

# GigaScience

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

<b>Manuscript Number:</b>	GIGA-D-21-00165R1														
<b>Full Title:</b>	An improved ovine reference genome assembly to facilitate in depth functional annotation of the sheep genome														
<b>Article Type:</b>	Data Note														
<b>Funding Information:</b>	<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														
<b>Abstract:</b>	<p><b>Background</b></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><b>Findings</b></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 PBSuitev15.8.24, 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><b>Conclusions</b></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>														
<b>Corresponding Author:</b>	Benjamin D Rosen  UNITED STATES														
<b>Corresponding Author Secondary</b>															

<b>Information:</b>	
<b>Corresponding Author's Institution:</b>	
<b>Corresponding Author's Secondary Institution:</b>	
<b>First Author:</b>	Kimberly M Davenport, M.S.
<b>First Author Secondary Information:</b>	
<b>Order of Authors:</b>	<p>Kimberly M Davenport, M.S.</p> <p>Derek M. Bickhart</p> <p>Kim C. Worley</p> <p>Shwetha C. Murali</p> <p>Mazdak Salavati</p> <p>Emily L. Clark</p> <p>Noelle E. Cockett</p> <p>Michael P. Heaton</p> <p>Timothy P.L. Smith</p> <p>Brenda M. Murdoch</p> <p>Benjamin D. Rosen</p>
<b>Order of Authors Secondary Information:</b>	
<b>Response to Reviewers:</b>	<p>Response to Reviewer Comments</p> <p>Reviewer #1: My comments are minimal as the paper is succinct. The authors present an improved genome, largely in the form of contiguity, and provide a number of statistics to support their argument. There are literally dozens upon dozens of different ways to assemble and polish a genome and I see no value in suggesting changes in this regard as the approach more-or-less reflects the state-of-the-art. Response: Thank you to Reviewer 1 for the thorough review of this manuscript and helpful suggestions for improvement.</p> <p>I also might question the description of "substantial improvement" as this really reflects the improvement in contiguity and less so the BUSCO, annotation. Response: The sentence in Line 76 was revised to specify the improvement in contiguity of the new genome. The sentence now reads, "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> <p>Also, scaffold L50 of the two available genomes is quite good, but not reported in Table 1, which I would suggest. The other Oar reference genomes were published 4 and 6 years ago, with this study offering the addition of nanopore sequence. Response: The scaffold L50 was added to the genome comparisons in Table 1.</p> <p>Minor L55 - long read vs long sequence? Do you mean contigs or scaffolds? Response: The sentence refers to the increasing length of raw sequence reads that can be generated and used in genome assembly. The sentence was revised to clarify this point, and now reads, "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>L233- does freebayes do polishing? This is what is suggested by the current wording</p>

	<p>Response: Freebayes was used in the polishing pipeline, we mention in the abstract that BCFtools is then used for consensus generation but failed to mention BCFtools in the methods section. This has been clarified “The final polishing with Illumina short read data consisted of two rounds of freebayes v1.3.1 [35] variant calling and BCFtools [36] consensus. Variants used for polishing with both Nanopolish and freebayes/BCFtools were screened with Merfin [37] which evaluates the k-mer consequences of variant calls and filters unsupported variants.”</p> <p>Reviewer #2: This paper is of high importance and will be well cited as people use the resource it is reporting on. It is well written and easy to read. Response: Thank you to Reviewer 2 for the thorough review of this manuscript and helpful comments.</p> <p>A mention of the genome length and NG50 (in addition to N50) in the abstract would be useful. Response: The genome length and NG50 were added to the abstract in Line 68 and reads, “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.” The manuscript has also been modified to refer to NG50 rather than N50 throughout.</p> <p>The authors might like to add the species name as a keyword (completely optional) Response: Species name (<i>Ovis aries</i>) has been added as a keyword.</p> <p>If not already done, the mito genome should have the start position matched to the older version. If this is already done then add it to the methods. Response: Yes, the start position was matched to the Oar_rambouillet_v1.0 mitochondrial genome. This was added to the methods in Line 233. The sentence reads, “The mitochondrial genome was identified by aligning the previously annotated mitochondrial sequence from Oar_rambouillet_v1.0 (RefSeq NC_001941.1) to the assembly contigs and the start positions were matched.”</p> <p>I think Figure 3B should be a table not a figure part. Also Figure 5A. Response: Figure 3B and Figure 5A were both separated and added as tables.</p>
<b>Additional Information:</b>	
<b>Question</b>	<b>Response</b>
Are you submitting this manuscript to a special series or article collection?	No
<p><b>Experimental design and statistics</b></p> <p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>. Information essential to interpreting the data presented should be made available in the figure legends.</p> <p>Have you included all the information</p>	Yes

<p>requested in your manuscript?</p>	
<p><b>Resources</b></p> <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 <a href="#">Research Resource Identifiers</a> (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>?</p>	<p>Yes</p>
<p><b>Availability of data and materials</b></p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in <a href="#">publicly available repositories</a> (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 <a href="#">Minimum Standards Reporting Checklist</a>?</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. Davenport<sup>1</sup>, Derek M. Bickhart<sup>2</sup>, Kim C. Worley<sup>3</sup>, Shwetha C. Murali<sup>4</sup>, Mazdak  
5 Salavati<sup>5</sup>, Emily L. Clark<sup>6</sup>, Noelle E. Cockett<sup>7</sup>, Michael P. Heaton<sup>8</sup>, Timothy P.L. Smith<sup>9</sup>, Brenda  
6 M. Murdoch<sup>10\*</sup>, and Benjamin D. Rosen<sup>11\*</sup>

7  
8 <sup>1</sup>Department of Animal, Veterinary, and Food Sciences, University of Idaho, 875 Perimeter Dr.,  
9 Moscow, ID, United States 83843. Email: [kmdavenport@uidaho.edu](mailto:kmdavenport@uidaho.edu)

10  
11 <sup>2</sup>US Dairy Forage Research Center, USDA-ARS, 1925 Linden Drive, Madison, WI, United  
12 States 53706. Email: [derek.bickhart@usda.gov](mailto:derek.bickhart@usda.gov)

13  
14 <sup>3</sup>Baylor College of Medicine, One Baylor Plaza, Houston, TX, United States 77030. Email:  
15 [kworley@bcm.edu](mailto:kworley@bcm.edu)

16  
17 <sup>4</sup>Baylor College of Medicine, One Baylor Plaza, Houston, TX, United States 77030.  
18 Email: [shwethac@gmail.com](mailto:shwethac@gmail.com)

19  
20 <sup>5</sup>The Roslin Institute, Royal (Dick) School of Veterinary Studies, The University of Edinburgh,  
21 Easter Bush Campus, Midlothian, United Kingdom, EH25 9RG, United Kingdom. Email:  
22 [mazdak.salavati@roslin.ed.ac.uk](mailto:mazdak.salavati@roslin.ed.ac.uk)

23  
24 <sup>6</sup>The Roslin Institute, Royal (Dick) School of Veterinary Studies, The University of Edinburgh,  
25 Easter Bush Campus, Midlothian, United Kingdom, EH25 9RG. Email:  
26 [emily.clark@roslin.ed.ac.uk](mailto:emily.clark@roslin.ed.ac.uk)

27  
28 <sup>7</sup>Utah State University, Old Main Hill, Logan, UT 84322. Email: [noelle.cockett@usu.edu](mailto:noelle.cockett@usu.edu)

29  
30 <sup>8</sup>US Meat Animal Research Center, USDA-ARS, State Spur 18D, Clay Center, NE 68933.  
31 Email: [mike.heaton@usda.gov](mailto:mike.heaton@usda.gov)

32  
33 <sup>9</sup>US Meat Animal Research Center, USDA-ARS, State Spur 18D, Clay Center, NE 68933.  
34 Email: [tim.smith2@usda.gov](mailto:tim.smith2@usda.gov)

35  
36 <sup>10</sup>Department of Animal, Veterinary, and Food Sciences, University of Idaho, 875 Perimeter Dr.,  
37 Moscow, ID 83843. Email: [bmurdoch@uidaho.edu](mailto:bmurdoch@uidaho.edu)

38  
39 <sup>11</sup>Animal Genomics and Improvement Laboratory, USDA-ARS, 10300 Baltimore Avenue,  
40 Beltsville, MD 20705. Email: [ben.rosen@usda.gov](mailto:ben.rosen@usda.gov)

41  
42 Correspondence:

43 Brenda M. Murdoch

44 [bmurdoch@uidaho.edu](mailto:bmurdoch@uidaho.edu)

45 Benjamin D. Rosen

46 [ben.rosen@usda.gov](mailto:ben.rosen@usda.gov)

47 **Abstract**

48

49 *Background*

50

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

58

59

60 *Findings*

61

62 Short-read Illumina (55x coverage), long-read PacBio (75x coverage), and Hi-C data from this  
63 ewe retrieved from public databases were combined with an additional 50x coverage of Oxford  
64 Nanopore data and assembled with canu v1.9. The assembled contigs were scaffolded using Hi-  
65 C data with Salsa v2.2, gaps filled with Pbsuite v15.8.24, and polished with Nanopolish v0.12.5.  
66 After duplicate contig removal with PurgeDups v1.0.1, chromosomes were oriented and polished  
67 with two rounds of a pipeline which consisted of freebayes v1.3.1 to call variants, Merfin to  
68 validate them, and BCFtools to generate the consensus fasta. The ARS-UI\_Ramb\_v2.0 assembly  
69 is 2.63 Gb in length and has improved continuity (contig NG50 of 43.18 Mb) with a 19-fold and

70 38-fold decrease in the number of scaffolds compared with Oar\_rambouillet\_v1.0 and Oar\_v4.0.  
71 ARS-UI\_Ramb\_v2.0 has greater per-base accuracy and fewer insertions and deletions identified  
72 from mapped RNA sequence than previous assemblies.

73

74

75 *Conclusions*

76

77 The ARS-UI\_Ramb\_v2.0 assembly is a substantial improvement in contiguity that will optimize  
78 the functional annotation of the sheep genome and facilitate improved mapping accuracy of  
79 genetic variant and expression data for traits in sheep.

80

81

82 **Keywords:** Rambouillet, genome assembly, reference genome, sheep, *Ovis aries*

83

84

85

86

87

88

89

90

91

92

93

94 **Context**

95

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

105

106 Genome research in sheep holds promise to improve efficiency and sustainability of production  
107 and reduce the environmental effects of animal agriculture [12]. The first sheep reference  
108 genome assembly was based on whole genome shotgun (WGS) short-read sequencing,  
109 scaffolded by genetic linkage and radiation hybrid maps. The sequence came from two unrelated  
110 Texel breed sheep, with the first assembly draft (Oar\_v3.1; International Sheep Genomics  
111 Consortium, 2010) having a contig NG50, based on a 2.6 gigabase (Gb) genome size, of 39  
112 kilobases (kb) and the update (Oar\_v4.0) [13] boosting the NG50 metric to 145 kb. More  
113 recently, the Ovine Functional Annotation of Animal Genomes (FAANG) project proposed to  
114 perform a variety of genome annotation assays for dozens of tissues from a single animal  
115 [14,15]. To maximize the success of assays that depend on mapping sequence data to a reference,



116 the FAANG project assembled the genome of that animal, a female of the Rambouillet breed.  
117 The assembly, released in 2017 (Oar\_rambouillet\_v1.0, GenBank accession GCF\_002742125;  
118 Worley et al., unpublished) is based on a combination of Pacific Biosciences RSII WGS long-  
119 read and Illumina short-read sequencing. It has an improved contig NG50 of 2.9 megabases (Mb)  
120 and is generally regarded as the official reference assembly for global sheep research.

121

122 The continued maturation of long read sequencing technologies provided an opportunity to  
123 improve upon the sheep reference genome assembly. Since most of the proposed FAANG  
124 annotation assays had already been performed on the Rambouillet ewe, lung tissue from the  
125 same animal was chosen for DNA extraction. This allowed the use of existing long read data to  
126 supplement new, longer-read, Oxford Nanopore PromethION sequencing. We report a *de novo*  
127 assembly of the same Rambouillet ewe used for Oar\_rambouillet\_v1.0, based on approximately  
128 50x coverage of nanopore reads (N50 47kb) and 75x coverage PacBio reads (N50 13kb). The  
129 new assembly, ARS-UI\_Ramb\_v2.0 offers a 15-fold improvement in contiguity and increased  
130 accuracy, providing a basis for regulatory element annotation in the FAANG project and  
131 facilitating the discovery of biological mechanisms that influence traits important in global sheep  
132 research and production.

133

134

## 135 **Methods**

136

### 137 *Sampling Strategy*

138

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

146

147

#### 148 *Sequencing and Data Acquisition*

149

150 DNA was extracted from approximately 50 mg of lung tissue using phenol:chloroform-based  
151 method as described (Logsdon 2019). Briefly, the frozen tissue was pulverized in a cryoPREP  
152 CP02 tissue disruption system (Covaris Inc., Woburn MA) as recommended by the  
153 manufacturer. The powdered tissue was transferred to a 50 mL conical tube and mixed in 200  
154  $\mu$ L of phosphate buffered saline (Sigma-Aldrich, St. Louis MO). The tissue was then diluted in  
155 10 mL of buffer TLB (100mM NaCl, 10mM Tris-HCl pH 8.0, 25 mM EDTA, 0.5% SDS) and  
156 mixed by vortexing, then incubated with 20  $\mu$ L 10 mg/mL RNase A at 37°C for one hour with  
157 gentle shaking. Protein digestion was performed with 100  $\mu$ L Proteinase K (20 mg/mL) at 50°C  
158 for 2 hours, with slow rotation of the tube to mix every 30 minutes. The lysate was distributed  
159 equally into two 15 mL Phase Lock tubes (Quantabio, Beverly MA) and each tube received 5  
160 mL of TE-saturated Phenol (Sigma-Aldrich, St. Louis MO) followed by mixing on a tube rotator  
161 at 20 RPM for 10 minutes at 22°C. The aqueous layer was collected after separating at 2300xg

162 for 10 minutes and transferred to another Phase Lock tube. A second extraction performed in the  
163 same way as the first was conducted using 2.5 mL phenol and 2.5 mL chloroform:isoamyl  
164 alcohol (Sigma). The final aqueous phase was transferred to a 50 mL conical tube and the DNA  
165 precipitated with 2 mL of 5M ammonium acetate and 15 mL of ice-cold 100% ethanol. The  
166 DNA was pulled from the alcohol using a Pasteur pipet “hook” and placed in 10 mL of cold 70%  
167 ethanol to wash the pellet. The ethanol was poured off and the DNA pellet dried for 20-30  
168 minutes, then dissolved in a dark drawer at room temperature for 48 hours in 1 mL of 10mM  
169 Tris-Cl pH 8.5. Library preparation for Oxford Nanopore long read sequencing was performed  
170 with an LSK-109 template preparation kit as recommended by the manufacturer (Oxford  
171 Nanopore, Oxford U.K.) with modifications as described by Logsdon  
172 ([https://www.protocols.io/view/hmw-gdna-purification-and-ont-ultra-long-read-data-  
173 bchhit36?comment\\_id=88927](https://www.protocols.io/view/hmw-gdna-purification-and-ont-ultra-long-read-data-bchhit36?comment_id=88927)). The ligated template was sequenced with a PromethION  
174 instrument using four R9.4 flow cells. (Oxford Nanopore Technologies, Oxford, United  
175 Kingdom). Output as fast5 files were basecalled with Guppy v3.1 [18].

176  
177 Sequence data used in the previous Oar\_rambouillet\_v1.0 assembly was retrieved from the  
178 Sequence Read Archive listed under project number PRJNA414087 [15]. PacBio RS II sequence  
179 generated from DNA extracted from whole blood was retrieved from SRX3445660,  
180 SRX3445661, SRX3445662, and SRX3445663. The Hi-C sequence data generated from liver  
181 using HindIII enzyme and sequenced at 150 bp paired end with an Illumina HiSeq X Ten was  
182 retrieved from SRX3399085 and SRX3399086. Short read whole genome sequencing from DNA  
183 extracted from whole blood collected from the Rambouillet ewe was performed with an Illumina  
184 HiSeq X Ten sequenced at 150 bp paired end and was retrieved from SRX3405602. Further

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

189

190

### 191 *Assembly*

192

193 Contigs were assembled with Oxford Nanopore and PacBio reads generated as described above  
194 using canu v1.8 through the trimmed reads stage of assembly. Parameters for contig construction  
195 were set as “batOptions=-dg 4 -db 4 -mo 1000” [19]. Canu v1.9 was used to complete the contig  
196 assembly because this update demonstrates better consensus generation of the overlapped contigs  
197 in the final step in the assembly process [20,21]. The corrected error rate option was set as  
198 “correctedErrorRate=0.105.”

199

200

### 201 *Scaffolding*

202

203 Two Hi-C datasets from liver tissue from two different library preparations were retrieved as  
204 described above. The Hi-C reads were first aligned to the polished contigs using the Arima  
205 Genomics mapping pipeline [22]. This pipeline first maps paired end reads individually with  
206 bwa-mem, then removes the 3' end of reads identified as chimeric and span ligation junctions.  
207 Reads were then paired, filtered by mapping quality with samtools [23], and PCR duplicates

208 removed with Picard [24]. The two Hi-C libraries were merged in the final step in the Arima  
209 pipeline to generate the merged BAM file. The BAM file was converted to a BED file for input  
210 into Salsa using the bedtools command `bamToBed` [25]. Salsa v2.2 was used for scaffolding by  
211 implementing “`python run_pipeline.py -a contigs.fasta -l contigs.fasta.fai -b alignment.bed -e`  
212 `HindIII -o scaffolds -m yes`” [26].

213  
214 The Hi-C reads were aligned to the scaffolded assembly with the Arima Genomics mapping  
215 pipeline and then processed with PretextView to visually evaluate the scaffolds as a contact map  
216 in PretextView [27]. The scaffolded assembly was also compared to *Oar\_rambouillet\_v1.0* by  
217 aligning the two genomes with “`minimap2 -cx asm5 Oar_rambouillet_v1.0_genomic.fasta`  
218 `scaffolds.fasta > alignment.paf`” [28]. A dotplot of the alignment was visualized with D-Genies  
219 [29]. Scaffolds were edited based on visual inspection of the contact map and dotplot, as well as  
220 the Hi-C alignment file. Scaffold joins and rearrangements were incorporated to the assembly  
221 using the *agp2fasta* mode of CombineFasta [30].

222

223

#### 224 *Gap Filling and Polishing*

225

226 Gap filling was completed with pbsuite v15.8.24 using both the PacBio and Oxford Nanopore  
227 reads. Nanopolish v0.12.5 [31] with the NanoGrid parallel wrapper [32] was employed with the  
228 raw fast5 files generated from the PromethION sequencing to polish the assembly. Duplicates  
229 were removed using PurgeDups v1.0.1 [33]. The chromosome orientation was confirmed in the  
230 polished assembly by identifying telomeres and centromeres using RepeatMasker v4.1.1 [34].

231 The mitochondrial genome was identified by aligning the previously annotated mitochondrial  
232 sequence from Oar\_rambouillet\_v1.0 (RefSeq NC\_001941.1) to the assembly contigs and the  
233 start positions were matched. Chromosomes were oriented centromere to telomere and placed in  
234 chromosome number order. The final polishing with Illumina short read data consisted of two  
235 rounds of freebayes v1.3.1 [35] variant calling and BCFtools [36] consensus. Variants used for  
236 polishing with both Nanopolish and freebayes/BCFtools were screened with Merfin [37] which  
237 evaluates the k-mer consequences of variant calls and filters unsupported variants.

238

239

#### 240 *RNA Sequencing*

241

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

251

252

#### 253 *Annotation*

254 The annotation for ARS-UI\_Ramb\_v2.0, NCBI Ovis aries Annotation Release 104, is available  
255 in RefSeq and other NCBI genome resources  
256 ([https://ftp.ncbi.nlm.nih.gov/genomes/all/annotation\\_releases/9940/104](https://ftp.ncbi.nlm.nih.gov/genomes/all/annotation_releases/9940/104)).

257  
258 Here we also provide a liftover of the annotation for Oar\_rambouillet\_v1.0 onto ARS-  
259 UI\_Ramb\_v2.0. The annotation used for the liftover was NCBI v103  
260 GCF\_002742125.1\_Oar\_rambouillet\_v1.0\_genomic.fna.gz. The GFF3 format gene annotation  
261 file was prepared for processing using liftOff v1.5.2 [41]. A set of matching chromosome names  
262 for Oar\_rambouillet\_v1.0 and ARS-UI\_Ramb\_v2.0 were generated according to the instructions  
263 for liftOff (*paste -d "," <(cut -d ' ' -f1 ramb1.chr) <(cut -d ' ' -f1 ramb2.chr) > chroms.txt*). The  
264 GFF file (annotation Ramb1LO2) generated by liftOff is included in Supplementary File 1  
265 (Ramb\_v1.0\_NCBI103\_lifted\_over\_ARS-UI\_Ramb\_v2.0.gff.gz).

266  
267 To compare the breakdown of transcripts captured by the three annotations  
268 (Oar\_Rambouillet\_v1.0, Ramb1LO2 (liftover) and ARS-UI\_Ramb\_v2.0), we generated  
269 transcript expression estimates using Kallisto v0.44.0 [42]. For the lifted over gene annotation  
270 the GFF file (Ramb\_v1.0\_NCBI103\_lifted\_over\_ARS-UI\_Ramb\_v2.0.gff.gz) was used to  
271 generate transcriptome sequence FASTA files, as a Kallisto index, for transcript expression  
272 estimation. Briefly, exonic blocks were extracted from the GFF3 file using the awk command  
273 (*awk '(\$3~/exon/)' input.gff*). The getfasta and groupby plugins from bedtools v2.30.0 [43] were  
274 used to extract the exonic sequences and group them by transcript name. Exonic sequences for  
275 each transcript were appended in the correct order, to produce the complete sequence for each  
276 transcript. The FASTA format file for the whole transcriptome was created using all of the

277 transcript level FASTA sequences for the liftover annotation Ramb1LO2 (Supplementary File 2;  
278 Ramb1LO2\_NCBI103\_geneBank\_rna.fa). The set of scripts used for this step are included in  
279 Supplementary File 3. The Kallisto indices for Oar\_Rambouillet\_v1.0  
280 (GCF\_002742125.1\_Oar\_rambouillet\_v1.0\_rna.fna.gz), Ramb1LO2 (liftover;  
281 Ramb1LO2\_NCBI103\_geneBank\_rna.fa) and ARS-UI\_Ramb\_v2.0 (GCF\_016772045.1\_ARS-  
282 UI\_Ramb\_v2.0\_rna.fna.gz) were then used with the RNA-Seq data from the 61 tissues from  
283 Benz2616 (GenBank BioProject PRJNA414087 and PRJEB35292) to estimate transcript level  
284 expression for every tissue as transcript per million mapped reads (TPM) and compared across  
285 the three annotations.

286

287

## 288 **Data Validation and Quality Control**

289

### 290 *Assembly Quality Statistics*

291

292 The four flow cells of PromethION data produced 136 Gb of WGS sequence (approximately 51x  
293 coverage) in reads having a read N50 of 47 kb. The initial generation of contigs used this data as  
294 well as 198.1 Gb of RSII data with a read N50 of 12.9 kb. The ARS-UI\_Ramb\_v2.0 assembly  
295 was submitted to NCBI GenBank under accession number GCA\_016772045.1, and statistics of  
296 contigs and scaffolds following initial polishing, scaffolding with Hi-C data and manual editing,  
297 gap-filling, and final polishing, are shown in Table 1. The assembly improved on the  
298 Oar\_v4.0/Oar\_rambouillet\_v1.0 sheep reference assemblies in all continuity measures (Table 1)  
299 including a 286/17-fold increase in contig N50 (the size of the shortest contig for which all larger



300 contigs contain half of the total assembly), a 214/33-fold reduction in the number of contigs in  
301 the assembly and concomitant 209/13-fold reduction of contig L50 (the number of contigs  
302 making up half of the total assembly), and 38/19-fold reduction in total number of scaffolds.  
303 Manual curation of scaffolds using Hi-C data improved scaffold continuity and led to  
304 chromosome length scaffolds (Figure 2).

305  
306 The Themis-ASM pipeline [44] was implemented to further assess assembly quality and  
307 compare sheep genome assemblies. Short read sequence from both the Rambouillet ewe used in  
308 this assembly and Texel sheep from previous sheep genome assemblies were used to compare  
309 ARS-UI\_Ramb\_v2.0 with Oar\_rambouillet\_v1.0 and Oar\_v4.0 assemblies.

310  
311 The k-mer based quality value and error rates improved with ARS-UI\_Ramb\_v2.0 compared  
312 with Oar\_rambouillet\_v1.0 and Oar\_v4.0. This is also reflected in the proportion of complete  
313 assembly based on k-mers (merCompleteness), which is similar between ARS-UI\_Ramb\_v2.0  
314 and Oar\_rambouillet\_v1.0 and both are higher than Oar\_v4.0. Further, the SNP and indel quality  
315 value (baseQV) were greatest overall in ARS-UI\_Ramb\_v2.0 (41.84), followed by  
316 Oar\_rambouillet\_v1.0 (40.69) and Oar\_v4.0 (32.40). The percentage of short reads not mapped  
317 to the genome was  $\leq 1\%$  in all three assemblies.

318  
319 The completeness of ARS-UI\_Ramb\_v2.0 was evaluated by examining the presence or absence  
320 of evolutionarily conserved genes in each assembly using Benchmarking Universal Single-Copy  
321 Ortholog (BUSCO) scores generated as an output of the Themis-ASM pipeline. The percent of  
322 single copy complete BUSCOs were higher (90.7%) in ARS-UI\_Ramb\_v2.0 when compared

323 with Oar\_rambouillet\_v1.0 (90.1%) and Oar\_v4.0 (86.1%). Complete duplicated BUSCO  
324 percentage was highest in Oar\_rambouillet\_v1.0 (1.6%) compared with ARS-UI\_Ramb\_v2.0  
325 (1.4%), and lowest in Oar\_v4.0 (1.0%). Further, ARS-UI\_Ramb\_v2.0 had the lowest percent of  
326 fragmented and missing BUSCOs (2.0% and 5.9%, respectively) compared with  
327 Oar\_rambouillet\_v1.0 (2.1% and 6.2%, respectively) and Oar\_v4.0 (3.7% and 9.2%,  
328 respectively).

329  
330 The three sheep genome assemblies were also compared with a feature response curve in which  
331 the quality of the assembly is analyzed as a function of the features, or maximum number of  
332 possible errors, allowed in the contigs (Figure 3) [45]. Both the ARS-UI\_Ramb\_v2.0 and  
333 Oar\_v4.0 feature response curves peak higher and to the left of Oar\_rambouillet\_v1.0, which  
334 indicate fewer errors in these assemblies (Figure 3A). The ARS-UI\_Ramb\_v2.0 genome also has  
335 fewer regions with either low or high coverage overall and for paired reads, suggesting fewer  
336 coverage issues, as well as fewer improperly paired or unmapped single reads when compared  
337 with other assemblies (Table 2). The number of high Comp/Expansion (CE) statistics in ARS-  
338 UI\_Ramb\_v2.0 was intermediate between Oar\_rambouillet\_v1.0 (higher) and Oar\_v4.0 (lower),  
339 however this latest assembly had the lowest number of regions with low CE statistics.

340  
341 Comparative alignment of ARS-UI\_Ramb\_v2.0 with previous assemblies Oar\_rambouillet\_v1.0  
342 and Oar\_v4.0 and visualization with a dotplot revealed a high amount of agreement between  
343 assemblies (Figure 4). Interestingly, chromosome 11 was improperly oriented in  
344 Oar\_rambouillet\_v1.0, and after confirming centromere and telomere locations on this  
345 chromosome, this was resolved in the ARS-UI\_Ramb\_v2.0 assembly. The percent identity

346 between ARS-UI\_Ramb\_v2.0 is very high when compared with Oar\_rambouillet\_v1.0 which  
347 was expected considering the same animal was used in both assemblies. However, Oar\_v4.0 was  
348 assembled from Texel sheep, which is apparent in the percent identity in the dotplot.

349

350 In summary, ARS-UI\_Ramb\_v2.0 offers greater contiguity, improved quality, more complete  
351 BUSCOs, and fewer assembly errors when compared with previous assemblies.

352

353

#### 354 *RNA sequencing alignment*

355

356 Insertions and deletions (indels) in the ARS-UI\_Ramb\_v2.0 assembly were characterized and  
357 compared with Oar\_rambouillet\_v1.0 by mapping 150 bp paired-end RNA-seq data from skin,  
358 thalamus, pituitary, lymph node (mesenteric), and abomasum pylorus generated from the same  
359 animal used to assemble the reference genome (Table 3). In all five tissues, ARS-UI\_Ramb\_v2.0  
360 had nearly half of the number of indels compared with Oar\_rambouillet\_v1.0. Most indels  
361 identified in both assemblies were 1bp in length. The ARS-UI\_Ramb\_v2.0 had a greater number  
362 of uniquely mapped reads in each tissue when compared with Oar\_rambouillet\_v1.0, leading to  
363 an approximate 2% increase in the percent of uniquely mapped reads in most tissues except  
364 pituitary, which saw an almost 13% improvement. The number of reads that mapped to multiple  
365 loci decreased in the new assembly by 12.59% in pituitary, and 1-2% in other tissues. Further,  
366 ARS-UI\_Ramb\_v2.0 had fewer unmapped reads than Oar\_rambouillet\_v1.0 across all five  
367 tissues by an average of 0.15%.

368

369 *Annotation*

370

371 The ARS-UI\_Ramb\_v2.0 annotation represents a substantial improvement over the annotation  
372 on Oar\_rambouillet\_v1.0. For example, for ARS-UI\_Ramb\_v2.0 16,500 coding genes have an  
373 ortholog to human (compared to 16,319 for Oar\_rambouillet\_v1.0), and the BUSCO scores  
374 demonstrate that 99.1% of the gene models (cetartiodactyla\_odb10) are complete in the new  
375 annotation versus 98.8% in the previous one. The annotation for ARS-UI\_Ramb\_v2.0 includes  
376 Iso-Sequencing for 8 tissues to improve contiguity of gene models, and CAGE sequencing for 56  
377 tissues to define TSS, that were not used to annotate Oar\_rambouillet\_v1.0. The full report for  
378 the annotation release is available at:

379 ([https://www.ncbi.nlm.nih.gov/genome/annotation\\_euk/Ovis\\_aries/104](https://www.ncbi.nlm.nih.gov/genome/annotation_euk/Ovis_aries/104)).

380

381 Using Kallisto we compared the number of expressed transcripts, for the RNA-Seq dataset of 61  
382 tissue samples from Benz2616, across the three annotations (Oar\_Rambouillet\_v1.0, Ramb1LO2  
383 (liftover) and ARS-UI\_Ramb\_v2.0). There was a considerable increase in the number of  
384 transcripts captured by the annotation for ARS-UI\_Ramb\_v2.0 (60,064) relative to  
385 Oar\_Rambouillet\_v1.0 (42,058) and the liftover annotation (Ramb1LO2) (40,910) (Table 4 and  
386 Figure 5). This equates to approximately 20,000 new annotated gene models for ARS-  
387 UI\_Ramb\_v2.0 and further reflects the substantial improvement over the annotation for  
388 Oar\_Rambouillet\_v1.0.

389 The lifted over annotation we have generated will provide a resource for those who wish to  
390 compare their results for ARS-UI\_Ramb\_v2.0 to previous work using

391 Oar\_Rambouillet\_v1.0. Only 2.7% of protein coding transcripts were lost (1148) lifting over the

392 annotation for Oar\_Rambouillet\_v1.0 onto ARS-UI\_Ramb\_v2.0. According to the annotation  
393 report provided by NCBI  
394 ([https://www.ncbi.nlm.nih.gov/genome/annotation\\_euk/Ovis\\_aries/104/](https://www.ncbi.nlm.nih.gov/genome/annotation_euk/Ovis_aries/104/)), 70% of the annotations  
395 were identical or had only minor changes between and Oar\_Rambouillet\_v1.0 and ARS-  
396 UI\_Ramb\_v2.0.

397

398

### 399 **Re-use potential**

400

401 The ARS-UI\_Ramb\_v2.0 genome assembly serves as a reference for genetic investigation of  
402 traits important in sheep research and production across the world. This genome is assembled  
403 from the same animal used in the Ovine FAANG Project, which provides a high-quality basis for  
404 epigenetic annotation to serve the international sheep genomics community and scientific  
405 community at large.

406

407

### 408 **Availability of supporting data**

409

410 The data sets supporting the results of this article are available in the GenBank repository,  
411 GCA\_016772045.1.

412

413

### 414 **Additional files**

415 Supplementary File 1 – Ramb\_v1.0\_NCBI103\_lifted\_over\_ARS-UI\_Ramb\_v2.0.gff.gz

416 Supplementary File 2 – Ramb1LO2\_NCBI103\_geneBank\_rna.fa

417 Supplementary File 3 – Supplementary\_File\_3\_scripts.txt

418

419

## 420 **Author contributions**

421

422 BMM, TPLS, DMB, and BDR conceptualized the study. BMM, NEC, MPH, and TPLS selected  
423 the animal and collected samples. KW and SCM facilitated the generation of RSII, short read,  
424 and Hi-C data. TPLS facilitated the nanopore long read data generation. KMD, DMB, TPLS,  
425 BMM, and BDR performed the genome assembly, scaffolding, RNA-sequencing alignment,  
426 polishing, and quality control. MS and ELC contributed the section describing the LiftOff  
427 annotation and comparative analysis of transcript expression estimates for the three annotations.  
428 KMD, DMB, TPLS, BMM, and BDR generated tables and figures and drafted the manuscript.  
429 KMD, DMB, KW, SCM, NEC, TPLS, BMM, and BDR edited the manuscript. All authors  
430 contributed to the article and approved the final version.

431

432

## 433 **Acknowledgements**

434

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

437

438

**439 Funding**

440

441 Funding was provided by Agriculture and Food Research Initiative Competitive grants from the  
442 USDA National Institute of Food and Agriculture supporting improvements of the sheep  
443 genomes (2013-67015-21228) and FAANG activities (2013-67015-21372, 2017-67016-26301).  
444 Additional funding was received from the International Sheep Genome Consortium  
445 (217201191442) and infrastructure support from a grant to R. Gibbs from the NIH NHGRI  
446 Large-Scale Sequencing Program (U54 HG003273).

447

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

453

454

**455 References**

456

- 457 1. Pedrosa S, Uzun M, Arranz JJ, Gutiérrez-Gil B, San Primitivo F, Bayón Y. Evidence of three  
458 maternal lineages in Near Eastern sheep supporting multiple domestication events. *Proc Biol*  
459 *Sci.* 2005;272:2211-7.

460

- 461 2. Zeder MA. Domestication and early agriculture in the Mediterranean Basin: origins,  
462 diffusion, and impact. *Proc Natl Acad Sci USA*. 2008;105:11597-604.  
463
- 464 3. Chessa B, Pereira F, Arnaud F, Amorim A, Goyache F, Mainland I, Kao RR, Pemberton JM,  
465 Beraldi D, Stear MJ, Alberti A, Pittau M, Iannuzzi L, Banabazi MH, Kazwala RR, Zhang  
466 YP, Arranz JJ, Ali BA, Wang Z, Uzun M, Dione MM, Olsaker I, Holm LE, Saarma U,  
467 Ahmad S, Marzanov N, Eythorsdottir E, Holland MJ, Ajmone-Marsan P, Bruford MW,  
468 Kantanen J, Spencer TE, Palmarini M. Revealing the history of sheep domestication using  
469 retrovirus integrations. *Science*. 2009;324:532-6.  
470
- 471 4. Kijas JW, Lenstra JA, Hayes B, Boitard S, Porto Neto LR, San Cristobal M, Servin B,  
472 McCulloch R, Whan V, Gietzen K, Paiva S, Barendse W, Ciani E, Raadsma H, McEwan J,  
473 Dalrymple B, International Sheep Genomics Consortium Members. Genome-wide analysis  
474 of the world's sheep breeds reveals high levels of historic mixture and strong recent selection.  
475 *PLoS Biol*. 2012;doi:10.1371/journal.pbio.1001258.  
476
- 477 5. Burke JM, Miller JE. Relative resistance of Dorper crossbred ewes to gastrointestinal  
478 nematode infection compared with St. Croix and Katahdin ewes in the southeastern United  
479 States. *Vet Parasitol*. 2002;109:265-75.  
480
- 481 6. Bowdridge SA, Zajac AM, Notter DR. St. Croix sheep produce a rapid and greater cellular  
482 immune response contributing to reduced establishment of *Haemonchus contortus*. *Vet*  
483 *Parasitol*. 2015;208:204-10.



- 484
- 485 7. Scherf BD. World watch list for domestic animal diversity. 3<sup>rd</sup> ed. Rome: Food and  
486 Agriculture Organization of the United Nations; 2000.
- 487
- 488 8. Lv FH, Agha S, Kantanen J, Colli L, Stucki S, Kijas JW, Joost S, Li MH, Ajmone Marsan P.  
489 Adaptations to climate-mediated selective pressures in sheep. *Mol Biol Evol.* 2014;31:3324-  
490 43.
- 491
- 492 9. Cao YH, Xu SS, Shen M, Chen ZH, Gao L, Lv FH, Xie XL, Wang XH, Yang H, Liu CB,  
493 Zhou P, Wan PC, Zhang YS, Yang JQ, Pi WH, Hehua E, Berry DP, Barbato M,  
494 Esmailizadeh A, Nosrati M, Salehian-Dehkordi H, Dehghani-Qanatqestani M, Dotsev AV,  
495 Deniskova TE, Zinovieva NA, Brem G, Štěpánek O, Ciani E, Weimann C, Erhardt G,  
496 Mwacharo JM, Ahbara A, Han JL, Hanotte O, Miller JM, Sim Z, Coltman D, Kantanen J,  
497 Bruford MW, Lenstra JA, Kijas J, Li MH. Historical Introgression from Wild Relatives  
498 Enhanced Climatic Adaptation and Resistance to Pneumonia in Sheep. *Mol Biol Evol.*  
499 2021;38:838-55.
- 500
- 501 10. Dickinson WF, Lush JL. Inbreeding and the genetic history of the Rambouillet sheep in  
502 America. *J Hered.* 1933;24:19-33.
- 503
- 504 11. Zhang L, Mousel MR, Wu X, Michal JJ, Zhou X, Ding B, Dodson MV, El-Halawany NK,  
505 Lewis GS, Jiang Z. Genome-wide genetic diversity and differentially selected regions among

- 506 Suffolk, Rambouillet, Columbia, Polypay, and Targhee sheep. PLoS One. 2013;doi:  
507 10.1371/journal.pone.0065942.  
508
- 509 12. Rexroad C, Vallet J, Matukumalli LK, Reecy J, Bickhart D, Blackburn H, Boggess M, Cheng  
510 H, Clutter A, Cockett N, Ernst C, Fulton JE, Liu J, Lunney J, Neibergs H, Purcell C, Smith  
511 TPL, Sonstegard T, Taylor J, Telugu B, Eenennaam AV, Tassell CPV, Wells K. Genome to  
512 Phenome: Improving Animal Health, Production, and Well-Being - A New USDA Blueprint  
513 for Animal Genome Research 2018-2027. Front Genet. 2019;10:327.  
514
- 515 13. Jiang Y, Xie M, Chen W, Talbot R, Maddox JF, Faraut T, Wu C, Muzny DM, Li Y, Zhang  
516 W, Stanton JA, Brauning R, Barris WC, Hourlier T, Aken BL, Searle SMJ, Adelson DL,  
517 Bian C, Cam GR, Chen Y, Cheng S, DeSilva U, Dixen K, Dong Y, Fan G, Franklin IR, Fu S,  
518 Guan R, Highland MA, Holder ME, Huang G, Ingham AB, Jhangiani SN, Kalra D, Kovar  
519 CL, Lee SL, Liu W, Liu X, Lu C, Lv T, Mathew T, McWilliam S, Menzies M, Pan S,  
520 Robelin D, Servin B, Townley D, Wang W, Wei B, White SN, Yang X, Ye C, Yue Y, Zeng  
521 P, Zhou Q, Hansen JB, Kristensen K, Gibbs RA, Flicek P, Warkup CC, Jones HE, Oddy VH,  
522 Nicholas FW, McEwan JC, Kijas J, Wang J, Worley KC, Archibald AL, Cockett N, Xu X,  
523 Wang W, Dalrymple BP. The sheep genome illuminates biology of the rumen and lipid  
524 metabolism. Science. 2014;344:1168-1173.  
525
- 526 14. Murdoch BM. The functional annotation of the sheep genome project. J Anim Sci.  
527 2019;97:16.  
528

- 529 15. Salavati M, Caulton A, Clark R, Gazova I, Smith TPL, Worley KC, Cockett NE, Archibald  
530 AL, Clarke SM, Murdoch BM, Clark EL. Global Analysis of Transcription Start Sites in the  
531 New Ovine Reference Genome (*Oar rambouillet v1.0*). *Front Genet.* 2020;11:580580.  
532
- 533 16. FAANG Data Coordination Center. 2016.  
534 [https://data.faang.org/api/fire\\_api/samples/USU\\_SOP\\_Ovine\\_Benz2616\\_Tissue](https://data.faang.org/api/fire_api/samples/USU_SOP_Ovine_Benz2616_Tissue)  
535 [\\_Collection\\_20160426.pdf](https://data.faang.org/api/fire_api/samples/USU_SOP_Ovine_Benz2616_Tissue).  
536
- 537 17. European Bioinformatics Institute, BioSample SAMEG329607. 2016.  
538 <https://www.ebi.ac.uk/biosamples/samples/SAMEG329607>.  
539
- 540 18. Guppy (2021). Guppy basecaller (Version 3.1) [www.nanoporetech.com](http://www.nanoporetech.com).  
541
- 542 19. Koren S, Walenz BP, Berlin K, Miller JR, Bergman NH, Phillippy AM. Canu: scalable and  
543 accurate long-read assembly via adaptive *k*-mer weighting and repeat separation. *Genome*  
544 *Res.* 2017;27:722-36.  
545
- 546 20. Heaton MP, Smith TPL, Bickhart DM, Vander Ley BL, Kuehn LA, Oppenheimer J, Shafer  
547 WR, Schuetze FT, Stroud B, McClure JC, Barfield JP, Blackburn HD, Kalbfleisch TS,  
548 Davenport KM, Kuhn KL, Green RE, Shapiro B, Rosen BD. A Reference Genome Assembly  
549 of Simmental Cattle, *Bos taurus taurus*. *J Hered.* 2021;112:184-91.  
550

- 551 21. Oppenheimer J, Rosen BD, Heaton MP, Vander Ley BL, Shafer WR, Schuetze FT, Stroud B,  
552 Kuehn LA, McClure JC, Barfield JP, Blackburn HD, Kalbfleisch TS, Bickhart DM,  
553 Davenport KM, Kuhn KL, Green RE, Shapiro B, Smith TPL. A Reference Genome  
554 Assembly of American Bison, *Bison bison bison*. *J Hered.* 2021;112:174-183.  
555
- 556 22. Arima Genomics Mapping Pipeline (2019). ArimaGenomics  
557 [https://github.com/ArimaGenomics/mapping\\_pipeline](https://github.com/ArimaGenomics/mapping_pipeline).  
558
- 559 23. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin  
560 R; 1000 Genome Project Data Processing Subgroup. The Sequence Alignment/Map format  
561 and SAMtools. *Bioinformatics.* 2009;25:2078-9.  
562
- 563 24. PicardTools (2019). Picard Toolkit, Broad Institute (Version 2.9.2)  
564 <http://broadinstitute.github.io/picard>.  
565
- 566 25. Quinlan AR. BEDTools: The Swiss-Army Tool for Genome Feature Analysis.  
567 *Bioinformatics.* 2014;47:1-34.  
568
- 569 26. Ghurye J, Pop M, Koren S, Bickhart D, Chin CS. Scaffolding of long read assemblies using  
570 long range contact information. *BMC Genomics.* 2017;18:527.  
571
- 572 27. Yardımcı GG, Noble W. Software tools for visualizing Hi-C data. *Genome Biol.* 2017;18:26.  
573

- 574 28. Heng L. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics*.  
575 2018;34:3094-100.  
576
- 577 29. D-Genies (2018). D-Genies (Version 1.2.0) [https://github.com/genotoul-](https://github.com/genotoul-bioinfo/dgenies/releases/tag/v1.2.0)  
578 [bioinfo/dgenies/releases/tag/v1.2.0](https://github.com/genotoul-bioinfo/dgenies/releases/tag/v1.2.0).  
579
- 580 30. CombineFasta agp2fasta (2020). CombineFasta (Version 0.0.17)  
581 <https://github.com/njdbickhart/CombineFasta>.  
582
- 583 31. Loman N, Quick J Simpson J. A complete bacterial genome assembled de novo using only  
584 nanopore sequencing data. *Nat Methods*. 2015;12:733-735.  
585
- 586 32. NanoGrid (2018). NanoGrid <https://github.com/skoren/NanoGrid>.  
587
- 588 33. Guan D, McCarthy SA, Wood J, Howe K, Wang Y, Durbin R. Identifying and removing  
589 haplotypic duplication in primary genome assemblies. *Bioinformatics*. 2020;36:2896-8.  
590
- 591 34. Smit AFA, Hubley R, Green P. RepeatMasker Open-4.0. 2013-2015.  
592 <https://www.repeatmasker.org>.  
593
- 594 35. Garrison E, Marth G. Haplotype-based variant detection from short-read sequencing. *arXiv*  
595 preprint. 2012:1207.3907.  
596

- 597 36. Danecek, P., Bonfield, J.K., Liddle, J., Marshall, J., Ohan, V., Pollard, M.O., Whitwham, A.,  
598 Keane, T., McCarthy, S.A., Davies, R.M., Li, H.. Twelve years of SAMtools and BCFtools.  
599 *Gigascience*. 2021;10(2):giab008. doi:10.1093/gigascience/giab008  
600
- 601 37. Merfin (2021). Merfin <https://github.com/arangrhie/merfin>.  
602
- 603 38. Trim Galore (2020). TrimGalore (Version 0.6.6)  
604 [https://www.bioinformatics.babraham.ac.uk/projects/trim\\_galore/](https://www.bioinformatics.babraham.ac.uk/projects/trim_galore/).  
605
- 606 39. Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, Batut P, Chaisson M, &  
607 Gingeras TR. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics*. 2013;29:15-21.  
608
- 609 40. Li H. A statistical framework for SNP calling, mutation discovery, association mapping and  
610 population genetical parameter estimation from sequencing data. *Bioinformatics*.  
611 2011;27:2987-93.  
612
- 613 41. Shumate, A., and Salzberg, S.L. (2020). Liftoff: accurate mapping of gene annotations.  
614 *Bioinformatics*. Doi:10.1093/bioinformatics/btaa1016.  
615
- 616 42. Bray, N. L., Pimentel, H., Melsted, P., and Pachter, L. (2016). Near-optimal probabilistic  
617 RNA-seq quantification. *Nat. Biotechnol.* 34, 525–527. doi:10.1038/nbt.3519.  
618
- 619 43. Quinlan, A. R., and Hall, I. M. (2010). BEDTools: A flexible suite of utilities for comparing  
620 genomic features. *Bioinformatics* 26, 841–842. doi:10.1093/bioinformatics/btq033.

- 621
- 622 44. Themis-ASM (2020). Themis-ASM pipeline <https://github.com/njdbickhart/Themis-ASM>.
- 623
- 624 45. Vezzi F, Narzisi G, Mishra B. Reevaluating Assembly Evaluations with Feature Response
- 625 Curves: GAGE and Assemblathons. PLoS ONE. 2012;7:e52210.
- 626
- 627 46. Rhie A, Walenz BP, Koren S, Phillippy AM. Merqury: reference-free quality, completeness,
- 628 and phasing assessment for genome assemblies. Genome Biol. 2020;21:245.
- 629
- 630 47. Koren S, Rhie A, Walenz BP, Dilthey AT, Bickhart DM, Kingan SB, Hiendleder S, Williams
- 631 JL, Smith TPL, Phillippy AM. De novo assembly of haplotype-resolved genomes with trio
- 632 binning. Nat Biotechnol. 2018;10.1038/nbt.4277.
- 633
- 634
- 635
- 636
- 637
- 638
- 639
- 640
- 641
- 642
- 643
- 644

645

646

647

648 **Tables**

649

650 Table 1: Assembly quality statistics comparison

Assembly Statistic	ARS-UI_Ramb_v2.0	Oar_rambouillet_v1.0	Oar_v4.0	Description
<b>Total Length (Mb)</b>	2628.15	2869.91	2615.52	Assembly length in Mbp
<b>Contig Number</b>	226	7,486	48,482	Total number of contigs
<b>Contig NG50 (bp)</b>	43,178,051	2,850,956	145,655	Half the length of the genome is in contigs of this size or greater, based on a 2600 Mb genome
<b>Contig LG50 (number of contigs)</b>	24	263	5,206	The smallest number of contigs whose length sum make up half of the assembly size
<b>Scaffold Number</b>	142	2,641	5,466	Total number of scaffolds and unplaced contigs in the assembly
<b>Scaffold L50 (number of scaffolds)</b>	8	8	8	The smallest number of scaffolds whose length sum make up half of the assembly size.
<b>merQV</b>	44.7721*	32.1705*	31.9131**	Kmer based quality from Merqury, which estimates the frequency of consensus errors in the assembly [46]
<b>merErrorRate</b>	0.000033327*	0.00060662*	0.000643714**	Kmer based error rate from Merqury, which estimates error rate of the assembly based on errors in kmers [46]
<b>merCompleteness</b>	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 [46]
<b>baseQV</b>	41.84*	40.69*	32.40**	SNP and INDEL quality value estimated from short read data mapped to the assembly [47]
<b>Unmap%</b>	0.96*	1.00*	0.73**	Percentage of short reads that are unmapped to each assembly [47]
<b>COMPLETESC</b>	90.7	90.1	86.1	Percent of complete, single copy BUSCOs
<b>COMPLETEDUP</b>	1.4	1.6	1.0	Percent of complete, duplicated BUSCOs
<b>FRAGMENT</b>	2.0	2.1	3.7	Percent of fragmented BUSCOs
<b>MISSING</b>	5.9	6.2	9.2	Percent of missing BUSCOs

651

652 \*Short read sequencing from the Rambouillet ewe used to assemble both ARS-UI\_Ramb\_v2.0

653 and Oar\_rambouillet\_v1.0 was used in these quality values.



654 \*\*Short read sequencing from the Texel animal used to assemble Oar\_v4.0 was used in these  
 655 quality values.

656

657 Table 2: Specific feature counts for each genome and descriptions.

Features	ARS-UI_Ramb_v2.0	Oar_rambouillet_v1.0	Oar_v4.0	Description
<b>LOW_COV_PE</b>	7212	95166	89103	Low read coverage areas
<b>LOW_NORM_COV_PE</b>	2990	24381	26860	Low coverage of normal paired end reads
<b>HIGH_SPAN_PE</b>	6522	22628	33232	Regions with high numbers of inter-contig paired end reads
<b>HIGH_COV_PE</b>	2051	3630	26276	Regions with high read coverage
<b>HIGH_NORM_COV_PE</b>	2366	2633	1875	Regions with high coverage of normal paired end reads
<b>HIGH_OUTIE_PE</b>	2514	28766	37495	Regions with high counts of improperly paired reads
<b>HIGH_SINGLE_PE</b>	0	0	0	Regions with high counts of single unmapped reads
<b>STRECH_PE</b>	74	84	67	Regions with high Comp/Expansion (CE) statistics
<b>COMPR_PE</b>	87	92	44	Regions with low Comp/Expansion (CE) statistics

658

659

660

661

662

663

664

665

666

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

667 Table 3: RNA-seq alignment statistics to ARS-UI\_Ramb\_v2.0 and Oar\_rambouillet\_v1.0 from  
668 five different tissues.

669

670 \* Genomes include v2.0 (ARS-UI\_Ramb\_v2.0) and v1.0 (Oar\_rambouillet\_v1.0) and the

671 difference (Δ).

672

673

674

675

676

677 Table 4: Expressed transcripts (TPM&gt;0) in Benz2616 tissues (n=61) based on

678 Oar\_rambouillet\_v1.0 and ARS-UI\_Ramb\_v2.0 and lift over (LO) (RefSeq v103 &amp; 104,

679 respectively).

Gene Biotype	Oar_rambouillet_v1.0	Oar_rambouillet_v1.0 LO	ARS-UI_Ramb_v2.0	LO vs. Oar_rambouillet_v1.0	LO vs. ARS-UI_Ramb_v2.0	Oar_rambouillet_v1.0 vs. ARS-UI_Ramb_v2.0
Guide RNA	30	29	30	-1	-1	0
lncRNA	3929	3752	6018	-177	-2266	-2089
Protein coding	42058	40910	60064	-1148	-19154	-18006
rRNA	272	17	22	-255	-5	250
snoRNA	644	590	593	-54	-3	51
snRNA	997	907	879	-90	28	118

680

681

682

683

684

685

686

687

688

689

690

691

692

693

**694 Figure Legends**

695

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

698

699 Figure 2: Hi-C contact map comparison of ARS-UI\_Ramb\_v2.0 A) directly after scaffolding and  
700 before manual curation and B) after manual curation with scaffold rearrangements and joins.

701

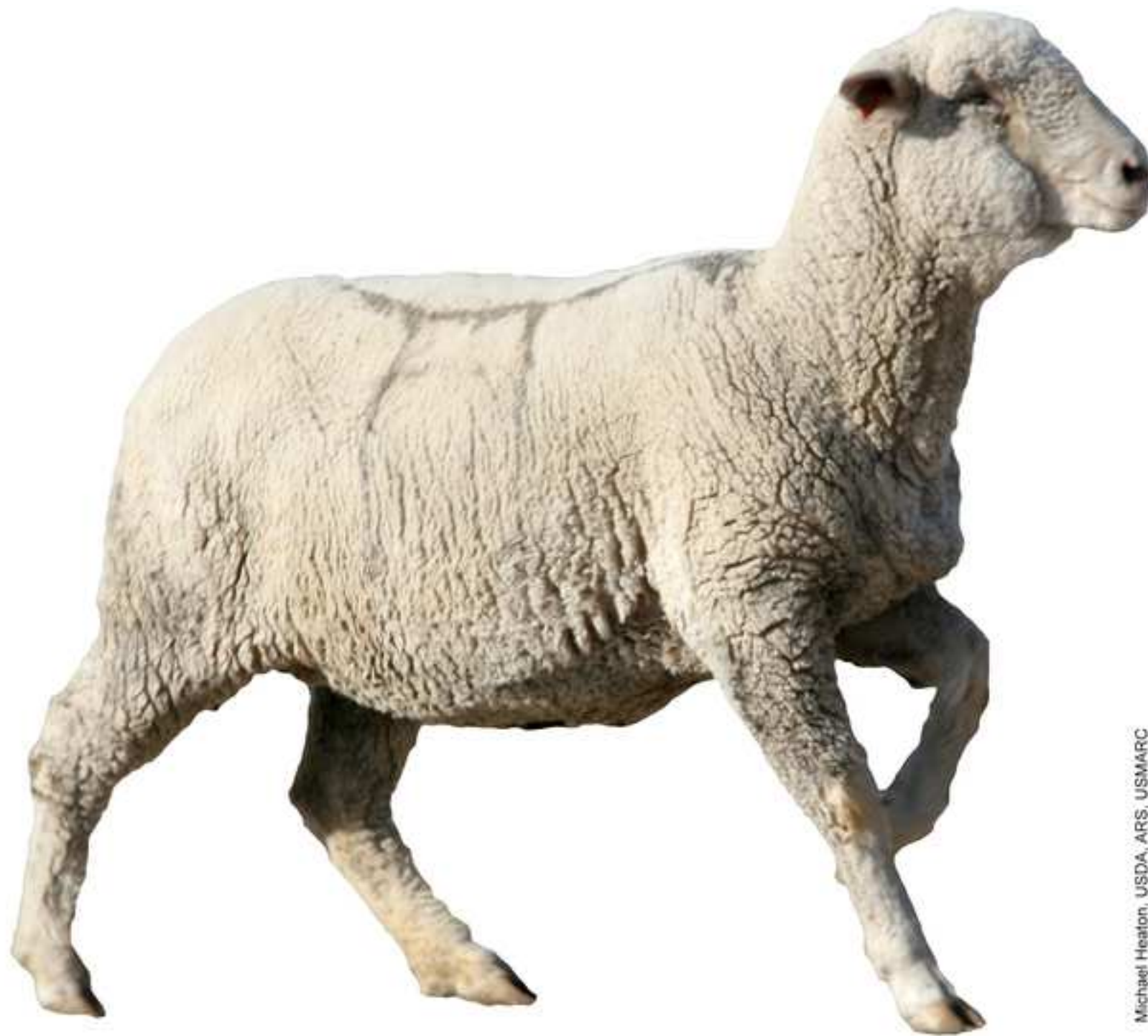
702 Figure 3: Assembly error comparison between ARS-UI\_Ramb\_v2.0, Oar\_rambouillet\_v1.0, and  
703 Oar\_v4.0 in a feature response curve displaying sorted lengths of the assemblies with the fewest  
704 errors.

705

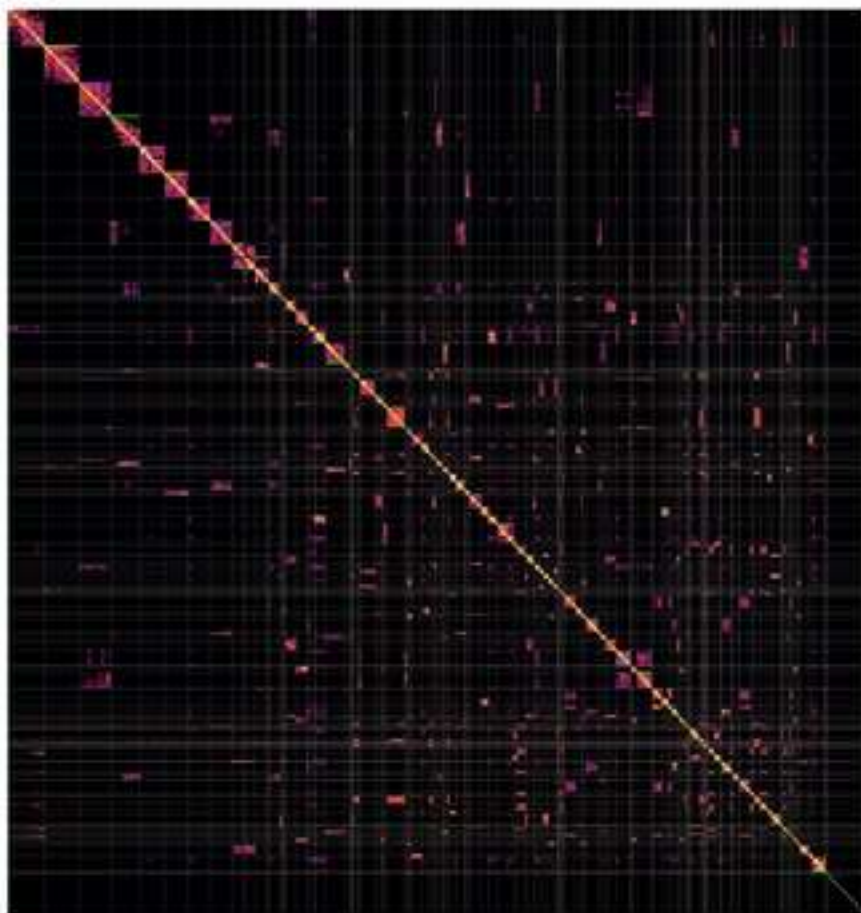
706 Figure 4: Dotplot comparison of genome assemblies between A) ARS-UI\_Ramb\_v2.0 and  
707 Oar\_rambouillet\_v1.0, and B) ARS-UI\_Ramb\_v2.0 and Oar\_v4.0.

708

709 Figure 5: Kallisto comparison of the number of expressed transcripts for the RNA-Seq dataset of  
710 61 tissue samples from Benz2616, across the three annotations (Oar\_rambouillet\_v1.0,  
711 Ramb1LO2 (liftover) and ARS-UI\_Ramb\_v2.0).



Michael Heaton, USDA, ARS, USMARC

**A****B**