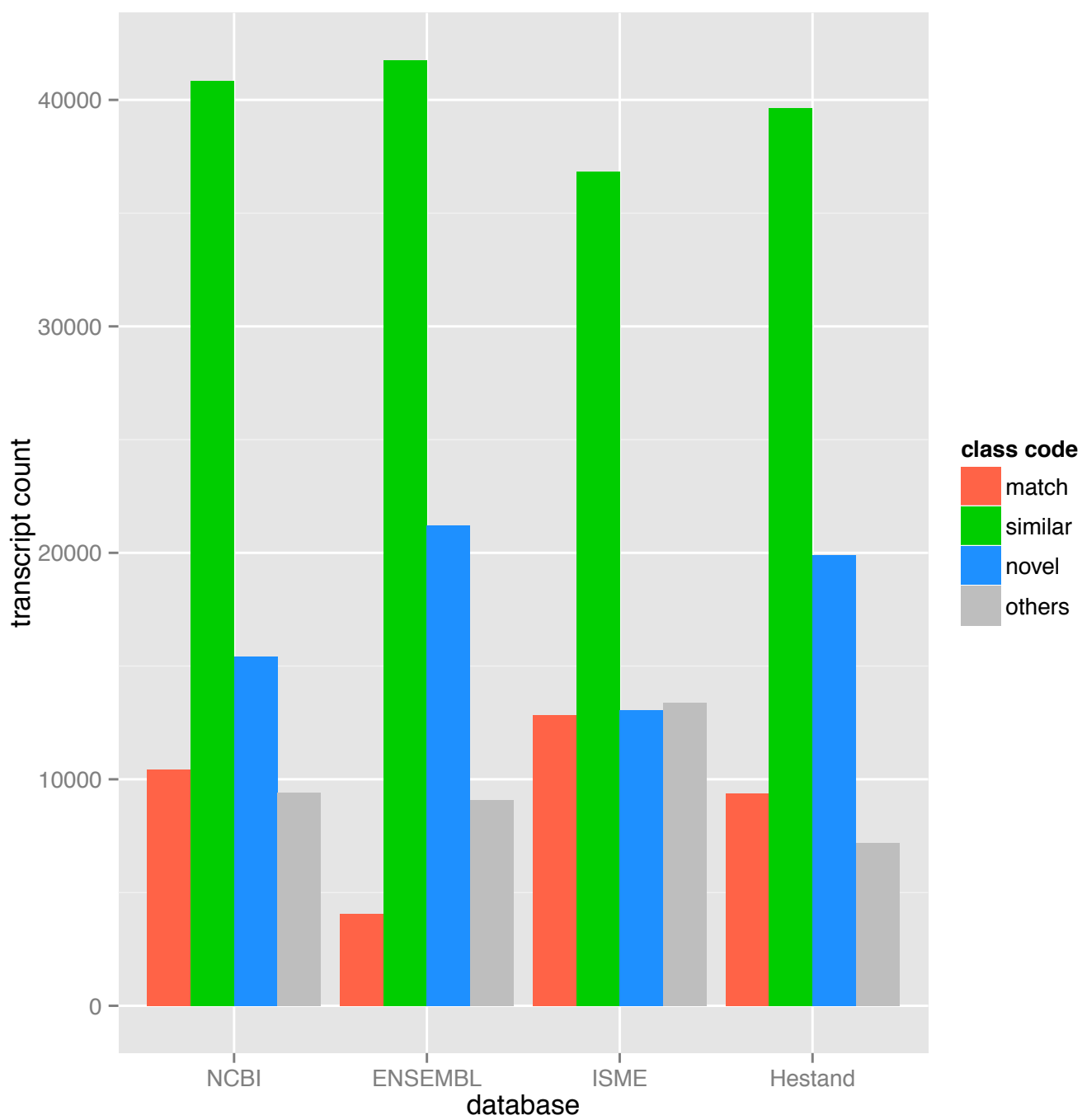570       **quanitification with Salmon**. *bioRxiv* 2015.

571  31.  Haas BJ, Papanicolaou A, Yassour M, Grabherr M, Blood PD, Bowden J, Couger MB, Eccles D,
572       Li B, Lieber M *et al*: **De novo transcript sequence reconstruction from RNA-seq using**
573       **the Trinity platform for reference generation and analysis**. *Nat Protoc* 2013,
574       **8**(8):1494-1512.
575  32.  Mi H, Muruganujan A, Casagrande JT, Thomas PD: **Large-scale gene function analysis**
576       **with the PANTHER classification system**. *Nat Protoc* 2013, **8**(8):1551-1566.
577  33.  Kent WJ, Zweig AS, Barber G, Hinrichs AS, Karolchik D: **BigWig and BigBed: enabling**
578       **browsing of large distributed datasets**. *Bioinformatics* 2010, **26**(17):2204-2207.
579  34.  Raney BJ, Dreszer TR, Barber GP, Clawson H, Fujita PA, Wang T, Nguyen N, Paten B, Zweig
580       AS, Karolchik D *et al*: **Track data hubs enable visualization of user-defined genome-**
581       **wide annotations on the UCSC Genome Browser**. *Bioinformatics* 2014, **30**(7):1003-1005.

582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616

**Figure 1**. An outline of the workflow used to generate each version of the transcriptome. Transcriptome products are in ovals. Programs used to perform various steps are indicated in parentheses. All transcriptome versions and the pipeline scripts are publically available.

A)

**Figure 2.** Comparison of our refined transcriptome to current equine annotations. (A) Our refined transcriptome compared to current annotations. (B) The annotation of *MUTYH* in the refined version of the transcriptome shows the addition of several isoforms, α, β, and γ, as seen in the human, of *MUTYH*. (C) The gene annotation of *CYP7A1* in the refined transcriptome also shows the inclusion of an extended alternative first exon not seen in other species.

A)

C)



**Figure 3**. Tissue-specific gene and isoform composition of the transcriptome. (A) A heatmap of genes with high expression and substantial expression differences across tissues. (B) A bar graph showing isoforms uniquely present (the bar outlined in red above the x-axis) or solely absent (the blue outlined bars extending below the x-axis). The green trendline corresponds to the cumulative TPM of the uniquely present transcripts. (C) A stacked bar graph showing the transcription percentage of mitochondrial genes versus nuclear encoded genes.

A)

B)

C)



**Figure 4**. Novel gene analysis and classification. (A) A bar graph showing the comparison of all the novel genes against the current equine annotations. The three categories of novel genes were supported novel genes (Category I), unsupported, but conserved, novel genes (Category II) and the unsupported, un-conserved, but novel genes with an ORF (Category III). (B) A stacked bar graph of transcript counts with all three categories of novel genes showing exonic composition and (C) their cumulative TPM in a tissue specific manner.

**Table 1.** Sample and library preparations used as input for our equine transcriptome.

| Tissue | Library Preparation | Library Characteristics | #Samples | #Frag(M) | #bp(Gb) | Reference |
|---|---|---|---|---|---|---|
| BrainStem | RiboRNA-depleted | PE100bp, stranded | 8* | 166.73 | 33.68 | Finno et al, 2016 |
| Cerebellum | RiboRNA-depleted | PE100bp, stranded | 12 | 411.48 | 82.3 | Scott et al, 2016 |
| Muscle | Poly-A capture | PE125bp, stranded | 12 | 301.94 | 76.08 | |
| Retina | Poly-A captured | PE80bp unstranded | 2 | 20.3 | 3.28 | Bellone et al, 2013 |
| SpinalCord | RiboRNA-depleted | PE100bp, stranded | 16* | 403 | 81.4 | Finno et al, 2016 |
| Skin | Poly-A captured | PE80bp, unstranded | 2 | 18.54 | 3 | Holl et al, 2016 |
| | Poly-A captured | SE80bp, unstranded | 2 | 16.57 | 1.34 | Holl et al, 2016 |
| | Poly-A captured | SE95bp unstranded | 3 | 105.51 | 10.02 | Bellone et al, 2013 |
| Embryo ICM | Ovation RNA-seq | PE100bp, unstranded | 3 | 126.32 | 25.26 | Iqbal et al, 2014 |
| | Ovation RNA-seq | SE100bp, unstranded | 3 | 115.21 | 11.52 | Iqbal et al, 2014 |
| Embryo TE | Ovation RNA-seq | PE100bp, unstranded | 3 | 129.84 | 25.96 | Iqbal et al, 2014 |
| | Ovation RNA-seq | SE100bp, unstranded | 3 | 102.26 | 10.23 | Iqbal et al, 2014 |
| Total | | | 69 | 1917.7 | 364.07 | |

*Seven individuals had both brainstem and spinal cord tissue collected from them. Seven of the skin samples were taken from 5 individuals and one individual had both retina and skin sampled, bringing our total number of individuals to 59.

**Table 2.** Comparison of current public equine annotations to six versions of our transcriptome (bolded and outline in red) in terms of gene numbers and composition

| | Unfiltered | Mature | High-exp | Supported | Refined | Merged | Hestand | ISME | NCBI | ENSEMBL |
|---|---|---|---|---|---|---|---|---|---|---|
| Genes (super-loci) | 117019 | 75102 | 75375 | 37062 | 36876 | 47760 | 56495 | 42654 | 24342 | 26962 |
| Transcripts | 211562 | 162261 | 114830 | 76323 | 76125 | 121997 | 68594 | 285538 | 43417 | 29196 |
| Multi-transcript loci | 17136 | 15430 | 14602 | 14511 | 14505 | 17835 | 8465 | 23833 | 7257 | 1592 |
| Multi-exon transcripts | 108985 | 108985 | 61570 | 60839 | 60782 | 97654 | 30949 | 259556 | 39272 | 19805 |
| Redundant transcripts | 0 | 0 | 0 | 0 | 0 | 2 | 3 | 12578 | 141 | 79 |
| Unstranded transcripts | 46928 | 35881 | 35872 | 6676 | 6618 | 5705 | 37673 | 4732 | 0 | 0 |
| Single exon transcripts | 102577 | 53276 | 53260 | 15484 | 15343 | 24341 | 37642 | 13404 | 3723 | 6862 |
| Two exons transcripts | 11114 | 11114 | 9449 | 8857 | 8808 | 11350 | 4410 | 13092 | 2425 | 2972 |
| Many exons transcripts | 97871 | 97871 | 52121 | 51982 | 51974 | 86306 | 26542 | 259042 | 37269 | 19362 |

**Table 3.** Tissue-specific splicing rate as calculated by Cuffcompare, with relevant number of multi-exonic transcripts and multi-transcript loci per tissue.

| | Embryo ICM | Embryo TE | Skin | Brainstem | Cerebellum | Retina | Spinal cord | Muscle |
|---|---|---|---|---|---|---|---|---|
| Genes | 33998 | 32050 | 30003 | 34792 | 36139 | 26733 | 34980 | 29549 |
| transcripts | 57400 | 54424 | 51995 | 62993 | 66364 | 47095 | 66001 | 52000 |
| multi-exon transcripts | 44069 | 42433 | 42432 | 49346 | 51640 | 39420 | 52175 | 42483 |
| multi-transcript loci | 11938 | 11461 | 11797 | 13066 | 13334 | 10866 | 13352 | 11560 |
| Splicing rate | 1.7 | 1.7 | 1.7 | 1.8 | 1.8 | 1.8 | 1.9 | 1.8 |