

Scotland's Rural College

Compendium of 4,941 rumen metagenome-assembled genomes for rumen microbiome biology and enzyme discovery

Stewart, Robert; Auffret, MD; Warr, Amanda; Walker, Alan; Roehe, R; Watson, Mick

Published in:
Nature Biotechnology

DOI:
[10.1038/s41587-019-0202-3](https://doi.org/10.1038/s41587-019-0202-3)

Print publication: 02/08/2019

Document Version
Peer reviewed version

[Link to publication](#)

Citation for published version (APA):

Stewart, R., Auffret, MD., Warr, A., Walker, A., Roehe, R., & Watson, M. (2019). Compendium of 4,941 rumen metagenome-assembled genomes for rumen microbiome biology and enzyme discovery: Comprehensive resource of cow rumen genomes and a database of predicted proteins. *Nature Biotechnology*, 37, 953-961. <https://doi.org/10.1038/s41587-019-0202-3>

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal ?

Take down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

1 **Editors summary**

2 **Comprehensive resource of cow rumen genomes and a database of predicted proteins.**

3
4
5 **Compendium of 4941 rumen metagenome-assembled genomes**
6 **for rumen microbiome biology and enzyme discovery**

7
8 Robert D. Stewart¹, Marc D. Auffret², Amanda Warr¹, Alan W. Walker³, Rainer Roehe² and Mick Watson^{1*}

9 ¹The Roslin Institute and R(D)SVS, University of Edinburgh, Easter Bush EH25 9RG, UK.

10 ²Scotland's Rural College, Edinburgh EH25 9RG, UK.

11 ³The Rowett Institute, University of Aberdeen, Aberdeen AB25 2ZD, UK

12 * Corresponding author: mick.watson@roslin.ed.ac.uk

13
14 **Abstract**

15 **Ruminants provide essential nutrition for billions of people worldwide. The rumen is a**
16 **specialised stomach that is adapted to the breakdown of plant-derived complex**
17 **polysaccharides. The genomes of the rumen microbiota encode thousands of enzymes adapted**
18 **to the digestion of plant matter which dominates the ruminant diet. We assembled 4941**
19 **rumen microbial metagenome-assembled genomes (MAGs) using ~ 6.5 terabytes of short- and**
20 **long-read sequence data from 283 ruminant cattle. We present a genome-resolved**
21 **metagenomics workflow that enabled assembly of bacterial and archaeal genomes that were**
22 **at least 80% complete. Of note, we obtained 3 single-contig, whole-chromosome assemblies of**
23 **rumen bacteria, two of which represent previously unknown rumen species, assembled from**
24 **long-read data. Using our rumen genome collection, we predicted and annotated the largest**
25 **set of rumen proteins to date. Our set of rumen MAGs increases the rate of mapping of**
26 **rumen metagenomic sequencing reads from 15% to 50-70%. These genomic and protein**
27 **resources will enable a better understanding of the structure and functions of the rumen**
28 **microbiota.**

29
30
31
32 **Introduction**

33 Ruminants convert human-inedible, low value plant biomass into products of high nutritional
34 value, such as meat and dairy products. The rumen, which is the first of four chambers of the
35 stomach, contains a mixture of bacteria, archaea, fungi and protozoa that ferment complex
36 carbohydrates e.g. lignocellulose, cellulose, to produce short-chain fatty acids (SCFAs) that the
37 ruminant uses for homeostasis and growth. Rumen microbes are a rich source of enzymes for
38 plant biomass degradation for use in biofuels production¹⁻³, and manipulation of the rumen
39 microbiome offers opportunities to reduce the cost of food production⁴.

40 Ruminants are important for both food security and climate change. For example, methane
41 is a by-product of ruminant fermentation, released by methanogenic archaea, and an estimated
42 14% of methane produced by humans has been attributed to ruminant livestock⁵. Methane
43 production has been directly linked to the abundance of methanogenic archaea in the rumen⁶,
44 offering possibilities for mitigating this issue through selection⁷ or manipulation of the
45 microbiome. Two studies have reported large collections of rumen microbial genomes. Stewart *et al*
46 assembled 913 draft metagenome-assembled genomes (MAGs) (named rumen-uncultured
47 genomes (RUGs)) from the rumen of 43 cattle raised in Scotland⁸ and Seshadri *et al* reported 410
48 reference archaeal and bacterial genomes from the Hungate collection⁹. As isolate genomes, the
49 Hungate genomes are generally higher quality, and crucially, exist in culture, so can be grown and
50 studied in the lab. However, we found that addition of the Hungate genomes only increased read
51 classification by 10%, compared to an increase of 50-70% when the RUGs are used, indicating
52 significant numbers of undiscovered microbes in the rumen.

53 We present a comprehensive analysis of more than 6.5Tb of sequence data from the
54 rumen of 283 cattle. Our catalog of rumen genomes (named RUG2) includes 4056 genomes that
55 were not present in Stewart *et al*⁸, and brings the number of rumen genomes assembled to date
56 to 5845. We also present a metagenomic assembly of nanopore (MinION) sequencing data (from
57 one rumen sample) that contains at least three whole bacterial chromosomes as single contigs,
58 and which represents the most continuous metagenomic assembly from the rumen to date. These
59 genomic and protein resources will underpin future studies on the structure and function of the
60 rumen microbiome.

61 Results

62 4941 metagenome-assembled genomes from the cow rumen

63 We sequenced DNA extracted from the rumen contents of 283 beef cattle (characteristics of
64 animals sequenced is in Supplementary Data 1), producing over 6.5Tb of Illumina sequence data.
65 We operate a continuous assembly-and-de-replication pipeline, which means that newer genomes
66 of the same strain (>99% average-nucleotide identity, ANI) can replace older genomes if their
67 completeness and contamination statistics are better. All of the 4941 RUGs we present here have
68 completeness \geq 80% and contamination \leq 10% (Supplementary Figure 1).

69 4941 RUGs were analysed using MAGpy¹⁰ and their assembly characteristics, putative
70 names and taxonomic classifications are in Supplementary data 2. Sourmash¹¹, DIAMOND¹² and
71 PhyloPhlAn¹³ outputs, which reveal genomic and proteomic similarity to existing public data, are in
72 Supplementary data 3. A phylogenetic tree of the 4941 RUGs, alongside 460 public genomes from
73 the Hungate collection, is presented in Figure 1 and Supplementary data 4. The tree is dominated
74 by large numbers of genomes from the *Firmicutes* and *Bacteroidetes* phyla (dominated by
75 *Clostridiales* and *Bacteroidales* respectively), but also contains many novel genomes from the
76 *Actinobacteria*, *Fibrobacteres* and *Proteobacteria* phyla. *Clostridiales* (2079) and *Bacteroidales*
77 (1081) are the dominant orders, with *Ruminococacae* (1111) and *Lachnospiraceae* (640) the
78 dominant families within the *Clostridiales*, and *Prevotellaceae* (521) the dominant family within
79 the *Bacteroidales*.

80 The genome taxonomy database (GTDB) proposed a new bacterial taxonomy based on
81 conserved concatenated protein sequences¹⁴, and we include the GTDB predicted taxa for all

82 RUGs (Supplementary data 3). 4763 RUGs have < 99% average-nucleotide-identity (ANI) with
83 existing genomes, and 3535 have < 95% ANI with existing genomes and therefore represent
84 potential novel species.

85 144 of 4941 genomes are classified to species level, 1092 of 4941 to genus level, 3188 of
86 4941 to family, 4084 to order, 4514 to class, 4801 to phylum and 4941 to kingdom. Of the
87 genomes classified at the species level, 43 represent genomes derived from uncultured strains of
88 *Ruminococcus flavefaciens*, 42 represent genomes from uncultured strains of *Fibrobacter*
89 *succinogenes*, 18 represent genomes from uncultured strains of *Sharpea azabuensis*, and 10
90 represent genomes from uncultured strains of *Selenomonas ruminantium*. These species belong to
91 genera known to play an important role in rumen homeostasis¹⁵.

92 We assembled 126 archaeal genomes, 111 of which are species of *Methanobrevibacter*.
93 There are two other members of the *Methanobacteriaceae* family, both predicted to be a member
94 of the *Methanosphaera* genus by GTDB. Nine of the archaeal RUGs have sourmash hits to
95 “*Candidatus Methanomethylophilus* sp. 1R26”; a further three have weak sourmash hits to
96 “Methanogenic archaeon ISO4-H5”; and the remaining archaeal genome has no sourmash hits,
97 and weaker DIAMOND hits to the same genome (“Methanogenic archaeon ISO4-H5”). All thirteen
98 are predicted to be members of the genus “*Candidatus Methanomethylophilus*” by GTDB, but this
99 is based on similarity to only two genomes, both of which have uncertain phylogenetic lineages. If
100 “*Candidatus Methanomethylophilus*” is a true genus, then our dataset increases the number of
101 sequenced genomes from 2 to 15.

102 Genome quality statistics were measured by analysing single copy core-genes (Supp Figure
103 1). There are different standards for the definition of MAG quality: Bowers *et al*¹⁶ describe high-
104 quality drafts as having $\geq 90\%$ completeness and $\leq 5\%$ contamination; 2417 of the RUGs meet
105 these criteria. Alternatively, Parks *et al*¹⁷ define a quality score as “Completeness – (5 *
106 contamination)” and exclude any MAG with a score less than 50; 4761 of the RUGs meet that
107 criterion; however, whilst the MAGs from Parks *et al* could be as low as 50% complete, the
108 genomes presented here are all $\geq 80\%$ complete. The RUGs range in size from 456kb to 6.6Mb
109 with N50s (50% of assembled bases are in contigs greater than the N50 value) ranging from 4.5kb
110 to 1.37Mb. The average number of tRNA genes per RUG is 16.9 and 446 of the RUGs have all 20.
111 As assemblies of Illumina metagenomes struggle to assemble repetitive regions, most of the RUGs
112 do not contain a 16S rRNA gene – 464 RUGs encode a fragment of the 16S rRNA gene, and 154
113 encode at least one full length 16S rRNA gene.

114 The coverage of each RUG in each sample is in Supplementary Data 5. Using a cut-off of 1X
115 coverage, most RUGs (4863) are present in more than one animal, 3937 are present in more than
116 10 animals, and 225 RUGs are present in more than 200 animals. One RUG is present in all
117 animals, RUG11026 a member of the *Prevotellaceae* family.

118

119

120 **A near-complete single-contig Proteobacteria genome**

121 Metagenomic assembly of Illumina data often results in highly fragmented assemblies but
122 RUG14498, an uncultured *Proteobacteria* species (genome completeness 87.91% and

123 contamination 0%) has 136 of 147 single-copy genes present with no duplications in a single contig
124 of just over 1Mb in size. *Proteobacteria* with small genomes (<1.5Mb size) are relatively common
125 (n=67) in our dataset and have also been found in other large metagenome assembly projects¹⁷.
126 The *Proteobacteria* genomes we present encode proteins with only 45% to 60% amino acid
127 identity with proteins in UniProt TREMBL¹⁸. We compared our single-contig *Proteobacteria*
128 assembly with nine *Proteobacteria* with similarly sized genomes assembled by Parks *et al*¹⁷ (see
129 whole-genome alignments in Supplementary Figure 2). Average nucleotide identity (ANI; often
130 used to delineate new strains and species) between the 9 UBA genomes and RUG14498 is
131 revealing. UBA2136, UBA1908, UBA3307, UBA3773 and UBA3768 have no detectable level of
132 identity with any other genome in the set; UBA4623, UBA6376, UBA6864, and UBA6830 all share
133 greater than 99.4% average nucleotide identity with one another, indicating that they are highly
134 similar strains of the same species. UBA4623, UBA6376, UBA6864 and UBA6830 also show around
135 77.8% ANI with RUG14498 suggesting that the single-contig RUG14498 is a high-quality, near-
136 complete whole genome of a novel *Proteobacteria* species. The single contig RUG14498 was
137 assembled by IDBA_ud from sample 10678_020. IDBA_ud exploits uneven depth in metagenomic
138 samples to improve assemblies. RUG14498 is the tenth most abundant genome in 10678_020, and
139 other genomes of similar depth in that sample are taxonomically unrelated, enabling IDBA_ud to
140 assemble almost the entire genome in a single contig.

141 RUG14498 has a single full length 16S rRNA gene (1507bp). The top hit in GenBank (97%
142 identity across 99% of the length) is accession AB824499.1, a sequence from an “uncultured
143 bacterium” from “the rumen of Thai native cattle and swamp buffaloes”. The top hit in SILVA¹⁹ is
144 to the same sequence, only this time annotated as an uncultured *Rhodospirillales*. Together these
145 results support the conclusion that RUG14498 represents a novel *Proteobacteria* species. Low
146 amino acid identity to known proteins limits our ability to predict function and metabolic activity;
147 nevertheless, RUG14498 encodes 73 predicted CAZymes, including 42 glycosyl transferases and 19
148 glycosyl hydrolases, suggesting a role in carbohydrate synthesis and metabolism.

149 **Novel microbial genomes from the rumen microbiome**

150 We compared 4941 RUGs to the Hungate collection and to our previous dataset⁸ (Figure 2).
151 149/4941 RUGs share > 95% protein identity with Hungate members; 271/4941 > than 90%; this
152 leaves 4670/4941 RUGs with < 90% protein identity with Hungate members. 2387/4941 RUGs
153 have < 90% protein identity with genomes in Stewart *et al*, and more than 1100 RUGs have < 70%
154 protein identity with Stewart *et al*. Many of the RUGs with the lowest protein identity to public
155 genomes could not be classified beyond Phylum level, and some are simply “uncultured
156 bacterium”.

157 We compiled a database comprising all RUG genomes, the Hungate collection genomes⁹,
158 and rumen MAGs from Hess *et al*¹, Parks *et al*¹⁷, Solden *et al*²⁰ and Svartström *et al*²¹ that we name
159 the “rumen superset”. The rumen superset was dereplicated at both 99% (strain-level) and 95%
160 (species level) average-nucleotide identity (ANI). At 95% ANI, the rumen superset was reduced to
161 2690 clusters, representing species-level bins. 2078 of these clusters contain only RUG genomes,
162 and therefore represent putative novel rumen microbial species identified in this study. 58 clusters
163 contain both Hungate and RUG genomes, and 268 clusters contained only Hungate genomes
164 (Supplementary Data 6). At 99% ANI, the rumen superset was reduced to 5574 clusters,
165 representing strain-level bins. 4845 of these clusters contain only RUG genomes, and may

166 represent putative novel rumen microbial strains (Supplementary Data 7). Supplementary Figure 3
167 shows how the various rumen MAG sets overlap at 95% ANI after de-replication.

168 We calculated an estimate of the completeness of the RUG2 dataset using the Chao 1
169 estimator²² (note we can only do this for our own dataset as it is based on the number of times
170 species are observed at different frequencies, and we do not have these values for other
171 datasets). De-replicating all RUG genomes at 95% gives us 2180 species-level bins. 948 of those are
172 singletons (i.e. observed exactly once), and 410 are doublets (i.e. observed exactly twice). Using
173 the Chao 1 formula, we predict 3276 species, which means we estimate that we have discovered
174 66.54% of the species present in our samples.

175 We assessed the impact of using rumen genomic data on the read classification rates of
176 several public datasets using three databases – the first, our custom rumen kraken database
177 consisting of RefSeq complete genomes and the Hungate collection (previously described^{23,24}); the
178 second was the same database plus only the RUGs; and the third was the same database plus the
179 rumen superset (which includes the RUGs). We classified five datasets – our own (Stewart *et al*), a
180 dataset we previously published (Wallace *et al*⁶), data from 14 cows from a study on niche
181 specialisation (Rubino *et al*²⁵), data from a methane emissions study of sheep (Shi *et al*²⁶) and a
182 recent metagenomic study of moose (Svartström *et al*²¹) (Supplementary Figure 4).

183 The classification rate is increased by using either the RUG or rumen superset databases,
184 though the rumen superset achieves only a marginal increase in most cases. We have improved
185 read classification rates from 15% to 70%, with more than a quarter of our samples achieving a
186 classification rate of 80% or higher. These are comparable with read classification rates for the
187 human microbiome as reported by Pasolli *et al*²⁷.

188

189 **Strain-level analysis of methane emissions in sheep**

190 Previously Shi *et al*²⁶ found no significant changes in community structure between low-
191 and high- methane emitting sheep, although there were differences in gene expression between
192 the two groups. We re-analysed the Shi *et al* dataset using our rumen metagenomic data;
193 specifically, we used our custom kraken database consisting of RefSeq genomes and the rumen
194 superset and used it to classify reads at the level of Kingdom, Phylum, Family, Genus and Species,
195 and tested differences between low methane-emitting (LME) and high-methane emitting (HME)
196 sheep. Whilst we found no significant differences at the level of Kingdom, we found significant and
197 profound differences at every other taxonomic level tested (Supplementary Tables 1-5 and
198 Supplementary Figures 5-9). At the Genus level, *Sharpea*, *Kandleria*, *Fibrobacter* and *Selenomonas*
199 are associated with LME sheep, and *Elusimicrobium* with HME sheep (Supplementary Table 4). At
200 the species level, we found that 340 species differ significantly between LME and HME emitting
201 sheep (Supplementary Table 5), including eleven species of *Bifidobacterium*, and six species of
202 *Olsenella*, all significantly more proportionally abundant in LME sheep, and nine species of
203 *Desulfovibrio* significantly more proportionally abundant in HME sheep. *Fibrobacter succinogenes*,
204 an important rumen microbe known to be heavily involved in plant fibre degradation, is also
205 significantly different between the two groups, and is associated with LME sheep. Some of these
206 microbes were previously identified as differentially proportionally abundant between LME and

207 HME sheep^{15,28} using marker-gene sequencing, though our results provide greater resolution and
208 for the first time reveal the genome sequences involved.

209 Kraken classifies data at different levels of the NCBI taxonomy; unfortunately, this does not
210 give us data on the RUGs which do not yet have specific NCBI taxonomy IDs. Therefore, to
211 estimate the abundance of individual strains, we aligned reads directly to the rumen superset, and
212 used the number of reads designated as primary alignments as a proxy for the relative abundance
213 of each genome. At $FDR \leq 0.05$, 1709 genomes show differentially proportional abundance
214 between low- and high-methane sheep (Supplementary Data 8, Supplementary Figure 10). In
215 supplementary figure 10, LME and HME sheep are clearly separated along principal component 1,
216 which explains 58% of the variance in the data. Supplementary data 8 lists the differentially
217 abundant genomes; of note are large numbers of previously uncharacterised *Lachnospiraceae*
218 species associated with LME sheep; and 22 strains of *Sharpea azabuensis* all higher proportional
219 abundance in LME sheep (all 18 *Sharpea azabuensis* RUGs and four *Sharpea azabuensis* strains
220 from the Hungate collection). These results agree with previous studies based on marker-genes¹⁵,
221 and our dataset increases the number of *Sharpea azabuensis* genomes publicly available from 4 to
222 22. Large numbers of uncharacterised *Ruminococcaceae* and *Bacteroidia* are also associated with
223 HME sheep. Multiple strains of uncharacterised *Proteobacteria*, including RUG14498 described
224 above, are more proportionally abundant in HME sheep; and *Fibrobacter* strains were almost all
225 associated with LME sheep.

226 The relationship between proportional abundance of Archaea and methane emissions is
227 not simple. Most archaeal strains are present at similar abundance in LME and HME sheep
228 (Supplementary Data 8). RUGs representing novel strains of *Methanobrevibacter* are often more
229 abundant in HME sheep. The RUG with the most striking proportional abundance is RUG12825,
230 which is likely a member of the *Methanosphaera* genus, and is more abundant in LME sheep. The
231 complex relationship between relative abundance of methanogens and methane emissions may
232 underlie our inability to find significant differences in overall archaeal proportional abundance.

233 That notwithstanding, these data represent the first strain-level view of methane emissions
234 in sheep to our knowledge, and support and confirm the hypothesis that there are major,
235 fundamental changes in rumen metagenomic relative abundance associated with extremes of low
236 and high methane emissions.

237 **Global rumen census updated**

238 The global rumen census attempted to determine the core rumen microbiome by using 16S rRNA
239 sequencing of rumen samples from 742 individual animals from around the world, comprising
240 eight ruminant species²⁹. *Prevotella*, *Butyrivibrio*, and *Ruminococcus*, as well as unclassified
241 *Lachnospiraceae*, *Ruminococcaceae*, *Bacteroidales*, and *Clostridiales* were the dominant rumen
242 bacteria and which may represent a core bacterial rumen microbiome. The same species are
243 abundant in our data (Supplementary Data 5). We also find that many *Proteobacteria* are highly
244 abundant, including *Succinivibrio* (Supplementary Data 5). This is noteworthy because
245 *Proteobacteria* were found to be highly abundant in many of the samples from the rumen census,
246 but were not highlighted as being part of the core rumen microbiome.

247 To further characterise the proportional abundance of *Proteobacteria* we used the rumen
248 superset database to classify data from this study, Wallace *et al*⁶, Rubino *et al*²⁵, Shi *et al*¹⁵ and

249 Svartström *et al*²¹ (Supplementary Figure 11). *Proteobacteria* are present in all datasets; abundant
250 in cattle datasets, but less so in moose and sheep. Given the high proportional abundance of
251 *Proteobacteria* in many samples, and their consistent presence in all of the samples we tested, we
252 suggest adding *Proteobacteria* to the core bacterial rumen microbiome that was proposed by
253 Henderson *et al*²⁹.

254 Long-read assembly of complete bacterial chromosomes

255 We analysed a single sample (10572_0012) using a MinION sequencer and compare Illumina and
256 MinION assembly statistics in Figure 3. Three flowcells produced 11.4Gb of data with a read N50 of
257 11,585bp. The mean read length was 6144bp, which is short compared to other reports^{30,31}. We
258 attribute this to short DNA fragments and nicks caused by bead-beating step during DNA
259 extraction. We assembled long reads using Canu³², to form an assembly 178Mb in length with an
260 N50 of 268kb. Regardless of length, Canu predicted 31 of the contigs to be circular. These circular
261 contigs might represent putative plasmids or other circular chromosomes.

262 One problem with single-molecule sequencing technologies is the presence of post-
263 assembly insertions and deletions (indels)³³. Canu can correct reads but not enough to remove all
264 indels. Detecting sequencing errors without a ground truth dataset is difficult so we hypothesized
265 that most indels would create premature stop-codons and that gene prediction tools (eg
266 Prodigal³⁴) would produce truncated proteins. We examined the ratio between the lengths of
267 predicted proteins and their top-hits in UniProt to estimate indels (Supplementary Figure 12).
268 Although these data indicate multiple errors compared with the Illumina short-read data, we
269 corrected errors by polishing with one round of Nanopolish and two rounds of Racon. We set-up a
270 software pipeline to calculate statistics and produce similar plots for any input genome or
271 metagenome called "IDEEL".

272 Statistics for all contigs \geq 500kb and all contigs predicted to be circular are in
273 Supplementary data 9. The Nanopore assembly contains several single contigs that we predict are
274 complete, or near-complete, circular whole chromosomes.

275 *Prevotella copri* nRUG14950 (tig00000032) is a single contig of 3.8Mb which most closely
276 resembles *Prevotella copri* DSM 18205, and which shows high similarity to RUG14032. *Prevotella*
277 *copri* nRUG14950 is predicted to be 98.48% complete by CheckM³⁵, with a contamination score of
278 2.03%; whereas RUG14032 is estimated to be 96.62% complete and 1.35% contaminated.
279 Comparative alignments between *Prevotella copri* nRUG14950, RUG14032 and *Prevotella copri*
280 DSM 18205 can be seen in Supplementary Figure 13. There is a clear and striking relationship
281 between *Prevotella copri* nRUG14950 and RUG14032. These two genomes, both estimated to be
282 near-complete, were assembled from different samples using different techniques, and sequenced
283 with different sequencing technologies. Our assembly of *Prevotella copri* nRUG14950, with only
284 one contig and estimated to be 98.48% complete, represents the most continuous chromosomal
285 assembly of *Prevotella copri* to date, despite having been assembled from a metagenome.

286 *Selenomonas* spp. nRUG14951 is a single contig of 3.1Mb in length, predicted to be
287 circular, and with completeness and contamination statistics of 98.13% and 0.16% respectively.
288 The most similar RUG is RUG10160, sharing a mean of 94% protein identity. RUG10160 is
289 estimated 97.66% complete and 0% contaminated. However, the closest public reference genome
290 is *Selenomonas ruminantium* GACV-9, part of the Hungate collection, which shares only ~64%

291 protein identity with *Selenomonas* spp. nRUG14951. There exists a good whole-genome alignment
292 between *Selenomonas* spp. nRUG14951 and RUG10160 (Supplementary Figure 14), albeit with
293 some evidence of re-arrangements, and some small sections of the genome that are only captured
294 by the Nanopore assembly.

295 We also identified *Lachnospiraceae bacterium* nRUG14952, which has a 2.5Mb circular,
296 near-complete genome (95.46%), a second RUG13141 (which has 96% protein identity to
297 nRUG14952) and a more distantly related public reference genome (*Lachnospiraceae bacterium*
298 KHCPX20, 63% protein identity to nRUG14952). The nanopore-assembled genome
299 *Lachnospiraceae bacterium* nRUG14952 contains several genome regions that are absent from
300 RUG13141 (Supplementary Figure 15).

301 nRUG14951 and nRUG14952 represent entire bacterial chromosomes assembled as single
302 contigs and are the first genome assemblies for these species. The remainder of the nanopore
303 assembly contains highly continuous contigs that represent large portions of previously
304 unsequenced bacterial chromosomes. These results taken together demonstrate the power of
305 long reads for assembling complete, whole chromosomes from complex metagenomes.

306 To assess the advantage of having complete chromosomal assemblies, we annotated the
307 three nanopore whole genomes and the 3 genomes of their closely related RUGs (Supplementary
308 Data 10). The three complete nanopore genomes contain 5, 7 or 3 full length 16S gene sequences
309 respectively, whereas all three RUGs contain none. In addition, the three nanopore genomes are
310 massively enriched for IS family transposase proteins compared to their RUG counterparts.
311 Transposases are associated with insertion sequences in bacterial genomes, and catalyse the
312 transposition of mobile elements³⁶. Finally, in all cases, the nanopore assemblies have more
313 annotated COGs (“clusters of orthologous genes”), suggesting that they have a more complete
314 functional annotation than their short-read counterparts.

315 **A protein database for rumen microbial proteomics**

316 We put together a non-redundant dataset of rumen proteins from the 4941 RUGs and 460 publicly
317 available Hungate collection genomes (10.69 million proteins), following the model of UniRef³⁷
318 and clustering the protein set at 100% (9.45 million clusters), 90% (5.69 million clusters) and 50%
319 (2.45 million clusters) identity to form RumiRef100, RumiRef90 and RumiRef50.

320 To assess the novelty of our dataset at the protein level as compared to other rumen MAG
321 datasets, we took RumiRef100 and added over 900,000 predicted proteins from the rumen
322 superset. We clustered these at 90% identity, which resulted in 6.24 million protein clusters. Of
323 these, 5 million clusters contain at least one RUG protein, 4.74 million contain only RUG proteins,
324 and 3.67 million are singletons containing only RUG proteins.

325 All 10.69 million predicted proteins from the RUGs have been compared to KEGG³⁸, 460
326 public genomes from the Hungate collection, UniRef100, UniRef90 and UniRef50. The mean
327 protein identity of the top hit for these databases is 55.88%, 63.58%, 67.52%, 67.25% and 59.97%
328 respectively. These data provide the most comprehensive and richly annotated protein dataset
329 from the rumen to date.

330 The RUG proteins were compared to the CAZy³⁹ database (31st July 2018) using dbCAN2⁴⁰.
331 442,917 are predicted to be involved in carbohydrate metabolism, including 235,001 glycoside

332 hydrolases, 120,494 glycosyltransferases, 55,523 carbohydrate esterases, 23,928 proteins with
333 carbohydrate binding modules, 6834 polysaccharide lyases, 907 proteins with predicted auxiliary
334 activities, 80 proteins with a predicted cohesin domain, and 150 proteins with an S-layer homology
335 module (SLH).

336 The similarity of the predicted CAZymes to the current CAZy database can be seen in Figure
337 4. None of the eight classes of carbohydrate active enzymes displays an average protein identity
338 greater than 60% indicating that CAZy poorly represents the diversity of CAZymes encoded in the
339 genomes of ruminant microbes. Of particular note is the class AA “auxiliary activities”, with an
340 average protein identity of less than 30% between CAZy and the RUG CAZymes. AA was created by
341 CAZy to classify ligninolytic enzymes and lytic polysaccharide mono-oxygenases (LPMOs).

342 The distribution of CAZymes across 12 different phyla and the group of “unknown” bacteria
343 can be seen in Figure 5. The *Bacteroidetes* (3.9 million) and *Firmicutes* (5.3 million) together
344 contribute the largest number of proteins to our dataset; however, whereas 5.7% of the proteome
345 of *Bacteroidetes* is devoted to CAZyme activity, in *Firmicutes* the figure is 3.2%. *Fibrobacteres*
346 devote the highest percentage of their proteome to carbohydrate metabolism (over 6.6%), as is
347 expected due to their fibre-attached, high cellulolytic activity. Only a few studies exist on the role
348 of *Planctomycetes* in the rumen^{24,41,42}, however whilst they contribute a relatively low number of
349 proteins in our dataset (30172), just over 5% of those proteins are predicted to be CAZymes,
350 suggesting a role in and adaptation to carbohydrate metabolism. 79 out of 80 cohesin-containing
351 proteins are encoded by the *Firmicutes* (the remaining one is encoded by an unknown bacterium),
352 as are 101 out of 149 SLH-domain containing proteins. Both are components of cellulosomes,
353 multi-enzyme complexes involved in fibre degradation, which are encoded by some members of
354 the *Clostridiales* family.

355 There are 1707 *Bacteroidetes* genomes in the RUGs, and additionally we have a whole
356 genome of *Prevotella copri* from the Nanopore assembly. These 1708 genomes were subject to
357 prediction of polysaccharide utilisation loci (PUL) using our pipeline PULpy⁴³. Of the 1708
358 genomes, 1469 are predicted to have at least one PUL, and in total there are 15,629 separate loci
359 involving 88260 proteins. The highest number of PUL per genome are 52 PUL for RUG13980 and
360 50 for RUG10279, both labelled uncultured *Prevotellaceae*; both of these genomes are closely
361 related to *Prevotella multisaccharivorax*, known to be able to utilise multiple carbohydrate
362 substrates⁴⁴.

363 Discussion

364 The rumen microbiome has a crucial role in food security and climate change. Recent studies have
365 released more than 1300 draft and complete rumen genomes. We add 4941 near-complete, de-
366 replicated metagenome-assembled genomes to these 1300 existing rumen genomes^{9,20,21}. By
367 combining our dataset with publicly available genomes, we assembled a “rumen superset” of 5845
368 public bacterial and archaeal genomes. This set contains 2690 unique species-level bins (95% ANI)
369 and 2078 of these 2690 putative species are RUG2 genomes discovered in this study. The RUG2
370 dataset and the rumen superset bring read classification rates up to 70% for our own data, and 45-
371 55% for other rumen metagenome datasets (some from non-cattle ruminants). The remaining
372 reads are likely to derive from low-abundance bacterial and archaeal species, difficult-to-assemble
373 genomes, and the fungal, protozoan and viral genomes that are not part of this study.

374 We estimate that we have discovered 65% of rumen species in our samples, representing 4
375 important beef cattle breeds, which suggests that there are over 1000 species yet to be sequenced
376 and assembled. Given that average read classification rates dip from 70% in our own data, to 50%
377 in the Rubino *et al* cattle data (Limousin x Friesian cross)²⁵ and Shi *et al* sheep data²⁶, and 45% in
378 moose²¹, there are many species yet to be discovered, and there are likely to be species- and
379 breed- specific rumen microbiomes. We note the high abundance of unclassified *Proteobacteria* in
380 our data, and in the rumen census data, and suggest that these may form part of a core rumen
381 microbiome. Our dataset contains 74 proteobacterial genomes, and we present one near-
382 complete genome in a single contig.

383 We apply our dataset to re-analyse data on methane emissions in sheep that was
384 published in 2014²⁶. Using a combined database of rumen microbial genomes we reveal
385 fundamental and large-scale changes in rumen metagenomic abundance between LME and HME
386 sheep. These differences occur at almost every taxonomic level tested, and the rumen superset
387 database enables us to analyse these data at previously unprecedented resolution. Whilst species-
388 and strain- level metagenomic data must always be interpreted with care – there remains a
389 possibility that strains that are not present in the database are driving the observed differences –
390 nonetheless we observe consistent patterns suggesting large changes in abundance for numerous
391 species. Our analysis confirms and strengthens subsequent studies of methane emissions in
392 sheep^{15,28} by identifying specific strains of bacteria and archaea involved and revealing their
393 genome sequence. Our analysis confirms that there is a complex relationship between archaeal
394 abundance and methane emissions, with archaeal species and strains both positively and
395 negatively associated with methane emissions. These insights into metagenomic species- and
396 strain- level aspects of methane emissions will form the basis of future studies.

397 The main rumen functions rely on the activity of proteins encoded in rumen microbe
398 genomes, and as researchers produce more proteomic data, it is vital that protein reference
399 datasets are available. We present the largest redundant and non-redundant rumen microbial
400 protein prediction datasets to date, and provide rich annotation using public protein, pathway and
401 enzyme databases. This resource will enable researchers to predict the function of each protein,
402 and better assess the functional consequences of changes in the rumen proteome.

403 Going forwards, it is vital that more rumen bacteria and archaea are brought into culture,
404 to better enable studying the functions of the rumen microbiome. In particular, if we are to design
405 rational interventions to manipulate rumen feed-conversion or methane emissions, we will need
406 to understand microbiome structure, which substrates are utilized by microbiota, and how they
407 interact with one another and the ruminant host. Sequencing and assembling rumen microbial
408 genomes is an important step towards improved culture collections and future manipulation of
409 the rumen microbiome for human benefit.

410 **Competing interests.**

411 The authors declare no competing interests.

412

413 **Acknowledgements**

414 The Rowett Institute and SRUC are core funded by the Rural and Environment Science and
415 Analytical Services Division (RESAS) of the Scottish Government. The Roslin Institute forms part of

416 the Royal (Dick) School of Veterinary Studies, University of Edinburgh. This project was supported
417 by the Biotechnology and Biological Sciences Research Council (BBSRC; BB/N016742/1,
418 BB/N01720X/1), including institute strategic programme and national capability awards to The
419 Roslin Institute (BBSRC: BB/P013759/1, BB/P013732/1, BB/J004235/1, BB/J004243/1); and by the
420 Scottish Government as part of the 2016–2021 commission.

421 **Data Availability**

422 Raw sequence reads for all samples are available under ENA project PRJEB31266, except for 10572
423 which are available under PRJEB21624. All metagenomic assemblies and RUGs are in the process
424 of being deposited in ENA under accession PRJEB31266. All protein predictions, clusters and
425 annotation are available at DOI: 10.7488/ds/2470.

426

427 **Code Availability**

428 Comparative genomic analysis was carried out using MAGpy¹⁰
429 (<https://github.com/WatsonLab/MAGpy>); analysis of PUL was carried out using PULpy⁴³
430 (<https://github.com/WatsonLab/PULpy>); analysis of indels in nanopore data was carried out using
431 IDEEL (<https://github.com/mw55309/ideel>)

432

433 **Figure 1** Phylogenetic tree of 4941 rumen uncultured genomes (RUGs) from the cow rumen,
434 additionally incorporating rumen genomes from the Hungate collection. The tree was produced
435 from concatenated protein sequences using PhyloPhlAn¹³ and subsequently drawn using
436 GraPhlAn⁴⁵. Labels show Hungate genome names, chosen to be informative but not overlap.

437 **Figure 2** A comparison of the 4941 RUG dataset with the Hungate collection (A) and our previously
438 published data from Stewart et al (B). Black line is average percentage protein identity with closest
439 match (right-hand y-axis), and blue dots are mash distance ($k=100,000$) between RUG and the
440 closest match (a measure of dissimilarity between two DNA sequences). As expected, a high
441 protein identity relates to a low mash distance, and vice versa. The RUGs are sorted independently
442 for figures A and B, by average protein identity. There is a clear inflection point in Figure 5B,
443 roughly half way along the x-axis, where the protein identity dips below 90% and the Mash
444 distance rises, neatly demonstrating the novelty represented by our new, larger dataset

445 **Figure 3** A comparison of Illumina and Nanopore metagenomic assembly statistics. The coloured
446 histograms show the distribution of statistics for 282 Illumina assemblies, and the single Nanopore
447 assembly is highlighted. A) N50; B) total length of the assembly; and C) length of the longest
448 contig. As can be seen, the Nanopore assembly N50 of 268kb is over 56-times longer than the
449 average Illumina assembly (4.7kb); whilst the Illumina assemblies are often longer (average
450 600Mb), the Nanopore assembly (at 178Mb in length) is not the shortest of the assemblies we
451 produced; and the Nanopore assembly produced the longest contig at 3.8Mb, seven-times longer
452 than the average for the Illumina assemblies (479kb) and 2.74-times longer than the longest single
453 Illumina contig (1.38Mb – one of thirteen contigs from the 99.19% complete “uncultured
454 *Bacteroidia* bacterium RUG14538”). In terms of a direct comparison, the Illumina-only assembly of
455 the same sample has an N50 of 12.2Kb, a total length of 247Mb and a longest contig of 358Kb

456 **Figure 4** Maximum percentage identity between CAZyme-predicted proteins from the RUGs and
457 the CAZy database. GH=glycoside hydrolase (n=235,001); GT=glycosyl transferase (n=120,494);
458 PL=polysaccharide lyase (n=6,834); CE=carbohydrate esterase (n=55,523); AA=auxiliary activities;
459 CBM=carbohydrate binding module (n=23,928); SLH=S-layer homology domain (n=150);
460 cohesin=cohesin domain (n=80). Centre line shows the median value; box shows the interquartile
461 range; whiskers extend to the most extreme data point which is no more than 1.5 times the
462 interquartile range from the box

463 **Figure 5** Top: Total number of proteins for 12 phyla and the group of unknown bacteria; Middle:
464 percentage of the proteome predicted to be CAZymes; Bottom: distribution of eight CAZyme
465 classes as a proportion of the total number of predicted CAZymes. GH=glycoside hydrolase;
466 GT=glycosyl transferase; PL=polysaccharide lyase; CE=carbohydrate esterase; AA=auxiliary
467 activities; CB=carbohydrate binding module; SL=S-layer homology domain; co=cohesin domain

468 References

- 469 1. Hess, M. *et al.* Metagenomic discovery of biomass-degrading genes and genomes from cow rumen.
470 *Science (80-.)*. **331**, 463–467 (2011).
- 471 2. Cowan, D. A. *et al.* Metagenomics, gene discovery and the ideal biocatalyst. *Biochem. Soc. Trans.* **32**,
472 298–302 (2004).
- 473 3. Roumpeka, D. D., Wallace, R. J., Escalettes, F., Fotheringham, I. & Watson, M. A Review of
474 Bioinformatics Tools for Bio-Prospecting from Metagenomic Sequence Data. *Front. Genet.* **8**, 23
475 (2017).
- 476 4. Huws, S. A. *et al.* Addressing Global Ruminant Agricultural Challenges Through Understanding the
477 Rumen Microbiome: Past, Present, and Future. *Front. Microbiol.* **9**, 2161 (2018).
- 478 5. Gerber, P. J. & Food and Agriculture Organization of the United Nations. *Tackling climate change
479 through livestock : a global assessment of emissions and mitigation opportunities*. (Rome: Food and
480 Agriculture Organization of the United Nations (FAO)., 2013).
- 481 6. Wallace, R. J. *et al.* The rumen microbial metagenome associated with high methane production in
482 cattle. *BMC Genomics* **16**, 839 (2015).
- 483 7. Roehe, R. *et al.* Bovine Host Genetic Variation Influences Rumen Microbial Methane Production with
484 Best Selection Criterion for Low Methane Emitting and Efficiently Feed Converting Hosts Based on
485 Metagenomic Gene Abundance. *PLoS Genet.* **12**, e1005846 (2016).
- 486 8. Stewart, R. D. *et al.* Assembly of 913 microbial genomes from metagenomic sequencing of the cow
487 rumen. *Nat. Commun.* **9**, 870 (2018).
- 488 9. Seshadri, R. *et al.* Cultivation and sequencing of rumen microbiome members from the
489 Hungate1000 Collection. *Nat. Biotechnol.* **36**, 359–367 (2018).
- 490 10. Stewart, R. D., Auffret, M., Snelling, T. J., Roehe, R. & Watson, M. MAGpy: a reproducible pipeline
491 for the downstream analysis of metagenome-assembled genomes (MAGs). *Bioinformatics* (2018).
492 doi:10.1093/bioinformatics/bty905
- 493 11. Brown, C. T. & Irber, L. sourmash: a library for MinHash sketching of DNA. *J. Open Source Softw.* **1**,
494 27 (2016).
- 495 12. Buchfink, B., Xie, C. & Huson, D. H. Fast and sensitive protein alignment using DIAMOND. *Nat.*
496 *Methods* **12**, 59–60 (2015).
- 497 13. Segata, N., Börnigen, D., Morgan, X. C. & Huttenhower, C. PhyloPhlAn is a new method for improved
498 phylogenetic and taxonomic placement of microbes. *Nat. Commun.* **4**, 2304 (2013).

- 499 14. Parks, D. H. *et al.* A standardized bacterial taxonomy based on genome phylogeny substantially
500 revises the tree of life. *Nat. Biotechnol.* **36**, 996 (2018).
- 501 15. Kamke, J. *et al.* Rumen metagenome and metatranscriptome analyses of low methane yield sheep
502 reveals a Sharpea-enriched microbiome characterised by lactic acid formation and utilisation.
503 *Microbiome* **4**, 56 (2016).
- 504 16. Bowers, R. M. *et al.* Minimum information about a single amplified genome (MISAG) and a
505 metagenome-assembled genome (MIMAG) of bacteria and archaea. *Nat. Biotechnol.* **35**, 725–731
506 (2017).
- 507 17. Parks, D. H. *et al.* Recovery of nearly 8,000 metagenome-assembled genomes substantially expands
508 the tree of life. *Nat. Microbiol.* **2**, 1533–1542 (2017).
- 509 18. The UniProt Consortium. UniProt: the universal protein knowledgebase. *Nucleic Acids Res.* **45**,
510 D158–D169 (2017).
- 511 19. Quast, C. *et al.* The SILVA ribosomal RNA gene database project: improved data processing and web-
512 based tools. *Nucleic Acids Res.* **41**, D590–6 (2013).
- 513 20. Solden, L. M. *et al.* Interspecies cross-feeding orchestrates carbon degradation in the rumen
514 ecosystem. *Nat. Microbiol.* **3**, 1274–1284 (2018).
- 515 21. Svartström, O. *et al.* Ninety-nine de novo assembled genomes from the moose (*Alces alces*) rumen
516 microbiome provide new insights into microbial plant biomass degradation. *ISME J.* **11**, 2538–2551
517 (2017).
- 518 22. Chao, A. Nonparametric Estimation of the Number of Classes in a Population. *Scandinavian Journal*
519 *of Statistics* **11**, 265–270 (1984).
- 520 23. Auffret, M. D. *et al.* Identification, Comparison, and Validation of Robust Rumen Microbial
521 Biomarkers for Methane Emissions Using Diverse *Bos Taurus* Breeds and Basal Diets. *Front.*
522 *Microbiol.* **8**, 2642 (2018).
- 523 24. Auffret, M. D. *et al.* The rumen microbiome as a reservoir of antimicrobial resistance and
524 pathogenicity genes is directly affected by diet in beef cattle. *Microbiome* **5**, 159 (2017).
- 525 25. Rubino, F. *et al.* Divergent functional isoforms drive niche specialisation for nutrient acquisition and
526 use in rumen microbiome. *ISME J.* **11**, 932–944 (2017).
- 527 26. Shi, W. *et al.* Methane yield phenotypes linked to differential gene expression in the sheep rumen
528 microbiome. *Genome Res.* **24**, 1517–1525 (2014).
- 529 27. Pasolli, E. *et al.* Extensive Unexplored Human Microbiome Diversity Revealed by Over 150,000
530 Genomes from Metagenomes Spanning Age, Geography, and Lifestyle. *Cell* **176**, 649–662.e20 (2019).
- 531 28. Kittelmann, S. *et al.* Two Different Bacterial Community Types Are Linked with the Low-Methane
532 Emission Trait in Sheep. *PLoS One* **9**, e103171 (2014).
- 533 29. Henderson, G. *et al.* Rumen microbial community composition varies with diet and host, but a core
534 microbiome is found across a wide geographical range. *Sci. Rep.* **5**, 14567 (2015).
- 535 30. Risse, J. *et al.* A single chromosome assembly of *Bacteroides fragilis* strain BE1 from Illumina and
536 MinION nanopore sequencing data. *Gigascience* **4**, 60 (2015).
- 537 31. Ip, C. L. C. *et al.* MinION Analysis and Reference Consortium: Phase 1 data release and analysis.
538 *F1000Research* **4**, 1075 (2015).
- 539 32. Koren, S. *et al.* Canu: scalable and accurate long-read assembly via adaptive *k*-mer weighting and
540 repeat separation. *Genome Res.* **27**, 722–736 (2017).

- 541 33. Watson, M. & Warr, A. Errors in long-read assemblies can critically affect protein prediction. *Nat.*
542 *Biotechnol.* **37**, 124–126 (2019).
- 543 34. Hyatt, D. *et al.* Prodigal: prokaryotic gene recognition and translation initiation site identification.
544 *BMC Bioinformatics* **11**, 119 (2010).
- 545 35. Parks, D. H., Imelfort, M., Skennerton, C. T., Hugenholtz, P. & Tyson, G. W. CheckM: assessing the
546 quality of microbial genomes recovered from isolates, single cells, and metagenomes. *Genome Res.*
547 **25**, 1043–55 (2015).
- 548 36. Siguier, P., Goubeyre, E. & Chandler, M. Bacterial insertion sequences: their genomic impact and
549 diversity. *FEMS Microbiol. Rev.* **38**, 865–891 (2014).
- 550 37. Suzek, B. E. *et al.* UniRef clusters: a comprehensive and scalable alternative for improving sequence
551 similarity searches. *Bioinformatics* **31**, 926–32 (2015).
- 552 38. Kanehisa, M. & Goto, S. KEGG: kyoto encyclopedia of genes and genomes. *Nucleic Acids Res.* **28**, 27–
553 30 (2000).
- 554 39. Cantarel, B. L. *et al.* The Carbohydrate-Active EnZymes database (CAZy): an expert resource for
555 Glycogenomics. *Nucleic Acids Res.* **37**, D233-8 (2009).
- 556 40. Zhang, H. *et al.* dbCAN2: a meta server for automated carbohydrate-active enzyme annotation.
557 *Nucleic Acids Res.* **46**, W95–W101 (2018).
- 558 41. Hook, S. E. *et al.* Impact of subacute ruminal acidosis (SARA) adaptation and recovery on the density
559 and diversity of bacteria in the rumen of dairy cows. *FEMS Microbiol. Ecol.* **78**, 275–284 (2011).
- 560 42. Kasparovska, J. *et al.* Effects of Isoflavone-Enriched Feed on the Rumen Microbiota in Dairy Cows.
561 *PLoS One* **11**, e0154642 (2016).
- 562 43. Stewart, R. D., Auffret, M., Roehe, R. & Watson, M. Open prediction of polysaccharide utilisation loci
563 (PUL) in 5414 public Bacteroidetes genomes using PULpy. *bioRxiv* 421024 (2018).
564 doi:10.1101/421024
- 565 44. Sakamoto, M., Umeda, M., Ishikawa, I. & Benno, Y. *Prevotella multisaccharivorax* sp. nov., isolated
566 from human subgingival plaque. *Int. J. Syst. Evol. Microbiol.* **55**, 1839–1843 (2005).
- 567 45. Asnicar, F., Weingart, G., Tickle, T. L., Huttenhower, C. & Segata, N. Compact graphical
568 representation of phylogenetic data and metadata with GraPhlAn. *PeerJ* **3**, e1029 (2015).

569 **Online Methods**

570 *Metagenomic samples.* Animal experiments were conducted at the Beef and Sheep Research
571 Centre of Scotland's Rural College (SRUC). The experiment was approved by the Animal
572 Experiment Committee of SRUC and was conducted in accordance with the requirements of the
573 UK Animals (Scientific Procedures) Act 1986.

574 The data were obtained from three cross breeds: Aberdeen Angus, Limousin and Charolais and
575 one pure breed: Luing (Supplementary Data 4). As previously described, the animals were
576 slaughtered in a commercial abattoir where two post-mortem digesta samples (approximately 50
577 mL) were taken immediately after the rumen was opened to be drained^{46,47}. DNA extraction was
578 carried out following the protocol of Yu and Morrison⁴⁸ and based on repeated bead beating plus
579 column filtration. Illumina TruSeq libraries were prepared from genomic DNA and sequenced on
580 an Illumina HiSeq 4000 by Edinburgh Genomics.

581 We experienced severe problems when using the MinION for rumen microbiome DNA when
582 following standard, recommended protocols, and we hope our adapted methods will be of

583 assistance to others. We found that the DNA did not meet the recommended purity for Nanopore
584 library prep following extraction, according to Nanodrop O/D ratios. RNase treatment using
585 Riboshredder and clean up with methods such as AMPure XP beads were sufficient to obtain O/D
586 ratios within the recommended range, but DNA from these methods typically led to poor or failed
587 sequencing runs. Successful clean-up reaching recommended O/D ratios and leading to successful
588 sequencing runs was carried out using RNase treatment with Riboshredder and a phenol-
589 chloroform purification. 1D libraries were prepared starting with 2ug of DNA per library following
590 Oxford Nanopore's SQK-LSK108 1D ligation protocol with modifications. The incubation in the end
591 prep stage of the protocol was extended to 30 minutes at 20°C and 30 minutes at 65°C and the
592 incubation in the ligation stage was extended to 15 minutes at room temperature. The optional
593 FFPE repair step was also carried out. Three sequencing runs were carried out using FLOMIN-106
594 flow cells on a MinION MK1b housed in the Watson lab, U. Edinburgh.

595 *Bioinformatics – metagenomic assembly and binning.* In total, 282 samples were sequenced for
596 this study generating between 24 and 140 million 150bp paired-end reads per sample. The
597 samples were sequenced in five batches of 48 samples and one batch of 42 samples (this 42-
598 sample batch was the sole basis of Stewart *et al*). An additional sample was used for Hi-C
599 sequencing in Stewart *et al*⁸, and the metagenome-assembled genomes from that sample are
600 included in the de-replicated set.

601 Unless otherwise stated, all parameters used were the default. Each sample was assembled and
602 binned individually using coverage and content as previously described⁸. Briefly, each sample was
603 assembled using *idba_ud*⁴⁹ (v1.1.3) with the options `--num_threads 16 --pre_correction --`
604 `min_contig 300`. BWA MEM⁵⁰ (v0.7.15) was used to map reads back to the filtered assembly and
605 *Samtools*⁵¹ (v 1.3.1) was used to convert to BAM format. Script
606 *jgi_summarize_bam_contig_depths* from the *MetaBAT2*⁵² (v 2.11.1) package was used to calculate
607 coverage from the resulting BAM files. A co-assembly was also produced for each of the 6 batches
608 of samples using *MEGAHIT*⁵³ (v1.1.1) with options `--kmin-1pass, -m 60e +10, --k-list`
609 `27,37,47,57,67,77,87, --min-contig-len 1000, -t 16`.

610 Metagenomic binning was applied to both single-sample assemblies and the coassemblies using
611 *MetaBAT2*⁵², with options `--minContigLength 2000, --minContigDepth 2`. Single-sample binning
612 produced a total of 37153 bins, and co-assembly binning produced a further 23335. All 60743 bins
613 were aggregated and then dereplicated using *dRep*⁵⁴ (v1.1.2). The *dRep* dereplication workflow
614 was used with options `dereplicate_wf -p 16 -comp 80 -con 10 -str 100 -strW 0`. Thus, in pre-
615 filtering, only bins assessed by *CheckM* (v1.0.5) as having both completeness $\geq 80\%$, and
616 contamination $\leq 10\%$ were retained for pairwise dereplication comparison (n=10586). Bin scores
617 were given as $\text{completeness} - 5 * \text{contamination} + 0.5 * \log(N50)$, and only the highest scoring RUG
618 from each secondary cluster was retained in the dereplicated set. For our dataset, 4941
619 dereplicated RUGs were obtained.

620 Note that we operate a continuous de-replication workflow. Therefore all 913 of the RUGs (both
621 *MetaBAT2* and Hi-C) we previously published have been merged with the newer RUGs and de-
622 replicated. Therefore, whilst some of the previously published RUGs exist in the newer dataset
623 published here, many have been replaced by newer RUGs of higher quality.

624 Supplementary Data 5 is the average depth for each genome in each sample as calculated by script
625 *jgi_summarize_bam_contig_depths* from the *MetaBAT2*⁵² (v 2.11.1) package.

626 *Bioinformatics – metagenomic assignment.* The output of metagenomic binning is simply a set of
627 DNA FASTA files containing putative genomes. These were all assessed for completeness and
628 contamination using CheckM³⁵ (v1.0.5). The 4941 best bins were analysed using MAGpy¹⁰, a
629 Snakemake⁵⁵ pipeline that runs a series of analyses on the bins, including: CheckM (v1.0.5);
630 prodigal³⁴ (v2.6.3) protein prediction; Pfam_Scan⁵⁶ (v1.6); DIMAOND¹² (v0.9.22.123) search against
631 UniProt TrEMBL; PhyloPhlAn¹³ (v0.99); and Sourmash (v2.0.0) search against all public bacterial
632 genomes. The MAGpy results were used to produce a putative taxonomic assignment for each bin
633 as follows:

- 634 • If the proportion of proteins assigned to a species is ≥ 0.9 and the average amino acid identity
635 ≥ 0.95 , assign to species based on DIAMOND results; else
- 636 • If sourmash score is ≥ 0.8 assign to species based on Sourmash results; else
- 637 • If PhyloPhlAn probability is high and the level is genus or species, then assign taxonomy based
638 on PhyloPhlAn results; else
- 639 • If the proportion of proteins assigned to a genus is ≥ 0.9 and the average amino acid identity
640 ≥ 0.9 , assign to genus based on DIAMOND results; else
- 641 • If PhyloPhlAn probability is high or medium and the level is genus, then assign to genus based
642 on PhyloPhlAn results; else
- 643 • If PhyloPhlAn probability is high or medium and the level is family, then assign to family based
644 on PhyloPhlAn results; else
- 645 • If the proportion of proteins assigned to a family is ≥ 0.8 and the average amino acid identity
646 ≥ 0.6 , assign to family based on DIAMOND results; else
- 647 • If PhyloPhlAn probability is high or medium and the level is order, then assign to order based
648 on PhyloPhlAn results; else
- 649 • If the proportion of proteins assigned to an order is ≥ 0.6 and the average amino acid identity
650 ≥ 0.6 , assign to order based on DIAMOND results; else
- 651 • If PhyloPhlAn probability is high or medium and the level is class, then assign to class based on
652 PhyloPhlAn results; else
- 653 • If PhyloPhlAn probability is high or medium and the level is phylum, then assign to phylum
654 based on PhyloPhlAn results; else
- 655 • Assign taxonomy based on CheckM lineage

656 Importantly, at this stage, these are only putative taxonomic assignments. Using these labels, a
657 phylogenetic tree consisting of the RUGs and genomes from the Hungate collection, produced
658 from concatenated protein subsequences using PhyloPhlAn¹³ (v0.99), was visually inspected using
659 FigTree (v.1.4.3), iTol⁵⁷ (v4.3.1) and GraPhlAn⁴⁵ (v0.9.7). Annotations were improved where they
660 could be - for example where MAGpy had only assigned a taxonomy at the Genus level, but that
661 genome clustered closely with a Hungate 1000 genome annotated at the species level, the
662 annotation was updated. The tree was also re-rooted manually at the Bacteria/Archaea branch
663 using FigTree.

664 *Bioinformatics – genome quality and comparative genomics.* Genome completeness and
665 contamination was assessed using CheckM (v1.0.5) (see above). tRNA genes were annotated using
666 tRNAscan-SE (v2.0.0) and 16S rRNA genes predicted using barrnap (v0.9). Whole-genome
667 alignments were calculated with MUMmer⁵⁸ (v 3.23) using promoter to find matches between
668 genomes. Average nucleotide identity was calculated using FastANI (v1.1). The RUGs were

669 compared to the Hungate collection and our previous dataset using DIAMOND blastp (v0.9.22.123)
670 and MASH⁵⁹ (v 2.0) with parameters -k 21 -s 100000.

671 The rumen superset was de-replicated using dRep as above, with -sa 0.99 for de-replication at
672 99% ANI and -sa 0.95 for de-replication at 95% ANI. Overlaps between sets were plotted with
673 UpSetR⁶⁰ (v1.3.3). Read classification rates were calculated using Kraken⁶¹ (v0.10.5) with
674 parameters --fastq-input --gzip-compressed --preload --paired.

675 *Bioinformatics – analysis of sheep methane data.* Reads from the low and high methane samples
676 from Shi *et al* were assigned to different taxonomic levels of the rumen superset database using
677 Kraken, as described above. The resulting read counts data was used as input into DESeq2
678 (v1.22.2) for differential analysis. Principal components analysis plots were created using the
679 plotPCA() function within DESeq2, and heatmaps were created using the heatmap.2() function
680 within the gplots package (v3.0.1.1). For strain-level analysis, reads from the low and high
681 methane samples from Shi *et al* were aligned directly to the rumen superset database using BWA
682 MEM (v0.7.15) and the number of primary alignments to each genome was used as input to
683 DESeq2. P-values for all comparisons were calculated by DESeq2 and adjusted for multiple
684 testing⁶².

685 *Bioinformatics – rumen census analysis.* The average and total depth for each genome in each
686 dataset (Supplementary Data 5) was used as a proxy for abundance in the dataset(s). Kraken (as
687 described above) was used with the rumen superset database to calculate the read abundance of
688 *Proteobacteria* in all samples.

689 *Bioinformatics – assembly and analysis of Nanopore sequence data.* The Nanopore reads were
690 extracted and QC-ed using poRe^{63,64} (v 0.24), and assembled using Canu³² (v1.8) with default
691 settings and genomeSize=150Mb. The resulting assembly was analysed using MAGpy¹⁰. The raw
692 assembly was corrected using both Nanopolish⁶⁵ (v 0.10.2) and Racon⁶⁶ (v 1.3.1) using Illumina
693 data aligned to the Nanopore assembly with Minimap2 (v 2.12) using short read settings (-x sr).
694 Query vs subject length data were extracted and plotted using ideel
695 (<https://github.com/mw55309/ideel>). Whole-genome alignments were calculated using
696 MUMmer79 (v 3.23) using promoter to find matches between genomes. The three complete
697 nanopore bacterial genomes and their Illumina counterparts were annotated using Prokka⁶⁷ (v
698 1.13.3). The Nanopore assembly was created with a minimum contig length of 1kb, therefore the
699 Illumina assemblies were similarly limited prior to comparison.

700 *Bioinformatics – proteome analysis.* Proteins were predicted using Prodigal (v2.6.3) with option -p
701 meta. Using DIAMOND, each protein was searched against KEGG (downloaded on Sept 15th 2018),
702 UniRef100, UniRef90 and UniRef50 (downloaded Oct 3rd 2018), and CAZy (dbCAN2 version,
703 31/07/2018). The protein predictions were clustered by CD-HIT⁶⁸ (v4.7) at 100%, 90% and 50%
704 identity, mirroring similar methods at UniRef.

705 All protein predictions were searched against the CAZy database using dbCAN2⁴⁰ and HMMER⁶⁹
706 (v3.1b2), and polysaccharide utilisation loci (PUL) were predicted for *Bacteroidetes* RUGs using
707 PULpy⁴³.

708 **References**

709 46. Duthie, C.-A. *et al.* Impact of adding nitrate or increasing the lipid content of two contrasting diets
710 on blood methaemoglobin and performance of two breeds of finishing beef steers. *animal* **10**, 786–

- 711 795 (2016).
- 712 47. Duthie, C.-A. *et al.* The impact of divergent breed types and diets on methane emissions, rumen
713 characteristics and performance of finishing beef cattle. *animal* **11**, 1762–1771 (2017).
- 714 48. Yu, Z. & Morrison, M. Improved extraction of PCR-quality community DNA from digesta and fecal
715 samples. *Biotechniques* **36**, 808–812 (2004).
- 716 49. Peng, Y., Leung, H. C. M., Yiu, S. M. & Chin, F. Y. L. IDBA-UD: a de novo assembler for single-cell and
717 metagenomic sequencing data with highly uneven depth. *Bioinformatics* **28**, 1420–1428 (2012).
- 718 50. Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. **3** (2013).
- 719 51. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078–9 (2009).
- 720 52. Kang, D. D., Froula, J., Egan, R. & Wang, Z. MetaBAT, an efficient tool for accurately reconstructing
721 single genomes from complex microbial communities. *PeerJ* **3**, e1165 (2015).
- 722 53. Li, D., Liu, C.-M., Luo, R., Sadakane, K. & Lam, T.-W. MEGAHIT: an ultra-fast single-node solution for
723 large and complex metagenomics assembly via succinct de Bruijn graph. *Bioinformatics* **31**, 1674–
724 1676 (2015).
- 725 54. Olm, M. R., Brown, C. T., Brooks, B. & Banfield, J. F. dRep: a tool for fast and accurate genomic
726 comparisons that enables improved genome recovery from metagenomes through de-replication.
727 *ISME J.* **11**, 2864–2868 (2017).
- 728 55. Koster, J. & Rahmann, S. Snakemake--a scalable bioinformatics workflow engine. *Bioinformatics* **28**,
729 2520–2522 (2012).
- 730 56. Finn, R. D. *et al.* Pfam: the protein families database. *Nucleic Acids Res.* **42**, D222–D230 (2014).
- 731 57. Letunic, I. & Bork, P. Interactive tree of life (iTOL) v3: an online tool for the display and annotation of
732 phylogenetic and other trees. *Nucleic Acids Res.* **44**, W242–W245 (2016).
- 733 58. Kurtz, S. *et al.* Versatile and open software for comparing large genomes. *Genome Biol.* **5**, R12
734 (2004).
- 735 59. Ondov, B. D. *et al.* Mash: fast genome and metagenome distance estimation using MinHash.
736 *Genome Biol.* **17**, 132 (2016).
- 737 60. Conway, J. R., Lex, A. & Gehlenborg, N. UpSetR: an R package for the visualization of intersecting
738 sets and their properties. *Bioinformatics* **33**, 2938–2940 (2017).
- 739 61. Wood, D. E. & Salzberg, S. L. Kraken: ultrafast metagenomic sequence classification using exact
740 alignments. *Genome Biol.* **15**, R46 (2014).
- 741 62. Benjamini, Y. & Hochberg, Y. *Controlling The False Discovery Rate - A Practical And Powerful*
742 *Approach To Multiple Testing.* *J. Royal Statist. Soc., Series B* **57**, (1995).
- 743 63. Watson, M. *et al.* poRe: an R package for the visualization and analysis of nanopore sequencing
744 data. *Bioinformatics* **31**, 114–5 (2015).
- 745 64. Stewart, R. D. & Watson, M. poRe GUIs for parallel and real-time processing of MinION sequence
746 data. *Bioinformatics* **33**, 2207–2208 (2017).
- 747 65. Loman, N. J., Quick, J. & Simpson, J. T. A complete bacterial genome assembled de novo using only
748 nanopore sequencing data. *Nat. Methods* **12**, 733–735 (2015).
- 749 66. Vaser, R., Sović, I., Nagarajan, N. & Šikić, M. Fast and accurate de novo genome assembly from long
750 uncorrected reads. *Genome Res.* **27**, 737–746 (2017).
- 751 67. Seemann, T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics* **30**, 2068–2069 (2014).

752 68. Fu, L., Niu, B., Zhu, Z., Wu, S. & Li, W. CD-HIT: accelerated for clustering the next-generation
753 sequencing data. *Bioinformatics* **28**, 3150–3152 (2012).

754 69. Mistry, J., Finn, R. D., Eddy, S. R., Bateman, A. & Punta, M. Challenges in homology search: HMMER3
755 and convergent evolution of coiled-coil regions. *Nucleic Acids Res.* **41**, e121–e121 (2013).

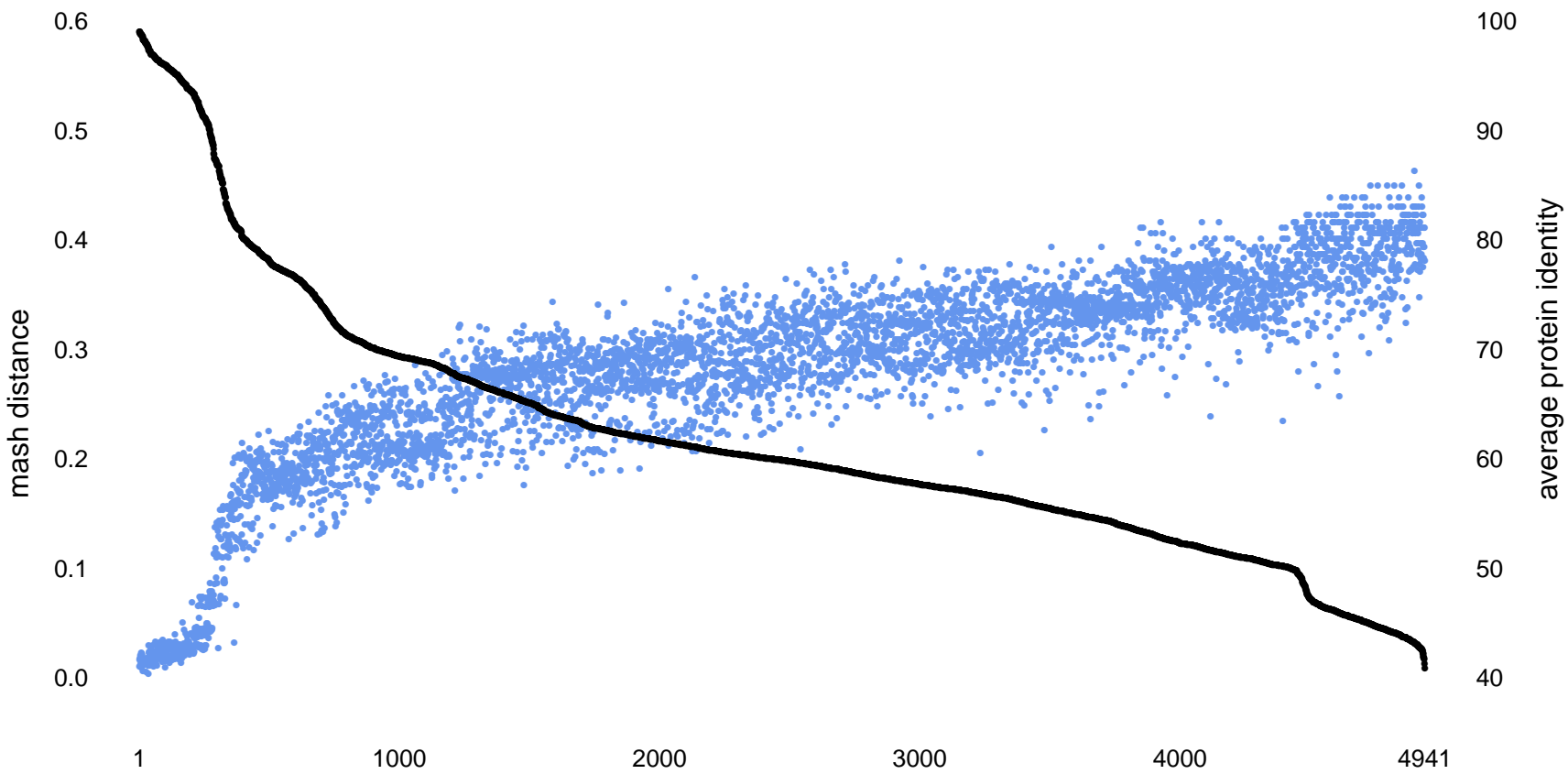
756

757 **Author contributions**

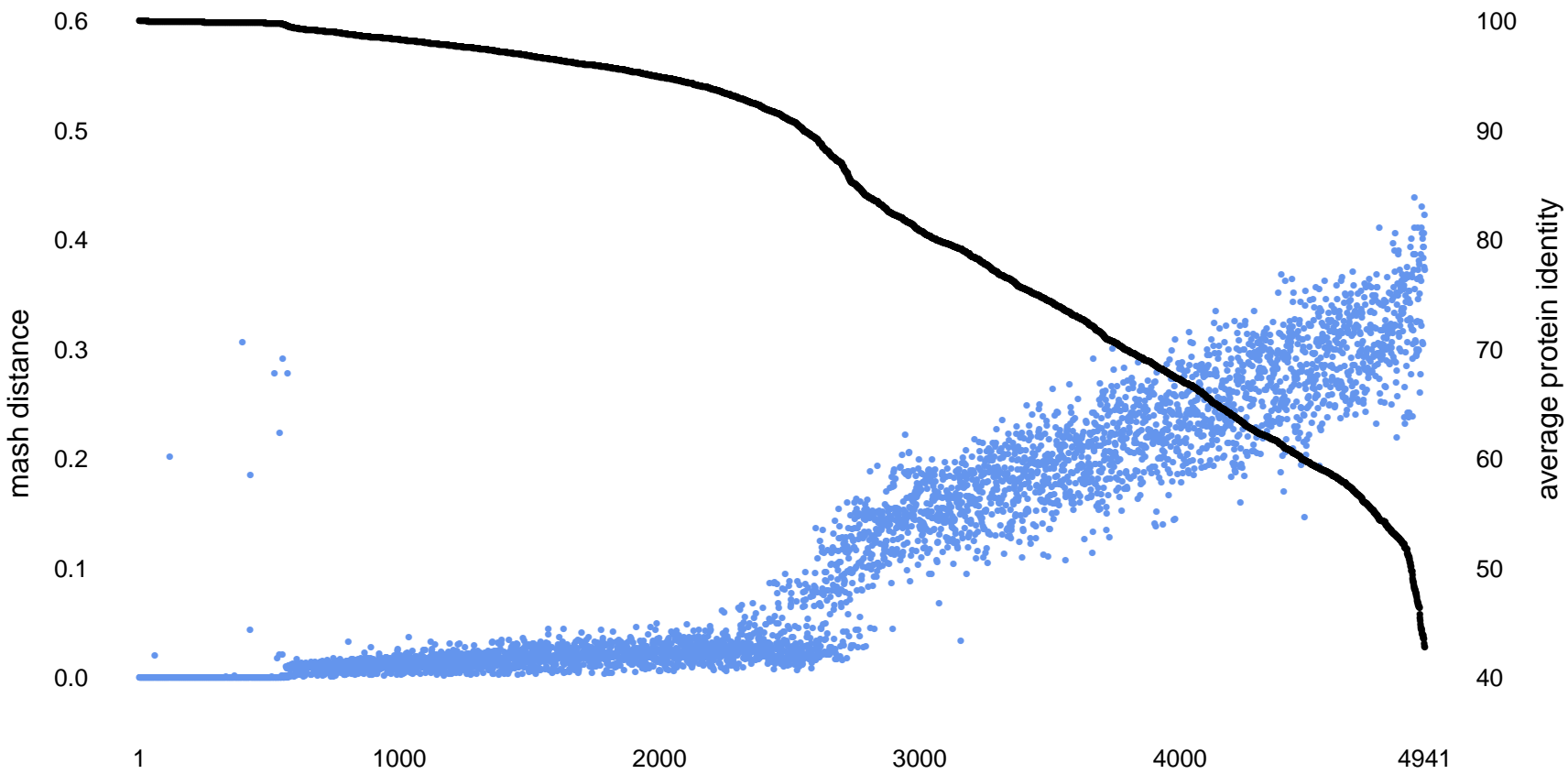
758 MW, RR and AWW conceived of the study and supervised the project. RDS and MW carried out all
759 bioinformatics work on the Illumina data and AW carried out all bioinformatics work on the
760 Nanopore data. MA carried out all laboratory work, except for DNA clean up, Nanopore library
761 prep and Nanopore sequencing which was done by AW. All of the authors contributed ideas, co-
762 wrote the paper, and reviewed and approved the manuscript.

763

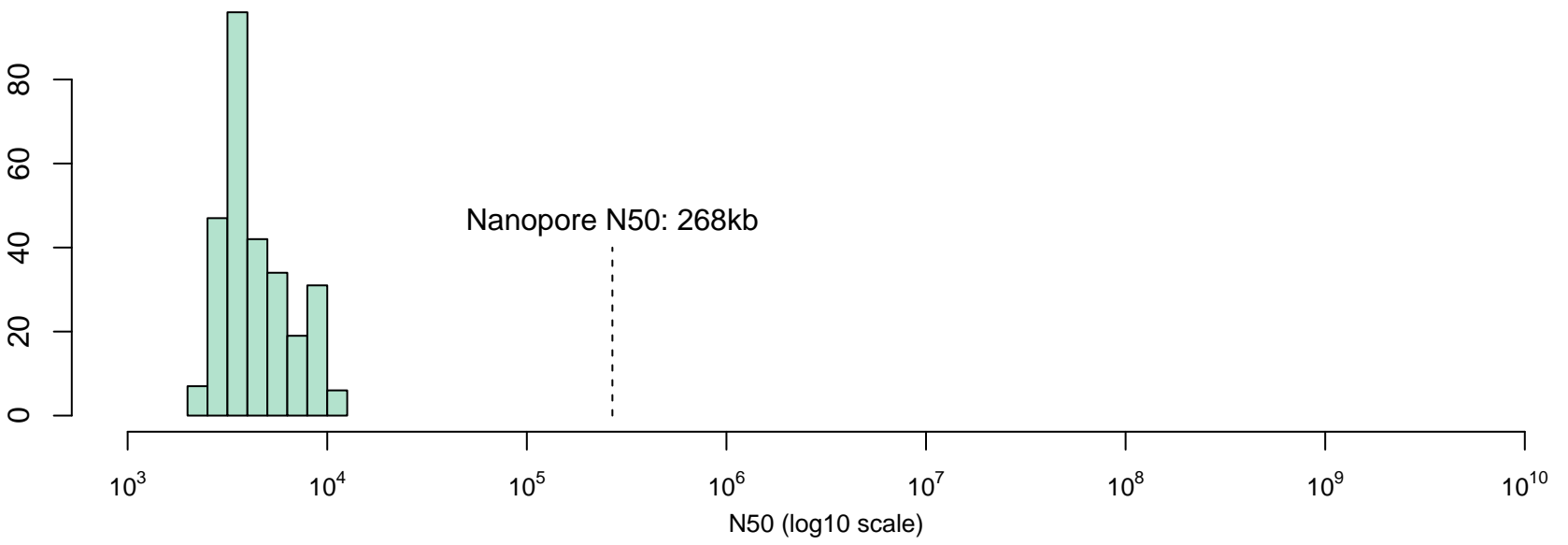
A: RUG vs Hungate



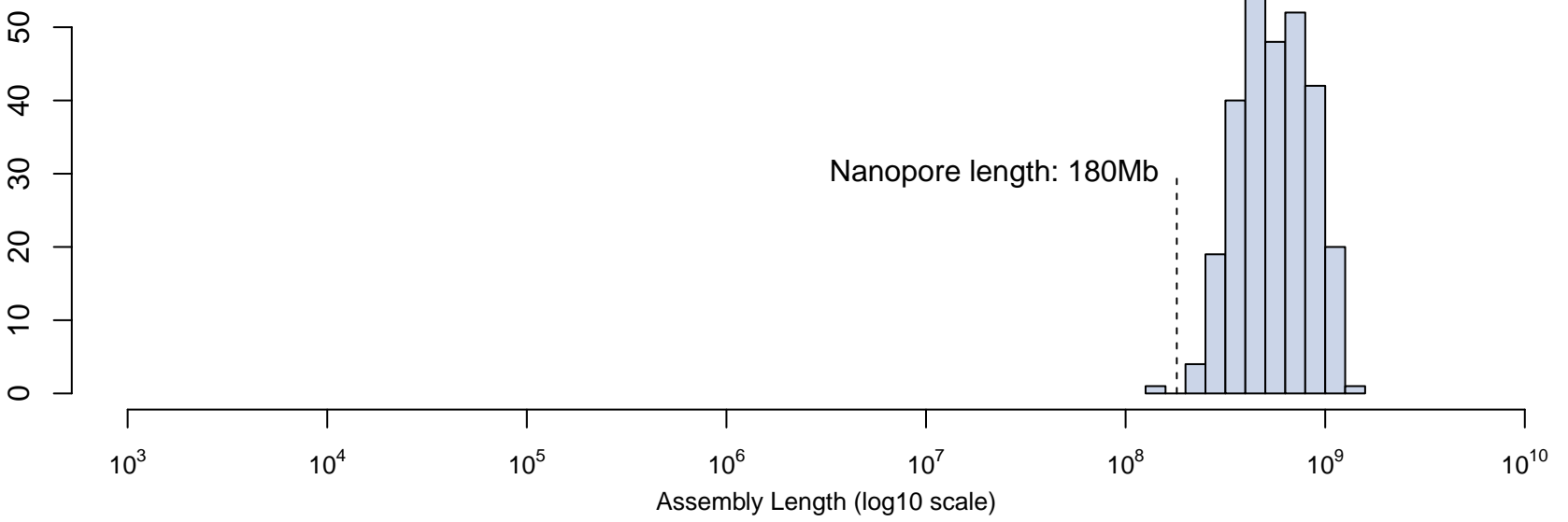
B: RUG vs Stewart et al



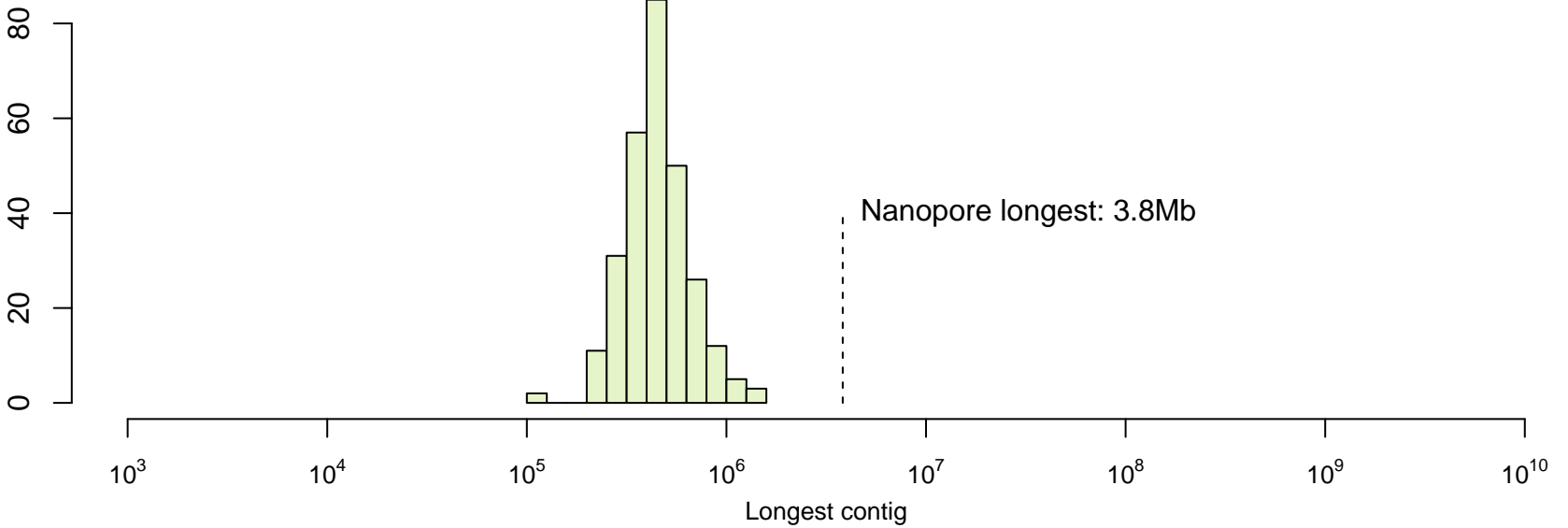
A: N50



B: Length of assembly



C: Longest contig



% identity against CAZy

