

Characterisation of the horse transcriptome from immunologically active tissues

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

1 Authors

2 J. Moreton^{1,2,3}, S. Malla², A. A. Aboobaker⁴, R. E. Tarlinton³ and R. D. Emes^{1,3}

3 1 Advanced Data Analysis Centre, University of Nottingham, Sutton Bonington Campus,
4 Loughborough, Leicestershire, LE12 5RD, UK

5 2 Deep Seq, Centre for Genetics and Genomics, University of Nottingham, Queen's Medical
6 Centre, NG7 2UH, UK

7 3 School of Veterinary Medicine and Science, University of Nottingham, Sutton Bonington
8 Campus, Loughborough, Leicestershire, LE12 5RD, UK

9 4 Department of Zoology, University of Oxford, Oxford, OX1 3PS, UK

10 Corresponding author:

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

13 0115 951 6359

14 Joanna.Moreton@nottingham.ac.uk

15 Introduction

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

25 Due to a lack of expressed sequence tag (EST) data, the current annotation of the protein
26 coding regions of the horse genome is largely derived from extrapolation from the genes of other
27 species ([Coleman *et al.* 2010](#)). Several recent papers have outlined transcription profiles using
28 digital gene analysis of a variety of horse tissues including muscle, leukocytes, cartilage, brain,
29 reproductive tissue, embryos, sperm and blood ([Coleman *et al.* 2010](#); [McGivney *et al.* 2010](#);
30 [Serteyn *et al.* 2010](#); [Park *et al.* 2012](#); [Capomaccio *et al.* 2013](#); [Das *et al.* 2013](#); [Iqbal *et al.* 2014](#)).
31 [Capomaccio *et al.* \(2013\)](#) identified new putative non-coding sequences within intergenic and
32 intronic regions whilst [Das *et al.* \(2013\)](#) suggested additions to the structural annotation of four
33 sperm genes. Two of the other studies ([Coleman *et al.* 2010](#); [Park *et al.* 2012](#)) detailed extensions
34 to the annotated gene catalogue in the horse based on transcriptome analysis of quite differing
35 tissue sets, methods and results to those used in this paper. They show that the actual expressed
36 transcription profile only partially overlaps the annotated gene set. A direct comparison of our
37 and these two studies is difficult due to the differing tissues, methodologies and the lack of
38 available locations of the predicted novel genes from these studies.

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

45 Materials and methods

46 *Samples, library preparation and sequencing*

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

53 Lymphocytes isolated by Ficoll Paque (GE healthcare) from a healthy 11 year old welsh
54 mountain pony gelding were kindly provided by Dr Julia Kydd (School of Veterinary Medicine
55 and Science, University of Nottingham) under the Home Office and local Ethical Approval
56 Committee (PPL 40/3354). RNA from lymphocytes isolated from a healthy Thoroughbred mare
57 (DNA the horse genome is derived from) was kindly provided by Donald Miller (Baker Institute
58 of Animal Health, Cornell University, USA). This horse was maintained at the Baker Institute for
59 Animal Health, Cornell University, Ithaca, N.Y., USA. Animal care and research activities were
60 performed in accordance with the guidelines set forth by the Institutional Animal Care and Use
61 Committee of Cornell University, protocol # 1986-0216, approved until March 2013.

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

65 *Read trimming and alignment*

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

73 TopHat 2.0.9 ([Trapnell et al. 2009](#)) was used to align the reads to the repeat masked
74 version of the horse genome (Ensembl v71) to enable non-redundant transcriptome analysis.
75 TopHat first aligns non-spliced reads using Bowtie 1.0.0 ([Langmead et al. 2009](#)) then identifies
76 splice junctions. Gapped alignments are then used by TopHat to map the reads not aligned by
77 Bowtie. In order to utilise the splice sites in all samples, two iterations of TopHat alignments
78 were carried out ([Cabili et al. 2011](#)). Firstly, the reads from each sample were aligned to the
79 repeat-masked horse genome with default parameters. The splice sites ("junctions") were
80 extracted from all of the output files and duplicates were removed leaving 216,007 sites. These
81 splice sites were pooled together with the non-redundant sites extracted from the Ensembl
82 annotation yielding 399,264 non-redundant splice sites. Each of the samples were then re-aligned
83 with TopHat using the pooled non-redundant splice sites file (with 'raw-juncs' and 'no-novel-
84 juncs' parameters) to the repeat-masked genome.

85 TopHat was used for the read alignment because it is part of the Tuxedo suite and is
86 therefore a natural input for the Cufflinks assembler ([Trapnell et al. 2010](#)). It is also the preferred
87 aligner for Scripture ([Guttman et al. 2010](#)). Cufflinks and Scripture are described in the
88 transcriptome assembly section.

89 *De novo transcriptome assembly*

90 Each of the samples were assembled into separate transcriptomes using two different tools;
91 Cufflinks v2.1.1 ([Trapnell et al. 2010](#)) and Scripture ([Guttman et al. 2010](#)) (beta2 version,
92 December 2010). These tools use different approaches for transcript assembly. A minimal set of

93 transcripts is assembled by Cufflinks using a probabilistic model. It performs a minimum cost
94 maximum matching in bipartite graphs ([Trapnell et al. 2010](#)). Scripture however creates a
95 connectivity graph which represents the adjacency that occurs in the RNA but that is broken in
96 the genome by an intron sequence. A statistical segmentation strategy is used to determine paths
97 with aligned read enrichment over background noise ([Guttman et al. 2010](#)).

98 Both Cufflinks and Scripture were run using default parameters, however due to
99 computational time Scripture was run on the named chromosomes only (not on the unanchored
100 contigs "chrUn"). The samples were assembled individually to reduce the complexity of isoforms
101 and hence reduce the chance of incorrectly assembled transcripts ([Trapnell et al. 2012](#)). The
102 Cufflinks and Scripture assembly files are provided as supplemental files 2 and 3.

103 The "Cuffmerge" program (included in the Cufflinks package) was used to merge the
104 Cufflinks and Scripture assemblies separately. Stranded transcripts from the two assemblies were
105 compared using the Cufflinks inclusive program "Cuffcompare" with the Cufflinks assembly as a
106 mock reference. The class codes in the Cuffcompare output were used to generate a consensus
107 assembly (University of Nottingham "UoN", supplemental file 4). This consensus assembly was
108 compared to the Ensembl annotations using Cuffcompare (supplemental file 5).

109 *Annotation*

110 The UoN cDNA sequences (supplemental file 6) were extracted from the consensus assembly
111 (*gtf) file and the longest open reading frames (ORFs) were determined. Gene annotation was
112 conducted by prediction of Pfam domains (PfamA.hmm library downloaded June 2013) ([Punta et al.](#)
113 [2012](#)) using HMMER ([Eddy 2011](#)). Associated gene ontology (GO) terms ([Ashburner et al.](#)
114 [2000](#)) were determined using the Pfam2GO database (version compiled 15/6/2013) of Interpro
115 ([Hunter et al. 2009](#)). The UoN transcripts were searched against the NCBI non-redundant (NR)
116 database (downloaded 14th November 2013) using BLASTX ([Altschul et al. 1997](#)), a cutoff
117 evaluate of 1e-10 was used to infer homology.

118 *Gene expression analyses*

119 The TopHat BAM files were filtered for unique alignments (SAM flag NH:i:1) and the number of
120 tags per Ensembl gene was calculated using htseq-count ([http://www-](http://www-huber.embl.de/users/anders/HTSeq/doc/count.html)
121 [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 per
122 Kb per million (RPKM) values ([Mortazavi et al. 2008](#)). A table of RPKM values for all Ensembl
123 genes is provided as supplemental file 7.

124 As the number of replicates was limiting, identification of genes differentially expressed
125 between samples was not attempted. However, genes enriched in each sample were identified
126 using the following criteria; RPKM > 5 for a sample and RPKM > 10 x the mean of RPKMs for
127 the other samples (supplemental file 8). The "hclust" command in R ([R-Core-Team 2013](#)) was
128 used for the hierarchical clustering analysis of gene expression values (RPKM). It was
129 performed using the default complete linkage method and Euclidean distance. Probability values
130 for each cluster were calculated using the "pvclust" R package ([Suzuki & Shimodaira 2006](#))
131 (bootstrap n = 1000).

132 *Comparison of horse and human gene families*

133 To identify orthologous and potential paralogous gene expansions in the horse evident in our
134 transcriptome data, translations of all horse transcripts were compared to proteins encoded by
135 known human genes (Ensembl build GRCh37.71). Both human and horse proteome sets were
136 first clustered to collapse within-species identical protein sequences generated from alternative
137 transcripts using CD-HIT ([Li & Godzik 2006](#)). This resulted in 64,231 human and 29,090 horse
138 sequences. These were compared using Inparanoid (version 4.1, overlap cutoff = 0.5, group

139 merging cutoff = 0.5, scoring matrix BLOSUM62) ([Remm et al. 2001](#)). Functional comparison of
140 gene sets was conducted using Ingenuity Pathway Analysis (Ingenuity Systems).

141 **Results**

142 *Transcriptome assemblies*

143 Around 184 million reads were generated and 159 million remained after trimming;
144 approximately 68.6 million of which were aligned to the reference genome EquCab2 (Table 1).
145 Scripture assembled 102,270 stranded transcripts (27,610 with >1 exon, supplemental file 3)
146 whereas Cufflinks reconstructed 58,182 (20,459 with >1 exon, supplemental file 2). There were
147 10,518 Cufflinks transcripts that completely matched the intron chain of the Scripture transcripts.
148 In addition to this 18,152 Cufflinks transcripts contained or covered at least one Scripture
149 transcript with the same compatible intron structure. The union of these two sets resulted in
150 28,230 transcripts, 14,762 of which contained more than one exon (supplemental file 4).

151 *Comparison of consensus assembly to Ensembl*

152 The similarities between the 28,230 consensus transcripts (henceforth referred to as “UoN”,
153 University Of Nottingham) and the 28,944 Ensembl transcripts on the named chromosomes were
154 compared (supplemental file 5). There were only 507 UoN transcripts which completely matched
155 the intron chain of an Ensembl transcript. The majority of transcripts (8763, 31%) were identified
156 as potentially novel isoforms of a predicted Ensembl transcript with at least one splice junction
157 shared.

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

163 Around 9,500 (34%) of the 28,230 UoN transcripts were annotated with a Pfam protein
164 domain, approximately 6,600 (23%) with at least one GO term and 16,166 (57%) had at least one
165 significant BLASTX hit against NCBI-NR (supplemental file 10). In total there were 16,305 UoN
166 transcripts with at least one annotation. The UoN annotated transcripts were split into
167 Cuffcompare categories based on the comparison to the Ensembl annotations (supplemental file
168 10). As expected, the transcripts matching the intron chain (“=”) or sharing at least one splice
169 junction (“j”) of the Ensembl annotations had the highest percentage of annotated transcripts (e.g.
170 97% and 99% with BLASTX hits respectively). There were 367 of the 16,166 UoN transcripts
171 with a BLASTX hit that showed homology to only a single species and just under half of these
172 (163) were to *Equus caballus*. The top hit was extracted for each transcript and as expected most
173 of these hits were also to the *Equus caballus* genome. Other mammals with a high number of top
174 hits were *Homo sapiens*, *Mus musculus*, *Ceratotherium simum simum*, *Tursiops truncatus* and
175 *Sus scrofa*. The full list is shown in supplemental file 11.

176 *Gene expression analyses*

177 The number of Ensembl genes specific to each sample is shown in Table 2 and supplemental file
178 8 (see also materials and methods). No genes were enriched in more than one sample. The
179 Lymphocyte A sample had many more specific genes than Lymphocyte B. This is possibly due to
180 sample A being taken from the same horse that the published genome is derived from, however
181 the read alignment rate between these two samples is similar suggesting this may not be the
182 major factor. Alternatively this may reflect the immune states of individual horses at the time of
183 sample collection.

184 The top ten gene ontology (GO) terms for the sample-enriched genes largely reflect the
185 known function of the tissues sampled (supplemental file 12). Hierarchical clustering analysis of
186 the RPKMs between tissues showed three clades (Figure 1). The branch values are the pvclust
187 approximately unbiased (AU) p-values (left) and bootstrap (BP) probability values (right) where
188 the p-values are expressed as percentages (95% is equivalent to p-value < 0.05) ([Beliakova-
189 Bethell et al. 2013](#)). For each of the nodes, the AU bootstraps are over 80% and these are reported
190 having superiority over the BP values ([Suzuki & Shimodaira 2006](#)). The lymphocyte samples
191 cluster most closely with the spleen sample which likely reflects the high number of lymphocytes
192 present in the spleen at the time of collection. Whilst the kidney and liver have general shared
193 roles in waste excretion suggesting a possible overlap of transcription profile, determining a
194 definitive reason for the separation of the clade containing lymph node, kidney and liver is not
195 clear. The jejunum sample forms an outgroup and this separation from the other immune-like
196 tissues likely reflects the relatively smaller proportion of lymphoid (Peyer's patch) tissue to non-
197 lymphoid material in this organ. It is also important to consider that only a limited number of
198 samples and animals are compared and so robustness of these relationships is not ensured.

199 Analysis of genes enriched in each sample identified related enriched canonical pathways.
200 The kidney sample is enriched in genes involved in the “ γ -glutamyl Cycle”, “Leukotriene
201 Biosynthesis”, “Glycine Cleavage Complex”, “ β -alanine Degradation I” and “4-hydroxyproline
202 Degradation I” pathways. Amino-acid catabolism pathways, possibly reflecting high-energy
203 consumption of the kidney, dominate these. The liver sample is enriched with genes involved in
204 the degradation of chemical products (e.g. nicotine and melatonin). Enzymes including members
205 of the CYP450 and UDP-Glucuronosyltransferase (UGT) gene families, which are known to be
206 highly expressed in the liver, are enriched. The spleen shows enrichment of genes involved in the
207 pathways “Autoimmune Thyroid Disease Signaling”, “Hematopoiesis from Pluripotent Stem
208 Cells”, “Primary Immunodeficiency Signaling”, “Dendritic Cell Maturation”, and “Agranulocyte
209 Adhesion and Diapedesis”. Largely these are due to the enrichment of genes encoding the
210 immunoglobulin heavy chain and Fc fragment of IgG. Enrichment of these pathways reflects the
211 role of the spleen as the primary site of white blood cell differentiation and storage. The lymph
212 node sample is enriched in the pathways, “Primary Immunodeficiency Signaling”,
213 “Hematopoiesis from Pluripotent Stem Cells”, “Autoimmune Thyroid Disease Signaling”,
214 “Allograft Rejection Signaling” and “Communication between Innate and Adaptive Immune
215 Cells”. As with the spleen these are predominantly due to the enrichment of genes encoding the
216 immunoglobulin heavy chain proteins and result from the contained white blood cell content.

217 *Identification of paralogous gene expansions in horse*

218 Previously the horse genome was described as containing lineage specific expansions of olfactory
219 and immune genes ([Wade et al. 2009](#)). The expansion of these families particularly immune
220 related genes is often seen in mammalian genome comparisons ([Emes et al. 2003](#)). [Wade et al.
221 \(2009\)](#) reported that there were 99 gene families expanded in the horse genome. Comparison of
222 the proteins encoded by the transcripts identified here identified 4,605 groups of horse:human
223 orthologs and 10,607 inparalogs. The majority of these represent expansions in human where a
224 single horse protein was encoded by the transcriptome data generated here. 91 families were
225 identified with a specific expansion in horses (many:1 relationship). Of these the large majority
226 (83/91) represent simple duplications in the horse transcriptome compared to human. Three
227 families have four non-identical encoded proteins orthologous to a single protein in humans.
228 Annotation of these genes identifies them as T cell receptor alpha constant (TRAC), heparin
229 sulfate proteoglycan 2 (HSPG2 and solute carrier family 23 (ascorbic acid transporter) member 1
230 (SLC23A1). An additional four gene families are identified with three encoded proteins in horse
231 compared to a single protein in human. These are GTPase, IMAP family member 7 (GIMAP7),

232 UDP glucuronosyltransferase 1 family polypeptide A6 (UGT1A6), solute carrier family 44
233 (SLC44A2), ATP-binding cassette, sub-family C member 8 (ABCC8 and sushi, nidogen and
234 EGF-like domains 1 (SNED1).

235 An additional 99 families were found with expansions in both human and horse
236 (many:many relationship). Reflecting the tissues used for RNA extraction, genes in this category
237 are highly enriched for immune functions. The most significantly populated pathways are “Role
238 of NFAT in regulation of the immune response”, “CD28 Signaling in T helper cells”, “iCOS-
239 iCOSL signaling in T helper cells”, “Natural killer cell signaling” and “PKC \square signaling in T
240 lymphocytes”.

241 Discussion

242 The analysis conducted here provided insight into the transcriptome of immune tissues from the
243 horse and made these analyses freely available (supplemental files). Whilst it is unclear why the
244 horse transcriptome should contain the specific expansions of gene families described, the
245 analysis provided insight into potential areas of T-cell biology which may underlie equine
246 specific immunobiology. The analysis conducted also allowed the identification of gene
247 expansions such as UGT1A6, part of a putative paralogous gene expansion in horse relative to
248 human. UGT1A6 is a member of the UDP-glucuronosyltransferases (UGTs), a gene family
249 essential for metabolism of both xenobiotic and endobiotic substances. In contrast to humans and
250 model organisms, there is currently little information regarding specific drug metabolism in
251 animals of veterinary importance. This is particularly true in the horse, despite it being potentially
252 exposed to extensive medical care. Due to the broad application of its mechanisms on xenobiotic
253 substances, the UGT enzyme group has important implications in pharmacokinetics, the
254 development of drugs and their associated elimination rates. Importantly, as many of the drugs
255 used in equids are adopted from those designed from human UGT research, understanding the
256 differences in genes encoding these proteins may provide a basis for investigation into the UGT
257 group of enzymes in horses and will open up further opportunities for specific pharmacokinetic
258 research into UGT related equine drug metabolism potentially reducing toxic drug interactions.

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

265 Acknowledgements

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

271 References

- 272 Altschul S.F., Madden T.L., Schaffer A.A., Zhang J., Zhang Z., Miller W. & Lipman D.J. (1997)
273 Gapped BLAST and PSI-BLAST: a new generation of protein database search programs.
274 *Nucleic Acids Res* **25**, 3389-402.
- 275 Ashburner M., Ball C.A., Blake J.A., Botstein D., Butler H., Cherry J.M., Davis A.P., Dolinski
276 K., Dwight S.S., Eppig J.T., Harris M.A., Hill D.P., Issel-Tarver L., Kasarskis A., Lewis
277 S., Matese J.C., Richardson J.E., Ringwald M., Rubin G.M. & Sherlock G. (2000) Gene

- 278 ontology: tool for the unification of biology. The Gene Ontology Consortium. *Nat Genet*
279 **25**, 25-9.
- 280 Beliakova-Bethell N., Massanella M., White C., Lada S.M., Du P., Vaida F., Blanco J., Spina
281 C.A. & Woelk C.H. (2013) The effect of cell subset isolation method on gene expression
282 in leukocytes. *Cytometry* **85**, 94-104.
- 283 Brown K., Moreton J., Malla S., Aboobaker A.A., Emes R.D. & Tarlinton R.E. (2012)
284 Characterisation of retroviruses in the horse genome and their transcriptional activity via
285 transcriptome sequencing. *Virology* **433**, 55-63.
- 286 Bryant C.E., Ouellette A., Lohmann K., Vandenplas M., Moore J.N., Maskell D.J. & Farnfield
287 B.A. (2007) The cellular Toll-like receptor 4 antagonist E5531 can act as an agonist in
288 horse whole blood. *Vet Immunol Immunopathol* **116**, 182-9.
- 289 Cabili M.N., Trapnell C., Goff L., Koziol M., Tazon-Vega B., Regev A. & Rinn J.L. (2011)
290 Integrative annotation of human large intergenic noncoding RNAs reveals global
291 properties and specific subclasses. *Genes & development* **25**, 1915-27.
- 292 Capomaccio S., Vitulo N., Verini-Supplizi A., Barcaccia G., Albiero A., D'Angelo M., Campagna
293 D., Valle G., Felicetti M. & Silvestrelli M. (2013) RNA Sequencing of the Exercise
294 Transcriptome in Equine Athletes. *PLoS one* **8**, e83504.
- 295 Coleman S.J., Zeng Z., Wang K., Luo S., Khrebtukova I., Mienaltowski M.J., Schroth G.P., Liu J.
296 & MacLeod J.N. (2010) Structural annotation of equine protein-coding genes determined
297 by mRNA sequencing. *Anim Genet* **41 Suppl 2**, 121-30.
- 298 Das P.J., McCarthy F., Vishnoi M., Paria N., Gresham C., Li G., Kachroo P., Sudderth A.K.,
299 Teague S. & Love C.C. (2013) Stallion sperm transcriptome comprises functionally
300 coherent coding and regulatory RNAs as revealed by microarray analysis and RNA-seq.
301 *PLoS one* **8**, e56535.
- 302 Eddy S.R. (2011) Accelerated Profile HMM Searches. *PLoS Comput Biol* **7**, e1002195.
- 303 Emes R.D., Goodstadt L., Winter E.E. & Ponting C.P. (2003) Comparison of the genomes of
304 human and mouse lays the foundation of genome zoology. *Hum Mol Genet* **12**, 701-9.
- 305 Guttman M., Garber M., Levin J.Z., Donaghey J., Robinson J., Adiconis X., Fan L., Koziol M.J.,
306 Gnirke A. & Nusbaum C. (2010) Ab initio reconstruction of cell type-specific
307 transcriptomes in mouse reveals the conserved multi-exonic structure of lincRNAs.
308 *Nature biotechnology* **28**, 503-10.
- 309 Hill E.W., Gu J., McGivney B.A. & MacHugh D.E. (2010) Targets of selection in the
310 Thoroughbred genome contain exercise-relevant gene SNPs associated with elite
311 racecourse performance. *Anim Genet* **41 Suppl 2**, 56-63.
- 312 Hunter S., Apweiler R., Attwood T.K., Bairoch A., Bateman A., Binns D., Bork P., Das U.,
313 Daugherty L., Duquenne L., Finn R.D., Gough J., Haft D., Hulo N., Kahn D., Kelly E.,
314 Laugraud A., Letunic I., Lonsdale D., Lopez R., Madera M., Maslen J., McAnulla C.,
315 McDowall J., Mistry J., Mitchell A., Mulder N., Natale D., Orengo C., Quinn A.F.,
316 Selengut J.D., Sigrist C.J., Thimma M., Thomas P.D., Valentin F., Wilson D., Wu C.H. &
317 Yeats C. (2009) InterPro: the integrative protein signature database. *Nucleic Acids Res* **37**,
318 D211-5.
- 319 Iqbal K., Chitwood J.L., Meyers-Brown G.A., Roser J.F. & Ross P.J. (2014) RNA-Seq
320 Transcriptome Profiling of Equine Inner Cell Mass and Trophectoderm. *Biology of*
321 *Reproduction*.
- 322 Langmead B., Trapnell C., Pop M. & Salzberg S.L. (2009) Ultrafast and memory-efficient
323 alignment of short DNA sequences to the human genome. *Genome Biol* **10**, R25.
- 324 Li W. & Godzik A. (2006) Cd-hit: a fast program for clustering and comparing large sets of
325 protein or nucleotide sequences. *Bioinformatics* **22**, 1658-9.

- 326 McGivney B.A., McGettigan P.A., Browne J.A., Evans A.C., Fonseca R.G., Loftus B.J., Lohan
327 A., MacHugh D.E., Murphy B.A., Katz L.M. & Hill E.W. (2010) Characterization of the
328 equine skeletal muscle transcriptome identifies novel functional responses to exercise
329 training. *BMC genomics* **11**, 398.
- 330 Mortazavi A., Williams B.A., McCue K., Schaeffer L. & Wold B. (2008) Mapping and
331 quantifying mammalian transcriptomes by RNA-Seq. *Nat Methods* **5**, 621-8.
- 332 Park K.D., Park J., Ko J., Kim B.C., Kim H.S., Ahn K., Do K.T., Choi H., Kim H.M. & Song S.
333 (2012) Whole transcriptome analyses of six thoroughbred horses before and after exercise
334 using RNA-Seq. *BMC genomics* **13**, 473.
- 335 Punta M., Coggill P.C., Eberhardt R.Y., Mistry J., Tate J., Boursnell C., Pang N., Forslund K.,
336 Ceric G. & Clements J. (2012) The Pfam protein families database. *Nucleic acids*
337 *research* **40**, D290-D301.
- 338 R-Core-Team (2013) R: A Language and Environment for Statistical Computing. R Foundation
339 for Statistical Computing. Vienna, Austria. <http://www.R-project.org>.
- 340 Remm M., Storm C.E. & Sonnhammer E.L. (2001) Automatic clustering of orthologs and in-
341 paralogs from pairwise species comparisons. *J Mol Biol* **314**, 1041-52.
- 342 Serteyn D., Piquemal D., Vanderheyden L., Lejeune J.P., Verwilghen D. & Sandersen C. (2010)
343 Gene expression profiling from leukocytes of horses affected by osteochondrosis. *J*
344 *Orthop Res* **28**, 965-70.
- 345 Suzuki R. & Shimodaira H. (2006) Pvcust: an R package for assessing the uncertainty in
346 hierarchical clustering. *Bioinformatics* **22**, 1540-2.
- 347 Trapnell C., Pachter L. & Salzberg S.L. (2009) TopHat: discovering splice junctions with RNA-
348 Seq. *Bioinformatics* **25**, 1105-11.
- 349 Trapnell C., Roberts A., Goff L., Pertea G., Kim D., Kelley D.R., Pimentel H., Salzberg S.L.,
350 Rinn J.L. & Pachter L. (2012) Differential gene and transcript expression analysis of
351 RNA-seq experiments with TopHat and Cufflinks. *Nature protocols* **7**, 562-78.
- 352 Trapnell C., Williams B.A., Pertea G., Mortazavi A., Kwan G., van Baren M.J., Salzberg S.L.,
353 Wold B.J. & Pachter L. (2010) Transcript assembly and quantification by RNA-Seq
354 reveals unannotated transcripts and isoform switching during cell differentiation. *Nat*
355 *Biotechnol* **28**, 511-5.
- 356 Wade C.M., Giulotto E., Sigurdsson S., Zoli M., Gnerre S., Imsland F., Lear T.L., Adelson D.L.,
357 Bailey E., Bellone R.R., Blocker H., Distl O., Edgar R.C., Garber M., Leeb T., Mauceli
358 E., MacLeod J.N., Penedo M.C., Raison J.M., Sharpe T., Vogel J., Andersson L., Antczak
359 D.F., Biagi T., Binns M.M., Chowdhary B.P., Coleman S.J., Della Valle G., Fryc S.,
360 Guerin G., Hasegawa T., Hill E.W., Jurka J., Kiialainen A., Lindgren G., Liu J., Magnani
361 E., Mickelson J.R., Murray J., Nergadze S.G., Onofrio R., Pedroni S., Piras M.F.,
362 Raudsepp T., Rocchi M., Roed K.H., Ryder O.A., Searle S., Skow L., Swinburne J.E.,
363 Syvanen A.C., Tozaki T., Valberg S.J., Vaudin M., White J.R., Zody M.C., Lander E.S. &
364 Lindblad-Toh K. (2009) Genome sequence, comparative analysis, and population genetics
365 of the domestic horse. *Science* **326**, 865-7.

Table 1 (on next page)

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

(A) "Twilight", healthy Thoroughbred (B) healthy castrated male welsh mountain pony (C) aged gelding euthanised for arthritis.

Table 2(on next page)

Number of sample-enriched genes

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

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

Figure 1

Hierarchical clustering of gene expression profiles in seven tissues

The R command “hclust” was used for the hierarchical clustering analysis. The branch values are the pvclust approximately unbiased (AU) p-values (left) and bootstrap (BP) probability values (right) where the p-values are expressed as percentages.

