

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

## 21 **Importance**

22 Considering colorectal cancer is the third leading cancer-related cause of death within the United  
23 States, there is a great need to detect colorectal tumors early without invasive colonoscopy  
24 procedures and to prevent the formation of tumors. Short-chain fatty acids (SCFAs) are often  
25 used as a surrogate for measuring gut health and for being anti-carcinogenic because of their  
26 anti-inflammatory properties. We evaluated the fecal SCFA concentration of a cohort of individuals  
27 with varying colonic tumor burden who were previously analyzed to identify microbiome-based  
28 biomarkers of tumors. We were unable to find an association between SCFA concentration and  
29 tumor burden or use SCFAs to improve our microbiome-based models of classifying people based  
30 on their tumor status. Furthermore, we were unable to find an association between the fecal  
31 community structure and SCFA concentrations. These data indicate that there is no conclusive link  
32 between the gut microbiome, SCFAs, and tumor burden.

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

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

59 **Fecal SCFA concentrations did not vary with diagnosis or treatment.** To quantify the  
60 associations between colorectal cancer, the microbiome, and SCFAs, we quantified the  
61 concentration of acetate, propionate, and butyrate in feces of previously characterized individuals

62 with normal colons (N=172) and those with colonic adenomas (N=198) or carcinomas (N=120)  
63 (16). We were unable to detect a significant difference in any SCFA concentration across the  
64 diagnoses groups (all  $P > 0.15$ ; Figure 1A). Among the individuals with adenomas and carcinomas,  
65 a subset ( $N_{\text{adenoma}}=41$ ,  $N_{\text{carcinoma}}=26$ ) were treated and sampled a year later (17). None of  
66 the SCFAs exhibited a significant change with treatment (all  $P > 0.058$ ; Figure 1B). For both the  
67 pre-treatment cross-sectional data and the pre/post treatment data, we also failed to detect any  
68 significant differences in the relative concentrations of any SCFAs ( $P > 0.16$ ). Finally, we pooled  
69 the SCFA concentrations on a total and per molecule of carbon basis and again failed to observe  
70 any significant differences ( $P > 0.077$ ). These results demonstrated that there were no significant  
71 associations between fecal SCFA concentration and diagnosis or treatment.

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

87 **Knowledge of microbial community structure does not predict SCFA concentrations.** We  
88 next asked whether the fecal community structure was predictive of fecal SCFA concentrations,  
89 regardless of a person's diagnosis. We trained random forest regression models using 16S rRNA  
90 gene sequence data binned into OTUs and genera to predict the concentration of the SCFAs (Figure

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

107 **Conclusion.** Our data indicate that fecal SCFA concentrations are not associated with the presence  
108 of adenomas or carcinomas and that they provide weak predictive power to improve the ability  
109 to diagnose someone with one of these lesions. Furthermore, knowledge of the taxonomic and  
110 genetic structure of gut microbiota was not predictive of SCFA concentrations. These results  
111 complement existing literature that suggest that fiber consumption and the production of SCFAs  
112 are unable to prevent the risk of developing colonic tumors. It is important to note that our analysis  
113 was based on characterizations of SCFA and microbiome profiles using fecal samples and that  
114 observations along the mucosa near the site of lesions may provide a stronger association. This  
115 may be a cautionary result to temper enthusiasm for SCFAs as a biomarker of gut health more  
116 generally. Going forward it is critical to develop additional hypotheses for how the microbiome and  
117 host interact to drive tumorigenesis so that we can better understand tumorigenesis and identify  
118 biomarkers that will allow early detection of lesions.

## 119 **Acknowledgements**

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

## 125 **Materials and Methods**

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

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

143 **16S rRNA gene sequence data analysis.** Sequence data from Baxter et al. (16) and Sze et  
144 al. (17) were obtained from the Sequence Read Archive (studies SRP062005 and SRP096978)  
145 and reprocessed using mothur v.1.42 (20). The original studies generated sequence data from  
146 V4 region of the 16S rRNA gene using paired 250 nt reads on an Illumina MiSeq sequencer. The  
147 resulting sequence data were assembled into contigs and screened to remove low quality contigs  
148 and chimeras. The curated sequences were then clustered into OTUs at a 97% similarity threshold  
149 and assigned to the closest possible genus with an 80% confidence threshold trained on the  
150 reference collection from the Ribosomal Database Project (v.16). We used PICRUSt (v.2.1.0-b)  
151 with the recommended standard operating protocol to generate imputed metagenomes based on  
152 the expected metabolic pathways and KEGG categories (21).

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

165 **random forest models.** The classification models were built to predict lesion type from microbiome  
166 information with or without SCFA concentrations. The regression models were built to predict the  
167 SCFA concentrations of acetate, butyrate, and propionate from microbiome information. For  
168 classification and regression models, we pre-processed the features by scaling them to vary  
169 between zero and one. Features with no variance in the training set were removed from both the  
170 training and testing sets. We randomly split the data into training and test sets so that the training  
171 set consisted of 80% of the full dataset while the test set was composed of the remaining data. The  
172 training set was used for hyperparameter selection and training the model and the test set was used  
173 for evaluating prediction performance. For each model, the best performing hyperparameter,  $mtry$ ,  
174 was selected in an internal five-fold cross-validation of the training set with 100 randomizations. Six  
175 values of  $mtry$  were tested and the value that provided the largest AUROC or  $R^2$  was selected. We  
176 trained the random forest model using the selected  $mtry$  value and predicted the held-out test set.  
177 The data-split, hyperparameter selection, training and testing steps were repeated 100 times to  
178 get a reliable and robust reading of model prediction performance. We used AUROC and  $R^2$  as  
179 the prediction performance metric for classification and regression models, respectively. We used  
180 randomForest package implemented to the caret package (version 4.6-14) in R statistical software  
181 (version 6.0-81) for our models.

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

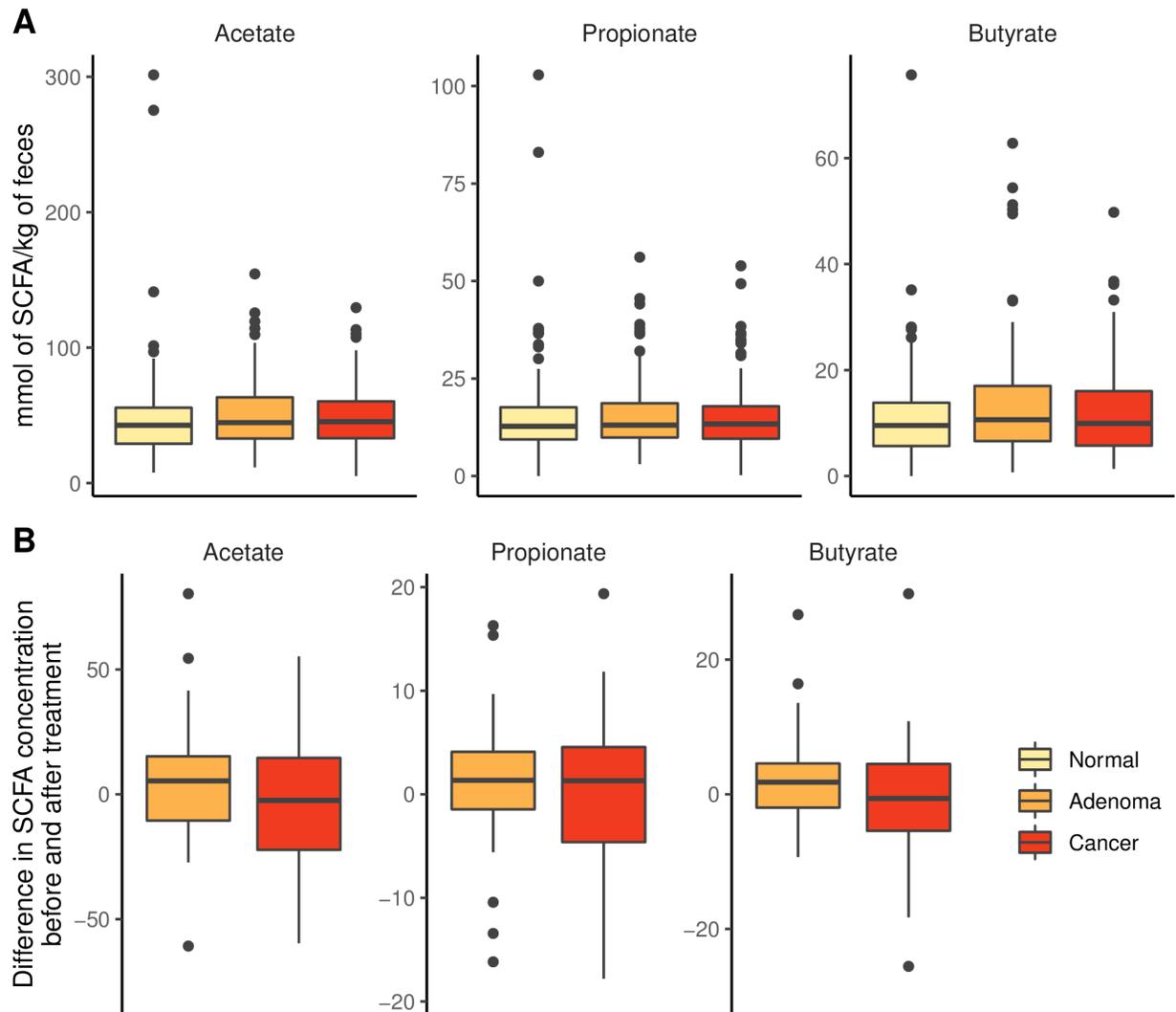
194 **Code availability.** The code for all sequence curation and analysis steps including an Rmarkdown  
195 version of this manuscript is available at [https://github.com/SchlossLab/Sze\\_SCFACRC\\_XXXX\\_](https://github.com/SchlossLab/Sze_SCFACRC_XXXX_2019/)  
196 2019/.

## 197 **References**

- 198 1. **Siegel RL, Miller KD, Jemal A.** 2016. Cancer statistics, 2016. CA: A Cancer Journal for  
199 Clinicians **66**:7–30. doi:10.3322/caac.21332.
- 200 2. **Fearon ER, Vogelstein B.** 1990. A genetic model for colorectal tumorigenesis. Cell **61**:759–767.  
201 doi:10.1016/0092-8674(90)90186-i.
- 202 3. **Fliss-Isakov N, Zelber-Sagi S, Webb M, Halpern Z, Kariv R.** 2017. Smoking habits are  
203 strongly associated with colorectal polyps in a population-based case-control study. Journal of  
204 Clinical Gastroenterology **1**. doi:10.1097/mcg.0000000000000935.
- 205 4. **Lee J, Jeon JY, Meyerhardt JA.** 2015. Diet and lifestyle in survivors of colorectal cancer.  
206 Hematology/Oncology Clinics of North America **29**:1–27. doi:10.1016/j.hoc.2014.09.005.
- 207 5. **Zackular JP, Baxter NT, Iverson KD, Sadler WD, Petrosino JF, Chen GY, Schloss PD.**  
208 2013. The gut microbiome modulates colon tumorigenesis. mBio **4**:e00692–13–e00692–13.  
209 doi:10.1128/mbio.00692-13.
- 210 6. **Shields CED, Meerbeke SWV, Housseau F, Wang H, Huso DL, Casero RA, O’Hagan**  
211 **HM, Sears CL.** 2016. Reduction of murine colon tumorigenesis driven by Enterotoxigenic  
212 Bacteroides fragilis using cefoxitin treatment. Journal of Infectious Diseases **214**:122–129.  
213 doi:10.1093/infdis/jiw069.
- 214 7. **Tomkovich S, Yang Y, Winglee K, Gauthier J, Mühlbauer M, Sun X, Mohamadzadeh**  
215 **M, Liu X, Martin P, Wang GP, Oswald E, Fodor AA, Jobin C.** 2017. Locoregional effects  
216 of microbiota in a preclinical model of colon carcinogenesis. Cancer Research **77**:2620–2632.  
217 doi:10.1158/0008-5472.can-16-3472.
- 218 8. **Sze MA, Schloss PD.** 2018. Leveraging existing 16S rRNA gene surveys to identify reproducible  
219 biomarkers in individuals with colorectal tumors. doi:10.1101/285486.
- 220 9. **Sanna S, Zuydam NR van, Mahajan A, Kurilshikov A, Vila AV, Vösa U, Mujagic Z, Masclee**  
221 **AAM, Jonkers DMAE, Oosting M, Joosten LAB, Netea MG, Franke L, Zhernakova A, Fu J,**

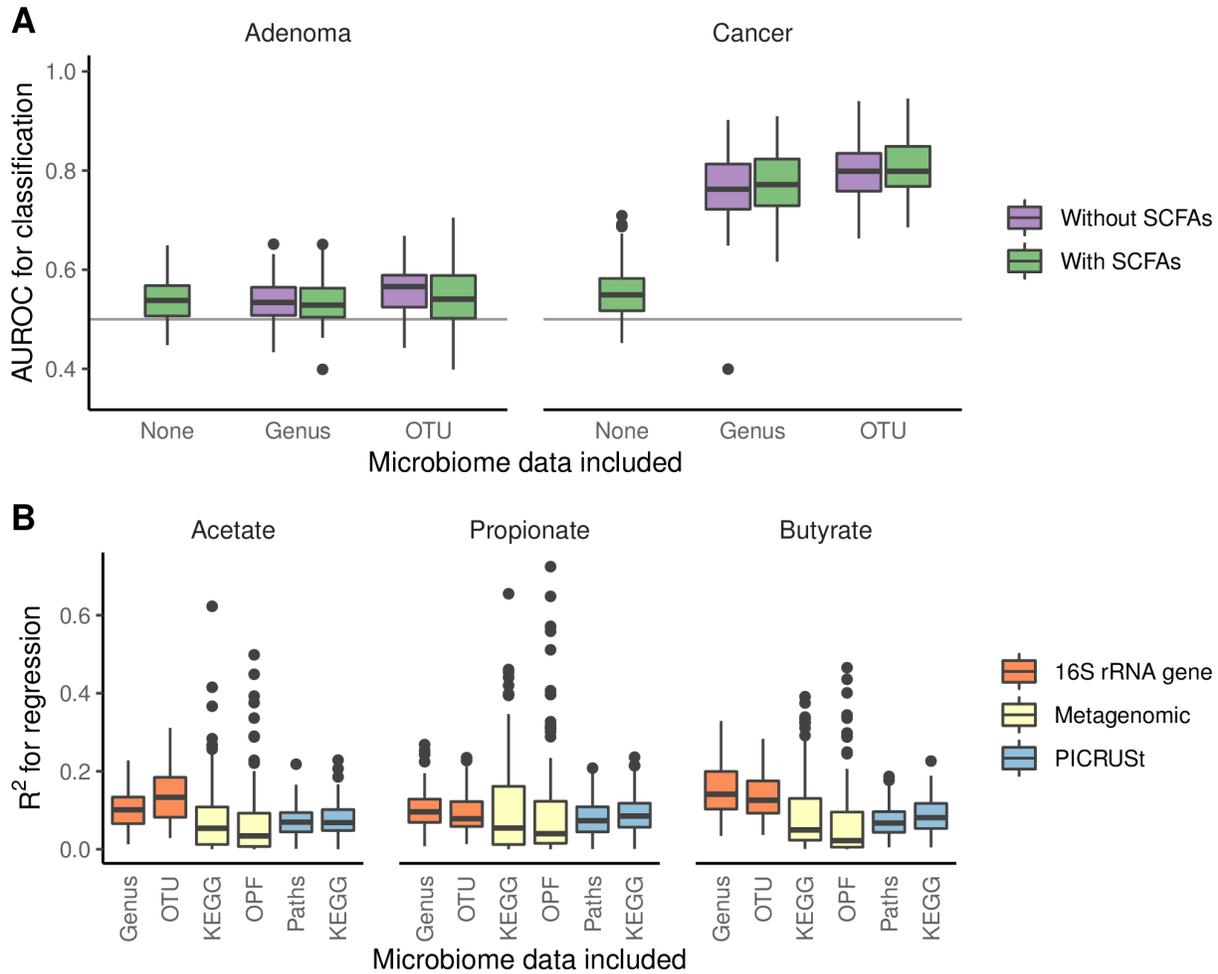
- 222 **Wijmenga C, McCarthy MI.** 2019. Causal relationships among the gut microbiome, short-chain  
223 fatty acids and metabolic diseases. *Nature Genetics*. doi:10.1038/s41588-019-0350-x.
- 224 10. **Meisel M, Mayassi T, Fehlner-Peach H, Koval JC, O'Brien SL, Hinterleitner R, Lesko**  
225 **K, Kim S, Bouziat R, Chen L, Weber CR, Mazmanian SK, Jabri B, Antonopoulos DA.** 2016.  
226 Interleukin-15 promotes intestinal dysbiosis with butyrate deficiency associated with increased  
227 susceptibility to colitis. *The ISME Journal* **11**:15–30. doi:10.1038/ismej.2016.114.
- 228 11. **O'Keefe SJD.** 2016. Diet, microorganisms and their metabolites and colon cancer. *Nature*  
229 *Reviews Gastroenterology & Hepatology* **13**:691–706. doi:10.1038/nrgastro.2016.165.
- 230 12. **Bishehsari F, Engen P, Preite N, Tuncil Y, Naqib A, Shaikh M, Rossi M, Wilber S, Green**  
231 **S, Hamaker B, Khazaie K, Voigt R, Forsyth C, Keshavarzian A.** 2018. Dietary fiber treatment  
232 corrects the composition of gut microbiota, promotes SCFA production, and suppresses colon  
233 carcinogenesis. *Genes* **9**:102. doi:10.3390/genes9020102.
- 234 13. **Weaver GA, Krause JA, Miller TL, Wolin MJ.** 1988. Short chain fatty acid distributions of  
235 enema samples from a sigmoidoscopy population: An association of high acetate and low butyrate  
236 ratios with adenomatous polyps and colon cancer. *Gut* **29**:1539–1543. doi:10.1136/gut.29.11.1539.
- 237 14. **Yao Y, Suo T, Andersson R, Cao Y, Wang C, Lu J, Chui E.** 2017. Dietary fibre for the  
238 prevention of recurrent colorectal adenomas and carcinomas. *Cochrane Database of Systematic*  
239 *Reviews*. doi:10.1002/14651858.cd003430.pub2.
- 240 15. **Gianfredi V, Salvatori T, Villarini M, Moretti M, Nucci D, Realdon S.** 2018. Is dietary fibre  
241 truly protective against colon cancer? A systematic review and meta-analysis. *International Journal*  
242 *of Food Sciences and Nutrition* **69**:904–915. doi:10.1080/09637486.2018.1446917.
- 243 16. **Baxter NT, Ruffin MT, Rogers MAM, Schloss PD.** 2016. Microbiota-based model improves  
244 the sensitivity of fecal immunochemical test for detecting colonic lesions. *Genome Medicine* **8**.  
245 doi:10.1186/s13073-016-0290-3.
- 246 17. **Sze MA, Baxter NT, Ruffin MT, Rogers MAM, Schloss PD.** 2017. Normalization of the  
247 microbiota in patients after treatment for colonic lesions. *Microbiome* **5**. doi:10.1186/s40168-017-0366-3.

- 248 **18. Hannigan GD, Duhaime MB, Ruffin MT, Koumpouras CC, Schloss PD.** 2017. Diagnostic  
249 potential & the interactive dynamics of the colorectal cancer virome. doi:10.1101/152868.
- 250 **19. Venkataraman A, Sieber JR, Schmidt AW, Waldron C, Theis KR, Schmidt TM.** 2016.  
251 Variable responses of human microbiomes to dietary supplementation with resistant starch.  
252 *Microbiome* **4**. doi:10.1186/s40168-016-0178-x.
- 253 **20. Schloss PD, Westcott SL, Ryabin T, Hall JR, Hartmann M, Hollister EB, Lesniewski RA,**  
254 **Oakley BB, Parks DH, Robinson CJ, Sahl JW, Stres B, Thallinger GG, Horn DJV, Weber CF.**  
255 2009. Introducing mothur: Open-source, platform-independent, community-supported software  
256 for describing and comparing microbial communities. *Applied and Environmental Microbiology*  
257 **75**:7537–7541. doi:10.1128/aem.01541-09.
- 258 **21. Langille MGI, Zaneveld J, Caporaso JG, McDonald D, Knights D, Reyes JA, Clemente JC,**  
259 **Burkepile DE, Thurber RLV, Knight R, Beiko RG, Huttenhower C.** 2013. Predictive functional  
260 profiling of microbial communities using 16S rRNA marker gene sequences. *Nature Biotechnology*  
261 **31**:814–821. doi:10.1038/nbt.2676.
- 262 **22. Li D, Liu C-M, Luo R, Sadakane K, Lam T-W.** 2015. MEGAHIT: An ultra-fast single-node  
263 solution for large and complex metagenomics assembly via succinct de bruijn graph. *Bioinformatics*  
264 **31**:1674–1676. doi:10.1093/bioinformatics/btv033.
- 265 **23. Hyatt D, Chen G-L, LoCascio PF, Land ML, Larimer FW, Hauser LJ.** 2010. Prodigal:  
266 Prokaryotic gene recognition and translation initiation site identification. *BMC Bioinformatics* **11**:119.  
267 doi:10.1186/1471-2105-11-119.
- 268 **24. Buchfink B, Xie C, Huson DH.** 2014. Fast and sensitive protein alignment using DIAMOND.  
269 *Nature Methods* **12**:59–60. doi:10.1038/nmeth.3176.
- 270 **25. Steinegger M, Söding J.** 2017. MMseqs2 enables sensitive protein sequence searching for  
271 the analysis of massive data sets. *Nature Biotechnology*. doi:10.1038/nbt.3988.



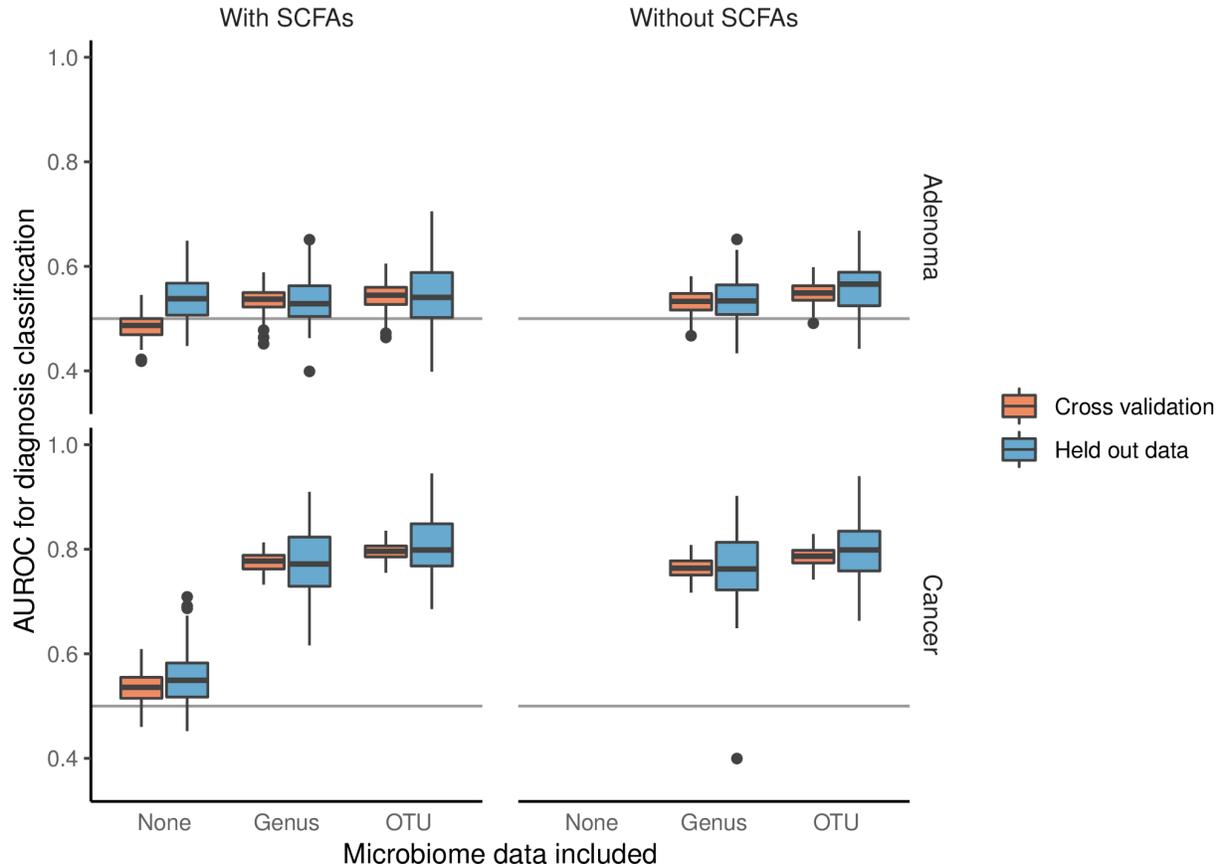
272

273 **Figure 1. SCFA concentrations did not vary meaningfully with diagnosis of colonic lesions**  
274 **or with treatment for adenomas or carcinomas.** (A) The concentration of fecal SCFAs from  
275 individuals with normal colons (N=172) or those with adenoma (N=198) or carcinomas (N=120). (B)  
276 A subset of individuals diagnosed with adenomas (N=41) or carcinomas (N=26) who underwent  
277 treatment were resampled a year after the initial sampling; one extreme propionate value (124.4  
278 mmol/kg) was included in the adenoma analysis but censored from the visualization for clarity.



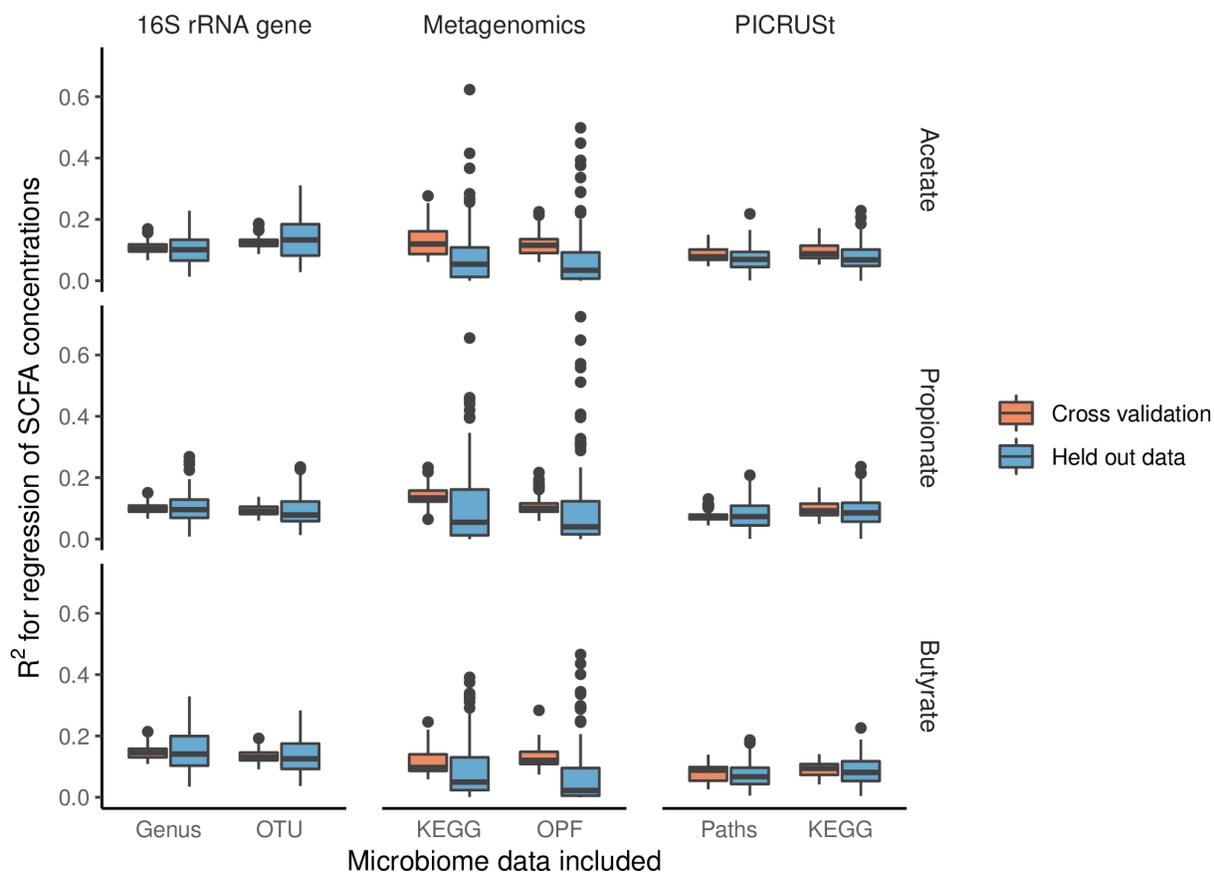
279

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



290

291 **Figure S1. Comparison of training and testing results for classification models shows that**  
292 **the models are robust and are not overfit.** random forest classification models were generated to  
293 differentiate between individuals with normal colons and those with adenomas or carcinomas using  
294 16S rRNA gene sequence data that were clustered into genera or OTUs with and without including  
295 the three SCFAs as additional features. random forest classification models were generated by  
296 partitioning the samples into a training set with 80% of the data and a testing set with the remaining  
297 samples for 100 randomizations.



298

299 **Figure S2. Comparison of training and testing results for regression models shows that**  
300 **the models are robust and are not overfit.** random forest regression models were generated  
301 to predict the concentration of each SCFA using each individuals' microbiome data generated  
302 using 16S rRNA gene sequence and metagenomic sequence data. These regression models were  
303 generated by partitioning the samples into a training set with 80% of the data and a testing set with  
304 the remaining samples for 100 randomizations.