

# Differences in the fecal Microbiome Before and After Colorectal Cancer Treatment

Running Title: Human Microbiome and Colorectal Cancer

Marc A Sze<sup>1</sup>, Nielson T Baxter<sup>2</sup>, Mack T Ruffin IV<sup>3</sup>, Mary AM Rogers<sup>2</sup>, and Patrick D Schloss<sup>1†</sup>

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

<sup>1</sup> Department of Microbiology and Immunology, University of Michigan, Ann Arbor, MI

<sup>2</sup> Department of Internal Medicine, University of Michigan, Ann Arbor, MI

<sup>3</sup> Department of Family Medicine and Community Medicine, Penn State Hershey Medical Center, Hershey, PA

## Abstract

**Background:** Colorectal cancer (CRC) continues to be a worldwide health problem with previous research suggesting that a link may exist between the fecal bacterial microbiome and CRC. The overall objective of our study was to test the hypothesis that changes in the bacterial microbiome occur after surgery in patients with lesions (i.e. adenoma or carcinoma). Specifically, we wanted to identify what within the community was different before and after surgical removal of said lesion.

**Results:** The bacterial microbiome in pre and post surgery samples of individuals with adenoma were more similar to each other than those with carcinoma (P-value < 0.05). There was no difference in the relative abundance of any OTU between the pre and post surgery samples (P-value > 0.05). A model with a total of 53 variables was able to classify lesion (AUC = 0.811 - 0.866) while a model built to classify pre surgery samples had 70 variables (AUC = 0.641 - 0.805). The post surgery sample for both models had a decrease in the positive probability for either lesion or pre surgery sample (P-value < 0.05). In total there were 23 OTUs that were common to both models and the majority of these classified to commensal bacteria (e.g. *Bacteroides*, *Clostridiales*, *Blautia*, and *Ruminococcaceae*).

**Conclusions:** Our data supports the hypothesis that there are differences in the bacterial microbiome between pre and post surgical samples. With individuals with carcinoma having more drastic differences to the overall community than those with adenoma. Changes to commensal bacteria were some of the most important variables for model classification, suggesting that these bacteria may be central to initial polyp formation and transition to carcinoma.

23 **Keywords**

24 bacterial microbiome; colorectal cancer; polyps; FIT; post surgery; risk factors

## Background

Colorectal cancer (CRC) continues to be a leading cause of cancer related deaths and is currently the third most common cause of cancer deaths [1,2]. Over the last few years death due to the disease has seen a significant decrease, thanks mainly to improvements in screening [1]. However, despite this improvement there are still approximately 50,000 deaths from the disease a year [2].

Recently, there has been promising work on the bacterial microbiome and it's ability to be able to complement existing screening methods such as Fecal Immunoglobulin Test (FIT) or act alone as a screening tool [3,4]. There has also been research into how this microbiome could be altered directly on tumor tissue itself [5]. A few studies have also shown how this microbiome [6] or specific members within it [7] could be directly involved with the pathogenesis of CRC. These studies have helped to provide a tantalizing link between the bacterial microbiome and CRC. Although these studies suggest that the bacterial microbiome might change after treatment there remains limited information on the bacterial microbiome before and after surgical removal of lesion (adenoma or carcinoma) and whether the community changes at all.

In this study we tested the hypothesis that the bacterial microbiome changes after surgery for individuals with a lesion. Our analysis included both alpha and beta diversity analysis along with investigation of individual operational taxonomic units (OTUs). We also utilized Random Forest models and observed how these models as well as specific OTUs within them performed pre (initial) and post (follow up) surgery. We also used these models to look for similar important OTUs to identify the crucial OTUs for not only classifying initial and follow up samples but also lesion or normal.

## Results

**Bacterial Community and FIT:** We first wanted to test whether there were any broad differences between initial and follow up samples based on lesion being either adenoma or carcinoma. What we found was that the bacterial community in those with carcinoma were more dissimilar (as measured by thetacyc) to their initial sample than those with adenoma (P-value < 0.001) [Figure 1a]. We also found that there were larger changes in fecal blood (measured by FIT) for those with carcinoma versus adenoma (P-value < 0.0001) [Figure 1b]. The broad shift in bacterial community structure before and after surgery was visualized using NMDS for both adenoma [Figure 1c] (PERMANOVA > 0.05) and carcinoma [Figure 1d] (PERMANOVA < 0.05). Interestingly, when initial and follow up samples were compared to each other, regardless of whether they were adenoma or carcinoma (lesion), there was no significant overall difference between them (PERMANOVA > 0.05). When investigating more broad alpha diversity metrics there was no difference found between initial and follow up samples for lesion, adenoma only, or carcinoma only for any metric tested [Table S1]. We also observed that there was no difference in OTU relative abundance between initial and follow up samples for lesion, adenoma only, or carcinoma only [Figure S1].

**Cancer Associated Bacteria:** Previous literature has suggested that a number of oral microbes may be important in CRC pathogenesis [8]. So we next examined whether there were differences in previously well described carcinoma associated OTUs. These included the OTUs that aligned with *Porphyromonas asaccharolytica* (Otu000202), *Fusobacterium nucleatum* (Otu000442), *Parvimonas micra* (Otu001273), and *Peptostreptococcus stomatis* (Otu001682). There was a difference in relative abundance in initial and follow up samples for lesion and carcinoma for *Parvimonas micra* (P-value < 0.05), and *Porphyromonas asaccharolytica* (P-value < 0.05). In contrast, there was no difference in relative abundance in any of these OTUs for those with adenoma [Figure 2]. We also observed that only a

small percentage of those with adenoma or carcinoma were positive or had an appreciable relative abundance of any of these respective OTUs [Figure 2].

**Full and Reduced Model:** We next wanted to identify if there were any common bacterial microbiome changes in individuals with adenoma or carcinoma. In order to investigate this we created two different models: one to classify lesion versus normal and one to classify pre (initial) versus post (follow up) samples based on the bacterial community and FIT measurements. The lesion model had an AUC range of 0.73 to 0.797 while the initial sample model had an AUC range of 0.485 to 0.686 after 100 iterations of 20 repeated 10-fold cross validations. By identifying the most important variables for each respective model and then reducing them to only these factors we were able to increase the AUC in the lesion model (0.811 - 0.866) and initial sample model (0.641 - 0.805).

The test set AUC range for the full and reduced lesion model were similar to that reported for the training set AUC ranges and the ROC curve ranges overlap with each other [Figure 3a]. The ROC curve for the final lesion model used was within the range of the reduced lesion model [Figure 3a]. Interestingly, the test set AUC range for the initial sample model performed much better than the training set AUCs. Both the full and reduced initial sample models overlapped with each other [Figure 3b] but there was a marked decrease in the ROC curve for the final before sample model used.

**Common OTUs to both Models:** The reduced models were built based on the most important variables to either classification model. For the lesion model there were a total of 53 variables [Figure S2] whereas for the initial sample model there were a total of 70 variables [Figure S3]. For both models FIT measurement resulted in one of the largest decreases in MDA [Figure S2a & S3a]. When we compared the two different reduced models with each other there were a total of 23 common OTUs. Some of the most common taxonomic identifications belonged to *Bacteroides*, *Clostridiales*, *Blautia*, and *Ruminococcaceae*. The vast majority of these OTUs had classifications to bacteria

typically thought of as commensal [Table S2].

**Positive Probability after Lesion Removal:** If there were common OTUs for individuals with adenoma and carcinoma, that were different versus normal controls, we would expect to find a decrease in the positive probability of the follow up sample to be either lesion or an initial sample. This is what we observed regardless of model used (lesion or initial sample) or whether it was built on the full or reduced variable data set [Figure 4 & S4] (P-value < 0.001).

When we separated individuals based on whether they had an adenoma or carcinoma there was only a decrease in positive probability for the carcinoma group (P-value < 0.001) and not for the adenoma group (P-value > 0.05). We also observed that there were no significant differences in whether the models classified the samples as having lesion between the predicted and actual (P-value > 0.05). This lack of difference between the predicted and actual classifications were also observed for the initial sample model (P-value > 0.05). Even though the lesion model was not as accurate in classifying samples as the initial sample model. It was able to correctly keep the one individual who still had a carcinoma on follow up above the cut off threshold for a positive call [Figure 4a & S4a] while the initial sample model did not [Figure 4b & S4b].

**Treatment and Time Differences:** After observing these changes in the bacterial community and positive probability we wanted to assess whether additional treatments, such as chemotherapy and radiation, could have an impact on the results that we observed. There was no difference in the amount of change in positive probability for either the full or reduced lesion model for either chemotherapy (P-value > 0.05) or radiation therapy (P-value > 0.05). In contrast, we observed for the the before sample model a significant difference in decreased positive probability for those treated with chemotherapy (P-value < 0.05). All other variables that were tested showed no difference based on whether chemotherapy or radiation was used [Table S3]. Finally, we wanted to know if the length of

126 time between the initial and follow up sample could be a possible confounder. Within our  
127 study there was a significant difference for the time elapsed in the collection of the follow  
128 up sample between adenoma and carcinoma (uncorrected P-value < 0.05), with time  
129 passed being less for adenoma (253 +/- 41.3 days) then carcinoma (351 +/- 102 days).



## Discussion

From our results there were some large observed differences in the bacterial microbiome between pre and post surgery samples based on whether the individual had an adenoma or carcinoma. There were much larger differences between initial and follow up samples based on the thetayc distance metric and in fecal blood as measured by FIT for individuals with carcinoma versus adenoma [Figure 1]. However, there were no differences between initial and follow up samples for Shannon Diversity, observed OTUs, or evenness regardless of whether the individual had an adenoma or carcinoma [Table S1]. There was also no differences in relative abundance of any specific OTU for lesion, adenoma only, or carcinoma only [Figure s1].

Although there were no differences when investigating all OTUs, when looking specifically at four OTUs that taxonomically classified to previously suggested cancer causing microbes we found that only 2/4 had a decrease in relative abundance between initial and follow up for those with carcinoma and 0/4 had differences for those with adenoma. This data would suggest that these specific OTUs may be important in the transition of an adenoma to a carcinoma but less so in the initiation of an adenoma from benign tissue.

We next created a model that incorporated FIT and the bacterial microbiome to either be able to classify lesions (adenoma or carcinoma) or initial samples in order to find common OTUs in the community that change for both adenoma and carcinoma. What we found was that the commonly associated CRC bacteria were not highly represented within our models but rather that OTUs that made up the most important variables overwhelmingly belonged to commensal bacteria. With only the lesion model having a single OTU from a previously associated cancer bacterium (*Porphyromonas asaccharolytica*). Using only these important OTUs and FIT both models (lesion and before sample) significantly decreased positive probability of either lesion or being an initial sample on follow up [Figure

4 & S4]. Further confirmation of the importance of the changes of commensal bacteria to these classifications was that a total of 23 OTUs were common to both models and the vast majority belonged to regular residents of our gut community.

For the majority of tests performed there were no differences in the bacterial microbiome based on whether chemotherapy or radiation was received [Table S3]. There was a difference in the length of time between initial and follow up sample between adenoma and carcinoma. These results would indicate that the findings described were specific to the surgical intervention and that some of the differences observed between carcinoma and adenoma samples could be due to differences in collection time between samples for the two different groups.

This study builds upon previous work from numerous labs that have looked into the bacterial microbiome as a potential screening tool [3,4] by exploring what happens to the bacterial community after surgical removal of a lesion. Based on previous work by Arthur, et al. [9] it may not be surprising to have E.coli as one of the most important OTUs and one that was common to both models. Interestingly, many of the most important OTUs had taxonomic identification for resident gut microbes. This could suggest that the bacterial community is one of the first components that could change during the pathogenesis of disease. These bacterial microbiome changes could be the first step in allowing more inflammatory bacterium to gain a foothold within the colon [8].

Curiously, we observed that the typical CRC associated bacteria were not predictive within our models. There are a number of reasons why this may have occurred. First, is that they were not present in enough individuals to be able to classify those with and without disease with a high degree of accuracy. Second, is that it is possible that our Random Forest models were able to gather the same information from measures such as FIT or other OTUs. It is also possible that both of these explanations could have played a role. Regardless, our observations would suggest that an individual's resident bacteria have a large role

181 to play in disease initiation and could change in a way that allows predictive models to  
182 lower the positive probability of a lesion after surgery [Figure 4]. It should be noted that  
183 our study does not argue against the importance of these CRC associated bacteria in the  
184 pathogenesis of disease but rather that the models do not utilize these specific bacteria  
185 for classification purposes (lesion or before sample). In fact, it is possible that these CRC  
186 associated bacteria are important in the transition from adenoma to carcinoma and would  
187 be one explanation as to why in our data we not only see high initial relative abundances, in  
188 certain individuals, but also large decreases in relative abundance in those with carcinoma  
189 but not in those with adenoma after surgery [Figure 2].

190 Many of the common OTUs between the different models used had many OTUs that  
191 taxonomically classified to potential butyrate producers [Table S2]. Another batch of OTUs  
192 classified to bacteria that can either degrade polyphenols or are inhibited by them. Both  
193 butyrate and polyphenols are thought to protective against cancer in part by reducing  
194 inflammation [10]. These protective compounds are derived from the breakdown of  
195 fiber, fruits, and vegetables by resident gut microbes. One example of this potential  
196 diet-microbiome-inflammation-polyp axis is that *Bacteroides*, which was highly prevalent  
197 in our models, are known to be increased in those with high non-meat based protein  
198 consumption [11]. High protein consumption in general has been linked with an increased  
199 CRC risk [12]. Conversely, *Bacteroides* are inhibited by polyphenols which are derived from  
200 fruits and vegetables [13]. Our data fits with the hypothesis that the microbial metabolites  
201 from breakdown products within our own diet could not only help to shape the exisiting  
202 community but also have an effect on CRC risk and disease progression.

203 One limitation of our study is that we do not know whether individuals who were still  
204 classified as positive by the lesion model eventually had a subsequent CRC diagnosis.  
205 This information would help to strengthen the case for our lesion model to have kept  
206 a number of individuals above the cutoff threshold even though at follow up they were

207 diagnosed as no longer having a lesion. Another limitation is that we do not know if  
208 adding modern tests such as the stool DNA test [14] could help improve our overall AUC.  
209 Another limitation is that this study drew heavily from those with Caucasian ancestry.  
210 The results may not be immediately representative of those with either Asian or African  
211 ancestry. Finally, although our training and test set are relatively large we still run the  
212 risk of over-fitting or having a model that may not be immediately extrapolate-able to  
213 other populations. We've done our best to safeguard against this by not only running  
214 10-fold cross validation but also having over 100 different 80/20 splits to try and mimic  
215 the type of variation that might be expected to occur. The time difference in collection of  
216 sample between adenoma and carcinoma could have affected our observed results for  
217 the differences between individuals with adenoma or carcinoma. This confounding though  
218 does not affect the observations based on the overall lesion results.

219 Another interesting outcome was that within figure 3 the before sample model showed  
220 better test AUC results than the training set AUC. This may have occurred because the  
221 training AUC that was determined from 20 repeated 10 fold cross validation removed  
222 samples at random and did not take into account that they were matched samples. Another  
223 potential reason is that the model itself may be over-fit since the total number of samples  
224 was not that large. However, the lesion model did not suffer from these discrepancies and  
225 similar conclusions can be drawn solely from this model. Regardless, further independent  
226 studies will need to be carried out to verify our findings since not only are we dealing with  
227 feces, which could be very different than the communities present on the actual tissue, but  
228 also are dealing with correlations that may not be representative of the true pathogenesis  
229 of disease.

230 Despite these limitations we think that these findings significantly add to the existing  
231 scientific knowledge on CRC and the bacterial microbiome: That there is a measurable  
232 difference in the bacterial community after surgical removal of lesion. Further, the ability

233 for machine learning algorithms to take bacterial microbiome data and successfully lower  
234 positive probability after either adenoma or carcinoma removal provides evidence that  
235 there are specific signatures, mostly attributable to commensal organisms, associated with  
236 these lesions. Our data provides evidence that commensal bacteria are important in the  
237 development of polyps and also potentially the transition from adenoma to carcinoma.

## Methods

**Study Design and Patient Sampling** The sampling and design of the study was similar to that reported in Baxter, et al [3]. In brief, study exclusion involved those who had already undergone surgery, radiation, or chemotherapy, had colorectal cancer before a baseline fecal sample could be obtained, had IBD, a known hereditary non-polyposis colorectal cancer, or Familial adenomatous polyposis. Samples used to build the models for prediction were collected either prior to a colonoscopy or between 1 - 2 weeks after. The bacterial microbiome has been shown to normalize within this time period [15]. Our follow up data set had a total of 67 individuals that not only had a sample as described but also a follow up sample between 188 - 546 days after surgery and treatment had been completed. This study was approved by the University of Michigan Institutional Review Board. All study participants provided informed consent and the study itself conformed to the guidelines set out by the Helsinki Declaration.

**FIT and 16S rRNA Gene Sequencing** FIT was analyzed as previously published using both OC FIT-CHEK and OC-Auto Micro 80 automated system (Polymedco Inc.) [16]. 16S rRNA gene sequencing was completed as previously described by Kozich, et al. [17]. In brief, DNA extraction used the 96 well Soil DNA isolation kit (MO BIO Laboratories) and an epMotion 5075 automated pipetting system (Eppendorf). The V4 variable region was amplified and the resulting product was split between three sequencing runs with normal, adenoma, and carcinoma evenly represented on each run. Each group was randomly assigned to avoid biases based on sample collection location.

**Sequence Processing** The mothur software package (v1.37.5) was used to process the 16S rRNA gene sequences. This process has been previously described [17]. The general processing workflow using mothur was as follows: Paired-end reads were first merged into contigs, quality filtered, aligned to the SILVA database, screened for chimeras,

classified with a naive Bayesian classifier using the Ribosomal Database Project (RDP), and clustered into Operational Taxonomic Units (OTUs) using a 97% similarity cutoff with an average neighbor clustering algorithm. The number of sequences for each sample was rarefied to 10523 in an attempt to minimize uneven sampling.

**Lesion Model Creation** The Random Forest [18] algorithm was used to create the model used for prediction of lesion (adenoma or carcinoma) with the main testing and training of the model completed on a data set of 490 individuals. This model was then applied to our follow up data set of 67 individuals. The model included data on FIT and the bacterial microbiome. Non-binary data was checked for near zero variance and OTUs that had near zero variance were removed. This pre-processing was performed with the R package caret (v6.0.73). Optimization of the mtry hyper-parameter involved taking the samples and making 100 different 80/20 (train/test) splits of the data where normal and lesion were represented in the same proportion within both the whole data set and the 80/20 split. Each of these splits were then run through 20 repeated 10-fold cross validations to optimize the mtry hyper-parameter by maximizing the AUC (Area Under the Curve of the Receiver Operator Characteristic). This resulting model was then tested on the 20% of the data that was originally held out from this overall process. Once the ideal mtry was found the entire 490 sample set was used to create the final Random Forest model on which classifications on the 67-person cohort was completed. The default cutoff of 0.5 was used as the threshold to classify individuals as positive or negative for lesion. The hyper-parameter, mtry, defines the number of variables to investigate at each split before a new division of the data is created with the Random Forest model.

**Before Sample Model Creation** We also investigated whether a model could be created that could identify before and after surgery samples. The main difference was that only the 67-person cohort was used at all stages of model building and classification. Other than this difference the creation of this model and optimization of the mtry hyper-parameter was

completed using the same procedure that was used to create the lesion model. Instead of classifying samples as positive or negative of lesion this model classified samples as positive or negative for being a before surgery sample.

**Selection of Important OTUs** In order to assess which variables were most important to all the models we counted the number of times a variable was present in the top 10% of mean decrease in accuracy (MDA) for each of the 100 different 80/20 split models and then filtered this list to variables that were only present more than 50% of the time. This final collated list of variables was what was considered the most important for the lesion or before sample models.

**Statistical Analysis** The R software package (v3.3.2) was used for all statistical analysis. Comparisons between bacterial community structure utilized PERMANOVA [19] in the vegan package (v2.4.1). Comparisons between probabilities as well as overall OTU differences between initial and follow up samples utilized a paired Wilcoxon ranked sum test. Where multiple comparison testing was needed a Benjamini-Hochberg (BH) correction was applied [20] and a corrected P-value of less than 0.05 was considered significant. Unless otherwise stated the P-values reported are those that were BH corrected.

**Analysis Overview** Initial and follow up samples were analyzed for differences in alpha and beta diversity. Next, differences in FIT between initial and follow ups for either adenoma or carcinoma were investigated. From here, all OTUs that were used in either model were then analyzed using a paired Wilcoxon test. We then investigated the relative abundance of specific previously associated CRC bacteria, specifically, OTUs that taxonomically classified to *Fusobacterium nucleatum*, *Parvimonas micra*, *Peptostreptococcus assacharolytica*, and *Porphyromonas stomatis*. We wanted to test if there were any differences based on whether the individual had an adenoma or carcinoma. From here the lesion model was then tested for accuracy in prediction and whether it reduced the positive probability of lesion after surgery. The most important OTUs for this



were used to build a reduced model and it was assessed for similarity to the original model. We then used the before sample model to assess whether it could classify samples better than the lesion model. The most important OTUs were then identified from this model and used to create a reduced feature before sample model. This reduced feature model, as was done with the lesion model, was compared to the full model for loss of accuracy. Finally, a list of common OTUs were found for the two different models used.

***Reproducible Methods.*** A detailed and reproducible description of how the data were processed and analyzed can be found at [https://github.com/SchlossLab/Size\\_followUps\\_2017](https://github.com/SchlossLab/Size_followUps_2017). Raw sequences have been deposited into the NCBI Sequence Read Archive (SRP062005 and SRP096978) and the necessary metadata can be found at <https://www.ncbi.nlm.nih.gov/Traces/study/> and searching the respective SRA study accession.

**Figure 1: General Differences between the Adenoma or Carcinoma Group.** A) A significant difference was found between the adenoma and carcinoma group for thetacyc (P-value = 0.000472). B) A significant difference was found between the adenoma and carcinoma group for change in FIT (P-value = 2.15e-05). C) NMDS of the initial and follow up samples for the Adenoma group. D) NMDS of the initial and follow up samples for the Carcinoma group. For C) and D) the teal represents initial samples and the pink represents follow up samples.

**Figure 2: Previously Associated CRC Bacteria in Initial and Follow up Samples.**

A) Carcinoma initial and follow up samples. There was a significant difference in initial and follow up sample for the OTUs classified as *Peptostreptococcus stomatis* (P-value = 0.0496) and *Porphyromonas asaccharolytica* (P-value = 0.00842). B) Adenoma initial and follow up samples. There were no significant differences between initial and follow up.

**Figure 3: Graph of the Receiver Operating Characteristic Curve for lesion and Before Sample Models.** The shaded areas represents the range of values of a 100 different 80/20 splits of the test set data using either all variables (grey) or reduced variable (red) models. The blue line represents the reduced variable model using 100% of the data set. A) Lesion model. B) Before sample model

**Figure 4: Breakdown by Carcinoma and Adenoma of Prediction Results for Lesion and Before Sample Reduced Variable Models** A) Lesion positive probability adjustment of those with carcinoma from initial to follow up sample B) Initial follow up positive probability adjustment of those with carcinoma from initial to follow up sample C) Lesion positive probability adjustment of those with adenoma as well as those with SRN and the probability adjustment from initial to follow up sample. D) Initial follow up positive probability adjustment of those with adenoma as well as those with SRN and the probability adjustment from initial to follow up sample. The dotted line represents the threshold used to make the decision of whether a sample was positive or not.

**Figure S1: Distribution of P-values from Paired Wilcoxon Analysis of OTUs in Initial versus Follow Up**

**Figure S2: Summary of Important Variables in the Lesion Model** A) MDA of the most important variables in the lesion model. The black point represents the median and the different colors are the different runs up to 100. B) The total number of appearances of each variable in the 100 different lesion models. The cutoff of 50% was used to assess importance.

**Figure S3: Summary of Important Variables in Before Sample Model** A) MDA of the most important variables in the lesion model. The black point represents the median and the different colors are the different runs up to 100. B) The total number of appearances of each variable in the 100 different lesion models. The cutoff of 50% was used to assess importance.

**Figure S4: Breakdown by Carcinoma and Adenoma of Prediction Results for Lesion and Before Sample Full Variable Models** A) Lesion positive probability adjustment of those with carcinoma from initial to follow up sample B) Initial follow up positive probability adjustment of those with carcinoma from initial to follow up sample C) Lesion positive probability adjustment of those with adenoma as well as those with SRN and the probability adjustment from initial to follow up sample. D) Initial follow up positive probability adjustment of those with adenoma as well as those with SRN and the probability adjustment from initial to follow up sample. The dotted line represents the threshold used to make the decision of whether a sample was positive or not.

**Figure S5: Thetayc Graphed Against Time of Follow up Sample from Initial**

## **Declarations**

### **Ethics approval and consent to participate**

### **Consent for publication**

### **Availability of data and material**

### **Competing Interests**

All authors declare that they do not have any relevant competing interests to report.

### **Funding**

This study was supported by funding from the National Institutes of Health to P. Schloss (R01GM099514, P30DK034933) and to the Early Detection Research Network (U01CA86400).

### **Authors' contributions**

All authors were involved in the conception and design of the study. MAS analyzed the data. NTB processed samples and analyzed the data. All authors interpreted the data. MAS and PDS wrote the manuscript. All authors reviewed and revised the manuscript. All authors read and approved the final manuscript.

## 389 **Acknowledgements**

390 The authors thank the Great Lakes-New England Early Detection Research Network for  
391 providing the fecal samples that were used in this study.

## References

1. Jemal A, Siegel R, Xu J, Ward E. Cancer statistics, 2010. *CA: a cancer journal for clinicians*. 2010;60:277–300.
2. Haggard FA, Boushey RP. Colorectal cancer epidemiology: Incidence, mortality, survival, and risk factors. *Clinics in Colon and Rectal Surgery*. 2009;22:191–7.
3. Baxter NT, Ruffin MT, Rogers MAM, Schloss PD. Microbiota-based model improves the sensitivity of fecal immunochemical test for detecting colonic lesions. *Genome Medicine*. 2016;8:37.
4. Zeller G, Tap J, Voigt AY, Sunagawa S, Kultima JR, Costea PI, et al. Potential of fecal microbiota for early-stage detection of colorectal cancer. *Molecular Systems Biology*. 2014;10:766.
5. Dejea CM, Wick EC, Hechenbleikner EM, White JR, Mark Welch JL, Rossetti BJ, et al. Microbiota organization is a distinct feature of proximal colorectal cancers. *Proceedings of the National Academy of Sciences of the United States of America*. 2014;111:18321–6.
6. Zackular JP, Baxter NT, Chen GY, Schloss PD. Manipulation of the Gut Microbiota Reveals Role in Colon Tumorigenesis. *mSphere*. 2016;1.
7. Arthur JC, Gharaibeh RZ, Mühlbauer M, Perez-Chanona E, Uronis JM, McCafferty J, et al. Microbial genomic analysis reveals the essential role of inflammation in bacteria-induced colorectal cancer. *Nature Communications*. 2014;5:4724.
8. Flynn KJ, Baxter NT, Schloss PD. Metabolic and Community Synergy of Oral Bacteria in Colorectal Cancer. *mSphere*. 2016;1.
9. Arthur JC, Perez-Chanona E, Mühlbauer M, Tomkovich S, Uronis JM, Fan T-J, et al.

Intestinal inflammation targets cancer-inducing activity of the microbiota. *Science* (New York, N.Y.). 2012;338:120–3.

10. Louis P, Hold GL, Flint HJ. The gut microbiota, bacterial metabolites and colorectal cancer. *Nature Reviews Microbiology* [Internet]. 2014 [cited 2017 Feb 14];12:661–72. Available from: <http://www.nature.com/doifinder/10.1038/nrmicro3344>

11. Zhu Y, Lin X, Li H, Li Y, Shi X, Zhao F, et al. Intake of Meat Proteins Substantially Increased the Relative Abundance of Genus *Lactobacillus* in Rat Feces. *PloS One*. 2016;11:e0152678.

12. Mu C, Yang Y, Luo Z, Guan L, Zhu W. The Colonic Microbiome and Epithelial Transcriptome Are Altered in Rats Fed a High-Protein Diet Compared with a Normal-Protein Diet. *The Journal of Nutrition*. 2016;146:474–83.

13. Ozdal T, Sela DA, Xiao J, Boyacioglu D, Chen F, Capanoglu E. The Reciprocal Interactions between Polyphenols and Gut Microbiota and Effects on Bioaccessibility. *Nutrients* [Internet]. 2016 [cited 2017 Feb 14];8:78. Available from: <http://www.mdpi.com/2072-6643/8/2/78>

14. Cotter TG, Burger KN, Devens ME, Simonson JA, Lowrie KL, Heigh RI, et al. Long-Term Follow-up of Patients Having False Positive Multi-target Stool DNA Tests after Negative Screening Colonoscopy: The LONG-HAUL Cohort Study. *Cancer Epidemiology, Biomarkers & Prevention: A Publication of the American Association for Cancer Research, Cosponsored by the American Society of Preventive Oncology*. 2016;

15. O'Brien CL, Allison GE, Grimpen F, Pavli P. Impact of colonoscopy bowel preparation on intestinal microbiota. *PloS One*. 2013;8:e62815.

16. Daly JM, Bay CP, Levy BT. Evaluation of fecal immunochemical tests for colorectal

- 437 cancer screening. *Journal of Primary Care & Community Health*. 2013;4:245–50.
- 438 17. Kozich JJ, Westcott SL, Baxter NT, Highlander SK, Schloss PD. Development of a  
439 dual-index sequencing strategy and curation pipeline for analyzing amplicon sequence  
440 data on the MiSeq Illumina sequencing platform. *Applied and Environmental Microbiology*.  
441 2013;79:5112–20.
- 442 18. Breiman L. Random Forests. *Machine Learning* [Internet]. 2001 [cited 2013 Feb  
443 7];45:5–32. Available from: <http://link.springer.com/article/10.1023/A%3A1010933404324>  
444 <http://link.springer.com/article/10.1023%2FA%3A1010933404324?LI=true>
- 445 19. Anderson MJ, Walsh DCI. PERMANOVA, ANOSIM, and the Mantel test in the face of  
446 heterogeneous dispersions: What null hypothesis are you testing? *Ecological Monographs*  
447 [Internet]. 2013 [cited 2017 Jan 5];83:557–74. Available from: [http://doi.wiley.com/10.1890/](http://doi.wiley.com/10.1890/12-2010.1)  
448 12-2010.1
- 449 20. Benjamini Y, Hochberg Y. Controlling the false discovery rate: A practical and  
450 powerful approach to multiple testing. *Journal of the Royal Statistical Society. Series B*  
451 (Methodological). 1995;57:289–300.