

Leveraging Existing 16S rRNA Gene Surveys to Identify Reproducible Biomarkers in Individuals with Colorectal Tumors

Marc A Sze¹ and Patrick D Schloss^{1†}

† To whom correspondence should be addressed: pschloss@umich.edu

¹ Department of Microbiology and Immunology, University of Michigan, Ann Arbor, MI

Co-author e-mails:

- marcsze@med.umich.edu

1 **Abstract**

2 An increasing body of literature suggests that both individual and collections of bacteria
3 are associated with the progression of colorectal cancer. As the number of studies
4 investigating these associations increases and the number of subjects in each study
5 increases, a meta-analysis to identify the associations that are the most predictive of
6 disease progression is warranted. We analyzed previously published 16S rRNA gene
7 sequencing data collected from feces and colon tissue. We quantified the odds ratios
8 (ORs) for individual bacterial taxa that were associated with an individual having tumors
9 relative to a normal colon. Among the fecal samples, there were no taxa that had significant
10 ORs associated with adenoma and there were 8 taxa with significant ORs associated with
11 carcinoma. Similarly, among the tissue samples, there were no taxa that had a significant
12 OR associated with adenoma and there were 3 taxa with significant ORs associated with
13 carcinoma. Among the significant ORs, the association between individual taxa and tumor
14 diagnosis was equal or below 7.11. Because individual taxa had limited association with
15 tumor diagnosis, we trained Random Forest classification models using only the taxa that
16 had significant ORs, using the entire collection of taxa found in each study, and using
17 operational taxonomic units defined based on a 97% similarity threshold. All training
18 approaches yielded similar classification success as measured using the Area Under the
19 Curve. The ability to correctly classify individuals with adenomas was poor and the ability
20 to classify individuals with carcinomas was considerably better using sequences from fecal
21 or tissue.

22 **Importance**

23 Colorectal cancer is a significant and growing health problem in which animal models and
24 epidemiological data suggest that the colonic microbiota have a role in tumorigenesis.
25 These observations indicate that the colonic microbiota is a reservoir of biomarkers that
26 may improve our ability to detect colonic tumors using non-invasive approaches. This
27 meta-analysis identifies and validates a set of 8 bacterial taxa that can be used within a
28 Random Forest modeling framework to differentiate individuals as having normal colons or
29 carcinomas. When models trained using one dataset were tested on other datasets, the
30 models performed well. These results lend support to the use of fecal biomarkers for the
31 detection of tumors. Furthermore, these biomarkers are plausible candidates for further
32 mechanistic studies into the role of the gut microbiota in tumorigenesis.

33 **Keywords**

34 microbiota; colorectal cancer; polyps; adenoma; tumor; meta-analysis.

35 **Background**

36 Colorectal cancer (CRC) is a growing world-wide health problem in which the microbiota
37 has been hypothesized to have a role in disease progression (1, 2). Numerous studies
38 using murine models of CRC have shown the importance of both individual microbes
39 (3–7) and the overall community (8–10) in tumorigenesis. Numerous case-control
40 studies have characterized the microbiota of individuals with colonic adenomas and
41 carcinomas in an attempt to identify biomarkers of disease progression (6, 11–17).
42 Because current CRC screening recommendations are poorly adhered to due to a
43 person's socioeconomic status, test invasiveness, and frequency of tests, development
44 and validation of microbiota-associated biomarkers for CRC progression could further
45 attempts to develop non-invasive diagnostics (18).

46 Recently, there has been an intense focus on identifying microbiota-based biomarkers
47 yielding a seemingly endless number of candidate taxa. Some studies point towards
48 mouth-associated genera such as *Fusobacterium*, *Peptostreptococcus*, *Parvimonas*, and
49 *Porphyromonas* that are enriched in people with carcinomas (6, 11–17). Other studies have
50 identified members of *Akkermansia*, *Bacteroides*, *Enterococcus*, *Escherichia*, *Klebsiella*,
51 *Mogibacterium*, *Streptococcus*, and *Providencia* (13–15). Additionally, *Roseburia* has been
52 found in some studies to be more abundant in people with tumors but in other studies it has
53 been found to be less abundant than what is found in subjects with normal colons (14, 17,
54 19, 20). There is support from mechanistic studies using tissue culture and murine models
55 that *Fusobacterium nucleatum*, pks-positive strains of *Escherichia coli*, *Streptococcus*
56 *gallolyticus*, and an enterotoxin-producing strain of *Bacteroides fragilis* are important in
57 tumorigenesis (5, 14, 21–24). These results point to a causative role for the microbiota in
58 tumorigenesis as well as their potential as diagnostic biomarkers.

59 Most studies have focused on identifying biomarkers in patients with carcinomas but

60 there is a clinical need to identify biomarkers associated with adenomas to facilitate
61 early detection of the tumors. Studies focusing on broad scale community metrics have
62 found that measures such as the total number of taxa (i.e. richness) are lower in those
63 with adenomas versus controls (25). Other studies have identified *Acidovorax*, *Bilophila*,
64 *Cloacibacterium*, *Desulfovibrio*, *Helicobacter*, *Lactobacillus*, *Lactococcus*, *Mogibacterium*,
65 and *Pseudomonas* to be enriched in those with adenomas (25–27). The ability to classify
66 individuals as having normal colons or adenomas based solely on the taxa within fecal
67 samples has been limited. However, when 16S rRNA gene sequence data was combined
68 with the results of a fecal immunochemical test (FIT), the ability to diagnose individuals
69 with adenomas was improved relative to using the FIT results alone (12).

70 A recent meta-analysis found that 16S rRNA gene sequences from members of
71 *Akkermansia*, *Fusobacterium*, and *Parvimonas* were fecal biomarkers for the presence of
72 carcinomas (28). Contrary to previous studies, they found sequences similar to members
73 of *Lactobacillus* and *Ruminococcus* to be enriched in patients with adenoma or carcinoma
74 relative to those with normal colons (12, 15, 16). In addition, they found that 16S rRNA
75 gene sequences from members of *Haemophilus*, *Methanosphaera*, *Prevotella*, and
76 *Succinovibrio* were enriched in patients with adenomas and *Pantoea* were enriched
77 in patients with carcinomas. Although this meta-analysis was helpful for distilling a
78 large number of possible biomarkers, the aggregate number of samples included in the
79 analysis (n=509) was smaller than several larger case-control studies that have since been
80 published (12, 27)

81 Here we provide an updated meta-analysis using 16S rRNA gene sequence data from
82 both feces (n=1737) and colon tissue (492 samples from 350 individuals) from 14 studies
83 (11–17, 19, 20, 23, 25–27, 29) [Table 1 & 2]. We expand both the breadth and scope
84 of the previous meta-analysis to investigate whether biomarkers describing the bacterial
85 community or specific members of the community can more accurately classify patients as

86 having adenoma or carcinoma. Our results suggest that the bacterial community changes
87 as disease severity worsens and that a subset of the microbial community can be used to
88 diagnose the presence of carcinoma.

89 Results

90 ***Lower bacterial diversity is associated with higher odds ratio (OR) of tumors.*** We
91 first assessed whether variation in broad community metrics like total number of operational
92 taxonomic units (OTUs) (i.e. richness), the evenness of their abundance, and the overall
93 diversity of the communities were associated with disease stage after controlling for
94 study and variable region differences. In fecal samples, both evenness and diversity
95 were significantly lower in successive disease severity categories (P-value=0.025 and
96 P-value=0.043, respectively) [Figure 1]; there was no significant difference for richness
97 (P-value=0.21). We next tested whether the lower value of these community metrics
98 translated into significant ORs for having an adenoma or carcinoma. For fecal samples,
99 the ORs for richness were not significantly greater than 1.0 for adenoma or carcinoma
100 (P-value=0.40) [Figure 2A]. The ORs for evenness were significantly higher than 1.0 for
101 adenoma (OR=1.3 (95% Confidence Interval: 1.02 - 1.65), P-value=0.035) and carcinoma
102 (OR=1.66 (1.2 - 2.3), P-value=0.0021) [Figure 2B]. The ORs for diversity were only
103 significantly greater than 1.0 for carcinoma (OR=1.61 (1.14 - 2.28), P-value=0.0069),
104 but not for adenoma (P-value=0.11) [Figure 2C]. Although these ORs are significantly
105 greater than 1.0, it is doubtful that they are clinically meaningful.

106 Similar to our analysis of sequences obtained from fecal samples, we repeated the analysis
107 using sequences obtained from colon tissue. There were no significant differences in
108 richness, evenness, or diversity as disease severity progressed from control to adenoma
109 to carcinoma (P-values > 0.05). We next analyzed the ORs, for matched (i.e. where
110 unaffected tissue and tumors were obtained from the same individual) and unmatched
111 (i.e. where unaffected tissue and tumor tissue were not obtained from the same individual)
112 tissue samples. The ORs for adenoma and carcinoma were not significantly different from
113 1.0 for any measure (P-values > 0.05) [Figure S1 & Table S1]. This is likely due to the
114 combination of a small effect size and the relatively small number of studies and size of

115 studies used in the analysis.

116 ***Disease progression is associated with changes in community structure.*** Based
117 on the differences in evenness and diversity, we next asked whether there were
118 community-wide differences in the structure of the communities associated with different
119 disease stages. We identified significant bacterial community differences in the feces of
120 patients with adenomas relative to those with normal colons in 1 of 4 studies and in patients
121 with carcinomas relative to those with normal colons in 6 of 7 studies (PERMANOVA;
122 P-value < 0.05) [Table S2]. Similar to the analyses using fecal samples, there were
123 significant differences in the bacterial community structure of subjects with normal colons
124 and those with adenomas (1 of 2 studies) and carcinomas (1 of 3 studies) [Table S2].
125 For studies that used matched samples, we did not observe any differences in bacterial
126 community structures [Table S2]. Combined, these results indicate that there were
127 consistent and significant community-wide changes in the fecal community structure of
128 subjects with carcinomas. However, the signal observed in subjects with adenomas or
129 when using tissue samples was not as consistent. This is likely due to a smaller effect
130 size or the relatively small sample sizes among the studies that characterized the tissue
131 microbiota.

132 ***Individual taxa are associated with significant ORs for carcinomas.*** We next
133 identified those taxa that had ORs that were significantly associated with having a
134 normal colon or the presence of adenomas or carcinomas. No taxa had a significant
135 OR for the presence of adenomas when we used data collected from fecal or tissue
136 samples (Table S3 & S4). In contrast, 8 taxa had significant ORs for the presence of
137 carcinomas using data from fecal samples. Of these, 4 are commonly associated with
138 the oral cavity: *Fusobacterium* (OR=2.74 (1.95 - 3.85)), *Parvimonas* (OR=3.07 (2.11
139 - 4.46)), *Porphyromonas* (OR=3.2 (2.26 - 4.54)), and *Peptostreptococcus* (OR=7.11
140 (3.84 - 13.17)) [Table S3]. The other 4 were *Clostridium XI* (OR=0.65 (0.49 - 0.86)),

141 *Enterobacteriaceae* (OR=1.79 (1.33 - 2.41)), *Escherichia* (OR=2.15 (1.57 - 2.95)), and
142 *Ruminococcus* (OR=0.63 (0.48 - 0.83)). Among the data collected from tissue samples,
143 only unmatched carcinoma samples had taxa with a significant OR. Those included *Dorea*
144 (OR=0.35 (0.22 - 0.55)), *Blautia* (OR=0.47 (0.3 - 0.73)), and *Weissella* (OR=5.15 (2.02 -
145 13.14)). Mouth-associated genera were not significantly associated with a higher OR for
146 carcinoma in tissue samples [Table S4]. For example, *Fusobacterium* had an OR of 3.98
147 (1.19 - 13.24); however, due to the small number of studies and considerable variation in
148 the data, the Benjimini-Hochberg corrected P-value was 0.93 [Table S4]. It is interesting
149 to note that *Ruminococcus* and members of *Clostridium XI* in fecal samples and *Dorea*
150 and *Blautia* in tissue had ORs that were significantly less than 1.0, which suggests that
151 these populations are protective against the development of carcinomas. Overall, there
152 was no overlap in the taxa with significant OR between fecal and tissue samples.

153 ***Individual taxa with a significant OR do a poor job of differentiating subjects with***
154 ***normal colons and those with carcinoma.*** We next asked whether those taxa that had
155 a significant OR associated with having a normal colon or carcinomas could be used
156 individually, to classify subjects as having a normal colon or carcinomas. OR values were
157 calculated based on whether the relative abundance for a taxon in a subject was above
158 or below the median relative abundance for that taxon across all subjects in a study. To
159 measure the ability of these taxa to classify individuals we instead generated receiver
160 operator characteristic (ROC) curves for each taxon in each study and calculated the area
161 under the curve (AUC). This allowed us to use a more fluid relative abundance threshold
162 for classifying individuals by their disease status. Using data from fecal samples, the 8 taxa
163 did no better at classifying the subjects than one would expect by chance (i.e. AUC=0.50)
164 [Figure 3A]. The taxa that performed the best included *Clostridium XI*, *Ruminococcus*,
165 and *Escherichia*. However, these had median AUC values less than 0.588 indicating
166 their limited value as biomarkers when used individually. Likewise, in unmatched tissue
167 samples the 3 taxa with significant OR taxa had AUC values that were marginally better

168 than one would expect by chance [Figure 3B]. The relative abundance of *Dorea* was the
169 best predictor of carcinomas and its median AUC was only 0.62. These results suggest that
170 although these taxa are associated with a significant OR for the presences of carcinomas,
171 they do a poor job of classifying a subject's disease status when used individually.

172 ***Combined taxa model classifies subjects better than using individual taxa.*** Instead
173 of attempting to classify subjects based on individual taxa, next we combined information
174 from the individual taxa and evaluated the ability to classify a subject's disease status
175 using Random Forest models. For data from fecal samples, the combined model had an
176 AUC of 0.75, which was significantly higher than any of the AUC values for the individual
177 taxa (P-value < 0.033). When this approach was used to train models using data from
178 each study, the most important taxa were *Ruminococcus* and *Clostridium XI* [Figure 4A].
179 Similarly, using data from the unmatched tissue samples, the combined model had an AUC
180 of 0.77, which was significantly higher than the AUC values for classifying based on the
181 relative abundances of *Blautia* and *Weissella* individually (P-value < 0.037). Both *Dorea*
182 and *Blautia* were the most important taxa in the tissue-based models [Figure 4B]. Pooling
183 the information from the taxa with significant ORs resulted in models that outperformed
184 classifications made using the same taxa individually.

185 ***Performance of models based on taxa relative abundance in full community is***
186 ***better than that of models based on taxa with significant ORs.*** Next, we asked
187 whether a Random Forest classification model built using all of the taxa found in the
188 communities would outperform the models generated using those taxa with a significant
189 OR. Similar to our inability to identify taxa associated with a significant OR for the presence
190 of adenomas, the median AUCs to classify subjects as having normal colons or having
191 adenomas using data from fecal or tissue samples were only marginally better than 0.5
192 for any study (median AUC=0.549 (range: 0.367 - 0.971)) [Figure 5A & S2A]. In contrast,
193 the models for classifying subjects as having normal colons or having carcinomas using

194 data from fecal or tissue samples yielded AUC values meaningfully higher than 0.5 [Figure
195 5B & S2B-C]. When we compared the models based on all of the taxa in a community to
196 models based on the taxa with significant ORs, the results were mixed. Using the data
197 from fecal samples, we found that the AUC for 6 of 7 studies were an average of 14.8%
198 higher and AUC for the Flemer study was 0.54% lower when using the relative abundance
199 data from all taxa relative to using the relative abundance of only the taxa with significant
200 ORs. The overall improvement in performance was statistically significant (mean=12.61%,
201 one-tailed paired T-test; P-value=0.005). Among the models trained using data from fecal
202 samples, *Bacteroides* and *Lachnospiraceae* were the most common taxa in the top 10%
203 mean decrease in accuracy across studies [Figure S3]. Using data from unmatched
204 tissue samples to train classification models, we found that the AUC of studies was an
205 average 19.11% higher when we used all of the taxa rather than the 3 taxa with significant
206 ORs (one-tailed paired T-test; P-value=0.03). For the models trained using data from
207 unmatched tissue samples, *Lachnospiraceae*, *Bacteroidaceae*, and *Ruminococcaceae*
208 were the most common taxa in the top 10% mean decrease in accuracy across studies
209 [Figure S4]. Although the models trained using those taxa with a significant OR perform
210 well for classifying individuals with and without carcinomas, models trained using data from
211 the full community perform better.

212 ***Performance of models based on OTU relative abundances are not significantly***
213 ***better than those based on taxa with significant ORs.*** The previous models were
214 based on relative abundance data where sequences were classified to coarse taxonomic
215 assignments (i.e. typically genus or family level). To determine whether model performance
216 improved with finer scale classification, we assigned sequences to operational taxonomic
217 units (OTUs) where the similarity among sequences within an OTU was more than 97%. We
218 again found that classification models built using all of the sequence data for a community
219 did a poor job of differentiating between subjects with normal colons and those with
220 adenomas (median AUC: 0.53 (0.37- 0.56)). However, they did a good job of differentiating

221 between subjects with normal colons and those with carcinomas (median AUC: 0.71 (0.50-
222 0.90)). The OTU-based models performed similarly to those constructed using the taxa
223 with significant ORs (one-tailed paired T-test; P-value=0.979) and those using all taxa
224 (one-tailed paired T-test; P-value=0.184) [Figure 4]. Among the OTUs that had the highest
225 mean decrease in accuracy for the OTU-based models, we found that OTUs that affiliated
226 with all of the 8 taxa that had a significant OR were within the top 10% for at least one study.
227 This result was surprising as it indicated that a finer scale classification of the sequences
228 and thus a larger number of features to select from, did not yield improved classification of
229 the subjects.

230 ***Generalizability of taxon-based models trained on one dataset to the other***
231 ***datasets.*** Considering the good performance of the Random Forest models trained using
232 the relative abundance of taxa with significant ORs and models trained using the relative
233 abundance of all taxa, we next asked how well the models would perform when given
234 data from a different cohort. For instance, if a model was trained using data from the
235 Ahn study, we wanted to know how well it would perform using the data from the Baxter
236 study. The models trained using the taxa with significant ORs all had a higher median AUC
237 than the models trained using all of the taxa when tested on the other datasets [Figure 6
238 & S5]. As might be expected, the difference between the performance of the modeling
239 approaches appeared to vary with the size of the training cohort ($R^2=0.66$) [Figure 6].
240 These data suggest that given a sufficient number of subjects with normal colons and
241 carcinomas, Random Forest models trained using a small number of taxa can accurately
242 classify individuals from a different cohort.

243 Discussion

244 We performed a meta-analysis to identify and validate microbiota-based biomarkers that
245 could be used to classify individuals as having normal colons or colonic tumors using fecal
246 or tissue samples. To our surprise, Random Forest classification models constructed to
247 differentiate individuals with normal colons from those with carcinomas using a subset of the
248 community performed well relative to models constructed using the full communities. When
249 we applied the models trained on each dataset to the other datasets in our study, we found
250 that the models trained using the subset of the communities performed better than those
251 using the full communities. These models were trained using data in which sequences were
252 assigned to bacterial taxa using a classifier that typically assigned sequences to the family
253 or genus level. When we attempted to improve the specificity of the classification by using
254 an OTU-based approach the resulting models performed as well as those constructed using
255 coarse taxonomic assignments. These results are significant because they strengthen the
256 growing literature indicating a role for the colonic microbiota in tumorigenesis, as a potential
257 tool as a non-invasive diagnostic, and for assessing risk of disease and recurrence (9, 12,
258 30).

259 Fine scale classification of sequences into OTUs did not improve our classification models.
260 This was also tested in earlier efforts to use shotgun metagenomic data to classify
261 individuals as having normal colons or tumors; however, it was shown that analyses
262 performed using shotgun metagenomic data did not perform better than using 16S
263 rRNA gene sequencing data (31). We hypothesize that fine scale classification may
264 not result in better classification because distribution of microbiota between individuals
265 is patchy. In contrast, models using coarser taxonomic assignments will pool the fine
266 scale diversity, resulting in less patchiness and better classification. Furthermore, the
267 ability of models trained using a subset of the community to outperform those using the
268 full community when testing the models on the other datasets may also be a product of

269 the patchiness of the human-associated microbiota. The models based on the 8 taxa that
270 had significant ORs used taxa that were found in every study and tended to have higher
271 relative abundances. Similar to the OTU-based models, those models based on the full
272 community taxonomy assignments were still sensitive to the patchy distribution of taxa.
273 Regardless, it is encouraging that a collection of 8 taxa could reliably classify individuals
274 as having carcinomas considering the differences in cohorts, DNA extraction procedures,
275 regions of the 16S rRNA gene, and sequencing methods.

276 When used to classify individuals with carcinomas, the taxa with significant ORs could
277 not reliably classify individuals on their own [Figure 3]. This result further supports the
278 hypothesis that carcinoma-associated microbiota have a patchy distribution. Two individuals
279 may have had the same classification, based on the relative abundance of different
280 populations within this group of 8 taxa. Although these results only reflect associations
281 with disease, it is tempting to hypothesize that the patchiness is indicative of distinct
282 mechanisms of exacerbating tumorigenesis or that multiple taxa have the same mechanism
283 of exacerbating tumorigenesis. For example, strains of *Escherichia coli* and *Fusobacterium*
284 *nucleatum* have been shown to worsen inflammation in mouse models of tumorigenesis
285 (5, 6, 21). In contrast to the patchiness of the taxa that were positively associated with
286 carcinomas, potentially beneficial taxa had a more consistent association [Figure 6]. This
287 result was particularly interesting because members of these taxa (i.e. *Ruminococcus*
288 and *Clostridium XI* in fecal samples and *Dorea* and *Blautia* in tissue) are thought to be
289 beneficial due to their involvement in production of anti-inflammatory short chain fatty acids
290 (32–34).

291 All of the adenoma classification models performed poorly, which is consistent with
292 previous studies (27, 30). However, the classification results are at odds with results
293 of the multitarget microbiota test (MMT) from Baxter, et al. (12) who observed an
294 AUC of 0.755 when the test was applied to individuals with adenomas. There are two

295 major differences between the models generated in this meta-analysis and that analysis.
296 The MMT attempted to classify individuals as having a normal colon or having colonic
297 lesions (i.e. adenomas or carcinomas) and not adenomas alone. Further, the MMT
298 incorporated fecal immunoglobulin test (FIT) data while our models only used 16S rRNA
299 gene sequencing data. Because FIT data were not available for the other studies in
300 our meta-analysis, it was not possible to validate the MMT approach. The ability to
301 differentiate between individuals with and without adenomas is an important problem since
302 early detection of tumors is critical to patient survivorship. However, it is possible that
303 we might have been able to detect differences in the bacterial community if individuals
304 with non-advanced and advanced adenomas were separated. This is a clinically relevant
305 distinction since advanced adenomas are at highest risk of progressing to carcinomas.
306 The initial changes of the microbiota during tumorigenesis could be focal to where the
307 initial adenoma develops and would not be easily assessed using fecal samples from an
308 individual with non-advanced adenomas. Unfortunately, distinguishing between individuals
309 with advanced and non-advanced adenomas was not possible in our meta-analysis since
310 the studies did not provide the clinical data needed to make that distinction.

311 Fecal samples represent a non-invasive approach to assess the structure of the gut
312 microbiota and are potentially useful for diagnosing individuals as having colonic tumors.
313 However, they do not reflect the structure of the mucosal microbiota (35). Regardless, the
314 taxa that were the most important in the feces-based models overlapped with those from
315 the models trained using the data from unmatched and matched colon tissue samples
316 [Figure S3]. Mucosal biopsies are preferred for focused mechanistic studies and have
317 offered researchers the opportunity to sample healthy and diseased tissue from the same
318 individuals (i.e. matched) using each individual as their own control or in a cross-sectional
319 design (i.e. unmatched). Because obtaining these samples is invasive, carries risks
320 to the individual, and is expensive, studies investigating the structure of the mucosal
321 microbiota generally have a limited number of participants. Thus, it was not surprising that

322 tissue-based studies did not provide clearer associations between the mucosal microbiota
323 and the presence of tumors. Interestingly, *Fusobacterium*, which has received increased
324 attention for its potential role in tumorigenesis (6) was not consistently identified across
325 the studies in our meta-analysis which is consistent with a recent replicability study (36).
326 This could be due to the relatively small number of individuals in the limited number of
327 studies. The classification models trained using the tissue-based data performed well when
328 tested with the training data (Figure S4), but performed poorly when tested on the other
329 tissue-associated datasets (Figure S5). Disturbingly, taxa that are commonly associated
330 with reagent contamination (e.g. *Novosphingobium*, *Acidobacteria Gp2*, *Sphingomonas*,
331 etc.) were detected within the tissue datasets. Such contamination is common in studies
332 where there is relatively low bacterial biomass (37). The lack of replication among the
333 tissue-based biomarkers may be a product of the relatively small number of studies and
334 individuals per study and possible reagent contamination.

335 Among the fecal sample data, we failed to identify several notable populations that are
336 commonly associated with carcinomas including an enterotoxigenic strain of *Bacteroides*
337 *fragilis* (ETBF) and *Streptococcus gallolyticus* subsp. *gallolyticus* (22, 24). ETBF have
338 been found in tumors in the proximal colon where they tend to form biofilms (20, 38).
339 Considering DNA from bacteria that are more prevalent in the proximal colon may be
340 degraded by the time it leaves the body, it is not surprising that we failed to identify a
341 significant OR for *Bacteroides* with carcinomas. In addition, since our approach could only
342 classify sequences to the genus level and there are likely multiple *Bacteroides* populations
343 in the colon, it is possible that sequences from ETBF and non-oncogenic *Bacteroides*
344 were pooled. This would then reduce the OR between *Bacteroides* and whether an
345 individual had carcinomas. It is also necessary to distinguish between populations that are
346 biomarkers for a disease and those that are known to cause disease. Although the latter
347 have been shown to have a causative role, they may appear at low relative abundance,
348 be found in specific locations, or may have a highly patchy distribution among affected

349 individuals.

350 Meta-analyses are a useful tool in microbiome research because they can demonstrate
351 whether a result can be replicated and facilitate new discoveries by pooling multiple
352 independent investigations. There have been several meta-analyses similar to this study
353 that have sought biomarkers for obesity (39–41), inflammatory bowel disease (40), and
354 colorectal cancer (28). Considering microbiome research is particularly prone to hype and
355 overgeneralization of results (42), these analyses are critical. Meta-analyses are difficult to
356 perform because the underlying 16S rRNA gene sequence data are not publicly available,
357 metadata are missing, incomplete, or vague, sequence data are of poor quality or derived
358 by non-standard approaches, and the original studies may be significantly underpowered.
359 Reluctance to publish negative results (i.e. the “file drawer effect”) is also likely to skew
360 our understanding of the relationship between microbiota and disease. Better attention to
361 these specific issues will increase the reproducibility and replicability of microbiota studies
362 and make it easier to perform these crucial meta-analyses. Moving forward, meta-analyses
363 will be important tools to help aggregate and find commonalities across studies when
364 investigating the microbiota in the context of a specific disease (28, 39–41).

365 Our meta-analysis suggests a strong association between the gut microbiota and colon
366 tumorigenesis. By aggregating the results from studies that sequenced the 16S rRNA
367 gene from fecal and tissue samples, we are able to provide evidence supporting the use of
368 microbial biomarkers to diagnose the presence of colonic tumors. Further development
369 of microbial biomarkers should focus on including other biomarkers (e.g. FIT), better
370 categorizing of people with adenomas, and expanding datasets to include larger numbers
371 of individuals. Based on prior research into the physiology of the biomarkers we identified,
372 it is likely that they have a causative role in tumorigenesis. Their patchy distribution across
373 individuals suggests that there are either multiple mechanisms causing disease or a single
374 mechanism (e.g. inflammation) that can be mediated by multiple, diverse bacteria.

375 **Methods**

376 **Datasets.** The studies used for this meta-analysis were identified through the review
377 articles written by Keku, et al. (43) and Vogtmann, et al. (44). Additional studies, not
378 mentioned in those reviews were obtained based on the authors' knowledge of the literature.
379 Studies were included that used tissue or feces as their sample source for 454 or Illumina
380 16S rRNA gene sequencing. A significant number of studies (N=12) were excluded from
381 the meta-analysis because they did not have publicly available sequences, did not use 454
382 or Illumina sequencing platforms, or did not have metadata that the authors were able to
383 share. We were able to obtain sequence data and metadata from the following studies:
384 Ahn, et al. (11), Baxter, et al. (12), Brim, et al. (29), Burns, et al. (15), Chen, et al. (13),
385 Dejea, et al. (20), Flemer, et al. (17), Geng, et al. (19), Hale, et al. (27), Kostic, et al. (45),
386 Lu, et al. (26), Sanapareddy, et al. (25), Wang, et al. (14), Weir, et al. (23), and Zeller,
387 et al. (16). The Zackular (46) study was excluded because the individuals studied were
388 included within the larger Baxter study (12). The Kostic study was excluded because after
389 we processed the sequences, all of the case samples had 100 or fewer sequences. The
390 final analysis included 14 studies (Tables 1 and 2). There were seven studies with only
391 fecal samples (Ahn, Baxter, Brim, Hale, Wang, Weir, and Zeller), five studies with only
392 tissue samples (Burns, Dejea, Geng, Lu, Sanapareddy), and two studies with both fecal
393 and tissue samples (Chen and Flemer). After curating the sequences, 1737 fecal samples
394 and 492 tissue samples remained in the analysis [Tables 1 and 2].

395 **Sequence Processing.** Raw sequence data and metadata were primarily obtained from
396 the Sequence Read Archive (SRA) and dbGaP. Other sequence and metadata were
397 obtained directly from the authors (n=4, (17, 23, 25, 27)). Each dataset was processed
398 separately using mothur (v1.39.3) using the default quality filtering methods for both 454
399 and Illumina sequence data (47). If it was not possible to use the defaults because the
400 trimmed sequences were too short, then the stated quality cut-offs from the original study

401 were used. Chimeric sequences were identified and removed using VSEARCH (48). The
402 curated sequences were assigned to OTUs at 97% similarity using the OptiClust algorithm
403 (49) and classified to the deepest taxonomic level that had 80% support using the naïve
404 Bayesian classifier trained on the RDP taxonomy outline (version 14, (50)).

405 **Community analysis.** We calculated alpha diversity metrics (i.e. OTU richness, evenness,
406 and Shannon diversity) for each sample. Within each dataset, we ensured that the data
407 followed a normal distribution using power transformations. Using the transformed data,
408 we tested the hypothesis that individuals with normal colons, adenomas, and carcinomas
409 had significantly different alpha diversity metrics using linear mixed-effect models. We
410 also calculated the OR for each study and metric by considering any value above the
411 median alpha diversity value to be positive. We measured the dissimilarity between
412 individuals by calculating the pairwise Bray-Curtis index and used PERMANOVA (51) to
413 test whether individuals with normal colons were significantly different from those with
414 adenomas or carcinomas. Finally, after binning sequences into the deepest taxa that
415 the naïve Bayesian classifier could classify the sequences, we quantified the ORs for
416 individuals having an adenoma or carcinoma and corrected for multiple comparisons using
417 the Benjamini-Hochberg method (52). Again, for each taxon, if the relative abundance was
418 greater than the median relative abundance for that taxon in the study, the individual was
419 considered to be positive.

420 **Random Forest classification analysis.** To classify individuals as having normal colons
421 or tumors, we built Random Forest classification models for each dataset and comparison
422 using taxa with significant ORs (after multiple comparison correction), all taxa, or OTUs.
423 Because no taxa were identified as having a significant OR associated with adenomas
424 using stool or tissue samples, classification models based on OR data were not constructed
425 to classify individuals as having normal colons or adenomas. For all models, the value of
426 trees included (i.e. ntree) was set to 500 and the number of variables that were randomly

427 tested (i.e. mtry) was set to the square root of the number of taxa or OTUs within the
428 model. Using the square root of the total number of features as the number of features
429 to test has been found to reliably approximate the optimum value after model tuning (53).
430 All fecal models were built using a 10-fold cross validation (CV) while tissue models were
431 built using 5-fold CV due to study sample size. One exception to this were the models
432 constructed using data from the Weir study, which was built using a 2-fold CV due to
433 the small number of samples. For models constructed based on the taxa that had a
434 significant OR or using all of the taxa, we trained the models using a single study and then
435 tested on the remaining studies with AUCs recorded during both train and testing phases.
436 For the models constructed using OTU data, 100 10-fold CVs were run to generate a
437 range of AUCs that could be reasonably expected to occur. The average AUC from these
438 100 repeats was reported. The Mean Decrease in Accuracy (MDA), a measure of the
439 importance of each taxon to the overall model, was used to rank the taxa used in each
440 model.

441 **Statistical Analysis.** All statistical analysis after sequence processing utilized the R
442 (v3.4.3) software package (54). For OTU richness, evenness, and Shannon diversity
443 analysis, values were power transformed using the rcompanion (v1.11.1) package (55)
444 and Z-score normalized using the car (v2.1.6) package (56). Testing for OTU richness,
445 evenness, and Shannon diversity differences utilized linear mixed-effect models to correct
446 for study, repeat sampling of individuals (tissue only), and 16S rRNA gene sequence
447 region used using the lme4 (v1.1.15) package (57). ORs were analyzed using both the
448 epiR (v0.9.93) and metafor (v2.0.0) packages (58, 59) by assessing how many individuals
449 with and without disease were above and below the overall median value within each
450 specific study. OR significance testing utilized the chi-squared test. Community structure
451 differences were calculated using the Bray-Curtis dissimilarity index and PERMANOVA was
452 used to test for tumor-associated differences in structure with the vegan (v2.4.5) package
453 (60). Random Forest models were built using both the caret (v6.0.78) and randomForest

454 (v4.6.12) packages (61, 62). All figures were created using both ggplot2 (v2.2.1) and
455 gridExtra (v2.3) packages (63, 64).

456 **Reproducible Methods.** The analysis code can be found at [https://github.com/](https://github.com/SchlossLab/Sze_CRCMetaAnalysis_mBio_2018)
457 SchlossLab/Sze_CRCMetaAnalysis_mBio_2018. Unless otherwise mentioned, the
458 accession number of raw sequences from the studies used in this analysis can be found
459 directly in the respective batch file in the GitHub repository or in the original manuscript.

460 **Acknowledgements**

461 The authors would like to thank all the study participants who were a part of each of the
462 individual studies analyzed. We would also like to thank each of the study authors for
463 making their sequencing reads and metadata available for use. Finally, we would like to
464 thank the members of the Schloss lab for their valuable feedback and proofreading during
465 the formulation of this manuscript.

466 **References**

- 467 1. **Siegel RL, Miller KD, Jemal A.** 2016. Cancer statistics, 2016. *CA: a cancer journal for*
468 *clinicians* **66**:7–30. doi:10.3322/caac.21332.
- 469 2. **Flynn KJ, Baxter NT, Schloss PD.** 2016. Metabolic and Community Synergy of Oral
470 Bacteria in Colorectal Cancer. *mSphere* **1**. doi:10.1128/mSphere.00102-16.
- 471 3. **Goodwin AC, Destefano Shields CE, Wu S, Huso DL, Wu X, Murray-Stewart TR,**
472 **Hacker-Prietz A, Rabizadeh S, Woster PM, Sears CL, Casero RA.** 2011. Polyamine
473 catabolism contributes to enterotoxigenic *Bacteroides fragilis*-induced colon tumorigenesis.
474 *Proceedings of the National Academy of Sciences of the United States of America*
475 **108**:15354–15359. doi:10.1073/pnas.1010203108.
- 476 4. **Abed J, Emgård JEM, Zamir G, Faroja M, Almogy G, Grenov A, Sol A, Naor R,**
477 **Pikarsky E, Atlan KA, Mellul A, Chaushu S, Manson AL, Earl AM, Ou N, Brennan CA,**
478 **Garrett WS, Bachrach G.** 2016. Fap2 Mediates *Fusobacterium nucleatum* Colorectal
479 Adenocarcinoma Enrichment by Binding to Tumor-Expressed Gal-GalNAc. *Cell Host &*
480 *Microbe* **20**:215–225. doi:10.1016/j.chom.2016.07.006.
- 481 5. **Arthur JC, Perez-Chanona E, Mühlbauer M, Tomkovich S, Uronis JM, Fan T-J,**
482 **Campbell BJ, Abujamel T, Dogan B, Rogers AB, Rhodes JM, Stintzi A, Simpson**
483 **KW, Hansen JJ, Keku TO, Fodor AA, Jobin C.** 2012. Intestinal inflammation targets
484 cancer-inducing activity of the microbiota. *Science (New York, NY)* **338**:120–123.
485 doi:10.1126/science.1224820.
- 486 6. **Kostic AD, Chun E, Robertson L, Glickman JN, Gallini CA, Michaud M, Clancy**
487 **TE, Chung DC, Lochhead P, Hold GL, El-Omar EM, Brenner D, Fuchs CS, Meyerson**
488 **M, Garrett WS.** 2013. *Fusobacterium nucleatum* potentiates intestinal tumorigenesis
489 and modulates the tumor-immune microenvironment. *Cell Host & Microbe* **14**:207–215.

490 doi:10.1016/j.chom.2013.07.007.

491 7. **Wu S, Rhee K-J, Albesiano E, Rabizadeh S, Wu X, Yen H-R, Huso DL, Brancati FL,**
492 **Wick E, McAllister F, Housseau F, Pardoll DM, Sears CL.** 2009. A human colonic
493 commensal promotes colon tumorigenesis via activation of T helper type 17 T cell
494 responses. *Nature Medicine* **15**:1016–1022. doi:10.1038/nm.2015.

495 8. **Zackular JP, Baxter NT, Chen GY, Schloss PD.** 2016. Manipulation of the Gut
496 Microbiota Reveals Role in Colon Tumorigenesis. *mSphere* **1**. doi:10.1128/mSphere.00001-15.

497 9. **Zackular JP, Baxter NT, Iverson KD, Sadler WD, Petrosino JF, Chen GY, Schloss**
498 **PD.** 2013. The gut microbiome modulates colon tumorigenesis. *mBio* **4**:e00692–00613.
499 doi:10.1128/mBio.00692-13.

500 10. **Baxter NT, Zackular JP, Chen GY, Schloss PD.** 2014. Structure of the gut
501 microbiome following colonization with human feces determines colonic tumor burden.
502 *Microbiome* **2**:20. doi:10.1186/2049-2618-2-20.

503 11. **Ahn J, Sinha R, Pei Z, Dominianni C, Wu J, Shi J, Goedert JJ, Hayes RB, Yang**
504 **L.** 2013. Human gut microbiome and risk for colorectal cancer. *Journal of the National*
505 *Cancer Institute* **105**:1907–1911. doi:10.1093/jnci/djt300.

506 12. **Baxter NT, Ruffin MT, Rogers MAM, Schloss PD.** 2016. Microbiota-based model
507 improves the sensitivity of fecal immunochemical test for detecting colonic lesions. *Genome*
508 *Medicine* **8**:37. doi:10.1186/s13073-016-0290-3.

509 13. **Chen W, Liu F, Ling Z, Tong X, Xiang C.** 2012. Human intestinal lumen and
510 mucosa-associated microbiota in patients with colorectal cancer. *PloS One* **7**:e39743.
511 doi:10.1371/journal.pone.0039743.

512 14. **Wang T, Cai G, Qiu Y, Fei N, Zhang M, Pang X, Jia W, Cai S, Zhao L.** 2012.

- 513 Structural segregation of gut microbiota between colorectal cancer patients and healthy
514 volunteers. *The ISME journal* **6**:320–329. doi:10.1038/ismej.2011.109.
- 515 **15. Burns MB, Lynch J, Starr TK, Knights D, Blehman R.** 2015. Virulence genes are
516 a signature of the microbiome in the colorectal tumor microenvironment. *Genome Medicine*
517 **7**:55. doi:10.1186/s13073-015-0177-8.
- 518 **16. Zeller G, Tap J, Voigt AY, Sunagawa S, Kultima JR, Costea PI, Amiot A, Böhm J,**
519 **Brunetti F, Habermann N, Hercog R, Koch M, Luciani A, Mende DR, Schneider MA,**
520 **Schrotz-King P, Tournigand C, Tran Van Nhieu J, Yamada T, Zimmermann J, Benes**
521 **V, Kloor M, Ulrich CM, Knebel Doeberitz M von, Sobhani I, Bork P.** 2014. Potential of
522 fecal microbiota for early-stage detection of colorectal cancer. *Molecular Systems Biology*
523 **10**:766.
- 524 **17. Flemer B, Lynch DB, Brown JMR, Jeffery IB, Ryan FJ, Claesson MJ, O’Riordain**
525 **M, Shanahan F, O’Toole PW.** 2017. Tumour-associated and non-tumour-associated
526 microbiota in colorectal cancer. *Gut* **66**:633–643. doi:10.1136/gutjnl-2015-309595.
- 527 **18. García AZG.** 2012. Factors influencing colorectal cancer screening participation.
528 *Gastroenterology Research and Practice* **2012**:1–8. doi:10.1155/2012/483417.
- 529 **19. Geng J, Fan H, Tang X, Zhai H, Zhang Z.** 2013. Diversified pattern of the human
530 colorectal cancer microbiome. *Gut Pathogens* **5**:2. doi:10.1186/1757-4749-5-2.
- 531 **20. Dejea CM, Wick EC, Hechenbleikner EM, White JR, Mark Welch JL, Rossetti**
532 **BJ, Peterson SN, Snedrud EC, Borisy GG, Lazarev M, Stein E, Vadivelu J, Roslani**
533 **AC, Malik AA, Wanyiri JW, Goh KL, Thevambiga I, Fu K, Wan F, Llosa N, Housseau**
534 **F, Romans K, Wu X, McAllister FM, Wu S, Vogelstein B, Kinzler KW, Pardoll DM,**
535 **Sears CL.** 2014. Microbiota organization is a distinct feature of proximal colorectal
536 cancers. *Proceedings of the National Academy of Sciences of the United States of*

537 America **111**:18321–18326. doi:10.1073/pnas.1406199111.

538 **21. Arthur JC, Gharaibeh RZ, Mühlbauer M, Perez-Chanona E, Uronis JM,**
539 **McCafferty J, Fodor AA, Jobin C.** 2014. Microbial genomic analysis reveals the essential
540 role of inflammation in bacteria-induced colorectal cancer. *Nature Communications* **5**:4724.
541 doi:10.1038/ncomms5724.

542 **22. Aymeric L, Donnadieu F, Mulet C, Merle L du, Nigro G, Saffarian A, Bérard**
543 **M, Poyart C, Robine S, Regnault B, Trieu-Cuot P, Sansonetti PJ, Dramsi S.**
544 2017. Colorectal cancer specific conditions promote *Streptococcus gallolyticus* gut
545 colonization. *Proceedings of the National Academy of Sciences* **115**:E283–E291.
546 doi:10.1073/pnas.1715112115.

547 **23. Weir TL, Manter DK, Sheflin AM, Barnett BA, Heuberger AL, Ryan EP.** 2013. Stool
548 microbiome and metabolome differences between colorectal cancer patients and healthy
549 adults. *PloS One* **8**:e70803. doi:10.1371/journal.pone.0070803.

550 **24. Boleij A, Hechenbleikner EM, Goodwin AC, Badani R, Stein EM, Lazarev MG,**
551 **Ellis B, Carroll KC, Albesiano E, Wick EC, Platz EA, Pardoll DM, Sears CL.** 2014. The
552 *bacteroides fragilis* toxin gene is prevalent in the colon mucosa of colorectal cancer patients.
553 *Clinical Infectious Diseases* **60**:208–215. doi:10.1093/cid/ciu787.

554 **25. Sanapareddy N, Legge RM, Jovov B, McCoy A, Burcal L, Araujo-Perez F, Randall**
555 **TA, Galanko J, Benson A, Sandler RS, Rawls JF, Abdo Z, Fodor AA, Keku TO.** 2012.
556 Increased rectal microbial richness is associated with the presence of colorectal adenomas
557 in humans. *The ISME journal* **6**:1858–1868. doi:10.1038/ismej.2012.43.

558 **26. Lu Y, Chen J, Zheng J, Hu G, Wang J, Huang C, Lou L, Wang X, Zeng Y.** 2016.
559 Mucosal adherent bacterial dysbiosis in patients with colorectal adenomas. *Scientific*

560 Reports **6**:26337. doi:10.1038/srep26337.

561 **27. Hale VL, Chen J, Johnson S, Harrington SC, Yab TC, Smyrk TC, Nelson H,**
562 **Boardman LA, Druliner BR, Levin TR, Rex DK, Ahnen DJ, Lance P, Ahlquist DA,**
563 **Chia N.** 2017. Shifts in the Fecal Microbiota Associated with Adenomatous Polyps. *Cancer*
564 *Epidemiology, Biomarkers & Prevention: A Publication of the American Association for*
565 *Cancer Research, Cosponsored by the American Society of Preventive Oncology* **26**:85–94.
566 doi:10.1158/1055-9965.EPI-16-0337.

567 **28. Shah MS, DeSantis TZ, Weinmaier T, McMurdie PJ, Cope JL, Altrichter**
568 **A, Yamal J-M, Hollister EB.** 2017. Leveraging sequence-based faecal microbial
569 community survey data to identify a composite biomarker for colorectal cancer. *Gut*.
570 doi:10.1136/gutjnl-2016-313189.

571 **29. Brim H, Yooseph S, Zoetendal EG, Lee E, Torralbo M, Laiyemo AO, Shokrani**
572 **B, Nelson K, Ashktorab H.** 2013. Microbiome analysis of stool samples from African
573 Americans with colon polyps. *PloS One* **8**:e81352. doi:10.1371/journal.pone.0081352.

574 **30. Sze MA, Baxter NT, Ruffin MT, Rogers MAM, Schloss PD.** 2017. Normalization
575 of the microbiota in patients after treatment for colonic lesions. *Microbiome* **5**.
576 doi:10.1186/s40168-017-0366-3.

577 **31. Hannigan GD, Duhaime MB, Ruffin MT, Koumpouras CC, Schloss PD.** 2017.
578 Diagnostic potential & the interactive dynamics of the colorectal cancer virome.
579 doi:10.1101/152868.

580 **32. Venkataraman A, Sieber JR, Schmidt AW, Waldron C, Theis KR, Schmidt TM.**
581 2016. Variable responses of human microbiomes to dietary supplementation with resistant
582 starch. *Microbiome* **4**. doi:10.1186/s40168-016-0178-x.

583 **33. Herrmann E, Young W, Reichert-Grimm V, Weis S, Riedel C, Rosendale D,**

- 584 **Stoklosinski H, Hunt M, Egert M.** 2018. In vivo assessment of resistant starch
585 degradation by the caecal microbiota of mice using RNA-based stable isotope probingA
586 proof-of-principle study. *Nutrients* **10**:179. doi:10.3390/nu10020179.
- 587 34. **Reichardt N, Vollmer M, Holtrop G, Farquharson FM, Wefers D, Bunzel M,**
588 **Duncan SH, Drew JE, Williams LM, Milligan G, Preston T, Morrison D, Flint HJ,**
589 **Louis P.** 2017. Specific substrate-driven changes in human faecal microbiota composition
590 contrast with functional redundancy in short-chain fatty acid production. *The ISME Journal*
591 **12**:610–622. doi:10.1038/ismej.2017.196.
- 592 35. **Flynn KJ, Ruffin MT, Turgeon DK, Schloss PD.** 2018. Spatial variation of the native
593 colon microbiota in healthy adults. *Cancer Prevention and Research* **In Press**.
- 594 36. **Repass J, Iorns E, Denis A, Williams SR, Perfito N, and TME.** 2018. Replication
595 study: *Fusobacterium nucleatum* infection is prevalent in human colorectal carcinoma.
596 *eLife* **7**. doi:10.7554/elife.25801.
- 597 37. **Salter SJ, Cox MJ, Turek EM, Calus ST, Cookson WO, Moffatt MF, Turner P,**
598 **Parkhill J, Loman NJ, Walker AW.** 2014. Reagent and laboratory contamination
599 can critically impact sequence-based microbiome analyses. *BMC Biology* **12**.
600 doi:10.1186/s12915-014-0087-z.
- 601 38. **Purcell RV, Pearson J, Aitchison A, Dixon L, Frizelle FA, Keenan JI.** 2017.
602 Colonization with enterotoxigenic *Bacteroides fragilis* is associated with early-stage
603 colorectal neoplasia. *PLOS ONE* **12**:e0171602. doi:10.1371/journal.pone.0171602.
- 604 39. **Sze MA, Schloss PD.** 2016. Looking for a signal in the noise: Revisiting obesity and
605 the microbiome. *mBio* **7**:e01018–16. doi:10.1128/mbio.01018-16.
- 606 40. **Walters WA, Xu Z, Knight R.** 2014. Meta-analyses of human gut microbes associated

607 with obesity and IBD. *FEBS Letters* **588**:4223–4233. doi:10.1016/j.febslet.2014.09.039.

608 **41. Finucane MM, Sharpton TJ, Laurent TJ, Pollard KS.** 2014. A taxonomic signature
609 of obesity in the microbiome? Getting to the guts of the matter. *PLoS ONE* **9**:e84689.
610 doi:10.1371/journal.pone.0084689.

611 **42. Hanage WP.** 2014. Microbiology: Microbiome science needs a healthy dose of
612 scepticism. *Nature* **512**:247–248. doi:10.1038/512247a.

613 **43. Keku TO, Dulal S, Deveaux A, Jovov B, Han X.** 2015. The gastrointestinal microbiota
614 and colorectal cancer. *American Journal of Physiology - Gastrointestinal and Liver*
615 *Physiology* **308**:G351–G363. doi:10.1152/ajpgi.00360.2012.

616 **44. Vogtmann E, Goedert JJ.** 2016. Epidemiologic studies of the human microbiome and
617 cancer. *British Journal of Cancer* **114**:237–242. doi:10.1038/bjc.2015.465.

618 **45. Kostic AD, Gevers D, Pedamallu CS, Michaud M, Duke F, Earl AM, Ojesina AI,**
619 **Jung J, Bass AJ, Tabernero J, Baselga J, Liu C, Shivdasani RA, Ogino S, Birren**
620 **BW, Huttenhower C, Garrett WS, Meyerson M.** 2012. Genomic analysis identifies
621 association of *Fusobacterium* with colorectal carcinoma. *Genome Research* **22**:292–298.
622 doi:10.1101/gr.126573.111.

623 **46. Zackular JP, Rogers MAM, Ruffin MT, Schloss PD.** 2014. The human gut
624 microbiome as a screening tool for colorectal cancer. *Cancer Prevention Research*
625 (Philadelphia, Pa) **7**:1112–1121. doi:10.1158/1940-6207.CAPR-14-0129.

626 **47. Schloss PD, Westcott SL, Ryabin T, Hall JR, Hartmann M, Hollister**
627 **EB, Lesniewski RA, Oakley BB, Parks DH, Robinson CJ, Sahl JW, Stres B,**
628 **Thallinger GG, Van Horn DJ, Weber CF.** 2009. Introducing mothur: Open-Source,
629 Platform-Independent, Community-Supported Software for Describing and Comparing

630 Microbial Communities. *Appl Environ Microbiol* **75**:7537–7541.

631 48. **Rognes T, Flouri T, Nichols B, Quince C, Mahé F.** 2016. VSEARCH: A versatile
632 open source tool for metagenomics. *PeerJ* **4**:e2584. doi:10.7717/peerj.2584.

633 49. **Westcott SL, Schloss PD.** 2017. OptiClust, an Improved Method for Assigning
634 Amplicon-Based Sequence Data to Operational Taxonomic Units. *mSphere* **2**.
635 doi:10.1128/mSphereDirect.00073-17.

636 50. **Wang Q, Garrity GM, Tiedje JM, Cole JR.** 2007. Naive bayesian classifier for
637 rapid assignment of rRNA sequences into the new bacterial taxonomy. *Applied and*
638 *Environmental Microbiology* **73**:5261–5267. doi:10.1128/aem.00062-07.

639 51. **Anderson MJ, Walsh DCI.** 2013. PERMANOVA, ANOSIM, and the mantel test in
640 the face of heterogeneous dispersions: What null hypothesis are you testing? *Ecological*
641 *Monographs* **83**:557–574. doi:10.1890/12-2010.1.

642 52. **Benjamini Y, Hochberg Y.** 1995. Controlling the false discovery rate: A practical and
643 powerful approach to multiple testing. *Journal of the Royal Statistical Society Series B*
644 (Methodological) **57**:289–300.

645 53. **Breiman L.** 2001. *Machine Learning* **45**:5–32. doi:10.1023/a:1010933404324.

646 54. **R Core Team.** 2017. R: A language and environment for statistical computing. R
647 Foundation for Statistical Computing, Vienna, Austria.

648 55. **Mangiafico S.** 2017. Rcompanion: Functions to support extension education program
649 evaluation.

650 56. **Fox J, Weisberg S.** 2011. *An R companion to applied regression* Second. Sage,

651 Thousand Oaks CA.

652 57. **Bates D, Mächler M, Bolker B, Walker S.** 2015. Fitting linear mixed-effects models
653 using lme4. *Journal of Statistical Software* **67**:1–48. doi:10.18637/jss.v067.i01.

654 58. **Telmo Nunes MS with contributions from, Heuer C, Marshall J, Sanchez J,**
655 **Thornton R, Reiczigel J, Robison-Cox J, Sebastiani P, Solymos P, Yoshida K, Jones**
656 **G, Pirikahu S, Firestone S, Kyle. R.** 2017. EpiR: Tools for the analysis of epidemiological
657 data.

658 59. **Viechtbauer W.** 2010. Conducting meta-analyses in R with the metafor package.
659 *Journal of Statistical Software* **36**:1–48.

660 60. **Oksanen J, Blanchet FG, Friendly M, Kindt R, Legendre P, McGlenn D, Minchin**
661 **PR, O’Hara RB, Simpson GL, Solymos P, Stevens MHH, Szoecs E, Wagner H.** 2017.
662 *Vegan: Community ecology package.*

663 61. **Jed Wing MKC from, Weston S, Williams A, Keefer C, Engelhardt A, Cooper**
664 **T, Mayer Z, Kenkel B, R Core Team, Benesty M, Lescarbeau R, Ziem A, Scrucca L,**
665 **Tang Y, Candan C, Hunt. T.** 2017. *Caret: Classification and regression training.*

666 62. **Liaw A, Wiener M.** 2002. Classification and regression by randomForest. *R News*
667 **2**:18–22.

668 63. **Wickham H.** 2009. *Ggplot2: Elegant graphics for data analysis.* Springer-Verlag New
669 York.

670 64. **Auguie B.** 2017. *GridExtra: Miscellaneous functions for “grid” graphics.*

671 **Table 1: Characteristics of the datasets included in the fecal-based analysis**

Study	Data Stored	Region	Control (n)	Adenoma (n)	Carcinoma (n)
Ahn	DBGap	V3-4	148	0	62
Baxter	SRA	V4	172	198	120
Brim	SRA	V1-3	6	6	0
Flemer	Author	V3-4	37	0	43
Hale	Author	V3-5	473	214	17
Wang	SRA	V3	56	0	46
Weir	Author	V4	4	0	7
Zeller	SRA	V4	50	37	41

672 **Table 2: Characteristics of the datasets included in the tissue-based analyses**

Study	Data Stored	Region	Control (n)	Adenoma (n)	Carcinoma (n)
Burns	SRA	V5-6	18	0	16
Chen	SRA	V1-3	9	0	9
Dejea	SRA	V3-5	31	0	32
Flemer	Author	V3-4	103	37	94
Geng	SRA	V1-2	16	0	16
Lu	SRA	V3-4	20	20	0
Sanapareddy	Author	V1-2	38	0	33

673 **Figure 1: Comparison of alpha diversity indices that were significant between**
674 **individuals with normal colons, and those with adenomas or carcinomas using**
675 **data collected from fecal samples** A) Comparison of evenness between individuals with
676 normal colons and adenomas. B) Comparison of evenness between individuals with
677 normal colons and carcinomas. C) Comparison of Shannon diversity between individuals
678 with normal colons and carcinomas. Blue points represent individuals with normal colons
679 and red points represent individuals with either adenomas (panel A) or carcinomas (panel
680 B and C). The black lines represent the median value for each group.

681 **Figure 2: Comparison of odds ratios calculated using alpha diversity community**
682 **metrics associated with the presence of adenomas (A) or carcinoma (B) relative to**
683 **those in individuals with normal colons using data collected from stool samples.**

684 **Figure 3: AUC values when classifying individuals as having normal colons or**
685 **carcinomas using taxa with significant ORs when using stool samples (A) and**
686 **unmatched tissue samples (B).** We did not identify any taxa as having a significant OR
687 to differentiate individuals with normal colons and adenomas or using matched tissue
688 samples. The large black circles represent the median AUC of all studies and the smaller
689 circles represent the individual AUC for a particular study. The dotted line denotes an AUC
690 of 0.5.

691 **Figure 4: Relative importance of taxa with significant ORs in Random Forest**
692 **models for differentiating between individuals with normal colons and carcinomas**
693 **using stool samples (A) or unmatched tissue samples (B).** The colors indicate the
694 z-transformed (i.e. mean of 0.0 and standard deviation of 1.0) mean decrease in accuracy
695 values calculated from the model for each study. The taxa are ranked by their mean
696 z-score-transformed mean decrease in accuracy.

697 **Figure 5: Comparison of Random Forest modeling approaches to classify**

698 **individuals as having normal colons or adenomas (A) or carcinomas (B) when**
699 **training the models using the taxa with significant ORs, all taxa in a community, or**
700 **all OTUs in a community when using stool samples.** No taxa had a significant OR
701 associated with the presence of adenomas using stool samples. The black line represents
702 the median AUC for the respective group. The dashed gray line indicates an AUC of 0.5.

703 **Figure 6: Testing of Random Forest models to classify individuals as having normal**
704 **colons or adenomas (A) or carcinomas (B) when using sequence data obtained**
705 **from stool samples.** Models were trained on data from each study (Figure 5) and tested
706 on the other studies. The black lines represent the median AUC of all test AUCs for a
707 specific study. The dashed gray line represents the AUC at 0.5.

708 **Figure S1: Comparison of Odds Ratios associated with normal colons or adenomas**
709 **(A) or carcinomas (B) calculated using alpha diversity indices with sequence data**
710 **generated from tissue samples.** The pooled results are from the aggregation of data
711 across all studies. The horizontal lines indicate the 95% confidence interval for the OR.

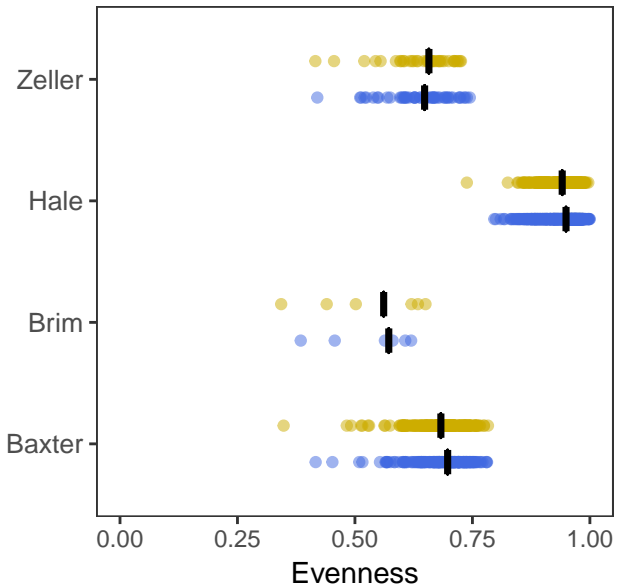
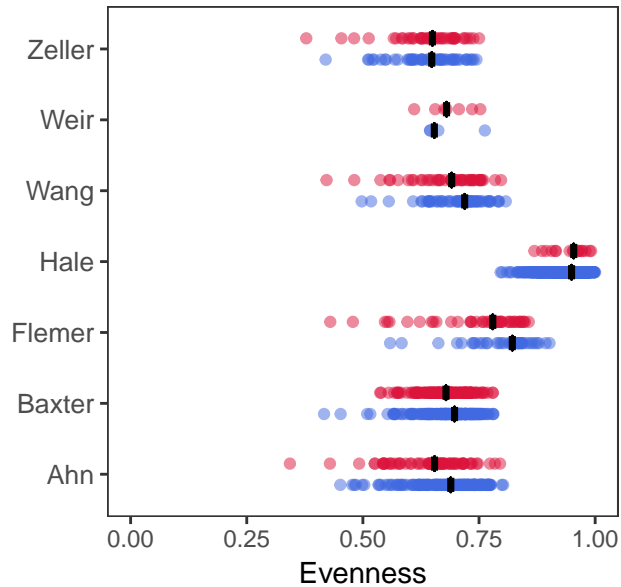
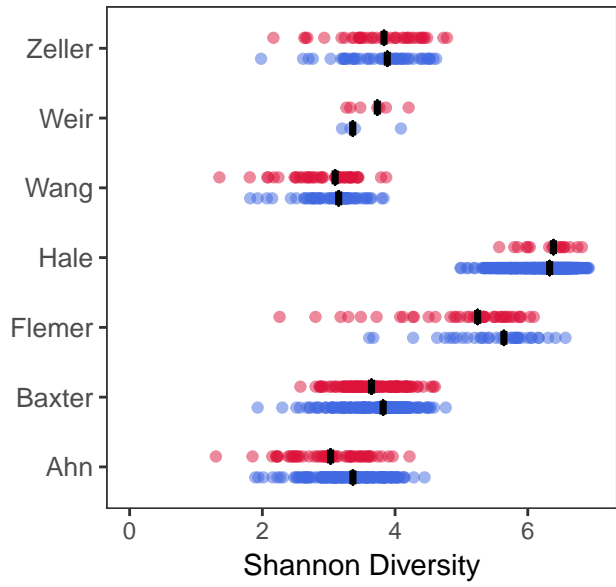
712 **Figure S2: Comparison of Random Forest modeling approaches to classify**
713 **individuals as having normal colons or adenomas (A) or carcinomas (B) when**
714 **training the models using the taxa with significant ORs, all taxa in a community,**
715 **or all OTUs in a community when using data from tissue samples.** No taxa had a
716 significant OR associated with the presence of adenomas using tissue samples. The black
717 line represents the median AUC for the respective group. The dashed gray line indicates
718 an AUC of 0.5.

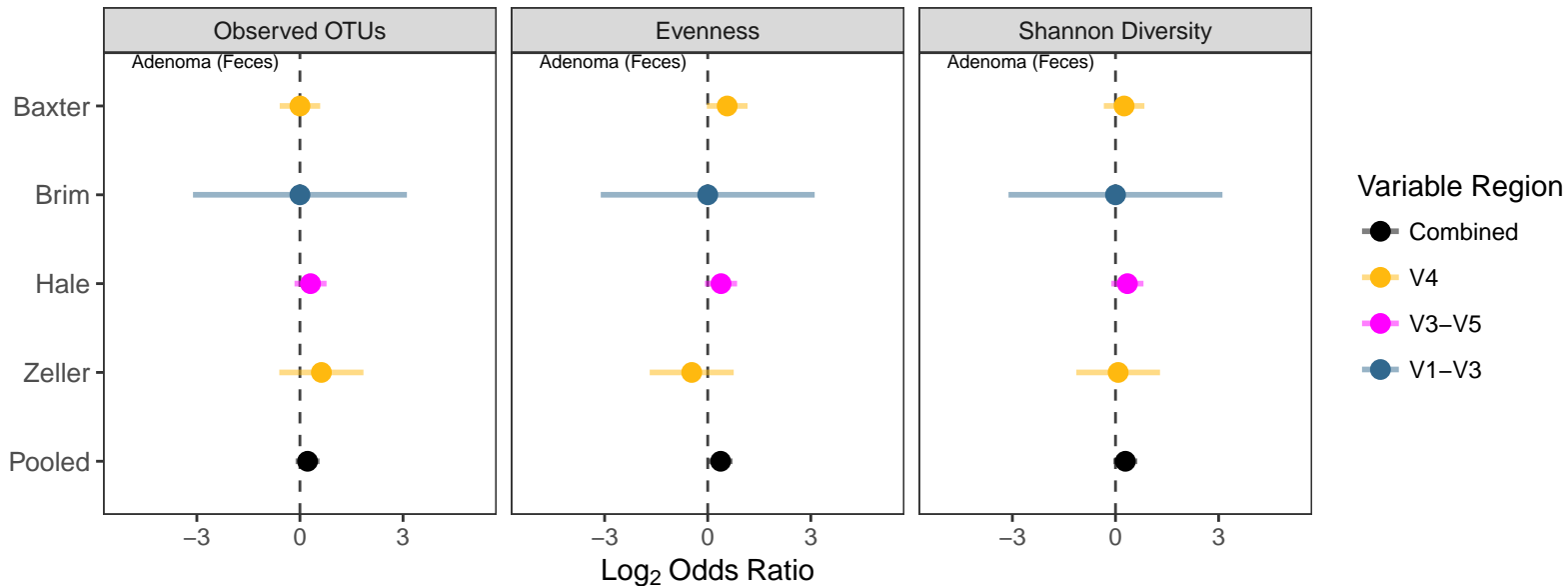
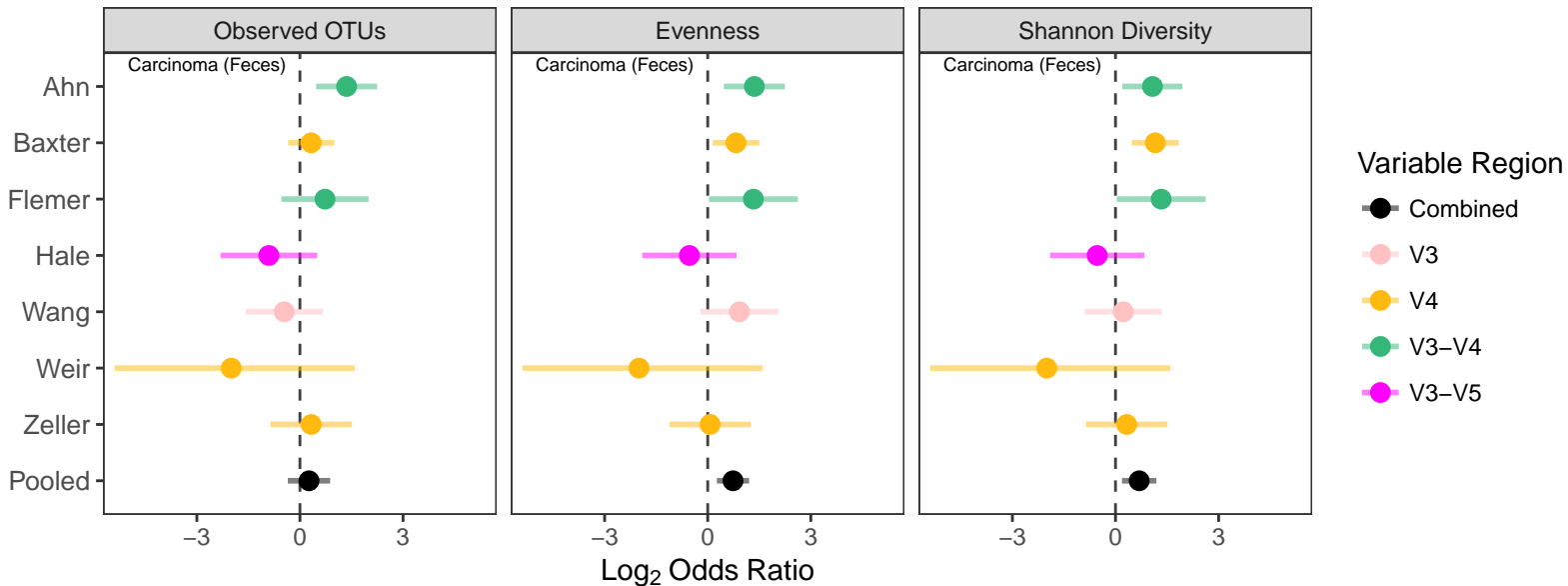
719 **Figure S3: Relative importance of taxa (A) and OTUs (B) in Random Forest models**
720 **for differentiating between individuals with normal colons and carcinomas using**
721 **stool samples.** These taxa and OTUs were among the top 10% most important features
722 in each model. The colors indicate the z-transformed (i.e. mean of 0.0 and standard
723 deviation of 1.0) mean decrease in accuracy values calculated from the model for each
724 study. The taxa are ranked by their mean z-score-transformed mean decrease in accuracy.

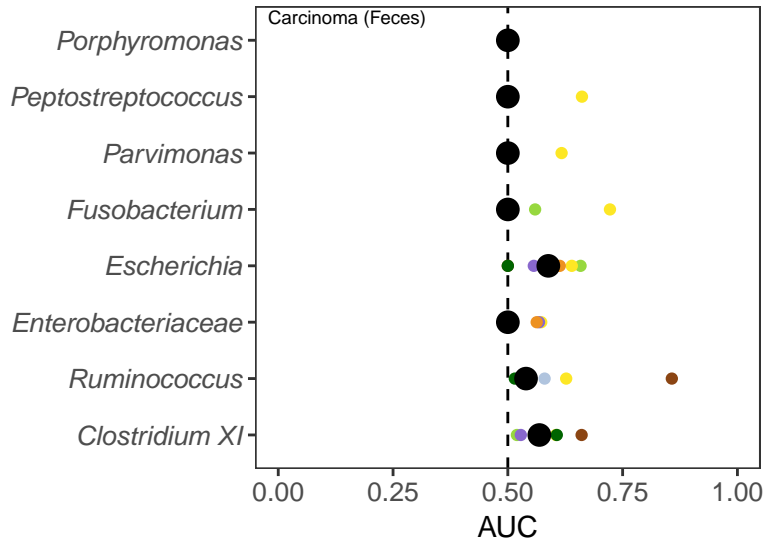
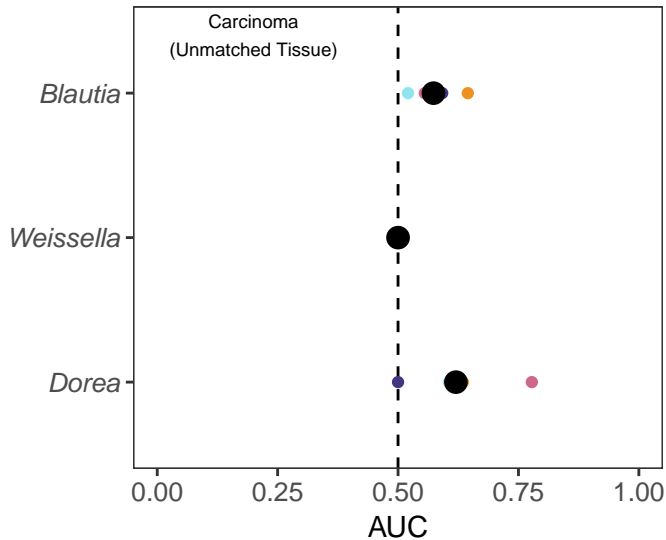
725 **Figure S4: Relative importance of taxa (A, B) and OTUs (C, D) in Random Forest**
726 **models for differentiating between individuals with normal colons and carcinomas**
727 **using matched (A, C) and unmatched (B, D) tissue samples.** These taxa and OTUs
728 were among the top 10% most important features in each model. The colors indicate the
729 z-transformed (i.e. mean of 0.0 and standard deviation of 1.0) mean decrease in accuracy
730 values calculated from the model for each study. The taxa are ranked by their mean
731 z-score-transformed mean decrease in accuracy.

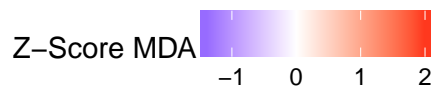
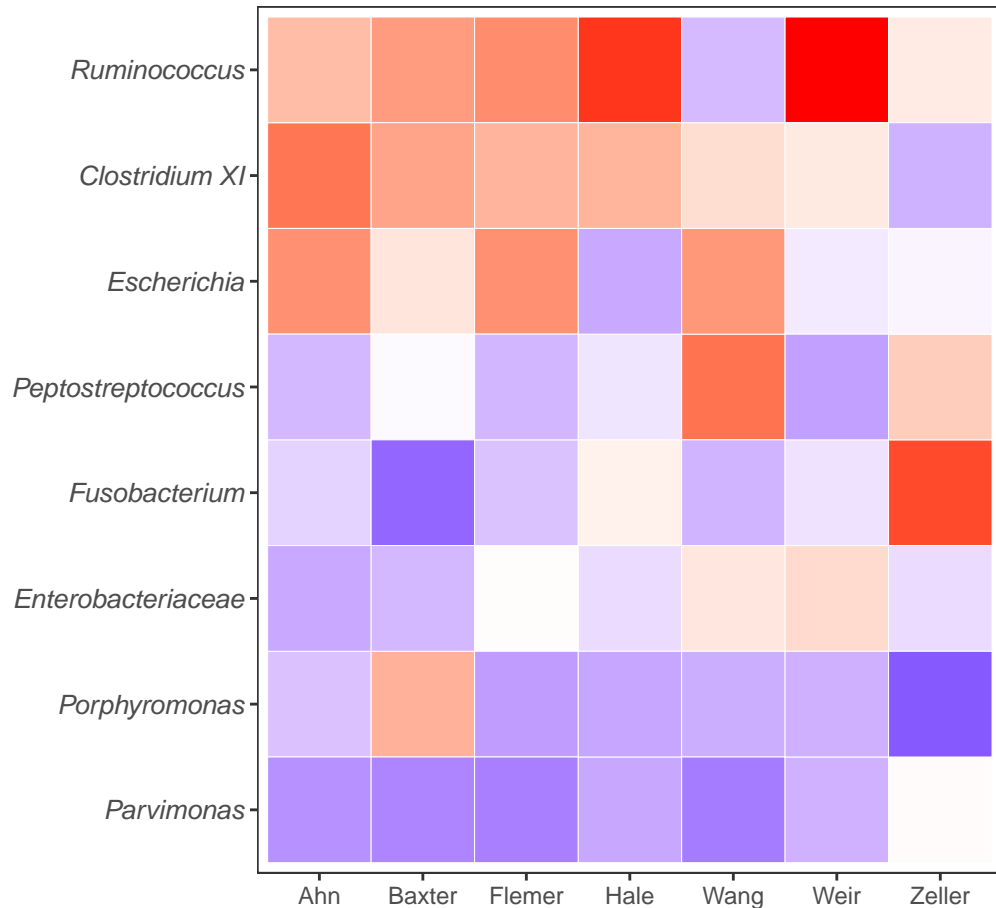
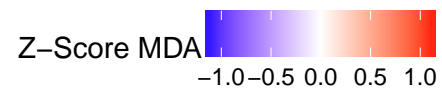
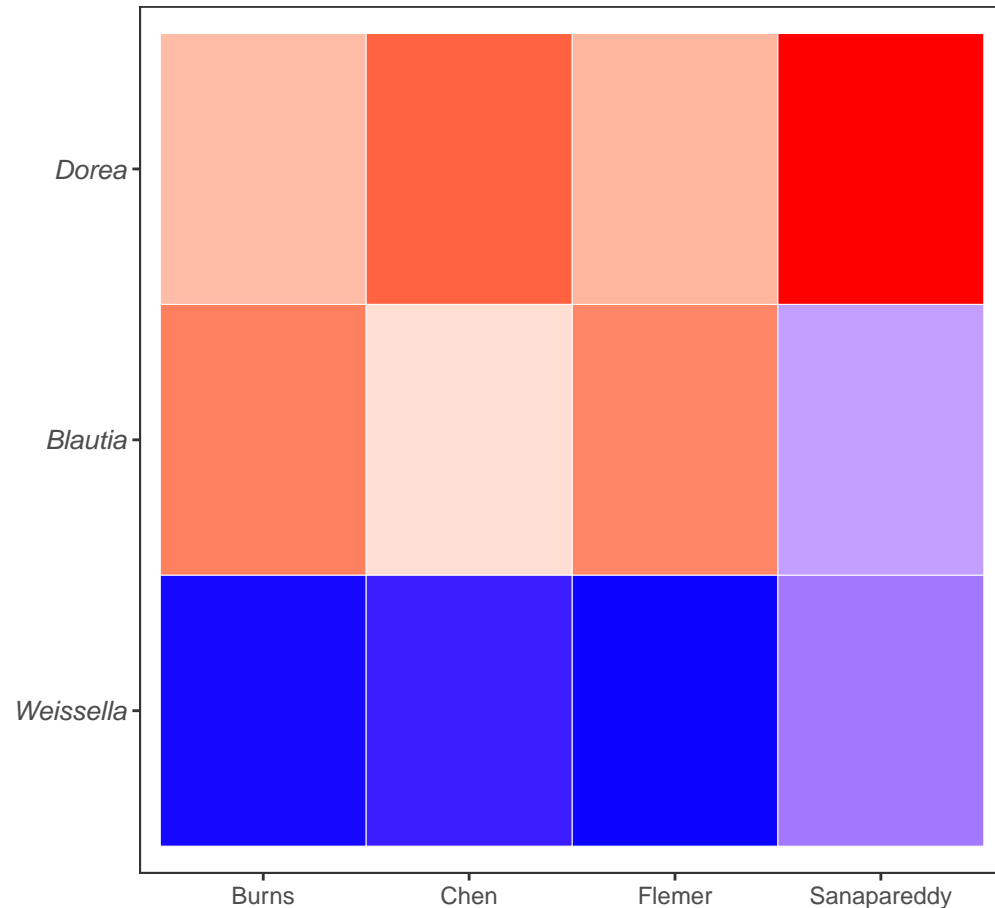
732 **Figure S5: Testing of Random Forest models to classify individuals as having**

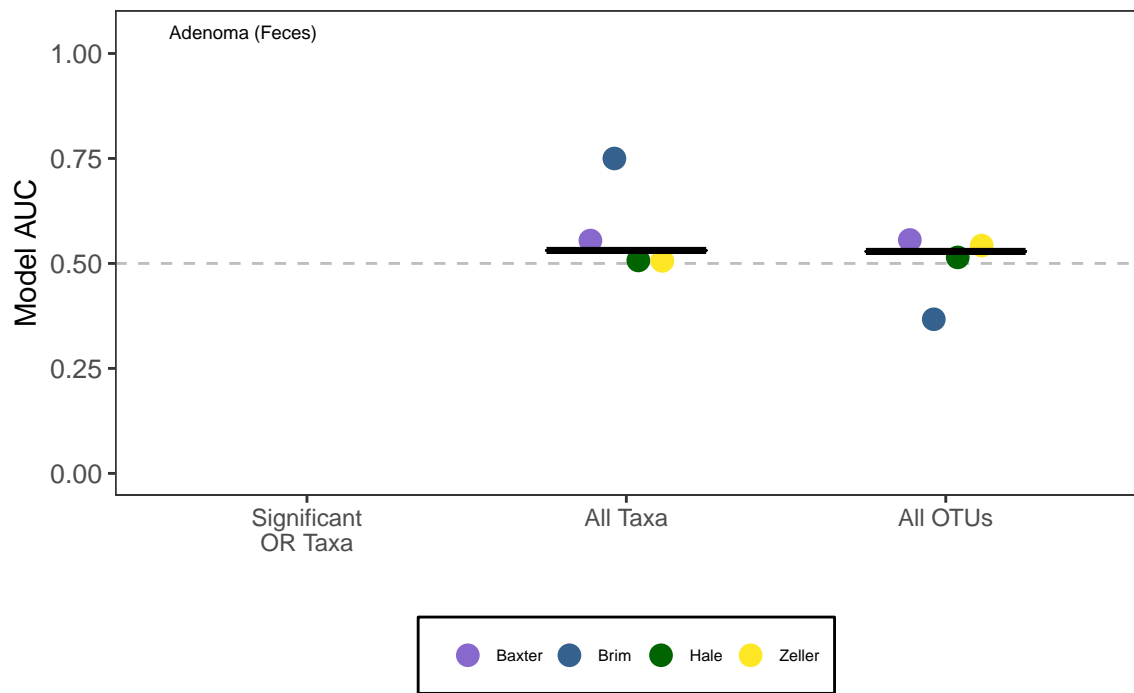
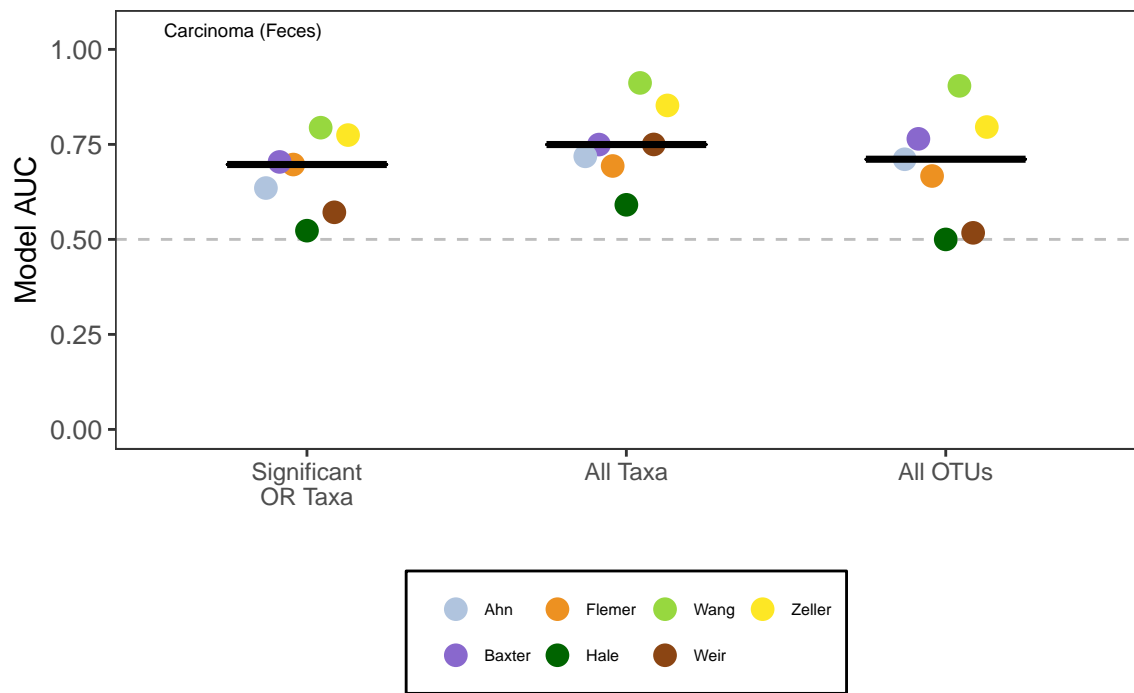
733 **normal colons or adenomas (A) or carcinomas (B, C) when using sequence data**
734 **obtained from tissue samples.** Models were trained on data from each study (Figure
735 S5) and tested on the other studies. The black lines represent the median AUC of all test
736 AUCs for a specific study. The dashed gray line represents the AUC at 0.5.

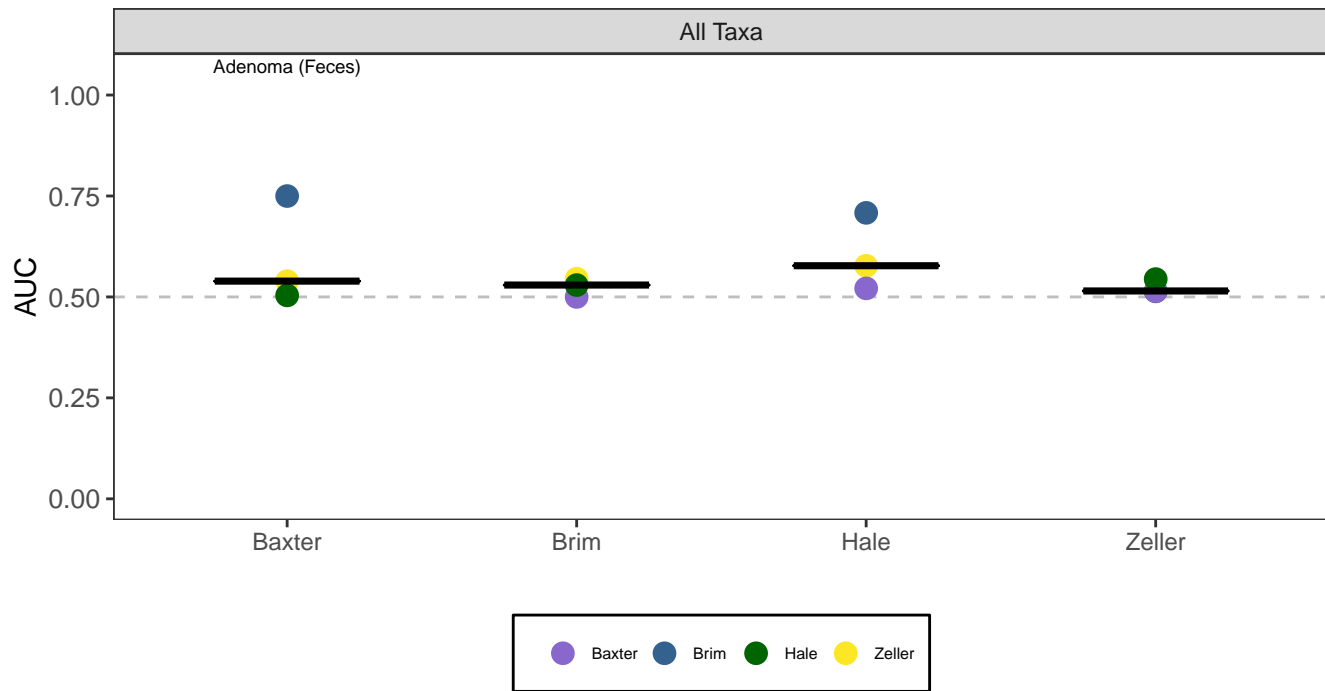
A**B****C**

A**B**

A**B**

A**B**

A**B**

A**B**