

**A peer-reviewed version of this preprint was published in PeerJ on 6 May 2014.**

[View the peer-reviewed version](https://doi.org/10.7717/peerj.382) (peerj.com/articles/382), which is the preferred citable publication unless you specifically need to cite this preprint.

Moreton J, Malla S, Aboobaker AA, Tarlinton RE, Emes RD. 2014. Characterisation of the horse transcriptome from immunologically active tissues. PeerJ 2:e382 <https://doi.org/10.7717/peerj.382>

1 **Title**

2 Characterisation of the horse transcriptome from immunologically active tissues

3  
4 **Authors**

5 J. Moreton<sup>1, 2, 3</sup>, S. Malla<sup>2</sup>, A. A. Aboobaker<sup>4</sup>, R. E. Tarlinton<sup>3</sup> and R. D. Emes<sup>1, 3</sup>

6  
7 1 Advanced Data Analysis Centre, University of Nottingham, Sutton Bonington Campus,  
8 Loughborough, Leicestershire, LE12 5RD, UK

9 2 Deep Seq, Centre for Genetics and Genomics, University of Nottingham, Queen's Medical  
10 Centre, NG7 2UH, UK

11 3 School of Veterinary Medicine and Science, University of Nottingham, Sutton Bonington  
12 Campus, Loughborough, Leicestershire, LE12 5RD, UK

13 4 Department of Zoology, University of Oxford, Oxford, OX1 3PS, UK

14  
15 Corresponding author:

16 Joanna Moreton, ADAC, School of Veterinary Medicine and Science, University of  
17 Nottingham, Sutton Bonington Campus, Loughborough, Leicestershire, LE12 5RD, UK

18 0115 951 6359

19 [Joanna.Moreton@nottingham.ac.uk](mailto:Joanna.Moreton@nottingham.ac.uk)

20  
21 **Abstract**

22 The immune system of the horse has not been well studied, despite the fact that the horse  
23 displays several features such as sensitivity to bacterial lipopolysaccharide that make them in  
24 many ways a more suitable model of some human disorders than the current rodent models.  
25 The difficulty of working with large animal models has however limited characterisation of  
26 gene expression in the horse immune system with current annotations for the equine genome  
27 restricted to predictions from other mammals and the few described horse proteins. This paper  
28 outlines sequencing of 184 million transcriptome short reads from immunologically active  
29 tissues of three horses including the genome reference "Twilight". In a comparison with the  
30 Ensembl horse genome annotation, we found 8,763 potentially novel isoforms.

## 51 Introduction

52 While no longer the principal means of transport in much of the world, the horse is still an  
53 economically important animal in agriculture, sport and gambling associated with horse  
54 racing. Individual stallions may be worth several millions of dollars and attract high stud fees  
55 creating considerable interest in the genetics of performance traits (Hill *et al.* 2010). The  
56 immune response of the horse has not been well characterised, largely due to the difficulties  
57 in working with large animals in experimental settings. There are however several  
58 components of the equine immune system that make them in many ways a better model of  
59 some human disorders than the current rodent models. These include, similarly to human, an  
60 exquisite sensitivity to the effects of lipopolysaccharide (LPS) with associated endotoxemia  
61 and sepsis (Bryant *et al.* 2007).

62 Due to a lack of expressed sequence tag (EST) data, the current annotation of the  
63 protein coding regions of the horse genome is largely derived from extrapolation from the  
64 genes of other species (Coleman *et al.* 2010). Several recent papers have outlined  
65 transcription profiles using digital gene analysis of a variety of horse tissues including muscle,  
66 leukocytes, cartilage, brain, reproductive tissue, embryos, sperm and blood (Coleman *et al.*  
67 2010; McGivney *et al.* 2010; SerTEYN *et al.* 2010; Park *et al.* 2012; Capomaccio *et al.* 2013;  
68 Das *et al.* 2013; Iqbal *et al.* 2014). Capomaccio *et al.* (2013) identified new putative non-  
69 coding sequences within intergenic and intronic regions whilst Das *et al.* (2013) suggested  
70 additions to the structural annotation of four sperm genes. Two of the other studies (Coleman  
71 *et al.* 2010; Park *et al.* 2012) detailed extensions to the annotated gene catalogue in the horse  
72 based on transcriptome analysis of quite differing tissue sets, methods and results to those  
73 used in this paper. They show that the actual expressed transcription profile only partially  
74 overlaps the annotated gene set. A direct comparison of our and these two studies is difficult  
75 due to the differing tissues, methodologies and the lack of available locations of the predicted  
76 novel genes from these studies.

77 This paper focuses on immunologically active tissues in the horse to further explore  
78 this issue. Uniquely we present data on the transcriptional profile from lymphocytes from  
79 Twilight, the animal that the published horse genome is derived from. Comparison of this  
80 animal with lymphocytes, core immunologically active tissues (lymph node and spleen) and  
81 other tissues (liver, kidney and jejunum) from two unrelated animals allows a unique insight  
82 into expression of genes with a functional role in the immune system.

## 84 Materials and methods

### 85 *Samples, library preparation and sequencing*

86 The methods are described fully in our previous work (Brown *et al.* 2012) but briefly, five  
87 tissue samples; kidney, jejunum, liver, spleen and mesenteric lymph node were collected from  
88 an aged gelding (castrated male horse) euthanised due to osteoarthritis. The tissue samples  
89 listed were collected from an animal euthanized for clinical reasons, by the veterinary  
90 surgeon, under the Veterinary Surgeons act of 1966. Full informed consent of the owner was  
91 obtained for use of the samples, taken from that animal post-mortem.

92 Lymphocytes isolated by Ficoll Paque (GE healthcare) from a healthy 11 year old  
93 welsh mountain pony gelding were kindly provided by Dr Julia Kydd (School of Veterinary  
94 Medicine and Science, University of Nottingham) under the Home Office and local Ethical  
95 Approval Committee (PPL 40/3354). RNA from lymphocytes isolated from a healthy  
96 Thoroughbred mare (DNA the horse genome is derived from) was kindly provided by Donald  
97 Miller (Baker Institute of Animal Health, Cornell University, USA). This horse was  
98 maintained at the Baker Institute for Animal Health, Cornell University, Ithaca, N.Y., USA.  
99 Animal care and research activities were performed in accordance with the guidelines set  
100 forth by the Institutional Animal Care and Use Committee of Cornell University, protocol #

101 1986-0216, approved until March 2013.

102 Sequencing was performed on a SOLiD 3 ABI sequencer generating 50bp reads  
103 according to the manufacturer's instructions. Read data are available at the EBI Sequence  
104 Read Archive (SRA) under the study accession number ERP001116.

105  
106 *Read trimming and alignment*

107 The horse genome assembly EquCab2 (Wade *et al.* 2009) was downloaded from Ensembl v71  
108 ([www.ensembl.org](http://www.ensembl.org)) and contained 26,991 genes and 29,196 transcripts. CLC Genomics  
109 Workbench version 6 (CLC Bio, Aarhus, Denmark, [www.clcbio.com](http://www.clcbio.com)) was used to apply  
110 quality, SOLiD adapter and Poly-N trimming to the read sequences (supplemental file 1). The  
111 limit for the removal of low quality sequences was set at 0.2 and a maximum of two  
112 ambiguous nucleotides were permitted in each sequence. Any reads less than 20bp were  
113 removed after trimming and the average read lengths were 47bp.

114 TopHat 2.0.9 (Trapnell *et al.* 2009) was used to align the reads to the repeat masked  
115 version of the horse genome (Ensembl v71) to enable non-redundant transcriptome analysis.  
116 TopHat first aligns non-spliced reads using Bowtie 1.0.0 (Langmead *et al.* 2009) then  
117 identifies splice junctions. Gapped alignments are then used by TopHat to map the reads not  
118 aligned by Bowtie. In order to utilise the splice sites in all samples, two iterations of TopHat  
119 alignments were carried out (Cabili *et al.* 2011). Firstly, the reads from each sample were  
120 aligned to the repeat-masked horse genome with default parameters. The splice sites  
121 ("junctions") were extracted from all of the output files and duplicates were removed leaving  
122 216,007 sites. These splice sites were pooled together with the non-redundant sites extracted  
123 from the Ensembl annotation yielding 399,264 non-redundant splice sites. Each of the  
124 samples were then re-aligned with TopHat using the pooled non-redundant splice sites file  
125 (with 'raw-juncs' and 'no-novel-juncs' parameters) to the repeat-masked genome.

126 TopHat was used for the read alignment because it is part of the Tuxedo suite and is  
127 therefore a natural input for the Cufflinks assembler (Trapnell *et al.* 2010). It is also the  
128 preferred aligner for Scripture (Guttman *et al.* 2010). Cufflinks and Scripture are described in  
129 the transcriptome assembly section.

130  
131 *De novo transcriptome assembly*

132 Each of the samples were assembled into separate transcriptomes using two different tools;  
133 Cufflinks v2.1.1 (Trapnell *et al.* 2010) and Scripture (Guttman *et al.* 2010) (beta2 version,  
134 December 2010). These tools use different approaches for transcript assembly. A minimal set  
135 of transcripts is assembled by Cufflinks using a probabilistic model. It performs a minimum  
136 cost maximum matching in bipartite graphs (Trapnell *et al.* 2010). Scripture however creates  
137 a connectivity graph which represents the adjacency that occurs in the RNA but that is broken  
138 in the genome by an intron sequence. A statistical segmentation strategy is used to determine  
139 paths with aligned read enrichment over background noise (Guttman *et al.* 2010).

140 Both Cufflinks and Scripture were run using default parameters, however due to  
141 computational time Scripture was run on the named chromosomes only (not on the  
142 unanchored contigs "chrUn"). The samples were assembled individually to reduce the  
143 complexity of isoforms and hence reduce the chance of incorrectly assembled transcripts  
144 (Trapnell *et al.* 2012). The Cufflinks and Scripture assembly files are provided as  
145 supplemental files 2 and 3.

146 The "Cuffmerge" program (included in the Cufflinks package) was used to merge the  
147 Cufflinks and Scripture assemblies separately. Stranded transcripts from the two assemblies  
148 were compared using the Cufflinks inclusive program "Cuffcompare" with the Cufflinks  
149 assembly as a mock reference. The class codes in the Cuffcompare output were used to  
150 generate a consensus assembly (University of Nottingham "UoN", supplemental file 4). This

151 consensus assembly was compared to the Ensembl annotations using Cuffcompare  
152 (supplemental file 5).

153

#### 154 *Annotation*

155 The UoN cDNA sequences (supplemental file 6) were extracted from the consensus assembly  
156 (\*gtf) file and the longest open reading frames (ORFs) were determined. Gene annotation was  
157 conducted by prediction of Pfam domains (PfamA.hmm library downloaded June 2013)  
158 (Punta *et al.* 2012) using HMMER (Eddy 2011). Associated gene ontology (GO) terms  
159 (Ashburner *et al.* 2000) were determined using the Pfam2GO database (version compiled  
160 15/6/2013) of Interpro (Hunter *et al.* 2009). The UoN transcripts were searched against the  
161 NCBI non-redundant (NR) database (downloaded 14<sup>th</sup> November 2013) using BLASTX  
162 (Altschul *et al.* 1997), a cutoff evalue of 1e-10 was used to infer homology.

163

#### 164 *Gene expression analyses*

165 The TopHat BAM files were filtered for unique alignments (SAM flag NH:i:1) and the  
166 number of tags per Ensembl gene was calculated using htseq-count ([http://www-](http://www-huber.embl.de/users/anders/HTSeq/doc/count.html)  
167 [huber.embl.de/users/anders/HTSeq/doc/count.html](http://www-huber.embl.de/users/anders/HTSeq/doc/count.html)). These counts were converted into Reads  
168 per Kb per million (RPKM) values (Mortazavi *et al.* 2008). A table of RPKM values for all  
169 Ensembl genes is provided as supplemental file 7.

170 As the number of replicates was limiting, identification of genes differentially  
171 expressed between samples was not attempted. However, genes enriched in each sample were  
172 identified using the following criteria; RPKM > 5 for a sample and RPKM > 10 x the mean of  
173 RPKMs for the other samples (supplemental file 8). The “hclust” command in R (R-Core-  
174 Team 2013) was used for the hierarchical clustering analysis of gene expression values  
175 (RPKMs). It was performed using the default complete linkage method and Euclidean  
176 distance. Probability values for each cluster were calculated using the “pvclust” R package  
177 (Suzuki & Shimodaira 2006) (bootstrap n = 1000).

178

#### 179 *Comparison of horse and human gene families*

180 To identify orthologous and potential paralogous gene expansions in the horse evident in our  
181 transcriptome data, translations of all horse transcripts were compared to proteins encoded by  
182 known human genes (Ensembl build GRCh37.71). Both human and horse proteome sets were  
183 first clustered to collapse within-species identical protein sequences generated from  
184 alternative transcripts using CD-HIT (Li & Godzik 2006). This resulted in 64,231 human and  
185 29,090 horse sequences. These were compared using Inparanoid (version 4.1, overlap cutoff =  
186 0.5, group merging cutoff = 0.5, scoring matrix BLOSUM62) (Remm *et al.* 2001). Functional  
187 comparison of gene sets was conducted using Ingenuity Pathway Analysis (Ingenuity  
188 Systems).

189

## 190 **Results**

### 191 *Transcriptome assemblies*

192 Around 184 million reads were generated and 159 million remained after trimming;  
193 approximately 68.6 million of which were aligned to the reference genome EquCab2 (Table  
194 1). Scripture assembled 102,270 stranded transcripts (27,610 with >1 exon, supplemental file  
195 3) whereas Cufflinks reconstructed 58,182 (20,459 with >1 exon, supplemental file 2). There  
196 were 10,518 Cufflinks transcripts that completely matched the intron chain of the Scripture  
197 transcripts. In addition to this 18,152 Cufflinks transcripts contained or covered at least one  
198 Scripture transcript with the same compatible intron structure. The union of these two sets  
199 resulted in 28,230 transcripts, 14,762 of which contained more than one exon (supplemental  
200 file 4).

201  
202 *Comparison of consensus assembly to Ensembl*  
203 The similarities between the 28,230 consensus transcripts (henceforth referred to as “UoN”,  
204 University Of Nottingham) and the 28,944 Ensembl transcripts on the named chromosomes  
205 were compared (supplemental file 5). There were only 507 UoN transcripts which completely  
206 matched the intron chain of an Ensembl transcript. The majority of transcripts (8763, 31%)  
207 were identified as potentially novel isoforms of a predicted Ensembl transcript with at least  
208 one splice junction shared.

209 The majority of Ensembl transcripts (18668, 65%) did not overlap with a UoN  
210 transcript (supplemental file 9). This could be due to the strict consensus approach used for  
211 the UoN assembly. Also, the specific tissues analysed would not be expected to reconstruct all  
212 the transcripts from Ensembl which are predicted from genomic DNA and hence all potential  
213 transcriptomes.

214 Around 9,500 (34%) of the 28,230 UoN transcripts were annotated with a Pfam  
215 protein domain, approximately 6,600 (23%) with at least one GO term and 16,166 (57%) had  
216 at least one significant BLASTX hit against NCBI-NR (supplemental file 10). In total there  
217 were 16,305 UoN transcripts with at least one annotation. The UoN annotated transcripts  
218 were split into Cuffcompare categories based on the comparison to the Ensembl annotations  
219 (supplemental file 10). As expected, the transcripts matching the intron chain (“=”) or sharing  
220 at least one splice junction (“j”) of the Ensembl annotations had the highest percentage of  
221 annotated transcripts (e.g. 97% and 99% with BLASTX hits respectively). There were 367 of  
222 the 16,166 UoN transcripts with a BLASTX hit that showed homology to only a single  
223 species and just under half of these (163) were to *Equus caballus*. The top hit was extracted  
224 for each transcript and as expected most of these hits were also to the *Equus caballus* genome.  
225 Other mammals with a high number of top hits were *Homo sapiens*, *Mus musculus*,  
226 *Ceratotherium simum simum*, *Tursiops truncatus* and *Sus scrofa*. The full list is shown in  
227 supplemental file 11.

228  
229 *Gene expression analyses*  
230 The number of Ensembl genes specific to each sample is shown in Table 2 and supplemental  
231 file 8 (see also materials and methods). No genes were enriched in more than one sample. The  
232 Lymphocyte A sample had many more specific genes than Lymphocyte B. This is possibly  
233 due to sample A being taken from the same horse that the published genome is derived from,  
234 however the read alignment rate between these two samples is similar suggesting this may not  
235 be the major factor. Alternatively this may reflect the immune states of individual horses at  
236 the time of sample collection.

237 The top ten gene ontology (GO) terms for the sample-enriched genes largely reflect  
238 the known function of the tissues sampled (supplemental file 12). Hierarchical clustering  
239 analysis of the RPKMs between tissues showed three clades (Figure 1). For each of the nodes,  
240 the approximately unbiased (au) bootstraps are over 80%. These are reported having  
241 superiority over the bootstrap probabilities (bp) (Suzuki & Shimodaira 2006). The  
242 lymphocyte samples cluster most closely with the spleen sample which likely reflects the high  
243 number of lymphocytes present in the spleen at the time of collection. Whilst the kidney and  
244 liver have general shared roles in waste excretion suggesting a possible overlap of  
245 transcription profile, determining a definitive reason for the separation of the clade containing  
246 lymph node, kidney and liver is not clear. The jejunum sample forms an outgroup and this  
247 separation from the other immune-like tissues likely reflects the relatively smaller proportion  
248 of lymphoid (Peyer’s patch) tissue to non-lymphoid material in this organ. It is also important  
249 to consider that only a limited number of samples and animals are compared and so  
250 robustness of these relationships is not ensured.



251 Analysis of genes enriched in each sample identified related enriched canonical  
252 pathways. The kidney sample is enriched in genes involved in the “ $\gamma$ -glutamyl Cycle”,  
253 “Leukotriene Biosynthesis”, “Glycine Cleavage Complex”, “ $\beta$ -alanine Degradation I” and “4-  
254 hydroxyproline Degradation I” pathways. Amino-acid catabolism pathways, possibly  
255 reflecting high-energy consumption of the kidney, dominate these. The liver sample is  
256 enriched with genes involved in the degradation of chemical products (e.g. nicotine and  
257 melatonin). Enzymes including members of the CYP450 and UDP-Glucuronosyltransferase  
258 (UGT) gene families, which are known to be highly expressed in the liver, are enriched. The  
259 spleen shows enrichment of genes involved in the pathways “Autoimmune Thyroid Disease  
260 Signaling”, “Hematopoiesis from Pluripotent Stem Cells”, “Primary Immunodeficiency  
261 Signaling”, “Dendritic Cell Maturation”, and “Agranulocyte Adhesion and Diapedesis”.  
262 Largely these are due to the enrichment of genes encoding the immunoglobulin heavy chain  
263 and Fc fragment of IgG. Enrichment of these pathways reflects the role of the spleen as the  
264 primary site of white blood cell differentiation and storage. The lymph node sample is  
265 enriched in the pathways, “Primary Immunodeficiency Signaling”, “Hematopoiesis from  
266 Pluripotent Stem Cells”, “Autoimmune Thyroid Disease Signaling”, “Allograft Rejection  
267 Signaling” and “Communication between Innate and Adaptive Immune Cells”. As with the  
268 spleen these are predominantly due to the enrichment of genes encoding the immunoglobulin  
269 heavy chain proteins and result from the contained white blood cell content.  
270

#### 271 *Identification of paralogous gene expansions in horse*

272 Previously the horse genome was described as containing lineage specific expansions of  
273 olfactory and immune genes (Wade *et al.* 2009). The expansion of these families particularly  
274 immune related genes is often seen in mammalian genome comparisons (Emes *et al.* 2003).  
275 Wade *et al.* (2009) reported that there were 99 gene families expanded in the horse genome.  
276 Comparison of the proteins encoded by the transcripts identified here identified 4,605 groups  
277 of horse:human orthologs and 10,607 inparalogs. The majority of these represent expansions  
278 in human where a single horse protein was encoded by the transcriptome data generated here.  
279 91 families were identified with a specific expansion in horses (many:1 relationship). Of these  
280 the large majority (83/91) represent simple duplications in the horse transcriptome compared  
281 to human. Three families have four non-identical encoded proteins orthologous to a single  
282 protein in humans. Annotation of these genes identifies them as T cell receptor alpha constant  
283 (TRAC), heparin sulfate proteoglycan 2 (HSPG2 and solute carrier family 23 (ascorbic acid  
284 transporter) member 1 (SLC23A1). An additional four gene families are identified with three  
285 encoded proteins in horse compared to a single protein in human. These are GTPase, IMAP  
286 family member 7 (GIMAP7), UDP glucuronosyltransferase 1 family polypeptide A6  
287 (UGT1A6), solute carrier family 44 (SLC44A2), ATP-binding cassette, sub-family C member  
288 8 (ABCC8 and sushi, nidogen and EGF-like domains 1 (SNED1).

289 An additional 99 families were found with expansions in both human and horse  
290 (many:many relationship). Reflecting the tissues used for RNA extraction, genes in this  
291 category are highly enriched for immune functions. The most significantly populated  
292 pathways are “Role of NFAT in regulation of the immune response”, “CD28 Signaling in T  
293 helper cells”, “iCOS-iCOSL signaling in T helper cells”, “Natural killer cell signaling” and  
294 “PKC $\theta$  signaling in T lymphocytes”.  
295

#### 296 **Discussion**

297 The analysis conducted here provided insight into the transcriptome of immune tissues from  
298 the horse and made these analyses freely available (supplemental files). Whilst it is unclear  
299 why the horse transcriptome should contain the specific expansions of gene families  
300 described, the analysis provided insight into potential areas of T-cell biology which may

301 underlie equine specific immunobiology. The analysis conducted also allowed the  
302 identification of gene expansions such as UGT1A6, part of a putative paralogous gene  
303 expansion in horse relative to human. UGT1A6 is a member of the UDP-  
304 glucuronosyltransferases (UGTs), a gene family essential for metabolism of both xenobiotic  
305 and endobiotic substances. In contrast to humans and model organisms, there is currently little  
306 information regarding specific drug metabolism in animals of veterinary importance. This is  
307 particularly true in the horse, despite it being potentially exposed to extensive medical care.  
308 Due to the broad application of its mechanisms on xenobiotic substances, the UGT enzyme  
309 group has important implications in pharmacokinetics, the development of drugs and their  
310 associated elimination rates. Importantly, as many of the drugs used in equids are adopted  
311 from those designed from human UGT research, understanding the differences in genes  
312 encoding these proteins may provide a basis for investigation into the UGT group of enzymes  
313 in horses and will open up further opportunities for specific pharmacokinetic research into  
314 UGT related equine drug metabolism potentially reducing toxic drug interactions.

315 The data presented here demonstrated the utility of second generation sequencing in  
316 significantly advancing knowledge of gene transcription in a poorly characterised species. A  
317 large number of potential novel genes were identified alongside some extensions to existing  
318 genes. The completeness of these predictions remains to be confirmed by traditional mRNA  
319 isolation and sequencing but the data presented provides a starting point for the study of  
320 whole groups of genes.

### 321 **Acknowledgements**

322 Dr. Julia Kydd (School of Veterinary Medicine and Science, University of Nottingham) and  
323 Dr. Donald Miller (Baker Institute of Animal Health, Cornell University, USA) for their kind  
324 donation of lymphocyte samples and RNA. We would also like to thank Dr Martin Blythe,  
325 Damian Kao, Victoria Wright and Katharine Rangeley (University of Nottingham) for useful  
326 discussions.  
327

### 328 **Supplemental Information**

#### 329 **Supplemental Files**

- 330 • *S1\_CLC\_SOLiD\_trim\_adapter\_list.xls*: The trim adapter list used in CLC with SOLiD  
331 adapter sequences and single-base-repeat sequences.
- 332 • *S2\_Cufflinks\_assembly.gtf.gz*: Annotation of the individual pre-consensus horse  
333 Cufflinks assembly.
- 334 • *S3\_Scripture\_assembly.gtf.gz*: Annotation of the individual pre-consensus horse  
335 Scripture assembly.
- 336 • *S4\_UoN\_horse\_consensus\_assembly.gtf.gz*: Annotation of the University of  
337 Nottingham (UoN) final consensus horse assembly.
- 338 • *S5\_Ensembl-vs-UoN\_Cuffcompare-results.xls*: Table of Cuffcompare results showing  
339 the similarities between the University of Nottingham (UoN) consensus assembly and  
340 the Ensembl annotations (using Ensembl as a Cuffcompare reference).
- 341 • *S6\_UoN\_horse\_cDNA\_sequences.fa.gz*: University of Nottingham (UoN) consensus  
342 horse cDNA sequences.
- 343 • *S7\_Ensembl71\_genes\_RPKMs\_UoN-reads.xls*: RPKM values for all Ensembl v71  
344 genes.
- 345 • *S8\_sample-enriched-genes.zip*: Tables containing the sample enriched Ensembl v71  
346 gene IDs.
- 347 • *S9\_Ensembl-vs-UoN\_Cuffcompare-results\_UoN-as-ref.xls*: Table of Cuffcompare  
348 results showing the similarities between the University of Nottingham (UoN)  
349 consensus assembly and the Ensembl annotations (using UoN as a Cuffcompare  
350



- 351 reference).
- 352 • *S10\_Ensembl-vs-UoN\_Cuffcompare-results\_PFAM-GO-BLASTX.xls*: Table of
- 353 Cuffcompare results showing the similarities between the University of Nottingham
- 354 (UoN) consensus assembly and Ensembl. The number of annotated UoN transcripts
- 355 for each of the Cuffcompare categories is also shown.
- 356 • *S11\_species\_UoN\_numTopHits\_BLASTX.xls*: Table showing the number of top hit
- 357 BLASTX hits for each species.
- 358 • *S12\_top10-GO-names\_sample-enriched-genes.xls*: Table showing the top ten gene
- 359 ontology (GO) term names for the sample-enriched genes.

360

361 **Figure Legends**

362 **Figure 1: Hierarchical clustering of gene expression profiles in 7 tissues**

363 The R command “hclust” was used for the hierarchical clustering analysis. The branch values

364 are the pvclust approximately unbiased (AU) p-values (left) and bootstrap (BP) probability

365 values (right) where the p-values are expressed as percentages (95% is equivalent to p-value <

366 0.05) (Beliakova-Bethell *et al.* 2013).

367

368 **Tables**

369 **Table 1: Read statistics for the seven samples**

Sample	Horse	Raw reads	Trimmed reads	Reads aligned
Lymphocyte A	A	20,853,992	18,243,283	7,856,017
Lymphocyte B	B	32,050,093	27,315,182	11,659,787
Jejunum	C	19,902,170	17,241,772	7,659,938
Kidney	C	33,158,285	27,746,321	10,937,750
Liver	C	23,176,545	19,982,256	8,565,159
Lymph node	C	24,671,029	21,444,476	9,221,340
Spleen	C	30,421,675	26,828,834	12,708,499

370 (A) "Twilight", healthy Thoroughbred (B) healthy castrated male welsh mountain pony (C)

371 aged gelding euthanised for arthritis.

372

373 **Table 2: Sample-enriched genes**

Sample	# Sample-enriched genes
Lymphocyte A	201
Lymphocyte B	23
Jejunum	228
Kidney	318
Liver	272
Lymph node	44
Spleen	79

374

375 **References**

- 376
- 377 Altschul S.F., Madden T.L., Schaffer A.A., Zhang J., Zhang Z., Miller W. & Lipman D.J.
- 378 (1997) Gapped BLAST and PSI-BLAST: a new generation of protein database search
- 379 programs. *Nucleic Acids Res* **25**, 3389-402.
- 380 Ashburner M., Ball C.A., Blake J.A., Botstein D., Butler H., Cherry J.M., Davis A.P.,
- 381 Dolinski K., Dwight S.S., Eppig J.T., Harris M.A., Hill D.P., Issel-Tarver L.,

- 382 Kasarskis A., Lewis S., Matese J.C., Richardson J.E., Ringwald M., Rubin G.M. &  
383 Sherlock G. (2000) Gene ontology: tool for the unification of biology. The Gene  
384 Ontology Consortium. *Nat Genet* **25**, 25-9.
- 385 Beliakova-Bethell N., Massanella M., White C., Lada S.M., Du P., Vaida F., Blanco J., Spina  
386 C.A. & Woelk C.H. (2013) The effect of cell subset isolation method on gene  
387 expression in leukocytes. *Cytometry* **85**, 94-104.
- 388 Brown K., Moreton J., Malla S., Aboobaker A.A., Emes R.D. & Tarlinton R.E. (2012)  
389 Characterisation of retroviruses in the horse genome and their transcriptional activity  
390 via transcriptome sequencing. *Virology* **433**, 55-63.
- 391 Bryant C.E., Ouellette A., Lohmann K., Vandenplas M., Moore J.N., Maskell D.J. &  
392 Farnfield B.A. (2007) The cellular Toll-like receptor 4 antagonist E5531 can act as an  
393 agonist in horse whole blood. *Vet Immunol Immunopathol* **116**, 182-9.
- 394 Cabili M.N., Trapnell C., Goff L., Koziol M., Tazon-Vega B., Regev A. & Rinn J.L. (2011)  
395 Integrative annotation of human large intergenic noncoding RNAs reveals global  
396 properties and specific subclasses. *Genes & development* **25**, 1915-27.
- 397 Capomaccio S., Vitulo N., Verini-Supplizi A., Barcaccia G., Albiero A., D'Angelo M.,  
398 Campagna D., Valle G., Felicetti M. & Silvestrelli M. (2013) RNA Sequencing of the  
399 Exercise Transcriptome in Equine Athletes. *PLoS one* **8**, e83504.
- 400 Coleman S.J., Zeng Z., Wang K., Luo S., Khrebtukova I., Mienaltowski M.J., Schroth G.P.,  
401 Liu J. & MacLeod J.N. (2010) Structural annotation of equine protein-coding genes  
402 determined by mRNA sequencing. *Anim Genet* **41 Suppl 2**, 121-30.
- 403 Das P.J., McCarthy F., Vishnoi M., Paria N., Gresham C., Li G., Kachroo P., Sudderth A.K.,  
404 Teague S. & Love C.C. (2013) Stallion sperm transcriptome comprises functionally  
405 coherent coding and regulatory RNAs as revealed by microarray analysis and RNA-  
406 seq. *PLoS one* **8**, e56535.
- 407 Eddy S.R. (2011) Accelerated Profile HMM Searches. *PLoS Comput Biol* **7**, e1002195.
- 408 Emes R.D., Goodstadt L., Winter E.E. & Ponting C.P. (2003) Comparison of the genomes of  
409 human and mouse lays the foundation of genome zoology. *Hum Mol Genet* **12**, 701-9.
- 410 Guttman M., Garber M., Levin J.Z., Donaghey J., Robinson J., Adiconis X., Fan L., Koziol  
411 M.J., Gnirke A. & Nusbaum C. (2010) Ab initio reconstruction of cell type-specific  
412 transcriptomes in mouse reveals the conserved multi-exonic structure of lincRNAs.  
413 *Nature biotechnology* **28**, 503-10.
- 414 Hill E.W., Gu J., McGivney B.A. & MacHugh D.E. (2010) Targets of selection in the  
415 Thoroughbred genome contain exercise-relevant gene SNPs associated with elite  
416 racecourse performance. *Anim Genet* **41 Suppl 2**, 56-63.
- 417 Hunter S., Apweiler R., Attwood T.K., Bairoch A., Bateman A., Binns D., Bork P., Das U.,  
418 Daugherty L., Duquenne L., Finn R.D., Gough J., Haft D., Hulo N., Kahn D., Kelly  
419 E., Laugraud A., Letunic I., Lonsdale D., Lopez R., Madera M., Maslen J., McAnulla  
420 C., McDowall J., Mistry J., Mitchell A., Mulder N., Natale D., Orengo C., Quinn A.F.,  
421 Selengut J.D., Sigrist C.J., Thimma M., Thomas P.D., Valentin F., Wilson D., Wu  
422 C.H. & Yeats C. (2009) InterPro: the integrative protein signature database. *Nucleic  
423 Acids Res* **37**, D211-5.
- 424 Iqbal K., Chitwood J.L., Meyers-Brown G.A., Roser J.F. & Ross P.J. (2014) RNA-Seq  
425 Transcriptome Profiling of Equine Inner Cell Mass and Trophectoderm. *Biology of  
426 Reproduction*.
- 427 Langmead B., Trapnell C., Pop M. & Salzberg S.L. (2009) Ultrafast and memory-efficient  
428 alignment of short DNA sequences to the human genome. *Genome Biol* **10**, R25.
- 429 Li W. & Godzik A. (2006) Cd-hit: a fast program for clustering and comparing large sets of  
430 protein or nucleotide sequences. *Bioinformatics* **22**, 1658-9.
- 431 McGivney B.A., McGettigan P.A., Browne J.A., Evans A.C., Fonseca R.G., Loftus B.J.,

- 432 Lohan A., MacHugh D.E., Murphy B.A., Katz L.M. & Hill E.W. (2010)  
433 Characterization of the equine skeletal muscle transcriptome identifies novel  
434 functional responses to exercise training. *BMC genomics* **11**, 398.
- 435 Mortazavi A., Williams B.A., McCue K., Schaeffer L. & Wold B. (2008) Mapping and  
436 quantifying mammalian transcriptomes by RNA-Seq. *Nat Methods* **5**, 621-8.
- 437 Park K.D., Park J., Ko J., Kim B.C., Kim H.S., Ahn K., Do K.T., Choi H., Kim H.M. & Song  
438 S. (2012) Whole transcriptome analyses of six thoroughbred horses before and after  
439 exercise using RNA-Seq. *BMC genomics* **13**, 473.
- 440 Punta M., Coghill P.C., Eberhardt R.Y., Mistry J., Tate J., Boursnell C., Pang N., Forslund K.,  
441 Ceric G. & Clements J. (2012) The Pfam protein families database. *Nucleic acids*  
442 *research* **40**, D290-D301.
- 443 R-Core-Team (2013) R: A Language and Environment for Statistical Computing. R  
444 Foundation for Statistical Computing. Vienna, Austria. <http://www.R-project.org>.
- 445 Remm M., Storm C.E. & Sonnhammer E.L. (2001) Automatic clustering of orthologs and in-  
446 paralogs from pairwise species comparisons. *J Mol Biol* **314**, 1041-52.
- 447 Serteyn D., Piquemal D., Vanderheyden L., Lejeune J.P., Verwilghen D. & Sandersen C.  
448 (2010) Gene expression profiling from leukocytes of horses affected by  
449 osteochondrosis. *J Orthop Res* **28**, 965-70.
- 450 Suzuki R. & Shimodaira H. (2006) Pvcust: an R package for assessing the uncertainty in  
451 hierarchical clustering. *Bioinformatics* **22**, 1540-2.
- 452 Trapnell C., Pachter L. & Salzberg S.L. (2009) TopHat: discovering splice junctions with  
453 RNA-Seq. *Bioinformatics* **25**, 1105-11.
- 454 Trapnell C., Roberts A., Goff L., Pertea G., Kim D., Kelley D.R., Pimentel H., Salzberg S.L.,  
455 Rinn J.L. & Pachter L. (2012) Differential gene and transcript expression analysis of  
456 RNA-seq experiments with TopHat and Cufflinks. *Nature protocols* **7**, 562-78.
- 457 Trapnell C., Williams B.A., Pertea G., Mortazavi A., Kwan G., van Baren M.J., Salzberg  
458 S.L., Wold B.J. & Pachter L. (2010) Transcript assembly and quantification by RNA-  
459 Seq reveals unannotated transcripts and isoform switching during cell differentiation.  
460 *Nat Biotechnol* **28**, 511-5.
- 461 Wade C.M., Giulotto E., Sigurdsson S., Zoli M., Gnerre S., Imsland F., Lear T.L., Adelson  
462 D.L., Bailey E., Bellone R.R., Blocker H., Distl O., Edgar R.C., Garber M., Leeb T.,  
463 Mauceli E., MacLeod J.N., Penedo M.C., Raison J.M., Sharpe T., Vogel J., Andersson  
464 L., Antczak D.F., Biagi T., Binns M.M., Chowdhary B.P., Coleman S.J., Della Valle  
465 G., Fryc S., Guerin G., Hasegawa T., Hill E.W., Jurka J., Kiialainen A., Lindgren G.,  
466 Liu J., Magnani E., Mickelson J.R., Murray J., Nergadze S.G., Onofrio R., Pedroni S.,  
467 Piras M.F., Raudsepp T., Rocchi M., Roed K.H., Ryder O.A., Searle S., Skow L.,  
468 Swinburne J.E., Syvanen A.C., Tozaki T., Valberg S.J., Vaudin M., White J.R., Zody  
469 M.C., Lander E.S. & Lindblad-Toh K. (2009) Genome sequence, comparative  
470 analysis, and population genetics of the domestic horse. *Science* **326**, 865-7.  
471  
472