

GigaScience

An improved ovine reference genome assembly to facilitate in depth functional annotation of the sheep genome --Manuscript Draft--

Manuscript Number:	GIGA-D-21-00165R2														
Full Title:	An improved ovine reference genome assembly to facilitate in depth functional annotation of the sheep genome														
Article Type:	Data Note														
Funding Information:	<table border="1"><tr><td>National Institute of Food and Agriculture (2013-67015-21228)</td><td>Dr. Kim C. Worley</td></tr><tr><td>National Institute of Food and Agriculture (2013-67015-21372)</td><td>Dr. Noelle E. Cockett</td></tr><tr><td>National Institute of Food and Agriculture (2017-67016-26301)</td><td>Dr. Brenda M. Murdoch</td></tr><tr><td>International Sheep Genomics Consortium (217201191442)</td><td>Dr. Kim C. Worley</td></tr><tr><td>Agricultural Research Service (5090-31000-026-00-D)</td><td>Dr. Derek M. Bickhart</td></tr><tr><td>Agricultural Research Service (3040-31000-100-00D)</td><td>Dr. Timothy P.L. Smith</td></tr><tr><td>Agricultural Research Service (8042-31000-001-00-D)</td><td>Dr. Benjamin D. Rosen</td></tr></table>	National Institute of Food and Agriculture (2013-67015-21228)	Dr. Kim C. Worley	National Institute of Food and Agriculture (2013-67015-21372)	Dr. Noelle E. Cockett	National Institute of Food and Agriculture (2017-67016-26301)	Dr. Brenda M. Murdoch	International Sheep Genomics Consortium (217201191442)	Dr. Kim C. Worley	Agricultural Research Service (5090-31000-026-00-D)	Dr. Derek M. Bickhart	Agricultural Research Service (3040-31000-100-00D)	Dr. Timothy P.L. Smith	Agricultural Research Service (8042-31000-001-00-D)	Dr. Benjamin D. Rosen
National Institute of Food and Agriculture (2013-67015-21228)	Dr. Kim C. Worley														
National Institute of Food and Agriculture (2013-67015-21372)	Dr. Noelle E. Cockett														
National Institute of Food and Agriculture (2017-67016-26301)	Dr. Brenda M. Murdoch														
International Sheep Genomics Consortium (217201191442)	Dr. Kim C. Worley														
Agricultural Research Service (5090-31000-026-00-D)	Dr. Derek M. Bickhart														
Agricultural Research Service (3040-31000-100-00D)	Dr. Timothy P.L. Smith														
Agricultural Research Service (8042-31000-001-00-D)	Dr. Benjamin D. Rosen														
Abstract:	<p>Background</p> <p>The domestic sheep (<i>Ovis aries</i>) is an important agricultural species raised for meat, wool, and milk across the world. A high-quality reference genome for this species enhances the ability to discover genetic mechanisms influencing biological traits. Further, a high-quality reference genome allows for precise functional annotation of gene regulatory elements. The rapid advances in genome assembly algorithms and emergence of sequencing technologies with increasingly long reads provide the opportunity for an improved de novo assembly of the sheep reference genome.</p> <p>Findings</p> <p>Short-read Illumina (55x coverage), long-read PacBio (75x coverage), and Hi-C data from this ewe retrieved from public databases were combined with an additional 50x coverage of Oxford Nanopore data and assembled with canu v1.9. The assembled contigs were scaffolded using Hi-C data with Salsa v2.2, gaps filled with PBSuite v15.8.24, and polished with Nanopolish v0.12.5. After duplicate contig removal with PurgeDups v1.0.1, chromosomes were oriented and polished with two rounds of a pipeline which consisted of freebayes v1.3.1 to call variants, Merfin to validate them, and BCFtools to generate the consensus fasta. The ARS-UI_Ramb_v2.0 assembly is 2.63 Gb in length and has improved continuity (contig NG50 of 43.18 Mb) with a 19-fold and 38-fold decrease in the number of scaffolds compared with Oar_rambouillet_v1.0 and Oar_v4.0. ARS-UI_Ramb_v2.0 has greater per-base accuracy and fewer insertions and deletions identified from mapped RNA sequence than previous assemblies.</p> <p>Conclusions</p> <p>The ARS-UI_Ramb_v2.0 assembly is a substantial improvement in contiguity that will optimize the functional annotation of the sheep genome and facilitate improved mapping accuracy of genetic variant and expression data for traits in sheep.</p>														
Corresponding Author:	Benjamin D Rosen UNITED STATES														
Corresponding Author Secondary															

Information:	
Corresponding Author's Institution:	
Corresponding Author's Secondary Institution:	
First Author:	Kimberly M Davenport, M.S.
First Author Secondary Information:	
Order of Authors:	Kimberly M Davenport, M.S.
	Derek M. Bickhart
	Kim C. Worley
	Shwetha C. Murali
	Mazdak Salavati
	Emily L. Clark
	Noelle E. Cockett
	Michael P. Heaton
	Timothy P.L. Smith
	Brenda M. Murdoch
	Benjamin D. Rosen
Order of Authors Secondary Information:	
Response to Reviewers:	Dear Hans, Here is the manuscript updated to include the Oxford Nanopore sequence accession numbers (line 175), the updated BUSCO results (lines 321-328), and the GigaDB reference in the "Availability of supporting data" section (line 410). Regards, Ben
Additional Information:	
Question	Response
Are you submitting this manuscript to a special series or article collection?	No
Experimental design and statistics	Yes
Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist . Information essential to interpreting the data presented should be made available in the figure legends.	
Have you included all the information requested in your manuscript?	
Resources	Yes

<p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist?</p>	
<p>Availability of data and materials</p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in publicly available repositories (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.</p> <p>Have you have met the above requirement as detailed in our Minimum Standards Reporting Checklist?</p>	<p>Yes</p>

1 **An improved ovine reference genome assembly to facilitate in depth functional annotation**
2 **of the sheep genome**

3
4 Kimberly M. Davenport¹ [0000-0003-2796-9252], Derek M. Bickhart² Bickhart [0000-0003-
5 2223-9285], Kim Worley³ [0000-0002-0282-1000], Shwetha C. Murali⁴, Mazdak Salavati⁵
6 [0000-0002-7349-2451], Emily L. Clark⁶ [0000-0002-9550-7407], Noelle E. Cockett⁷, Michael
7 P. Heaton⁸ [0000-0003-1386-1208], Timothy P.L. Smith⁹ [0000-0003-1611-6828], Brenda M.
8 Murdoch^{10*} [0000-0001-8675-3473], and Benjamin D. Rosen^{11*} [0000-0001-9395-8346]
9

10 ¹Department of Animal, Veterinary, and Food Sciences, University of Idaho, 875 Perimeter Dr.,
11 Moscow, ID, United States 83843. Email: kmdavenport@uidaho.edu
12

13 ²US Dairy Forage Research Center, USDA-ARS, 1925 Linden Drive, Madison, WI, United
14 States 53706. Email: derek.bickhart@usda.gov
15

16 ³Baylor College of Medicine, One Baylor Plaza, Houston, TX, United States 77030. Email:
17 kworley@bcm.edu
18

19 ⁴Baylor College of Medicine, One Baylor Plaza, Houston, TX, United States 77030.
20 Email: shwethac@gmail.com
21

22 ⁵The Roslin Institute, Royal (Dick) School of Veterinary Studies, The University of Edinburgh,
23 Easter Bush Campus, Midlothian, United Kingdom, EH25 9RG, United Kingdom. Email:
24 mazdak.salavati@roslin.ed.ac.uk
25

26 ⁶The Roslin Institute, Royal (Dick) School of Veterinary Studies, The University of Edinburgh,
27 Easter Bush Campus, Midlothian, United Kingdom, EH25 9RG. Email:
28 emily.clark@roslin.ed.ac.uk
29

30 ⁷Utah State University, Old Main Hill, Logan, UT 84322. Email: noelle.cockett@usu.edu
31

32 ⁸US Meat Animal Research Center, USDA-ARS, State Spur 18D, Clay Center, NE 68933.
33 Email: mike.heaton@usda.gov
34

35 ⁹US Meat Animal Research Center, USDA-ARS, State Spur 18D, Clay Center, NE 68933.
36 Email: tim.smith2@usda.gov
37

38 ¹⁰Department of Animal, Veterinary, and Food Sciences, University of Idaho, 875 Perimeter Dr.,
39 Moscow, ID 83843. Email: bmurdoch@uidaho.edu
40

41 ¹¹Animal Genomics and Improvement Laboratory, USDA-ARS, 10300 Baltimore Avenue,
42 Beltsville, MD 20705. Email: ben.rosen@usda.gov
43

44 Correspondence:

45 Brenda M. Murdoch

46 bmurdoch@uidaho.edu

47 Benjamin D. Rosen
48 ben.rosen@usda.gov
49 **Abstract**

50

51 *Background*

52

53 The domestic sheep (*Ovis aries*) is an important agricultural species raised for meat, wool, and
54 milk across the world. A high-quality reference genome for this species enhances the ability to
55 discover genetic mechanisms influencing biological traits. Further, a high-quality reference
56 genome allows for precise functional annotation of gene regulatory elements. The rapid advances
57 in genome assembly algorithms and emergence of increasingly long sequence read length
58 provide the opportunity for an improved *de novo* assembly of the sheep reference genome.

59

60

61 *Findings*

62

63 Short-read Illumina (55x coverage), long-read PacBio (75x coverage), and Hi-C data from this
64 ewe retrieved from public databases were combined with an additional 50x coverage of Oxford
65 Nanopore data and assembled with canu v1.9. The assembled contigs were scaffolded using Hi-
66 C data with Salsa v2.2, gaps filled with PBSuitev15.8.24, and polished with Nanopolish v0.12.5.
67 After duplicate contig removal with PurgeDups v1.0.1, chromosomes were oriented and polished
68 with two rounds of a pipeline which consisted of freebayes v1.3.1 to call variants, Merfin to
69 validate them, and BCFtools to generate the consensus fasta. The ARS-UI_Ramb_v2.0 assembly
70 has improved continuity (contig N50 of 43.18 Mb) with a 19-fold and 38-fold decrease in the

71 number of scaffolds compared with Oar_rambouillet_v1.0 and Oar_v4.0. ARS-UI_Ramb_v2.0
72 has greater per-base accuracy and fewer insertions and deletions identified from mapped RNA
73 sequence than previous assemblies.

74

75

76 *Conclusions*

77

78 The ARS-UI_Ramb_v2.0 assembly is a substantial improvement that will optimize the
79 functional annotation of the sheep genome and facilitate improved mapping accuracy of genetic
80 variant and expression data for traits in sheep.

81

82

83 **Keywords:** Rambouillet, genome assembly, reference genome, sheep

84

85

86

87

88

89

90

91

92

93

94

95 **Context**

96

97 The domestic sheep (*Ovis aries*) is a globally important livestock species raised for a variety of
98 purposes including meat, wool, and milk. Domestication likely occurred in multiple events
99 approximately 11,000 years ago [1-4]. Selection for desirable traits including meat, wool, and
100 milk began approximately 4,000-5,000 years ago [2,4]. Modern sheep breeds exhibit a wide
101 variety of phenotypes and adaptations to specific environments, for example the enhanced
102 parasite tolerance evident in hair sheep [5,6]. As many as 1,400 breeds of sheep exist today [7-9]
103 including the Rambouillet breed developed in France from a Merino fine wool lineage that is
104 regarded for its ability to produce high quality wool as well as meat products in production
105 systems across the world [10,11].

106

107 Genome research in sheep holds promise to improve efficiency and sustainability of production
108 and reduce the environmental effects of animal agriculture [12]. The first sheep reference
109 genome assembly was based on whole genome shotgun (WGS) short-read sequencing,
110 scaffolded by genetic linkage and radiation hybrid maps. The sequence came from two unrelated
111 Texel breed sheep, with the first assembly draft (Oar_v3.1; International Sheep Genomics
112 Consortium, 2010) having a contig N50 of 40 kilobases (kb) and the update (Oar_v4.0) [13]
113 boosting the N50 metric to 150 kb. More recently, the Ovine Functional Annotation of Animal
114 Genomes (FAANG) project proposed to perform a variety of genome annotation assays for
115 dozens of tissues from a single animal [14,15]. To maximize the success of assays that depend on
116 mapping sequence data to a reference, the FAANG project assembled the genome of that animal,

117 a female of the Rambouillet breed. The assembly, released in 2017 (Oar_rambouillet_v1.0,
118 GenBank accession GCF_002742125; Worley et al., unpublished) is based on a combination of
119 Pacific Biosciences RSII WGS long-read and Illumina short-read sequencing. It has an improved
120 contig N50 of 2.6 megabases (Mb) and is generally regarded as the official reference assembly
121 for global sheep research.

122

123 The continued maturation of long read sequencing technologies provided an opportunity to
124 improve upon the sheep reference genome assembly. Since most of the proposed FAANG
125 annotation assays had already been performed on the Rambouillet ewe, lung tissue from the
126 same animal was chosen for DNA extraction. This allowed the use of existing long read data to
127 supplement new, longer-read, Oxford Nanopore PromethION sequencing. We report a *de novo*
128 assembly of the same Rambouillet ewe used for Oar_rambouillet_v1.0, based on approximately
129 50x coverage of nanopore reads (N50 47kb) and 75x coverage PacBio reads (N50 13kb). The
130 new assembly, ARS-UI_Ramb_v2.0 offers a 20-fold improvement in contiguity and increased
131 accuracy, providing a basis for regulatory element annotation in the FAANG project and
132 facilitating the discovery of biological mechanisms that influence traits important in global sheep
133 research and production.

134

135

136 **Methods**

137

138 *Sampling Strategy*

139

140 The fullblood Rambouillet ewe used for this genome assembly (Benz 2616, USMARC ID
141 200935900) (Figure 1) was selected by the Ovine Functional Annotation of Animal Genomes
142 project and acquired from the USDA. Tissues were collected postmortem from the healthy six-
143 year-old ewe as approved by the Utah State University Institutional Animal Care and Use
144 Committee. A full description of the tissue collection strategy is available in the FAANG Data
145 Coordination Center [15,16]. Details regarding the tissues collected from the animal are available
146 under BioSample number SAMEG329607 [17].

147

148

149 *Sequencing and Data Acquisition*

150

151 DNA was extracted from approximately 50 mg of lung tissue using phenol:chloroform-based
152 method as described [18]. Briefly, the frozen tissue was pulverized in a cryoPREP CP02 tissue
153 disruption system (Covaris Inc., Woburn MA) as recommended by the manufacturer. The
154 powdered tissue was transferred to a 50 mL conical tube and mixed in 200 μ L of phosphate
155 buffered saline (Sigma-Aldrich, St. Louis MO). The tissue was then diluted in 10 mL of buffer
156 TLB (100mM NaCl, 10mM Tris-HCl pH 8.0, 25 mM EDTA, 0.5% SDS) and mixed by
157 vortexing, then incubated with 20 μ L 10 mg/mL RNase A at 37°C for one hour with gentle
158 shaking. Protein digestion was performed with 100 μ L Proteinase K (20 mg/mL) at 50°C for 2
159 hours, with slow rotation of the tube to mix every 30 minutes. The lysate was distributed equally
160 into two 15 mL Phase Lock tubes (Quantabio, Beverly MA) and each tube received 5 mL of TE-
161 saturated Phenol (Sigma-Aldrich, St. Louis MO) followed by mixing on a tube rotator at 20
162 RPM for 10 minutes at 22°C. The aqueous layer was collected after separating at 2300xg for 10

Commented [HZ1]: = Logsdon 2019 ?

163 minutes and transferred to another Phase Lock tube. A second extraction performed in the same
164 way as the first was conducted using 2.5 mL phenol and 2.5 mL chloroform:isoamyl alcohol
165 (Sigma). The final aqueous phase was transferred to a 50 mL conical tube and the DNA
166 precipitated with 2 mL of 5M ammonium acetate and 15 mL of ice-cold 100% ethanol. The
167 DNA was pulled from the alcohol using a Pasteur pipet “hook” and placed in 10 mL of cold 70%
168 ethanol to wash the pellet. The ethanol was poured off and the DNA pellet dried for 20-30
169 minutes, then dissolved in a dark drawer at room temperature for 48 hours in 1 mL of 10mM
170 Tris-Cl pH 8.5. Library preparation for Oxford Nanopore long read sequencing was performed
171 with an LSK-109 template preparation kit as recommended by the manufacturer (Oxford
172 Nanopore, Oxford U.K.) with modifications as described by Logsdon [18]. The ligated template
173 was sequenced with a PromethION instrument using four R9.4 flow cells. (Oxford Nanopore
174 Technologies, Oxford, United Kingdom). Output as fast5 files were basecalled with Guppy v3.1
175 [19]. Fastq files are available under the Sequence Read Archive (SRA) accessions
176 SRR17080040-SRR17080043.

177
178 Sequence data used in the previous Oar_rambouillet_v1.0 assembly was retrieved from the SRA
179 listed under project number PRJNA414087 [15]. PacBio RS II sequence generated from DNA
180 extracted from whole blood was retrieved from SRX3445660, SRX3445661, SRX3445662, and
181 SRX3445663. The Hi-C sequence data generated from liver using HindIII enzyme and
182 sequenced at 150 bp paired end with an Illumina HiSeq X Ten was retrieved from SRX3399085
183 and SRX3399086. Short read whole genome sequencing from DNA extracted from whole blood
184 collected from the Rambouillet ewe was performed with an Illumina HiSeq X Ten sequenced at
185 150 bp paired end and was retrieved from SRX3405602. Further details about these sequences

186 can be found under the umbrella project number PRJNA414087. Short read 45 bp paired end
187 whole genome sequence from an Illumina Genome Analyzer II generated from Texel sheep used
188 in previous genome assemblies were retrieved from the SRA under accessions SRX511533-
189 SRX511565 (BioProject PRJNA169880).

190

191

192 *Assembly*

193

194 Contigs were assembled with Oxford Nanopore and PacBio reads generated as described above
195 using canu v1.8 (Canu, RRID:SCR_015880) through the trimmed reads stage of assembly.

196 Parameters for contig construction were set as “batOptions=-dg 4 -db 4 -mo 1000” [20]. Canu

197 v1.9 was used to complete the contig assembly because this update demonstrates better

198 consensus generation of the overlapped contigs in the final step in the assembly process [21,22].

199 The corrected error rate option was set as “correctedErrorRate=0.105.”

200

201

202 *Scaffolding*

203

204 Two Hi-C datasets from liver tissue from two different library preparations were retrieved as

205 described above. The Hi-C reads were first aligned to the polished contigs using the Arima

206 Genomics mapping pipeline [23]. This pipeline first maps paired end reads individually with

207 bwa-mem, then removes the 3' end of reads identified as chimeric and span ligation junctions.

208 Reads were then paired, filtered by mapping quality with samtools [24], and PCR duplicates

209 removed with Picard [25]. The two Hi-C libraries were merged in the final step in the Arima
210 pipeline to generate the merged BAM file. The BAM file was converted to a BED file for input
211 into Salsa using the bedtools command bamToBed (BEDTools, RRID:SCR_006646) [26]. Salsa
212 v2.2 was used for scaffolding by implementing “python run_pipeline.py -a contigs.fasta -l
213 contigs.fasta.fai -b alignment.bed -e HindIII -o scaffolds -m yes” [27].

214
215 The Hi-C reads were aligned to the scaffolded assembly with the Arima Genomics mapping
216 pipeline and then processed with PretextView to visually evaluate the scaffolds as a contact map
217 in PretextView [28]. The scaffolded assembly was also compared to *Oar_rambouillet_v1.0* by
218 aligning the two genomes with “minimap2 -cx asm5 Oar_rambouillet_v1.0_genomic.fasta
219 scaffolds.fasta > alignment.paf” [29]. A dotplot of the alignment was visualized with D-Genies
220 [30]. Scaffolds were edited based on visual inspection of the contact map and dotplot, as well as
221 the Hi-C alignment file. Scaffold joins and rearrangements were incorporated to the assembly
222 using the *agp2fasta* mode of CombineFasta [31].

223

224

225 *Gap Filling and Polishing*

226

227 Gap filling was completed with pbsuite v15.8.24 using both the PacBio and Oxford Nanopore
228 reads. Nanopolish v0.12.5 (Nanopolish, RRID:SCR_016157) [32] with the NanoGrid parallel
229 wrapper [33] was employed with the raw fast5 files generated from the PromethION sequencing
230 to polish the assembly. Duplicates were removed using PurgeDups v1.0.1 [34]. The chromosome
231 orientation was confirmed in the polished assembly by identifying telomeres and centromeres

232 using RepeatMasker v4.1.1 (**RepeatMasker**, RRID:SCR_012954) [35]. The mitochondrial genome
233 was identified by aligning the previously annotated mitochondrial sequence from
234 Oar_rambouillet_v1.0 (RefSeq NC_001941.1) to the assembly contigs. Chromosomes were
235 oriented centromere to telomere and placed in chromosome number order. The final polishing
236 was performed with two rounds of freebayes v1.3.1 (**FreeBayes**, RRID:SCR_010761) using the
237 Illumina short read data after final chromosome orientations and mitochondrial genome were
238 confirmed [36]. Variants used for polishing with both Nanopolish and freebayes were screened
239 with Merfin [37] which evaluates the k-mer consequences of variant calls and filters unsupported
240 variants.

241

242

243 *RNA Sequencing*

244

245 RNA sequencing data was generated from five tissues including skin, thalamus, pituitary, lymph
246 node (mesenteric), and abomasum pylorus collected from the animal used to assemble the
247 reference genome. Details regarding the RNA isolation protocol, library preparation, and
248 sequencing as well as the raw data can be found in GenBank under BioProject PRJEB35292,
249 specifically under SRA run numbers ERR3665717 (skin), ERR3728435 (thalamus),
250 ERR3650379 (pituitary), ERR3665711 (lymph node mesenteric), and ERR3650373 (abomasum
251 pylorus). Reads were trimmed with Trim Galore v0.6.4 [38] and alignment to both Rambouillet
252 genomes was performed with STAR v2.7 using default parameters [39]. Indels were identified
253 with bcftools mpileup, filtering allele depth (AD) at > 5 [40].

254

255

256 **Annotation**

257 The annotation for ARS-UI_Ramb_v2.0, NCBI Ovis aries Annotation Release 104, is available
258 in RefSeq and other NCBI genome resources [41].

259

260 Here we also provide a liftover of the annotation for Oar_rambouillet_v1.0 onto ARS-
261 UI_Ramb_v2.0. The annotation used for the liftover was NCBI v103
262 GCF_002742125.1_Oar_rambouillet_v1.0_genomic.fna.gz. The GFF3 format gene annotation
263 file was prepared for processing using liftOff v1.5.2 [42]. A set of matching chromosome names
264 for Oar_rambouillet_v1.0 and ARS-UI_Ramb_v2.0 were generated according to the instructions
265 for liftOff (*paste -d " " <(cut -d ' ' -f1 ramb1.chr) <(cut -d ' ' -f1 ramb2.chr) > chroms.txt*). The
266 GFF file (annotation Ramb1LO2) generated by liftOff is included in Supplementary File 1
267 (Ramb_v1.0_NCBI103_lifted_over_ARS-UI_Ramb_v2.0.gff.gz).

268

269 To compare the breakdown of transcripts captured by the three annotations
270 (Oar_Rambouillet_v1.0, Ramb1LO2 (liftover) and ARS-UI_Ramb_v2.0), we generated
271 transcript expression estimates using Kallisto v0.44.0 (kallisto, RRID:SCR_016582) [43]. For
272 the lifted over gene annotation the GFF file (Ramb_v1.0_NCBI103_lifted_over_ARS-
273 UI_Ramb_v2.0.gff.gz) was used to generate transcriptome sequence FASTA files, as a Kallisto
274 index, for transcript expression estimation. Briefly, exonic blocks were extracted from the GFF3
275 file using the awk command (*awk '(\$3~/exon/)' input.gff*). The getfasta and groupby plugins
276 from bedtools v2.30.0 [44] were used to extract the exonic sequences and group them by
277 transcript name. Exonic sequences for each transcript were appended in the correct order, to

278 produce the complete sequence for each transcript. The FASTA format file for the whole
279 transcriptome was created using all of the transcript level FASTA sequences for the liftover
280 annotation Ramb1LO2 (Supplementary File 2; Ramb1LO2_NCBI103_geneBank_rna.fa). The
281 set of scripts used for this step are included in Supplementary File 3. The Kallisto indices for
282 Oar_Rambouillet_v1.0 (GCF_002742125.1_Oar_rambouillet_v1.0_rna.fna.gz), Ramb1LO2
283 (liftover; Ramb1LO2_NCBI103_geneBank_rna.fa) and ARS-UI_Ramb_v2.0
284 (GCF_016772045.1_ARS-UI_Ramb_v2.0_rna.fna.gz) were then used with the RNA-Seq data
285 from the 61 tissues from Benz2616 (GenBank BioProject PRJNA414087 and PRJEB35292) to
286 estimate transcript level expression for every tissue as transcript per million mapped reads
287 (TPM) and compared across the three annotations.

288

289

290 **Data Validation and Quality Control**

291

292 *Assembly Quality Statistics*

293

294 The four flow cells of PromethION data produced 136 gigabases (Gb) of WGS sequence
295 (approximately 51x coverage) in reads having a read N50 of 47 kb. The initial generation of
296 contigs used this data as well as 198.1 Gb of RSII data with a read N50 of 12.9 kb. The ARS-
297 UI_Ramb_v2.0 assembly was submitted to NCBI GenBank under accession number
298 GCF_016772045.1, and statistics of contigs and scaffolds following initial polishing, scaffolding
299 with Hi-C data and manual editing, gap-filling, and final polishing, are shown in Table 1. The
300 assembly improved on the Oar_v4.0/Oar_rambouillet_v1.0 sheep reference assemblies in all

301 continuity measures (Table 1) including a 286/17-fold increase in contig N50 (the size of the
302 shortest contig for which all larger contigs contain half of the total assembly), a 214/33-fold
303 reduction in the number of contigs in the assembly and concomitant 209/13-fold reduction of
304 contig L50 (the number of contigs making up half of the total assembly), and 38/19-fold
305 reduction in total number of scaffolds. Manual curation of scaffolds using Hi-C data improved
306 scaffold continuity and led to chromosome length scaffolds (Figure 2).

307
308 The Themis-ASM pipeline [45] was implemented to further assess assembly quality and
309 compare sheep genome assemblies. Short read sequence from both the Rambouillet ewe used in
310 this assembly and Texel sheep from previous sheep genome assemblies were used to compare
311 ARS-UI_Ramb_v2.0 with Oar_rambouillet_v1.0 and Oar_v4.0 assemblies.

312
313 The k-mer based quality value and error rates improved with ARS-UI_Ramb_v2.0 compared
314 with Oar_rambouillet_v1.0 and Oar_v4.0. This is also reflected in the proportion of complete
315 assembly based on k-mers (merCompleteness), which is similar between ARS-UI_Ramb_v2.0
316 and Oar_rambouillet_v1.0 and both are higher than Oar_v4.0. Further, the SNP and indel quality
317 value (baseQV) were greatest overall in ARS-UI_Ramb_v2.0 (41.84), followed by
318 Oar_rambouillet_v1.0 (40.69) and Oar_v4.0 (32.40). The percentage of short reads not mapped
319 to the genome was $\leq 1\%$ in all three assemblies.

320
321 The completeness of ARS-UI_Ramb_v2.0 was evaluated by examining the presence or absence
322 of evolutionarily conserved genes in each assembly using Benchmarking Universal Single-Copy
323 Ortholog (BUSCO, RRID:SCR_015008)) v5.2.2 scores with the cetartiodactyla_odb10 dataset

324 and metaeuk gene predictor [46]. The percent of single copy complete BUSCOs were higher
325 (93.9%) in ARS-UI_Ramb_v2.0 when compared with Oar_rambouillet_v1.0 (93.0%) and
326 Oar_v4.0 (91.2%). Complete duplicated BUSCO percentage was highest in
327 Oar_rambouillet_v1.0 (2.6%) compared with ARS-UI_Ramb_v2.0 (2.1%), and lowest in
328 Oar_v4.0 (1.6%). Further, ARS-UI_Ramb_v2.0 had the lowest percent of fragmented and
329 missing BUSCOs (0.9% and 3.1%, respectively) compared with Oar_rambouillet_v1.0 (1.1%
330 and 3.3%, respectively) and Oar_v4.0 (2.4% and 4.8%, respectively).

331
332 The three sheep genome assemblies were also compared with a feature response curve in which
333 the quality of the assembly is analyzed as a function of the features, or maximum number of
334 possible errors, allowed in the contigs (Figure 3) [47]. Both the ARS-UI_Ramb_v2.0 and
335 Oar_v4.0 feature response curves peak higher and to the left of Oar_rambouillet_v1.0, which
336 indicate fewer errors in these assemblies (Figure 3A). The ARS-UI_Ramb_v2.0 genome also has
337 fewer regions with either low or high coverage overall and for paired reads, suggesting fewer
338 coverage issues, as well as fewer improperly paired or unmapped single reads when compared
339 with other assemblies (Figure 3B). The number of high Comp/Expansion (CE) statistics in ARS-
340 UI_Ramb_v2.0 was intermediate between Oar_rambouillet_v1.0 (higher) and Oar_v4.0 (lower),
341 however this latest assembly had the lowest number of regions with low CE statistics.

342
343 Comparative alignment of ARS-UI_Ramb_v2.0 with previous assemblies Oar_rambouillet_v1.0
344 and Oar_v4.0 and visualization with a dotplot revealed a high amount of agreement between
345 assemblies (Figure 4). Interestingly, chromosome 11 was improperly oriented in
346 Oar_rambouillet_v1.0, and after confirming centromere and telomere locations on this

347 chromosome, this was resolved in the ARS-UI_Ramb_v2.0 assembly. The percent identity
348 between ARS-UI_Ramb_v2.0 is very high when compared with Oar_rambouillet_v1.0 which
349 was expected considering the same animal was used in both assemblies. However, Oar_v4.0 was
350 assembled from Texel sheep, which is apparent in the percent identity in the dotplot.

351

352 In summary, ARS-UI_Ramb_v2.0 offers greater contiguity, improved quality, more complete
353 BUSCOs, and fewer assembly errors when compared with previous assemblies.

354

355

356 *RNA sequencing alignment*

357

358 Insertions and deletions (indels) in the ARS-UI_Ramb_v2.0 assembly were characterized and
359 compared with Oar_rambouillet_v1.0 by mapping 150 bp paired-end RNA-seq data from skin,
360 thalamus, pituitary, lymph node (mesenteric), and abomasum pylorus generated from the same
361 animal used to assemble the reference genome. In all five tissues, ARS-UI_Ramb_v2.0 had
362 nearly half of the number of indels compared with Oar_rambouillet_v1.0. Most indels identified
363 in both assemblies were 1bp in length. The ARS-UI_Ramb_v2.0 had a greater number of
364 uniquely mapped reads in each tissue when compared with Oar_rambouillet_v1.0, leading to an
365 approximate 2% increase in the percent of uniquely mapped reads in most tissues except
366 pituitary, which saw an almost 13% improvement. The number of reads that mapped to multiple
367 loci decreased in the new assembly by 12.59% in pituitary, and 1-2% in other tissues. Further,
368 ARS-UI_Ramb_v2.0 had fewer unmapped reads than Oar_rambouillet_v1.0 across all five
369 tissues by an average of 0.15%.

370

371 *Annotation*

372

373 The ARS-UI_Ramb_v2.0 annotation represents a substantial improvement over the annotation
374 on Oar_rambouillet_v1.0. For example, for ARS-UI_Ramb_v2.0 16,500 coding genes have an
375 ortholog to human (compared to 16,319 for Oar_rambouillet_v1.0), and the BUSCO scores
376 demonstrate that 99.1% of the gene models (cetartiodactyla_odb10) are complete in the new
377 annotation versus 98.8% in the previous one. The annotation for ARS-UI_Ramb_v2.0 includes
378 Iso-Sequencing for 8 tissues to improve contiguity of gene models, and CAGE sequencing for 56
379 tissues to define TSS, that were not used to annotate Oar_rambouillet_v1.0. The full report for
380 the annotation release is available at:

381 (https://www.ncbi.nlm.nih.gov/genome/annotation_euk/Ovis_aries/104).

382

383 Using Kallisto we compared the number of expressed transcripts, for the RNA-Seq dataset of 61
384 tissue samples from Benz2616, across the three annotations (Oar_Rambouillet_v1.0, RambILO2
385 (liftover) and ARS-UI_Ramb_v2.0). There was a considerable increase in the number of
386 transcripts captured by the annotation for ARS-UI_Ramb_v2.0 (60,064) relative to
387 Oar_Rambouillet_v1.0 (42,058) and the liftover annotation (RambILO2) (40,910) (Figure 5).
388 This equates to approximately 20,000 new annotated gene models for ARS-UI_Ramb_v2.0 and
389 further reflects the substantial improvement over the annotation for Oar_Rambouillet_v1.0.

390 The lifted over annotation we have generated will provide a resource for those who wish to
391 compare their results for ARS-UI_Ramb_v2.0 to previous work using

392 Oar_Rambouillet_v1.0. Only 2.7% of protein coding transcripts were lost (1148) lifting over the

393 annotation for Oar_Rambouillet_v1.0 onto ARS-UI_Ramb_v2.0. According to the annotation
394 report provided by NCBI
395 (https://www.ncbi.nlm.nih.gov/genome/annotation_euk/Ovis_aries/104/), 70% of the annotations
396 were identical or had only minor changes between and Oar_Rambouillet_v1.0 and ARS-
397 UI_Ramb_v2.0.

398

399

400 **Re-use potential**

401

402 The ARS-UI_Ramb_v2.0 genome assembly serves as a reference for genetic investigation of
403 traits important in sheep research and production across the world. This genome is assembled
404 from the same animal used in the Ovine FAANG Project, which provides a high-quality basis for
405 epigenetic annotation to serve the international sheep genomics community and scientific
406 community at large.

407

408

409 **Availability of supporting data**

410

411 The data sets supporting the results of this article are available in the RefSeq repository,
412 GCF_016772045.1, and in the GigaScience Database [50]. RNA sequencing data is available
413 under BioProject PRJEB35292. Ovis aries Annotation Release 104 is also available in RefSeq
414 and other NCBI genome resources [41].

415

416

417 Additional files

418 Supplementary File 1 – Ramb_v1.0_NCBI103_lifted_over_ARS-UI_Ramb_v2.0.gff.gz

419 Supplementary File 2 – Ramb1LO2_NCBI103_geneBank_rna.fa

420 Supplementary File 3 – Supplementary_File_3_scripts.txt

421

422

423 Author contributions

424

425 BMM, TPLS, DMB, and BDR conceptualized the study. BMM, NEC, MPH, and TPLS selected

426 the animal and collected samples. KW and SCM facilitated the generation of RSII, short read,

427 and Hi-C data. TPLS facilitated the nanopore long read data generation. KMD, DMB, TPLS,

428 BMM, and BDR performed the genome assembly, scaffolding, RNA-sequencing alignment,

429 polishing, and quality control. MS and ELC contributed the section describing the LiftOff

430 annotation and comparative analysis of transcript expression estimates for the three annotations.

431 KMD, DMB, TPLS, BMM, and BDR generated tables and figures and drafted the manuscript.

432 KMD, DMB, KW, SCM, NEC, TPLS, BMM, and BDR edited the manuscript. All authors

433 contributed to the article and approved the final version.

434

435

436 Acknowledgements

437

438 The authors thank Dr. Kristen Kuhn for technical support and Dr. Kreg Leymaster for overseeing
439 the acquisition, animal care and housing, and interstate transportation of the Rambouillet ewe.

440

441

442 **Funding**

443

444 Funding was provided by Agriculture and Food Research Initiative Competitive grants from the
445 USDA National Institute of Food and Agriculture supporting improvements of the sheep
446 genomes (2013-67015-21228) and FAANG activities (2013-67015-21372, 2017-67016-26301).

447 Additional funding was received from the International Sheep Genome Consortium
448 (217201191442) and infrastructure support from a grant to R. Gibbs from the NIH NHGRI
449 Large-Scale Sequencing Program (U54 HG003273).

450

451 DMB was supported by appropriated USDA CRIS project 5090-31000-026-00-D. TPLS was
452 supported by appropriated USDA CRIS Project 3040-31000-100-00D. BDR was supported by
453 appropriated USDA CRIS Project 8042-31000-001-00-D. The USDA does not endorse any
454 products or services. Mentioning of trade names is for information purposes only. The USDA is
455 an equal opportunity employer.

456

457

458 **References**

459

- 460 1. Pedrosa S, Uzun M, Arranz JJ, Gutiérrez-Gil B, San Primitivo F, Bayón Y. Evidence of three
461 maternal lineages in Near Eastern sheep supporting multiple domestication events. *Proc Biol*
462 *Sci.* 2005;272:2211-7.
- 463
- 464 2. Zeder MA. Domestication and early agriculture in the Mediterranean Basin: origins,
465 diffusion, and impact. *Proc Natl Acad Sci USA.* 2008;105:11597-604.
- 466
- 467 3. Chessa B, Pereira F, Arnaud F, Amorim A, Goyache F, Mainland I, Kao RR, Pemberton JM,
468 Beraldi D, Stear MJ, Alberti A, Pittau M, Iannuzzi L, Banabazi MH, Kazwala RR, Zhang
469 YP, Arranz JJ, Ali BA, Wang Z, Uzun M, Dione MM, Olsaker I, Holm LE, Saarma U,
470 Ahmad S, Marzanov N, Eythorsdottir E, Holland MJ, Ajmone-Marsan P, Bruford MW,
471 Kantanen J, Spencer TE, Palmarini M. Revealing the history of sheep domestication using
472 retrovirus integrations. *Science.* 2009;324:532-6.
- 473
- 474 4. Kijas JW, Lenstra JA, Hayes B, Boitard S, Porto Neto LR, San Cristobal M, Servin B,
475 McCulloch R, Whan V, Gietzen K, Paiva S, Barendse W, Ciani E, Raadsma H, McEwan J,
476 Dalrymple B, International Sheep Genomics Consortium Members. Genome-wide analysis
477 of the world's sheep breeds reveals high levels of historic mixture and strong recent selection.
478 *PLoS Biol.* 2012;doi:10.1371/journal.pbio.1001258.
- 479
- 480 5. Burke JM, Miller JE. Relative resistance of Dorper crossbred ewes to gastrointestinal
481 nematode infection compared with St. Croix and Katahdin ewes in the southeastern United
482 States. *Vet Parasitol.* 2002;109:265-75.

- 483
- 484 6. Bowdridge SA, Zajac AM, Notter DR. St. Croix sheep produce a rapid and greater cellular
485 immune response contributing to reduced establishment of *Haemonchus contortus*. *Vet*
486 *Parasitol.* 2015;208:204-10.
- 487
- 488 7. Scherf BD. World watch list for domestic animal diversity. 3rd ed. Rome: Food and
489 Agriculture Organization of the United Nations; 2000.
- 490
- 491 8. Lv FH, Agha S, Kantanen J, Colli L, Stucki S, Kijas JW, Joost S, Li MH, Ajmone Marsan P.
492 Adaptations to climate-mediated selective pressures in sheep. *Mol Biol Evol.* 2014;31:3324-
493 43.
- 494
- 495 9. Cao YH, Xu SS, Shen M, Chen ZH, Gao L, Lv FH, Xie XL, Wang XH, Yang H, Liu CB,
496 Zhou P, Wan PC, Zhang YS, Yang JQ, Pi WH, Hehua E, Berry DP, Barbato M,
497 Esmailizadeh A, Nosrati M, Salehian-Dehkordi H, Dehghani-Qanatqestani M, Dotsev AV,
498 Deniskova TE, Zinovieva NA, Brem G, Štěpánek O, Ciani E, Weimann C, Erhardt G,
499 Mwacharo JM, Ahbara A, Han JL, Hanotte O, Miller JM, Sim Z, Coltman D, Kantanen J,
500 Bruford MW, Lenstra JA, Kijas J, Li MH. Historical Introgression from Wild Relatives
501 Enhanced Climatic Adaptation and Resistance to Pneumonia in Sheep. *Mol Biol Evol.*
502 2021;38:838-55.
- 503
- 504 10. Dickinson WF, Lush JL. Inbreeding and the genetic history of the Rambouillet sheep in
505 America. *J Hered.* 1933;24:19-33.

- 506
- 507 11. Zhang L, Mousel MR, Wu X, Michal JJ, Zhou X, Ding B, Dodson MV, El-Halawany NK,
508 Lewis GS, Jiang Z. Genome-wide genetic diversity and differentially selected regions among
509 Suffolk, Rambouillet, Columbia, Polypay, and Targhee sheep. *PLoS One*. 2013;doi:
510 10.1371/journal.pone.0065942.
- 511
- 512 12. Rexroad C, Vallet J, Matukumalli LK, Reecy J, Bickhart D, Blackburn H, Boggess M, Cheng
513 H, Clutter A, Cockett N, Ernst C, Fulton JE, Liu J, Lunney J, Neiberghs H, Purcell C, Smith
514 TPL, Sonstegard T, Taylor J, Telugu B, Eenennaam AV, Tassell CPV, Wells K. Genome to
515 Phenome: Improving Animal Health, Production, and Well-Being - A New USDA Blueprint
516 for Animal Genome Research 2018-2027. *Front Genet*. 2019;10:327.
- 517
- 518 13. Jiang Y, Xie M, Chen W, Talbot R, Maddox JF, Faraut T, Wu C, Muzny DM, Li Y, Zhang
519 W, Stanton JA, Brauning R, Barris WC, Hourlier T, Aken BL, Searle SMJ, Adelson DL,
520 Bian C, Cam GR, Chen Y, Cheng S, DeSilva U, Dixen K, Dong Y, Fan G, Franklin IR, Fu S,
521 Guan R, Highland MA, Holder ME, Huang G, Ingham AB, Jhangiani SN, Kalra D, Kovar
522 CL, Lee SL, Liu W, Liu X, Lu C, Lv T, Mathew T, McWilliam S, Menzies M, Pan S,
523 Robelin D, Servin B, Townley D, Wang W, Wei B, White SN, Yang X, Ye C, Yue Y, Zeng
524 P, Zhou Q, Hansen JB, Kristensen K, Gibbs RA, Flicek P, Warkup CC, Jones HE, Oddy VH,
525 Nicholas FW, McEwan JC, Kijas J, Wang J, Worley KC, Archibald AL, Cockett N, Xu X,
526 Wang W, Dalrymple BP. The sheep genome illuminates biology of the rumen and lipid
527 metabolism. *Science*. 2014;344:1168-1173.
- 528

- 529 14. Murdoch BM. The functional annotation of the sheep genome project. *J Anim Sci.*
530 2019;97:16.
531
- 532 15. Salavati M, Caulton A, Clark R, Gazova I, Smith TPL, Worley KC, Cockett NE, Archibald
533 AL, Clarke SM, Murdoch BM, Clark EL. Global Analysis of Transcription Start Sites in the
534 New Ovine Reference Genome (*Oar rambouillet v1.0*). *Front Genet.* 2020;11:580580.
535
- 536 16. FAANG Data Coordination Center. 2016.
537 https://data.faang.org/api/fire_api/samples/USU_SOP_Ovine_Benz2616_Tissue
538 [_Collection_20160426.pdf](https://data.faang.org/api/fire_api/samples/USU_SOP_Ovine_Benz2616_Tissue).
539
- 540 17. European Bioinformatics Institute, BioSample SAMEG329607. 2016.
541 <https://www.ebi.ac.uk/biosamples/samples/SAMEG329607>.
542
- 543 18. Logsdon G. HMW gDNA purification and ONT ultra-long-read data generation
544 <https://dx.doi.org/10.17504/protocols.io.bchhit36>
545
- 546 19. Guppy (2021). Guppy basecaller (Version 3.1) www.nanoporetech.com.
547
- 548 20. Koren S, Walenz BP, Berlin K, Miller JR, Bergman NH, Phillippy AM. Canu: scalable and
549 accurate long-read assembly via adaptive *k*-mer weighting and repeat separation. *Genome*
550 *Res.* 2017;27:722-36.
551

- 552 21. Heaton MP, Smith TPL, Bickhart DM, Vander Ley BL, Kuehn LA, Oppenheimer J, Shafer
553 WR, Schuetze FT, Stroud B, McClure JC, Barfield JP, Blackburn HD, Kalbfleisch TS,
554 Davenport KM, Kuhn KL, Green RE, Shapiro B, Rosen BD. A Reference Genome Assembly
555 of Simmental Cattle, *Bos taurus taurus*. *J Hered.* 2021;112:184-91.
556
- 557 22. Oppenheimer J, Rosen BD, Heaton MP, Vander Ley BL, Shafer WR, Schuetze FT, Stroud B,
558 Kuehn LA, McClure JC, Barfield JP, Blackburn HD, Kalbfleisch TS, Bickhart DM,
559 Davenport KM, Kuhn KL, Green RE, Shapiro B, Smith TPL. A Reference Genome
560 Assembly of American Bison, *Bison bison bison*. *J Hered.* 2021;112:174-183.
561
- 562 23. Arima Genomics Mapping Pipeline (2019). ArimaGenomics
563 https://github.com/ArimaGenomics/mapping_pipeline.
564
- 565 24. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin
566 R; 1000 Genome Project Data Processing Subgroup. The Sequence Alignment/Map format
567 and SAMtools. *Bioinformatics.* 2009;25:2078-9.
568
- 569 25. PicardTools (2019). Picard Toolkit, Broad Institute (Version 2.9.2)
570 <http://broadinstitute.github.io/picard>.
571
- 572 26. Quinlan AR. BEDTools: The Swiss-Army Tool for Genome Feature Analysis.
573 *Bioinformatics.* 2014;47:1-34.
574

- 575 27. Ghurye J, Pop M, Koren S, Bickhart D, Chin CS. Scaffolding of long read assemblies using
576 long range contact information. *BMC Genomics*. 2017;18:527.
577
- 578 28. Yardımcı GG, Noble W. Software tools for visualizing Hi-C data. *Genome Biol*. 2017;18:26.
579
- 580 29. Heng L. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics*.
581 2018;34:3094-100.
582
- 583 30. D-Genies (2018). D-Genies (Version 1.2.0) [https://github.com/genotoul-](https://github.com/genotoul-bioinfo/dgenies/releases/tag/v1.2.0)
584 [bioinfo/dgenies/releases/tag/v1.2.0](https://github.com/genotoul-bioinfo/dgenies/releases/tag/v1.2.0).
585
- 586 31. CombineFasta agp2fasta (2020). CombineFasta (Version 0.0.17)
587 <https://github.com/njdbickhart/CombineFasta>.
588
- 589 32. Loman N, Quick J Simpson J. A complete bacterial genome assembled de novo using only
590 nanopore sequencing data. *Nat Methods*. 2015;12:733-735.
591
- 592 33. NanoGrid (2018). NanoGrid <https://github.com/skoren/NanoGrid>.
593
- 594 34. Guan D, McCarthy SA, Wood J, Howe K, Wang Y, Durbin R. Identifying and removing
595 haplotypic duplication in primary genome assemblies. *Bioinformatics*. 2020;36:2896-8.
596

- 597 35. Smit AFA, Hubley R, Green P. RepeatMasker Open-4.0. 2013-2015.
598 <https://www.repeatmasker.org>.
599
- 600 36. Garrison E, Marth G. Haplotype-based variant detection from short-read sequencing. arXiv
601 preprint. 2012:1207.3907.
602
- 603 37. Merfin (2021). Merfin <https://github.com/arangrhie/merfin>.
604
- 605 38. Trim Galore (2020). TrimGalore (Version 0.6.6)
606 https://www.bioinformatics.babraham.ac.uk/projects/trim_galore/.
607
- 608 39. Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, Batut P, Chaisson M, &
609 Gingeras TR. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics*. 2013;29:15-21.
610
- 611 40. Li H. A statistical framework for SNP calling, mutation discovery, association mapping and
612 population genetical parameter estimation from sequencing data. *Bioinformatics*.
613 2011;27:2987-93.
614
- 615 41. Oves aries Annotation Release 104. NCBI. 2021.
616 https://ftp.ncbi.nlm.nih.gov/genomes/all/annotation_releases/9940/104
617
- 618 42. Shumate, A., and Salzberg, S.L. (2020). Liftoff: accurate mapping of gene annotations.
619 *Bioinformatics*. Doi:10.1093/bioinformatics/btaa1016.
620

- 621 43. Bray, N. L., Pimentel, H., Melsted, P., and Pachter, L. (2016). Near-optimal probabilistic
622 RNA-seq quantification. *Nat. Biotechnol.* 34, 525–527. doi:10.1038/nbt.3519.
623
- 624 44. Quinlan, A. R., and Hall, I. M. (2010). BEDTools: A flexible suite of utilities for comparing
625 genomic features. *Bioinformatics* 26, 841–842. doi:10.1093/bioinformatics/btq033.
626
- 627 45. Themis-ASM (2020). Themis-ASM pipeline <https://github.com/njdbickhart/Themis-ASM>.
628
- 629 46. Manni, M, Berkeley, M.R., Seppey, M., Simão, F.A., Zdobnov, E.M., BUSCO Update:
630 Novel and Streamlined Workflows along with Broader and Deeper Phylogenetic Coverage
631 for Scoring of Eukaryotic, Prokaryotic, and Viral Genomes, *Molecular Biology and*
632 *Evolution*, 2021;38:4647–54, <https://doi.org/10.1093/molbev/msab199>
633
- 634 47. Vezzi, F., Narzisi, G., Mishra, B. Reevaluating Assembly Evaluations with Feature Response
635 Curves: GAGE and Assemblathons. *PLoS ONE*. 2012;7:e52210.
636
- 637 48. Rhie A, Walenz BP, Koren S, Phillippy AM. Merqury: reference-free quality, completeness,
638 and phasing assessment for genome assemblies. *Genome Biol.* 2020;21:245.
639
- 640 49. Koren S, Rhie A, Walenz BP, Dilthey AT, Bickhart DM, Kingan SB, Hiendleder S, Williams
641 JL, Smith TPL, Phillippy AM. De novo assembly of haplotype-resolved genomes with trio
642 binning. *Nat Biotechnol.* 2018;10.1038/nbt.4277.
643

644 50. Davenport KM; Bickhart DM; Worley KC; Murali SC; Salavati M; Clark EL; Cockett NE;
 645 Heaton MP; Smith TPL; Murdoch BM; Rosen BD (2021): Supporting data for "An improved
 646 ovine reference genome assembly to facilitate in depth functional annotation of the sheep
 647 genome" GigaScience Database. <http://doi.org/10.5524/100944>

648

649 **Tables**

650

651 Table 1: Assembly quality statistics comparison

Assembly Statistic	ARS-UI_Ramb_v2.0	Oar_rambouillet_v1.0	Oar_v4.0	Description
Total Length (Mb)	2628.15	2869.91	2615.52	Assembly length in Mbp
Contig Number	226	7,486	48,482	Total number of contigs
Contig N50 (bp)	43,178,051	2,572,683	150,472	Half the length of the assembly is in contigs of this size or greater
Contig L50 (number of contigs)	24	313	5,008	The smallest number of contigs whose length sum make up half of the assembly size
Scaffold Number	142	2,641	5,466	Total number of scaffolds and unplaced contigs in the assembly
merQV	44.7721*	32.1705*	31.9131**	Kmer based quality from Merqury, which estimates the frequency of consensus errors in the assembly [48]
merErrorRate	0.000033327*	0.00060662*	0.000643714**	Kmer based error rate from Merqury, which estimates error rate of the assembly based on errors in kmers [48]
merCompleteness	93.0479*	93.4711*	92.2182**	Proportion of complete assembly estimated by Merqury based on "reliable" kmers, or kmers unlikely to be caused by sequencing error [48]
baseQV	41.84*	40.69*	32.40**	SNP and INDEL quality value estimated from short read data mapped to the assembly [49]
Unmap%	0.96*	1.00*	0.73**	Percentage of short reads that are unmapped to each assembly [49]
COMPLETESC	93.9	93.0	91.2	Percent of complete, single copy BUSCOs
COMPLETEDUP	2.1	2.6	1.6	Percent of complete, duplicated BUSCOs
FRAGMENT	0.9	1.1	2.4	Percent of fragmented BUSCOs
MISSING	3.1	3.3	4.8	Percent of missing BUSCOs

652

653 *Short read sequencing from the Rambouillet ewe used to assemble both ARS-UI_Ramb_v2.0
 654 and Oar_rambouillet_v1.0 was used in these quality values.

Tissue	Genome*	# input reads	# reads uniquely mapped	% of reads uniquely mapped	# reads multi-mapped	% reads multi-mapped	# reads unmapped	% reads unmapped	# indels
Skin	v2.0	62,630,134	53,990,480	86.20%	6,684,213	10.67%	1,955,441	3.12%	962
	v1.0		52,523,732	83.86%	8,114,599	12.96%	1,991,803	3.18%	2,512
	Δ	N/A	1,466,748	2.34%	-1,430,386	-2.29%	-36,362	-0.06%	-1,550
Thalamus	v2.0	54,655,873	45,721,452	83.65%	5,414,620	9.91%	3,519,801	6.44%	649
	v1.0		44,904,096	82.16%	6,126,363	11.21%	3,625,414	6.63%	1,054
	Δ	N/A	817,356	1.49%	-711,743	-1.30%	-105,613	-0.19%	-405
Pituitary	v2.0	43,368,663	39,710,031	91.56%	2,405,103	5.55%	1,253,529	2.89%	604
	v1.0		34,115,417	78.66%	7,866,251	18.14%	1,386,995	3.20%	960
	Δ	N/A	5,594,614	12.90%	-5,461,148	-12.59%	-133,466	-0.31%	-356
Lymph node – mesenteric	v2.0	43,673,576	38,819,419	88.88%	3,562,121	8.16%	1,292,036	2.96%	684
	v1.0		38,296,065	87.69%	4,057,915	9.29%	1,319,596	3.02%	999
	Δ	N/A	523,354	1.19%	-495,794	-1.13%	-27,560	-0.06%	-315
Abomasum pylorus	v2.0	45,977,534	41,018,529	89.21%	2,978,042	6.48%	1,980,963	4.31%	512
	v1.0		40,403,981	87.88%	3,533,015	7.68%	2,040,538	4.44%	846
	Δ	N/A	614,548	1.33%	-554,973	-1.20%	-59,575	-0.13%	-334

655 **Short read sequencing from the Texel animal used to assemble Oar_v4.0 was used in these
 656 quality values.

657

658

659 Table 2: RNA-seq alignment statistics to ARS-UI_Ramb_v2.0 and Oar_rambouillet_v1.0 from
 660 five different tissues.

661

662 * Genomes include v2.0 (ARS-UI_Ramb_v2.0) and v1.0 (Oar_rambouillet_v1.0) and the
663 difference (Δ).

664

665

666

667

668

669 **Figure Legends**

670

671 Figure 1: Image of Benz 2616 Rambouillet ewe selected for the ovine reference genome
672 assembly.

673

674 Figure 2: Hi-C contact map comparison of ARS-UI_Ramb_v2.0 A) directly after scaffolding and
675 before manual curation and B) after manual curation with scaffold rearrangements and joins.

676

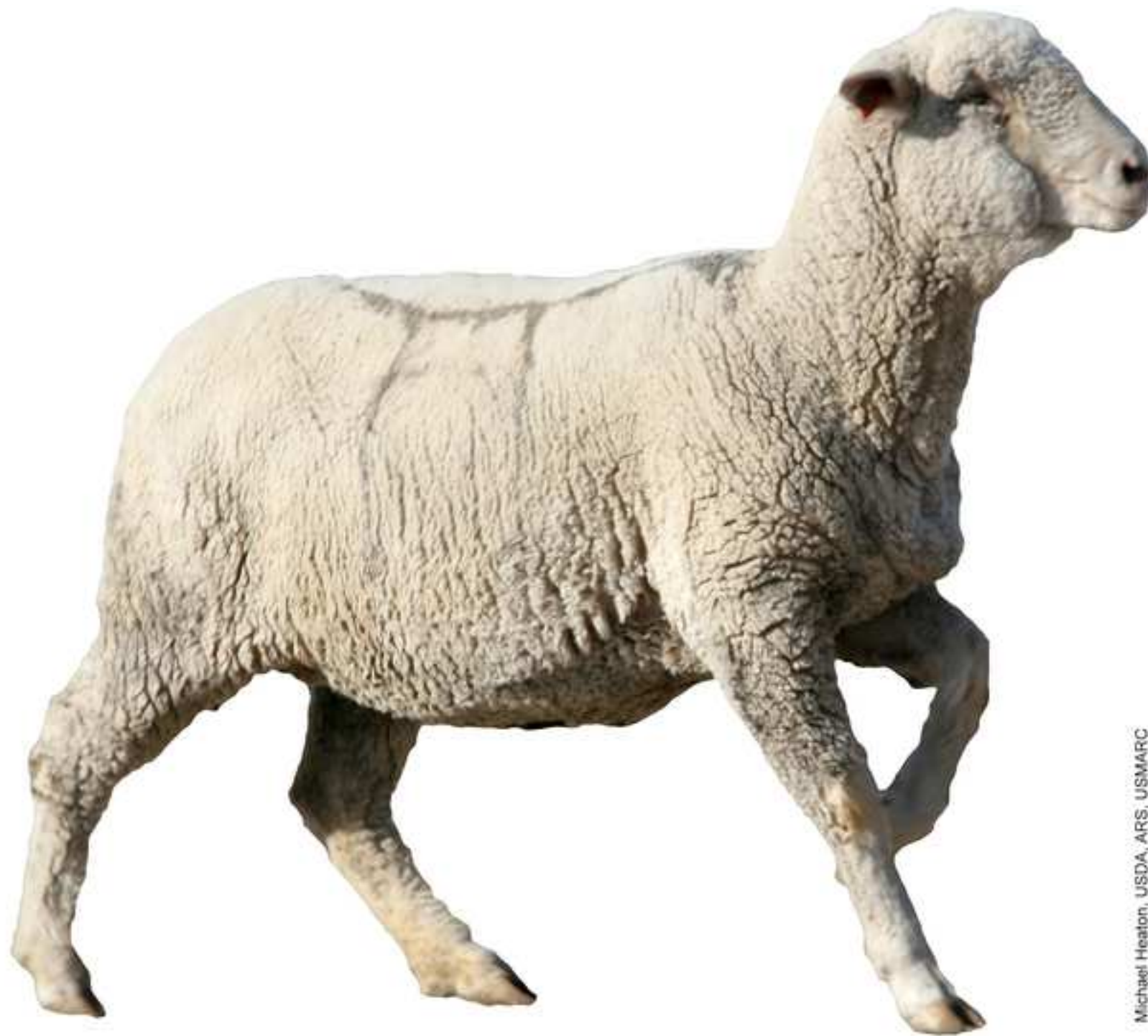
677 Figure 3: Assembly error comparison between ARS-UI_Ramb_v2.0, Oar_rambouillet_v1.0, and
678 Oar_v4.0 in A) a feature response curve displaying sorted lengths of the assemblies with the
679 fewest errors and B) specific feature counts for each genome and descriptions.

680

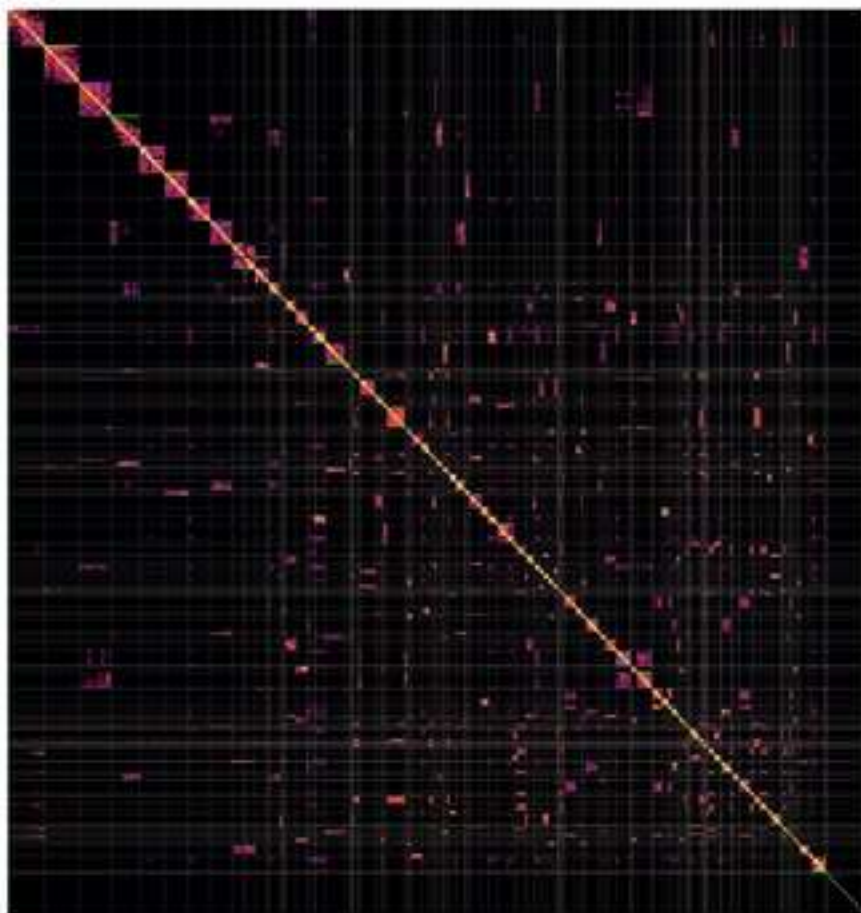
681 Figure 4: Dotplot comparison of genome assemblies between A) ARS-UI_Ramb_v2.0 and
682 Oar_rambouillet_v1.0, and B) ARS-UI_Ramb_v2.0 and Oar_v4.0.

683

684 Figure 5: Kallisto comparison of the number of expressed transcripts for the RNA-Seq dataset of
685 61 tissue samples from Benz2616, across the three annotations (Oar_Rambouillet_v1.0,
686 Ramb1LO2 (liftover) and ARS-UI_Ramb_v2.0).



Michael Heaton, USDA, ARS, USMARC

A**B**