# Revisiting Short-Chain Fatty Acids and the Microbiota in Colorectal Cancer

Running title: SCFAs and Colorectal Cancer

Marc A Sze[1], Nicholas A Lesniak[1], Mack T Ruffin IV[2], 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 48109

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

# 1 Abstract

## 2 Introduction

# 3 Results

# 4 Discussion

## 5 Conclusions

## Materials and Methods

**Study design and sampling.** The overall protocol has been described in detail previously (1, 2). In brief, this study used fecal samples obtained at either a single cross-sectional time point (n=490) or from before and after treatment for their tumor (n=67). For patients undergoing treatment for their tumor the length of time between their initial and follow up sample ranged from 188 - 546 days. Our use of treatment has been previously defined as encompassing removal of a tumor with or without chemotherapy and radiation (1). Diagnosis of tumor was made by colonoscopic examination and histopathological review of biopsies obtained (1, 2). The University of Michigan Institutional Review Board approved the study and informed consent was obtained from all participants in accordance to the guidelines set out by the Helsinki Decleration.

**Measuring specific SCFAs.** Our protocol for the measurement of acetate, butyrate, and propionate followed a previously published protocol (3). The following changes to this protocol included the use of frozen fecal samples suspended in 1ml of PBS instead of fecal suspensions in DNA Genotek OmniGut tubes, and the use of the acutal weight of fecal samples instead of the average weight for SCFA concentration normalizations. These changes did not affect the overall median concentrations of these SCFAs between the two studies (see Table 1 (3) and Figure 1 in this report).

**16s rRNA gene sequencing.** The workflow and processing have been described previously (1, 4, 5). The major differences from these previous reports include: the use of version 1.39.5 of the mothur software package and clustering Operational Taxonomic Units (OTUs) at 97% similarity used the OptClust algorithm (6).

**Generating imputed metagenomes.** The use of PICRUSt version 1.1.2 with the recommended standard operating protocol (7) was used. Briefly, the mothur shared file and metadata was converted into a biom formated table using the biom convert function, the subsequent biom file was processed with the normalize_by_copy_number.py function, and subsequent imputed metagenomes created using the predict_metagenomes.py function.

**Obtaining OPFs from metagenomes.** A subset of the cross-sectional group (n=490) containing a total of 85 individuals (normal n=29 normal, adenoma n=28, and carcinoma n=28) was shotgun

7

sequenced on an Illumina HiSeq with 125 bp paired end reads using a previously described method (8). Briefly, the sequences were quality filtered and sequences aligning to the human genome were removed prior to contig assembly with MEGAHIT (9). Open Reading Frames (ORFs) were identified using Prodigal (10), counts generated using Diamond (11), subsequent clustering into Operational Protein Families (OPFs) used mmseq2 (12), and OPF gene alignment used the KEGG database (13).

**Pulling genes involved with SCFA synthesis.** Specific genes located near the end of the pathways involved in the synthesis of acetate, butyrate, and propionate were analyzed for any differences between individuals with normal colons and those with tumors. These genes were based on pathways from KEGG and can be found in the supplemental material [Table S1].

**Statistical analysis workflow.** All analysis was performed using the statistical language R (14). Generally, differences between the different disease groups used a Kruskal-Walis rank sum test with a Dunn's post-hoc test. Models to classify high or low SCFA concentration based on 16S rRNA gene sequencing data were created using Random Forest (15). Regression models to classify the exact SCFA concentration based on 16S rRNA gene sequencing data also were built using the Random Forest algorithm. The measured SCFA concentrations were first tested for differences between groups. The ability of 16S rRNA gene sequencing to classify these concentrations were then assessed. Next, the imputed gene counts of important mediators of SCFA creation were tested. Finally, the counts generated for OPFs that matched important genes involved with SCFA creation were analyzed.

8

## Acknowledgements

## References

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

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

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

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

5. **Kozich JJ**, **Westcott SL**, **Baxter NT**, **Highlander SK**, **Schloss PD**. 2013. Development of a dual-index sequencing strategy and curation pipeline for analyzing amplicon sequence data on the MiSeq illumina sequencing platform. Applied and Environmental Microbiology **79**:5112–5120. doi:10.1128/aem.01043-13.

6. **Westcott SL**, **Schloss PD**. 2017. OptiClust, an improved method for assigning amplicon-based sequence data to operational taxonomic units. mSphere **2**:e00073–17. doi:10.1128/mspheredirect.00073-17.

7. **Langille MGI**, **Zaneveld J**, **Caporaso JG**, **McDonald D**, **Knights D**, **Reyes JA**, **Clemente JC**, **Burkepile DE**, **Thurber RLV**, **Knight R**, **Beiko RG**, **Huttenhower C**. 2013. Predictive functional profiling of microbial communities using 16S rRNA marker gene sequences. Nature Biotechnology

84 **31**:814–821. doi:10.1038/nbt.2676.

85 8. **Hannigan GD**, **Duhaime MB**, **Ruffin MT**, **Koumpouras CC**, **Schloss PD**. 2017. Diagnostic

86 potential & the interactive dynamics of the colorectal cancer virome. doi:10.1101/152868.

87 9. **Li D**, **Liu C-M**, **Luo R**, **Sadakane K**, **Lam T-W**. 2015. MEGAHIT: An ultra-fast single-node

88 solution for large and complex metagenomics assembly via succinct de bruijn graph. Bioinformatics

89 **31**:1674–1676. doi:10.1093/bioinformatics/btv033.

90 10. **Hyatt D**, **Chen G-L**, **LoCascio PF**, **Land ML**, **Larimer FW**, **Hauser LJ**. 2010. Prodigal:

91 Prokaryotic gene recognition and translation initiation site identification. BMC Bioinformatics **11**:119.

92 doi:10.1186/1471-2105-11-119.

93 11. **Buchfink B**, **Xie C**, **Huson DH**. 2014. Fast and sensitive protein alignment using DIAMOND.

94 Nature Methods **12**:59–60. doi:10.1038/nmeth.3176.

95 12. **Steinegger M**, **Söding J**. 2017. MMseqs2 enables sensitive protein sequence searching for

96 the analysis of massive data sets. Nature Biotechnology. doi:10.1038/nbt.3988.

97 13. **Kanehisa M**, **Sato Y**, **Kawashima M**, **Furumichi M**, **Tanabe M**. 2015. KEGG as a

98 reference resource for gene and protein annotation. Nucleic Acids Research **44**:D457–D462.

99 doi:10.1093/nar/gkv1070.

100 14. **R Core Team**. 2017. R: A language and environment for statistical computing. R Foundation

101 for Statistical Computing, Vienna, Austria.

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

103  Insert figure legends with the first sentence in bold, for example:

104  **Figure 1.  Number of OTUs sampled among bacterial and archaeal 16S rRNA gene**

105  **sequences for different OTU definitions and level of sequencing effort.** Rarefaction curves

106  for different OTU definitions of Bacteria (A) and Archaea (B). Rarefaction curves for the coarse

107  environments in Table 1 for Bacteria (C) and Archaea (D).