

# **Fecal short-chain fatty acids are not predictive of colonic tumor status and cannot be predicted based on bacterial community structure**

Marc A. Sze<sup>1</sup>, Begüm D. Topçuoğlu<sup>1</sup>, Nicholas A. Lesniak<sup>1</sup>, Mack T. Ruffin IV<sup>2</sup>, Patrick D. Schloss<sup>1†</sup>

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

1 Department of Microbiology and Immunology, University of Michigan, Ann Arbor, MI 48109

2 Department of Family Medicine and Community Medicine, Penn State Hershey Medical Center, Hershey, PA

## **Observation format**

## 1 Abstract

2 Colonic bacterial populations are thought to have a role in the development of colorectal cancer  
3 with some protecting against inflammation and others exacerbating inflammation. Short-chain  
4 fatty acids (SCFAs) have been shown to have anti-inflammatory properties and are produced  
5 in large quantities by colonic bacteria which produce SCFAs by fermenting fiber. We assessed  
6 whether there was an association between fecal SCFA concentrations and the presence of colonic  
7 adenomas or carcinomas in a cohort of individuals using 16S rRNA gene and metagenomic shotgun  
8 sequence data. We measured the fecal concentrations of acetate, propionate, and butyrate within  
9 the cohort and found that there were no significant associations between SCFA concentration and  
10 tumor status. When we incorporated these concentrations into random forest classification models  
11 trained to differentiate between people with normal colons and those with adenomas or carcinomas,  
12 we found that they did not significantly improve the ability of 16S rRNA gene or metagenomic gene  
13 sequence-based models to classify individuals. Finally, we generated random forest regression  
14 models trained to predict the concentration of each SCFA based on 16S rRNA gene or metagenomic  
15 gene sequence data from the same samples. These models performed poorly and were able to  
16 explain at most 14% of the observed variation in the SCFA concentrations. These results support  
17 the broader epidemiological data that questions the value of fiber consumption for reducing the  
18 risks of colorectal cancer. Although other bacterial metabolites may serve as biomarkers to detect  
19 adenomas or carcinomas, fecal SCFA concentrations have limited predictive power.

20 **Importance**

21 Considering colorectal cancer is the third leading cancer-related cause of death within the United  
22 States, it is important to detect colorectal tumors early and to prevent the formation of tumors.  
23 Short-chain fatty acids (SCFAs) are often used as a surrogate for measuring gut health and for  
24 being anti-carcinogenic because of their anti-inflammatory properties. We evaluated the fecal SCFA  
25 concentration of a cohort of individuals with varying colonic tumor burden who were previously  
26 analyzed to identify microbiome-based biomarkers of tumors. We were unable to find an association  
27 between SCFA concentration and tumor burden or use SCFAs to improve our microbiome-based  
28 models of classifying people based on their tumor status. Furthermore, we were unable to find an  
29 association between the fecal community structure and SCFA concentrations. Our results indicate  
30 that the association between fecal SCFAs, the gut microbiome, and tumor burden is weak.

31 Colorectal cancer is the third leading cancer-related cause of death within the United States (1).  
32 Less than 10% of cases can be attributed to genetic risk factors (2). This leaves a significant  
33 role for environmental, behavioral, and dietary factors (3, 4). Colorectal cancer is thought to be  
34 initiated by a series of mutations that accumulate as the mutated cells begin to proliferate leading  
35 to adenomatous lesions, which are succeeded by carcinomas (2). Throughout this progression,  
36 there are ample opportunities for bacterial populations to have a role as some bacteria are known  
37 to cause mutations, induce inflammation, and accelerate tumorigenesis (5–7). Additional cross  
38 sectional studies in humans have identified microbiome-based biomarkers of disease (8). These  
39 studies suggest that in some cases, it is the loss of bacterial populations that produce short-chain  
40 fatty acids (SCFAs) that results in increased inflammation and tumorigenesis.

41 Many microbiome studies use the concentrations of SCFAs and the presence of 16S rRNA gene  
42 sequences from organisms and the genes involved in producing them as a biomarker of a healthy  
43 microbiota (9, 10). Depending on the concentrations, SCFAs can have proliferative activities  
44 at low concentrations or anti-proliferative activities at higher concentrations; they can also have  
45 anti-inflammatory activities (11). Direct supplementation of SCFAs or feeding of fiber caused  
46 an overall reduction in tumor burden in mouse models of colorectal cancer (12). These results  
47 suggest that supplementation with fiber, which many colonic bacteria ferment to produce SCFAs,  
48 may confer beneficial effects against colorectal cancer. Regardless, there is a lack of consistent  
49 evidence that increasing SCFA concentrations can protect against colorectal cancer in humans.  
50 Case-control studies that have investigated possible associations between SCFAs and colon tumor  
51 status have been plagued by relatively small numbers of subjects, but have reported increased total  
52 and relative fecal acetate levels and decreased relative fecal butyrate concentrations in subjects with  
53 colonic lesions (13). In randomized controlled trials fiber supplementation has been inconsistently  
54 associated with protection against tumor formation and recurrence (14, 15). Such studies are  
55 confounded by difficulties ensuring subjects took the proper dose and using subjects with prior  
56 polyp history who may be beyond a point of benefiting from fiber supplementation. Together, these  
57 findings temper enthusiasm for treatments that target the production of SCFAs or for using them as  
58 biomarkers for protection against tumorigenesis.

59 **Fecal SCFA concentrations did not vary with diagnosis or treatment.** To test for a significant

60 association between colorectal cancer and SCFAs, we quantified the concentration of acetate,  
61 propionate, and butyrate in feces of previously characterized individuals with normal colons (N=172)  
62 and those with colonic adenomas (N=198) or carcinomas (N=120) (16). We were unable to detect  
63 a significant difference in any SCFA concentration across the diagnoses groups (all  $P>0.15$ ; Figure  
64 1A). Among the individuals with adenomas and carcinomas, a subset ( $N_{adenoma}=41$ ,  $N_{carcinoma}=26$ )  
65 were treated and sampled a year later (17). None of the individuals showed signs of recurrence  
66 and yet none of the SCFAs exhibited a significant change with treatment (all  $P>0.058$ ; Figure 1B).  
67 For both the pre-treatment cross-sectional data and the pre/post treatment data, we also failed to  
68 detect any significant differences in the relative concentrations of any SCFAs ( $P>0.16$ ). Finally, we  
69 pooled the SCFA concentrations on a total and per molecule of carbon basis and again failed to  
70 observe any significant differences ( $P>0.077$ ). Although some of the P-values from our analyses  
71 were close to 0.05, the effect sizes were all relatively small and inconsistent given the disease  
72 progression (Figure 1). These results demonstrated that there were no significant associations  
73 between fecal SCFA concentration and diagnosis or treatment.

74 **Combining SCFA and microbiome data does not improve the ability to diagnose individual**  
75 **as having adenomas or carcinomas using a random forest model.** We previously found that  
76 binning 16S rRNA gene sequence data into operational taxonomic units (OTUs) based on 97%  
77 similarity or into genera enabled us to classify individuals as having adenomas or carcinomas  
78 using random forest machine learning models (8, 16). We repeated that analysis but added the  
79 concentration of the SCFAs as possible features to train the models (Figure S1). Models trained  
80 using SCFAs to classify individuals as having adenomas or carcinomas rather than normal colons  
81 had median areas under the receiver operator characteristic curve (AUROC) that were significantly  
82 greater than 0.5 ( $P_{adenoma}<0.001$  and  $P_{carcinoma}<0.001$ ). However, the AUROC values to detect  
83 the presence of adenomas or carcinomas were only 0.54 and 0.55, respectively, indicating that  
84 SCFAs had poor predictive power on their own (Figure 2A). When we trained the models with  
85 the SCFAs concentrations and OTU or genus-level relative abundances the AUROC values were  
86 not significantly different from the same models trained without the SCFA concentrations ( $P>0.15$ ;  
87 Figure 2A). These data demonstrate that knowledge of the SCFA profile from a subject's fecal  
88 sample did not improve the ability to diagnose a colonic lesion.

89 **Knowledge of microbial community structure does not predict SCFA concentrations using**  
90 **a random forest model.** We next asked whether the fecal community structure was predictive  
91 of fecal SCFA concentrations, regardless of a person's diagnosis. We trained random forest  
92 regression models using 16S rRNA gene sequence data binned into OTUs and genera to predict the  
93 concentration of the SCFAs (Figure S2). The largest  $R^2$  between the observed SCFA concentrations  
94 and the modeled concentrations was 0.14, which was observed when using genus data to predict  
95 butyrate concentrations (Figure 2B). We also used a smaller dataset of shotgun metagenomic  
96 sequencing data generated from a subset of our cohort ( $N_{\text{normal}}=27$ ,  $N_{\text{adenoma}}=25$ , and  $N_{\text{cancer}}=26$ )  
97 (18). We binned genes extracted from the assembled metagenomes into operational protein families  
98 (OPFs) or KEGG categories and trained random forest regression models using metagenomic  
99 sequence data to predict the concentration of the SCFAs (Figure S2). Similar to the analysis using  
100 16S rRNA gene sequence data, the metagenomic data was not predictive of SCFA concentration.  
101 The largest  $R^2$  was 0.055, which was observed when using KEGG data to predict propionate  
102 concentrations (Figure 2B). Because of the limited number of samples that we were able to  
103 generate metagenomic sequence data from, we used our 16S rRNA gene sequence data to impute  
104 metagenomes that were binned into metabolic pathways or KEGG categories using PICRUSt  
105 (Figure S2). SCFA concentrations could not be predicted based on the imputed metagenomic  
106 data. The largest  $R^2$  was 0.085, which was observed when using KEGG data to predict propionate  
107 concentrations (Figure 2B). The inability to model SCFA concentrations from microbiome data  
108 indicates that the knowledge of the abundance of organisms and their genes was insufficient to  
109 predict fecal SCFA concentrations.

110 **Conclusion.** Our data indicate that fecal SCFA concentrations are not associated with the presence  
111 of adenomas or carcinomas and that they provide weak predictive power to improve the ability  
112 to diagnose someone with one of these lesions. Furthermore, knowledge of the taxonomic and  
113 genetic structure of gut microbiota was not meaningfully predictive of SCFA concentrations. These  
114 results complement existing literature that suggest that fiber consumption and the production of  
115 SCFAs are unable to prevent the risk of developing colonic tumors. It is important to note that our  
116 analysis was based on characterizations of SCFA and microbiome profiles using fecal samples at a  
117 single time point. Furthermore, observations along the mucosa near the site of lesions may provide

118 a stronger association. This may be a cautionary result to temper enthusiasm for SCFAs as a  
119 biomarker of gut health more generally. Going forward it is critical to develop additional hypotheses  
120 for how the microbiome and host interact to drive tumorigenesis so that we can better understand  
121 tumorigenesis and identify biomarkers that will allow early detection of lesions.

## 122 **Acknowledgements**

123 The authors thank the Great Lakes-New England Early Detection Research Network for providing  
124 the fecal samples that were used in this study. We would thank the University of Michigan Center for  
125 Microbial Systems for enabling our short-chain fatty acid analysis. Support for MAS came from the  
126 Canadian Institute of Health Research and the National Institutes of Health (UL1TR002240). This  
127 work was also supported by the National Institutes of Health (P30DK034933 and R01CA215574).

128 **Materials and Methods**

129 **Study design and sampling.** The overall study design and the resulting sequence data have  
130 been previously described (16, 17). In brief, fecal samples were obtained from 172 individuals  
131 with normal colons, 198 individuals with colonic adenomas, and 120 individuals with carcinomas.  
132 Of the individuals diagnosed as having adenomas or carcinomas, a subset ( $N_{adenoma}=41$   
133 and  $N_{carcinoma}=26$ ) were sampled after treatment of the lesion (median=255 days between  
134 sampling, IQR=233 to 334 days). Tumor diagnosis was made by colonoscopic examination and  
135 histopathological review of the biopsies (16). The University of Michigan Institutional Review Board  
136 approved the studies that generated the samples and informed consent was obtained from all  
137 participants in accordance to the guidelines set out by the Helsinki Declaration.

138 **Measuring specific SCFAs.** The measurement of acetate, propionate, isobutyrate, and butyrate  
139 used a previously published protocol that used High-Performance Liquid Chromatography (HPLC)  
140 (19). Two changes were made to the protocol. First, instead of using fecal samples suspended  
141 in DNA Genotek OmniGut tubes, we suspended frozen fecal samples in 1 mL of PBS. Second,  
142 instead of using the average weight of fecal sample aliquots to normalize SCFA concentrations, we  
143 used the actual weight of the fecal samples. These methodological changes did not affect the range  
144 of concentrations of these SCFAs between the two studies. The concentrations of isobutyrate were  
145 consistently at or below the limit of detection and were not included in our analysis.

146 **16S rRNA gene sequence data analysis.** Sequence data from Baxter et al. (16) and Sze et  
147 al. (17) were obtained from the Sequence Read Archive (studies SRP062005 and SRP096978)  
148 and reprocessed using mothur v.1.42 (20). The original studies generated sequence data from  
149 V4 region of the 16S rRNA gene using paired 250 nt reads on an Illumina MiSeq sequencer. The  
150 resulting sequence data were assembled into contigs and screened to remove low quality contigs  
151 and chimeras. The curated sequences were then clustered into OTUs at a 97% similarity threshold  
152 and assigned to the closest possible genus with an 80% confidence threshold trained on the  
153 reference collection from the Ribosomal Database Project (v.16). We used PICRUSt (v.2.1.0-b)  
154 with the recommended standard operating protocol to generate imputed metagenomes based on  
155 the expected metabolic pathways and KEGG categories (21).

156 **Metagenomic DNA sequence analysis.** A subset of the samples from the samples described by  
157 Baxter et al. (16) were used to generate metagenomic sequence data ( $N_{\text{normal}}=27$ ,  $N_{\text{adenoma}}=25$ ,  
158 and  $N_{\text{cancer}}=26$ ). These data were generated by Hannigan et al. (18) and deposited into the  
159 Sequence Read Archive (study SRP108915). Fecal DNA was subjected to shotgun sequencing on  
160 an Illumina HiSeq using 125 bp paired end reads. The archived sequences were already quality  
161 filtered and aligned to the human genome to remove contaminating sequence data. We downloaded  
162 the sequences and assembled them into contigs using MEGAHIT (22), which were used to identify  
163 open reading frames (ORFs) using Prodigal (23). We determined the abundance of each ORF  
164 by mapping the raw reads back to the ORFs using Diamond (24). We clustered the ORFs into  
165 operational protein families (OPFs) in which the clustered ORFs were more than 40% identical to  
166 each other using mmseq2 (25). We also used mmseq2 to map the ORFs to the KEGG database  
167 and clustered the ORFs according to which category the ORFs mapped.

168 **Random forest models.** The classification models were built to predict lesion type from microbiome  
169 information with or without SCFA concentrations. The regression models were built to predict the  
170 SCFA concentrations of acetate, butyrate, and propionate from microbiome information. For  
171 classification and regression models, we pre-processed the features by scaling them to vary  
172 between zero and one. Features with no variance in the training set were removed from both the  
173 training and testing sets. We randomly split the data into training and test sets so that the training  
174 set consisted of 80% of the full dataset while the test set was composed of the remaining data. The  
175 training set was used for hyperparameter selection and training the model and the test set was used  
176 for evaluating prediction performance. For each model, the best performing hyperparameter, mtry,  
177 was selected in an internal five-fold cross-validation of the training set with 100 randomizations. The  
178 mtry parameter represents the number of features randomly sampled from the available features at  
179 a question point in the classification tree (i.e. called splits of nodes) that, when answered, lead to the  
180 greatest improvement in classification. Six values of mtry were tested and the value that provided  
181 the largest AUROC or  $R^2$  was selected. We trained the random forest model using the selected  
182 mtry value and predicted the held-out test set. The data-split, hyperparameter selection, training  
183 and testing steps were repeated 100 times to get a reliable and robust reading of model prediction  
184 performance. We used AUROC and  $R^2$  as the prediction performance metric for classification

185 and regression models, respectively. We used the randomForest R package (version 4.6-14) as  
186 implemented in the caret R package (version 6.0-81) for developing and testing our models.

187 **Statistical analysis workflow.** Data summaries, statistical analysis, and data visualizations were  
188 performed using R (v.3.5.1) with the tidyverse package (v.1.2.1). To assess differences in SCFA  
189 concentrations between individuals normal colons and those with adenomas or carcinomas, we  
190 used the Kruskal-Wallis rank sum test. If a test had a P-value below 0.05, we then applied a  
191 pairwise Wilcoxon rank sum test with a Benjamini-Hochberg correction for multiple comparisons. To  
192 assess differences in SCFA concentrations between individuals samples before and after treatment  
193 we used paired Wilcoxon rank sum tests to test for significance. To compare the median AUCROC  
194 for the held out data for the model generated using only the SCFAs, we compared the distribution of  
195 the data to the expected median of 0.5 using the Wilcoxon rank sum test to test whether the model  
196 performed better than would be achieved by randomly assigning the data to each diagnosis. When  
197 we compared the random forest models generated without and with SCFA data included, we used  
198 Wilcoxon rank sum tests to determine whether the models with the SCFA data included did better.

199 **Code availability.** The code for all sequence curation and analysis steps including an Rmarkdown  
200 version of this manuscript is available at [https://github.com/SchlossLab/Sze\\_SCFACRC\\_mBio\\_](https://github.com/SchlossLab/Sze_SCFACRC_mBio_)  
201 2019/.

202 **References**

203 1. **Siegel RL, Miller KD, Jemal A.** 2016. Cancer statistics, 2016. CA: A Cancer Journal for  
204 Clinicians **66**:7–30. doi:10.3322/caac.21332.

205 2. **Fearon ER, Vogelstein B.** 1990. A genetic model for colorectal tumorigenesis. Cell **61**:759–767.  
206 doi:10.1016/0092-8674(90)90186-i.

207 3. **Fliss-Isakov N, Zelber-Sagi S, Webb M, Halpern Z, Kariv R.** 2017. Smoking habits are  
208 strongly associated with colorectal polyps in a population-based case-control study. Journal of  
209 Clinical Gastroenterology **1**. doi:10.1097/mcg.0000000000000935.

210 4. **Lee J, Jeon JY, Meyerhardt JA.** 2015. Diet and lifestyle in survivors of colorectal cancer.  
211 Hematology/Oncology Clinics of North America **29**:1–27. doi:10.1016/j.hoc.2014.09.005.

212 5. **Zackular JP, Baxter NT, Iverson KD, Sadler WD, Petrosino JF, Chen GY, Schloss PD.**  
213 2013. The gut microbiome modulates colon tumorigenesis. mBio **4**:e00692–13–e00692–13.  
214 doi:10.1128/mbio.00692-13.

215 6. **Shields CED, Meerbeke SWV, Housseau F, Wang H, Huso DL, Casero RA, O'Hagan  
HM, Sears CL.** 2016. Reduction of murine colon tumorigenesis driven by Enterotoxigenic  
217 Bacteroides fragilis using cefoxitin treatment. Journal of Infectious Diseases **214**:122–129.  
218 doi:10.1093/infdis/jiw069.

219 7. **Tomkovich S, Yang Y, Winglee K, Gauthier J, Mühlbauer M, Sun X, Mohamadzadeh  
M, Liu X, Martin P, Wang GP, Oswald E, Fodor AA, Jobin C.** 2017. Locoregional effects  
221 of microbiota in a preclinical model of colon carcinogenesis. Cancer Research **77**:2620–2632.  
222 doi:10.1158/0008-5472.can-16-3472.

223 8. **Sze MA, Schloss PD.** 2018. Leveraging existing 16S rRNA gene surveys to identify reproducible  
224 biomarkers in individuals with colorectal tumors. doi:10.1101/285486.

225 9. **Sanna S, Zuydam NR van, Mahajan A, Kurilshikov A, Vila AV, Võsa U, Mujagic Z, Masclee  
AAM, Jonkers DMAE, Oosting M, Joosten LAB, Netea MG, Franke L, Zhernakova A, Fu J,**

227 **Wijmenga C, McCarthy MI.** 2019. Causal relationships among the gut microbiome, short-chain  
228 fatty acids and metabolic diseases. *Nature Genetics*. doi:10.1038/s41588-019-0350-x.

229 **10. Meisel M, Mayassi T, Fehlner-Peach H, Koval JC, O'Brien SL, Hinterleitner R, Lesko  
230 K, Kim S, Bouziat R, Chen L, Weber CR, Mazmanian SK, Jabri B, Antonopoulos DA.** 2016.  
231 Interleukin-15 promotes intestinal dysbiosis with butyrate deficiency associated with increased  
232 susceptibility to colitis. *The ISME Journal* **11**:15–30. doi:10.1038/ismej.2016.114.

233 **11. O'Keefe SJD.** 2016. Diet, microorganisms and their metabolites and colon cancer. *Nature  
234 Reviews Gastroenterology & Hepatology* **13**:691–706. doi:10.1038/nrgastro.2016.165.

235 **12. Bishehsari F, Engen P, Preite N, Tuncil Y, Naqib A, Shaikh M, Rossi M, Wilber S, Green  
236 S, Hamaker B, Khazaie K, Voigt R, Forsyth C, Keshavarzian A.** 2018. Dietary fiber treatment  
237 corrects the composition of gut microbiota, promotes SCFA production, and suppresses colon  
238 carcinogenesis. *Genes* **9**:102. doi:10.3390/genes9020102.

239 **13. Weaver GA, Krause JA, Miller TL, Wolin MJ.** 1988. Short chain fatty acid distributions of  
240 enema samples from a sigmoidoscopy population: An association of high acetate and low butyrate  
241 ratios with adenomatous polyps and colon cancer. *Gut* **29**:1539–1543. doi:10.1136/gut.29.11.1539.

242 **14. Yao Y, Suo T, Andersson R, Cao Y, Wang C, Lu J, Chui E.** 2017. Dietary fibre for the  
243 prevention of recurrent colorectal adenomas and carcinomas. *Cochrane Database of Systematic  
244 Reviews*. doi:10.1002/14651858.cd003430.pub2.

245 **15. Gianfredi V, Salvatori T, Villarini M, Moretti M, Nucci D, Realdon S.** 2018. Is dietary fibre  
246 truly protective against colon cancer? A systematic review and meta-analysis. *International Journal  
247 of Food Sciences and Nutrition* **69**:904–915. doi:10.1080/09637486.2018.1446917.

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

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

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

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

258 20. **Schloss PD, Westcott SL, Ryabin T, Hall JR, Hartmann M, Hollister EB, Lesniewski RA,**  
259 **Oakley BB, Parks DH, Robinson CJ, Sahl JW, Stres B, Thallinger GG, Horn DJV, Weber CF.**  
260 2009. Introducing mothur: Open-source, platform-independent, community-supported software  
261 for describing and comparing microbial communities. Applied and Environmental Microbiology  
262 75:7537–7541. doi:10.1128/aem.01541-09.

263 21. **Langille MGI, Zaneveld J, Caporaso JG, McDonald D, Knights D, Reyes JA, Clemente JC,**  
264 **Burkepile DE, Thurber RLV, Knight R, Beiko RG, Huttenhower C.** 2013. Predictive functional  
265 profiling of microbial communities using 16S rRNA marker gene sequences. Nature Biotechnology  
266 31:814–821. doi:10.1038/nbt.2676.

267 22. **Li D, Liu C-M, Luo R, Sadakane K, Lam T-W.** 2015. MEGAHIT: An ultra-fast single-node  
268 solution for large and complex metagenomics assembly via succinct de bruijn graph. Bioinformatics  
269 31:1674–1676. doi:10.1093/bioinformatics/btv033.

270 23. **Hyatt D, Chen G-L, LoCascio PF, Land ML, Larimer FW, Hauser LJ.** 2010. Prodigal:  
271 Prokaryotic gene recognition and translation initiation site identification. BMC Bioinformatics 11:119.  
272 doi:10.1186/1471-2105-11-119.

273 24. **Buchfink B, Xie C, Huson DH.** 2014. Fast and sensitive protein alignment using DIAMOND.  
274 Nature Methods 12:59–60. doi:10.1038/nmeth.3176.

275 25. **Steinegger M, Söding J.** 2017. MMseqs2 enables sensitive protein sequence searching for  
276 the analysis of massive data sets. Nature Biotechnology. doi:10.1038/nbt.3988.

277 **Figures**

278 **Figure 1. SCFA concentrations did not vary meaningfully with diagnosis of colonic lesions**

279 **or with treatment for adenomas or carcinomas.** (A) The concentration of fecal SCFAs from  
280 individuals with normal colons (N=172) or those with adenoma (N=198) or carcinomas (N=120). (B)  
281 A subset of individuals diagnosed with adenomas (N=41) or carcinomas (N=26) who underwent  
282 treatment were resampled a year after the initial sampling; one extreme propionate value (124.4  
283 mmol/kg) was included in the adenoma analysis but censored from the visualization for clarity.

284 **Figure 2. SCFA concentrations do not improve models for diagnosing the presence of**

285 **adenomas, carcinomas, or all lesions and cannot be reliably predicted from 16S rRNA**  
286 **gene or metagenomic sequence data.** (A) The median AUROC for diagnosing individuals as  
287 having adenomas or carcinomas using SCFAs was slightly better than chance (depicted by  
288 horizontal line at 0.50), but did not improve performance of the models generated using 16S rRNA  
289 gene sequence data. (B) Regression models that were trained using 16S rRNA gene sequence,  
290 metagenomic, and PICRUSt data to predict the concentrations of SCFAs performed poorly (all  
291 median  $R^2$  values  $< 0.14$ ). Regression models generated using 16S rRNA gene sequence and  
292 PICRUSt data included data from 490 samples and those generated using metagenomic data  
293 included data from 78 samples.

294 **Figure S1. Comparison of training and testing results for classification models shows that**

295 **the models are robust and are not overfit.** random forest classification models were generated to  
296 differentiate between individuals with normal colons and those with adenomas or carcinomas using  
297 16S rRNA gene sequence data that were clustered into genera or OTUs with and without including  
298 the three SCFAs as additional features. random forest classification models were generated by  
299 partitioning the samples into a training set with 80% of the data and a testing set with the remaining  
300 samples for 100 randomizations.

301 **Figure S2. Comparison of training and testing results for regression models shows that**

302 **the models are robust and are not overfit.** random forest regression models were generated  
303 to predict the concentration of each SCFA using each individuals' microbiome data generated

304 using 16S rRNA gene sequence and metagenomic sequence data. These regression models were  
305 generated by partitioning the samples into a training set with 80% of the data and a testing set with  
306 the remaining samples for 100 randomizations.



