

1 **Global analysis of transcription start sites in the new ovine reference genome**
2 **(*Oar rambouillet v1.0*)**

3

4 Mazdak Salavati^{1*}, Alex Caulton^{2,3*}, Richard Clark⁴, Iveta Gazova^{1,5}, Timothy P. L.
5 Smith⁶, Kim C. Worley⁷, Noelle E. Cockett⁸, Alan L. Archibald¹, Shannon M. Clarke²,
6 Brenda M. Murdoch⁹, Emily L. Clark^{1§} and The Ovine FAANG Project Consortium

7

8 ¹The Roslin Institute and Royal (Dick) School of Veterinary Studies, The University of
9 Edinburgh, Edinburgh, UK

10 ²AgResearch, Invermay Agricultural Centre, Mosgiel, New Zealand

11 ³University of Otago, Dunedin, New Zealand

12 ⁴Genetics Core, Edinburgh Clinical Research Facility, The University of Edinburgh,
13 Edinburgh, UK

14 ⁵MRC Human Genetics Unit, The University of Edinburgh, Edinburgh, UK

15 ⁶USDA, Agricultural Research Service, USMARC, Clay Center, Nebraska, USA

16 ⁷Baylor College of Medicine, Houston, Texas, USA

17 ⁸Utah State University, Logan, Utah, USA

18 ⁹University of Idaho, Moscow, Idaho, USA

19

20 *These two authors contributed equally to the work

21 [§]Corresponding author: emily.clark@roslin.ed.ac.uk

22

23 **Abstract**

24

25 The overall aim of the Ovine FAANG project is to provide a comprehensive
26 annotation of the new highly contiguous sheep reference genome sequence (*Oar*
27 *rambouillet v1.0*). Mapping of transcription start sites (TSS) is a key first step in
28 understanding transcript regulation and diversity. Using 56 tissue samples collected
29 from the reference ewe Benz2616 we have performed a global analysis of TSS and
30 TSS-Enhancer clusters using Cap Analysis Gene Expression (CAGE) sequencing.
31 CAGE measures RNA expression by 5' cap-trapping and has been specifically
32 designed to allow the characterization of TSS within promoters to single-nucleotide
33 resolution. We have adapted an analysis pipeline that uses TagDust2 for clean-up
34 and trimming, Bowtie2 for mapping, CAGEfightR for clustering and the Integrative
35 Genomics Viewer (IGV) for visualization. Mapping of CAGE tags indicated that the
36 expression levels of CAGE tag clusters varied across tissues. Expression profiles
37 across tissues were validated using corresponding polyA+ mRNA-Seq data from the
38 same samples. After removal of CAGE tags with < 10 read counts, 39.3% of TSS
39 overlapped with 5' ends of transcripts, as annotated previously by NCBI. A further
40 14.7% mapped to within 50bp of annotated promoter regions. Intersecting these
41 predicted TSS regions with annotated promoter regions (± 50 bp) revealed 46% of the
42 predicted TSS were 'novel' and previously un-annotated. Using whole genome
43 bisulphite sequencing data from the same tissues we were able to determine that a
44 proportion of these 'novel' TSS were hypo-methylated (32.2%) indicating that they
45 are likely to be reproducible rather than 'noise'. This global analysis of TSS in sheep
46 will significantly enhance the annotation of gene models in the new ovine reference
47 assembly. Our analyses provide one of the highest resolution annotations of
48 transcript regulation and diversity in a livestock species to date.

49

50 **Key words: Ovine, TSS, CAGE-Seq, WGBS, promotor, enhancer, RNA,**
51 **transcriptome, FAANG, methylation, mRNA-Seq**

52

53

54

55

56 Introduction

57

58 The Functional Annotation of Animal Genomes (FAANG) consortium is a concerted
59 international effort to use molecular assays, developed during the Human ENCODE
60 project (Birney et al., 2007), to annotate the majority of functional elements in the
61 genomes of domesticated animals (Andersson et al., 2015; Giuffra and Tuggle,
62 2019). Towards this aim the overarching goal of the Ovine FAANG project (Murdoch,
63 2019) is to provide a comprehensive annotation of the new highly contiguous
64 reference genome for sheep, *Oar rambouillet v1.0*
65 (https://www.ncbi.nlm.nih.gov/assembly/GCF_002742125.1/). The Ovine FAANG
66 project is developing a deep and robust dataset of expressed elements and
67 regulatory features in the sheep genome as a resource for the livestock genomics
68 community. Here we describe a global analysis of transcription start sites (TSS)
69 using Cap Analysis Gene Expression (CAGE) sequencing.

70 CAGE measures RNA expression by 5' cap-trapping to identify the 5' ends of
71 non-polyadenylated RNAs including lncRNAs and miRNAs, and has been specifically
72 designed to allow the characterization of TSS within promoters to single-nucleotide
73 resolution (Takahashi et al., 2012). This level of resolution allows investigation of the
74 regulatory inputs driving transcript expression, and construction of transcriptional
75 networks to study, for example, the genetic basis for disease susceptibility (Baillie et
76 al., 2017) or for systematic analysis of transcription start sites through development
77 (Lizio et al., 2017). Using CAGE sequencing technology, the FANTOM5 consortium
78 generated a comprehensive annotation of TSS for the human genome, which
79 included the major primary cell and tissue types (Forrest et al., 2014).

80 The goal of this study was to generate a comprehensive annotation of TSS
81 and TSS-Enhancer clusters for the ovine genome. Our approach was to perform
82 CAGE analysis on 55 tissues and one type of primary immune cell (alveolar
83 macrophages). Tissues representing all the major organ systems were collected from
84 Benz2616, the Rambouillet ewe used to generate the *Oar rambouillet v1.0* reference
85 assembly. CAGE tags for each tissue sample clustered with a high level of specificity
86 according to their expression profiles as measured by RNA-Seq. Mapping of CAGE
87 tags indicated that a large proportion of detected TSS did not overlap with the current
88 annotated 5' end of transcripts. The reproducibility of these 'novel' TSS was tested
89 using whole genome DNA methylation profiles from a subset of the same tissues.

90 DNA methylation plays a key role in the regulation of gene expression and the
91 maintenance of genome stability (Ibeagha-Awemu and Zhao, 2015), and is the most
92 highly studied epigenetic mark. In mammalian species, DNA methylation occurs

93 primarily at cytosine-phosphate-guanine dinucleotides (CpG) and to a lesser extent
94 at CHH and CHG sites (where C = Cytosine; H = Adenine, Guanine, or Thymine; and
95 G = Guanine) (An et al., 2018). Generally, DNA methylation in the promoter region of
96 genes represses transcription, inhibiting elongation by transcriptional machinery.
97 Methylation over TSS blocks transcription initiation; while, conversely, methylation
98 within gene bodies stimulates elongation and influences alternative splicing of
99 transcripts (Jones, 2012; Lev Maor et al., 2015; An et al., 2018). Using DNA
100 methylation profiles, we were able to determine the proportion of 'novel' TSS in our
101 dataset that were likely true signals of transcription initiation based on a hypo-
102 methylated state rather than being an artefact of CAGE-sequencing.

103 We provide the annotation of TSS in the ovine genome as tracks in a genome
104 browser via the Track Hub Registry and visualise these in the R package GViz,
105 ensuring the data is accessible and useable to the livestock genomics community.
106 The global analysis of TSS we present here will significantly enhance the annotation
107 of gene models in the new ovine reference assembly demonstrating the utility of the
108 datasets generated by the Ovine FAANG project and providing a foundation for
109 future work.

110

111 **Methods**

112

113 **Animals**

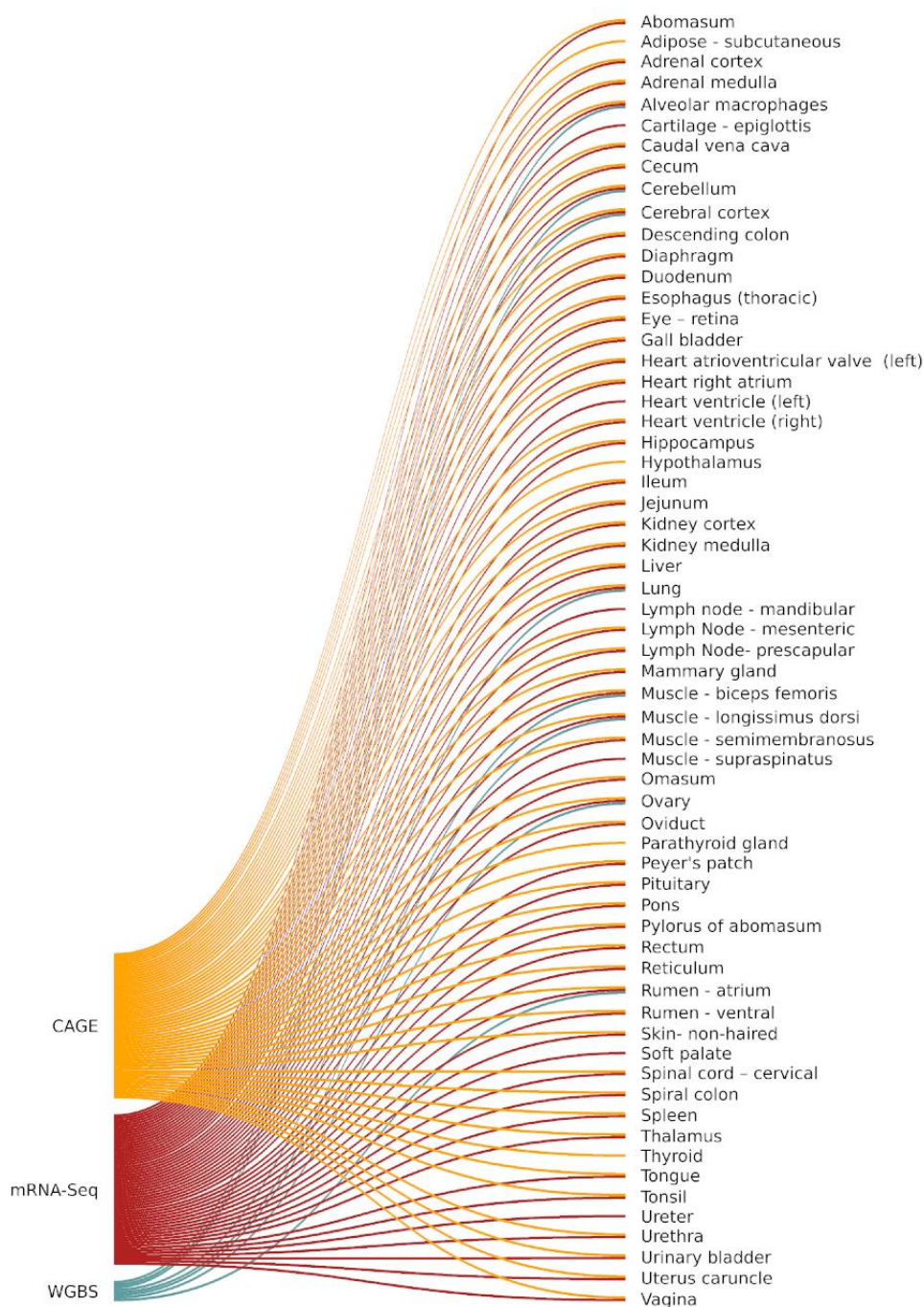
114 Tissues were collected from an adult female Rambouillet sheep at the Utah
115 Veterinary Diagnostic Laboratory on April 29, 2016. At the time of sample collection
116 Benz2616 was approximately 6 years of age and after a thorough veterinary
117 examination confirmed to be healthy. Benz 2616 was donated to the project by the
118 USDA. Sample collection methods were planned and tested over 15 months in 2015
119 to 2016, a description of these is available via the FAANG Data Coordination Centre
120 [https://data.faang.org/api/fire_api/samples/USU_SOP_Ovine_Benz2616_Tissue_Coll](https://data.faang.org/api/fire_api/samples/USU_SOP_Ovine_Benz2616_Tissue_Collection_20160426.pdf)
121 [ection_20160426.pdf](https://data.faang.org/api/fire_api/samples/USU_SOP_Ovine_Benz2616_Tissue_Collection_20160426.pdf) .

122

123 **Sample collection**

124 Necropsy of Benz2616 was performed by a veterinarian to ensure proper
125 identification of tissues, and a team of scientists on hand provided efficient and rapid
126 transfer of tissue sections to containers which were snap frozen in liquid nitrogen
127 prior to transfer to -80C for long-term storage. Alveolar macrophages were collected
128 by bronchoalveolar lavage as described in (Cordier et al., 1990). Details of all 100
129 samples collected from Benz2616 are included in the BioSamples database under

130 submission GSB-7268, group accession number SAMEG329607
131 (<https://www.ebi.ac.uk/biosamples/samples/SAMEG329607>) and associated
132 information is recorded according to FAANG metadata specifications (Harrison et al.,
133 2018). The FAANG assays, as described below, were generated from a subset of
134 tissues for CAGE (56 tissues), polyA+ mRNA-Seq (58 tissues) and WGBS (8 tissues)
135 (Figure 1).



136
137 *Figure 1. FAANG assays (CAGE, WGBS and mRNA-Seq) performed on each tissue*
138 *from Benz2616.*
139

140 **CAGE Library Preparation and Analysis**

141

142 **RNA Isolation for CAGE library preparation**

143 Frozen tissues (60-100mg per sample) were homogenised by grinding with a mortar
144 and pestle on dry ice and RNA was isolated using TRIzol Reagent (Invitrogen)
145 according to the manufacturer's instructions. After RNA isolation 10ug of RNA per
146 sample was treated with DNase I (NEB) then column purified using a RNeasy
147 MinElute kit (Qiagen), according to the manufacturer's instructions. Full details of the
148 RNA extraction protocol are available via the FAANG Data Coordination
149 Centre [https://data.faang.org/api/fire_api/assays/USDA_SOP_RNA_Extraction_Fro](https://data.faang.org/api/fire_api/assays/USDA_SOP_RNA_Extraction_From_Tissue_20180626.pdf)
150 [m_Tissue_20180626.pdf](https://data.faang.org/api/fire_api/assays/USDA_SOP_RNA_Extraction_From_Tissue_20180626.pdf) . Each RNA sample was run on an Agilent BioAnalyzer to
151 ensure RNA integrity was sufficiently high (RIN^e>6). Details of RNA purity metrics for
152 each sample are included in Supplementary Table 1. RNA samples were then stored
153 at -80°C for downstream analysis.

154

155 **CAGE library preparation and sequencing**

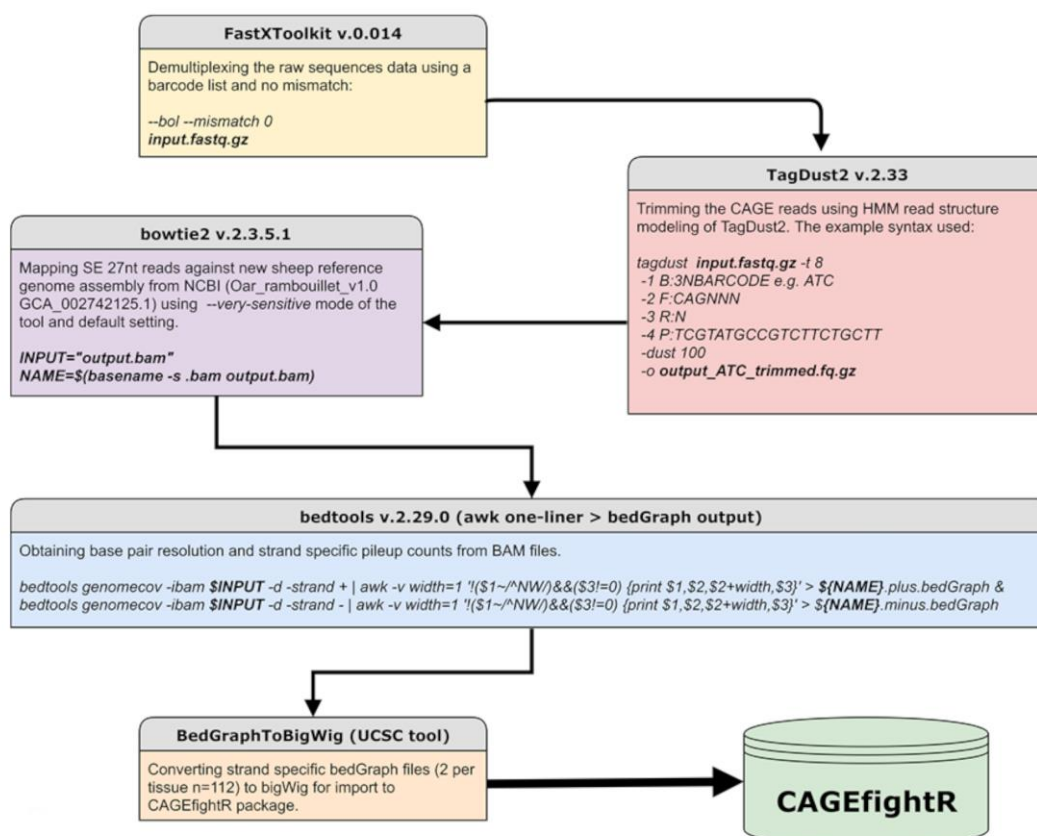
156 CAGE libraries were prepared for each sample as described in (Takahashi et al.,
157 2012) from a starting quantity of 5ug of DNase treated total RNA. Random primers
158 were used to ensure conversion of all 5' cap-trapping RNAs according to (Takahashi
159 et al., 2012). The full protocol is available via the FAANG Data Coordination Centre
160 [https://data.faang.org/api/fire_api/assays/ROSLIN_SOP_CAGE-library-](https://data.faang.org/api/fire_api/assays/ROSLIN_SOP_CAGE-library-preparation_20190903.pdf)
161 [preparation_20190903.pdf](https://data.faang.org/api/fire_api/assays/ROSLIN_SOP_CAGE-library-preparation_20190903.pdf) . Libraries were prepared in batches of eight and pooled.
162 Sequencing was performed on the Illumina HiSeq 2500 platform by multiplexing 8
163 samples on one lane to generate approximately 20 million 50bp single-end reads per
164 sample. Eight of the available fifteen 5' linker barcodes from (Takahashi et al., 2012)
165 were used for multiplexing: ACG, GAT, CTT, ATG, GTA, GCC, TAG and TGG. In
166 total 8 separate library pools were generated and spread across two HiSeq 2500 flow
167 cells. Details of barcodes assigned to each sample and pool IDs are included in
168 Supplementary Table 1.

169

170 **Processing and mapping of CAGE libraries**

171 All sequence data were processed using in house scripting (bash and R) on the
172 University of Edinburgh high performance computing facility (Edinburgh, 2020). The
173 analysis protocol for CAGE is available
174 via [https://data.faang.org/api/fire_api/analysis/ROSLIN_SOP_CAGE_analysis_pipeli](https://data.faang.org/api/fire_api/analysis/ROSLIN_SOP_CAGE_analysis_pipeline_20191029.pdf)
175 [ne_20191029.pdf](https://data.faang.org/api/fire_api/analysis/ROSLIN_SOP_CAGE_analysis_pipeline_20191029.pdf) and summarised in Figure 2. To de-multiplex the data we used the
176 FastX toolkit version 0.014 (Hannon Lab, 2017) for short read pre-processing. We

177 then used TagDust2 v.2.33 (Lassmann, 2015) to extract mappable reads from the
 178 raw data and for read clean-up to remove the *EcoP1* site and barcode, according to
 179 the recommendations of the FANTOM5 consortium e.g. (Bertin et al., 2017). This
 180 process resulted in cleaned approximately 27bp reads (hereafter referred to as
 181 CAGE tags) which were mapped to the Rambouillet Benz2616 genome available
 182 from NCBI (*Oar rambouillet v1.0* GCA_002742125.1) using Bowtie2 v.2.3.5.1 in --
 183 very-sensitive mode equivalent to options `-D 20 -R 3 -N 0 -L 20 -i S,1,0.50`
 184 (Langmead and Salzberg, 2012). The mapped BAM files were then processed for
 185 base pair resolution strand specific read counts using bedtools v.2.29.0 (Quinlan and
 186 Hall, 2010). In order for the bedGraph files to be used in the CAGEfightR package
 187 they were converted to bigWig format using UCSCs tool BedGraphToBigWig (Kent et
 188 al., 2010).



189
 190 *Figure 2. Workflow of the analysis pipeline and respective tools used for CAGE*
 191 *sequence data analysis*

192
 193 **Normalisation and mapping of CAGE tags**

194 For normalisation and clustering of CAGE tags (as CAGE Tags-Per-Million Mapped:
 195 CTPM) we used the software package CAGEfightR v.1.5.1 (Thodberg and Sandelin,
 196 2019). The normalisation was performed via dividing CAGE tag counts in each

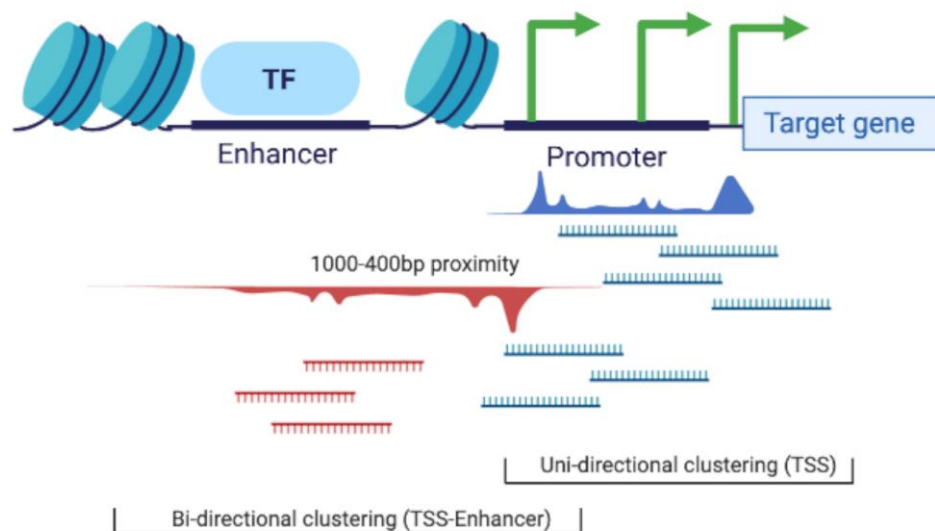
197 predicted cluster by the total mapped CAGE tags in the sample, multiplied by 1.0e6.
198 To perform these analyses we created a custom BSgenome object (a container of
199 the genomic sequence) for sheep from *Oar rambouillet v1.0* using the BSgenome
200 Bioconductor package v.1.53.1 (Pages, 2020). Distribution metrics of CAGE tags
201 across the genome were annotated and analysed using the TxDB transcript ID
202 assignment and Genomic Features package v.1.36.4 (Lawrence et al., 2013). The
203 TxDB object was created using the NCBI gff3 gene annotation file from NCBI *Oar*
204 *rambouillet v1.0* GCA_002742125.1 (*GCF_002742125.1_Oar rambouillet v1.0*
205 *_genomic.gff release 103*).

206

207 **Clustering of CAGE tags**

208 To annotate TSS in the *Oar rambouillet v1.0* genome assembly we first generated
209 expression read counts for each tag (bp resolution). Any tags with read counts < 10
210 (in each tissue) that were not present in at least 37/56 tissues (i.e. two thirds of the
211 tissues) were removed. This conservative representation threshold was introduced to
212 ensure CAGE tags included in downstream analysis were reproducible. In the
213 absence of additional biological replicates we based this on the assumption that a
214 CAGE tag was more likely to be reproducible if it was shared across multiple tissues,
215 although sensitivity to putative highly tissue-specific TSS is reduced (see below).
216 Gene annotation from the NCBI's GTF file was used to validate the coordinates of
217 predicted CAGE clusters (i.e. residing within or outside the promoter of annotated
218 genes). Five thresholds for representation, of CAGE tags across tissues, were
219 compared (1 tissue, 1/3rd of the tissues, half of the tissues, 2/3rd of the tissues and all
220 of the tissues). The proportion of CAGE tag clusters within (tagged by unique gene
221 IDs) or outside the promoter region (untagged) was used to compare each threshold.
222 Including all of the tissues (56/56 representation) resulted in 90.4% loss of genes
223 tagged by any CAGE tag cluster i.e. only 2974 genes from a total of 30,862 genes in
224 *Oar rambouillet v1.0* were tagged by any CAGE tag cluster in the transcriptome.
225 Reducing the threshold further to 1/3rd of tissues resulted in a high proportion of
226 untagged CAGE tag clusters (40.9%) and 18,244 tagged genes. A less stringent
227 reduction of the threshold to 2/3rd (37/56 tissues) resulted in 14,105 genes tagged by
228 any CAGE tag cluster. The 2/3rd representation threshold was therefore chosen to
229 maximise the number of annotated genes with expressed CAGE tags that were
230 shared across tissues and minimise the number of CAGE tags mapped to outside the
231 promoter region. Further details of this comparison are included in Supplementary
232 File 1 Section 3. TSS expression profiles (as CTPM) were then generated for each
233 tissue using the CAGEfightR v. 1.5.1 quickTSS, quickEnhancers and findLinks

234 functions (Thodberg and Sandelin, 2019). The CAGE tags clustered A) uni-
235 directionally (according to the sense or anti-sense flag of the mapped CAGE tag) into
236 predicted TSS and B) bi-directionally, using the TSS-Enhancer detection algorithm
237 from CAGEfightR (Thodberg and Sandelin, 2019), into correlated TSS and enhancer
238 (TSS-Enhancer) clusters. Bi-directional (TSS-Enhancer) clusters are defined as
239 clusters of CAGE tags that are located on the opposing strand within 400 bp-1 Kbp
240 proximity of the centre of a promoter (Thodberg and Sandelin, 2019). The bi-
241 directional clusters outside of this range were excluded from this analysis according
242 to the previously described method in (Thodberg et al., 2019). The concept of uni-
243 directional and bi-directional clustering is illustrated in Figure 3.



244
245 *Figure 3. Schematic representation of the two clustering algorithms used in the*
246 *CAGEfightR package for TSS (uni-directional) and TSS-Enhancer (bi-directional)*
247 *clustering.*

248

249 **Identification of shared TSS or TSS-enhancer clusters across tissues**

250 TSS or TSS-Enhancer clusters that were shared across tissues, were identified by
251 investigating the CTPM expression profile of each of the tissues using correlation
252 based and mutual information (MI) distance matrices (Priness et al., 2007; Reshef et
253 al., 2018). This method of MI based clustering tolerates missingness and outlier-
254 induced grouping errors in gene expression profiles (Priness et al., 2007). Using this
255 method, we assumed that the CTPM expression profile, for each cluster, could vary
256 across tissues, but for a predicted TSS or TSS-Enhancer cluster to be considered
257 reproducible it must be present in at least two thirds of the tissues (37/56) in the
258 dataset.

259

260 **Identification of tissue-specific TSS or TSS-enhancer clusters**

261 The 2/3rd representation threshold applied above would remove all tissue-specific
262 CAGE tag clusters. To overcome this, a rerun of the clustering algorithm was
263 performed with the representation threshold reduced. Tissue-specific uni-directional
264 TSS clusters that were only present in 1/56 tissues were identified by filtering for
265 CAGE tags with >10 expressed counts to create a data frame. The data frame was
266 then filtered tissue-by-tissue to only retain uni-directional TSS clusters present in
267 each tissue separately. This process was then repeated for the TSS-Enhancer
268 clusters.

269

270 **Annotation of 'novel' TSS in the ovine genome**

271 We expected given the diversity of tissues sampled that we would detect a significant
272 number of 'novel', previously unannotated TSS. The CAGE tag uni-directional
273 clusters (TSS) were annotated using the `mergeByOverlay` function of the
274 `GenomicFeatures` package in R and the custom `TxDB` object as following:

275 *mergeByOverlaps(subject = TSS, query = promoters(txdb, upstream = 25,*
276 *downstream = 25, use.names = T,c("tx_name", "GENEID")), maxgap = 25, type =*
277 *"any")*. The `TxDB` object calculates the range of the promoter based on the 5'UTR
278 and first CDS codon coordinates. In each tissue any TSS region within 50bp range of
279 the promoter coordinate of a gene model was considered 'annotated'. In addition, we
280 expanded this range to 400bp to determine whether this would identify significantly
281 more unannotated TSS further from the promoter. A reverse sub setting of the 50bp
282 window region was performed as follows: *subsetByOverlaps(x = TSS ,ranges =*
283 *annotated, invert = TRUE)*. These regions were considered 'novel' TSS previously
284 unannotated in the assembly. This process was repeated for every tissue separately
285 (n=56).

286

287 **Comparative analysis of WGBS and CAGE Data**

288

289 **Preparation of genomic DNA from tissue**

290 Extraction of DNA for bisulphite sequencing was performed using
291 phenol:chloroform:isoamyl alcohol method. Briefly, approximately 1 g frozen tissue
292 was pulverized and resuspended in 2.26 ml of digestion buffer (10 mM Tris-HCl, 400
293 mM NaCl, 2 mM EDTA, pH 8.0) with 200 µl of SDS 10% and 60 µl RnaseA
294 (10mg/ml) (Sigma-Aldrich, St. Louis, MO, USA) RNA degradation proceeded for one
295 hour at 37°C with gentle shaking. Next, 25 µl of proteinase K (20mg/ml) (Sigma-

296 Aldrich) was added to the suspension and incubated overnight (approximately 16
297 hours) at 37 °C with gentle shaking. The viscous lysate was transferred to a 2 mL
298 Phase Lock tube (VWR, Radnor, PA) and extracted twice with Tris-HCl-saturated
299 phenol:chloroform:isoamyl alcohol (25:24:1) pH 8.0, followed by extraction with 2.5
300 ml chloroform. The DNA was precipitated by addition of 5.5 ml of 100% ethanol and
301 250 µl of 3M Sodium Acetate to the aqueous phase in a 15 mL conical tube, mixed
302 by gentle inversion until the DNA became visible. The DNA was removed with a bent
303 Pastuer pipette hook, washed in 5 ml 70% cold ethanol, air dried then resuspended
304 in 250 µl – 1 ml of 1X TE and stored at –20 °C until use. DNA concentration was
305 quantified fluorometrically on the Qubit® 3.0 Fluorometer (Thermo Fisher Scientific,
306 Waltham, MA, United States) using the Qubit dsDNA HS Assay Kit. The purity of the
307 extractions was determined via 260/280 and 260/230 ratios measured on the
308 NanoDrop 8000 (Thermo Fisher Scientific) and DNA integrity was assessed by 1%
309 agarose gel electrophoresis. The protocol is available via the FAANG Data
310 Coordination Centre
311 [https://data.faang.org/api/fire_api/assays/USDA_SOP_DNA_Extraction_From_Whole](https://data.faang.org/api/fire_api/assays/USDA_SOP_DNA_Extraction_From_Whole_BloodandLiver_20200611.pdf)
312 [BloodandLiver_20200611.pdf](https://data.faang.org/api/fire_api/assays/USDA_SOP_DNA_Extraction_From_Whole_BloodandLiver_20200611.pdf).

313

314 **Whole Genome Bisulphite Conversion and Sequencing**

315 Library preparation and sequencing of seven tissues and 1 cell type (Figure 1),
316 selected to include a representative from all major organ systems, were performed by
317 The Garvan Institute of Medical Research, Darlinghurst, Sydney, New South Wales.
318 Un-methylated lambda DNA was added at 0.5% of the total sample DNA
319 concentration prior to bisulphite conversion as a conversion efficiency control. DNA
320 conversion was carried out using the EZ DNA Methylation-Gold Kit (Zymo Research,
321 CA, USA) following the manufacturer's instructions. The Accel-NGS Methyl-seq DNA
322 kit (Swift Biosciences, MI, USA) for single indexing, was used to prepare the libraries,
323 following the manufacturer's instructions. Libraries were pooled together and
324 sequenced across 6 lanes of a flow-cell on an Illumina HiSeq X platform using paired
325 end chemistry for 150 bp reads (min 10X coverage). The protocol is available via
326 FAANG Data Coordination Centre
327 [https://data.faang.org/api/fire_api/assays/AGR_SOP_WGBS_AgR_Library_prep_202](https://data.faang.org/api/fire_api/assays/AGR_SOP_WGBS_AgR_Library_prep_20200610.pdf)
328 [00610.pdf](https://data.faang.org/api/fire_api/assays/AGR_SOP_WGBS_AgR_Library_prep_20200610.pdf)

329

330 **WGBS data processing**

331 Paired end Illumina WGBS sequence data was processed and analysed using in
332 house scripting (bash and R) and a range of purpose-built bioinformatics tools on the

333 AgResearch and University of Edinburgh high performance computing facilities. The
334 analysis protocol for WGBS is available via the FAANG Data Coordination
335 Centre [https://data.faang.org/api/fire_api/analysis/AGR_SOP_WGBS_AgR_data_an](https://data.faang.org/api/fire_api/analysis/AGR_SOP_WGBS_AgR_data_analysis_20200610.pdf)
336 [alysis_20200610.pdf](https://data.faang.org/api/fire_api/analysis/AGR_SOP_WGBS_AgR_data_analysis_20200610.pdf) and summarised below.

337 Briefly, FASTQ files for each sample, run across multiple lanes were merged
338 together. TrimGalore v. 0.5.0. (<https://github.com/FelixKrueger/TrimGalore>) was used
339 to trim raw reads to remove adapter oligos, poor quality bases (phred score less than
340 20) and the low complexity sequence tag introduced during Accel-NGS Methyl-seq
341 DNA kit library preparation as follows: `trim_galore -q 20 --fastqc --paired --clip_R2`
342 `18 --three_prime_clip_R1 18 --retain_unpaired -o Trim_out INPUT_R1.fq.gz`
343 `INPUT_R2.fq.gz`

344 A bisulphite-sequencing amenable reference genome was built using the *Oar*
345 *rambouillet v1.0*, GenBank Accession number: GCA_002742125.1 genome with the
346 BSSeeker2 script `bs_seeker2-build.py` using bowtie v2.3.4.3 (Langmead and
347 Salzberg, 2012). and default parameters. The Enterobacteria phage lambda genome
348 available from NCBI (Accession number: NC_001416) was added to the
349 *Oar_rambouillet v1.0* genome as an extra chromosome to enable alignment of the
350 unmethylated lambda DNA conversion control reads. Paired-end, trimmed reads
351 were aligned to the reference genome using the BSSeeker2 script `bs_seeker2-`
352 `align.py` and bowtie v2.3.4.3 (Langmead and Salzberg, 2012) allowing four
353 mismatches (-m 4). Aligned bam files were sorted with samtools v1.6 (Li et al., 2009)
354 and duplicate reads were removed with picard tools v2.17.11
355 (<https://broadinstitute.github.io/picard/>) MarkDuplicates function.

356 Deduplicated bam files were used to call DNA methylation levels using the
357 “bam2cgmap” function within CGmaptools (Guo et al., 2018) with default options to
358 generate ATCGmap and CGmap files for each sample. The ATCGmap file format
359 summarises mapping information for all covered nucleotides on both strands, and is
360 specifically designed for BS-seq data; whilst the CGmap format is a more condensed
361 summary providing sequence context and estimated methylation levels at any
362 covered cytosine in the reference genome.

363 Hyper-methylated and hypo-methylated regions were determined for each
364 sample using methpipe v3.4.3 (Song et al., 2013). Specifically, CGmap files for each
365 sample were reformatted for the methpipe v3.4.3 workflow using custom awk scripts.
366 The methpipe symmetric-cpgs program was used to merge individual methylation
367 levels at symmetric CpG pairs. Hypo-methylated and hyper-methylated regions were
368 determined using the hmr program within methpipe, which uses a hidden Markov

369 model (HMM) using a Beta-Binomial distribution to describe methylation levels at
370 individual CpG sites, accounting for the read coverage at each site.

371 Visualisation of the individual CpG site methylation levels with a minimum
372 read depth cut-off of 10x coverage was done using Gviz package v.1.28.3 (Hahne
373 and Ivanek, 2016).

374

375 **Comparative analysis of annotated and ‘novel’ TSS with WGBS methylation** 376 **information**

377 We expected that reproducible TSS, either annotated or novel, would overlap with
378 hypo-methylated regions of the genome (Yamashita et al., 2005; Yagi et al., 2008).
379 To test whether this was true for those identified in our analysis, both annotated and
380 novel TSS from the CAGE BED tracks were intersected with WGBS hypo
381 methylation profiles using bedtools v.2.29.2 (Quinlan and Hall, 2010) and the
382 following script: `bedtools intersect -b WGBS_HypoCpG.bed -a Novel_or_`
383 `Annotated.bed > Novel_or_annotated_HypoCpG.bed`. Any annotated and novel TSS
384 (within a ± 50 bp window of the promoter) that intersected hypomethylated regions of
385 DNA in each tissue, were verified as reproducible TSS and the remainder as ‘noise’.
386 The overlay of these regions was visualised as a genomic track using the Gviz
387 package v.1.28.3 (Hahne and Ivanek, 2016).

388

389 **Visualisation of the annotated TSS, mRNA-Seq and WGBS tracks in the ovine** 390 **genome**

391 In order to confirm the simultaneous expression of mRNA, CAGE tags corresponding
392 to an active TSS and a hypomethylated region of DNA, a genomic track on which all
393 three datasets could be visualised was generated. This visualisation consists of the
394 following tracks: 1) Uni-directional CAGE tag clusters (TSS) 2) Bi-directional CAGE
395 tag clusters (TSS-Enhancers) 3) WGBS hypomethylation score (bp resolution) 4)
396 Transcript level expression (mRNA-Seq [TPM]) 5) The transcript models and 6) The
397 gene model. Areas of the genome where TSS or TSS-Enhancer regions overlapped
398 regions with a high hypomethylation score, within 5' end of an actively expressing
399 transcript (TPM score), were considered reproducible TSS for that tissue. This
400 process was performed using eight tissues with matching mRNA-Seq, CAGE and
401 WGBS sequence data. The Gviz package v.1.28.3 was used to visualise these tracks
402 (Hahne and Ivanek, 2016).

403

404 **Validation of tissue-specific expression profiles**

405

406 **mRNA-Sequencing**

407 Total RNA for mRNA-Seq from 32 tissues (Figure 1) was prepared, as above for the
408 CAGE samples, by USMARC, and for 26 tissues by Baylor College of Medicine
409 (BCM) using the MagMAX mirVana total RNA isolation kit (Thermo Fisher Scientific,
410 Waltham, MA, United States) according to the manufacturer's instructions. Paired
411 end polyA selected mRNA-Seq libraries were prepared and sequenced on an
412 Illumina NextSeq500 at USMARC or the Illumina HiSeq2000 at BCM using the
413 Illumina Tru-Seq Stranded mRNA Library Preparation Kit. For each tissue a set of
414 expression estimates, as transcripts per million (TPM), were obtained using the
415 transcript quantification tool Kallisto v0.43.0 (Bray et al., 2016). The mRNA-Seq
416 analysis pipeline is accessible via the FAANG Data Coordination Centre
417 [https://data.faang.org/api/fire_api/analysis/ROSLIN_SOP_RNA-](https://data.faang.org/api/fire_api/analysis/ROSLIN_SOP_RNA-Seq_analysis_pipeline_20200610.pdf)
418 [Seq_analysis_pipeline_20200610.pdf](https://data.faang.org/api/fire_api/analysis/ROSLIN_SOP_RNA-Seq_analysis_pipeline_20200610.pdf). A pairwise distance matrix (multiple
419 correlation coefficient based) was produced using MI values for all tissues and a
420 dendrogram of tissues was created in order to visualise grouping patterns of tissues
421 with similar mRNA expression profiles, and for comparison with the CAGE dataset.

422

423 **Comparative analysis of tissue-specific expression profiles using information** 424 **from CAGE and mRNA-Seq**

425 We assessed whether TSS expression profiles from the CAGE dataset were
426 biologically meaningful using the mutual information (MI) sharing algorithm (Joe,
427 1989). Tissues with the same function and physiology should have similar TSS
428 expression profiles. The CTPM expression level was binned (n=10) using the bioDist
429 package v.1.56.0 (Ding et al., 2012) and mutual information (MI) for each pair of
430 tissue samples was calculated as in (Joe 1989). :

431

$$432 \quad \delta = (1 - \exp(-2 \times \delta))^{0.5}$$

433

$$434 \quad MI \text{ distance} = 1 - \delta$$

435 A pairwise distance matrix (multiple correlation coefficient based) was produced
436 using MI values for all tissues and a dendrogram of tissues created to visualise
437 grouping patterns of tissues with similar TSS expression profiles. If the expression
438 profiles were meaningful then tissues with similar function and physiology would
439 group together in clades within the dendrogram. These tissue specific groupings
440 were then further validated by comparison with mRNA-Seq data for the same
441 samples, using the MI sharing algorithm and dendrogram approach.

442

443 **Results**

444

445 **Library size and annotation metrics**

446 The mean CAGE library depth based on raw CAGE tags was 4,862,957 tags. Library
447 depth varied across tissues. Tissues with low depth were not related to any specific
448 barcodes and were evenly spread over the two sequencing runs (Supplemental
449 Table S1), suggesting random variation rather than systematic differences due to
450 specific barcodes or sequencing runs. The RIN^e values were also consistently >7 for
451 all tissues with low counts, indicating RNA integrity was also unlikely to be affecting
452 library depth. Differences in tag numbers are therefore more likely to relate to
453 variation in efficiency between individual libraries or tissue-specific differences
454 related to the physiology of the tissue.

455

456 **CAGE tag clustering and annotation by genomic regions**

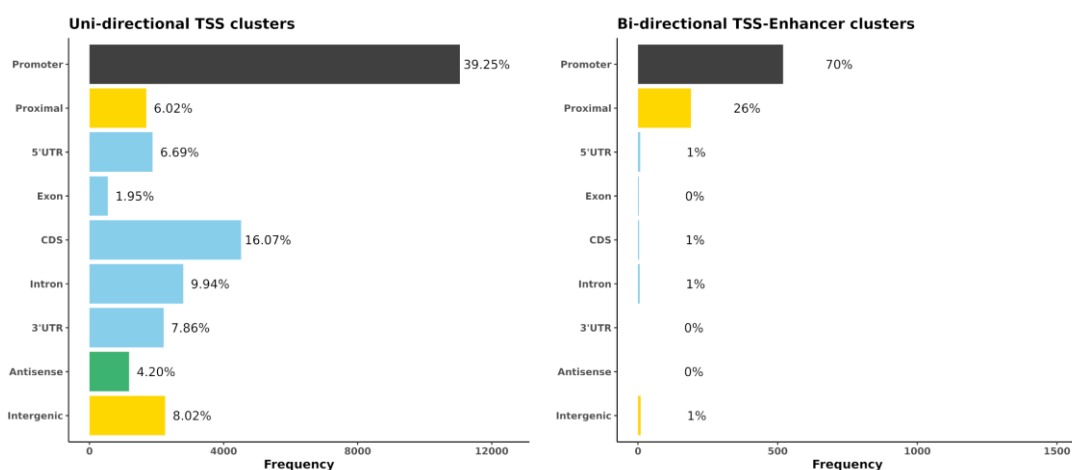
457 We used a newly developed software package to annotate TSS in the Rambouillet
458 Benz2616 genome (Thodberg and Sandelin, 2019; Thodberg et al., 2019) which
459 clustered the CAGE tags as A) uni-directionally into predicted TSS or B) bi-
460 directionally into correlated TSS and enhancer (TSS-Enhancer) clusters (Figure 3).
461 The clustered CAGE tags were filtered to remove any clusters with a minimum
462 expression level of <10 tag counts.

463 In order to reduce 'noise' for downstream analysis (i.e. large proportions of
464 un-annotated CAGE tags that were likely to be spurious) we applied a conservative
465 2/3rd representation criteria, i.e. a minimum of 37/56 tissues had to express the tag
466 cluster with >10 expressed counts. Application of 2/3rd representation criteria resulted
467 in 28,148 uni-directional TSS clusters, from a total of 5,450,864, for downstream
468 analysis. The mean (\pm SD) and median number of tissues per cluster was 3.68 ± 4.78
469 and 2, respectively. This level of noise in CAGE sequencing datasets (0.5% retained
470 clusters) is somewhat lower than reported for other mammalian promoter-level
471 expression atlas projects, e.g. by the FANTOM consortium, using less conservative
472 criteria, where approximately 5% of clusters were retained (Forrest et al., 2014).

473 Bi-directional TSS-enhancer clusters were far fewer in number, although
474 retention was higher with over 23% meeting the same 2/3rd representation criteria
475 741 from a total of 3,131. Though fewer in number these bi-directional (or TSS-
476 enhancer) clusters are functionally important in the regulation of expression of their
477 target genes (Andersson et al., 2014; Thodberg and Sandelin, 2019), consistent with
478 finding them in over 2/3rd of tissues. The co-expression of leading enhancer RNA

479 (eRNA) which is captured by CAGE sequencing can provide a map to enhancer
480 families in the genome and the genes under their regulation (Andersson et al., 2014).

481 The locations of both uni-directional TSS and bi-directional TSS-enhancer
482 clusters were identified in *Oar rambouillet v1.0* and the proportion of TSS clusters
483 located within or near annotated gene features was estimated (Figure 4). The custom
484 BSgenome and TxDB objects created from the GFF3 file format provide detailed
485 calculated coordinates for the following sections: intergenic (>1000bp before 5'UTR
486 or after the end of 3'UTR), proximal (1000bp upstream of the 5'UTR), promoter
487 (± 100 bp from 5'UTR) and the standard gene model (5'UTR, exon, intron and 3'UTR).
488 The genomic region class with the highest number of unidirectional clusters (39.25%)
489 was the promoter regions (± 100 bp from 5'UTR) (Figure 4A), with a relatively even
490 distribution within the other regions of the genome, including 6% mapping proximally
491 to the 5'UTR. The majority of bi-directional TSS-enhancer clusters were also located
492 in promoter regions (70.1%) with a smaller proportion (25.6%) located in proximal
493 regions (Figure 4B). The lack of bi-directional TSS-enhancer clusters in other regions
494 is a consequence of the operation of the CAGEfightR algorithm, which only
495 considers bi-directional clusters within a 400-1000bp window of a TSS CAGE tag
496 cluster (Thodberg and Sandelin, 2019; Thodberg et al., 2019). This approach also
497 reduced the total count compared to unidirectional clusters (28,148 uni-directional
498 clusters relative to 741 bidirectional TSS-enhancer clusters across tissues)
499 (Thodberg et al., 2019).



500

501 *Figure 4. The genomic region distribution of CAGE tag clusters mapped against Oar*
502 *rambouillet v1.0 assembly and gene annotation. The counts were averaged across*
503 *tissues. A) Uni-directional TSS clusters with the highest proportion in promoter*
504 *region (± 100 bp of the 5'UTR beginning at the [TSS]). B) Bi-directional TSS-enhancer*
505 *clusters with the highest proportion in the proximal region (1000bp upstream of the*
506 *5'UTR beginning at the [TSS]).*

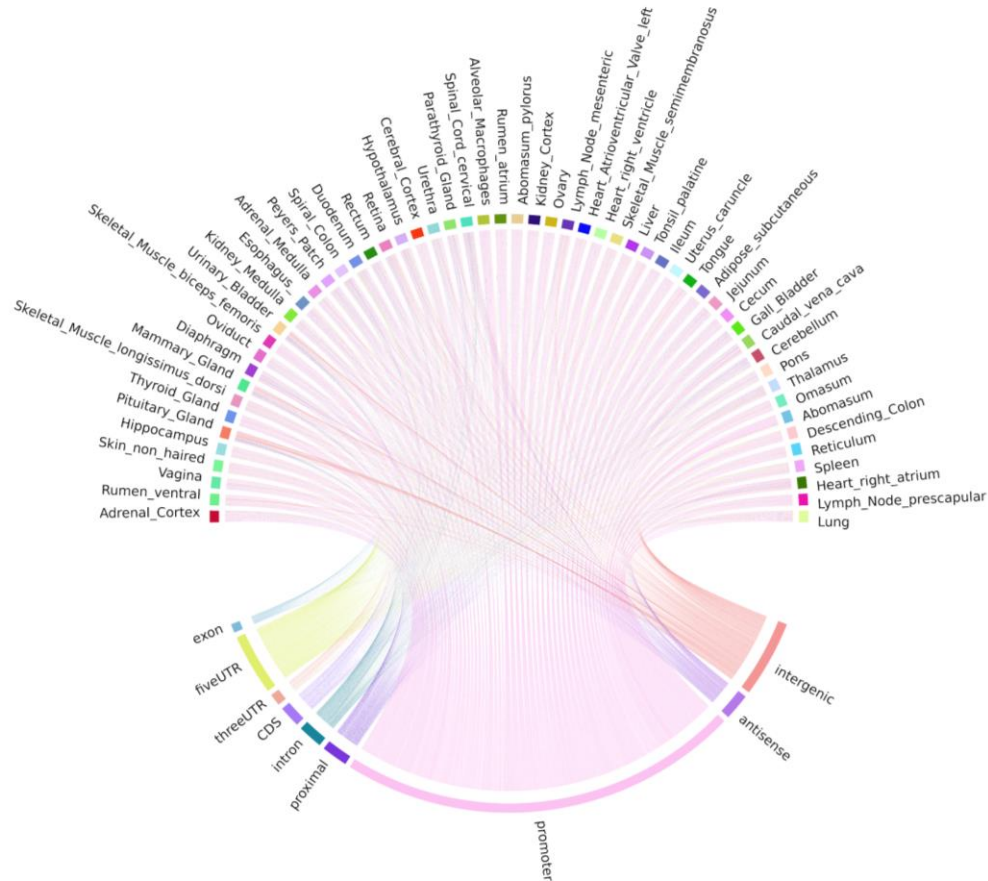
507 **Distribution of CAGE tag clusters in *Oar rambouillet v1.0* relative to *Oar_v3.1***

508 As a proxy for improvement in the accuracy of the gene models in *Oar rambouillet*
509 *v1.0* we investigated how mapped CAGE tag clusters were distributed across
510 genomic features when compared with *Oar_v3.1* (Jiang et al., 2014), the reference
511 genome assembly it superseded (Supplemental Figure S2). The percentage of uni-
512 directional CAGE tag clusters mapping to intergenic regions, which usually occurs
513 due to missing gene model information, was greater for *Oar_v3.1* (33.9%) relative to
514 *Oar rambouillet v1.0* (8%). The percentage of uni-directional CAGE tag clusters
515 mapping to annotated promoter regions was greater for *Oar rambouillet v1.0*
516 (39.25%) compared to *Oar_v3.1* (14.94%), indicating the proportion of accurate gene
517 models in *Oar rambouillet v1.0* was greater. Ensembl annotated 20,921 protein
518 coding genes, 5,843 non-coding genes and 29,118 transcripts on *Oar_v3.1*. NCBI
519 annotated 20,883 protein-coding genes, 7,533 non-coding and 62,535 transcripts on
520 *Oar rambouillet v1.0*. Of the 28,148 unidirectional TSS clusters mapped to *Oar*
521 *rambouillet v1.0*, 87.74% mapped to 13,868 unique genes (31,729 transcripts). In
522 comparison, of the 23,829 unidirectional TSS clusters mapped to *Oar_v3.1*, 49.1%
523 mapped to 6,549 genes (9,914 transcripts). A larger number of TSS-Enhancer CAGE
524 clusters were detected in *Oar_v3.1* (1121) in comparison to *Oar rambouillet v1.0*
525 (741) mapping to 1371 and 2598 unique genes, respectively. A detailed comparison
526 of mapping of the CAGE tags to the two reference assemblies is included in
527 Supplementary File 1, Sections 1 and 2.

528

529 **Mapping of CAGE tags shared across all tissue samples**

530 Correlation-based and mutual information (MI) distance matrices were used to
531 evaluate the occurrence of TSS and enhancer TSS across tissues. The mean \pm SD
532 number of tissues in which each cluster passing the 2/3rd criteria (expressed in 37/56
533 tissues) was (47.73 \pm 6.03). Uni-directional TSS clusters (n=28,148 TSS regions)
534 that were shared across tissues and detected in at least 37/56 tissues are visualised
535 in Figure 5. Each chord in Figure 5 represents the presence of an expressed
536 unidirectional TSS cluster shared across tissues. The majority of the unidirectional
537 TSS that were shared across tissues mapped to promoters (39.25%) and were
538 shared evenly across the tissues sampled (Figure 5). Some tissues e.g. mammary
539 gland, pituitary gland and urinary bladder had more uni-directional TSS mapping to
540 intergenic regions, which might indicate evidence of alternative splicing or differential
541 TSS usage across tissues (Figure 5). Alternative splicing events and differential TSS
542 usage, captured by CAGE, are often not included in the reference gene prediction
543 models (Berger et al., 2019).



544

545 *Figure 5. Chord diagram of expression level (TPM) of CAGE tag clusters (uni-*
546 *directional TSS) across all the tissues collected from Benz2616. Shared CAGE tag*
547 *clusters are common to at least 2/3rd of the tissues (37/56).*

548

549 Bi-directional TSS-Enhancer CAGE clusters were far fewer in number but
550 were shared in a similar pattern across tissues as the uni-directional TSS clusters
551 (Figure 6). The majority (70.1%) of the TSS-Enhancer clusters mapped to promoters
552 (n=520) while 25.6% mapped to 'proximal' regions as expected according to the
553 400bp-1Kbp detection window for TSS-Enhancer clusters from the centre of the
554 promoter (Figure 6). For some tissues including abomasum, spleen and heart right
555 atrium the proportion of bi-directional TSS-Enhancer clusters mapping to proximal
556 regions was greater indicating more enhancer families could be present within these
557 tissues (Figure 6).

558

559

560

561

562

582 tag clusters was very low (<2 CTPM), which combined with the small sample size
 583 (n=1) for each tissue, meant that analysis of tissue-specific TSS was not particularly
 584 meaningful using this dataset. The analysis was repeated for tissue specific TSS-
 585 Enhancer clusters which is detailed in Supplementary Figure S3B.

586

587 **Proportion of ‘novel’ TSS within the CAGE dataset for each tissue**

588 CAGE tag clusters were annotated initially using the *Oar rambouillet v1.0* gene
 589 models from NCBI. A tissue-by-tissue annotation was performed using the same
 590 gene models to identify any CAGE tag clusters within a 50bp window of the promoter
 591 boundaries of every gene. From a total of 23,837 TSS (the average number of TSS
 592 per tissue) we found 11,328 (49.6%) were located within 50bp of the promoter. The
 593 CAGE tag clusters were annotated using the NCBI *Oar rambouillet v1.0* GFF3 gene
 594 track file (version 103) and a TxDB object created in the GenomicFeatures package
 595 (version 1.36.4) in R. CAGE tag clusters within 50bp (short range) or 400bp (long
 596 range) of the promoter were defined as annotated. Supplementary File 2 includes
 597 BED files for these CAGE tag clusters. The percentage of ‘novel’ previously un-
 598 annotated, but likely to be reproducible, CAGE tag clusters for each tissue within
 599 50bp (short range) and 400bp (long range) from the promoter are detailed in Table 1.

600

601 Table 1: The total number and percentage of ‘novel’ CAGE tag clusters for each
 602 tissue within 50bp (short range) and 400bp (long range) from the promoter.

603

<i>Tissue</i>	%	<i>Tags within</i>	<i>Tags within</i>	<i>Total</i>
	<i>Novel</i>	<i>50bp</i>	<i>400bp</i>	
<i>Abomasum</i>	49.38	8,161	8,688	16,584
<i>Abomasum pylorus</i>	49.89	12,339	13,074	25,963
<i>Adipose subcutaneous</i>	51.74	12,336	13,074	26,970
<i>Adrenal cortex</i>	51.79	12,285	13,019	26,859
<i>Adrenal medulla</i>	52.86	11,520	12,210	25,604
<i>Alveolar macrophages</i>	51.19	12,008	12,731	25,845
<i>Caudal vena cava</i>	37.75	10,937	11,578	18,179
<i>Cecum</i>	51.68	12,070	12,801	26,261
<i>Cerebellum</i>	48.59	8,393	8,936	16,796
<i>Cerebral cortex</i>	51.19	12,199	12,917	26,327
<i>Descending colon</i>	51.80	11,830	12,539	25,810
<i>Diaphragm</i>	52.53	10,367	11,016	22,733
<i>Duodenum</i>	52.34	11,243	11,932	24,620
<i>Esophagus</i>	49.90	10,016	10,625	20,741

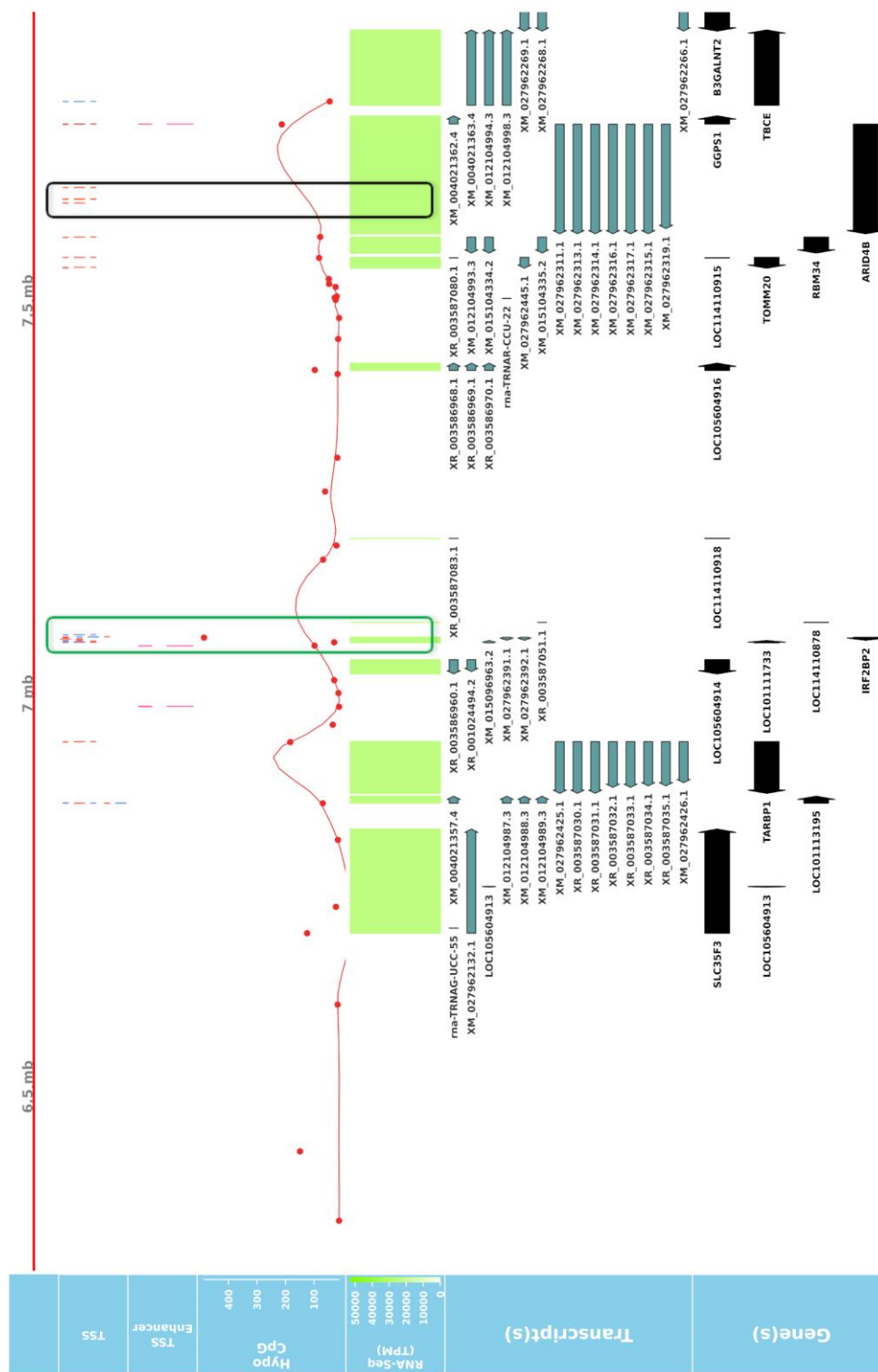
<i>Gall bladder</i>	47.78	11,870	12,578	23,852
<i>Heart atrioventricular valve left</i>	50.90	12,268	13,000	26,330
<i>Heart right atrium</i>	52.96	10,996	11,666	24,444
<i>Heart right ventricle</i>	50.47	12,260	12,987	26,082
<i>Hippocampus</i>	53.40	12,142	12,878	27,451
<i>Ileum</i>	52.45	12,352	13,094	27,411
<i>Jejunum</i>	31.67	10,810	11,418	16,361
<i>Kidney cortex</i>	52.04	12,317	13,057	27,076
<i>Kidney medulla</i>	51.07	10,946	11,618	23,365
<i>Liver</i>	49.35	12,255	12,981	25,459
<i>Lung</i>	52.91	11,644	12,339	25,995
<i>Lymph node mesenteric</i>	46.34	12,132	12,838	23,742
<i>Lymph node prescapular</i>	53.56	11,533	12,228	26,096
<i>Mammary gland</i>	49.75	10,048	10,688	20,774
<i>Omasum</i>	39.89	9,167	9,708	15,692
<i>Ovary</i>	50.79	12,334	13,073	26,434
<i>Oviduct</i>	53.29	11,563	12,260	25,957
<i>Payer's patch</i>	52.41	11,881	12,578	26,240
<i>Pituitary gland</i>	47.25	6,918	7,362	13,400
<i>Pons</i>	40.69	11,622	12,296	20,506
<i>Rectum</i>	53.55	12,002	12,723	27,192
<i>Reticulum</i>	53.39	12,185	12,911	27,589
<i>Retina</i>	53.54	11,805	12,537	26,691
<i>Rumen atrium</i>	50.69	12,335	13,077	26,363
<i>Rumen ventral</i>	40.20	7,109	7,567	12,165
<i>Skeletal muscle biceps femoris</i>	50.23	12,151	12,872	25,715
<i>Skeletal muscle longissimus dorsi</i>	53.67	11,356	12,060	25,748
<i>Skeletal Muscle semimembranosus</i>	51.15	12,262	12,993	26,471
<i>Spinal cord cervical</i>	51.47	11,376	12,050	24,508
<i>Spiral colon</i>	53.25	11,937	12,662	26,813
<i>Spleen</i>	53.46	12,161	12,892	27,568
<i>Thalamus</i>	41.61	11,426	12,079	20,404
<i>Tongue</i>	39.57	9,639	10,244	16,512
<i>Tonsil palatine</i>	46.57	12,178	12,875	23,978
<i>Urethra</i>	52.76	11,387	12,087	25,292
<i>Urinary bladder</i>	51.68	11,163	11,840	24,174
<i>Uterus caruncle</i>	48.33	12,199	12,917	24,857
<i>Vagina</i>	52.30	11,600	12,300	25,543

Average	49.60	11,328	12,009	23,837
----------------	--------------	---------------	---------------	---------------

604

605 **Comparative analysis of CAGE and WGBS to validate 'novel' TSS**

606 True TSS and TSS enhancer elements are very likely to be associated with areas of
607 hypomethylation (Yamashita et al., 2005, Yagi et al, 2008). The assessment of
608 hypomethylation of regions where “novel” TSS were identified thus provides a means
609 to support or refute their designation as true TSS. The methylation status of putative
610 TSS regions for eight of the tissues used for CAGE analysis was examined at single
611 nucleotide resolution using WGBS. Each WGBS library was pooled prior to
612 sequencing and multiplexed across eight lanes of the HiSeq X 10 platform. Following
613 trimming of the raw reads, the sequenced libraries produced an average of 103 Gbp
614 of clean data. The average mapping rate of the reads was 78.8%. A small proportion
615 (8.5%) of reads were identified as PCR or optical duplicates and were removed prior
616 to downstream analysis. The average read depth of the filtered libraries was 20x
617 coverage (Supplementary Table S3). Only cytosines with a minimum of ten reads
618 were retained for the subsequent comparative analysis with CAGE data to ensure a
619 high level of confidence in the methylation level estimates, as per published
620 recommendations (Doherty and Couldrey, 2014; Ziller et al., 2015). This work
621 represents one of the most comprehensive and high-quality methylation profiling
622 datasets in livestock to date. We would expect that reproducible TSS, either
623 annotated or novel, would overlap with hypo-methylated regions of the genome
624 (Yamashita et al., 2005; Yagi et al., 2008). Comparative analysis of the CAGE data
625 with the WGBS methylation levels from eight tissues from Benz2616 was used to
626 investigate methylation levels at the TSS in comparison to gene body and UTR
627 regions. For the majority of genes, the methylation level was much lower around the
628 transcriptionally active TSS or regulatory enhancer candidate regions compared to
629 the gene body (e.g. for gene *IRF2BP2* Figure 7). We overlaid the WGBS
630 hypomethylated regions and the CAGE uni-directional TSS clusters (annotated and
631 'novel') within 50bp of the promoter. For the eight matching tissues 88.7% of the
632 annotated TSS clusters and 32.2% of the 'novel' TSS were hypomethylated (Figure
633 8). The combined evidence of the hypomethylation and TSS support the conclusion
634 that 32.2% are in fact novel TSS clusters, whereas 67.8% of the novel TSS clusters
635 lack this confirmation.



636

637 *Figure 7. Overlay of CAGE, RNA-Seq and WGBS data tracks centred using the*
 638 *genomic coordinates of gene IRF2BP2. The green box shows a hypomethylated area*
 639 *overlapping multiple uni and bi-directional CAGE tag clusters. The black box*
 640 *represents predicted CAGE tag clusters with no verifying hypomethylation island,*
 641 *which are likely to be 'noise'.*

642

643



B

Tissue	Annotated + HypoCpG	Annotated w/o	Novel + HypoCpG	Novel w/o
Alveolar_Macrophages	90%	10%	33%	67%
Cerebellum	91%	9%	35%	65%
Cerebral_Cortex	89%	11%	35%	65%
Lung	88%	12%	28%	72%
Ovary	89%	11%	35%	65%
Rumen_atrium	87%	13%	32%	68%
Skeletal_Muscle_biceps_femoris	89%	11%	32%	68%
Skeletal_Muscle_longissimus_dorsi	86%	14%	27%	73%

644

645 *Figure 8: Numbers of CAGE TSS that were hypomethylated according to the WGBS*
 646 *data to distinguish between 'novel' reproducible (+HypoCpG) TSS and 'noise' (w/o).*
 647 *A) Shows the distribution of CAGE clusters as novel and annotated with or without*
 648 *HypoCpG. B) Percentage of CAGE clusters in each categories for each of the eight*
 649 *tissues.*

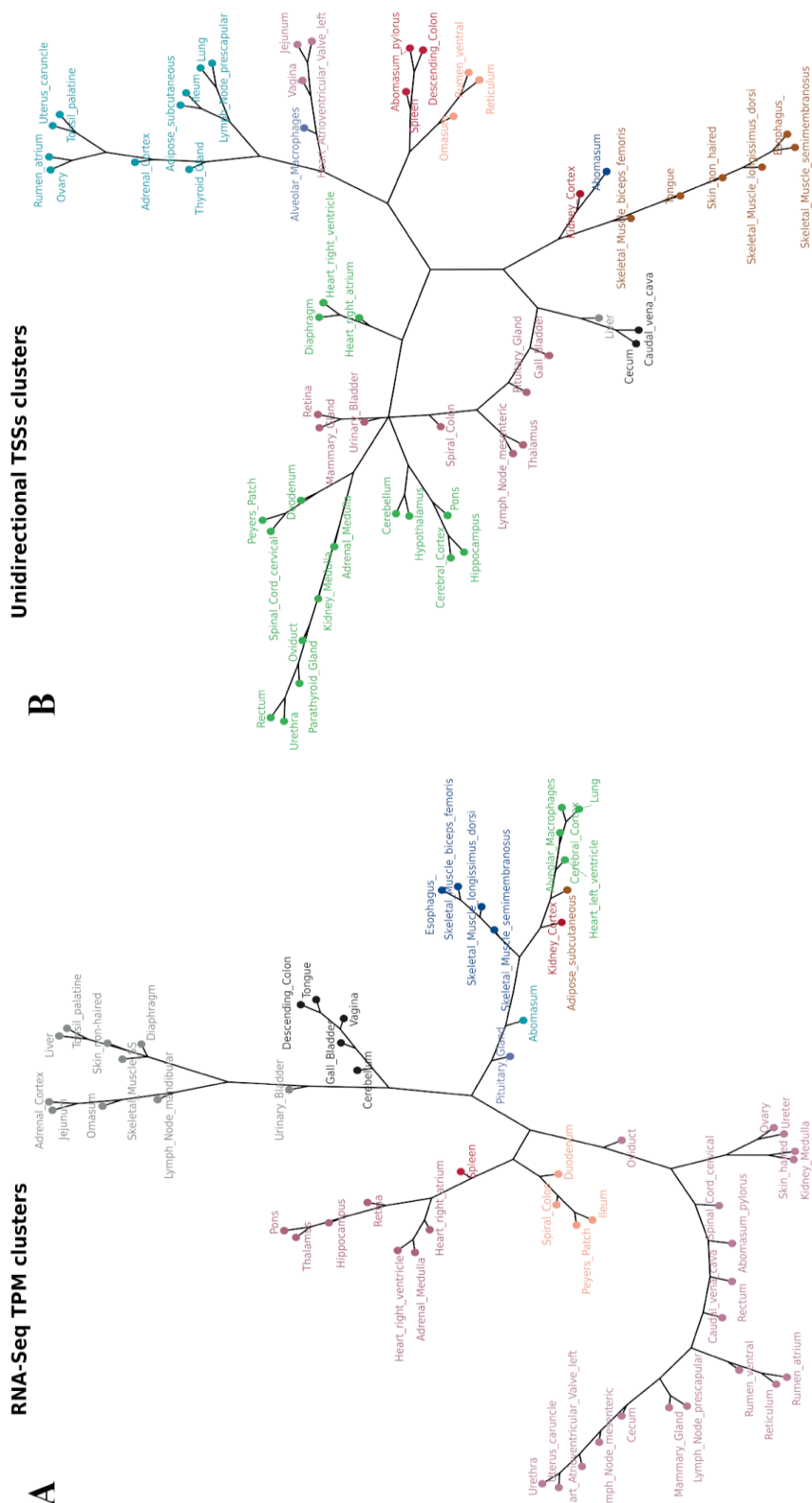
650

651 Validation of tissue expression profiles using mRNA-Seq

652 The tissue samples from Benz2616 were collected for the purpose of annotating her
 653 genome and as such N=1 in all cases. As an alternative strategy to having multiple
 654 biological replicates we validated the expression profiles for each tissue by
 655 comparing the CAGE data (CTPM) and mRNA-Seq (TPM) in 52 matching tissues.
 656 The transcript expression TPM was significantly correlated with the CAGE tag cluster
 657 CTPM values (correlation coefficient 0.19, Pearson p value $< 1.0e-08$) and visualised

658 as a heatmap (Supplementary Figure S4). A subset of house-keeping transcripts that
659 exhibit consistent expression for both the CAGE and mRNA-Seq datasets across all
660 tissues sampled, are visible from the heatmap (Supplementary Figure S4).

661 The similarity of tissue expression profiles for the uni-directional TSS clusters
662 was estimated in order to determine if tissues with similar physiology and function
663 formed distinct groups as expected. Similarity (distance) analysis showed a partial
664 grouping based on tissue type and organ system as shown in Figure 9A.
665 Physiologically similar tissues including nervous system and muscle tissues grouped
666 closely together. This grouping was also present in the mRNA-Seq data from tissue
667 matched samples Figure 9B, indicating good correlation between the two datasets.
668



669

670 *Figure 9. The network analysis of tissue TSS and gene expression profiles in 52*
 671 *matched samples from Benz1626. The clustering algorithm was based on MI*
 672 *distance of each tissue given the expressed A) mRNA-Seq transcript level TPM and*
 673 *B) CAGE tag clusters (TSSs).*

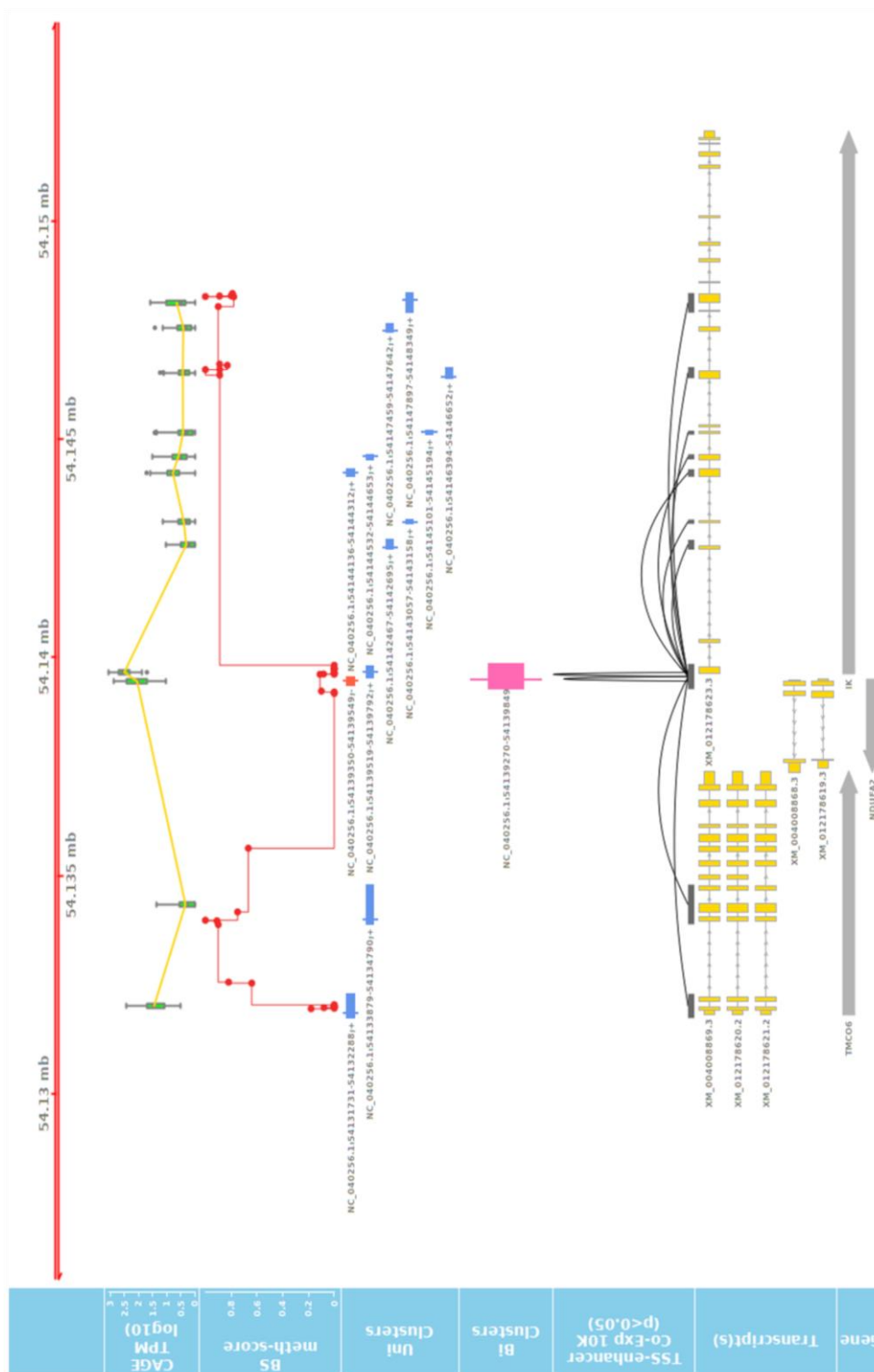
674 **Comparative visualisation of the datasets**

675 An interactive visualisation interface was developed in order to make these datasets
676 accessible and useable for the livestock genomics community. The genomic browser
677 incorporates the bp resolution hypomethylation data, the CTPM expression of TSS
678 and TSS-Enhancer regions and the mRNA-Seq TPM expression at transcript level.
679 These tracks are also overlaid using the coordinates provided by the TxDB objects
680 for transcripts and gene models as shown in Figure 10. This form of overlaid view
681 allows for confirmation of transcript expression and the exact coordinate of the
682 corresponding TSS in each tissue. For validation purposes the promoter region
683 should be under a hypo-methylated CpG island on the DNA track for any actively
684 transcribed gene in each tissue. The detailed bigBED format tracks for all the tissues
685 are available at https://data.faang.org/api/fire_api/trackhubregistry/hub.txt.

686 These visualisation tools were used to identify any co-expressed enhancers
687 within the proximity of a TSS. We were able to identify 741 TSS-Enhancer clusters
688 across the 56 tissues. An example of these bi-directional clusters is shown in Figure
689 10 as a pink line. The pairwise CTPM levels of co-expression of the bi-directional
690 clusters and those of the uni-directional TSS clusters were compared using the
691 Kendal correlation function in CAGEfightR (Thodberg and Sandelin, 2019). There
692 were 5383 significant co-expression pairs between uni-directional clusters (28,148)
693 and bi-directional clusters (741). An example of a co-expressed TSS-enhancer is
694 shown in Figure 10 as a black line connecting the significant start positions of the co-
695 expression pairs.

696 The co-expression range of bi-directional clusters, in some cases, can span
697 beyond the 10Kbp distance, as shown in the *IK* gene example Figure 10. The
698 expression of enhancer RNA (eRNA) with the promoter expression level of their
699 target genes has been reported before (Tippens et al., 2018). This layer of annotation
700 provides a foundation for enhancer target mapping in the sheep genome. The
701 detailed list and annotated target transcripts of these co-expression clusters can be
702 found in Supplementary File 2.

703



704

705 *Figure 10. Long range correlation of single enhancer site with multiple promoters of*
 706 *several genes. The track shows the significant correlation of a leading/primary*
 707 *enhancer site highly co-expressed with several TSS sites of different genes in a*
 708 *relatively long coding frame ($\pm 10,000$ Kb). The 3rd track from the top also shows the*
 709 *level of methylation at CpG sites at DNA level of Benz2616 overlaying the same*
 710 *coordinates of the IK gene and ± 10 Kbp.*

711

712 Discussion

713

714 High-quality reference genomes are now available for many farmed animal species
715 including domestic sheep (*Ovis aries*). The earlier draft genome sequence (Jiang et
716 al., 2014) has been superseded by a more contiguous genome assembly (*Oar*
717 *rambouillet v1.0* https://www.ncbi.nlm.nih.gov/assembly/GCF_002742125.1/).
718 Annotation of this genome sequence, however, is currently limited to gene and
719 transcript models. There is a lack of information on regulatory sequences and the
720 complexity of the transcriptome is underestimated. For example, promoters and TSS
721 are not well-annotated and alternative promoters and transcripts are poorly
722 characterised. The overall aim of the Ovine FAANG project was to provide a
723 comprehensive annotation of *Oar rambouillet v1.0*. To contribute to this aim we
724 generated a high-resolution global annotation of transcription start sites (TSS) for
725 sheep. After removal, of CAGE tags with < 10 read counts, 39.3% of TSS overlapped
726 with 5' ends of transcripts, as annotated previously by NCBI. A further 14.7%
727 mapped to within 50bp of annotated promoter regions. Intersecting these predicted
728 TSS regions with annotated promoter regions (± 50 bp) revealed 46% of the predicted
729 TSS were 'novel' and previously un-annotated. Using whole genome bisulphite
730 sequencing data from the same tissues we were able to determine that a proportion
731 of these 'novel' TSS were hypo-methylated (32.2%) indicating that they are likely to
732 be reproducible rather than 'noise'. This global annotation of TSS in sheep will
733 significantly enhance the annotation of gene models in the new ovine reference
734 assembly (*Oar rambouillet v1.0*).

735 The quality of the annotation of reference genomes for livestock species is
736 improving rapidly with reductions in the cost of sequencing and generation of new
737 datasets from multiple different functional assays (Giuffra and Tuggle, 2019). *Oar*
738 *rambouillet v1.0* superseded the Texel reference assembly (*Oar_v3.1*), which was
739 released in 2014 (Jiang et al., 2014). *Oar_v3.1* is still widely utilised by the sheep
740 genomics community and the Ensembl annotation
741 (https://www.ensembl.org/Ovis_aries/Info/Index) also includes sequence variation
742 information. We compared how mapped CAGE tag clusters were distributed across
743 genomic features in *Oar rambouillet v1.0* and *Oar_v3.1* (Jiang et al., 2014) and found
744 that the proportion of CAGE tag clusters mapping to promoter regions was greater for
745 *Oar rambouillet v1.0* (39%) than *Oar_v3.1* (15%). This may be because *Oar_v3.1*
746 was built using short read technology (Jiang et al., 2014), which had a significant bias
747 to GC rich regions, and therefore did not robustly capture the 5' ends of many genes
748 (Chen et al., 2013). In comparison, the *Oar rambouillet v1.0* assembly was generated

749 using long read technology, that dramatically improves the ease of assembly
750 resulting in increased contiguity (Contig N50: *Oar_v3.1* 0.07Mb and *Oar rambouillet*
751 *v1.0* 2.57Mb). Other recent high quality reference genome assemblies for livestock,
752 e.g. goat (Bickhart et al., 2017; Worley, 2017) and water buffalo (Low et al., 2019),
753 have been built using long read sequencing technology in combination with optical
754 mapping for scaffolding.

755 Highly annotated genomes are powerful tools that can help us to understand
756 the mechanisms underlying complex traits in livestock (Georges et al., 2018; Giuffra
757 and Tuggle, 2019) and mitigate future challenges to food production (Rexroad et al.,
758 2019). GWAS results, for example, can be integrated with functional annotation
759 information to identify causal variants enriched in trait-linked tissues or cell types
760 (reviewed in (Cano-Gamez and Trynka, 2020)). Using enrichment analysis (Finucane
761 et al., 2018) showed that heritable disease associated variants from GWAS were
762 enriched in enhancer regions in relevant tissues and cell types in humans. The TSS
763 and TSS-enhancer clusters identified in this study could be utilised in a similar way
764 for SNP enrichment analysis of GWAS variants in sheep. Using ChIP-Seq data
765 (Naval-Sanchez et al., 2018) found that selective sweeps were significantly enriched
766 for proximal regulatory elements to protein coding genes and genome features
767 associated with active transcription. A high quality set of variants for sheep,
768 generated using whole genome sequencing information for hundreds of animals
769 across multiple breeds, is available through (Sheep Genomes Database, 2020). This
770 dataset could be used to identify functional SNPs enriched in the TSS and TSS-
771 enhancer clusters for multiple tissues and cell types that we have annotated in the
772 *Oar rambouillet v1.0* assembly. High throughput functional screens using gene
773 editing technologies, are now possible to validate these functional variants (reviewed
774 in (Tait-Burkard et al., 2018)). New iPSC lines for livestock species also now offer the
775 potential to do this in relevant cell types (Ogorevc et al., 2016).

776 Our high-resolution atlas of TSS complements other available large-scale
777 RNA-Seq datasets for sheep e.g. (Clark et al., 2017). The analysis we present here
778 includes tissues representing all major organ systems. However, we were unable to
779 generate CAGE libraries for a small number of difficult to collect or problematic
780 tissues, and as such may have missed transcripts specific to these tissues. We were
781 also only able to generate CAGE libraries from one isolated cell type, alveolar
782 macrophages. As demonstrated by the FANTOM5 (Forrest et al., 2014) and
783 ENCODE (Birney et al., 2007) and FragENCODE (Foissac et al., 2019) projects,
784 including a diversity of immune cell types, in both activated and inactivated states, in
785 future work would capture additional transcriptional diversity. New technologies, such

786 as single cell sequencing, will allow annotation of cell-specific expressed and
787 regulatory regions of the genome at unprecedented resolution (Papatheodorou et al.,
788 2019). C1 CAGE now offers the opportunity to detect TSS and enhancer activity at
789 single-cell resolution (Kouno et al., 2019).

790 We have also generated full-length transcript information using the Iso-Seq
791 method, for a small subset of tissues from Benz2616. Integrating mRNA-Seq and
792 Iso-Seq datasets has been used successfully to improve the annotation of the pig
793 genome (Beiki et al., 2019). By merging the Iso-Seq data with the CAGE and mRNA-
794 Seq datasets we will be able to measure differential transcript usage across tissues
795 and improve the resolution of the *Oar rambouillet v1.0* transcriptome further. As such
796 the study we present here represents just the first step in demonstrating the power
797 and utility of the different datasets generated for the Ovine FAANG project, which will
798 provide one of the highest resolution annotations of transcript regulation and diversity
799 in a livestock species to date.

800

801 **Acknowledgements**

802 The authors would like to thank the Human Genome Sequencing Center staff for the
803 *Oar rambouillet v1.0* reference genome assembly and for long-read and short read
804 mRNA and miRNA sequencing. HGSC contributors include the genome assembly
805 and analysis team of Y. Liu, R.A. Harris, X. Qin, led by K. Worley and production
806 team members including M-C Gingras and L. Perez for RNA and DNA preparation,
807 V. Vee, Y. Han. V. Korchina, S. Dugan-Perez for sequencing, Q. Meng, H.
808 Doddapaneni, M. Wang for library production, the system support and LIMs teams,
809 D. M. Muzny the HGSC Director of Operations, R. A. Gibbs the HGSC Director. We
810 thank S. Sullivan and I. Liachko of Phase Genomics for PGA genome scaffolding.

811 We would also like to thank William Thompson for isolation of RNA and Jacky
812 Carnahan for isolation of DNA at USMARC. The authors are grateful to Lucas
813 Lefevre, Rachel Young and Heather Finlayson for performing the initial optimisation
814 to establish the CAGE protocol at the Roslin Institute and Sara Clohisey and Lee
815 Murphy for advice on data analysis and assistance in establishing the protocol at the
816 Edinburgh Clinical Research Facility. The authors are also grateful for the support of
817 the FAANG Data Coordination Centre (<http://data.fang.org>) in the upload and
818 archiving of the sample data and metadata, and hosting of the genome tracks.

819 Furthermore, we would like to thank the tissue collection design team; Noelle
820 Cockett and Kim Worley, who coordinated the team, James Kijas, Brian Dalrymple,
821 Tracy Hadfield, Kara Thornton, Tom Baldwin, Shuna Jones, Bob Lee, Sue Hauver,
822 Christy Kelley, Jessica Eisenhauer, Mike Heaton, William Thompson, Timothy Smith,

823 Stephen White, Michelle Mousel, Alisha Massa; Brian Sayre and Brenda Murdoch
824 who all contributed in planning and designing the collection. Additionally, we would
825 like to thank all those that contributed to the FAANG tissue collection at USU: Noelle
826 Cockett, Tracy Hadfield, Tom Baldwin, Rusty Stott, Arnaud Van Wettere, Holly
827 Mason, Jaqueline LaRose, Dave Forrester, Corey Wareham, Sarah Behunin, Kara
828 Thornton, Gordon Hullinger (VMES), Alisha Massa, Maria Herndon, Brenda
829 Murdoch, Brian Sayre, Caylee Birge, Codie Durfee, Michelle R. Mousel, Rachael
830 Christianson, Nicole Ineck, Angie Robinson, Dallin Wengert, Kerry Rood, Erica
831 Moscoso, Rickie Warr, Dustin Kinney, Abbey Benninghoff, Sumira Phatak, Kevin
832 Contreras, Braden Abercrombie, Misha Regouski.

833

834 **Ethics Statement**

835 All protocols were approved by the animal care and use in accordance with Utah
836 State University. IACUC approval: #2826, expiration date 21st of February 2021.

837

838 **Data Availability**

839 All the raw sequence data and analysis BAM files for this study are publicly available
840 via the OAR_USU_Benz2616 NCBI BioProject:
841 <https://www.ncbi.nlm.nih.gov/bioproject/PRJNA414087> and via the European
842 Nucleotide Archive (ENA): <https://www.ebi.ac.uk/ena/browser/view/PRJEB34864>
843 (CAGE), <https://www.ebi.ac.uk/ena/data/view/PRJEB35292> (mRNA-Seq) and
844 <http://www.ebi.ac.uk/ena/data/view/PRJEB39178> (WGBS). Details of all 100 samples
845 collected from Benz 2616 are included in the BioSamples database under
846 submission GSB-7268, group accession number SAMEG329607
847 (<https://www.ebi.ac.uk/biosamples/samples/SAMEG329607>). The datasets are
848 accessible via the FAANG data portal and were submitted according to FAANG
849 sample and experimental metadata requirements (Harrison et al., 2018). Oar
850 rambouillet v1.0 is now available on the Ensembl Rapid Release site
851 https://rapid.ensembl.org/Ovis_aries_rambouillet/Info/Index.

852

853 **Author Contributions**

854 RC and IG performed CAGE library, optimisation, preparation and sequencing. MS
855 performed all bioinformatic and data analyses, with the exception of the WGBS data,
856 which was analysed by AC. MS and AC generated the GViz tracks. TPLS
857 coordinated generation of the mRNA-Seq data at US-MARC. KCW coordinated
858 generation of the Oar Rambouillet v1.0 reference assembly and mRNA-Seq data at
859 BCM. SMC coordinated the generation and analysis of the WGBS with AC. ELC and

860 ALA coordinated the CAGE components of the study. NEC and KCW planned and
861 coordinated the sample collection at USU. BMM is coordinator of the Ovine FAANG
862 project with ALA, SNW, BPD, JWK, RB, NEC, SMC, KCW, ELC and TPLS who
863 designed the overall project and acquired the funding to support the work. ELC wrote
864 the manuscript with MS and AC. All authors contributed to editing and approved the
865 final version of the manuscript.

866

867 **Conflict of Interest**

868 The authors declare that the research was conducted in the absence of any
869 commercial or financial relationships that could be construed as a potential conflict of
870 interest.

871

872 **Code Availability**

873 All the code base for the analytical pipeline in this study are available at
874 https://msalavat@bitbucket.org/msalavat/rnaseqwrap_public.git for RNA-Seq
875 analysis, https://msalavat@bitbucket.org/msalavat/cagewrap_public.git for the CAGE
876 mapping, annotation and metrics pipeline and
877 https://msalavat@bitbucket.org/caultona/wgbswrap_public.git for WGBS pipeline.

878

879 **Funding**

880 This work was supported by National Institute of Food and Agriculture, U.S.
881 Department of Agriculture awards USDA-NIFA-2017-67016-26301 and USDA-NIFA-
882 2013-67015-21228. Sample collection was funded by NRSP-8 Sheep Genome
883 Coordinator Funding, Project UTA-1172. ELC, ALA and MS were partially supported
884 by Institute Strategic Program grants awarded to the Roslin Institute 'Farm Animal
885 Genomics' (BBS/E/D/2021550) and 'Prediction of genes and regulatory elements in
886 farm animal genomes' (BBS/E/D/10002070). ELC is supported by a University of
887 Edinburgh Chancellors Fellowship. The Edinburgh Clinical Research Facility is
888 funded by the Wellcome Trust. AC was supported by the University of Otago PhD
889 scholarship and granted the Marjorie McCallum travel award to visit the Roslin
890 Institute, her PhD research is funded by NZ MBIE C10X1906. TPLS was supported
891 by USDA-ARS Project No. 3040-31000-100-00D. The funders had no role in study
892 design, data collection and analysis, decision to publish, or preparation of the
893 manuscript.

894

895 **Members of the Ovine FAANG Project Consortium (listed by institution)**

896 Brenda Murdoch (University of Idaho)
897 Kimberly M Davenport (University of Idaho)
898 Stephen White (USDA, ARS, Washington State University)
899 Michelle Mousel (USDA, ARS ADRC)
900 Alisha Massa (Washington State University)
901 Kim Worley (Baylor College of Medicine)
902 Alan Archibald (The Roslin Institute, University of Edinburgh)
903 Emily Clark (The Roslin Institute, University of Edinburgh)
904 Brian Dalrymple (University of Western Australia)
905 James Kijas (CSIRO)
906 Shannon Clarke (AgResearch)
907 Rudiger Brauning (AgResearch)
908 Timothy Smith (USDA, ARS MARC)
909 Tracey Hadfield (Utah State University)
910 Noelle Cockett (Utah State University)

911

912 **References**

913 An, X., Ma, H., Han, P., Zhu, C., Cao, B., and Bai, Y. (2018). Genome-wide
914 differences in DNA methylation changes in caprine ovaries between oestrous
915 and dioestrous phases. *J. Anim. Sci. Biotechnol.* doi:10.1186/s40104-018-0301-
916 x.
917 Andersson, L., Archibald, A. L., Bottema, C. D., Brauning, R., Burgess, S. C., Burt, D.
918 W., et al. (2015). Coordinated international action to accelerate genome-to-
919 phenome with FAANG, the Functional Annotation of Animal Genomes project.
920 *Genome Biol.* 16, 57. doi:10.1186/s13059-015-0622-4.
921 Andersson, R., Gebhard, C., Miguel-Escalada, I., Hoof, I., Bornholdt, J., Boyd, M., et
922 al. (2014). An atlas of active enhancers across human cell types and tissues.
923 *Nature* 507, 455–461. doi:10.1038/nature12787.
924 Baillie, J. K., Arner, E., Daub, C., De Hoon, M., Itoh, M., Kawaji, H., et al. (2017).
925 Analysis of the human monocyte-derived macrophage transcriptome and
926 response to lipopolysaccharide provides new insights into genetic aetiology of
927 inflammatory bowel disease. *PLOS Genet.* 13, e1006641.
928 Beiki, H., Liu, H., Huang, J., Manchanda, N., Nonneman, D., Smith, T. P. L., et al.
929 (2019). Improved annotation of the domestic pig genome through integration of
930 Iso-Seq and RNA-seq data. *BMC Genomics* 20, 344. doi:10.1186/s12864-019-
931 5709-y.
932 Berger, M. R., Alvarado, R., and Kiss, D. L. (2019). mRNA 5' ends targeted by

- 933 cytoplasmic recapping cluster at CAGE tags and select transcripts are
934 alternatively spliced. *FEBS Lett.* 593, 670–679. doi:10.1002/1873-3468.13349.
- 935 Bertin, N., Mendez, M., Hasegawa, A., Lizio, M., Abugessaisa, I., Severin, J., et al.
936 (2017). Linking FANTOM5 CAGE peaks to annotations with CAGEscan. *Sci.*
937 *Data* 4, 170147. doi:10.1038/sdata.2017.147.
- 938 Bickhart, D. M., Rosen, B. D., Koren, S., Sayre, B. L., Hastie, A. R., Chan, S., et al.
939 (2017). Single-molecule sequencing and chromatin conformation capture enable
940 de novo reference assembly of the domestic goat genome. *Nat. Genet.* 49,
941 643–650. doi:10.1038/ng.3802.
- 942 Birney, E., Stamatoyannopoulos, J. A., Dutta, A., Guigó, R., Gingeras, T. R.,
943 Margulies, E. H., et al. (2007). Identification and analysis of functional elements
944 in 1% of the human genome by the ENCODE pilot project. *Nature* 447, 799–
945 816. doi:10.1038/nature05874.
- 946 Bray, N. L., Pimentel, H., Melsted, P., and Pachter, L. (2016). Near-optimal
947 probabilistic RNA-seq quantification. *Nat. Biotechnol.* 34, 525–527.
948 doi:10.1038/nbt.3519.
- 949 Cano-Gamez, E., and Trynka, G. (2020). From GWAS to Function: Using Functional
950 Genomics to Identify the Mechanisms Underlying Complex Diseases. *Front.*
951 *Genet.* 11, 424.
- 952 Chen, Y.-C., Liu, T., Yu, C.-H., Chiang, T.-Y., and Hwang, C.-C. (2013). Effects of
953 GC Bias in Next-Generation-Sequencing Data on De Novo Genome Assembly.
954 *PLoS One* 8, e62856.
- 955 Clark, E. L., Bush, S. J., McCulloch, M. E. B., Farquhar, I. L., Young, R., Lefevre, L.,
956 et al. (2017). A high resolution atlas of gene expression in the domestic sheep
957 (*Ovis aries*). *PLOS Genet.* 13, e1006997.
- 958 Cordier, G., Cozon, G., Greenland, T., Rocher, F., Guiguen, F., Guerret, S., et al.
959 (1990). In vivo activation of alveolar macrophages in ovine lentivirus infection.
960 *Clin. Immunol. Immunopathol.* 55, 355–367. doi:[https://doi.org/10.1016/0090-](https://doi.org/10.1016/0090-1229(90)90124-9)
961 [1229\(90\)90124-9](https://doi.org/10.1016/0090-1229(90)90124-9).
- 962 Ding, B., Gentleman, R., and Carey, V. (2012). bioDist: Different distance measures.
963 R package version 1.28.0. Available at:
964 <https://www.bioconductor.org/packages/release/bioc/html/bioDist.html>.
- 965 Doherty, R., and Couldrey, C. (2014). Exploring genome wide bisulfite sequencing for
966 DNA methylation analysis in livestock: A technical assessment. *Front. Genet.*
967 doi:10.3389/fgene.2014.00126.
- 968 Edinburgh, U. of (2020). Edinburgh Compute and Data Facility. Available at:
969 <https://www.ed.ac.uk/is/research-computing-service> [Accessed July 6, 2020].

- 970 Finucane, H. K., Reshef, Y. A., Anttila, V., Slowikowski, K., Gusev, A., Byrnes, A., et
971 al. (2018). Heritability enrichment of specifically expressed genes identifies
972 disease-relevant tissues and cell types. *Nat. Genet.* 50, 621–629.
973 doi:10.1038/s41588-018-0081-4.
- 974 Foissac, S., Djebali, S., Munyard, K., Vialaneix, N., Rau, A., Muret, K., et al. (2019).
975 Multi-species annotation of transcriptome and chromatin structure in
976 domesticated animals. *BMC Biol.* 17, 108. doi:10.1186/s12915-019-0726-5.
- 977 Forrest, A. R. R., Kawaji, H., Rehli, M., Baillie, J. K., De Hoon, M. J. L., Haberle, V.,
978 et al. (2014). A promoter-level mammalian expression atlas. *Nature* 507, 462–
979 470. doi:10.1038/nature13182.
- 980 Georges, M., Charlier, C., and Hayes, B. (2018). Harnessing genomic information for
981 livestock improvement. *Nat. Rev. Genet.* 20, 1. doi:10.1038/s41576-018-0082-2.
- 982 Giuffra, E., and Tuggle, C. K. (2019). Functional Annotation of Animal Genomes
983 (FAANG): Current Achievements and Roadmap. *Annu. Rev. Anim. Biosci.* 7,
984 65–88. doi:10.1146/annurev-animal-020518-114913.
- 985 Guo, W., Zhu, P., Pellegrini, M., Zhang, M. Q., Wang, X., and Ni, Z. (2018).
986 CGmapTools improves the precision of heterozygous SNV calls and supports
987 allele-specific methylation detection and visualization in bisulfite-sequencing
988 data. *Bioinformatics* 34, 381–387. doi:10.1093/bioinformatics/btx595.
- 989 Hahne, F., and Ivanek, R. (2016). “Visualizing Genomic Data Using Gviz and
990 Bioconductor BT - Statistical Genomics: Methods and Protocols,” in, eds. E.
991 Mathé and S. Davis (New York, NY: Springer New York), 335–351.
992 doi:10.1007/978-1-4939-3578-9_16.
- 993 Hannon Lab (2017). FASTX-Toolkit FASTQ/A short reads pre-processing tools.
994 Available at: http://hannonlab.cshl.edu/fastx_toolkit/ [Accessed July 6, 2020].
- 995 Harrison, P. W., Fan, J., Richardson, D., Clarke, L., Zerbino, D., Cochrane, G., et al.
996 (2018). FAANG, establishing metadata standards, validation and best practices
997 for the farmed and companion animal community. *Anim. Genet.* 49, 520–526.
998 doi:10.1111/age.12736.
- 999 Ibeagha-Awemu, E. M., and Zhao, X. (2015). Epigenetic marks: Regulators of
1000 livestock phenotypes and conceivable sources of missing variation in livestock
1001 improvement programs. *Front. Genet.* doi:10.3389/fgene.2015.00302.
- 1002 Jiang, Y., Xie, M., Chen, W., Talbot, R., Maddox, J. F., Faraut, T., et al. (2014). The
1003 sheep genome illuminates biology of the rumen and lipid metabolism. *Science*
1004 344, 1168–1173. doi:10.1126/science.1252806.
- 1005 Joe, H. (1989). Relative Entropy Measures of Multivariate Dependence. *J. Am. Stat.*
1006 *Assoc.* 84, 157–164. doi:10.1080/01621459.1989.10478751.

- 1007 Jones, P. A. (2012). Functions of DNA methylation: Islands, start sites, gene bodies
1008 and beyond. *Nat. Rev. Genet.* doi:10.1038/nrg3230.
- 1009 Kent, W. J., Zweig, A. S., Barber, G., Hinrichs, A. S., and Karolchik, D. (2010).
1010 BigWig and BigBed: enabling browsing of large distributed datasets.
1011 *Bioinformatics* 26, 2204–2207. doi:10.1093/bioinformatics/btq351.
- 1012 Kouno, T., Moody, J., Kwon, A. T.-J., Shibayama, Y., Kato, S., Huang, Y., et al.
1013 (2019). C1 CAGE detects transcription start sites and enhancer activity at
1014 single-cell resolution. *Nat. Commun.* 10, 360. doi:10.1038/s41467-018-08126-5.
- 1015 Langmead, B., and Salzberg, S. L. (2012). Fast gapped-read alignment with Bowtie
1016 2. *Nat. Methods* 9, 357–359. doi:10.1038/nmeth.1923.
- 1017 Lassmann, T. (2015). TagDust2: a generic method to extract reads from sequencing
1018 data. *BMC Bioinformatics* 16, 24. doi:10.1186/s12859-015-0454-y.
- 1019 Lawrence, M., Huber, W., Pagès, H., Aboyoun, P., Carlson, M., Gentleman, R., et al.
1020 (2013). Software for computing and annotating genomic ranges. *PLoS Comput.*
1021 *Biol.* 9, e1003118–e1003118. doi:10.1371/journal.pcbi.1003118.
- 1022 Lev Maor, G., Yearim, A., and Ast, G. (2015). The alternative role of DNA methylation
1023 in splicing regulation. *Trends Genet.* doi:10.1016/j.tig.2015.03.002.
- 1024 Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., et al. (2009).
1025 The Sequence Alignment/Map format and SAMtools. *Bioinformatics* 25, 2078–
1026 2079. doi:10.1093/bioinformatics/btp352.
- 1027 Lizio, M., Deviatiiarov, R., Nagai, H., Galan, L., Arner, E., Itoh, M., et al. (2017).
1028 Systematic analysis of transcription start sites in avian development. *PLOS Biol.*
1029 15, e2002887.
- 1030 Low, W. Y., Tearle, R., Bickhart, D. M., Rosen, B. D., Kingan, S. B., Swale, T., et al.
1031 (2019). Chromosome-level assembly of the water buffalo genome surpasses
1032 human and goat genomes in sequence contiguity. *Nat. Commun.* 10, 260.
1033 doi:10.1038/s41467-018-08260-0.
- 1034 Murdoch, B. M. (2019). The functional annotation of the sheep genome project. *J.*
1035 *Anim. Sci.* 97, 16.
- 1036 Naval-Sanchez, M., Nguyen, Q., McWilliam, S., Porto-Neto, L. R., Tellam, R.,
1037 Vuocolo, T., et al. (2018). Sheep genome functional annotation reveals proximal
1038 regulatory elements contributed to the evolution of modern breeds. *Nat.*
1039 *Commun.* 9, 859. doi:10.1038/s41467-017-02809-1.
- 1040 Ogorevc, J., Orehek, S., and Dovč, P. (2016). Cellular reprogramming in farm
1041 animals: an overview of iPSC generation in the mammalian farm animal
1042 species. *J. Anim. Sci. Biotechnol.* 7, 10. doi:10.1186/s40104-016-0070-3.
- 1043 Pages, H. (2020). BSgenome: Software infrastructure for efficient representation of

- 1044 full genomes and their SNPs. R package version 1.56.0.
- 1045 Papatheodorou, I., Moreno, P., Manning, J., Fuentes, A. M.-P., George, N., Fexova,
1046 S., et al. (2019). Expression Atlas update: from tissues to single cells. *Nucleic
1047 Acids Res.* 48, D77–D83. doi:10.1093/nar/gkz947.
- 1048 Priness, I., Maimon, O., and Ben-Gal, I. (2007). Evaluation of gene-expression
1049 clustering via mutual information distance measure. *BMC Bioinformatics* 8, 111.
1050 doi:10.1186/1471-2105-8-111.
- 1051 Quinlan, A. R., and Hall, I. M. (2010). BEDTools: A flexible suite of utilities for
1052 comparing genomic features. *Bioinformatics* 26, 841–842.
1053 doi:10.1093/bioinformatics/btq033.
- 1054 Reshef, D. N., Reshef, Y. A., Sabeti, P. C., and Mitzenmacher, M. (2018). An
1055 empirical study of the maximal and total information coefficients and leading
1056 measures of dependence. *Ann. Appl. Stat.* 12, 123–155. doi:10.1214/17-
1057 AOAS1093.
- 1058 Rexroad, C., Vallet, J., Matukumalli, L. K., Reecy, J., Bickhart, D., Blackburn, H., et
1059 al. (2019). Genome to Phenome: Improving Animal Health, Production, and
1060 Well-Being – A New USDA Blueprint for Animal Genome Research 2018–2027
1061 . *Front. Genet.* 10, 327.
- 1062 Sheep Genomes Database (2020). Available at: <https://sheepgenomesdb.org/>
1063 [Accessed July 6, 2020].
- 1064 Song, Q., Decato, B., Hong, E. E., Zhou, M., Fang, F., Qu, J., et al. (2013). A
1065 reference methylome database and analysis pipeline to facilitate integrative and
1066 comparative epigenomics. *PLoS One* 8, e81148–e81148.
1067 doi:10.1371/journal.pone.0081148.
- 1068 Tait-Burkard, C., Doeschl-Wilson, A., McGrew, M. J., Archibald, A. L., Sang, H. M.,
1069 Houston, R. D., et al. (2018). Livestock 2.0 – genome editing for fitter, healthier,
1070 and more productive farmed animals. *Genome Biol.* 19, 204.
1071 doi:10.1186/s13059-018-1583-1.
- 1072 Takahashi, H., Kato, S., Murata, M., and Carninci, P. (2012). “CAGE (Cap Analysis of
1073 Gene Expression): A Protocol for the Detection of Promoter and Transcriptional
1074 Networks,” in *Methods in Molecular Biology*, 181–200. doi:10.1007/978-1-
1075 61779-292-2_11.
- 1076 Thodberg, M., and Sandelin, A. (2019). A step-by-step guide to analyzing CAGE data
1077 using R/Bioconductor. *F1000Research* 8, 886.
1078 doi:10.12688/f1000research.18456.1.
- 1079 Thodberg, M., Thieffry, A., Vitting-Seerup, K., Andersson, R., and Sandelin, A.
1080 (2019). CAGEfightR: analysis of 5'-end data using R/Bioconductor. *BMC*

1081 *Bioinformatics* 20, 487. doi:10.1186/s12859-019-3029-5.
1082 Tippens, N. D., Vihervaara, A., and Lis, J. T. (2018). Enhancer transcription: what,
1083 where, when, and why? *Genes Dev.* 32, 1–3. doi:10.1101/GAD.311605.118.
1084 Worley, K. C. (2017). A golden goat genome. *Nat. Genet.* 49, 485–486.
1085 doi:10.1038/ng.3824.
1086 Yagi, S., Hirabayashi, K., Sato, S., Li, W., Takahashi, Y., Hirakawa, T., et al. (2008).
1087 DNA methylation profile of tissue-dependent and differentially methylated
1088 regions (T-DMRs) in mouse promoter regions demonstrating tissue-specific
1089 gene expression. *Genome Res.* 18, 1969–1978. doi:10.1101/gr.074070.107.
1090 Yamashita, R., Suzuki, Y., Sugano, S., and Nakai, K. (2005). Genome-wide analysis
1091 reveals strong correlation between CpG islands with nearby transcription start
1092 sites of genes and their tissue specificity. *Gene* 350, 129–136.
1093 doi:<https://doi.org/10.1016/j.gene.2005.01.012>.
1094 Ziller, M. J., Hansen, K. D., Meissner, A., and Aryee, M. J. (2015). Coverage
1095 recommendations for methylation analysis by whole-genome bisulfite
1096 sequencing. *Nat. Methods.* doi:10.1038/nmeth.3152.

1097

1098 **Figure Legends**

1099

1100 Figure 1. FAANG assays (CAGE, WGBS and mRNA-Seq) performed on each tissue
1101 from Benz2616.

1102

1103 Figure 2. Workflow of the analysis pipeline and respective tools used for CAGE
1104 sequence data analysis

1105

1106 Figure 3. Schematic representation of the two clustering algorithms used in the
1107 CAGEfightR package for TSS (uni-directional) and TSS-Enhancer (bi-directional)
1108 clustering.

1109

1110 Figure 4. The genomic region distribution of CAGE tag clusters mapped against *Oar*
1111 *rambouillet v1.0* assembly and gene annotation. The counts were averaged across
1112 tissues. A) Uni-directional TSS clusters with the highest proportion in promoter
1113 region (± 100 bp of the 5'UTR beginning at the [TSS]). B) Bi-directional TSS-enhancer
1114 clusters with the highest proportion in the proximal region (1000bp upstream of the
1115 5'UTR beginning at the [TSS]).

1116

1117 Figure 5. Chord diagram of expression level (TPM) of CAGE tag clusters (uni-
1118 directional TSS) across all the tissues collected from Benz2616. Shared CAGE tag
1119 clusters are common to at least 2/3rd of the tissues (37/56).

1120

1121 Figure 6. Chord diagram of expression level (TPM) of CAGE tag clusters (bi-
1122 directional TSS-Enhancer) across all the tissues collected from Benz2616. CAGE tag
1123 clusters expressed (>10 CTPM) by at least 2/3rd of the tissues (37/56).

1124

1125 Figure 7. Overlay of CAGE, RNA-Seq and WGBS data tracks centred using the
1126 genomic coordinates of gene *IRF2BP2*. The green box shows a hypomethylated area
1127 overlapping multiple uni and bi-directional CAGE tag clusters. The black box
1128 represents predicted CAGE tag clusters with no verifying hypomethylation island,
1129 which are likely to be 'noise'.

1130

1131 Figure 8: Numbers of CAGE TSS that were hypomethylated according to the WGBS
1132 data to distinguish between 'novel' reproducible (+HypoCpG) TSS and 'noise' (w/o).
1133 A) Shows the distribution of CAGE clusters as novel and annotated with or without
1134 HypoCpG. B) Percentage of CAGE clusters in each categories for each of the eight
1135 tissues.

1136

1137 Figure 9. The network analysis of tissue TSS and gene expression profiles in 52
1138 matched samples from Benz1626. The clustering algorithm was based on MI
1139 distance of each tissue given the expressed A) mRNA-Seq transcript level TPM and
1140 B) CAGE tag clusters (TSSs).

1141

1142 Figure 10. Long range correlation of single enhancer site with multiple promoters of
1143 several genes. The track shows the significant correlation of a leading/primary
1144 enhancer site highly co-expressed with several TSS sites of different genes in a
1145 relatively long coding frame ($\pm 10,000$ Kb). The 3rd track from the top also shows the
1146 level of methylation at CpG sites at DNA level of Benz2616 overlaying the same
1147 coordinates of the IK gene and ± 10 Kbp.

1148

1149

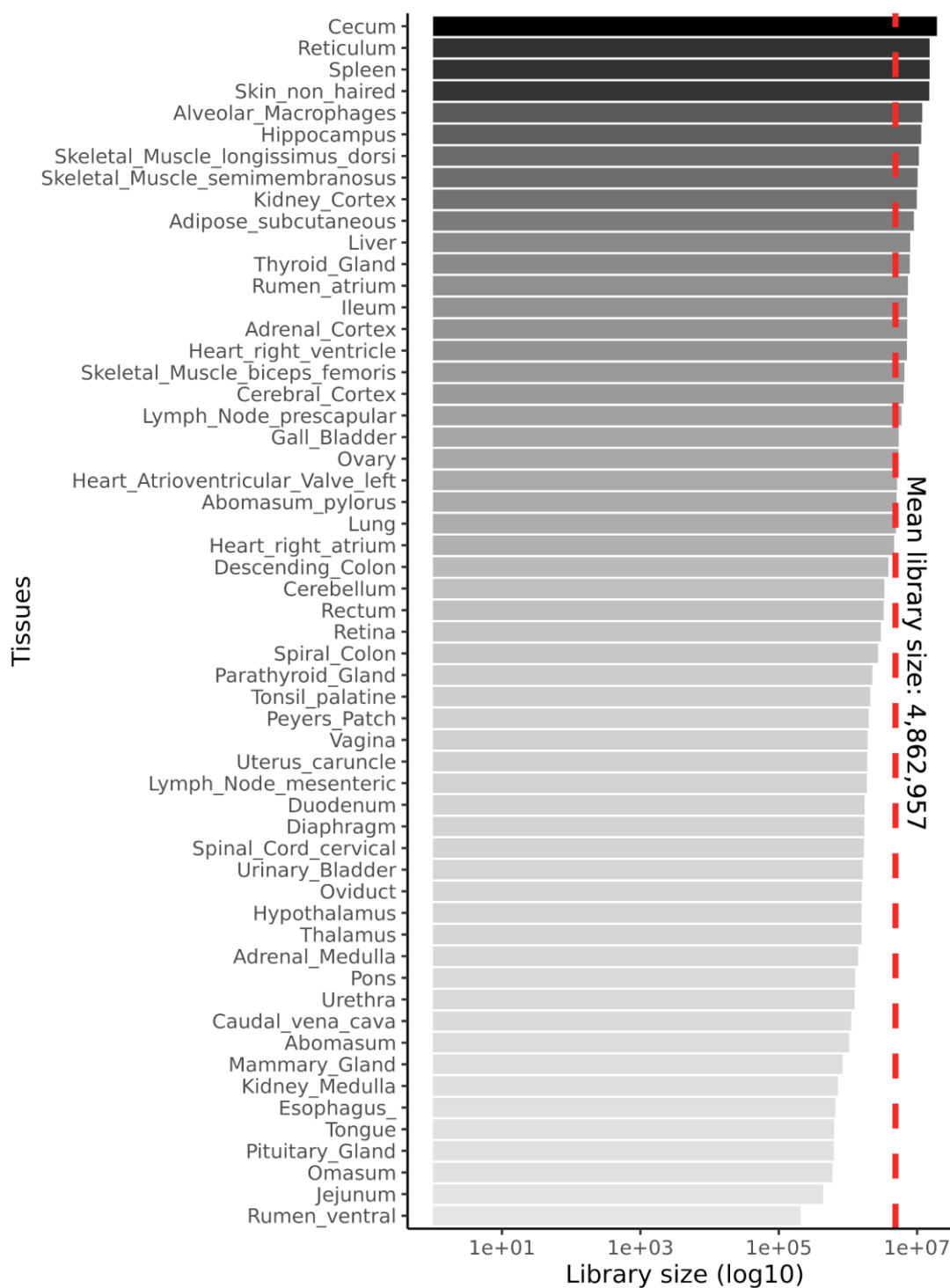
1150

1151

1152

1153

1154 **Supplemental Material**



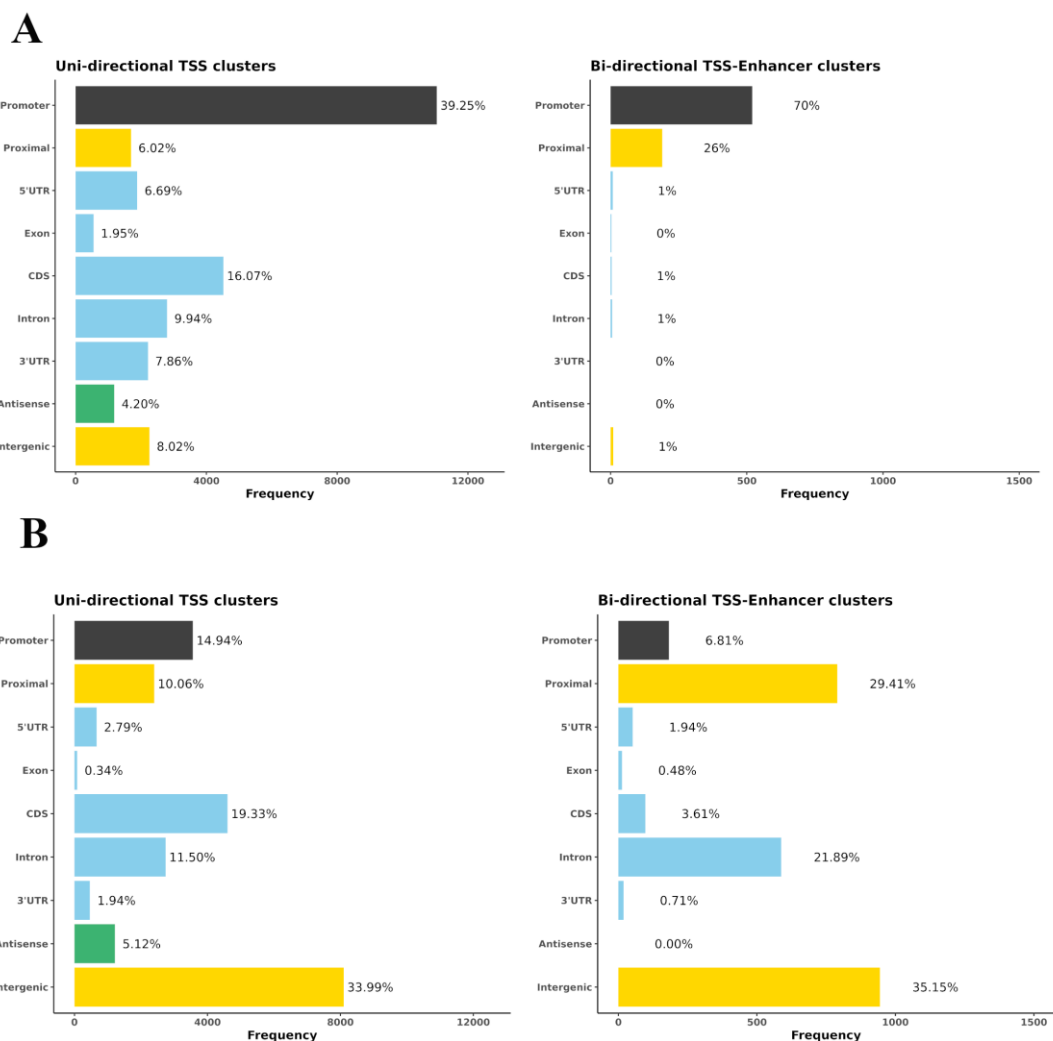
1155

1156 Supplementary Figure S1. CAGE library size for each of the 56 tissues analysed.

1157

1158

1159

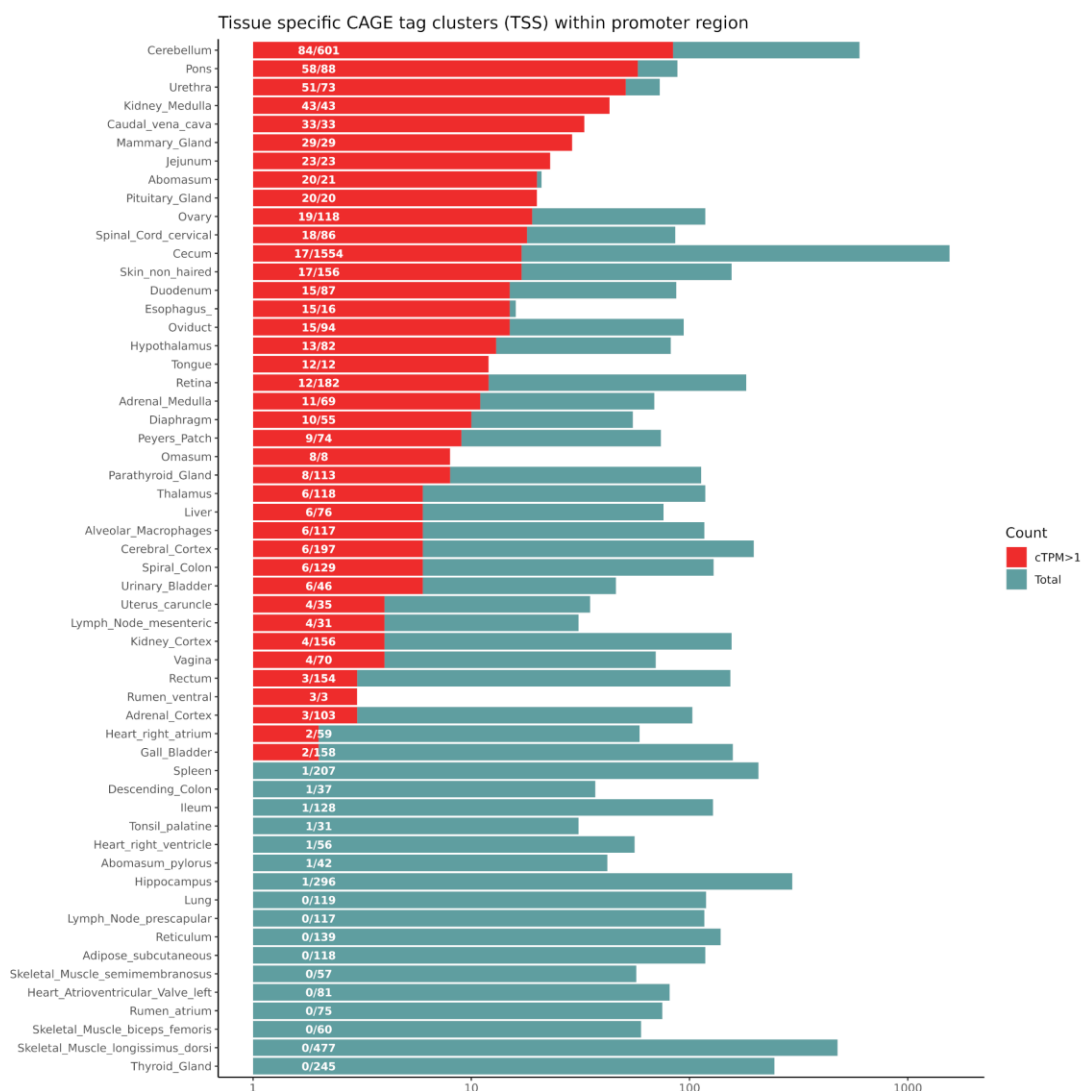


1160

1161 Supplementary Figure S2. The percentage of CAGE tags mapped to each genomic
 1162 region for Oar_rambouillet_v1.0 (A) and Oar_v3.1 (B) reference genome assemblies.

1163 The counts were average across tissues prior to annotation.

1164

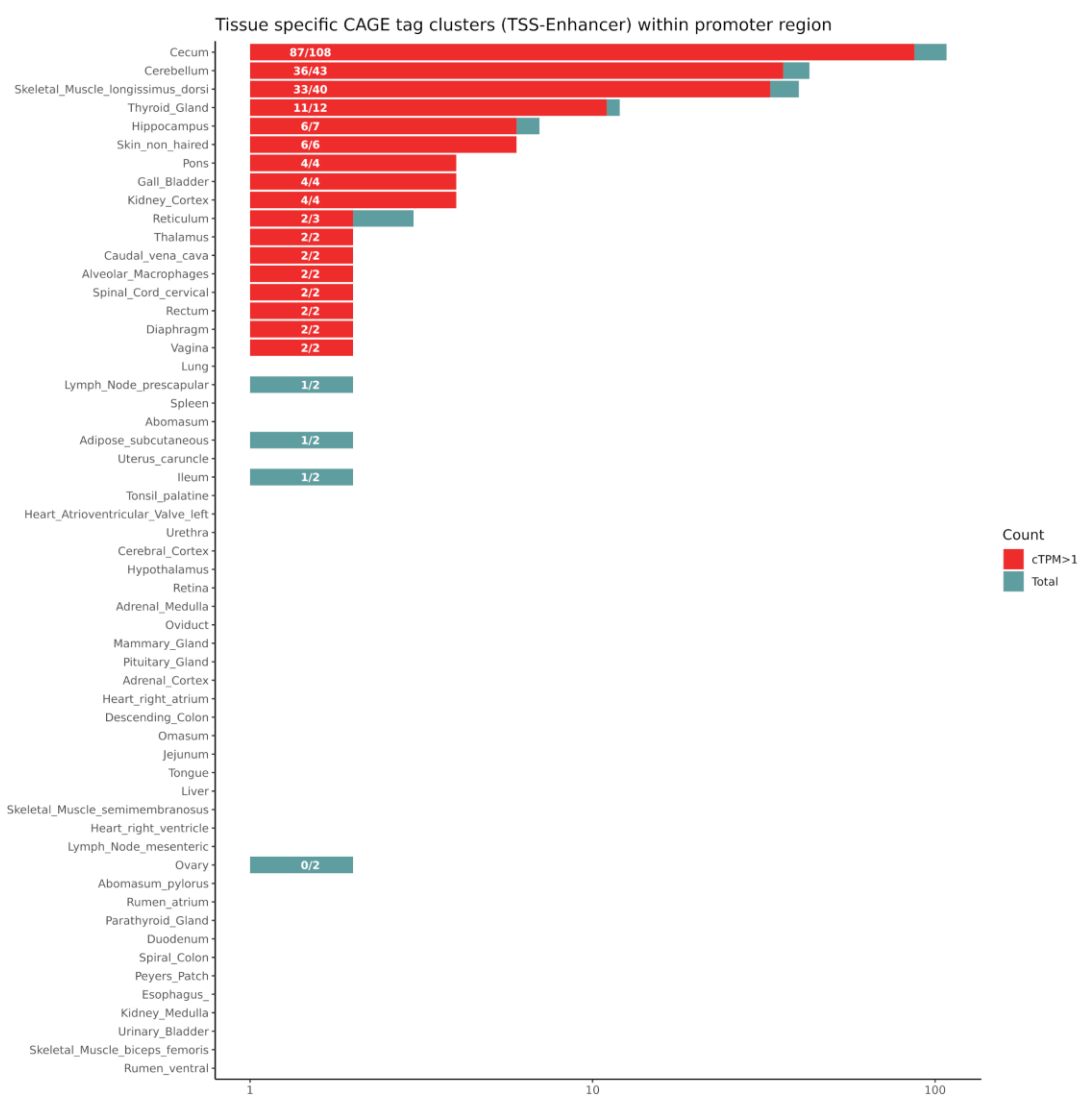


1165

1166

1167 Supplementary Figure S3A. The distribution of tissue specific TSS in 56 tissues of
 1168 Benz2616. The bar shows the count of tissue specific TSS in each tissue with the
 1169 proportion being expressed with CTPM > 1 coloured in red.

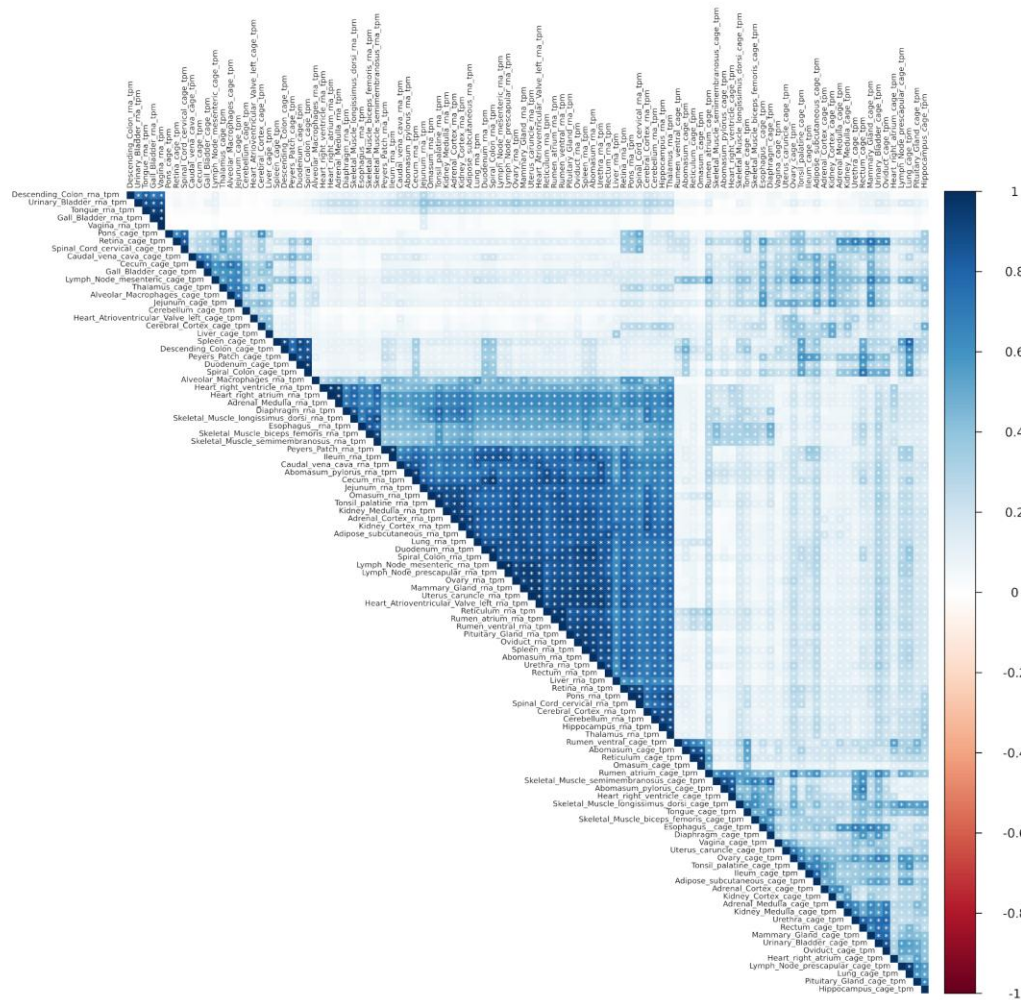
1170



1171

1172 Supplementary Figure S3B. The distribution of tissue specific TSS-Enhancers across
 1173 the 56 tissues from Benz2616. The bars show the count of tissue specific TSS in
 1174 each tissue with the proportion being expressed with CTPM >1 coloured in red.

1175



1176

1177 Supplementary Figure S4. Heatmap of mRNA-Seq and CAGE expression profiles
1178 (TPM and CTPM). The correlation was calculated over 52 matched tissues and 5732
1179 transcripts-TSS expressed in all tissues.

1180

1181 Supplementary Table 1. Details of 5' linker barcodes and pool ID assigned to each
1182 tissue sample.

1183

1184 Supplementary Table 2. Percentage of tissue-specific CAGE tags mapping to
1185 genomic features.

1186

1187 Supplementary Table .3. Summary of WGBS sequencing and mapping results.

1188

1189 Supplementary File 1. A detailed comparison of mapping of the CAGE tags to the
1190 two reference assemblies Oar_v3.1 and Oar_rambouillet_v1.0 and rationale for
1191 selecting the 2/3rd representation threshold for mapped CAGE tags.

1192

1193 Supplementary File 2. Expression data frames from Uni-, Bi-directional, long range
1194 Linked co-expression clustering and transcript level mRNA-Seq from all 56 tissues
1195 (2/3 representation rule applied).

1196