

1

2

3 A restriction enzyme reduced representation sequencing
4 approach for low-cost, high-throughput metagenome profiling

5

6 Melanie K. Hess^{1*}, Suzanne J. Rowe¹, Tracey C. Van Stijn¹, Hannah M. Henry¹, Sharon M.
7 Hickey², Rudiger Brauning¹, Alan F. McCulloch¹, Andrew S. Hess¹, Michelle R. Kirk³, Sandra
8 Kittelmann³, Graham R. Wood¹, Peter H. Janssen³ and John C. McEwan¹

9

10

11 ¹AgResearch Limited, Invermay Agricultural Centre, Puddle Alley, Mosgiel 9092, New Zealand

12 ²AgResearch Limited, Ruakura Agricultural Centre, Bisley Road, Hamilton 3214, New Zealand

13 ³AgResearch Limited, Grasslands Research Centre, Tennent Drive, Palmerston North 4442, New
14 Zealand

15

16 *Corresponding author

17 Email: melanie.hess@agresearch.co.nz (MKH)

18

19 **Abstract**

20 Microbial community profiles have been associated with a variety of traits, including
21 methane emissions in livestock, however, these profiles can be difficult and expensive to obtain
22 for thousands of samples. The objective of this work was to develop a low-cost, high-throughput
23 approach to capture the diversity of the rumen microbiome. Restriction enzyme reduced
24 representation sequencing (RE-RRS) using *ApeKI* or *PstI*, and two bioinformatic pipelines
25 (reference-based and reference-free) were compared to 16S rRNA gene sequencing using repeated
26 samples collected two weeks apart from 118 sheep that were phenotypically extreme (60 high and
27 58 low) for methane emitted per kg dry matter intake (n=236). DNA was extracted from freeze-
28 dried rumen samples using a phenol chloroform and bead-beating protocol prior to sequencing.
29 The resulting sequences were used to investigate the repeatability of the rumen microbial
30 community profiles, the effect of host genetics, laboratory and analytical method, and the genetic
31 and phenotypic correlations with methane production. The results suggested that the best method
32 was *PstI* RE-RRS analyzed with the reference-free approach via a correspondence analysis, with
33 estimates for repeatability of 0.62 ± 0.06 , heritability 0.31 ± 0.29 , and genetic and phenotypic
34 correlation with methane emissions of 0.88 ± 0.25 and 0.64 ± 0.05 respectively for the first
35 component of correspondence analysis. The reference-free approach assigned $62.0\pm 5.7\%$ of reads
36 to common 65 bp tags, much higher than the reference-based approach of $6.8\pm 1.8\%$ of reads
37 assigned. Sensitivity studies suggested approximately 2000 samples could be sequenced in a single
38 lane on an Illumina HiSeq 2500, therefore the current work of 118 samples/lane and future
39 proposed 384 samples/lane are well within that threshold. Our approach is now being used to
40 investigate host factors affecting the rumen and its association with a variety of production and

41 environmental traits. With minor adaptations, our approach could be used to obtain microbial
42 profiles from other metagenomic samples.

43

44

45 **Introduction**

46 Metagenomics is the study of genetic material recovered directly from environmental
47 samples and captures the myriad of organisms present in that environment. Samples of soil and
48 water are obvious examples of environmental samples; however, the gut can also be considered an
49 environment due to the presence of microbes that interact with the host, e.g. during digestion.
50 Metagenomic studies have gained popularity in recent years, primarily in human health, e.g.
51 Irritable Bowel Disease (1) and Coeliac Disease (2).

52 In an agricultural setting, rumen microbial community (RMC) profiles have been
53 associated with environmentally and economically important traits, such as methane emissions (3,
54 4) and feed efficiency (5, 6). The RMC breaks down ingested feed to produce volatile or short
55 chain fatty acids, which are a source of energy for the host. Hydrogen produced by this process is
56 metabolized into methane by methanogenic archaea. Methane production is not solely dependent
57 on the abundance of methanogens present, but also on the substrate available to them (7). Different
58 “ruminotypes”, generalized microbial community types, can be found in sheep with high and low
59 methane production, with at least two ruminotypes present in low-methane sheep fed lucerne
60 pellets (3). Furthermore, RMC profiles from rumen samples are moderately heritable (8),
61 suggesting that selection on microbiomes is likely to result in changes in offspring microbiomes.
62 Given that traits such as methane emissions and feed efficiency are difficult and expensive to

63 measure, selection on RMC profiles may facilitate progress in these traits, provided costs are low
64 enough and the method is high-throughput.

65 Historically, there have been two approaches used for sequencing metagenome samples:
66 targeted sequencing and whole genome shotgun (WGS) sequencing. Targeted sequencing
67 amplifies specified phylogenetically informative genes from a sample, such as the 16S rRNA gene
68 (16S) of microbes, which typically distinguishes taxonomic groups well due to large,
69 comprehensive databases of 16S rRNA sequences that include both culturable and unculturable
70 organisms (9, 10). This approach usually relies on having long sequence reads (11), only captures
71 phylogenetic variation at one gene, and is subject to PCR primer bias due to mismatches in the
72 flanking regions where the primers bind (12). WGS can capture any part of the microbial, host or
73 feed genome; but to be informative, a reference database of genome assemblies with known
74 taxonomies is needed to obtain taxonomic information on WGS sequences, e.g. the Hungate1000
75 Collection (13). Whole genome assemblies are currently difficult to obtain on unculturable
76 microbes, so these are usually missing from reference databases. Hundreds of millions of reads are
77 generated per sample for WGS, making it an expensive and time-consuming method that
78 additionally requires significant computation resources.

79 Restriction Enzyme-Reduced Representation Sequencing (RE-RRS, also known as
80 Genotyping-by-Sequencing or GBS) is a next-generation sequencing technique that reduces
81 genome complexity by digestion of genomic DNA by restriction enzymes, followed by the
82 sequencing of fragments within a given size range (14). RE-RRS is used to obtain genotypes for
83 parentage identification or genomic selection (to identify the individuals with the most favorable
84 phenotypes) in a variety of species across livestock, plants and aquaculture (15-18), as well as
85 population diversity studies, e.g. for conservation (19). RE-RRS holds promise as a technique for

86 rapid, high-throughput and cost-effective sequencing of metagenome samples at a fraction of the
87 cost of WGS. Underlying the RE-RRS method is the assumption that sequencing only a specific
88 fraction, typically 0.5-1% of any microbial genome as defined by restriction site and fragment size,
89 captures the majority of information on composition and diversity of the microbial community at
90 a fraction of the sequencing cost. This study used sheep rumen samples to show the potential of
91 RE-RRS as a low-cost, high-throughput approach for obtaining metagenome profiles on thousands
92 of samples, and describes pipelines for obtaining profiles both with and without a reference
93 database.

94

95

96 **Materials and methods**

97 **Rumen sampling and associated methane yields**

98 The sheep rumen samples and methane yield data used for this study were those for which
99 rumen microbial community structure was analysed using 16S rRNA gene sequencing in
100 Kittelmann et al. (3) and part of a larger experiment described in Pinares-Patiño et al. (20). Briefly,
101 respiration chambers were used to measure methane yield (g CH₄/ kg DMI) on 340 sheep at two
102 independent measuring rounds two weeks apart, each over two days in 4 separate cohorts of
103 animals. The rumen sample was collected via stomach tubing at the end of each measuring round
104 and immediately stored at -20°C. Two rumen samples from a subsample of 118 sheep,
105 representing the ~17% highest and lowest emitters (60 high-methane and 58 low-methane based
106 on methane yield phenotype), were previously freeze dried, homogenized and stored at -85°C.
107 Subsamples were used for analysis of RMC by sequencing amplified bacterial 16S rRNA genes,
108 and these sequences are available in the EMBL database under the study accession number

109 ERP003779 (3). The 16S rRNA gene profiles used in our study used the recently updated
110 taxonomic classification of these sequences by Kumar (21). These same freeze-dried and
111 homogenised samples were used in our study to evaluate the potential of using RE-RRS for RMC
112 profiling, as described below.

113

114 **DNA extraction and Restriction Enzyme-Reduced Representation**

115 **Sequencing**

116 DNA was extracted from rumen samples using a combined bead-beating, phenol and
117 column purification protocol using the QIAquick 96 PCR purification kit (Qiagen, Hilden,
118 Germany), as described in Text S1 of Kittelmann et al. (3), to provide high quality nucleic acids
119 for RE-RRS. *ApeKI* and *PstI* restriction enzymes were used separately to test whether RE-RRS is
120 a suitable approach for rumen metagenome profiling. These two enzymes were selected because
121 an in-silico digestion and size filtering of rumen microbial genome assemblies from the
122 Hungate1000 Collection (13) showed that RE-RRS using either *ApeKI* (G|CWGC) or *PstI*
123 (CTGCA|G) captured microbial sequence from all species present in the collection with an average
124 of 8.6% and 0.3% of each genome, respectively (22).

125 After digestion of DNA by either *ApeKI* or *PstI*, barcodes were ligated to link sequences
126 to samples, as described by Elshire et al. (14), and samples were grouped into two libraries, one
127 library for each restriction enzyme used, and PCR amplified. Amplified sequences between 193
128 and 318 bp (equivalent to 65 – 195 bp inserts) were selected using a Pippin Prep (SAGE Science,
129 Massachusetts, USA) and each library was run on two lanes on the same flow cell on an Illumina
130 HiSeq2500 machine, generating 101 bp single end reads using version 4 chemistry. One plate of

131 94 samples for *PstI* were re-run (in a single lane) because barcodes did not ligate in the initial run.
132 FastQ files were deposited in the NCBI SRA.

133

134 **Bioinformatic pipeline**

135 Sequenced reads were demultiplexed using GBSX (23), and trimmed using trim_galore
136 (24) for single reads with a minimum length of 40 base pairs. Samples with fewer than 100,000
137 reads across both lanes of sequencing for a single restriction enzyme were removed from all further
138 analyses, consisting of one sample for *ApeKI* and two samples for *PstI*. The trimmed sequences
139 from the remaining samples were run through both the reference-based and reference-free
140 pipelines, described below.

141

142 *Reference-based approach*

143 The reference-based (RB) approach used nucleotide BLAST (BLASTN) in BLAST
144 v2.2.28+ (25) with default parameters to compare sequenced reads against the 410 rumen
145 microbial genome assemblies from the Hungate1000 Collection (13). This was found to be the
146 optimal approach for aligning query sequences to the Hungate1000 Collection by Hess et al. (22).
147 Reads were assigned to a taxonomic node using the algorithm from MEGAN (26) implemented in
148 R with default parameters: a minimum bitscore of 50 and considering only hits within 10% of the
149 maximum bitscore for a query read. This approach was evaluated by Hess, Rowe (22) and found
150 to assign reads at the genus level with high accuracy (>96%). The RMC profile was defined as the
151 number of sequences assigned to each of the ~60 genera represented in the Hungate1000
152 Collection and associated analyses will be denoted by *ApeKI*_RB and *PstI*_RB for analyses of
153 *ApeKI* and *PstI* profiles, respectively.

154

155 *Reference-free approach*

156 The reference-free approach involved collating a set of “tags”, i.e. non-redundant 65 bp-
157 long DNA sequences (evaluated across all samples) commencing at the initial cut site. Tags were
158 required to be present in 10% and 25% of samples in *ApeKI* and *PstI*, respectively. A comparison
159 of performance for other tag lengths and a variety of prevalence thresholds can be found in File
160 S1– requiring 65 bp tags to be present in at least 10% of samples for *ApeKI* and 25% of samples
161 for *PstI* gave high estimates of heritability and repeatability with low standard errors, so these
162 parameters were selected for further analysis. The rumen community profile for each sample was
163 generated by counting the abundance of each tag from the sequenced reads; these profiles were
164 collated into a count matrix with samples as rows and tags as columns, obtained using an in-house
165 Unix script. Reference-free analyses will be denoted as *ApeKI*_RF and *PstI*_RF for profiles from
166 the *ApeKI* and *PstI* restriction enzymes, respectively.

167

168 **Comparison of methods for obtaining RMC profiles**

169 *Correspondence Analysis*

170 A correspondence analysis (27) was used to reduce the dimensionality of the dataset and
171 facilitate comparisons between the different methods for generating RMC profiles. In a similar
172 approach to that used by Rowe et al. (8), the coordinates of the first dimension of the
173 correspondence analysis (CA1) were used as the RMC profile phenotype for parameter estimation,
174 described below. The first dimension of a correspondence analysis is the dimension that explains
175 the largest proportion of variation in the data.

176 *Parameter Estimation*

177 Heritability, repeatability, and genetic and phenotypic correlations with scaled methane
178 yield were estimated in ASReml 4.1 (28). Heritability and repeatability of CA1 were estimated
179 using a univariate mixed linear model, and correlation between CA1 and scaled methane yield was
180 estimated with a bivariate mixed linear model. Scaled methane yield was obtained by dividing
181 methane yield by the contemporary group mean and multiplying by the overall mean, where
182 contemporary group included recording year, lot, group and round and the overall mean was 16.0
183 kg, as described in Pinares-Patiño et al. (20). In both univariate and bivariate models, sex and
184 cohort (lot and round) were fitted as fixed class effects, random animal genetic effects were
185 estimated based on pedigree relationships, and a random permanent environmental effect linked
186 duplicate samples from the same animal. In some cases, the model was unable to separate animal
187 and permanent environmental effects, in which case animal was dropped from the model and only
188 repeatability reported.

189

190 **Sensitivity to sequencing depth**

191 Sequencing more samples per lane would lower the cost of RE-RRS profiling but would
192 consequently reduce the sequencing depth. At low depths the profiling would not accurately
193 capture the proportion of each microbe in the sample, particularly microbes that are in low
194 abundance. Therefore, a sensitivity analysis was performed evaluate the impact of reducing the
195 sequencing depth in our approach. Reads were subsampled with probability 0.5, 0.25, 0.1, 0.05,
196 0.01, 0.005, 0.002 or 0.001; representing sequencing 2, 4, 10, 20, 100, 200, 500 or 1000 times the
197 number of samples per lane, respectively. The set of sampled reads at a given simulated sequencing
198 depth were then used to calculate repeatability, as above, or compression efficiency (29).
199 Compression efficiency compares the size of a compressed file to its original size as (original –

200 compressed)/original and is a measure of the non-redundant information present in the file. In our
201 study, the original file contained the reads for a given sample without their identifiers. This file
202 was compressed using gzip 1.3.12 (30), which uses the DEFLATE algorithm (31). The value of
203 compression efficiency was the mean across all samples for the simulated sequencing depth.
204 Standard errors for repeatability and compression efficiency were the standard deviation across
205 five replicates at that sequencing depth.

206

207

208 **Results and Discussion**

209 **Sequencing Results**

210 Sequence read quality was high for all lanes sequenced (Figure S1). A greater average
211 number of reads per sample was observed for samples digested with *PstI* (Table 1), likely partially
212 due to re-running of samples – only 94 samples were run in that lane rather than 118 (i.e. 236
213 samples across 2 lanes). Sequences from the *ApeKI* digest were slightly shorter than from the *PstI*
214 digest, however this difference was not significant based on a t-test with $\alpha = 0.05$.

215 **Table 1: Average Number of Reads per Sample and Average Read Length of RE-RRS Reads.**

Restriction enzyme	Reads per sample (sd)	Read length (sd)¹
<i>ApeKI</i>	2.4M (870k)	71 (17)
<i>PstI</i>	2.7M (680k)	84 (15)

216 1. Trimmed read length in base pairs after barcode removed

217 **Reference-based approach**

218 Using the MEGAN algorithm on nucleotide BLAST results, $5.3 \pm 1.7\%$ and $6.8 \pm 1.8\%$ of
219 reads were assigned at the genus level for *ApeKI* and *PstI*, respectively. This assignment rate is

220 consistent with querying the Hungate1000 Collection with reads from WGS (22, data from Shi et
 221 al. (7)). Comparing against a protein database is one method that could potentially improve the
 222 proportion of sequences assigned (hit rate) at the genus level. However, Hess et al. (22) found only
 223 a small increase in hit rate at the genus level when BLASTX was used. They also found the time
 224 taken to perform the BLAST query and analyze the results was much longer when BLASTX was
 225 used compared to BLASTN and therefore was not desirable for a high-throughput pipeline.

226 A significant difference in hit rate between high- and low- methane animals was found for
 227 both *ApeKI* ($p = 1.4 \times 10^{-8}$) and *PstI* ($p = 1.4 \times 10^{-8}$; Table 2). This may be attributed to the presence
 228 or absence of some species associated with methane yield in the Hungate1000 Collection. For
 229 example, Kittelmann et al. (3) identified the genera *Fibrobacter*, *Kandleria*, *Olsenella* and *Sharpea*
 230 to be in higher prevalence in low-methane yield animals. These genera are all present within the
 231 Hungate1000 Collection and have equal or significantly higher abundance in samples from low-
 232 methane animals. The Hungate1000 Collection also has poor or no representation of other genera
 233 that were found by Kittelmann et al. (3) to be in higher abundance in high-methane yield animals
 234 e.g. *Coprococcus*. This shows that using a method that is reliant on a reference database is limited
 235 by the genomes present within that database.

236 **Table 2: Hit Rates by Taxonomic Level from RE-RRS Samples using One of Two Restriction**
 237 **Enzymes.**

Restriction enzyme	Sample ¹	Hit rate by taxonomic level (%)						
		Kingdom	Phylum	Class	Order	Family	Genus	Species
<i>ApeKI</i>	High	4.9	4.9	4.8	4.8	4.7	4.7	1.4
	Low	6.2	6.1	6.1	6.1	5.9	5.9	2.2
<i>PstI</i>	High	6.5	6.4	6.4	6.4	6.3	6.3	1.7
	Low	7.4	7.4	7.4	7.4	7.3	7.3	2.3

238 1. Methane yield classification (high- or low-methane yield) of the sheep the sample came from.

239

240 A major gap in microbial genome assemblies is the inability, at least historically, to
241 sequence the unculturable microbes that make up a large proportion of any environment.
242 Technological advances, such as single-cell sequencing (32) and the ability to assemble genomes
243 from metagenomic datasets (33), offer alternative solutions to sequence and assemble microbial
244 genomes and will provide opportunities to improve reference databases. Judicious addition of new
245 microbial genome assemblies as they become available will improve hit rates, however, any
246 additional sequences added to the database will also increase the time to complete the analysis,
247 which may not be desirable for a high-throughput approach if there are time constraints. If
248 additional genomes were to be added to the Hungate1000 Collection (or another reference
249 database), expert curation would be needed to ensure the quality of genome assemblies and balance
250 the resource across taxa to maximise coverage and minimise duplication.

251

252 **Reference-free approach**

253 The reference-free approach is not subject to the biases of the species represented in the
254 Hungate1000 Collection. We explored the use of different tag lengths and filtering thresholds and
255 showed that there were different optimal filtering levels when using *ApeKI* (10%) or *PstI* (25%)
256 restriction enzymes (File S1). When *ApeKI* was used, there were ~1.2M 65 bp tags that were
257 present in at least 10% of the samples and these tags accounted for $20.0 \pm 3.5\%$ of reads
258 (corresponding to $33.6 \pm 5.8\%$ of reads at least 65 bp long). When *PstI* was used, there were ~500k
259 65 bp tags present in at least 25% of samples and $53.3 \pm 5.9\%$ of reads were accounted for
260 (corresponding to $61.9 \pm 5.7\%$ of reads at least 65 bp long). Although these proportions are given
261 at different filtering levels, File S1 shows that *PstI* captures a greater proportion of reads than
262 *ApeKI* at all filtering levels and tag lengths. These differences can be explained by the proportion

263 of each microbial genome that is expected to be captured using each restriction enzyme: *ApeKI*
264 captures 8.6% of Hungate1000 Collection genomes on average, while *PstI* captures 0.3% (22).
265 This means that, for a given number of sequences (e.g. one lane of sequencing), the fewer regions
266 captured by *PstI* reads will be at higher depth, whereas the greater number of regions captured by
267 *ApeKI* will be at lower depth. This is shown by the much larger number of unique tags present
268 when using *ApeKI* compared to *PstI* (File S1) despite a slightly higher number of reads per sample
269 for *PstI* (Table 1).

270 The reference-free pipeline is particularly useful for prediction of a trait that is correlated
271 with a microbial profile because knowledge of which taxonomic group a sequence belongs to is of
272 less importance than its predictive ability, so sequences that don't align to a reference genome can
273 still be utilized. Subsequently, if taxonomic information is desired, tags can be searched against a
274 relevant database; this process would be computationally inexpensive because there are fewer
275 search terms, i.e. fewer tags (hundreds of thousands) than the full set of reads in the original dataset
276 (tens or hundreds of millions). Given the large number of tags generated using the reference-free
277 approach it may be desirable to cluster these into groups (e.g. through sequence similarity,
278 taxonomic assignment, high positive correlations between tag abundances). Tags that come from
279 the same organism will be highly correlated, however, a high correlation (positive or negative)
280 could also come about due to interactions between the microbes, or by chance.

281

282 **Comparison of methods for obtaining rumen microbial profiles**

283 *Variance components of RMC profiles*

284 The first component of the correspondence analysis (CA1) was analysed as a trait for the
285 four RE-RRS approaches and the 16S rRNA gene taxonomic classifications from Kumar (21) on

286 the same samples. The variance of CA1 from the 16S rRNA gene profile was the highest, capturing
 287 24.3% of the variance in that profile (Table 3). The reference-based approaches had the smallest
 288 variances, but they explained 40-50% of the variation in the profiles. Lastly, the reference-free
 289 approaches both had variances just over 0.2 and explained less than 5% of the variance in these
 290 profiles. The percent variance explained by CA1 was negatively correlated with the number of
 291 tags or taxa: the reference-based approaches assigned reads to only 60 genera; the 16S rRNA gene
 292 approach assigned reads to ~250 genera and the reference-free approach assigned reads to ~1.2M
 293 and ~500k tags for *ApeKI* and *PstI*, respectively. This indicates that a large proportion of the
 294 variation in the rumen community profile is not accounted for by analysing only CA1 for these
 295 profiles, particularly for the reference-free approach. Nevertheless, evaluating CA1 allowed us to
 296 easily compare the different approaches and below we have discussed statistical approaches that
 297 will utilize more of the information contained within the profile.

298 **Table 3: Variance, heritability and repeatability of the first component of a correspondence**
 299 **analysis of rumen community profiles**

Method	%Var	Var	Gen Var	Heritability	Repeatability
16S	24.3	0.30	0.06 ± 0.05	0.26 ± 0.23	0.45 ± 0.08
<i>ApeKI</i> _RB	47.4	0.07	0.03 ± 0.02	0.58 ± 0.32	0.61 ± 0.06
<i>PstI</i> _RB	41.0	0.06	NE	NE	0.60 ± 0.06
<i>ApeKI</i> _RF	3.3	0.22	0.03 ± 0.03	0.18 ± 0.25	0.60 ± 0.06
<i>PstI</i> _RF	3.9	0.22	0.03 ± 0.04	0.24 ± 0.27	0.62 ± 0.06

300 %Var = Percent of the overall variance explained by first component of correspondence analysis

301 Var = Variance of first component of correspondence analysis

302 Gen Var = Genetic variance of first component of correspondence analysis

303 16S = 16S rRNA gene sequencing approach

304 RB = Reference-based approach

305 RF = Reference-free approach

306 NE = Not Estimable

307 Despite the differences in the variances of the first component of the correspondence
 308 analysis, the genetic variance was very similar for each of the methods (Table 3). The combination
 309 of similar genetic variances but different phenotypic variances produced heritability estimates
 310 ranging from 0.18 (*ApeKI*_RF) to 0.58 (*ApeKI*_RB), with the 16S rRNA gene approach giving an

311 intermediate heritability estimate of 0.26. Despite coming from animals with extreme phenotypes
312 for methane yield, these heritability estimates are consistent with other studies on livestock and
313 human intestinal microbial communities (4, 8); however, all estimates of heritability in our study
314 had large standard errors due to the relatively small number of samples and were not significantly
315 different from each other. The *PstI*_RB approach was unable to separate genetic and permanent
316 environmental effects, with the animal effect being bound at zero, so the animal effect was
317 removed from the model. The inability of the *PstI*_RB approach to separate these effects, as well
318 as the low variance of that component, indicates that it may be a less powerful approach. Increasing
319 the number of samples would allow the model to more accurately separate the genetic and
320 permanent environmental effects for all methods.

321 Repeatability is a measure of the similarity of two samples from the same individual and
322 is calculated as the genetic plus permanent environmental effect as a proportion of the phenotypic
323 variance. CA1 of the RMC profiles from RE-RRS were more consistent across time than 16S
324 rRNA gene profiles, as evidenced by repeatability estimates almost 1.5 times higher (Table 3).

325

326 *Correlations with methane yield*

327 The genetic correlation of the first component of the correspondence analysis with methane
328 yield was highest for the reference-free approach using *PstI*, followed by the 16S rRNA gene
329 approach and the reference-based approach using *ApeKI* (Table 4). Estimates were not available
330 for the other two approaches due to an inability to separate genetic and permanent environmental
331 effects. Phenotypic correlations between the first component of a correspondence analysis and
332 methane yield were highest for the two reference-free approaches, followed by the two reference-
333 based and finally the 16S rRNA gene results. These phenotypic and genetic correlations suggest

334 that RE-RRS may produce rumen community profiles with a greater predictive ability than profiles
 335 produced from 16S rRNA gene sequencing; however, more samples are needed to test this
 336 accurately. The *ApeKI*_RB method has the same genetic correlation with methane yield as 16S
 337 rRNA gene sequencing but a higher phenotypic correlation, suggesting that *ApeKI* captures more
 338 of the metagenomic variation that is phenotypically correlated with methane yield. There is likely
 339 to be some non-microbial DNA that is being captured and influencing the correlations. These
 340 results suggest that *ApeKI* will predict the individual's methane production better than 16S rRNA
 341 gene sequencing but not necessarily the genetic potential of that individual as a parent.

342 **Table 4: Genetic and phenotypic correlations of the first component of a correspondence**
 343 **analysis of rumen community profiles with methane yield**

Method	r_g (CH ₄ Yield)	r_p (CH ₄ Yield)
16S	0.63 ± 0.49	0.26 ± 0.07
<i>ApeKI</i> _RB	0.63 ± 0.31	0.48 ± 0.06
<i>PstI</i> _RB	NE	0.33 ± 0.07
<i>ApeKI</i> _RF	NE	0.66 ± 0.04
<i>PstI</i> _RF	0.88 ± 0.25	0.64 ± 0.05

344 16S = 16S rRNA gene sequencing approach
 345 RB = Reference-based approach
 346 RF = Reference-free approach
 347 NE = Not Estimable
 348

349 The higher correlations between methane yield and the reference free approaches (Table
 350 4) suggest that this approach might be capturing components of the rumen microbiome that are not
 351 being captured by 16S rRNA gene sequencing or the reference-based approaches. Both 16S rRNA
 352 gene and the reference-based approaches are focused on capturing microbial community profiles,
 353 while the reference-free approach may capture DNA from a much wider taxonomic range, e.g.
 354 host, feed, viruses. If the aim is to obtain the most accurate predictions, then this information from
 355 a wider taxonomic range is beneficial to include in the analysis. The reference-free approaches
 356 also have the lowest proportion of the variance captured within the first component of the

357 correspondence analysis (Table 3), suggesting that there is potential for the reference-free rumen
358 community profiles to more accurately predict methane yield if additional information is
359 considered. The additional information that is not contained within the first component will be a
360 mixture of “signal” that is associated with methane yield, and “noise” that is not. Too much noise
361 will have a negative impact on prediction accuracies; approaches to separate signal from noise are
362 suggested in the statistical modelling section below.

363

364 **Sensitivity to sequencing depth**

365 The number of samples per lane influences the cost per sample of sequencing as well as
366 the average number of sequenced reads per sample. A sensitivity analysis was performed by
367 subsampling reads from our RE-RRS samples and evaluating the repeatability of the first
368 component of a reference-based correspondence analysis, as well as the compression efficiency of
369 the dataset. Both analyses showed that sampling 5% of reads, corresponding to 20 times the
370 number of samples per lane, i.e. 2000 samples, is the lower bound because compression efficiency
371 drops, and the standard error of repeatability increases when sequencing depth is lowered beyond
372 this point (Figure 1).

373 **Figure 1: Repeatability (A) and Compression Efficiency (B) of RE-RRS Data as the Percent 374 of Reads Sampled Decreases.**

375 The standard error of the estimate of repeatability of the first component of a reference-based
376 correspondence analysis increases (A), and the compression efficiency of sequence data decreases
377 (B), when less than 5% of reads are sampled, a sequencing depth that corresponds to 20 times the
378 number of samples sequenced per lane. This number is consistent across both restriction enzymes
379 used for this study. Standard errors for compression efficiency are negligible and therefore not
380 visible (B).

381

382 We have determined a potential cut-off for how much we can increase the throughput
383 without losing crucial information using two separate approaches (compression efficiency and
384 repeatability, Figure 1). It is worth mentioning, however, that these samples are from a relatively

385 small set of individuals that have extreme phenotypes, so in practice more sequences may be
386 needed than the observed cut-off. This analysis shows that the depth with which sequencing has
387 occurred in this study is well within reasonable bounds for capturing metagenomics data, and that
388 the throughput could be safely increased 2-4× over what was done in this study. This will reduce
389 costs and allow faster turn-around times for obtaining sequencing data when large numbers of
390 samples are analysed.

391
392

393 **Utility of a high-throughput Metagenomics Method**

394 *RE-RRS vs. 16S rRNA gene sequencing and WGS*

395 Our RE-RRS approach to sequencing rumen samples is likely to perform as well or better
396 than 16S rRNA gene sequencing in terms of the variation in sequence reads that is accounted for,
397 and the predictive ability of rumen community profiles (Table 2). Although RE-RRS can capture
398 taxonomic information, like 16S rRNA gene sequencing it cannot directly quantify the abundance
399 of particular genes within a sample (16S rRNA gene sequencing only captures the relative
400 abundance of variants of the 16S rRNA gene); most will be missed because it is a reduced
401 representation sequencing approach that only captures a small percentage of each microbial
402 genome. WGS can capture information on the abundances of these genes; however, sequencing
403 must be done at high depth to capture this information accurately, which is very expensive.

404 With any metagenomics sequencing approach it is important to try to avoid biases in what
405 is captured. The 16S rRNA gene is present in all bacteria and archaea and contains some highly
406 conserved regions, suitable for universal primer design but also contains highly variable regions
407 which allows the sequences to be assigned to taxonomies. However, sometimes the highly
408 conserved regions have some variation, which leads to primer bias (e.g. Sim et al. (12)). RE-RRS

409 was shown to capture part of each genome in the Hungate1000 Collection (22), suggesting it is not
410 prone to similar biases; however, the Hungate1000 Collection captures only culturable microbes
411 so it is possible that some other microbes are not captured by our approach. This can be evaluated
412 further as more microbial genome assemblies become available.

413 16S rRNA gene sequencing, WGS and the reference-based RE-RRS approach all require
414 a reference database to assign taxonomic information to the sequences. 16S rRNA gene reference
415 databases tend to be more comprehensive because only a single gene needs to be sequenced,
416 enabling both culturable and unculturable microbes to be captured (10). WGS and the reference-
417 based RE-RRS approach both need reference databases containing genome assemblies. These
418 databases are (currently) less complete and capture a smaller range of taxa. All reference databases
419 such as this are prone to sequencing errors present in the reference database, which may result in
420 incorrect taxonomic assignment.

421

422 *Reference-based vs reference-free*

423 The reference-free approach was developed to overcome the reliance on a reference-
424 database to generate the RMC profile. If the goal of the analysis is prediction of a trait or disease
425 status, it is not critical to know the taxonomic origin of a sequence, as long as it can positively
426 influence prediction accuracy. Our results in Tables 3 and 4 indicate that the reference-free
427 approach will be a valuable approach for predicting methane yield. The reference-free approach is
428 likely capturing more than just microbial variation e.g. host, feed and microbial eukaryotes, which
429 may all play a role in methane emissions.

430

431

432 *Impact of different restriction enzymes*

433 The choice of restriction enzyme will impact sample throughput. *ApeKI* has a lower
434 compression efficiency than *PstI*, which is consistent with *ApeKI* capturing more of the genome
435 (Figure 1). This means that two sequences from the same genome with *PstI* are more likely to be
436 identical (i.e. less of the genome but at a higher depth), resulting in improved compression
437 efficiency. This can also be inferred from the filtering threshold required to achieve a similar
438 repeatability estimate for *ApeKI* (tags must be present in at least 10% of samples, ~1.2M tags)
439 compared to *PstI* (tags must be present in 25% of samples, ~500k tags) (File S1).

440 The *PstI*_RF approach performed the best of all the models in terms of heritability,
441 repeatability, and correlations with methane yield. File S1 shows that heritability and repeatability
442 estimates changed very little when different filtering parameters and tag lengths were applied to
443 the *PstI*_RF data. This indicates that the depth of sequencing is appropriate, which is supported by
444 the higher compression efficiency of *PstI* compared to *ApeKI* (Fig 1b).

445 *ApeKI* had a lower compression efficiency than *PstI* (Fig 1b), and when the reference-free
446 approach was used tags needed to be present in 10% of samples before heritability and repeatability
447 estimates were in line with *PstI* (File S1). This is because *ApeKI* captures more regions of the
448 microbial genome at lower depth than *PstI* so is not as powerful for the reference-free approach.
449 By extension, the reference-free approach developed here would not be suitable for WGS data
450 because any part of the metagenome could be captured, this could explain the poor results using a
451 k-mer approach in Ross et al. (34). Therefore, the profiling pipeline needs to be chosen based on
452 the sequencing approach used and the intended analysis.

453 Reducing tag length will join tags with similar sequences (e.g. two 65 bp tags that differ at
454 only base 40 will be merged into a single 32 bp tag). The shorter tag might represent a higher

455 taxonomic order than the larger one (possibly being less biologically informative, given the higher
456 specificity of a longer sequence) but this is outweighed by the increased power to accurately
457 capture the abundance of that taxonomic group.

458 These results show the importance of selecting an appropriate restriction enzyme for the
459 purpose of the study when using RE-RRS for metagenome profiling.

460

461 *Application of RE-RRS in livestock*

462 Most metagenome studies in livestock have used small sample sizes and many used
463 animals with extreme phenotypes, which is valuable for identifying whether there is a relationship
464 between microbes and traits of interest. However, knowledge of the RMC of thousands of animals
465 has the potential to reduce the carbon footprint of farming through selection of individuals with a
466 rumen microbiome genetically associated with lower methane emissions. Traits aimed at reducing
467 the carbon footprint of livestock animals, e.g. methane emissions or feed efficiency, are often
468 difficult and expensive to measure. Therefore, provided the costs can be reduced sufficiently and
469 high-throughput profiling is possible, large volumes of samples can be processed and the data
470 analyzed quickly and cheaply. In this situation, metagenome profiling could provide an alternative
471 solution to reducing the carbon footprint that circumvents the need to continually measure
472 expensive methane yield phenotypes on thousands of animals.

473

474 *Other sample types*

475 Much research has been done into sequencing microbial samples from humans (35, 36),
476 particularly samples related to the digestive tract and their association with a variety of health
477 issues (1, 2). A cheap, high-throughput metagenome sequencing approach has the potential to

478 make screening of these samples more accessible to those that require them. High-throughput
479 metagenome profiling has the potential to improve monitoring of other environmental samples as
480 well. This could range from identifying pathogens in water samples, evaluating the quality of water
481 in different environments, to identifying favorable and unfavorable soil environments for the
482 growth of particular crops. Further research is required to evaluate the potential of RE-RRS in each
483 of these situations.

484

485 *Statistical modelling*

486 Biological data has increased in complexity in recent years, as technological advances have
487 enabled the collection of many different parameters on large numbers of samples. Metagenomics
488 provides a rich source of information to include in predictions of associated traits, for example
489 rumen metagenomics samples to predict methane emissions. Statistical methodologies have
490 advanced along with data complexity, enabling the integration of many different –omics in an
491 attempt to improve prediction accuracies. However, it is important to appropriately model this
492 information to obtain robust predictions that extend to a wide diversity of situations.

493 One important aspect is to appropriately represent the data. Some taxa are present in all
494 samples while others appear in only some samples. For those present in most or all samples, a
495 correlation between the abundance of each tag/genus within each sample is the most appropriate
496 measure; however, for those present in only a subset of samples, a presence/absence coding may
497 be more appropriate. The abundance of each genera within a sample can vary by orders of
498 magnitude, therefore for prediction purposes it is common to transform the counts or proportions
499 of each taxonomic group, often a log transformation (37).

500 We summarized the RMC profile using the first component of a correspondence analysis;
501 however, other approaches can also be used, e.g. heritability and repeatability could be calculated

502 for each genus or tag. This information can be used to remove genera/tags, or a “metagenome wide
503 association study” (38) – similar to a genome wide association study – can be used to identify tags
504 that are associated with the trait of interest for inclusion in prediction models. This approach will
505 reduce the amount of noise in the profile, hopefully improving prediction accuracy. Other ways to
506 reduce the dimensionality, particularly for the reference-free approach, would be to calculate
507 correlations between each tag and merge the counts for those with high correlations. This approach
508 may be more powerful than using all tags independently because if a tag is not sequenced by
509 chance even though it is present in the sample then a correlated tag may be picked up instead; this
510 will be most important for genera at low abundance.

511 The metagenomics information can be included in prediction models in a variety of ways
512 depending on the circumstances. Microbial Relationship Matrices (MRMs) have been used to
513 include microbial information to predict traits of interest, particularly in livestock (37). This
514 approach will work when there are more taxa/tags than samples but runs into problems when the
515 opposite is true (39). When there are more samples than taxa/tags, fitting each taxa/tag as a
516 covariate, such as in Bayesian random regression model, is an alternative approach (39). Machine
517 learning algorithms may be useful for integrating host genetic information with microbial data for
518 trait predictions, while random forests hold potential for classifying samples into groups. Our
519 development of a high-throughput, low-cost approach to sequence rumen metagenomics samples
520 has set the stage for generating a large dataset where different modelling approaches can be
521 compared in the future.

522

523

524 **Conclusions**

525 We have shown that RE-RRS is a promising method for obtaining low-cost, high-throughput
526 metagenomic profiles that performs at least as well as 16S rRNA gene sequencing. Metagenomic
527 profiles can be generated either with or without a reference database (reference-based or reference-
528 free, respectively) depending on the purpose of the analysis. Gathering metagenomic information
529 on a large number of animals can be a useful addition to genomic information for the prediction
530 of traits in livestock production and human health. The next steps are to use this approach to
531 sequence thousands of environmental samples and develop appropriate statistical models for
532 prediction purposes.

533

534

535 **Acknowledgements**

536 Our thanks to Drs Graeme Attwood, Kathryn McRae and Ken Dodds (all from
537 AgResearch) for critically reviewing this manuscript, and to Dr Graeme Attwood for insightful
538 discussions throughout the process.

539

540

541 **References**

- 542 1. Young W, Jester T, Stoll ML, Izcue A. Inflammatory Bowel Disease. In: Ragab G,
543 Atkinson T, Stoll M, editors. *The Microbiome in Rheumatic Disease and Infection*. Cham:
544 Springer; 2018.
- 545 2. Lebwohl B, Sanders DS, Green PH. Coeliac disease. *The Lancet*. 2018;391(10115):70-81.

- 546 3. Kittelmann S, Pinares-Patiño CS, Seedorf H, Kirk MR, Ganesh S, McEwan JC, et al. Two
547 different bacterial community types are linked with the low-methane emission trait in sheep. *PLoS*
548 *ONE*. 2014;9(7):e103171.
- 549 4. Difford GF, Plichta DR, Løvendahl P, Lassen J, Noel SJ, Højberg O, et al. Host genetics
550 and the rumen microbiome jointly associate with methane emissions in dairy cows. *PLoS Genet*.
551 2018;14(10):e1007580.
- 552 5. Shabat SKB, Sasson G, Doron-Faigenboim A, Durman T, Yaacoby S, Miller MEB, et al.
553 Specific microbiome-dependent mechanisms underlie the energy harvest efficiency of ruminants.
554 *ISME J*. 2016;10(12):2958.
- 555 6. Sasson G, Ben-Shabat SK, Seroussi E, Doron-Faigenboim A, Shterzer N, Yaacoby S, et
556 al. Heritable bovine rumen bacteria are phylogenetically related and correlated with the cow's
557 capacity to harvest energy from its feed. *MBio*. 2017;8(4):e00703-17.
- 558 7. Shi W, Moon CD, Leahy SC, Kang D, Froula J, Kittelmann S, et al. Methane yield
559 phenotypes linked to differential gene expression in the sheep rumen microbiome. *Genome Res*.
560 2014;24(9):1517-25.
- 561 8. Rowe SJ, Kittelmann S, Pinares-Patiño CS, Wood G, Dodds KG, Kirk MR, et al., editors.
562 BRIEF COMMUNICATION: Genetic control of the rumen microbiome in sheep. *Proceedings of*
563 *the New Zealand Society of Animal Production*. 2015;75:67-69.
- 564 9. DeSantis TZ, Hugenholtz P, Larsen N, Rojas M, Brodie EL, Keller K, et al. Greengenes, a
565 chimera-checked 16S rRNA gene database and workbench compatible with ARB. *Appl. Environ.*
566 *Microbiol*. 2006;72(7):5069-72.

- 567 10. Henderson G, Yilmaz P, Kumar S, Forster RJ, Kelly WJ, Leahy SC, et al. Improved
568 taxonomic assignment of rumen bacterial 16S rRNA sequences using a revised SILVA taxonomic
569 framework. *PeerJ*. 2019;7:e6496.
- 570 11. Franzén O, Hu J, Bao X, Itzkowitz SH, Peter I, Bashir A. Improved OTU-picking using
571 long-read 16S rRNA gene amplicon sequencing and generic hierarchical clustering. *Microbiome*.
572 2015;3(1):43.
- 573 12. Sim K, Cox MJ, Wopereis H, Martin R, Knol J, Li M-S, et al. Improved detection of
574 bifidobacteria with optimized 16S rRNA-gene based pyrosequencing. *PloS ONE*.
575 2012;7(3):e32543.
- 576 13. Seshadri R, Leahy SC, Attwood GT, Teh KH, Lambie SC, Cookson AL, et al. Cultivation
577 and sequencing of rumen microbiome members from the Hungate1000 Collection. *Nat*.
578 *Biotechnol*. 2018;36(4):359.
- 579 14. Elshire RJ, Glaubitz JC, Sun Q, Poland JA, Kawamoto K, Buckler ES, et al. A robust,
580 simple genotyping-by-sequencing (GBS) approach for high diversity species. *PLoS ONE*.
581 2011;6(5):e19379.
- 582 15. Dodds KG, McEwan JC, Brauning R, Anderson RM, Stijn TC, Kristjánsson T, et al.
583 Construction of relatedness matrices using genotyping-by-sequencing data. *BMC Genomics*.
584 2015;16(1):1047.
- 585 16. Faville MJ, Ganesh S, Cao M, Jahufer MZ, Bilton TP, Easton HS, et al. Predictive ability
586 of genomic selection models in a multi-population perennial ryegrass training set using
587 genotyping-by-sequencing. *Theor. Appl. Genet*. 2018;131(3):703-20.

- 588 17. Rowe S, Dodds K, Ward J, Asher G, Brauning R, McEwan J, et al. Using low-depth
589 genotyping-by-sequencing for genomic analyses in farmed New Zealand deer. World Congress
590 on Genetics Applied to Livestock Production; Auckland, New Zealand. 2018.
- 591 18. Kim C, Guo H, Kong W, Chandnani R, Shuang L-S, Paterson AH. Application of
592 genotyping by sequencing technology to a variety of crop breeding programs. *Plant Sci.*
593 2016;242:14-22.
- 594 19. Dussex N, Taylor HR, Stovall WR, Rutherford K, Dodds KG, Clarke SM, et al. Reduced
595 representation sequencing detects only subtle regional structure in a heavily exploited and rapidly
596 recolonizing marine mammal species. *Ecol Evol.* 2018;8(17):8736-49.
- 597 20. Pinares-Patiño C, Hickey S, Young E, Dodds K, MacLean S, Molano G, et al. Heritability
598 estimates of methane emissions from sheep. *Animal.* 2013;7(s2):316-21.
- 599 21. Kumar S. Physiology of rumen bacteria associated with low methane emitting sheep.
600 Doctoral Dissertation, Massey University. 2017. Available from: <http://hdl.handle.net/10179/13403>
- 601 22. Hess MK, Rowe SJ, Van Stijn TC, Brauning R, Hess AS, Kirk MR, et al. High-throughput
602 rumen microbial profiling using genotyping-by-sequencing. World Congress on Genetics Applied
603 to Livestock Production; Auckland, New Zealand. 2018.
- 604 23. Herten K, Hestand MS, Vermeesch JR, Van Houdt JK. GBSX: a toolkit for experimental
605 design and demultiplexing genotyping by sequencing experiments. *BMC Bioinformatics.*
606 2015;16(1):73.
- 607 24. Krueger F. Trim Galore: A wrapper tool around Cutadapt and FastQC to consistently apply
608 quality and adapter trimming to FastQ files. 2015. Available from:
609 https://www.bioinformatics.babraham.ac.uk/projects/trim_galore/

- 610 25. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, et al. BLAST+:
611 architecture and applications. *BMC Bioinformatics*. 2009;10(1):421.
- 612 26. Huson DH, Auch AF, Qi J, Schuster SC. MEGAN analysis of metagenomic data. *Genome*
613 *Res*. 2007;17(3):377-86.
- 614 27. Greenacre MJ. *Theory and applications of correspondence analysis*. 1984.
- 615 28. Gilmour A, Gogel B, Cullis B, Welham S, Thompson R. *ASReml user guide release 4.1*
616 *structural specification*. Hemel hempstead: VSN international ltd. 2015.
- 617 29. Hudson NJ, Porto-Neto LR, Kijas J, McWilliam S, Taft RJ, Reverter A. Information
618 compression exploits patterns of genome composition to discriminate populations and highlight
619 regions of evolutionary interest. *BMC Bioinformatics*. 2014;15(1):66.
- 620 30. Gailly J-l. *gzip*. 1.3.12 ed. <http://www.gzip.org/2007>.
- 621 31. Ziv J, Lempel A. A universal algorithm for sequential data compression. *IEEE Trans Inf*
622 *Theory*. 1977;23(3):337-43.
- 623 32. Hwang B, Lee JH, Bang D. Single-cell RNA sequencing technologies and bioinformatics
624 pipelines. *Exp Mol Med*. 2018;50(8):96.
- 625 33. Parks DH, Rinke C, Chuvochina M, Chaumeil P-A, Woodcroft BJ, Evans PN, et al.
626 Recovery of nearly 8,000 metagenome-assembled genomes substantially expands the tree of life.
627 *Nat Microbiol*. 2017;2(11):1533.
- 628 34. Ross EM, Moate PJ, Marett LC, Cocks BG, Hayes BJ. Metagenomic predictions: from
629 microbiome to complex health and environmental phenotypes in humans and cattle. *PLoS ONE*.
630 2013;8(9):e73056.
- 631 35. Methé BA, Nelson KE, Pop M, Creasy HH, Giglio MG, Huttenhower C, et al. A framework
632 for human microbiome research. *Nature*. 2012;486(7402):215-21.

- 633 36. Qin J, Li R, Raes J, Arumugam M, Burgdorf KS, Manichanh C, et al. A human gut
634 microbial gene catalogue established by metagenomic sequencing. *Nature*. 2010;464(7285):59.
- 635 37. Ross E, Moate P, Maret L, Cocks B, Hayes B. Investigating the effect of two methane-
636 mitigating diets on the rumen microbiome using massively parallel sequencing. *J Dairy Sci*.
637 2013;96(9):6030-46.
- 638 38. Wang J, Jia H. Metagenome-wide association studies: fine-mining the microbiome. *Nat*
639 *Rev Microbiol*. 2016;14(8):508.
- 640 39. Fernando RL, Dekkers JCM, Garrick DJ. A class of Bayesian methods to combine large
641 numbers of genotyped and non-genotyped animals for whole-genome analyses. *Genet Sel Evol*.
642 2014;46(1):50.

643
644

645 **Figure S1: Sequence Quality per Base Pair for all Lanes of Sequencing.**

646 Box and whisker plots of sequence quality (Phred Score) at positions along the sequenced read.
647 Red, orange and green signify low, medium and high quality bases, respectively. Sequence quality
648 was high throughout the entire read; however, it did drop slightly towards the end of the read.
649 Sequence quality for *ApeKI* was more variable than for *PstI*. The third plot for *PstI* represents the
650 94 samples that were re-sequenced due to barcodes not ligating in the initial run.
651

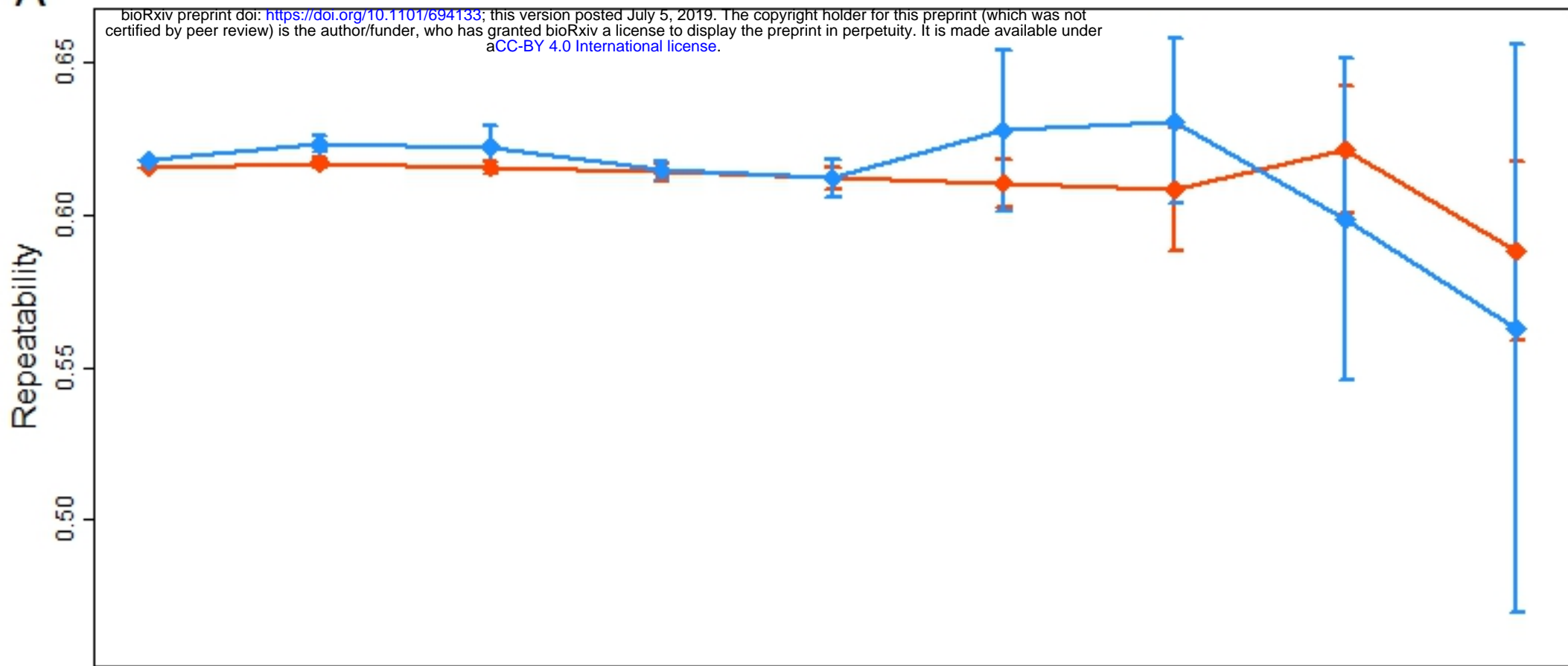
652 **File S1: Tag Filtering Comparison.**

653 This file contains an investigation into suitable tag lengths (16, 32 or 65 bp) and tag prevalence
654 thresholds (present in 10, 25, 50 or 100% of samples) for use with *ApeKI* and *PstI* restriction
655 enzymes.
656

657

■ ApeKI ■ PstI

A



B

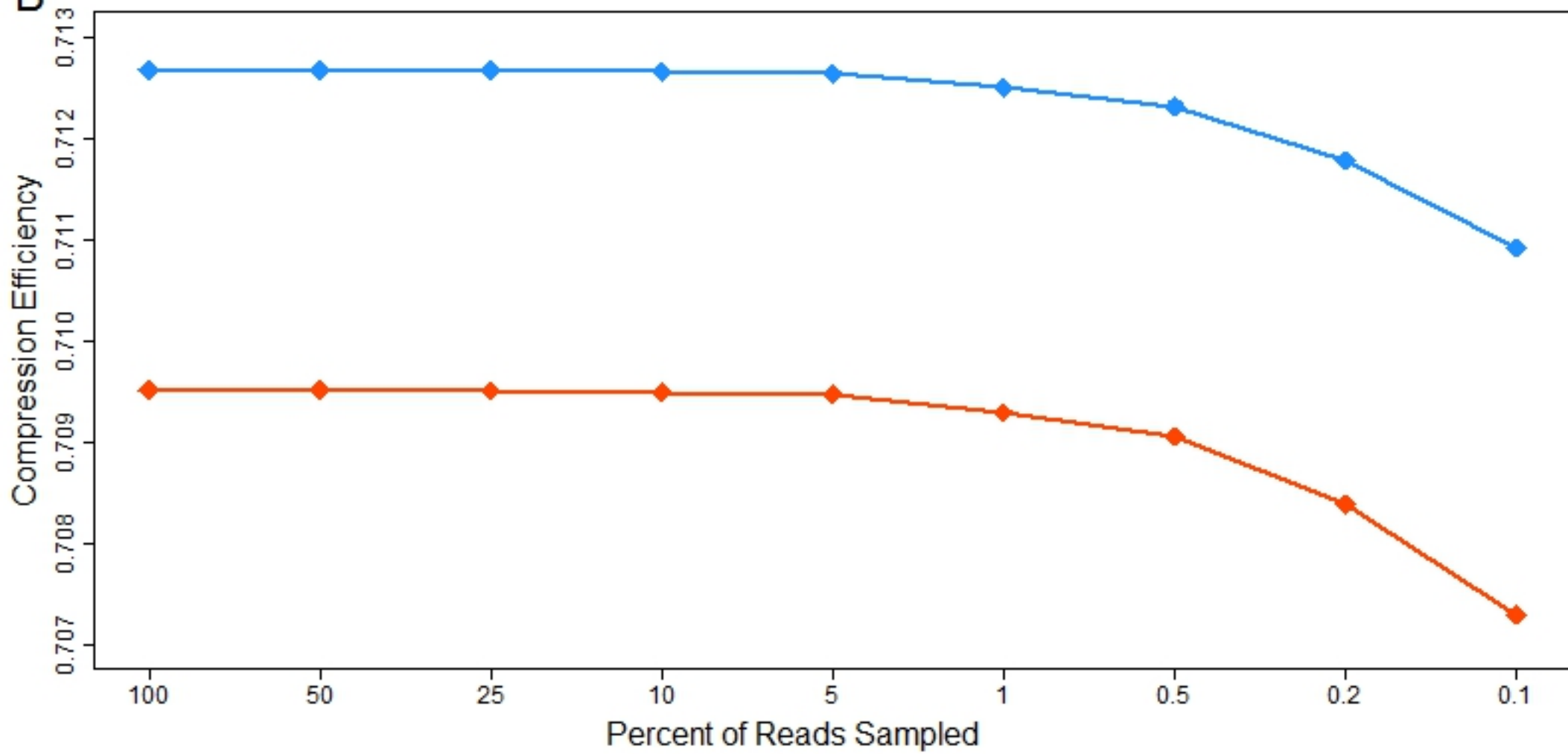


Figure 1