1   **Validating the use of bovine buccal sampling as a proxy for the rumen**

2   **microbiota using a time course and random forest classification approach**

3

4   Juliana Young,[a] Joseph H. Skarlupka,[b] Rafael Tassinari Resende,[c] Amelie Fischer,[d] Kenneth F.

5   Kalscheur,[a] Jennifer C. McClure,[a] John B. Cole,[e] Garret Suen,[b] Derek M. Bickhart[a#]

6   [#]Address correspondence to Derek Bickhart, derek.bickhart@usda.gov.

7   [a]US Dairy Forage Research Center, USDA-Agricultural Research Service, Madison,

8   Wisconsin, USA

9   [b]Department of Bacteriology, University of Wisconsin, Madison, Wisconsin, USA

10  [c]School of Agronomy, Universidade Federal de Goiás (UFG), Goiânia, GO, Brazil

11  [d]Institut de l'élevage, Beaucouse, France

12  [e]Animal Genomics and Improvement Laboratory, USDA-Agricultural Research Service,

13  Beltsville, Maryland, USA

14

15  **Running head:** Buccal swabbing as a noninvasive proxy for rumen bacteria

16

17

18

19

20

21

22

23

## ABSTRACT

Analysis of the cow microbiome, as well as host genetic influences on the establishment and colonization of the rumen microbiota, is critical for development of strategies to manipulate ruminal function toward more efficient and environmentally friendly milk production. To this end, the development and validation of noninvasive methods to sample the rumen microbiota at a large-scale is required. Here, we further optimized the analysis of buccal swab samples as a proxy for direct microbial samples of the rumen of dairy cows. To identify an optimal time for sampling, we collected buccal swab and rumen samples at six different time points relative to animal feeding. We then evaluated several biases in these samples using a machine learning classifier (random forest) to select taxa that discriminate between buccal swab and rumen samples. Differences in the Simpson's diversity, Shannon's evenness and Bray-Curtis dissimilarities between methods were significantly less apparent when sampling was performed prior to morning feeding ($P<0.05$), suggesting that this time point was optimal for representative sampling. In addition, the random forest classifier was able to accurately identify non-rumen taxa, including 10 oral and feed-associated taxa. Two highly prevalent (> 60%) taxa in buccal and rumen samples had significant variance in absolute abundance between sampling methods, but could be qualitatively assessed via regular buccal swab sampling. This work not only provides new insights into the oral community of ruminants, but further validates and refines buccal swabbing as a method to assess the rumen microbiota in large herds.

48 **IMPORTANCE**

49 The gastrointestinal tract of ruminants harbors a diverse microbial community that coevolved

50 symbiotically with the host, influencing its nutrition, health and performance. While the

51 influence of environmental factors on rumen microbes is well-documented, the process by which

52 host genetics influences the establishment and colonization of the rumen microbiota still needs

53 to be elucidated. This knowledge gap is due largely to our inability to easily sample the rumen

54 microbiota. There are three common methods for rumen sampling but all of them present at least

55 one disadvantage, including animal welfare, sample quality, labor, and scalability. The

56 development and validation of non-invasive methods, such as buccal swabbing, for large-scale

57 rumen sampling is needed to support studies that require large sample sizes to generate reliable

58 results. The validation of buccal swabbing will also support the development of molecular tools

59 for the early diagnosis of metabolic disorders associated with microbial changes in large herds.

60 **KEYWORDS**

61 Bacteria, oral community, rumen microbiota, buccal swab, machine learning, random forest

62

63

64

65

66

67

68

69

## INTRODUCTION

71      The rumen is a specialized organ found in cattle that hosts a wide diversity of

72      microorganisms from all three super kingdoms (for a review see (1, 2)). Essential to the digestion

73      of complex plant polymers by the host, the rumen microbiota consists of several species of

74      specialized fibrolytic bacteria capable of degrading lignocellulose (3). Microbial changes

75      following total rumen exchanges (4) and some preliminary genome-wide association data (5, 6)

76      suggest that the microbial community composition is unique to each individual cow and that the

77      genetics of the host animal may influence community development/maintenance in the rumen.

78      Unfortunately, the statistical determination of the extent of host-animal control over this

79      phenomenon requires a large amount of input data and rumen microbial samples are often quite

80      laborious to obtain.

81      Methods that directly sample the rumen contents of cattle are the rate-limiting step for

82      generating a population-scale metric of the rumen microbiome. The gold-standard method for

83      assessing rumen microbial contents is via rumen cannulation; however, this requires invasive

84      surgery and cannot be performed on hundreds of cows in a herd. Stomach tubing is another

85      method of sampling that provides direct access to rumen contents, but this method is labor-

86      intensive and is uncomfortable for the cow (7, 8) (7). Given the requirements for surgery or

87      labor-intensive sample collection, respectively, neither method is suitable for the development

88      of a scalable industrial product. In light of the deficiencies of these methods, buccal swabbing

89      has been proposed as a proxy for the rumen microbiota (9, 10). The ease of this method,

90      combined with high-throughput sequencing of the 16S rRNA gene and its lower cost of

91      implementation, make it a tantalizing option for obtaining population-scale rumen microbial

92      samples.

93      Buccal swabbing is a noninvasive method that takes advantage of cattle rumination, an

94      innate behavioral process that characterizes the ruminant clade of mammals (11, 12). During this

95    process, the cow regurgitates, masticates, moistens, and swallows a bolus from the rumen, which

96    is a mixture of previously ingested plant material that is resistant to prolonged chemical

97    degradation. This process exposes additional surface area of the digesting plant matter to

98    continued microbial fermentation (12). However, rumen microbes are not effaced from the

99    surface of the bolus prior to mastication, and microbial DNA in the oral cavity may constitute a

100   representative proxy of the rumen microbiota.

101   Indeed, the oral cavity has its own resident microbiota that contains both transient

102   facultative anaerobes and feed-associated microbes (13, 14) that can be concurrently sampled

103   during buccal swabbing. The identification and exclusion of these contaminants constitute a pre-

104   requisite for the use of buccal swabs as proxy for the rumen microbiota (9). In previous studies,

105   the depletion of these contaminants was performed with mathematical filtering based on the

106   comparison of the relative abundances of a given taxa between rumen and buccal swab samples

107   (9, 10). However, these approaches noted the need for further statistical and qualitative validation

108   for wide-spread adoption of the technique due to confounding factors that could impact microbial

109   taxa counts (9). This is a necessary step towards the use of buccal swabbing as an independent

110   method, as future surveys may not always have access to paired rumen samples for calibration.

111   Previous surveys have also not considered sampling time as a potential confounding

112   factor for interrogating rumen microbial counts via buccal swabbing (15–18). In the case of

113   sample time, salivary dilution and contamination with feed silage communities could impact

114   measured community composition and abundance. It is possible that there is a specific window

115   of time in which buccal swab samples best mirror the rumen contents of the sampled cow. Prior

116   to its widespread adoption as a suitable proxy for rumen sampling, buccal swabbing data must

117   be compared in a modeling experiment to identify the magnitude of these biases.

118   In this study, we apply statistical learning methods to buccal swab data obtained from 21

119   cannulated Holstein cows to identify microbial taxa that are specific to the oral cavity. We

5

120 hypothesize that the presence of non-rumen bacterial communities and the eventual salivary

121 dilution of rumen microbial DNA impacts the comparability of buccal swab samples with in-situ

122 rumen samples. We also tested if buccal swab OTU abundances can be used in regression models

123 to determine the approximate abundance of rumen microbial genera in individual animals. Our

124 analysis reveals an additional complexity in the diversity of microbes that colonize the ruminant

125 gastrointestinal tract, and we expand the future use of buccal swabs in population-scale surveys

126 of the rumen microbial community.

127 **MATERIAL AND METHODS**

128 **Animal care and use.** All animal procedures were conducted according to Research Animal

129 Resource Center (RARC) protocol A005902-A02 approved on 07/28/2017 by the University of

130 Wisconsin-Madison College of Agriculture and Life Sciences Institutional Animal Care and Use

131 Committee. This work was carried out at the US Dairy Forage Research Center Farm, Prairie du

132 Sac, WI, from 11/2017 to 06/2019 using a cohort of 21 cannulated lactating Holstein dairy cows

133 (~2.5 years old) fed a total mixed ration in a free stall barn.

134 **Sampling.** To identify the sampling time at which oral microbiota would best represent the

135 rumen microbiota, paired oral (Buccal Swab, BS) and ruminal samples (Rumen Anterior Liquid,

136 RAL; Rumen Anterior Solid, RAS; Rumen Ventral Liquid, RVL; Rumen Ventral Solid, RVS)

137 were collected from 8 cannulated Holstein cows every 2 hours over the course of 10 hours,

138 starting 1 hour prior to morning feeding (~ 9 AM) and ending just prior to evening feeding (~ 7

139 PM), totaling six time points (T1-T6). This dataset is hereafter referred to in the text as the

140 summer time course (STC; see Table 1).

141 Two other surveys of paired buccal swab and rumen content samplings were conducted on

142 different animals in the same herd at two other timepoints separated by at least three months

143 (Table 1). These datasets consist of a spring sampling (SPS; 5 cows) and a summer sampling

144 (SUS; 8 cows) taken a year prior to the STC dataset. Swabs and rumen contents were processed

145 in the same manner as listed for the time course survey, but samples were collected from animals

146 four hours after feeding (all cows in SPS) or prior to feeding (all cows in SUS), representing

147 equivalents to T4 and T1 from the time course trial, respectively. These samples were collected

148 to provide additional power for training and testing regression models (see Table 1).

149      In all trials, two swabs (Puritan PurFlock Ultra sterile flocked swab with an 80 mm break

150 point, Puritan Medical Products, Guilford, ME) were inserted in the buccal cavity of each cow

151 and were gently scraped across the inner side of the right cheek for approximately 10 seconds.

152 The buccal swabs were placed in a sterile conical tube (15 mL) containing 1 mL of sterile

153 phosphobuffer saline and stored on ice during sampling. Immediately after buccal swabbing,

154 rumen contents were collected via the rumen cannula and squeezed through double layers of

155 cheesecloth to obtain an aliquot of 40mL of rumen liquids and 50 mL of a loosely packed rumen

156 solid fraction. The solid fraction was squeezed once more to remove all liquids and the residual

157 solid material was transferred to another container. All samples were stored and transported on

158 wet ice and stored at -80 °C until processing and DNA extraction.

159 **DNA extraction and sequencing.** Total genomic DNA was extracted from buccal swab, rumen

160 liquid, and rumen solid samples as previously described (19). Sequencing was performed at the

161 UW-Madison Biotechnology Center using the $2 \times 250$ bp paired-end method on an Illumina

162 MiSeq following manufacturer's guidelines (Illumina, Inc., San Diego, CA, USA). Detailed

163 methods about and the library preparation and sequencing can be found in Skarlupka et al. (20).

164 **Bioinformatics analysis.** DNA sequences were analyzed using mothur (v1.39.0) (21) as

165 described previously (22). Coverage was assessed by Good's index (23) and samples that

166 displayed coverage less than 93% were discarded prior to normalization. To address differences

167 in sequencing depths, the operational taxonomic unit (OTU) table was normalized by

168 subsampling sequences to the sample with the smallest number of sequences and then

169 normalizing across samples to produce equal sequence counts (3,000 sequences per sample). The

170    normalized OTU table was used in further analyses as well as to calculate alpha diversity indices

171    (i.e., Chao1 (24), Shannon (25), and Simpson (26)), Bray-Curtis dissimilarity index (27) as well

172    as the relative abundance (reads/total reads in a sample x 100) of OTUs in each sample. Alpha

173    diversity indices were calculated in mothur (v1.39.0) (21) whereas Bray-Curtis dissimilarity

174    index was calculated using function vegdist available at R package vegan (v2.5-6) (28)

175    **Statistical analysis.** All statistical analyses were performed in R (v3.6.1) and source code to

176    reproduce these analyses is available in Supplementary Materials. Measurements of α-diversity

177    (Chao1's richness, Shannon's evenness and Simpson's diversity) and absolute abundance (i.e.,

178    sequence read counts) of OTUs detected in at least 80% of all samples, were assessed for

179    normality and were found to follow a non-normal distribution. Differences in the alpha diversity

180    indices and OTU absolute abundance values were analyzed, respectively, under Gamma and

181    Poisson distributions, using a repeated-measure generalized linear mixed model estimated via

182    penalized quasi-likelihood (29):

183    $$Y_i^* = X_i\beta + Z_i b + \varepsilon_i$$

184    where $Y_i^* = (y_{i,1,1}^*, \ldots, y_{i,n_i,1}^*, \ldots, y_{i,n_i,m_i}^*)$ is a vector of Gamma- or Poisson-transformed of alpha

185    diversity indices or OTU counts; $X_i$ is a design matrix relating individual observations to levels

186    of fixed effects, $\beta$ is a vector of fixed effects (i.e., sampling time, sample type, and their

187    interaction), $Z_i$ is the incidence-matrix on random effects, $b$ is the vector of random animal

188    effects; $\varepsilon_i$ is a vector of random error terms. The resulting ANOVA P-values were adjusted for

189    false discovery rate (FDR) using the Benjamini-Hochberg method, and values ≤0.05 were

190    considered significant. Pairwise comparisons among the Least Squares Means (LSMEANS)

191    were performed using Tukey's Honest Significant Difference (Tukey HSD) method. In the

192    presence of significant interaction effects, the LSMEANS of the sample types were compared

193    within each sampling time. These analyses were performed using functions available at R

194    package fitdistrplus (v1.0-14), MASS (v7.3-51.5), lsmeans (v2.30-0), and ggplot2 (v3.2.1) (30–

195    33).

196    To visually explore the degree of dissimilarity between bacterial composition of oral and

197    rumen samples collected at six distinct sampling times, Principal Coordinates Analysis (PCoA)

198    was conducted on the Bray-Curtis distance matrix (27). In addition, Permutational Multivariate

199    Analysis of Variance (PERMANOVA, nperm=1000) (34) with *post hoc* test using Benjamini-

200    Hochberg correction was performed to assess differences in the composition of bacterial

201    communities according to sample type, time points and their interaction. These analyses were

202    performed using functions available in the R packages ggplot2 (v3.2.1), vegan (v2.5-6), and

203    EcolUtils (v0.1) (28, 35, 36) .

204    To identify taxa that discriminate between oral and rumen samples, a Random Forest

205    classifier was trained on a random selection of 70% (162 samples) of the database composed of

206    232 samples and 2,031 OTUs and validated using the remaining 30% (70 samples). Only OTUs

207    with relative abundance ≥ 0.05% present in at least one sample were included as input. The

208    number of trees was set to 500, while the number of variables available for splitting at each tree

209    node (mtry) was tuned and accuracy was used to select the optimal model using the largest value.

210    In addition, to evaluate the capability of our model to predict on independent dataset, we adopted

211    a repeated k-fold cross validation method (10-fold repeated 3 times). Prediction performance

212    metrics (i.e., accuracy, sensitivity, specificity, precision and recall) and a confusion matrix were

213    calculated and summarized by sample type.  Finally, the Mean Decrease in Gini (i.e., Gini index)

214    was used to calculate the variable importance score (VIMP) and select bacterial OTUs that were

215    most predictive of sample types. To that end, we used the function varImp ((37)) that

216    automatically scales the importance scores to be between 0 and 100. These results were plotted

217    to show the most important sample type-associated bacterial OTUs with VIMP score ≥50%.

218 These analyses were performed using the R packages randomForest (v4.6-14) and caret (v6.0-

219 85) (37, 38).

220 In order to evaluate if abundance of oral microbiota can be used to predict the abundance of

221 rumen microbiota, we tested distinct regressions models (i.e., random forest, Random

222 generalized linear model, GLMM zero-inflated quasi-Poisson). These analyses were performed

223 using the R packages MASS (v7.3-51.5), caret (v6.0-85), randomForest (v4.6-14), and

224 randomGLM (v1.02-1) (37–39).

225 **Data Availability.** The raw sequence reads from all samples analyzed in this study are available

226 on the NCBI Sequence Read Archive (https://www.ncbi.nlm.nih.gov/sra/) under the Bioproject

227 accession number: PRJNA623113.

228 **RESULTS**

229 **Amplicon sequencing and quality control.** To provide metrics for quality control and optimal

230 parameter selection, we sampled buccal and rumen contents from several cohorts of cannulated

231 cattle (Table 1). To test if a difference in rumen sampling site had closer resemblance to swab

232 samples, rumen strata (solids and liquids) from the anterior and ventral side of the rumen lumen

233 were simultaneously collected. Samples are hereafter referred to by acronyms that denote their

234 sample type (BS and R for buccal swab and rumen, respectively), and their location and content

235 in the case of rumen samples (A, V, S, and L, for anterior, ventral, solid, and liquid, respectively).

236 For example, the acronym RAL refers to a rumen anterior liquid sample. All samples were

237 sequenced using the same methods and resulting data were processed using the same pipeline.

238 After sequence quality filtering and normalization, a total of 1,392,036 reads (mean

239 $6,000.155 \pm 132.615$ SD per sample) and 196,258 OTUs (mean $845.94 \pm 199.411$ per sample)

240 were obtained from 232 buccal, rumen solid, and rumen liquid samples in total. Good's coverage

241 estimation prior to normalization ($0.969 \pm 0.034$ per sample) was deemed adequate and indicated

242    that sequences sufficiently covered the diversity of the bacterial communities in our study. A full

243    summary of sequencing statistics as well as rarefaction curves divided by sample type and time

244    point is shown in Fig. S1 and Table S1.

245        Taxonomic composition analysis of the bacterial communities revealed a total of 2,031

246    OTUs (mean 112.46 ± 32.91 SD) present at relative abundances ≥0.05% and representing 20

247    phyla, 116 families and 279 genera. The average percentage of sequences unassigned to any

248    phylum, family, or genus were $0.19 \pm 0.15$, $1.15 \pm 0.45$, and $10.49 \pm 2.69$, respectively. The most

249    abundant OTUs, summarized at the phylum, family and genus levels according to sampling time

250    and type are shown in Fig. S2.

251    **Time course analysis and sampling method comparability.** We first sought to identify the

252    effects of sampling method on the composition of observed microbial communities in the rumen.

253    For this analysis, we used paired rumen strata (solid and liquid) and buccal swab samples taken

254    from the STC cohort (see Table 1) in 2-hour intervals, with the first time point (T1) taken 1 hour

255    prior to feeding. Rather than seeking a singular optimal time for sampling, we investigated the

256    possibility that there are periods where the buccal microbial community may be less

257    representative in terms of species prevalence and relative abundance of the rumen community.

258        Sampling type (i.e., buccal swabbing vs. rumen cannula sampling) had the largest effect

259    on observed microbial content, as expected. Alpha diversity analysis revealed that Chao1

260    richness (number of species) varied significantly with sample type ($P = 0.014$) but not sampling

261    time ($P = 0.208$) or the interaction of these two factors ($P = 0.091$). Shannon's evenness

262    (population density) and Simpson's diversity (richness and abundance) varied with sample type

263    ($P < 0.001$; $P < 0.001$), sampling time ($P = 0.021$; $P = 0.047$), and the interaction of these factors

264    was significant ($P < 0.001$; $P < 0.001$). Regardless of sampling time, buccal swab samples

265    displayed lower richness (i.e., Chao1) and evenness (i.e., Shannon), but higher diversity (i.e.,

266    Simpson) when compared to all types of rumen samples (Tukey HSD < 0.05). Regardless of

11

267    sample type, bacterial communities sampled at T3 and T4 displayed the lowest and highest

268    Shannon's evenness, respectively (Tukey HSD < 0.05). Significant differences in Shannon's

269    evenness and Simpson's diversity were not observed between others timepoints (Tukey

270    HSD < 0.05; Table S2). In regard to interaction terms, we observed that buccal swabs collected

271    at T1 and T4 displayed similar evenness and diversity to all types of rumen samples. In contrast,

272    buccal swab samples from other time points (T2, T3, T5 and T6) displayed lower evenness but

273    higher diversity, relative to rumen samples (Tukey HSD < 0.05; Table S2).

274    We used PCoA to visually inspect the similarity of buccal swab samples to contemporary

275    rumen cannula samples. In general, rumen samples grouped by phase (i.e., L vs S) rather than

276    location (i.e., A vs V). Additionally, we found that bacterial communities from buccal swab

277    samples obtained just prior to morning feeding (T1) grouped most closely to rumen solid samples

278    (RAS + RVS) (Fig. 1). Moreover, ordination plots showed that T3 had the most pronounced

279    differences between swab and rumen samples. The presence of higher OTU counts of silage-

280    associated microbes belonging to the *Lactobacilli* in T3 suggest that feed contamination was a

281    major contributor to this discrepancy (Figs. 2 and S2\).

282    PERMANOVA showed that Bray-Curtis dissimilarities in the composition of bacterial

283    communities were significantly driven by sampling time (R squared= 0.044, P< 0.001), sample

284    type (R squared= 0.284, P< 0.001), as well as by the interaction of these two factors (R squared=

285    0.106, P< 0.001). Pairwise comparisons between sample types showed that the composition of

286    BS samples differs significantly from all types of rumen samples (P=0.010). In addition, we

287    found that bacterial composition at sampling time T1 was significantly different from T3 (P =

288    0.015) and T5 (P = 0.045). Lastly, comparisons between sample types within each sampling time

289    indicates that the composition of bacterial communities in BS samples is similar to those

290    observed in the RAS samples only at T1 (P = 0.054), confirming the clustering observed in the

291    PCoA (Fig.1 and Table S3).

292   In addition to compositional dissimilarity, we assessed differences in the absolute

293   abundance (i.e., read counts) of 277 bacterial OTUs (prevalence of at least 80% of all samples)

294   in response to sampling time, sample type and the interaction of these two factors (Figs. 5, 6,

295   and Table S7).  Overall, most of the variance in the absolute abundance of bacterial communities

296   in our study was ascribed to interaction terms given that 240 OTUs varied simultaneously with

297   sampling time and sample type. Meanwhile, the differences ascribed to main effects were far

298   less apparent, given the abundance of only 38 and 20 OTUs that varied independently in response

299   to sample type and sampling time, respectively (Table S7).

300   Comparisons between sample types within each sampling time showed that fewer OTUs

301   had significantly different absolute abundance between buccal and rumen samples taken at T1

302   followed by T4 and T6 (Fig. 5A). At these particular time points, the significant differences in

303   the absolute abundance of OTUs between buccal swab and rumen samples were less pronounced

304   than observed at other sampling times (Tukey HSD $\leq$ 0.05; Fig. 5B). In contrast, greater

305   significant differences in the absolute abundance of OTUs between BS and all rumen samples

306   were observed at T3 followed by T5 and T2 (Figs. 5A, 5B, and Table S8).

307   In addition, no significant differences in the absolute abundance of OTUs between RAS

308   and RVS were observed at T1, T2 and T4. However, some OTUs varied in absolute abundance

309   between RAL and RVL at others sampling times, mainly at T3 followed by T5 and T6 (Fig. 6A

310   and Table S8). Pronounced differences in the absolute abundance of several OTUs between

311   liquids and solids contents were observed at all time points. Specifically, the majority of the

312   OTUs sampled at T1 and T2 displayed higher absolute abundance in rumen liquids than in rumen

313   solids (i.e., RAL vs. RAS and RVL vs. RVS) while the opposite was observed at other time

314   points (Fig. 6A, 6B, and Table S7).

315   Regardless of sample type, comparisons performed between sampling times showed that

316   the absolute abundance of bacterial OTUs were significantly lower at T3 and T5 in comparison

317  to the other time points, particularly with T4 and T1 (Figure S3 and Table S8). Finally,

318  comparisons performed between sample types showed that absolute abundance of bacterial

319  OTUs were significantly lower in buccal swabs than all types of rumen samples (Tukey HSD $\leq$

320  0.05), regardless of sampling time. These differences were less apparent when buccal swab and

321  rumen solids were compared (see Figure S3 and Table S8). However, a few exceptions were

322  observed for OTUs assigned to Prevotellaceae_Ga6A1_group and Succinivibrionaceae_UCG-

323  002, whose absolute abundance were significantly higher in BS in comparison to rumen liquids

324  (RAL or RVL; Tukey HSD < 0.05) (Table S8).

325  **Random forest classifier analysis.** We next sought to identify key microbial taxa present in the

326  oral microbial community that contributed to discrepancies observed in our ordination plots. To

327  statistically distinguish between taxa that had differences in relative abundance in each sample

328  type, we trained a random forest classifier model using the STC cohort samples. Random forest

329  is a supervised learning algorithm which uses ensemble learning method (i.e., combine several

330  trees base algorithms) to construct better predictive performance (for a review see (38, 40) and

331  has been widely and successfully employed for classification and regression purposes. In a

332  classification problem, the algorithm returns a list of predictor variables (i.e., bacterial OTUs)

333  that can be ranked according to their individual importance (i.e., VIMP score) in classifying the

334  data.

335  Our preliminary analyses showed that the overall performance of the random forest classifier

336  using five classification categories for sample type (BS, RAL, RAS, RVL, and RVS) was quite

337  low (Accuracy 58.6% and Kappa 48.2%), even after estimation and tuning of model hyper-

338  parameters (Table S4). This result supports the observation of high similarity between bacterial

339  communities from rumen solids (RAS and RVS) and liquids samples (RAL and RVL) from

340  different rumen lumen areas as observed in the PCoA (Fig. 1). We found improved classifier

341  accuracy when rumen samples were merged based on rumen content strata (liquids and solids)

342 into a single type in the training and testing sets (collectively referred to as RL and RS,

343 respectively). This merger unbalanced our training set by providing a two-fold increase in rumen

344 categories (RL and RS = 95 samples each), and we thus implemented a re-sampling method for

345 future model training to prevent misclassification of our minority class (BS = 42 samples). We

346 tested three additional re-sampling methods (i.e., under-sampling, over-sampling, and Synthetic

347 Minority Over-sampling Technique, SMOTE) to prevent classification bias towards the majority

348 classes (41, 42). The results showed that random forest trained with additional re-sampling using

349 the SMOTE had higher performance metrics than the other methods (Table S5).

350 Our final model was able to predict sample type-associated bacterial features with high accuracy

351 (97.78% ± 3.7%) and Cohen's kappa values (96.3% ± 5.4%). Cohen's kappa is a frequently used

352 statistic to assess the performance of machine learning models under a multi-class classification

353 problem and or unbalanced data (43, 44). Other performance metrics such as sensitivity,

354 specificity, precision and recall were also calculated for each sample type and are presented in

355 Table S5. Additionally, our classifier returned the variable importance score (VIMP), as a

356 function of the Mean Decrease in Gini, of each bacterial OTU, which can be used to discriminate

357 between oral and rumen samples (Table S6). Thus, higher values of VIMP score expressed as a

358 percentage indicate higher feature importance (i.e., bacterial OTU) in discerning between classes

359 and, in our case, between sample types.

360 **OTU categorization based on variable importance estimates.**

361     Bacterial OTUs with high VIMP scores (≥ 50% mean decrease Gini) displayed patterns

362 that allowed for manual categorization. Based on average taxon prevalence per sample type and

363 sampling time, we categorized these OTUs into three categories: core, oral, and rumen (Table 2,

364 Fig. 3 and see Table S6 in the supplementary material). The remaining OTUs whose VIMP score

365 was lower than 50% were also categorized for the sake of completeness but were not analyzed

366 further (Table S6). The core category consisted of OTUs that displayed moderate to high

367    prevalence (≥60 to 100%) in all sample types (both rumen and buccal) consistently across

368    timepoints. The rumen category was defined as the community well represented (prevalence

369    ≥75%) in rumen liquids and/or solids, and was underrepresented in buccal swab samples

370    (prevalence<60%) at all time points (Fig. 3, Table 2 and Table S6). Finally, the oral group

371    consisted of OTUs well represented in buccal swab samples (prevalence≥60%) but were either

372    absent or underrepresented in the rumen samples (<60% prevalence) across time points. The oral

373    group was found to contain silage community microbes (i.e., *Lactobacilli*) at time points where

374    feed was provided to the animals (e.g., T3, see Fig. 4), further supporting our classification and

375    the model's accuracy.

376    In the core group, we identified two OTUs (VIMP>80%) assigned to the genus

377    *Prevotella*_1 (Fig. 3 and Table 2) that displayed high prevalence in both buccal swab and rumen

378    (liquid and solid) samples. The absolute abundance of these taxa was significantly lower (Tukey

379    HSD≤0.05) in buccal swabs than in rumen samples (Tables S6 and S7). This suggests that these

380    taxa can be reliably sampled via swabbing but that their absolute abundances are greatly biased

381    compared to the paired rumen samples.

382    We also identified taxa in the families Neisseriaceae, Pasteurellacea, Micrococcaceae,

383    and Planococcacea, as well as in the genera *Streptococcus*, *Jeotgalicoccus*, and *Bibersteinia*,

384    which displayed moderate to high VIMP scores (≥ 50%) and were assigned to the oral category.

385    These taxa were overrepresented in terms of prevalence and abundance in buccal swab samples

386    and displayed very low or zero abundance in rumen liquid and solid samples (Fig. 3 and Tables

387    2 and S6). In addition, we observed that several oral taxa (i.e., *Oceanobacillus*, *Lactobacilli*,

388    *Lachonoclostrium*, *Leuconostoc*, *Rothia*, and *Proteus*) were underrepresented in terms of

389    abundance and prevalence at specific time points, including T1, T4 and T6, relative to time points

390    T2, T3 and T5 (Fig. 4 and Table S6).

391    Finally, the classifier also selected rumen strata OTUs that have lower relative abundance

392    in the buccal swab samples (rumen category). Several were specific to rumen liquids (0405-p-

393    1088-a5_gut_group, *Howardella*, Ruminococcaceaa_ge, *Synergistes*, Prevotellaceae_UCG-001,

394    Rikenellaceae_RC9_gut_group) and others were derived from the rumen solids

395    (*Ruminoccocus*_1, Prevotellaceae_UCG-001 and *Oribacterium*) whose overall importance was

396    ≥33% (Fig. 3 and Tables 2 and S6).

397    **Random forest regression analysis.** We next sought to test whether the abundance of OTUs

398    found in buccal swab samples could be used to predict the abundance of rumen OTUs. We tested

399    the ability of four linear models (random forest regression, three log-linear models with either a

400    Poisson distribution, zero inflated, or random generalized linear model (RGLM)) to characterize

401    the relationship between bacterial OTUs of paired buccal swab and rumen liquid samples. In

402    order to provide additional data for our training regression models, we incorporated data from

403    21 cows sampled in two other surveys (Table 1) processed with the same methods used for the

404    time course study. It is important to note that random forest regression was performed using

405    sequence relative abundances whereas log-linear models use sequence absolute abundance (i.e.,

406    number of reads) for each OTU, assuming a Poisson distribution of read counts. Our random

407    forest and Poisson regression model converged, but exhibited low accuracy in cross-validation

408    studies as shown by a low coefficient of determination (R-Squared = 0.39 ± 0.05) and high Root

409    Mean Square Error (RMSE = 0.28 ± 0.09). We attempted to tune additional parameters in the

410    random forest model, but were unable to achieve an accuracy R-Squared above of 0.42 ± 0.07

411    on a per-OTU basis. Conversely, zero inflated and RGLM trials failed to converge, despite

412    several attempts to filter the OTU tables and tune model parameters. These results may be related

413    to our use of a small dataset as well as the non-linear relationship between the buccal swab and

414    rumen OTU abundance/counts on a per-sample basis.

415

416    **DISCUSSION**

417    In this study we evaluated the ability of the buccal swabbing method to describe bacterial

418    communities found in two types of rumen samples taken at six distinct sampling times over the

419    course of ten hours. Buccal swab samples are an attractive alternative to more labor-intensive

420    methods of sampling the rumen microbial community, but may suffer from bias due to

421    contamination by the surrounding oral community (9, 10). We first sought to identify the effect

422    of sampling time on buccal swab community composition as we hypothesized that animal

423    rumination patterns and salivary flow may change the relative abundance of key members of the

424    rumen community.

425    Our time course analysis suggests that there is a small, but statistically significant, effect

426    of sampling time on the comparisons of several buccal swab microbial taxa with contemporary

427    rumen samples from the same animal. After dividing sampling times into two-hour intervals, we

428    sampled buccal contents from each animal just prior to the start of morning feeding (T1), within

429    regular intervals during and after feeding (T2, T3, T4, and T5), and prior to evening feeding

430    (T6). We found that the only major outlier was at time point 3 (T3), where the greatest

431    dissimilarities in the bacterial communities between buccal swabs and rumen samples were

432    observed. It is possible that additional contamination by the silage microbial community and

433    increased salivary flow induced by feeding changed the relative abundance of key rumen taxa in

434    the oral samples of cows sampled at T3. This is evidenced by the presence of *Lactobacilli* from

435    silage communities in the buccal swabs, but not in the rumen contents (Fig. 2 and 4). Our results

436    support a hypothesis that there are brief windows of time in which buccal swab data best

437    represent contemporary rumen microbial data. This means that future surveys will need to record

438    time of sampling relative to animal feeding in order to standardize results.

439    We also tested the possibility that buccal swab samples may be compositionally similar

440    to rumen content fractions taken from different positions in the rumen (i.e., Anterior vs Ventral).

18

441    Our comparisons of sampling time and sample types found no differences between the bacterial

442    communities of the anterior and ventral rumen microbial communities, which prevented us from

443    finding such an association (Fig. 1). This result is likely associated with the constant mixing of

444    rumen contents due to the contractions of the reticulorumen, which would result in

445    indistinguishable variation in our observed rumen microbial OTU counts (12). This finding

446    contrasts from previously published work that identified noticeable differences in sample

447    composition from five different locations of the rumen lumen via PCR DGGE surveys (45). We

448    therefore cannot rule out the possibility that our sampling and analysis methods could not

449    identify the small effects that these locations have on the community.

450    We also found greater similarity between bacterial taxa present in buccal swabs and

451    rumen solids than in rumen liquids (Fig. 1). We suspect that this reflects a key stage of the

452    rumination process whereby, immediately after regurgitation, the liquid fraction of the bolus is

453    swallowed (12). It is possible that the bacterial taxa that are predominant in the liquid-phase of

454    the rumen contents are evacuated from the oral cavity early in the process of rumination. During

455    mastication of the bolus, bacteria from solid-phase of rumen contents are more likely to adhere

456    to oral mucosal surfaces and are more likely to be sampled during buccal swabbing.

457    In order to identify non-rumen taxa in buccal swab samples, we employed a machine

458    learning classifier to assist in the filtering of oral and silage microbial communities in buccal

459    swab samples. As has been noted previously (9), the presence of the commensal oral microbial

460    community in buccal swab samples prevents direct comparisons between rumen content samples

461    and buccal swabs and must be filtered from buccal swab samples prior to analysis using manual

462    and mathematical methods (9, 10). By using a random forest classifier, we were able to assign

463    importance estimates to individual microbial taxa based on their use as a feature in our

464    classification models, as has been done previously (46, 47). The top OTUs, after variable

465    importance analysis, consisted of microbes that were oral-specific (oral, n = 10), rumen-biased

466     (rumen, n = 12), and those with high prevalence regardless of sample type but varied based on

467     relative abundance (core, n = 2). These findings support our observations of the influence of

468     sample type on OTU ~~relative~~ abundance, and also identified members of the oral-microbial

469     community that were prevalent only in buccal swab samples. In addition, the top OTUs identified

470     by our VIMP analysis included two members of the *Prevotella*, which were found to vary

471     substantially between buccal and rumen samples (Table S7). These two OTUs were prevalent in

472     all samples and at all time points; however, their ~~relative~~ abundance in buccal swabs was lower

473     than in the rumen samples. These differences were far less apparent at T1, which as just prior to

474     feeding, than at any other sampling time. This observation of similarity at only one time point

475     implies that sampling time had a large effect on the estimated ~~relative~~ abundance of this clade,

476     as confirmed by our ANOVA.

477         The OTUs present within the oral category represent taxa that are poorly represented in

478     buccal swab samples. Indeed, we identified commensal oral microbes from the genus *Rothia* that

479     were present only in the buccal swab samples (the oral category). These taxa can be safely

480     removed from future buccal swab surveys. We also identified several oral taxa (i.e.,

481     *Lactobacillus*, Chryseobacterium, *Burkholderiaceae*, *Oceanobacillus*) that were prevalent at

482     some time points, and underrepresented or even absent at others (Fig. 4) showing that sampling

483     time is a critical factor to be considered in future studies. The higher prevalence of these taxa

484     during (T2) and immediately after (T3) feeding suggests that these sampling times will result in

485     buccal swab data that is least representative of the rumen contents of the animal.

486         Our use of random forest classifiers suggests that machine-learning methods can be used

487     to approximate the rumen microbial community at the time of sampling. More accurate

488     estimation of these communities will be beneficial to rumen microbial ecology experiments that

489     suffer from low sample counts. However, we were unable to achieve an acceptable rate of error

490     (measured via residual error of observed and predicted OTU counts) from our regression

491    analysis. We found that multicollinearity of predictors and weak linear association between oral

492    and rumen OTUs prevented accurate regression. We suspect that other factors (i.e., sampling

493    time, herd, diet) must be controlled for in the modelling of these data, as evidenced by

494    significance of sampling time and interaction terms in our PERMANOVA and ANOVA.

495    Moreover, it is possible that the taxonomic affiliation of our OTU counts could be masking

496    individual species level abundances that provide far more variance than expected for the

497    regression model. Similarly, our genus-level assignments could also contain inaccuracies due to

498    strain abundance differences in the oral cavity vs. the rumen contents.

499    Finally, we cannot rule out the possibility that several OTUs are metabolically active

500    (i.e., facultative aerobes) in both locations and can proliferate in the oral cavity, thereby creating

501    a non-linear relationship between their abundance estimates in buccal swabs and rumen contents.

502    While this presents an impediment to the use of buccal swabs for classical microbial ecology

503    experiments, we note that buccal swab data is still useful for other associative analysis. The

504    ability to collect large numbers of samples from a diverse cohort of animals can present an

505    opportunity for associations of microbial profiles with animal production and performance

506    metrics including milk production, health and even fertility phenotypes. Such experiments would

507    benefit from the removal of biases that we identified in this survey.

508    In summary, we have identified significant effects of sampling time and sample type on

509    the composition of rumen microbial OTU counts derived from buccal swabs and rumen samples.

510    The buccal swab samples were prone to significant bias based on the time of sampling, with

511    specific time points showing higher prevalence of the oral- or feed-associated microbial

512    community than others. For future surveys using buccal swabs as a proxy for rumen microbial

513    counts, we recommend buccal sampling at least 2 hours prior or four hours after feeding. Our

514    data also suggests that a portion of the rumen microbial community will remain inaccessible to

515 buccal swab samples; however, this bias may not necessarily impede future association studies

516 with host animal phenotypic traits.

517

518

## ACKNOWLEDGMENTS

529 Mention of trade names or commercial products in this article is solely for the purpose of

530 providing specific information and does not imply recommendation or endorsement by the

531 USDA. The USDA is an equal opportunity provider and employer.

535

536

537

538

539

540

541

## REFERENCES

1. Weimer PJ. 2015. Redundancy, resilience, and host specificity of the ruminal microbiota: Implications for engineering improved ruminal fermentations. Front Microbiol 6:296. https://doi.org/10.3389/fmicb.2015.00296.

2. Bickhart DM, Weimer PJ. 2018. Symposium review: Host–rumen microbe interactions may be leveraged to improve the productivity of dairy cows. J Dairy Sci101:7680–7689. https://doi.org/10.3168/jds.2017-13328.

3. Neumann AP, Suen G. 2018. The Phylogenomic Diversity of Herbivore-Associated Fibrobacter spp. Is Correlated to Lignocellulose-Degrading Potential. mSphere 3(6): e00593-18. https://doi.org/10.1128/mSphere.00593-18.

4. Weimer PJ, Stevenson DM, Mantovani HC, Man SLC. 2010. Host specificity of the ruminal bacterial community in the dairy cow following near-total exchange of ruminal contents. J Dairy Sci93:5902–5912. https://doi.org/10.3168/jds.2010-3500.

5. Li F, Li C, Chen Y, Liu J, Zhang C, Irving B, Fitzsimmons C, Plastow G, Guan LL. 2019. Host genetics influence the rumen microbiota and heritable rumen microbial features associate with feed efficiency in cattle. Microbiome 7(1):92. https://doi.org/10.1186/s40168-019-0699-1.

6. Wallace RJ, Rooke JA, McKain N, Duthie CA, Hyslop JJ, Ross DW, Waterhouse A, Watson M, Roehe R. 2015. The rumen microbial metagenome associated with high methane production in cattle. BMC Genomics 16:1–14. https://doi.org/10.1186/s12864-015-2032-0.

7. Henderson G, Cox F, Kittelmann S, Miri VH, Zethof M, Noel SJ, Waghorn GC, Janssen PH. 2013. Effect of DNA extraction methods and sampling techniques on the apparent structure of cow and sheep rumen microbial communities. PloS One 8(9). https://doi.org/10.1371/journal.pone.0074787.

8. Paz HA, Anderson CL, Muller MJ, Kononoff PJ, Fernando SC. 2016. Rumen Bacterial Community Composition in Holstein and Jersey Cows Is Different under Same Dietary Condition and Is Not Affected by Sampling Method. Front Microbiol 07:1206. https://doi.org/10.3389/fmicb.2016.01206.

578    9.  Kittelmann S, Kirk MR, Jonker A, McCulloch A, Janssen PH. 2015. Buccal swabbing as a
579        noninvasive method to determine bacterial, archaeal, and eukaryotic microbial community
580        structures in the rumen. Appl Environ Microbiol 81:7470–7483.
581        https://doi.org/10.1128/AEM.02385-15.

582

583   10. Tapio I, Shingfield KJ, McKain N, Bonin A, Fischer D, Bayat AR, Vilkki J, Taberlet P,
584        Snelling TJ, Wallace RJ. 2016. Oral Samples as Non-Invasive Proxies for Assessing the
585        Composition of the Rumen Microbial Community. PloS One 11:e0151220.
586        https://doi.org/10.1371/journal.pone.0151220.

587   11. Lindström T, Redbo I. 2000. Effect of feeding duration and rumen fill on behaviour in dairy
588        cows. Appl. Anim 70:83–97. https://doi.org/10.1016/s0168-1591(00)00148-9.

589

590   12. Beauchemin KA. 2018. Invited review: Current perspectives on eating and rumination
591        activity in dairy cows. J Dairy Sci 101:4762–4784. https://doi.org/10.3168/jds.2017-13706.

592

593   13. RamÅ¡ak A, Peterka M, Tajima K, Martin JC, Wood J, Johnston MEA, Aminov RI, Flint
594        HJ, AvguÅ¡tin G. 2000. Unravelling the genetic diversity of ruminal bacteria belonging to
595        the CFB phylum. FEMS Microbiol Ecol 33:69–79. https://doi.org/10.1111/j.1574-
596        6941.2000.tb00728.x.

597

598   14. Creevey CJ, Kelly WJ, Henderson G, Leahy SC. 2014. Determining the culturability of the
599        rumen bacterial microbiome. Microb Biotechnol 7:467–479. https://doi.org/10.1111/1751-
600        7915.12141.

601

602   15. Duffield T, Plaizier JC, Fairfield A, Bagg R, Vessie G, Dick P, Wilson J, Aramini J, McBride
603        B. 2004. Comparison of techniques for measurement of rumen pH in lactating dairy cows. J
604        Dairy Sci 87:59–66. https://doi.org/10.3168/jds.S0022-0302(04)73142-2.

605

606   16. Jewell KA, McCormick CA, Odt CL, Weimer PJ, Suen G. 2015. Ruminal bacterial
607        community composition in dairy cows is dynamic over the course of two lactations and
608        correlates with feed efficiency. Appl Environ Microbiol 81:4697–4710.
609        https://doi.org/10.1128/AEM.00720-15.

610

611   17. de Mulder T, Goossens K, Peiren N, Vandaele L, Haegeman A, de Tender C, Ruttink T, de
612        Wiele T van, de Campeneere S. 2016. Exploring the methanogen and bacterial communities
613        of rumen environments: solid adherent, fluid and epimural. FEMS Microbiol Ecol
614        93(3):fiw251. https://doi.org/:10.1093/femsec/fiw251.

615

616   18. Ji S, Zhang H, Yan H, Azarfar A, Shi H, Alugongo G, Li S, Cao Z, Wang Y. 2017.
617        Comparison of rumen bacteria distribution in original rumen digesta, rumen liquid and solid
618        fractions in lactating Holstein cows. J ANIM SCI BIOTECHNO 8(1):16.
619        https://doi.org/doi:10.1186/s40104-017-0142-z.

620

621   19. Stevenson DM, Weimer PJ. 2007. Dominance of Prevotella and low abundance of classical
622        ruminal bacterial species in the bovine rumen revealed by relative quantification real-time
623        PCR. Appl Microbiol Biotechnol 75:165–174. http://dx.doi.org/10.1007/s00253-006-0802-
624        y.

625

20. Skarlupka JH, Kamenetsky ME, Jewell KA, Suen G. 2019. The ruminal bacterial community in lactating dairy cows has limited variation on a day-to-day basis. J ANIM SCI BIOTECHNO 10:66. https://doi.org/10.1186/s40104-019-0375-0.

21. Schloss PD, Westcott SL, Ryabin T, Hall JR, Hartmann M, Hollister EB, Lesniewski RA, Oakley BB, Parks DH, Robinson CJ, Sahl JW, Stres B, Thallinger GG, van Horn DJ, Weber CF. 2009. Introducing mothur: Open-source, platform-independent, community-supported software for describing and comparing microbial communities. Appl Environ Microbiol 75:7537–7541. https://doi.org/10.1128/AEM.01541-09.

22. Dill-Mcfarland KA, Breaker JD, Suen G. 2017. Microbial succession in the gastrointestinal tract of dairy cows from 2 weeks to first lactation. Sci. Rep. p 7:40864. https://doi.org/10.1038/srep40864.

23. Good IJ. 1953. The population frequencies of species and the estimation of population parameters. Biometrika 40:237–264. https://doi.org/10.2307/2333344.

24. Chao A. 1984. Nonparametric estimation of the number of classes in a population. SCAND J STAT 1:265–270.

25. Shannon CE. 1948. A Mathematical Theory of Communication. Bell System Technical Journal 27:379–423. https://doi.org/10.1002/j.1538-7305.1948.tb01338.x

26. Simpson EH. 1949. Measurement of Diversity. Nature 163:688–688.

27. Bray JR, Curtis JT. 1957. An Ordination of the Upland Forest Communities of Southern Wisconsin. Ecol Monogr 27:325–349. https://doi.org/10.2307/1942268.

28. Oksanen J, Blanchet FG, Friendly M, Kindt R, Legendre P, McGlinn D, Minchin PR, O'Hara RB, Simpson GL, Solymos P, Stevens MHH, Szoecs E, Wagner H. 2016. Vegan: an introduction to ordination. vegan: Community Ecology Package R package version 24-1. https://CRAN.R-project.org/package=vegan.

29. Breslow NE, Clayton DG. 1993. Approximate Inference in Generalized Linear Mixed Models. Journal of the American Statistical Association 88:9. https://doi.org/10.2307/2290687.

30. Delignette-Muller ML, Dutang C. 2015. fitdistrplus: An R Package for Fitting Distributions. J Stat Softw 64(4), 1-34. https://doi.org/10.18637/jss.v064.i04.

31. Venables WN, Ripley BD. 2002. Modern applied statistics with S, 4th ed. Springer, New York, NY.

32. Lenth R v. 2016. Least-Squares Means: The R Package lsmeans. Journal of Statistical Software, 69(1), 1-33. https://doi.org/10.18637/jss.v069.i01.

33. Wickham H, Chang W. 2016. ggplot2: Elegant Graphics for Data Analysis. Springer-Verlag New York, 2016. https://cran.r-project.org/web/packages/ggplot2/index.html.

674

675    34. Anderson MJ. 2001. A new method for non-parametric multivariate analysis of variance.
676        Austral Ecol 26:32–46. https://doi.org/10.1111/j.1442-9993.2001.01070.pp.x

677

678    35. Guillem Salazar. 2018. EcolUtils: Utilities for community ecology analysis. R package
679        version 0.1. https://github.com/GuillemSalazar/EcolUtils.

680

681    36. Liaw A, Wiener M. 2002. Classification and regression by randomForest. R news.2(3):18-
682        22. https://cran.r-project.org/web/packages/randomForest.

683

684    37. Kuhn M. 2020. caret: Classification and Regression Training. R package version 6.0-86.
685        https://CRAN.R-project.org/package=caret.

686

687    38. Breiman L. 2001. Random forests. Mach Learn 45:5–32.
688        https://doi.org/10.1023/A:1010933404324.

689

690    39. Song L, Maintainer PL, Langfelder P. 2015. randomGLM: Random General Linear Model
691        Prediction. R package version 1.02-1. https://CRAN.R-project.org/package=randomGLM.

692

693    40. Chen X, Ishwaran H. 2012. Random forests for genomic data analysis. Genomics 99:323–
694        329. https://doi.org/10.1016/j.ygeno.2012.04.003.

695

696    41. Chawla N v, Bowyer KW, Hall LO, Kegelmeyer WP. 2002. SMOTE: Synthetic Minority
697        Over-sampling Technique. J ARTIF INTELL RES. https://doi.org/10.1613/jair.953.

698

699    42. Blagus R, Lusa L. 2013. SMOTE for high-dimensional class-imbalanced data. BMC
700        Bioinformatics 14:106. https://doi.org/10.1186/1471-2105-14-106.

701

702    43. Landis JR, Koch GG. 1977. The Measurement of Observer Agreement for Categorical Data.
703        Biometrics 33:159-174. https://doi.org/10.2307/2529310.

704

705    44. Landis JR, Koch GG. 1977. An Application of Hierarchical Kappa-type Statistics in the
706        Assessment of Majority Agreement among Multiple Observers. Biometrics 33:363-374.
707        https://doi.org/10.2307/2529786.

708

709    45. Li M, Penner GB, Hernandez-Sanabria E, Oba M, Guan LL. 2009. Effects of sampling
710        location and time, and host animal on assessment of bacterial diversity and fermentation
711        parameters in the bovine rumen. J Appl Microbiol 107:1924–1934.
712        https://doi.org/10.1111/j.1365-2672.2009.04376.x.

713

714    46. Lv X, Chai J, Diao Q, Huang W, Zhuang Y, Zhang N. 2019. The Signature Microbiota Drive
715        Rumen Function Shifts in Goat Kids Introduced to Solid Diet Regimes. Microorganisms
716        7:516. https://doi.org/10.3390/microorganisms7110516.

717

718    47. Clemmons BA, Martino C, Powers JB, Campagna SR, Voy BH, Donohoe DR, Gaffney J,
719        Embree MM, Myer PR. 2019. Rumen Bacteria and Serum Metabolites Predictive of Feed
720        Efficiency Phenotypes in Beef Cattle. Sci Rep 9:19265. https://doi.org/10.1038/s41598-019-
721        55978-y.

722

723

724

725

726     **Table 1.** Samples and experimental design.

| Sample set | Description | Sample count | Used in classification? | Used to train regression model? |
|---|---|---|---|---|
| Summer, Time course, Farm 1 (STC) | Six timepoints of sampling paired buccal and rumen contents. | 8 animals | Yes | Yes |
| Spring sampling, Farm 1 (SPS) | Paired rumen and buccal contents; taken 4 hours after feeding | 5 animals | No | Yes |
| Summer sampling, Farm 2 (SUS) | Paired rumen and buccal contents; taken 2 hours prior to feeding | 8 animals | No | Yes |

727

**Table 2.** Variable importance analysis from the random forest classifier showing the most important bacterial OTUs (importance: scaled Mean Decrease in Gini≥50%) that discriminate between buccal swab and rumen samples.

| Taxa | Importance | Sample[1] | T1 Mean[2] | T1 Prev[3]. | T2 Mean | T2 Prev. | T3 Mean | T3 Prev. | T4 Mean | T4 Prev. | T5 Mean | T5 Prev. | T6 Mean | T6 Prev. | Group |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| OTU0003-Prevotella_1* | 100 | BS | 1.38 | 100.0 | 0.73 | 75.0 | 0.08 | 62.5 | 1.07 | 100.0 | 0.40 | 87.5 | 0.68 | 100.0 | CORE |
| | | RL | 3.17 | 100.0 | 3.03 | 100.0 | 4.09 | 100.0 | 3.49 | 100.0 | 3.63 | 100.0 | 2.95 | 100.0 | |
| | | RS | 2.12 | 100.0 | 2.06 | 100.0 | 2.13 | 100.0 | 2.41 | 100.0 | 2.36 | 100.0 | 2.61 | 100.0 | |
| Otu0405-p-1088-a5_gut_group | 96.8 | RS | 0.00 | 18.8 | 0.00 | 0.0 | 0.00 | 18.8 | 0.01 | 31.3 | 0.01 | 37.5 | 0.00 | 13.3 | RUMEN |
| | | RL | 0.07 | 93.3 | 0.09 | 100.0 | 0.11 | 100.0 | 0.08 | 100.0 | 0.09 | 100.0 | 0.04 | 93.8 | |
| | | BS | 0.01 | 33.3 | 0.00 | 12.5 | 0.00 | 12.5 | 0.01 | 37.5 | 0.00 | 0.0 | 0.00 | 25.0 | |
| Otu0001-Prevotella_1* | 87.4 | BS | 3.13 | 100.0 | 1.16 | 87.5 | 0.17 | 62.5 | 2.78 | 100.0 | 0.90 | 100.0 | 2.21 | 100.0 | CORE |
| | | RL | 8.08 | 100.0 | 9.27 | 100.0 | 12.13 | 100.0 | 9.15 | 100.0 | 9.83 | 100.0 | 7.97 | 100.0 | |
| | | RS | 5.36 | 100.0 | 5.35 | 100.0 | 5.84 | 100.0 | 6.79 | 100.0 | 6.42 | 100.0 | 6.54 | 100.0 | |
| Otu0241-Neisseriaceae | 86.5 | BS | 0.97 | 33.3 | 1.13 | 87.5 | 0.16 | 100.0 | 0.18 | 50.0 | 0.18 | 75.0 | 0.08 | 100.0 | ORAL |
| | | RL | 0.00 | 20.0 | 0.00 | 6.3 | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 12.5 | |
| | | RS | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | |
| Otu0113-Streptococcus | 86.1 | BS | 0.19 | 50.0 | 0.75 | 75.0 | 0.31 | 100.0 | 0.54 | 62.5 | 0.38 | 87.5 | 10.05 | 100.0 | ORAL |
| | | RL | 0.00 | 6.7 | 0.00 | 6.3 | 0.00 | 6.3 | 0.00 | 0.0 | 0.00 | 6.3 | 0.00 | 6.3 | |
| | | RS | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | |
| Otu0401-Streptococcus | 84.9 | BS | 0.63 | 50.0 | 0.19 | 87.5 | 0.17 | 100.0 | 0.10 | 37.5 | 0.26 | 75.0 | 0.13 | 100.0 | ORAL |
| | | RL | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | |
| | | RS | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | |
| Otu0434-Howardella | 83.3 | BS | 0.01 | 66.7 | 0.00 | 12.5 | 0.00 | 12.5 | 0.01 | 50.0 | 0.01 | 50.0 | 0.00 | 0.0 | RUMEN |
| | | RL | 0.07 | 93.3 | 0.09 | 100.0 | 0.12 | 100.0 | 0.06 | 93.8 | 0.05 | 93.8 | 0.04 | 93.8 | |
| | | RS | 0.00 | 0.0 | 0.00 | 18.8 | 0.00 | 12.5 | 0.01 | 31.3 | 0.00 | 25.0 | 0.00 | 13.3 | |
| Otu0424-Ruminococcaceae_ge | 81.3 | BS | 0.01 | 50.0 | 0.00 | 12.5 | 0.00 | 0.0 | 0.01 | 37.5 | 0.00 | 12.5 | 0.00 | 0.0 | RUMEN |
| | | RL | 0.06 | 93.3 | 0.06 | 87.5 | 0.14 | 93.8 | 0.06 | 81.3 | 0.08 | 100.0 | 0.07 | 93.8 | |
| | | RS | 0.00 | 0.0 | 0.01 | 31.3 | 0.00 | 25.0 | 0.00 | 18.8 | 0.01 | 31.3 | 0.01 | 33.3 | |
| Otu0838-Micrococcaceae | 79 | BS | 0.19 | 33.3 | 0.06 | 62.5 | 0.12 | 100.0 | 0.04 | 37.5 | 0.13 | 75.0 | 0.03 | 100.0 | ORAL |
| | | RL | 0.00 | 6.7 | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | |

| | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | RS | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | |
| Otu0184-Pasteurellaceae | 76.3 | BS | 0.97 | 50.0 | 2.45 | 75.0 | 0.09 | 100.0 | 0.08 | 37.5 | 0.16 | 75.0 | 0.12 | 100.0 | |
| | | RL | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | ORAL |
| | | RS | 0.00 | 0.0 | 0.00 | 6.3 | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 6.7 | |
| Otu0720-Jeotgalicoccus | 75 | BS | 0.24 | 50.0 | 0.03 | 87.5 | 0.15 | 100.0 | 0.07 | 50.0 | 0.13 | 75.0 | 0.03 | 100.0 | |
| | | RL | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 6.3 | ORAL |
| | | RS | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | |
| Otu0042-Ruminococcaceae_NK4A214_group* | 70.7 | BS | 0.16 | 100.0 | 0.06 | 50.0 | 0.01 | 25.0 | 0.10 | 87.5 | 0.05 | 50.0 | 0.10 | 100.0 | |
| | | RL | 0.59 | 100.0 | 0.80 | 100.0 | 0.99 | 100.0 | 0.72 | 100.0 | 0.81 | 100.0 | 0.55 | 100.0 | RUMEN |
| | | RS | 0.09 | 100.0 | 0.15 | 100.0 | 0.18 | 100.0 | 0.15 | 100.0 | 0.18 | 100.0 | 0.15 | 100.0 | |
| Otu0322-Streptococcus | 66.4 | BS | 0.19 | 50.0 | 0.05 | 62.5 | 0.60 | 75.0 | 0.03 | 62.5 | 0.83 | 75.0 | 1.67 | 75.0 | |
| | | RL | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 6.3 | 0.00 | 0.0 | ORAL |
| | | RS | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | |
| Otu0115-Bacteroidales_RF16_group_ge* | 63.8 | BS | 0.12 | 100.0 | 0.12 | 50.0 | 0.00 | 12.5 | 0.11 | 75.0 | 0.06 | 37.5 | 0.12 | 100.0 | |
| | | RL | 0.21 | 100.0 | 0.27 | 100.0 | 0.33 | 100.0 | 0.33 | 100.0 | 0.32 | 100.0 | 0.39 | 100.0 | RUMEN |
| | | RS | 0.02 | 81.3 | 0.02 | 87.5 | 0.02 | 75.0 | 0.03 | 81.3 | 0.01 | 68.8 | 0.02 | 66.7 | |
| Otu1233-Planococcaceae | 62.3 | BS | 0.03 | 66.7 | 0.03 | 50.0 | 0.03 | 75.0 | 0.06 | 25.0 | 0.07 | 75.0 | 0.05 | 75.0 | |
| | | RL | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | ORAL |
| | | RS | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | |
| Otu0780-Synergistes | 59.8 | BS | 0.00 | 16.7 | 0.00 | 12.5 | 0.00 | 0.0 | 0.01 | 37.5 | 0.00 | 12.5 | 0.00 | 0.0 | |
| | | RL | 0.02 | 80.0 | 0.02 | 87.5 | 0.06 | 93.8 | 0.03 | 75.0 | 0.03 | 87.5 | 0.03 | 75.0 | RUMEN |
| | | RS | 0.00 | 6.3 | 0.00 | 6.3 | 0.00 | 6.3 | 0.00 | 0.0 | 0.00 | 12.5 | 0.00 | 20.0 | |
| Otu0239-Streptococcus | 58.3 | BS | 0.06 | 66.7 | 0.17 | 62.5 | 1.07 | 75.0 | 0.40 | 75.0 | 0.49 | 75.0 | 2.48 | 75.0 | |
| | | RL | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 12.5 | 0.00 | 0.0 | 0.00 | 0.0 | ORAL |
| | | RS | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | |
| Otu0443-Prevotellaceae_UCG-001 | 55.4 | BS | 0.02 | 100.0 | 0.01 | 25.0 | 0.00 | 12.5 | 0.02 | 75.0 | 0.01 | 37.5 | 0.01 | 50.0 | |
| | | RL | 0.07 | 100.0 | 0.06 | 100.0 | 0.06 | 100.0 | 0.07 | 100.0 | 0.05 | 87.5 | 0.07 | 100.0 | RUMEN |
| | | RS | 0.01 | 62.5 | 0.01 | 56.3 | 0.01 | 37.5 | 0.01 | 62.5 | 0.00 | 25.0 | 0.01 | 60.0 | |
| Otu0056-Bibersteinia | 53.3 | BS | 1.17 | 83.3 | 2.55 | 87.5 | 1.52 | 100.0 | 1.55 | 87.5 | 0.93 | 87.5 | 8.06 | 100.0 | |
| | | RL | 0.00 | 6.7 | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 12.5 | 0.00 | 18.8 | ORAL |
| | | RS | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 6.3 | 0.00 | 6.7 | |

| Taxa | Importance | Type | Imp | Prev | Imp | Prev | Imp | Prev | Imp | Prev | Imp | Prev | Imp | Prev | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Otu0788-Rikenellaceae_RC9_gut_group | 52.7 | BS | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 0.0 | 0.00 | 25.0 | 0.00 | 12.5 | 0.00 | 25.0 | |
| | | RL | 0.03 | 86.7 | 0.03 | 87.5 | 0.04 | 87.5 | 0.03 | 87.5 | 0.03 | 87.5 | 0.03 | 62.5 | RUMEN |
| | | RS | 0.00 | 0.0 | 0.00 | 6.3 | 0.00 | 12.5 | 0.00 | 12.5 | 0.00 | 0.0 | 0.00 | 0.0 | |
| Otu0356-Rikenellaceae_RC9_gut_group | 52.5 | BS | 0.01 | 33.3 | 0.00 | 12.5 | 0.00 | 12.5 | 0.01 | 37.5 | 0.01 | 25.0 | 0.01 | 50.0 | |
| | | RL | 0.09 | 100.0 | 0.09 | 100.0 | 0.13 | 100.0 | 0.11 | 100.0 | 0.09 | 100.0 | 0.07 | 87.5 | RUMEN |
| | | RS | 0.01 | 37.5 | 0.01 | 43.8 | 0.01 | 31.3 | 0.01 | 43.8 | 0.01 | 25.0 | 0.01 | 33.3 | |
| Otu0120-Succiniclasticum* | 52.4 | BS | 0.07 | 100.0 | 0.02 | 50.0 | 0.00 | 0.0 | 0.05 | 75.0 | 0.02 | 37.5 | 0.03 | 50.0 | |
| | | RL | 0.18 | 100.0 | 0.20 | 100.0 | 0.22 | 100.0 | 0.19 | 100.0 | 0.17 | 100.0 | 0.26 | 100.0 | RUMEN |
| | | RS | 0.16 | 100.0 | 0.15 | 100.0 | 0.15 | 100.0 | 0.12 | 100.0 | 0.15 | 100.0 | 0.17 | 100.0 | |
| Otu0096-Ruminococcus_1* | 50.2 | BS | 0.11 | 100.0 | 0.06 | 50.0 | 0.01 | 25.0 | 0.17 | 87.5 | 0.07 | 37.5 | 0.13 | 75.0 | |
| | | RL | 0.05 | 93.3 | 0.04 | 75.0 | 0.01 | 50.0 | 0.04 | 100.0 | 0.04 | 81.3 | 0.09 | 93.8 | RUMEN |
| | | RS | 0.40 | 100.0 | 0.44 | 100.0 | 0.36 | 100.0 | 0.24 | 100.0 | 0.32 | 100.0 | 0.38 | 100.0 | |
| Otu0094-CPla-4_termite_group | 49.7 | BS | 0.02 | 66.7 | 0.01 | 25.0 | 0.00 | 12.5 | 0.03 | 87.5 | 0.01 | 25.0 | 0.00 | 0.0 | |
| | | RL | 0.25 | 100.0 | 0.31 | 100.0 | 0.56 | 100.0 | 0.37 | 100.0 | 0.45 | 100.0 | 0.26 | 100.0 | RUMEN |
| | | RS | 0.01 | 43.8 | 0.02 | 62.5 | 0.02 | 68.8 | 0.02 | 62.5 | 0.02 | 68.8 | 0.02 | 46.7 | |

*Taxa that varied with interaction of sampling time and sample type (Table S7); Importance and Prevalence are both expressed as percentages; [1]BS= buccal swab, rumen samples were merged based on rumen content strata: RL=rumen liquids (RAL +RVL) and RS=rumen solids (RAS+RVS);

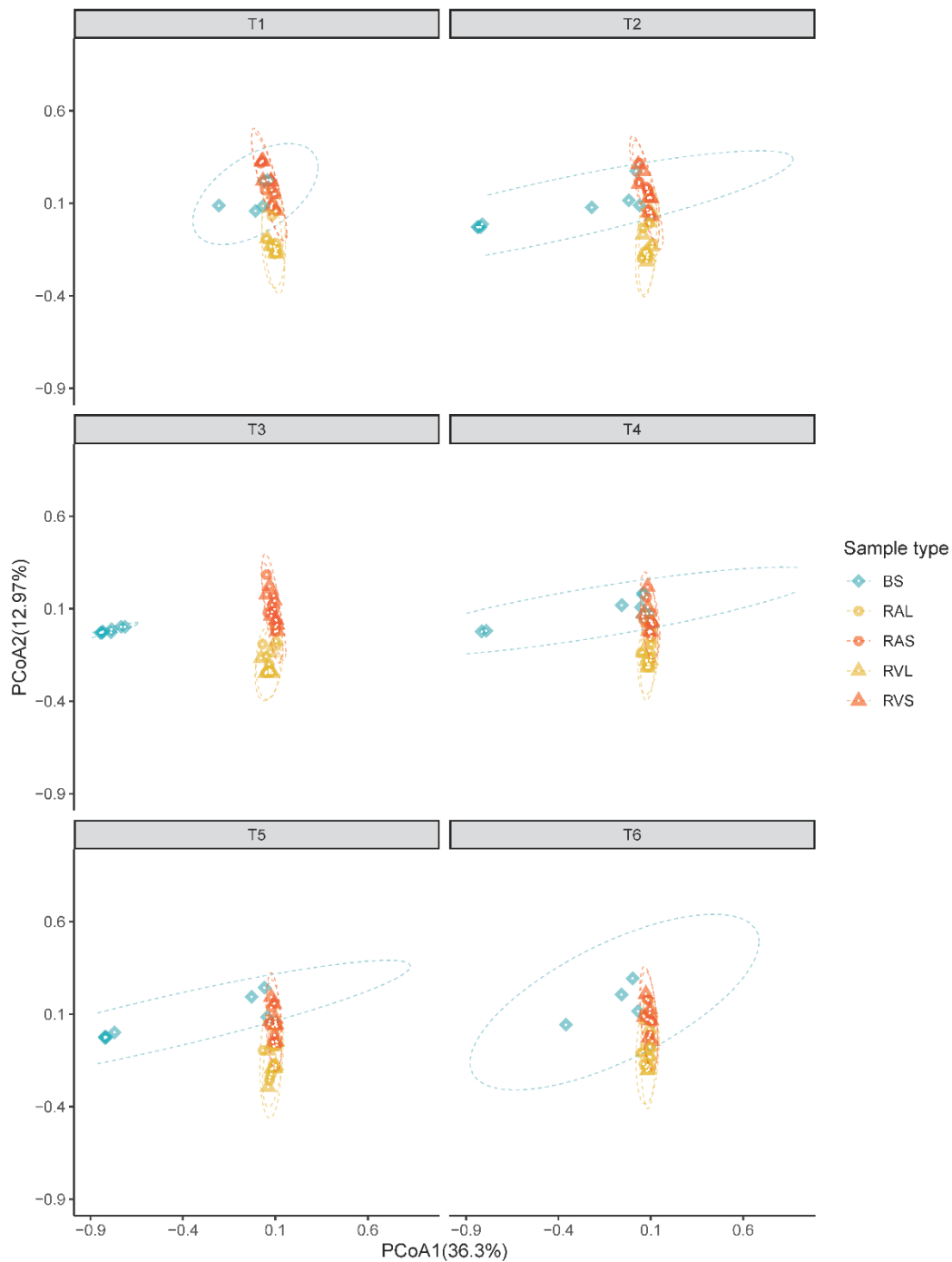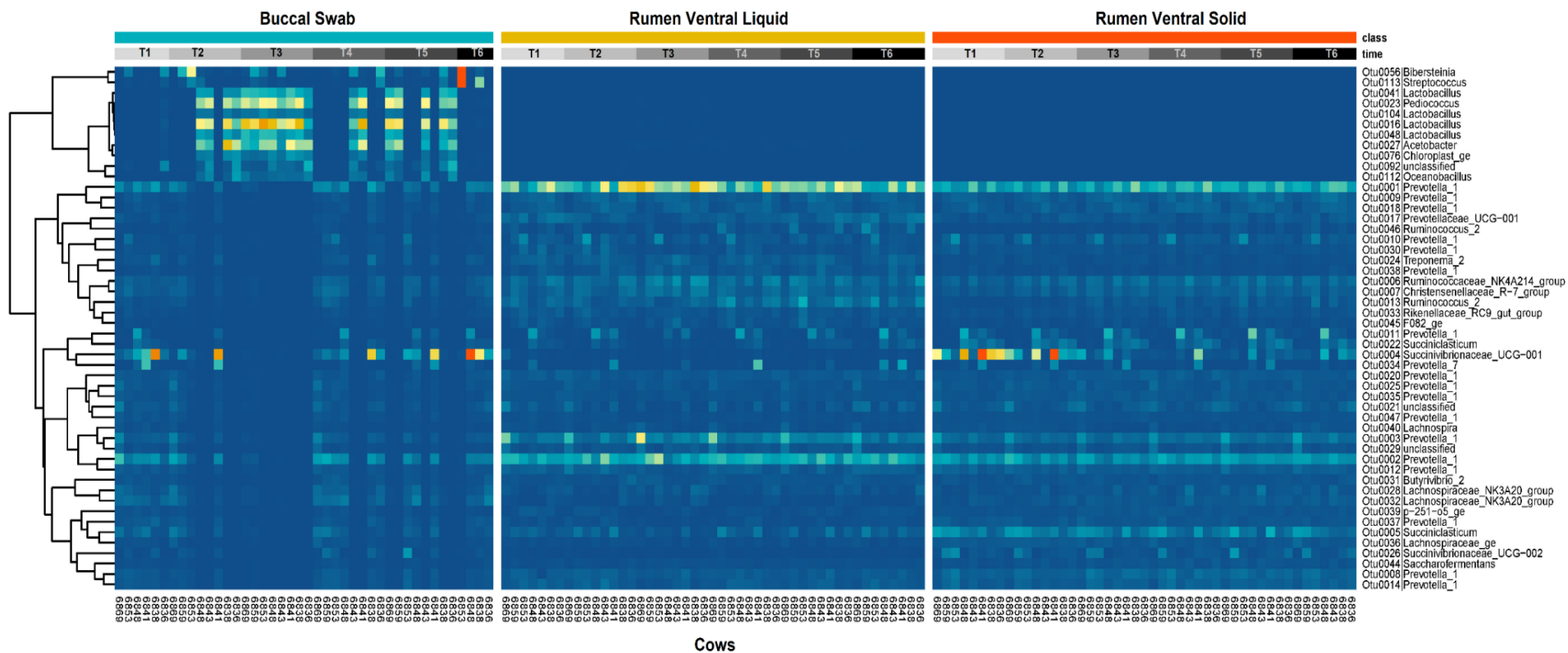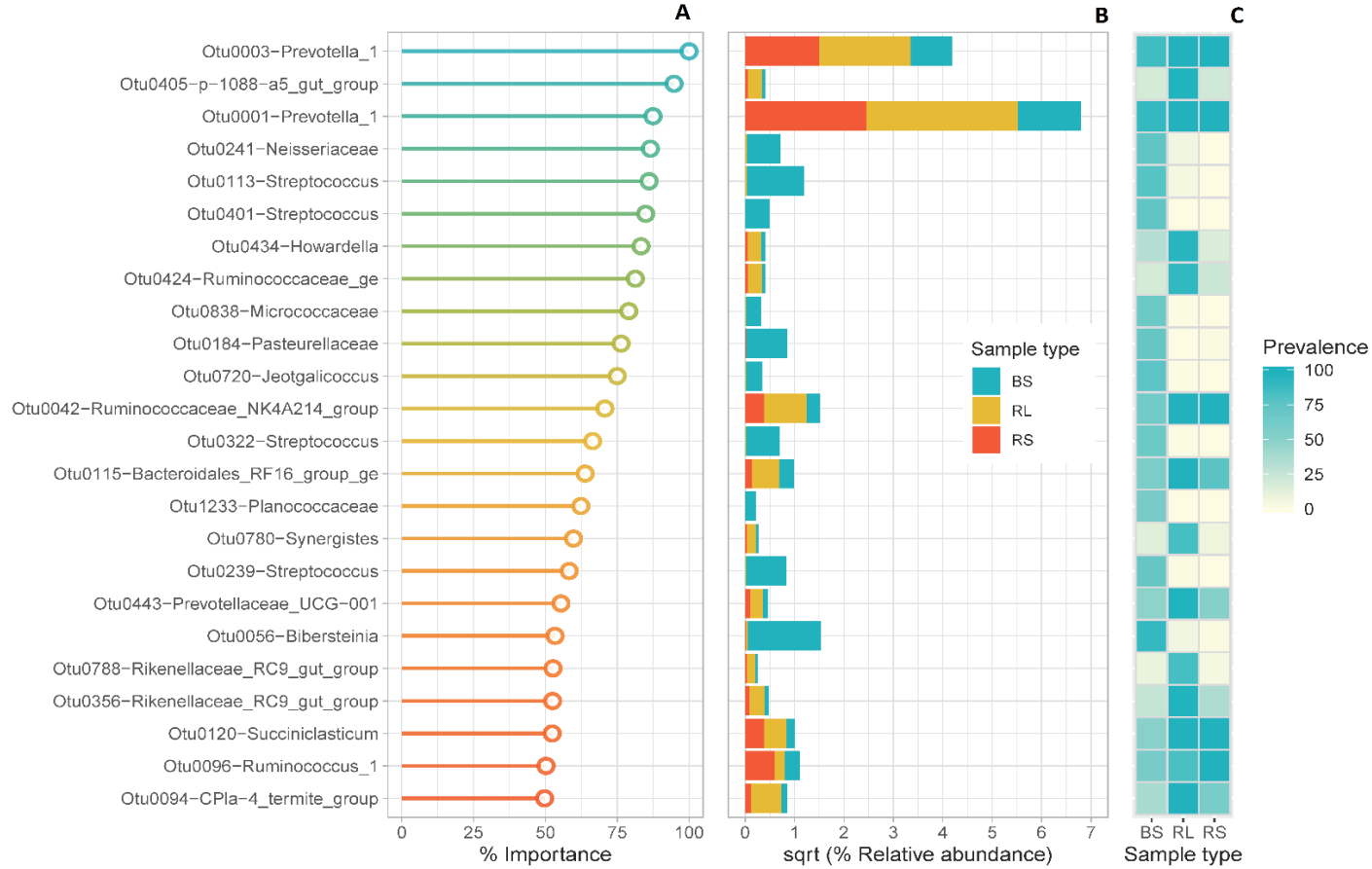[2]average relative abundance; [3]average prevalence.

**FIG 1.** Principal coordinate analysis (PCoA) showing Bray-Curtis dissimilarities in the composition of bacterial communities between sample types within each sampling time. Individual points in each plot represent a dairy cow, different colors and shapes represent a sample type (BS: buccal swab, RAL: rumen anterior liquid, RAS: rumen anterior solid, RVL: rumen ventral liquid and RVS: rumen ventral solid), and each facet represents a time point (T1 to T6). Percentages showed along the axes represent, respectively, the proportion of dissimilarities captured by PCoA in 2D coordinate space.

**FIG 2.** Distribution of the most abundant bacterial taxa among individual dairy cows according to sample type (BS: buccal swab, RVL: rumen ventral liquid and RVS: rumen ventral solid) and sampling time (T1:T6). The color-key represents the relative abundance at gradient of color from dark blue (low abundance) to dark orange (high abundance). The hierarchical dendrogram was established using Pearson product-moment correlations as the distance measure and "complete" as a clustering method.
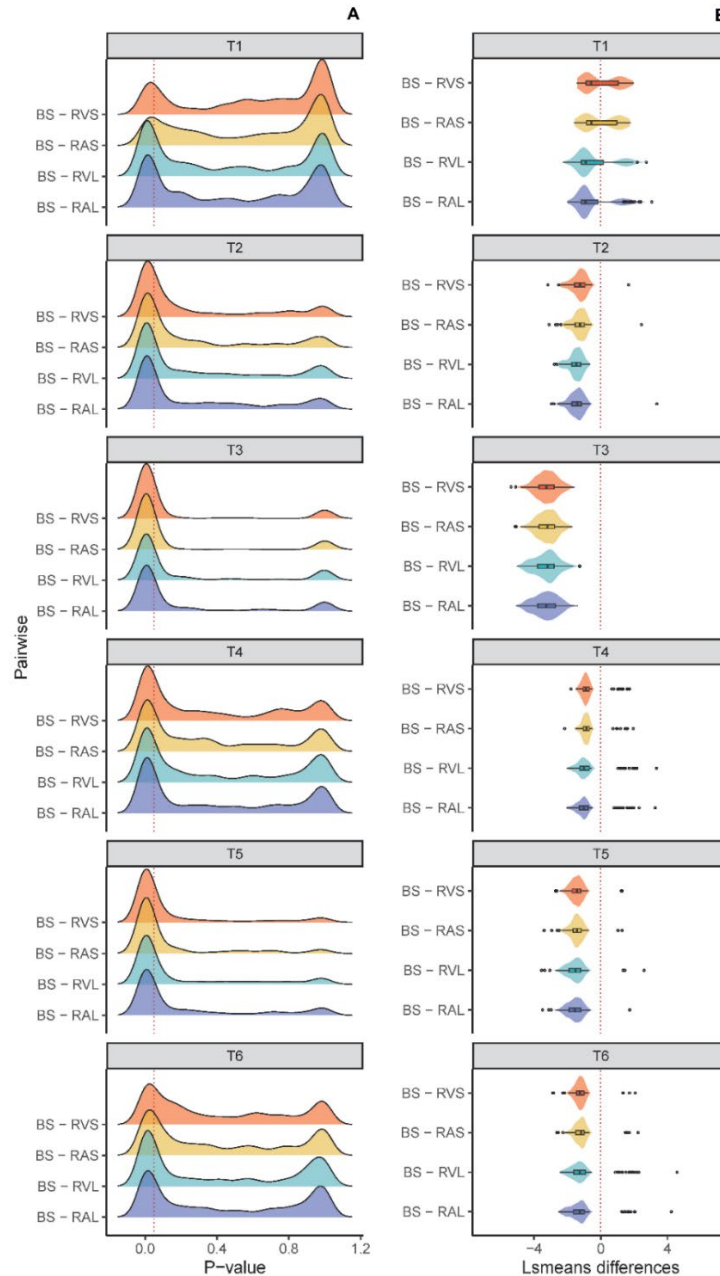
**FIG 3.** Variable importance (VIMP) plot from the random forest classifier. A) Lollipop chart showing the most important bacterial signatures that displayed importance (% Mean Decrease in Gini≥50) and that discriminate between buccal swab (BS), rumen liquids (RL) and rumen solids (RS) samples. B) Bar-plots of sqrt-relative abundance of OTUs according to sample type; C) Heat map of prevalence of OTUs in each sample type.
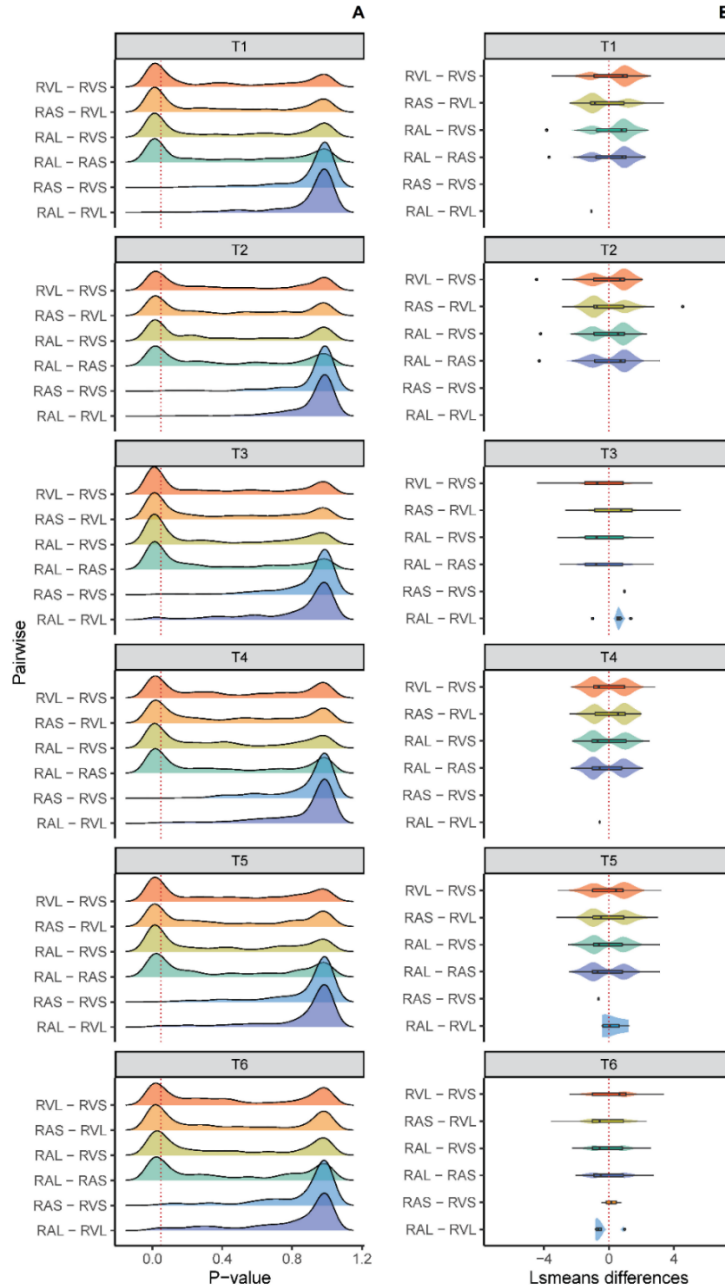
**FIG 4.** Bubble chart showing the prevalence and relative abundance of the oral OTUs assigned to higher taxa (phylum, family or genus level) according to sampling time (T1:T6).

**FIG 5. A)** Ridgeline plots showing the distribution of bacterial OTUs whose abundance varied significantly (red line = P-value ≤ 0.05) in pairwise comparisons between buccal swab (BS) and all types of rumen samples (RAL, RVL, RAS, and RVS) within each sampling time (T1:T6). **B)** Violin plot showing the Least Squares Means (LSmeans) differences of significant pairwise comparisons (Tukey HSD ≤ 0.05) between buccal swab and all types of rumen samples within each sampling time.

**FIG 6. A)** Ridgeline plot showing the distribution of bacterial OTUs whose abundance varied

significantly (red line=P-value≤0.05) in pairwise comparisons between all types of rumen samples

(RAL, RVL, RAS, and RVS) within each sampling time (T1:T6). **B)** Violin plot showing the Least

Square Means (LSMEANS) differences of significant pairwise comparisons (Tukey HSD ≤0.05)

between all types of rumen samples within each sampling time.