

Making Sense of the Noise: Leveraging Existing 16S rRNA Gene Surveys to Identify Key Community Members in Colorectal Tumors

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

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

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

Co-author e-mails:

- marcsze@med.umich.edu

Abstract

Background. An increasing body of literature suggests that both individual and collections of bacteria are associated with the progression of colorectal cancer. As the number of studies investigating these associations increases and the number of subjects in each study increases, a meta-analysis to identify the associations that are the most predictive of disease progression is warranted. For our meta-analysis, we analyzed previously published 16S rRNA gene sequencing data collected from feces (1737 individuals from 8 studies) and colon tissue (492 total samples from 350 individuals from 7 studies).

Results. We quantified the odds ratios for individual bacterial genera that were associated with an individual having tumors relative to a normal colon. Among the stool samples, there were no genera that had a significant odds ratio associated with adenoma and there were 8 genera with significant odds ratios associated with carcinoma. Similarly, among the tissue samples, there were no genera that had a significant odds ratio associated with adenoma and there were 3 genera with significant odds ratios associated with carcinoma. Among the significant odds ratios, the association between individual taxa and tumor diagnosis was equal or below 7.11. Because individual taxa had limited association with tumor diagnosis, we trained Random Forest classification models using the genera with the five highest and lowest odds ratios, using the entire collection of genera found in each study, and using operational taxonomic units defined based on a 97% similarity threshold. All training approaches yielded similar classification success as measured using the Area Under the Curve. The ability to correctly classify individuals with adenomas was poor and the ability to classify individuals with carcinomas was considerably better using sequences from stool or tissue.

Conclusions. This meta-analysis confirms previous results indicating that individuals with adenomas cannot be readily classified based on their bacterial community, but that those

26 with carcinomas can. Regardless of the dataset, we found a subset of the fecal community
27 that was associated with carcinomas was as predictive as the full community.

28 **Keywords**

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

Background

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

Recently, there has been an intense focus on identifying microbiota-based biomarker yielding a seemingly endless number of candidate taxa. Some studies point towards mouth-associated genera such as *Fusobacterium*, *Peptostreptococcus*, *Parvimonas*, and *Porphyromonas* that are enriched in people with carcinomas [6,11–17]. Other studies have identified members of *Akkermansia*, *Bacteroides*, *Enterococcus*, *Escherichia*, *Klebsiella*, *Mogibacterium*, *Streptococcus*, and *Providencia* are also associated with carcinomas [13–15]. Additionally, *Roseburia* has been found in some studies to be more abundant in people with tumors but in other studies it has been found to be either less abundant or no different than what is found in subjects with normal colons [14,17,19,20]. There are strong results from tissue culture and murine models that *Fusobacterium nucleatum*, pks-positive strains of *Escherichia coli*, *Streptococcus gallolyticus*, and an enterotoxin-producing strain of *Bacteroides fragilis* are important in the pathogenesis of CRC [5,14,21–24]. These results point to a causative role for the microbiota in CRC pathogenesis as well as their potential as diagnostic biomarkers.

Most studies have focused on identifying biomarkers in patients with carcinomas but there is a greater clinical need to identify biomarkers associated with adenomas. Studies focusing on broad scale community metrics have found that measures such as the total number of Operational Taxonomic Units (OTUs) are decreased in those with adenomas versus controls [25]. Other studies have identified *Acidovorax*, *Bilophila*, *Cloacibacterium*, *Desulfovibrio*, *Helicobacter*, *Lactobacillus*, *Lactococcus*, *Mogibacterium*, and *Pseudomonas* to be enriched in those with adenomas [25–27]. There are few genera that are enriched in patients with adenoma or carcinoma tumors.

Confirming some of these previous findings, a recent meta-analysis found that 16S rRNA gene sequences from members of the *Akkermansia*, *Fusobacterium*, and *Parvimonas* were fecal biomarkers for the presence of carcinomas [28]. Contrary to previous studies they found sequences similar to members of *Lactobacillus* and *Ruminococcus* to be enriched in patients with adenoma or carcinoma relative to those with normal colons [12,15,16]. In addition, they found 16S rRNA gene sequences from members of *Haemophilus*, *Methanospaera*, *Prevotella*, *Succinivibrio* were enriched in patients with adenoma and *Pantoea* were enriched in patients with carcinomas. Although this meta-analysis was helpful for distilling a large number of possible biomarkers, the aggregate number of samples included in the analysis (n = 509) was smaller than several larger case-control studies that have since been published [12,27]

Here we provide an updated meta-analysis using 16S rRNA gene sequence data from both feces (n = 1737) and colon tissue (492 samples from 350 individuals) from 14 studies [11–17,19,20,23,25–27,29] [Table 1 & 2]. We expand both the breadth and scope of the previous meta-analysis to investigate whether biomarkers describing the bacterial community or specific members of the community can more accurately classify patients as having adenoma or carcinoma. Our results suggest that the bacterial community changes as disease severity worsens and that a subset of the microbial community can be

81 used to diagnose the presence of carcinoma.

Results

Lower Bacterial Diversity is Associated with Increased OR of Tumors: To assess differences in broad scale community metrics as disease severity worsens Operational Taxonomic Unit (OTU) richness, evenness, and Shannon diversity measurements were power transformed and Z-score normalized. These metrics are commonly used to assess the total number of OTUs, the equality of their abundance, and the overall diversity, respectively. Using linear mixed-effect models to control for study and variable region we assessed whether OTU richness, evenness, or Shannon diversity changed in a step-wise manner with disease severity. In stool, there was a significant decrease in both evenness and Shannon diversity as disease severity moved from control to adenoma to carcinoma (P-value = 0.025 and 0.043, respectively). We next tested whether the detectable differences in community significantly increased in OR of having an adenoma or carcinoma. For fecal samples, a decrease versus the overall median in evenness resulted in a significantly increased OR for carcinoma (OR = 1.66 (1.2 - 2.3), P-value = 0.0021) and adenoma (OR = 1.3 (1.02 - 1.65), P-value = 0.035) while a decrease versus the overall median in Shannon diversity only increased the OR for carcinoma (OR = 1.61 (1.14 - 2.28), P-value = 0.0069) [Figure 1]. Using the Bray-Curtis distance metric and PERMANOVA, it was also possible to identify significant bacterial community changes, in specific studies, for both carcinoma-associated and adenoma-associated microbiota versus control [Table S1].

Using similar transformations for tissue samples, linear mixed-effect models were used on the transformed combined data to control for study, re-sampling of the same individual, and 16S variable region to test whether OTU richness, evenness, or Shannon diversity changed in a step-wise manner as disease severity increased. For colon tissue, there were no significant changes in OTU richness, evenness, or Shannon diversity as disease severity progressed from control to adenoma to carcinoma (P-value > 0.05). We next

analyzed the OR, for matched (unaffected tissue and an adenoma or carcinoma from the same individual) and unmatched (control and adenoma or carcinoma tissue not from the same individual) colon tissue samples. For individuals at either an adenoma or carcinoma stage of disease there was no significant change in OR based on lower than median values for OTU richness, evenness, and Shannon diversity [Figure S1 & Table S2]. Similar to stool samples, significant differences in bacterial community, assessed by PERMANOVA, were identified in unmatched tissue samples, for those at either adenoma or carcinoma stage of CRC [Table S1]. For studies with matched samples no differences in bacterial community were observed when assessed with PERMANOVA [Table S1]. These tissue results suggest that the microbiota within an individual are similar to each other regardless of disease status.

Mouth-Associated Genera are Associated with an Increased OR of Tumor: Next, we asked if being higher than the median relative abundance, for any specific genera, resulted in an altered OR for adenoma or carcinoma, in stool and colon tissue, due to our previous observations of small increases in OR using OTU richness and Shannon diversity. To investigate this we analyzed all common genera across each study, in colon tissue or stool, and assessed whether a relative abundance higher than the median results in an increase or decrease in OR. For both tissue and stool samples only ORs associated with an increase or decrease in carcinoma tumors were significant after multiple comparison correction [Table S3 & S4]. Out of the 8 taxa that had significant ORs in stool samples 4 were mouth-associated microbes. These mouth-associated genera significantly increased the ORs of carcinoma for stool samples and included *Fusobacterium*, *Parvimonas*, *Porphyromonas*, and *Peptostreptococcus* [Table S3]. Conversely, mouth-associated genera were not significantly associated with an increased OR for carcinoma in tissue samples [Table S4]. Only unmatched tissue samples had significant OR taxa and these were *Dorea* (OR = 0.35 (0.22 - 0.55)), *Blautia* (OR = 0.47 (0.3 - 0.73)), and *Weissella* (OR = 5.15 (2.02 - 13.14)). Overall, there was little direct overlap in increased or decreased OR

taxa between both tumor type and sample site.

Select Community Models can Recapitulate Whole Community Models: Since specific genera increased the OR for carcinoma over diversity metrics we assessed whether the bacterial community was better at classifying disease versus only a select group of genera. We selected these genera based on significance after multiple comparison correction. If no taxa were significant after multiple comparison correction then no model for that specific grouping (i.e. adenoma stool) was analyzed. We first tested three model AUCs. These models were created using Random Forest and where either all genera, all OTUs, or significant OR taxa only. Next, the all genera models and any significant OR taxa models were tested across all studies that were not used to train the model. For stool, all models used had similar AUCs [Figure 2]. Although for adenoma and unmatched carcinoma this trend held, there were large differences in matched tissue based on whether all genera or OTUs were utilized [Figure S2]. When analyzing the tests sets that were comprised of genera data from other studies, both the all genera and significant OR taxa only models had a similar ability to detect individuals with carcinomas, for both stool and tissue samples [Figure 3-S3].

In stool, the most common genera in the top 10 most important variables, in the full community models using all genera-based models, were *Ruminococcus*, *Bacteroides*, and *Roseburia* [Figure 4]. Regardless of sample type, mouth-associated genera were present in models for carcinomas [Figure 4B & D]. Yet, none were present in the majority of studies and none were present in the adenoma models [Figure 4A & B]. For the full community OTU-based models, *Ruminococcaceae* was present in the top 10 consistently for both adenoma and carcinoma models while *Roseburia* was only present in many adenoma models and *Bacteroides* was present in the overwhelming majority of the carcinoma models [Figure 4C & 4D].

Unlike the stool-based Random Forest models, the tissue-based models, for the full

genera from the first approach, showed no consistent representation of *Ruminococcaceae*, *Ruminococcus*, *Bacteroides*, and *Roseburia* in the top 10 most important model variables across study [Figure S4]. The vast majority of the top 10 model variables for the full community genera- and OTU-based models using colon tissue tended to be study specific. Further, there was very little overlap in the top 10 important variables between adenoma and carcinoma stage models, regardless of whether colon tissue or stool was used [Figure S4]. This discordance between stool and colon tissue samples also applies to the mouth-associated genera with one noticeable skew being that *Fusobacterium* and *Fusobacteriaceae* occur more often in the top 10 of matched versus unmatched colon tissue Random Forest models [Figure S4B-C & S4E-F]. This suggests that either the colon tissue microbiota is study or person dependent, that kit and/or other types of contamination associated with low biomass samples may be skewing the results, or that multiple microbes could act as the inflammatory stimulus needed to exacerbate mutations.

CRC Studies are Underpowered for Detecting Small Effect Sizes: Next, we assessed how much confidence should be placed in the reported outcomes from each individual study by calculating the ability to detect a difference (power) and sample size needed for small, medium, and large effect size differences between cases and controls. When assessing the power of each study at different effect sizes the majority of studies achieved 80% power to detect a 30% or greater difference between groups [Figure 5A & B]. No study that we analyzed had the standard 80% power to detect an effect size difference equal to or below 10% [Figure 5A & B]. In order to achieve a power of 80%, for small effect sizes, studies used in our meta-analysis would need to recruit over 1000 individuals for both the case and control arms [Figure 5C]

Discussion

Targeting the identification of tumor microbial biomarkers within stool seems logical since it offers an easy and cost-effective way to stratify risk of disease. The current gold standard for diagnosis, a colonoscopy, can be time-consuming and is not without risk of complications. Although stool represents an easy and less invasive way to assess risk, it is not clear how well this sample reflects adenoma- and carcinoma- associated microbial communities. Some studies have tried to assess this in health and disease but are limited by their sample size [17,30]. Sampling the microbiota directly associated with colon tissue may provide clearer answers but is not without their own limitations. After the colonoscopy bowel prep the bacterial community sampled may reflect the better adhered microbiota versus the resident community. Additionally, these samples contain more host DNA, potentially limiting the types of analysis that can be done. It is well known that low biomass samples can be very difficult to work with and results can be study dependent due to the randomness of contamination [31].

Our study identifies clear but small differences in diversity at the community level and larger differences for individual genera, present in patients with tumors versus controls [Figure 1-3]. Although there was a step-wise decrease in diversity as disease progressed from control to adenoma to carcinoma, this did not translate into large effect size increases in OR for either adenoma or carcinoma tumors. Even though mouth-associated genera increased individual's OR of having a carcinoma for certain sample types, they did not consistently increase the OR of having an adenoma. By using these taxa that had significant ORs after multiple comparison correction we found that we could classify individuals with either adenoma or carcinoma as well as models that use either all genera or all OTUs. Finally, many studies were individually under powered to be able to reject the null hypothesis and this could one reason only the comparison between control and carcinoma individuals for stool samples had reliable detectable differences.

The data presented herein support the importance of specific taxa for carcinoma, but not necessarily adenoma, tumor formation. The results that we have presented show that the significant OR taxa model and both the full genera and OTU models, for individuals with carcinoma, had similar AUCs [Figure 2 & 3]. This suggests that an interplay between a select number of potentially protective and exacerbating microbes within the GI community could be crucial for carcinoma formation. Importantly, it suggests that there may be key members of the GI community that should be studied further to potentially help reduce the risk of carcinoma tumor formation. Conversely, using the present data, it is clear that new approaches may be needed to identify members of the community associated with adenoma tumors. Regardless of sample type and whether a full genera- or OTU-based model was used, our Random Forest models consistently performed poorly. Yet, the step-wise decrease in diversity suggests that the adenoma-associated community is not normal but has changed subtly. This change in diversity, at this early stage of disease, could be focal to the adenoma itself. How the host interacts with these subtle changes at early stages of the disease could be what leads to a thoroughly dysfunctional community that is supportive of tumorigenesis.

For the full genera- and OTU-based models within stool, common GI microbes were most consistently present in the top 10 genera or OTUs across studies [Figure 4]. Changes in *Bacteroides*, *Ruminococcaceae*, *Ruminococcus*, and *Roseburia* were consistently found to be in the top 10 most important variables across the different studies for both individuals with adenoma and carcinoma [Figure 4]. These data suggest that whether the non-resident bacterium is *Fusobacteria* or *Peptostreptococcus* may not be as important as how these bacteria interact with the changing resident community. Based on these observations, it is possible to hypothesize that small changes in community structure lead to new niches in which any one of the mouth-associated or general inflammatory genera can gain a foothold, exacerbating the initial changes in community and facilitating the transition from adenoma to carcinoma stage of disease.

The colon tissue-based studies did not provide a clearer understanding of how the microbiota may be associated with tumors. Generally, the full OTU-based models of unmatched and matched colon tissue samples were concordant with stool samples showing that GI resident microbes were the most prevalent in the top 10 most important variables across study [Figure S4E & F]. Unlike in stool, *Fusobacterium* was the only mouth-associated bacteria consistently present in the top 10 most important variables of the full carcinoma stage models [Figure S4B-C & E-F]. The majority of the colon tissue-based results seem to be study specific with many of the top 10 taxa being present only in a single study. Additionally, the presence of genera associated with contamination, within the top 10 most important variables for the genera and OTU models is worrying. The low bacterial biomass of tissue samples coupled with potential contamination could explain why these results seem to be more sporadic than the stool results.

One important caveat to this study is that even though genera associated with certain species such as *Bacteroides fragilis* and *Streptococcus gallolyticus* subsp. *gallolyticus* were not identified, it does not necessarily mean that these specific species are not important in human CRC [22,24]. Since we are limited in our aggregation of the data to the genus level, it is not possible to clearly delineate which species are contributing to overall disease progression. Our observations are not inconsistent with the previous literature on either *Bacteroides fragilis* or *Streptococcus gallolyticus* subsp. *gallolyticus*. As an example, the stool-based full community models consistently identified the genus *Bacteroides*, as well as OTUs that classified as *Bacteroides*, to be important model components across studies. This suggests that even though *Bacteroides* may not increase the OR of individuals having an adenoma or carcinoma and may not vary in relative abundance, like *Fusobacterium*, it is still important in CRC. Additionally, *Streptococcus gallolyticus* subsp. *gallolyticus* is a mouth-associated microbe, and the results from this study suggest that regardless of sample type, mouth-associated genera are commonly associated with an increased OR for individuals to have a carcinoma tumor.

The associations between the microbiota and individuals with adenoma tumors are inconclusive, in part, because many studies may not be powered effectively to observe small effect sizes. None of the studies analyzed were properly powered to detect a 10% or lower change between cases and controls. The results within our meta-analysis suggest that a small effect size may well be the scope in which differences consistently occur between controls and those with adenomas. Future studies investigating adenoma tumors and the microbiota need to take power into consideration to reproducibly study whether the microbiota contributes to adenoma formation. In contrast to adenoma stage of disease, our observations suggest that most studies analyzed have sufficient power to detect many changes in the carcinoma-associated microbiota because of large effect size differences between cases and controls [Figure 5].

Conclusion

By aggregating together a large collection of studies analyzing both fecal and colon tissue samples, we are able to provide evidence supporting the importance of the bacterial community in colorectal tumors. The data presented here suggests that mouth-associated microbes can gain a foothold within the colon and are commonly associated with the greatest OR of individuals having a carcinoma. Conversely, no conclusive signal with these mouth-associated microbes could be detected for individuals with an adenoma. Our observations also highlight the importance of power and sample number considerations when investigating the microbiota and adenoma tumors due to possible subtle changes in the community. Overall, the associations between the microbiota and individuals with carcinomas were much stronger than with those with adenomas.

Methods

Obtaining Data Sets: The studies used for this meta-analysis were identified through the review articles written by Keku, *et al.* and Vogtmann, *et al.* [32,33] and additional studies not mentioned in the reviews were obtained based on the authors' knowledge of the literature. Studies that used tissue or feces as their sample source for 454 or Illumina 16S rRNA gene sequencing analysis and had data sets with sequences available for analysis were included. Some studies were excluded because they did not have publicly available sequences or did not have metadata in which the authors were able to share. After these filtering steps, the following studies remained: Ahn, *et al.* [11], Baxter, *et al.* [12], Brim, *et al.* [29], Burns, *et al.* [15], Chen, *et al.* [13], Dejea, *et al.* [20], Flemer, *et al.* [17], Geng, *et al.* [19], Hale, *et al.* [27], Kostic, *et al.* [34], Lu, *et al.* [26], Sanapareddy, *et al.* [25], Wang, *et al.* [14], Weir, *et al.* [23], and Zeller, *et al.* [16]. The Zackular [35] study was not included because the 90 individuals analyzed within the study are contained within the larger Baxter study [12]. After sequence processing, all the case samples for the Kostic study had 100 or less sequences remaining and was excluded, leaving a total of 14 studies that analysis could be completed on.

Data Set Breakdown: In total, there were seven studies with only fecal samples (Ahn, Baxter, Brim, Hale, Wang, Weir, and Zeller), five studies with only tissue samples (Burns, Dejea, Geng, Lu, Sanapareddy), and two studies with both fecal and tissue samples (Chen and Flemer). The total number of individuals analyzed after sequence processing for feces was 1737 [Table 1]. The total number of matched and unmatched tissue samples that were analyzed after sequence processing was 492 [Table 2].

Sequence Processing: For the majority of studies, raw sequences were downloaded from the Sequence Read Archive (SRA) (<ftp://ftp-trace.ncbi.nih.gov/sra/sra-instant/reads/ByStudy/sra/SRP/>) and metadata were obtained by searching the respective accession

number of the study at the following website: <http://www.ncbi.nlm.nih.gov/Traces/study/>. Of the studies that did not have sequences and metadata on the SRA, data was obtained from DBGap (n = 1, [11]) and directly from the authors (n = 4, [17,23,25,27]). Each study was processed using the mothur (v1.39.3) software program [36] and quality filtering utilized the default methods for both 454 and Illumina based sequencing. If it was not possible to use the defaults, the stated quality cut-offs, from the study itself, were used instead. Sequences that were made up of an artificial combination of two or more different sequences and commonly known as chimeras were identified and removed using VSEARCH [37] before *de novo* OTU clustering at 97% similarity was completed using the OptiClust algorithm [38].

Statistical Analysis: All statistical analysis after sequence processing utilized the R (v3.4.3) software package [39]. For OTU richness, evenness, and Shannon diversity analysis, values were power transformed using the rcompanion (v1.11.1) package [40] and then Z-score normalized using the car (v2.1.6) package [41]. Testing for OTU richness, evenness, and Shannon diversity differences utilized linear mixed-effect models created using the lme4 (v1.1.15) package [42] to correct for study, repeat sampling of individuals (tissue only), and 16S hyper-variable region used. Odds ratios (OR) were analyzed using both the epiR (v0.9.93) and metafor (v2.0.0) packages [43,44] by assessing how many individuals with and without disease were above and below the overall median value within each specific study. OR significance testing utilized the chi-squared test. Diversity differences measured by the Bray-Curtis index utilized the creation of distance matrix and testing with PERMANOVA executed with the vegan (v2.4.5) package [45]. Random Forest models were built using both the caret (v6.0.78) and randomForest (v4.6.12) packages [46,47]. Power analysis and estimations were made using the pwr (v1.2.1) and statmod (v1.4.30) packages [48,49]. All figures were created using both ggplot2 (v2.2.1) and gridExtra (v2.3) packages [50,51].

Study Analysis Overview: OTU richness, evenness, and Shannon diversity were first assessed for differences between controls, adenoma tumors, and carcinoma tumors using both linear mixed-effect models and ORs. For each individual study the Bray-Curtis index was used to assess differences between control-adenoma and control-carcinoma individuals. Next, all common genera were assessed for differences in ORs for individuals having an adenoma or carcinoma and corrected for multiple comparisons using the Benjamini-Hochberg method [52]. We then built Random Forest models based on all genera, all OTUs, or significant OR taxa (if any were present after multiple comparison correction). For both the full genera and significant OR taxa, models were trained on one study then tested on the remaining studies using genera-based relative abundances. The OTU-based models were built using OTU level data and a 10-fold CV over 100 different iterations, based on random 80/20 splitting of the data, was used to generate a range of expected AUCs. This process was repeated for every study in the meta-analysis. Comparisons of the initial trained model AUCs for the full genera and significant OR taxa were made to the mean AUC generated from the 100 different 10-fold CV runs of the respective OTU-based model. Finally, the power of each study was assessed for an effect size ranging from 1% to 30% and an estimated sample size, for these effect sizes, was generated based on 80% power. For comparisons in which only control versus adenoma individuals were made, the carcinoma individuals were excluded from each respective study. Similarly, for comparisons in which control versus carcinoma individuals were made the adenoma individuals were excluded from each respective study. For all analysis completed fecal and tissue samples were kept separate. Within the tissue groups the data were further divided between samples from the same individual (matched) and those from different individuals (unmatched).

Obtaining Genera Relative Abundance and Significant OR Taxa Models: For the genera analysis of the ORs, OTUs were added together based on the genus or lowest available taxonomic classification level and the total average counts, for 100 different

subsamplings was obtained. The significant OR taxa models for the Random Forest models utilized all taxa that had significant ORs after multiple comparison correction. This meant only models for the carcinoma stool (8 variables) and carcinoma unmatched (3 variables) samples were possible to be created and analyzed.

Matched versus Unmatched Tissue Samples: In general, tissue samples with control and tumor samples from different individuals were classified as unmatched while samples that belonged to the same individual were classified as matched. Studies with matched data included Burns, Dejea, Geng, and Lu while those with unmatched data were from Burns, Flemer, Chen, and Sanapareddy. For some studies samples became unmatched when a corresponding matched sample did not make it through sequence processing. All samples, from both matched and unmatched tissue samples, were analyzed together for the linear mixed-effect models with samples from the same individual being corrected for. All other analysis, where it is not specified explicitly, matched and unmatched samples were analyzed separately using the statistical approaches mentioned in the Statistical Analysis section.

Assessing Important Random Forest Model Variables: Using Mean Decrease in Accuracy (MDA) the top 10 most important variables to the Random Forest model were obtained for the full models of the two different approaches used. For the first approach utilizing genus-based models, the number of times that a specific taxa showed up in the top 10 of the training set across each study was counted. For the second approach, that utilized the OTU-based models, the medians for each OTU across 100 different 80/20 splits of the data was generated and the top 10 OTUs then counted for each study. Common taxa were then identified by using the lowest classification for each of the specific OTUs obtained from these counts and the number of times this classification occurred across the top 10 of each study was recorded. Finally, the two studies that had adenoma tumor tissue (Lu and Flemer) were equally divided between matched and unmatched studies and were

390 grouped together for the counting of the top 10 genera and OTUs for both Random Forest
391 approaches.

392 ***Reproducible Methods:*** The code and analysis can be found at [https://github.com/](https://github.com/SchlossLab/Size_CRCMetaAnalysis_Microbiome_2017)
393 SchlossLab/Size_CRCMetaAnalysis_Microbiome_2017. Unless otherwise mentioned, the
394 accession number of raw sequences from the studies used in this analysis can be found
395 directly in the respective batch file in the GitHub repository or in the original manuscript.

Declarations

Ethics approval and consent to participate

Ethics approval and informed consent for each of the studies used is mentioned in the respective manuscripts used in this meta-analysis.

Consent for publication

Not applicable.

Availability of data and material

A detailed and reproducible description of how the data were processed and analyzed for each study can be found at https://github.com/SchlossLab/Size_CRCMetaAnalysis_Microbiome_2017. Raw sequences can be downloaded from the SRA in most cases and can be found in the respective study batch file in the GitHub repository or within the original publication. For instances when sequences are not publicly available, they may be accessed by contacting the corresponding authors from whence the data came.

Competing Interests

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

Funding

MAS is supported by a Canadian Institute of Health Research fellowship and a University of Michigan Postdoctoral Translational Scholar Program grant.

Authors' contributions

All authors helped to design and conceptualize the study. MAS identified and analyzed the data. MAS and PDS interpreted the data. MAS wrote the first draft of the manuscript and both he and PDS reviewed and revised updated versions. All authors approved the final manuscript.

Acknowledgements

The authors would like to thank all the study participants who were a part of each of the individual studies utilized. We would also like to thank each of the study authors for making their data available for use. Finally, we would like to thank the members of the Schloss lab for valuable feed back and proof reading during the formulation of this manuscript.

References

1. Siegel RL, Miller KD, Jemal A. Cancer statistics, 2016. CA: a cancer journal for clinicians. 2016;66:7–30.
2. Flynn KJ, Baxter NT, Schloss PD. Metabolic and Community Synergy of Oral Bacteria in Colorectal Cancer. mSphere. 2016;1.
3. Goodwin AC, Destefano Shields CE, Wu S, Huso DL, Wu X, Murray-Stewart TR, et al. Polyamine catabolism contributes to enterotoxigenic *Bacteroides fragilis*-induced colon tumorigenesis. Proceedings of the National Academy of Sciences of the United States of America. 2011;108:15354–9.
4. Abed J, Emgård JEM, Zamir G, Faroja M, Almogy G, Grenov A, et al. Fap2 Mediates *Fusobacterium nucleatum* Colorectal Adenocarcinoma Enrichment by Binding to Tumor-Expressed Gal-GalNAc. Cell Host & Microbe. 2016;20:215–25.
5. 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, NY). 2012;338:120–3.
6. Kostic AD, Chun E, Robertson L, Glickman JN, Gallini CA, Michaud M, et al. *Fusobacterium nucleatum* potentiates intestinal tumorigenesis and modulates the tumor-immune microenvironment. Cell Host & Microbe. 2013;14:207–15.
7. Wu S, Rhee K-J, Albesiano E, Rabizadeh S, Wu X, Yen H-R, et al. A human colonic commensal promotes colon tumorigenesis via activation of T helper type 17 T cell responses. Nature Medicine. 2009;15:1016–22.
8. Zackular JP, Baxter NT, Chen GY, Schloss PD. Manipulation of the Gut Microbiota

Reveals Role in Colon Tumorigenesis. *mSphere*. 2016;1.

9. Zackular JP, Baxter NT, Iverson KD, Sadler WD, Petrosino JF, Chen GY, et al. The gut microbiome modulates colon tumorigenesis. *mBio*. 2013;4:e00692–00613.

10. Baxter NT, Zackular JP, Chen GY, Schloss PD. Structure of the gut microbiome following colonization with human feces determines colonic tumor burden. *Microbiome*. 2014;2:20.

11. Ahn J, Sinha R, Pei Z, Dominianni C, Wu J, Shi J, et al. Human gut microbiome and risk for colorectal cancer. *Journal of the National Cancer Institute*. 2013;105:1907–11.

12. 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.

13. Chen W, Liu F, Ling Z, Tong X, Xiang C. Human intestinal lumen and mucosa-associated microbiota in patients with colorectal cancer. *PloS One*. 2012;7:e39743.

14. Wang T, Cai G, Qiu Y, Fei N, Zhang M, Pang X, et al. Structural segregation of gut microbiota between colorectal cancer patients and healthy volunteers. *The ISME journal*. 2012;6:320–9.

15. Burns MB, Lynch J, Starr TK, Knights D, Blekhman R. Virulence genes are a signature of the microbiome in the colorectal tumor microenvironment. *Genome Medicine*. 2015;7:55.

16. 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.

17. Flemer B, Lynch DB, Brown JMR, Jeffery IB, Ryan FJ, Claesson MJ, et al. Tumour-associated and non-tumour-associated microbiota in colorectal cancer. *Gut*.

2017;66:633–43.

18. García AZG. Factors influencing colorectal cancer screening participation. Gastroenterology Research and Practice [Internet]. Hindawi Limited; 2012;2012:1–8. Available from: <https://doi.org/10.1155/2012/483417>

19. Geng J, Fan H, Tang X, Zhai H, Zhang Z. Diversified pattern of the human colorectal cancer microbiome. Gut Pathogens. 2013;5:2.

20. 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.

21. 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 [Internet]. Springer Nature; 2014;5:4724. Available from: <https://doi.org/10.1038/ncomms5724>

22. Aymeric L, Donnadieu F, Mulet C, Merle L du, Nigro G, Saffarian A, et al. Colorectal cancer specific conditions promote *Streptococcus gallolyticus* gut colonization. Proceedings of the National Academy of Sciences [Internet]. Proceedings of the National Academy of Sciences; 2017;115:E283–91. Available from: <https://doi.org/10.1073/pnas.1715112115>

23. Weir TL, Manter DK, Sheflin AM, Barnett BA, Heuberger AL, Ryan EP. Stool microbiome and metabolome differences between colorectal cancer patients and healthy adults. PLoS One. 2013;8:e70803.

24. Boleij A, Hechenbleikner EM, Goodwin AC, Badani R, Stein EM, Lazarev MG, et al. The bacteroides fragilis toxin gene is prevalent in the colon mucosa of colorectal cancer patients. Clinical Infectious Diseases [Internet]. Oxford University Press (OUP);

2014;60:208–15. Available from: <https://doi.org/10.1093/cid/ciu787>

25. Sanapareddy N, Legge RM, Jovov B, McCoy A, Burcal L, Araujo-Perez F, et al. Increased rectal microbial richness is associated with the presence of colorectal adenomas in humans. *The ISME journal*. 2012;6:1858–68.

26. Lu Y, Chen J, Zheng J, Hu G, Wang J, Huang C, et al. Mucosal adherent bacterial dysbiosis in patients with colorectal adenomas. *Scientific Reports*. 2016;6:26337.

27. Hale VL, Chen J, Johnson S, Harrington SC, Yab TC, Smyrk TC, et al. Shifts in the Fecal Microbiota Associated with Adenomatous Polyps. *Cancer Epidemiology, Biomarkers & Prevention: A Publication of the American Association for Cancer Research, Cosponsored by the American Society of Preventive Oncology*. 2017;26:85–94.

28. Shah MS, DeSantis TZ, Weinmaier T, McMurdie PJ, Cope JL, Altrichter A, et al. Leveraging sequence-based faecal microbial community survey data to identify a composite biomarker for colorectal cancer. *Gut*. 2017;

29. Brim H, Yooseph S, Zoetendal EG, Lee E, Torralbo M, Laiyemo AO, et al. Microbiome analysis of stool samples from African Americans with colon polyps. *PloS One*. 2013;8:e81352.

30. Flynn KJ, Ruffin MT, Turgeon DK, Schloss PD. Spatial variation of the native colon microbiota in healthy adults. Cold Spring Harbor Laboratory; 2017; Available from: <https://doi.org/10.1101/189886>

31. Salter SJ, Cox MJ, Turek EM, Calus ST, Cookson WO, Moffatt MF, et al. Reagent and laboratory contamination can critically impact sequence-based microbiome analyses. *BMC Biology [Internet]*. Springer Nature; 2014;12. Available from: <https://doi.org/10.1186/>

- 514 32. Keku TO, Dulal S, Deveau A, Jovov B, Han X. The gastrointestinal microbiota and
515 colorectal cancer. *American Journal of Physiology - Gastrointestinal and Liver Physiology*
516 [Internet]. 2015 [cited 2017 Oct 30];308:G351–63. Available from: [http://ajpgi.physiology.](http://ajpgi.physiology.org/lookup/doi/10.1152/ajpgi.00360.2012)
517 [org/lookup/doi/10.1152/ajpgi.00360.2012](http://ajpgi.physiology.org/lookup/doi/10.1152/ajpgi.00360.2012)
- 518 33. Vogtmann E, Goedert JJ. Epidemiologic studies of the human microbiome and cancer.
519 *British Journal of Cancer* [Internet]. 2016 [cited 2017 Oct 30];114:237–42. Available from:
520 <http://www.nature.com/doi/10.1038/bjc.2015.465>
- 521 34. Kostic AD, Gevers D, Pedamallu CS, Michaud M, Duke F, Earl AM, et al. Genomic
522 analysis identifies association of *Fusobacterium* with colorectal carcinoma. *Genome*
523 *Research*. 2012;22:292–8.
- 524 35. Zackular JP, Rogers MAM, Ruffin MT, Schloss PD. The human gut microbiome as
525 a screening tool for colorectal cancer. *Cancer Prevention Research (Philadelphia, Pa)*.
526 2014;7:1112–21.
- 527 36. Schloss PD, Westcott SL, Ryabin T, Hall JR, Hartmann M, Hollister EB, et al.
528 Introducing mothur: Open-Source, Platform-Independent, Community-Supported Software
529 for Describing and Comparing Microbial Communities. *Appl Environ Microbiol* [Internet].
530 2009 [cited 12AD Jan 1];75:7537–41. Available from: [http://aem.asm.org/cgi/content/](http://aem.asm.org/cgi/content/abstract/75/23/7537)
531 [abstract/75/23/7537](http://aem.asm.org/cgi/content/abstract/75/23/7537)
- 532 37. Rognes T, Flouri T, Nichols B, Quince C, Mahé F. VSEARCH: A versatile open source
533 tool for metagenomics. *PeerJ*. 2016;4:e2584.
- 534 38. Westcott SL, Schloss PD. OptiClust, an Improved Method for Assigning

- 535 Amplicon-Based Sequence Data to Operational Taxonomic Units. mSphere. 2017;2.
- 536 39. R Core Team. R: A language and environment for statistical computing [Internet].
537 Vienna, Austria: R Foundation for Statistical Computing; 2017. Available from: <https://www.R-project.org/>
538
- 539 40. Mangiafico S. Rcompanion: Functions to support extension education program
540 evaluation [Internet]. 2017. Available from: [https://CRAN.R-project.org/package=](https://CRAN.R-project.org/package=rcompanion)
541 [rcompanion](https://CRAN.R-project.org/package=rcompanion)
- 542 41. Fox J, Weisberg S. An R companion to applied regression [Internet]. Second. Thousand
543 Oaks CA: Sage; 2011. Available from: [http://socserv.socsci.mcmaster.ca/jfox/Books/](http://socserv.socsci.mcmaster.ca/jfox/Books/Companion)
544 [Companion](http://socserv.socsci.mcmaster.ca/jfox/Books/Companion)
- 545 42. Bates D, Mächler M, Bolker B, Walker S. Fitting linear mixed-effects models using lme4.
546 Journal of Statistical Software. 2015;67:1–48.
- 547 43. Telmo Nunes MS with contributions from, Heuer C, Marshall J, Sanchez J, Thornton
548 R, Reiczigel J, et al. EpiR: Tools for the analysis of epidemiological data [Internet]. 2017.
549 Available from: <https://CRAN.R-project.org/package=epiR>
- 550 44. Viechtbauer W. Conducting meta-analyses in R with the metafor package. Journal of
551 Statistical Software [Internet]. 2010;36:1–48. Available from: [http://www.jstatsoft.org/v36/](http://www.jstatsoft.org/v36/i03/)
552 [i03/](http://www.jstatsoft.org/v36/i03/)
- 553 45. Oksanen J, Blanchet FG, Friendly M, Kindt R, Legendre P, McGlinn D, et al. Vegan:
554 Community ecology package [Internet]. 2017. Available from: [https://CRAN.R-project.org/](https://CRAN.R-project.org/package=vegan)
555 [package=vegan](https://CRAN.R-project.org/package=vegan)
- 556 46. Jed Wing MKC from, Weston S, Williams A, Keefer C, Engelhardt A, Cooper T,
557 et al. caret: Classification and regression training [Internet]. 2017. Available from:

558 <https://CRAN.R-project.org/package=caret>

559 47. Liaw A, Wiener M. Classification and regression by randomForest. R News [Internet].
560 2002;2:18–22. Available from: <http://CRAN.R-project.org/doc/Rnews/>

561 48. Champely S. Pwr: Basic functions for power analysis [Internet]. 2017. Available from:
562 <https://CRAN.R-project.org/package=pwr>

563 49. Giner G, Smyth GK. Statmod: Probability calculations for the inverse gaussian
564 distribution. R Journal. 2016;8:339–51.

565 50. Wickham H. Ggplot2: Elegant graphics for data analysis [Internet]. Springer-Verlag
566 New York; 2009. Available from: <http://ggplot2.org>

567 51. Auguie B. GridExtra: Miscellaneous functions for “grid” graphics [Internet]. 2017.
568 Available from: <https://CRAN.R-project.org/package=gridExtra>

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

Table 1: Total Individuals in each Study Included in the Stool Analysis

Study	Data Stored	16S Region	Control (n)	Adenoma (n)	Carcinoma (n)
Ahn	DBGap	V3-4	148	0	62
Baxter	SRA	V4	172	198	120
Brim	SRA	V1-3	6	6	0
Flemer	Author	V3-4	37	0	43
Hale	Author	V3-5	473	214	17
Wang	SRA	V3	56	0	46
Weir	Author	V4	4	0	7
Zeller	SRA	V4	50	37	41

Table 2: Studies with Tissue Samples Included in the Analysis

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

Figure 1: Odds Ratio for Adenoma or Carcinoma based on Bacterial Community Metrics in Stool. A) Community-based odds ratio for adenoma. B) Community-based odds ratio for carcinoma. Colors represent the different variable regions used within the respective study.

Figure 2: Stool Random Forest Model Train AUCs. A) Adenoma random forest model AUCs between all genera, all OTU, and select model based on significant OR taxa. B) Carcinoma random forest model AUCs between all genera, all OTU, and select model based on significant OR taxa. The black line represents the median AUC for the respective group. If no values are present in the significant OR taxa group then there were no significant taxa identified and no model was tested.

Figure 3: Stool Random Forest Genus-Based Model Test AUCs. A) Test AUCs of adenoma models using all genera across study. B) Test AUCs of carcinoma models using all genera or significant OR taxa only. The black line represents the AUC at 0.5. The red lines represent the median AUC of all test AUCs for a specific study.

Figure 4: Most Common Taxa Across Full Community Stool Study Models. A) Common taxa in the top 10 for adenoma Random Forest all genera-based models. B) Common taxa in the top 10 for carcinoma Random Forest all genera-based models. C) Common taxa in the top 10 for adenoma Random Forest all OTU-based models. D) Common genera in the top 10 for carcinoma Random Forest all OTU-based models.

Figure 5: Power and Effect Size Analysis of Studies Included. A) Power based on effect size for studies with adenoma individuals. B) Power based on effect size for studies with carcinoma individuals. C) The estimated sample number needed for each arm of each study to detect an effect size of 1-30%. The dotted red lines in A) and B) represent a power of 0.8.

Figure S1: Odds Ratio for Adenoma or Carcinoma based on Bacterial Community Metrics in Tissue. A) Community-based odds ratio for adenoma. B) Community-based odds ratio for carcinoma. Colors represent the different variable regions used within the respective study.

Figure S2: Tissue Random Forest Model Train AUCs. A) Adenoma random forest model AUCs between all genera, all OTU, and select model based on significant OR taxa in unmatched and matched tissue. B) Carcinoma random forest model AUCs between all genera, all OTU, and select model based on significant OR taxa in unmatched and matched tissue. The black line represents the median AUC for the respective group. If no values are present in the significant OR taxa group then there were no significant taxa identified and no model was tested.

Figure S3: Tissue Random Forest Genus-Based Model Test AUCs. A) Test AUCs of adenoma models using all genera across study. B) Test AUCs of carcinoma models using all genera for matched tissue studies. C) Test AUCs of carcinoma models using all genera or significant OR taxa only for unmatched tissue studies. The black line represents the AUC at 0.5. The red lines represent the median AUC of all test AUCs for a specific study.

Figure S4: Most Common Genera Across Full Community Tissue Study Models.

A) Common genera in the top 10 for adenoma Random Forest for all genera-based models. B) Common genera in the top 10 for unmatched carcinoma Random Forest for all genera-based models. C) Common genera in the top 10 for matched carcinoma Random Forest for all genera-based models. D) Common genera in the top 10 for adenoma Random Forest for all OTU-based models. E) Common genera in the top 10 for unmatched carcinoma Random Forest for all OTU-based models. F) Common genera in the top 10 for matched carcinoma Random Forest for all OTU-based models.