

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

Marc A Sze<sup>1</sup> and Patrick D Schloss<sup>1†</sup>

† To whom correspondence should be addressed: [pschloss@umich.edu](mailto:pschloss@umich.edu)

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

Co-author e-mails:

- [marcsze@med.umich.edu](mailto: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*, *Methanosphaera*, *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 Odds Ratio (OR) of Tumors:***

We first assessed whether variation in broad community metrics like total number of operational taxonomic units (OTUs) (i.e. richness), the evenness of their abundance, and the overall diversity was associated with disease stage after controlling for study and variable region differences. In stool, there was a significant decrease in both evenness and diversity as disease severity progressed from normal to adenoma to carcinoma (P-value = 0.025 and 0.043, respectively) [Figure 1]; there was not a significant difference for richness (P-value = 0.21). We next tested whether the decrease in these community metrics translated into significant ORs for having an adenoma or carcinoma. For fecal samples, the ORs for richness were not significantly greater than 1.0 for adenoma or carcinoma (P-value = 0.40) [Figure 2A]. The ORs for evenness were significantly higher than 1.0 for adenoma (OR = 1.3 (1.02 - 1.65), P-value = 0.035) and carcinoma (OR = 1.66 (1.2 - 2.3), P-value = 0.0021) [Figure 2B]. The ORs for diversity were only significantly greater than 1.0 for carcinoma (OR = 1.61 (1.14 - 2.28), P-value = 0.0069), but not for adenoma (P-value = 0.11) [Figure 2C]. Although these OR are significantly greater than 1.0, it is doubtful that these are clinically meaningful values.

Similar to our analysis of sequences obtained from stool samples, we repeated the analysis using sequences obtained from colon tissue. There were no significant changes in richness, evenness, or diversity as disease severity progressed from control to adenoma to carcinoma (P-value > 0.05). We next analyzed the OR, for matched (i.e. where unaffected tissue and tumors were obtained from the same individual) and unmatched (i.e. where unaffected tissue and tumor tissue were not obtained from the same individual) tissue samples. The ORs for adenoma and carcinoma by any measure were not significantly different from 1.0 (P-value > 0.05) [Figure S1 & Table S1]. This is likely due to the combination of a small effect size, as suggested from the results using stool, and the

relatively small number of studies and size of studies used in the analysis.

***Disease Progression is Associated with Community-Wide Changes in Composition***

***and Abundance:*** Based on the differences in evenness and diversity, we next asked whether there were community-wide differences in the structure of the communities associated with different disease stages. We identified significant bacterial community differences in the stool of patients with adenomas relative to those with normal colons in 1 of 4 studies and in patients with carcinomas relative to those with normal colons in 6 of 7 studies (PERMANOVA; P-value < 0.05) [Table S2]. Similar to the analyses using stool samples, there were significant differences in bacterial community structure between subjects with normal colons and those with adenoma (1 of 2 studies) and carcinoma (1 of 3 studies) [Table S2]. For studies that used matched samples no differences in bacterial community structures were observed [Table S2]. Combined, these results indicate that there consistent and significant community-wide changes in the fecal community structure of subjects with carcinomas. However, the signal observed in subjects with adenomas or when using tissue samples was not as consistent. This is likely due to a smaller effect size or the relatively small sample sizes among the studies that characterized the tissue microbiota.

***Individual Taxa are Associated with Significant ORs for Carcinomas:***

Next we identified those taxa were associated with ORs that were significantly associated with having a normal colon or the presence of adenomas or carcinomas. No taxa had a significant OR for the presence of adenomas when we used data collected from stool or tissue samples (Table S3 & S4). In contrast, 8 taxa had significant ORs for the presence of carcinomas using data from stool samples. Of these, 4 are commonly associated with the oral cavity: *Fusobacterium* (OR = 2.74 (1.95 - 3.85)), *Parvimonas* (OR = 3.07 (2.11 - 4.46)), *Porphyromonas* (OR = 3.2 (2.26 - 4.54)), and *Peptostreptococcus* (OR = 7.11 (3.84 - 13.17)) [Table S3]. The other 4 were *Clostridium* XI (OR = 0.65 (0.49 - 0.86)),



Enterobacteriaceae (OR = 1.79 (1.33 - 2.41)), Escherichia (OR = 2.15 (1.57 - 2.95)), and Ruminococcus (OR = 0.63 (0.48 - 0.83)). Among the data collected from tissue samples, only unmatched carcinoma samples had taxa with a significant OR. Those included 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)). Mouth-associated genera were not significantly associated with an increased OR for carcinoma in tissue samples [Table S4]. For example, Fusobacterium had an OR of 3.98 (1.19 - 13.24; however, due to the small number of studies and considerable variation in the data, the Benjamini-Hochberg-corrected P-value was 0.93 [Table S4]. It is interesting to note that Ruminococcus and members of Clostridium group XI in stool and Dorea and Blautia in tissue had ORs that were significantly less than 1.0, which suggests that these populations are protective against the development of carcinomas. Overall, there was no overlap in the taxa with significant OR between stool and tissue samples.

***Individual taxa with a significant OR do a poor job of differentiating subjects with normal colons and those with carcinoma:*** We next asked whether those taxa that had a significant OR associated with having a normal colon or carcinomas could be used individually, to classify subjects as having a normal colon or carcinomas. Whereas the OR was defined based on whether the relative abundance for a taxon in a subject was above or below the median relative abundance for that taxon across all subjects in a study, we generated receiver operator characteristic (ROC) curves for each taxon in each study and calculated the area under the curve (AUC). This allowed us to use a more fluid relative abundance threshold for defining disease status. Using data from stool samples, the 8 taxa did no better at classifying the subjects than one would expect by chance (i.e. AUC=0.50) [Figure 3A]. The taxa that performed the best included Clostridium XI, Ruminococcus, and Escherichia and even these had median AUC values less than 0.588. Likewise, in unmatched tissue samples the 8 taxa with significant ORs taxa were marginally better than one would expect by chance [Figure 3B]. The relative abundance of Dorea was the best predictor of carcinomas and its median AUC was only 0.62. These results suggest that

although these taxa are associated with a decreased or increased OR for the presences of carcinomas, individually, they do a poor job of classifying a subject's disease status.

**Combined taxa model classifies subjects better than using individual taxa:** Instead of attempting to classify subjects based on individual taxa, next we generated Random Forest models that combined the individual taxa and evaluated the ability to classify as subject's disease status. For data from stool samples, the combined model had an AUC of 0.75, which was significantly higher than any of the AUC values for the individual taxa (P-value < 0.033). For the full taxa models using stool, *Bacteroides* and *Lachnospiraceae* were the most common taxa in the top 10% mean decrease in accuracy (MDA) across studies [Figure S2]. Similarly, using data from the unmatched tissue samples, the combined model had an AUC of 0.77, which was significantly higher than the AUC values for *Blautia* and *Weissella* (P-value < 0.037). For the full taxa models using unmatched tissue, *Lachnospiraceae*, *Bacteroidaceae*, and *Ruminococcaceae* were the most common taxa in the top 10% mean decrease in accuracy across studies [Figure S3]. Clearly, pooling the information from the taxa with significant ORs results in a model that outperforms classifications made using individual taxa.

**Performance of models based on taxa relative abundance in full community are better than those based on taxa with significant ORs:** Next, we asked whether a Random Forest classification model built using all of the taxa found in the communities would outperform the models generated using those taxa with a significant OR. Similar to our inability to identify taxa associated with a significant OR for the presence of adenomas, the median AUCs to classify subjects as having normal colons or having adenomas using data from stool or tissue samples were marginally better than 0.5 for any study [Figure 4A & S4A]. In contrast, the models for classifying subjects as having normal colons or having carcinomas using data from stool or tissue samples yielded AUC values meaningfully higher than 0.5 [Figure 4B & S4B-C]. When we compared the models based on all of

the taxa in a community to models based on the taxa with significant ORs, the results were mixed. Using the data from stool samples we found that although the AUC for 6 of 7 studies increased (mean decrease = 9.53%), the more expansive models performed worse for 1 of the studies (decrease = 0.38%). The overall improvement in performance was statistically significant (one-tailed paired T-test; P-value = 0.005). Of the 8 taxa with significant ORs, all 8 were among the top 10% most important taxa as measured by mean decrease in accuracy, in at least one study. Similarly, using the data from unmatched tissue samples we found that the AUC for 4 out of 4 studies decreased between full versus select OR models (mean decrease = 19.11%, one-tailed paired T-test; P-value = 0.03). Of the 3 taxa with significant ORs, all 3 were among the top 10% most important taxa as measured by mean decrease in accuracy, in at least one study. The most important taxa across study within the significant OR taxa only models for stool were *Ruminococcus* and *Clostridium XI* [Figure 5A]. For the significant OR taxa unmatched tissue models both *Dorea* and *Blautia* were the important based on mean decrease in accuracy [Figure 5B]. These results were surprising because it demonstrated that the ability to classify subjects could be done based on a limited characterization of the communities.

***Performance of models based on OTU relative abundance in full community are not significantly better than those based on taxa with significant ORs:*** The previous models were based on relative abundance data where sequences were assigned to coarse taxonomic assignments (i.e. typically genus or family level). To determine whether model performance improved with a more fine scale classification, we assigned sequences to operational taxonomic units (OTUs) where the similarity among sequences within an OTU was more than 97%. We again found that classification models built using all of the sequence data for a community did a poor job of differentiating between subjects with normal colons and those with adenomas (median AUC: 0.53 [0.37- 0.56]), but did a good job of differentiating between subjects with normal colons and those with carcinomas (median AUC: 0.71 [0.5- 0.9]). The OTU-based models performed similarly to those

constructed using the taxa with significant ORs (one-tailed paired T-test; P-value = 0.966) and those using all taxa (one-tailed paired T-test; P-value = 0.146). Among the OTUs that had the highest mean decrease in accuracy for the OTU-based models, we found that OTUs that affiliated with all of the 8 taxa that had a significant OR were within the top 10% for at least one study. Again, this result was surprising as it indicated that a finer scale classification of the sequence data and thus a larger number of features to select from, did not yield improved classification of the subjects.

***Generalizability of taxon-based models trained on one dataset to the other datasets:*** Considering the good performance of the Random Forest models using taxa with a significant OR and using all of the taxa, we next asked how well the models would perform when given data from a different subject cohort. For instance, if a model was trained using data from the Ahn study, we wanted to know how well it would perform using the data from the Baxter study. We found the models trained using the taxa with a significant OR all had a higher median AUC than the models trained using all of the taxa when tested on the other datasets [Figure 6 & S5]. As might be expected, the difference between the performance of the modelling approaches appeared to vary with the size of the training cohort [Figure 6]. These data suggest that given a sufficient number of subjects with normal colons and carcinomas, Random Forest models trained using a small number of taxa can accurately classify individuals from a different cohort.

## Discussion

Although we expected that the full OTU models would perform the best at classifying individuals with and without carcinomas, our observations suggest that both the full and significant OR taxa models performed equally well. These results suggest that lower level classification to species and strain may not add extra useful information with respect to prediction models. This has been suggested in previous literature where metagenomics did not perform better than 16S rRNA gene sequencing data at classifying individuals with normal colons and those with carcinomas (30). One possible reason as to why lower level classification may not result in better models is that the communities are patchy and higher level taxonomic information pools some of this patchiness allowing for better prediction models. There may also be a fair bit of data redundancy within models that utilize more of the community. An example of this redundancy would be when we trained models on one study and tested it on the other studies and the AUCs of the models created with the select OR taxa performed as well as full taxa models [Figure 6B].

Our observations also suggest that a small collection of taxa can classify disease as well as full OTU-based models but that these taxa individually perform quite poorly [Figure 3]. This result supports the contention that there might be redundancy of function even amongst the taxa included in the significant OR models. As an example multiple different microbes could be similarly stimulating the activation of inflammatory pathways and by doing so exacerbate disease progression. Multiple reports within the literature have found that different bacteria, such as *Escherichia coli* and *Fusobacterium nucleatum*, can worsen similarly worsen inflammation in mouse models of tumorigenesis (5, 6, 21). Although the inflammatory taxa were patchy in their importance and presence across studies those that were not typically associated with inflammation were consistently important for every study [Figure 6]. The loss of these taxa (*Ruminococcus* and *Clostridium XI* in stool and *Dorea* and *Blautia* in unmatched tissue) is particularly interesting because many are commonly

thought to be beneficial due to their involvement in production of short chain fatty acids (31–33).

The adenoma models as a whole performed poorly in classifying individuals with and without adenomas. This outcome is not inconsistent with what has been published previously (27, 34). However, the modeling results are at odds with results obtained in Baxter, *et al.* (12). There are some major differences between the models generated in this meta-analysis and what was used in this previous report. First, the previous report's models investigated the classification of lesions (individuals with adenoma or carcinoma) and not adenoma alone. The Baxter, *et al.* models also contained Fecal Immunoglobulin Test data while our meta-analysis models only contained 16S rRNA gene sequencing data. Although being able to classify individuals with adenomas is important, the classification of advanced adenomas is a more clinically meaningful diagnostic tool (i.e. those that are at high risk of progressing to a carcinoma). It is possible that we might have been able to detect differences in the bacterial community if advanced adenomas were separated from adenomas but that data was not available for the majority of studies analyzed. It is also possible that the initial changes to the bacterial community are focal to where the initial adenoma develops and would not be easily assessed with a fecal sample.

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, 35). Sampling the microbiota directly associated with colon tissue may provide clearer answers but 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 S3]. *Fusobacterium* was not identified

consistently across studies and this could be due to both a small number of studies and a small sample size within these studies. Additionally, the majority of the colon tissue-based results were study specific with many of the top 10 taxa being present only in a single study. The presence of genera associated with contamination (36), within the top 10 most important variables for the genera and OTU models is worrying (e.g. *Novosphingobium*, *Acidobacteria Gp2*, *Sphingomonas*, etc.). The low bacterial biomass of tissue samples coupled with potential contamination and small sample sizes 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). There are reports that *Bacteroides fragilis*, positive for the enterotoxigenic gene, are found at specific locations along the colon but the samples we were able to use in this meta-analysis could not identify these types of differences (37). Additionally, 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.

311 Despite these limitations the findings that we present here would not be possible without  
312 performing a meta-analysis. These types of studies can be a useful tool in microbiota  
313 research because they can both validate existing research and make new discoveries  
314 by pooling many independent investigations together. Yet, it is still difficult to perform  
315 these studies because of inaccessible 16S sequencing data, missing or vague metadata  
316 (e.g. which samples are carcinoma and which are not), varying sequence quality, and  
317 multiple small data sets. Better attention to these specific problems could help to increase  
318 the reproducibility and replicability of microbiota studies and make it easier to perform  
319 these crucial meta-analyses. Moving forward, meta-analyses will be important tools to help  
320 aggregate and find commonalities across studies when investigating the microbiota in the  
321 context of a specific disease and more are needed (28, 38–40).

322 By aggregating together a large collection of studies analyzing both fecal and colon tissue  
323 samples, we are able to provide evidence supporting the importance of the bacterial  
324 community in carcinoma tumors. Although further validation of the biomarkers presented  
325 here need to be undertaken, the replicability of the AUC of a specific collection of taxa  
326 across multiple studies suggests a strong potential for the use of the microbiota as a risk  
327 stratification tool for individuals with carcinomas.



## 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.* (41, 42) 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.* (43), Lu, *et al.* (26), Sanapareddy, *et al.* (25), Wang, *et al.* (14), Weir, *et al.* (23), and Zeller, *et al.* (16). The Zackular (44) 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 (45) 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 (46) before *de novo* OTU clustering at 97% similarity was completed using the OptiClust algorithm (47).

**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 (48). We then built Random Forest models based on all genera, all OTUs, or significant OR taxa (only using taxa still significant 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. 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 grouped together for the counting of the top 10 genera and OTUs for both Random Forest approaches.

**Statistical Analysis:** All statistical analysis after sequence processing utilized the R (v3.4.3) software package (49). For OTU richness, evenness, and Shannon diversity analysis, values were power transformed using the rcompanion (v1.11.1) package (50) and then Z-score normalized using the car (v2.1.6) package (51). Testing for OTU richness, evenness, and Shannon diversity differences utilized linear mixed-effect models created using the lme4 (v1.1.15) package (52) 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 (53, 54) 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 (55). Random Forest models were built using both the caret (v6.0.78) and randomForest (v4.6.12) packages (56, 57). All figures were created using both ggplot2 (v2.2.1) and gridExtra (v2.3) packages (58, 59).

**Reproducible Methods:** The code and analysis can be found at <https://github.com/>

432 SchlossLab/Size\_CRCMetaAnalysis\_Microbiome\_2017. Unless otherwise mentioned, the  
433 accession number of raw sequences from the studies used in this analysis can be found  
434 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](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, R. L., K. D. Miller, and A. Jemal.** 2016. Cancer statistics, 2016. *CA: a cancer journal for clinicians* **66**:7–30.
2. **Flynn, K. J., N. T. Baxter, and P. D. Schloss.** 2016. Metabolic and Community Synergy of Oral Bacteria in Colorectal Cancer. *mSphere* **1**.
3. **Goodwin, A. C., C. E. Destefano Shields, S. Wu, D. L. Huso, X. Wu, T. R. Murray-Stewart, A. Hacker-Prietz, S. Rabizadeh, P. M. Woster, C. L. Sears, and R. A. Casero.** 2011. Polyamine catabolism contributes to enterotoxigenic *Bacteroides fragilis*-induced colon tumorigenesis. *Proceedings of the National Academy of Sciences of the United States of America* **108**:15354–15359.
4. **Abed, J., J. E. M. Emgård, G. Zamir, M. Faroja, G. Almogy, A. Grenov, A. Sol, R. Naor, E. Pikarsky, K. A. Atlan, A. Mellul, S. Chaushu, A. L. Manson, A. M. Earl, N. Ou, C. A. Brennan, W. S. Garrett, and G. Bachrach.** 2016. Fap2 Mediates *Fusobacterium nucleatum* Colorectal Adenocarcinoma Enrichment by Binding to Tumor-Expressed Gal-GalNAc. *Cell Host & Microbe* **20**:215–225.
5. **Arthur, J. C., E. Perez-Chanona, M. Mühlbauer, S. Tomkovich, J. M. Uronis, T.-J. Fan, B. J. Campbell, T. Abujamel, B. Dogan, A. B. Rogers, J. M. Rhodes, A. Stintzi, K. W. Simpson, J. J. Hansen, T. O. Keku, A. A. Fodor, and C. Jobin.** 2012. Intestinal inflammation targets cancer-inducing activity of the microbiota. *Science (New York, N.Y.)* **338**:120–123.
6. **Kostic, A. D., E. Chun, L. Robertson, J. N. Glickman, C. A. Gallini, M. Michaud, T. E. Clancy, D. C. Chung, P. Lochhead, G. L. Hold, E. M. El-Omar, D. Brenner, C. S. Fuchs, M. Meyerson, and W. S. Garrett.** 2013. *Fusobacterium nucleatum* potentiates intestinal tumorigenesis and modulates the tumor-immune microenvironment. *Cell Host &*



487 Microbe **14**:207–215.

488 7. **Wu, S., K.-J. Rhee, E. Albesiano, S. Rabizadeh, X. Wu, H.-R. Yen, D. L. Huso, F. L.**  
489 **Brancati, E. Wick, F. McAllister, F. Housseau, D. M. Pardoll, and C. L. Sears.** 2009. A  
490 human colonic commensal promotes colon tumorigenesis via activation of T helper type  
491 17 T cell responses. *Nature Medicine* **15**:1016–1022.

492 8. **Zackular, J. P., N. T. Baxter, G. Y. Chen, and P. D. Schloss.** 2016. Manipulation of  
493 the Gut Microbiota Reveals Role in Colon Tumorigenesis. *mSphere* **1**.

494 9. **Zackular, J. P., N. T. Baxter, K. D. Iverson, W. D. Sadler, J. F. Petrosino, G. Y. Chen,**  
495 **and P. D. Schloss.** 2013. The gut microbiome modulates colon tumorigenesis. *mBio*  
496 **4**:e00692–00613.

497 10. **Baxter, N. T., J. P. Zackular, G. Y. Chen, and P. D. Schloss.** 2014. Structure of the  
498 gut microbiome following colonization with human feces determines colonic tumor burden.  
499 *Microbiome* **2**:20.

500 11. **Ahn, J., R. Sinha, Z. Pei, C. Dominianni, J. Wu, J. Shi, J. J. Goedert, R. B. Hayes,**  
501 **and L. Yang.** 2013. Human gut microbiome and risk for colorectal cancer. *Journal of the*  
502 *National Cancer Institute* **105**:1907–1911.

503 12. **Baxter, N. T., M. T. Ruffin, M. A. M. Rogers, and P. D. Schloss.** 2016.  
504 Microbiota-based model improves the sensitivity of fecal immunochemical test for detecting  
505 colonic lesions. *Genome Medicine* **8**:37.

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

508 14. **Wang, T., G. Cai, Y. Qiu, N. Fei, M. Zhang, X. Pang, W. Jia, S. Cai, and L. Zhao.**  
509 2012. Structural segregation of gut microbiota between colorectal cancer patients and

510 healthy volunteers. The ISME journal **6**:320–329.

511 **15. Burns, M. B., J. Lynch, T. K. Starr, D. Knights, and R. Blehman.** 2015. Virulence  
512 genes are a signature of the microbiome in the colorectal tumor microenvironment.  
513 *Genome Medicine* **7**:55.

514 **16. Zeller, G., J. Tap, A. Y. Voigt, S. Sunagawa, J. R. Kultima, P. I. Costea, A. Amiot,**  
515 **J. Böhm, F. Brunetti, N. Habermann, R. Hercog, M. Koch, A. Luciani, D. R. Mende,**  
516 **M. A. Schneider, P. Schrotz-King, C. Tournigand, J. Tran Van Nhieu, T. Yamada, J.**  
517 **Zimmermann, V. Benes, M. Kloor, C. M. Ulrich, M. von Knebel Doeberitz, I. Sobhani,**  
518 **and P. Bork.** 2014. Potential of fecal microbiota for early-stage detection of colorectal  
519 cancer. *Molecular Systems Biology* **10**:766.

520 **17. Flemer, B., D. B. Lynch, J. M. R. Brown, I. B. Jeffery, F. J. Ryan, M. J. Claesson,**  
521 **M. O’Riordain, F. Shanahan, and P. W. O’Toole.** 2017. Tumour-associated and  
522 non-tumour-associated microbiota in colorectal cancer. *Gut* **66**:633–643.

523 **18. García, A. Z. G.** 2012. Factors influencing colorectal cancer screening participation.  
524 *Gastroenterology Research and Practice*. Hindawi Limited **2012**:1–8.

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

527 **20. Dejea, C. M., E. C. Wick, E. M. Hechenbleikner, J. R. White, J. L. Mark Welch,**  
528 **B. J. Rossetti, S. N. Peterson, E. C. Snesrud, G. G. Borisy, M. Lazarev, E. Stein,**  
529 **J. Vadivelu, A. C. Roslani, A. A. Malik, J. W. Wanyiri, K. L. Goh, I. Thevambiga, K.**  
530 **Fu, F. Wan, N. Llosa, F. Housseau, K. Romans, X. Wu, F. M. McAllister, S. Wu, B.**  
531 **Vogelstein, K. W. Kinzler, D. M. Pardoll, and C. L. Sears.** 2014. Microbiota organization  
532 is a distinct feature of proximal colorectal cancers. *Proceedings of the National Academy*

of Sciences of the United States of America **111**:18321–18326.

21. **Arthur, J. C., R. Z. Gharaibeh, M. Mühlbauer, E. Perez-Chanona, J. M. Uronis, J. McCafferty, A. A. Fodor, and C. Jobin.** 2014. Microbial genomic analysis reveals the essential role of inflammation in bacteria-induced colorectal cancer. *Nature Communications*. Springer Nature **5**:4724.

22. **Aymeric, L., F. Donnadieu, C. Mulet, L. du Merle, G. Nigro, A. Saffarian, M. Bérard, C. Poyart, S. Robine, B. Regnault, P. Trieu-Cuot, P. J. Sansonetti, and S. Dramsi.** 2017. Colorectal cancer specific conditions promote *Streptococcus gallolyticus* gut colonization. *Proceedings of the National Academy of Sciences*. *Proceedings of the National Academy of Sciences* **115**:E283–E291.

23. **Weir, T. L., D. K. Manter, A. M. Sheflin, B. A. Barnett, A. L. Heuberger, and E. P. Ryan.** 2013. Stool microbiome and metabolome differences between colorectal cancer patients and healthy adults. *PloS One* **8**:e70803.

24. **Boleij, A., E. M. Hechenbleikner, A. C. Goodwin, R. Badani, E. M. Stein, M. G. Lazarev, B. Ellis, K. C. Carroll, E. Albesiano, E. C. Wick, E. A. Platz, D. M. Pardoll, and C. L. Sears.** 2014. The bacteroides fragilis toxin gene is prevalent in the colon mucosa of colorectal cancer patients. *Clinical Infectious Diseases*. Oxford University Press (OUP) **60**:208–215.

25. **Sanapareddy, N., R. M. Legge, B. Jovov, A. McCoy, L. Burcal, F. Araujo-Perez, T. A. Randall, J. Galanko, A. Benson, R. S. Sandler, J. F. Rawls, Z. Abdo, A. A. Fodor, and T. O. Keku.** 2012. Increased rectal microbial richness is associated with the presence of colorectal adenomas in humans. *The ISME journal* **6**:1858–1868.

26. **Lu, Y., J. Chen, J. Zheng, G. Hu, J. Wang, C. Huang, L. Lou, X. Wang, and Y. Zeng.** 2016. Mucosal adherent bacterial dysbiosis in patients with colorectal adenomas.

Scientific Reports **6**:26337.

27. **Hale, V. L., J. Chen, S. Johnson, S. C. Harrington, T. C. Yab, T. C. Smyrk, H. Nelson, L. A. Boardman, B. R. Druliner, T. R. Levin, D. K. Rex, D. J. Ahnen, P. Lance, D. A. Ahlquist, and N. Chia.** 2017. 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* **26**:85–94.

28. **Shah, M. S., T. Z. DeSantis, T. Weinmaier, P. J. McMurdie, J. L. Cope, A. Altrichter, J.-M. Yamal, and E. B. Hollister.** 2017. Leveraging sequence-based faecal microbial community survey data to identify a composite biomarker for colorectal cancer. *Gut*.

29. **Brim, H., S. Yooseph, E. G. Zoetendal, E. Lee, M. Torralbo, A. O. Laiyemo, B. Shokrani, K. Nelson, and H. Ashktorab.** 2013. Microbiome analysis of stool samples from African Americans with colon polyps. *PloS One* **8**:e81352.

30. **Hannigan, G. D., M. B. Duhaime, M. T. Ruffin, C. C. Koumpouras, and P. D. Schloss.** 2017. Diagnostic potential & the interactive dynamics of the colorectal cancer virome. Cold Spring Harbor Laboratory.

31. **Venkataraman, A., J. R. Sieber, A. W. Schmidt, C. Waldron, K. R. Theis, and T. M. Schmidt.** 2016. Variable responses of human microbiomes to dietary supplementation with resistant starch. *Microbiome. Springer Nature* **4**.

32. **Herrmann, E., W. Young, V. Reichert-Grimm, S. Weis, C. Riedel, D. Rosendale, H. Stoklosinski, M. Hunt, and M. Egert.** 2018. In vivo assessment of resistant starch degradation by the caecal microbiota of mice using RNA-based stable isotope probingA proof-of-principle study. *Nutrients. MDPI AG* **10**:179.

33. **Reichardt, N., M. Vollmer, G. Holtrop, F. M. Farquharson, D. Wefers, M. Bunzel,**

581 **S. H. Duncan, J. E. Drew, L. M. Williams, G. Milligan, T. Preston, D. Morrison, H. J.**  
582 **Flint, and P. Louis.** 2017. Specific substrate-driven changes in human faecal microbiota  
583 composition contrast with functional redundancy in short-chain fatty acid production. The  
584 ISME Journal. Springer Nature **12**:610–622.

585 34. **Sze, M. A., N. T. Baxter, M. T. Ruffin, M. A. M. Rogers, and P. D. Schloss.** 2017.  
586 Normalization of the microbiota in patients after treatment for colonic lesions. Microbiome.  
587 Springer Nature **5**.

588 35. **Flynn, K. J., M. T. Ruffin, D. K. Turgeon, and P. D. Schloss.** 2017. Spatial variation  
589 of the native colon microbiota in healthy adults. Cold Spring Harbor Laboratory.

590 36. **Salter, S. J., M. J. Cox, E. M. Turek, S. T. Calus, W. O. Cookson, M. F. Moffatt,**  
591 **P. Turner, J. Parkhill, N. J. Loman, and A. W. Walker.** 2014. Reagent and laboratory  
592 contamination can critically impact sequence-based microbiome analyses. BMC Biology.  
593 Springer Nature **12**.

594 37. **Purcell, R. V., J. Pearson, A. Aitchison, L. Dixon, F. A. Frizelle, and J. I. Keenan.**  
595 2017. Colonization with enterotoxigenic bacteroides fragilis is associated with early-stage  
596 colorectal neoplasia. PLOS ONE. Public Library of Science (PLoS) **12**:e0171602.

597 38. **Sze, M. A., and P. D. Schloss.** 2016. Looking for a signal in the noise: Revisiting  
598 obesity and the microbiome. mBio. American Society for Microbiology **7**:e01018–16.

599 39. **Walters, W. A., Z. Xu, and R. Knight.** 2014. Meta-analyses of human gut microbes  
600 associated with obesity and IBD. FEBS Letters. Wiley-Blackwell **588**:4223–4233.

601 40. **Finucane, M. M., T. J. Sharpton, T. J. Laurent, and K. S. Pollard.** 2014. A taxonomic  
602 signature of obesity in the microbiome? Getting to the guts of the matter. PLoS ONE.

603 Public Library of Science (PLoS) **9**:e84689.

604 41. **Keku, T. O., S. Dulal, A. Deveaux, B. Jovov, and X. Han.** 2015. The gastrointestinal  
605 microbiota and colorectal cancer. *American Journal of Physiology - Gastrointestinal and*  
606 *Liver Physiology* **308**:G351–G363.

607 42. **Vogtmann, E., and J. J. Goedert.** 2016. Epidemiologic studies of the human  
608 microbiome and cancer. *British Journal of Cancer* **114**:237–242.

609 43. **Kostic, A. D., D. Gevers, C. S. Pedomallu, M. Michaud, F. Duke, A. M. Earl, A. I.**  
610 **Ojesina, J. Jung, A. J. Bass, J. Tabernero, J. Baselga, C. Liu, R. A. Shivdasani, S.**  
611 **Ogino, B. W. Birren, C. Huttenhower, W. S. Garrett, and M. Meyerson.** 2012. Genomic  
612 analysis identifies association of *Fusobacterium* with colorectal carcinoma. *Genome*  
613 *Research* **22**:292–298.

614 44. **Zackular, J. P., M. A. M. Rogers, M. T. Ruffin, and P. D. Schloss.** 2014. The human  
615 gut microbiome as a screening tool for colorectal cancer. *Cancer Prevention Research*  
616 *(Philadelphia, Pa.)* **7**:1112–1121.

617 45. **Schloss, P. D., S. L. Westcott, T. Ryabin, J. R. Hall, M. Hartmann, E. B. Hollister,**  
618 **R. A. Lesniewski, B. B. Oakley, D. H. Parks, C. J. Robinson, J. W. Sahl, B. Stres, G.**  
619 **G. Thallinger, D. J. Van Horn, and C. F. Weber.** 2009. Introducing mothur: Open-Source,  
620 Platform-Independent, Community-Supported Software for Describing and Comparing  
621 Microbial Communities. *Appl. Environ. Microbiol.* **75**:7537–7541.

622 46. **Rognes, T., T. Flouri, B. Nichols, C. Quince, and F. Mahé.** 2016. VSEARCH: A  
623 versatile open source tool for metagenomics. *PeerJ* **4**:e2584.

624 47. **Westcott, S. L., and P. D. Schloss.** 2017. OptiClust, an Improved Method for

Assigning Amplicon-Based Sequence Data to Operational Taxonomic Units. *mSphere* **2**.

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

49. **R Core Team.** 2017. R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria.

50. **Mangiafico, S.** 2017. Rcompanion: Functions to support extension education program evaluation.

51. **Fox, J., and S. Weisberg.** 2011. An R companion to applied regressionSecond. Sage, Thousand Oaks CA.

52. **Bates, D., M. Mächler, B. Bolker, and S. Walker.** 2015. Fitting linear mixed-effects models using lme4. *Journal of Statistical Software* **67**:1–48.

53. **Telmo Nunes, M. S. with contributions from, C. Heuer, J. Marshall, J. Sanchez, R. Thornton, J. Reiczigel, J. Robison-Cox, P. Sebastiani, P. Solymos, K. Yoshida, G. Jones, S. Pirikahu, S. Firestone, and R. Kyle.** 2017. EpiR: Tools for the analysis of epidemiological data.

54. **Viechtbauer, W.** 2010. Conducting meta-analyses in R with the metafor package. *Journal of Statistical Software* **36**:1–48.

55. **Oksanen, J., F. G. Blanchet, M. Friendly, R. Kindt, P. Legendre, D. McGlinn, P. R. Minchin, R. B. O'Hara, G. L. Simpson, P. Solymos, M. H. H. Stevens, E. Szoecs, and H. Wagner.** 2017. Vegan: Community ecology package.

56. **Jed Wing, M. K. C. from, S. Weston, A. Williams, C. Keefer, A. Engelhardt, T. Cooper, Z. Mayer, B. Kenkel, the R Core Team, M. Benesty, R. Lescarbeau, A. Ziem,**

- 648 **L. Scrucca, Y. Tang, C. Candan, and T. Hunt.** 2017. Caret: Classification and regression  
649 training.
- 650 57. **Liaw, A., and M. Wiener.** 2002. Classification and regression by randomForest. R  
651 News **2**:18–22.
- 652 58. **Wickham, H.** 2009. Ggplot2: Elegant graphics for data analysis. Springer-Verlag New  
653 York.
- 654 59. **Auguie, B.** 2017. GridExtra: Miscellaneous functions for “grid” graphics.



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

Study	Data Stored	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	Region	Control (n)	Adenoma (n)	Carcinoma (n)
Burns	SRA	V5-6	18	0	16
Chen	SRA	V1-3	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: Significant Bacterial Community Metrics for Adenoma or Carcinoma in Stool.** A) Adenoma evenness. B) Carcinoma evenness. C) Carcinoma Shannon diversity. Blue represents controls and red represents either adenoma (panel A) or carcinoma (panel B and C). The black lines represent the median value for each respective group.

**Figure 2: 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 3: The AUC of Individual Significant OR Taxa to classify Carcinoma.** A) Stool samples. B) Unmatched tissue samples. The larger circle represents the median AUC of all studies and the smaller circles represent the individual AUC for a particular study. The dotted line denotes an AUC of 0.5.

**Figure 4: 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 5: Most Important Members in Significant OR Taxa Carcinoma Models.** A) Common taxa in the top 10 percent for carcinoma Random Forest stool-based models. B) Common taxa in the top 10 percent for carcinoma Random Forest unmatched tissue-based models. Blue represents less important and red represents more important to the Random Forest Model. White means that particular taxa was not in the top 10%.

**Figure 6: 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 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: Most Common Taxa Across Carcinoma Full Community Stool Study Models.** A) Common taxa in the top 10 percent for carcinoma Random Forest all taxa-based models. B) Common taxa in the top 10 percent for carcinoma Random Forest all OTU-based models. Blue represents less important and red represents more important to the Random Forest Model. White means that particular taxa was not in the top 10%.

**Figure S3: Most Common Genera Across Full Community Tissue Study Models.** A) Common genera in the top 10 percent for matched carcinoma Random Forest all genera-based models. B) Common genera in the top 10 percent for unmatched carcinoma Random Forest all genera-based models. C) Common genera in the top 10 percent for matched carcinoma Random Forest all OTU-based models. D) Common genera in the top 10 percent for unmatched carcinoma Random Forest all OTU-based models. Blue represents less important and red represents more important to the Random Forest Model. White means that particular taxa was not in the top 10%.

**Figure S4: 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 S5: 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.