

Investigating the microbial community of the human gut in colorectal cancer

Nicole Bowers¹, Brittany Hicks², Christina Kang-Yun², Katelyn Polemi³, Kelly Sovacool¹

¹Department of Computational Medicine and Bioinformatics

²Department of Civil and Environmental Engineering

³Department of Environmental Health Sciences

Contents

Abstract	2
Introduction	2
Methods	4
Study design and data collection	4
16S rRNA gene sequence processing	4
Metagenome and virome quality control	4
Virome assembly	5
Classification modeling	5
Metagenome taxonomic profiling and gene annotation	5
Results	6
Bacterial diversity in colorectal cancer	6
Classification modeling discriminates disease states	7
Taxonomic profiles of metagenomes	7
Metagenome annotation and quantification	8
Discussion	9
Acknowledgements and project contributions	11
References	12

Abstract

Colorectal cancer (CRC) is the third leading cause of cancer related deaths globally (1). CRC initiates in the large intestine emerging from glandular epithelial cells. Due to a selective advantage, obtained from a series of genetic or epigenetic mutations, these unregulated cells grow and can potentially develop into colorectal adenocarcinoma (2). Many studies have used the cancer microbiome to investigate how these mutations occur, however, they have almost exclusively focused on bacteria. Human viruses have been implicated in the development of many cancers such as HPV, HTLV-1 and HIV (3).

Through mutagenic and manipulative abilities, viruses cause unregulated growth by disrupting the normal function of cells. In addition, some bacteriophages have been shown to influence cancer processes through the immunological response (4). The association of CRC and the gut virome remains unknown, therefore, in this study we propose to investigate the differences in the bacterial and viral community composition and their effect on CRC development.

We performed operational taxonomic unit (OTU) clustering of 16S sequences for bacterial taxonomy and operational viromic unit (OVU) clustering of metaviromes for viral taxonomy to investigate the differences in human colorectal bacterial and viral community composition. Using machine learning classification of OTUs and OVUs from healthy or cancerous samples, we aim to compare performance of the random forest model vs logistic regression. In addition, gene abundance was calculated from metagenomes in healthy vs adenoma or CRC patients and were mapped to pathways to identify key genes and/or pathways that may have a role in CRC. Overall, the goal of the current study is to classify whether samples are healthy or cancerous based on bacterial and viral taxa. These results provide evidence that the gut microbiome of patients with CRC differ from healthy patients.

Introduction

The leading cancer-related death in the United States is colorectal cancer (CRC) (5). Due to screening techniques and improvements in treatment the rate of CRC has fallen. The primary screening technique is colonoscopy. A new technique gaining popularity is Exact Sciences Cologuard, a combination test which looks for both the genetic material found in colon cancer and some more advanced colon polyps and it detects hemoglobin in the stool using fecal immunochemical testing (Cologuard). While colon cancer screening techniques are advancing, we aim to understand the environment of CRC.

Gut dysbiosis is any change to the normal microbiome that could change the relationship between the host and associated microbes (6). This is well documented gastrointestinal (GI) disorders including irritable bowel syndrome, peptic ulcers, and even gastric and colon cancers (6). Gut dysbiosis is documented in patients with CRC (7). Though the underlying mechanisms of these

microbes is yet unknown, some bacteria were found to be associated with CRC (8). These include *Fusobacterium*, *Peptostreptococcus*, *Porphyromonas*, *Prevotella*, *Parvimonas*, *Bacteroides*, and *Gemella*, and they may be able to affect the disease state via metabolite secretion, host tissue invasion, and inducing host immune response (8).

At present, the majority of cancer microbiome studies focus on the bacteria (6). CRC was previously linked to changes to the colonic bacterial composition, yet the gut virome is vastly unexplored (7).

Virome is a viral community associated with a particular ecosystem or holobiont. In mammals, the viruses infect host cells and the variety of organisms that inhabit us. Identifying and understanding the virome within a host requires genetic and transcriptional identification of the mammal. This virome may be able to shed light on the host's genetics in states of health and states of disease. The virome is composed of both nucleic acids (DNA and RNA).

Human viruses are associated with many cancers, due to their manipulative and mutagenic abilities (5). Viral metagenomic taxa can distinguish CRC patients from control subjects (7). Furthermore, a subset of markers were identified as high-risk from a subset of patients with CRC.

The composition of the microbiome is further affected by the presence and population of bacteriophages (9). Hsu *et al.* (9) introduces lytic phages colonized with human gut bacteria. The phages showed direct effects, that phage predation decreases the susceptible bacterial species, while resistant populations flourish. Indirect effects include cascading effects on other bacterial species through interbacterial interactions. Shifts in the microbiome caused by phage predation have a direct effect on the gut metabolome (9).

Here we address the knowledge gap with respect to whether bacterial and viral community composition differ between healthy patients and those diagnosed with adenoma and carcinoma.

We used the data from Hannigan *et al.* (5) and NCBI BioProject ID: PRJA389927 to try to understand the differences in the bacterial and viral community composition between human stool samples from patients diagnosed as normal, adenoma, and carcinoma. We also try to illuminate if there differences in bacterial & viral community composition between human stool samples from patients of the three groups. In order to do achieve these objectives, we use the same approach as Hannigan *et al.* (5) with a random forest model with OTUs or OVUs as features. We also seek to understand differences in functional potential between metagenomes of the three samples groups by examining the relative gene abundance and pathways. This would be demonstrated by differences in OTU levels between the three conditions: cancer, adenoma, and healthy control.

Methods

Study design and data collection

The data used for our evaluation have been previously reported (5, 10). Whole bacterial metagenomes and viromes were isolated from whole evacuated stool samples collected from patients. Patient requirements included a minimum age of 18 years of age, informed consent, a histologically confirmed colonic disease status from a previous colonoscopy, no history of surgery or chemotherapy, and free of known comorbidities. These comorbidities included chronic viral hepatitis, inflammatory bowel disease, hereditary nonpolyposis colorectal cancer (HNPCC), HIV, and familial adenomatous polyposis (FAP). Ninety patients were recruited to the study. Three disease classes were designated: adenoma, carcinoma, and healthy. Thirty patients were classified into each of these classes. Samples were collected from four geographic locations: Ann Arbor (MI, USA), Boston (MA, USA), Houston (TX, USA), and Toronto (Ontario, Canada).

The data are available in the NCBI Sequence Read Archive as [PRJNA389927](https://www.ncbi.nlm.nih.gov/sra/PRJNA389927). All code used in preparation of this report is available in the following GitHub repository: <https://github.com/kelly-sovacool/bioinf545-group3-project>.

16S rRNA gene sequence processing

The 16S rRNA gene sequences associated with this study were previously reported by (10). The fastq sequence and metadata files were downloaded from the NCBI sequence read archive with the sra toolkit. The 16S rRNA gene sequences were analyzed with the mothur software package (v1.43.0) (11) according to the MiSeq standard operating procedure (12). Mothur allows contigs from forward and reverse reads to be assembled, pre-processing steps for primer removal and quality screening, and the identification of unique sequences. After these steps were taken to reduce sequencing and PCR errors, sequences were aligned to the SILVA database (13), screened for chimeras using VSEARCH, and clustered into operational taxonomic units (OTU) using the *de novo* OptiClust algorithm (14) at a 97% sequence similarity threshold. Alpha and beta diversity metrics were calculated with the vegan R package (15).

Metagenome and virome quality control

Trimmomatic (v.0.39) (16) was used to remove adapter sequences and low-quality reads in metagenome and virome sample sets. The read quality was assessed using FastQC (v.0.11.9) (17) before and after adapter trimming to confirm removal of adapters and low quality reads. Unpaired reads were dropped using the repair function in BBTools (v.37.62) (18). Then, reads mapping to the human GRCh38 reference genome were removed using the BWA-MEM algorithm (BWA v.0.7.12) (19). The paired unmapped reads from metagenome samples were used for taxonomic profiling and gene annotation, while those from virome samples were used for virome assembly.

Virome assembly

After the quality control steps, viral reads were assembled into contigs using megahit (20), then reads were mapped back to the assembled contigs using BWA-MEM. Next, contigs were binned using the CONCOCT algorithm (v0.4.0) (21), which applies a variation-based Bayesian approach to bin related contigs by similar tetramer and coabundance profiles within samples. These bins defined operational viral units (OVU) in the absence of taxonomic identity. The sample-wise table of OVU abundances derived from the CONCOCT output was used for classification modeling of disease states. Due to a lack of time, the OVUs were not aligned to a virome reference database for taxonomic identification.

Classification modeling

The caret R package (22) was used to train random forest models to classify samples as healthy or cancerous based on taxonomic abundance data. Only binary classifier models were trained, rather than three-class models, to reduce the time and computing resources required. Two models were trained; one with OTU (bacterial) abundances as features and the other with OVU (viral) abundances as features. For both models, the data were split into training and test sets with 65% of the data for training. The models were validated with five-fold cross validation and tuned for mtry values to maximize the area under the receiver operating curve (AUC).

Metagenome taxonomic profiling and gene annotation

The human genome-free reads were used to profile the metagenomic taxonomy at the species, genus, family, and phylum levels using MetaPhlAn2 (23). The results were visualized using the R package ggplot2 (v.3.3.0) (24).

The paired reads that did not map to the human genome were aligned to the Integrated Gene Catalog (IGC) (25) of 1,267 gut microbiome samples consisting of approximately 10 million genes with the BWA MEM algorithm for metagenome annotation (19). Annotated genes were extracted from the alignment results using functions geneList and countKegg modified from the MGS-Fast pipeline (26). The differences in gene abundance between healthy, adenoma, and cancer groups were assessed using the R package edgeR (27). Up to 480 top KEGG numbers of genes that were significantly different between healthy and other groups were selected and the relevant KEGG pathway was determined using the KEGG mapper tool (28).

Results

Bacterial diversity in colorectal cancer

The influence of CRC on bacterial diversity was evaluated amongst the three disease classes. A total of 3,904 OTUs were identified. Alpha and beta diversity metrics are used to characterize and distinguish ecological communities in broad strokes. As a metric of alpha diversity, the Shannon diversity index describes diversity within communities in terms of richness and evenness, taking into account the abundance of the species present (29). As seen in Fig. 1A, individuals with healthy colons do not appear to have a distinctly different gut microbiome from a species richness and evenness perspective. Although the mean Shannon diversity decreased along the progression from healthy, to adenoma, to cancer (Fig. 1A), these perceived differences were not statistically significant among the disease groups (analysis of variance [ANOVA] p -value = 0.3). This result suggests that individuals with cancerous and healthy colons may have a select group of OTUs that dominate the environment since their index values are not relatively high.

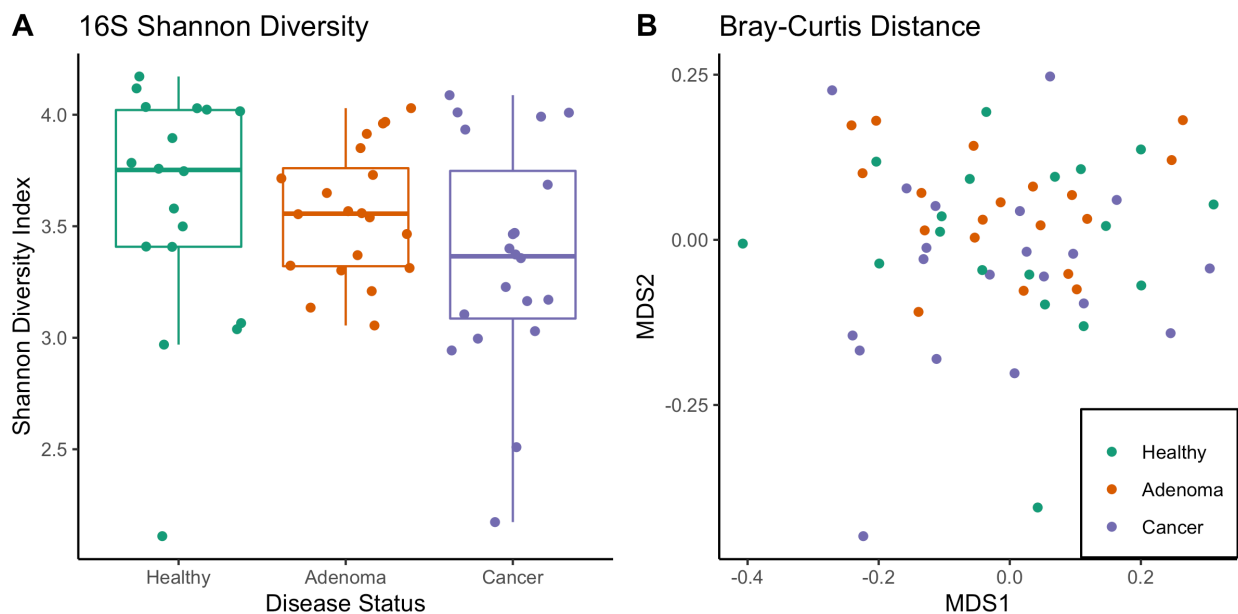


Figure 1: (A) Shannon diversity index across disease states and (B) multi-dimensional scaling plot of Bray-Curtis Distances of pairwise OTU abundance.

In contrast to Shannon diversity, which quantifies diversity within communities, Bray-Curtis dissimilarity is a beta diversity metric that quantifies diversity *between* communities (29). The pairwise Bray-Curtis distances of sample-wise OTU abundances were calculated and plotted as a multi-dimensional scaling plot (Fig. 1B). There does not appear to be distinct clustering of the disease classes. Instead, each disease class is fairly dispersed. There was no statistically significant clustering of the disease groups ([ANOVA] p -value > 0.23).

Classification modeling discriminates disease states

Random forest models were trained on OTU abundances from 16S bacterial sequence data and on OVU abundances from viral metagenomes to classify samples as healthy or cancerous. Both models had modest performance with an AUROC of 0.765 for the bacterial model and an AUROC of 0.73 for the viral model (Fig. 2). The bacterial abundance data consisted of 3,704 OTUs, while the viral abundance data contained only 70 OVUs. This discrepancy in the number of features as inputs to the random forest models may have contributed to the bacterial model outperforming the viral model for classifying samples as healthy or cancerous.

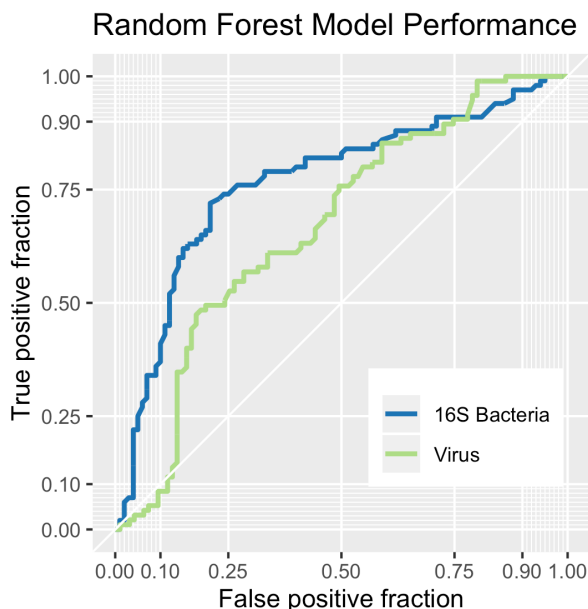


Figure 2: Model performance as measured by the receiver operating characteristic (ROC) curve for models trained on bacterial abundance or viral abundance data.

An OTU containing members of the genus *Fusobacterium* was by far the most important feature of the bacterial model. The viral metagenomes were not aligned to a reference database, so OVUs were not able to be taxonomically identified for feature importance determination. Previous work by Hannigan *et al.* additionally included models trained on bacterial whole metagenomes and models combining bacterial and viral abundance data.

Taxonomic profiles of metagenomes

The taxonomic profiles were constructed based on the bacterial metagenomes extracted from the gut microbiomes of healthy, adenoma, and cancer patients using MetaPhlAn2 (Fig. 3). There are subtle differences of community composition between healthy patients and those who were diagnosed with adenoma or cancer.

At the phylum level, some healthy patients have more members of the phyla Verrucomicrobia, Bacteroidetes, and Actinobacteria as compared to other patient groups. In addition, some patients with adenoma or cancer have a greater abundance of species belonging to the Proteobacteria phylum. The gut microbiome of healthy patients are relatively more abundant with members of the families *Ruminococcaceae*, *Lachnospiraceae*, *Bifidobacteriaceae*, and *Akkermansiaceae*. Conversely, the gut microbiome of patients with adenoma or cancer seem to have greater abundance of members of the families *Enterobacteriaceae* and *Coriobacteriaceae*.

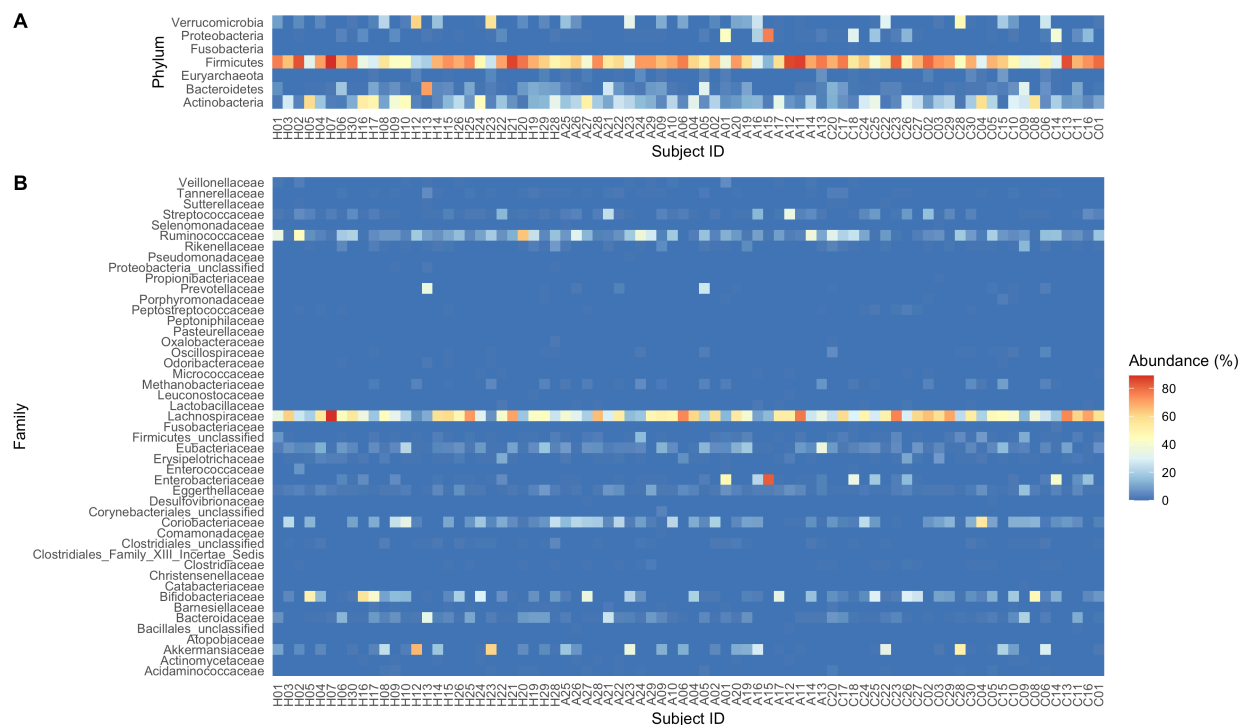


Figure 3: Taxonomic profile constructed using MetaPhlAn2 at the (A) phylum and (B) family level based on metagenome of gut microbiome of patients that are healthy, H, or diagnosed with adenoma, A, or cancer, C.

Metagenome annotation and quantification

The metagenomes isolated from stool samples of patients that were healthy or diagnosed with adenoma or cancer were aligned to the IGC database to annotate the genes with KEGG numbers. KEGG number counts were calculated for all samples, then were used to determine whether specific genes were more or less abundant in patients with adenoma or cancer as compared to those who are healthy using edgeR. There were seven samples including the negative control with extremely low read counts which were excluded from subsequent analyses, as they strongly skewed the results. Multidimensional scaling (MDS) analysis of the KEGG number counts between different groups did not show any distinct clustering within groups (Fig. 4).

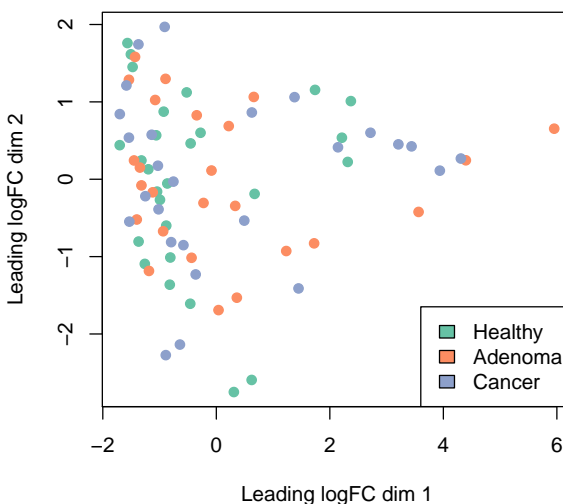


Figure 4: Multidimensional scaling (MDS) abundance plot based on log2 fold-change (logFC) of KEGG annotation counts derived from gut microbiome of patients that are healthy, or diagnosed with adenoma or carcinoma.

In the subsequent differential gene abundance analysis, the abundances of 6,104 distinct genes were compared between healthy vs adenoma and healthy vs cancer patients. In both cases, there were approximately 130 genes that were less abundant in adenoma and cancer patients, whereas around 600 to 1,000 genes were more abundant compared to healthy patients. Differentially abundant genes that were found in both adenoma and cancer patient groups were selected and mapped to metabolic pathways using the online KEGG mapper tool (see Table 1).

The differentially-abundant genes found in adenoma and cancer patients mapped to approximately 20 and 550 genes in the KEGG pathways, respectively. The number of genes mapped may be different from the query list due to lack of mapping or genes involved in multiple pathways. Less abundant genes in adenoma and cancer patients were slightly enriched in pathways in metabolism, genetic information processing, signal transduction, transport and catabolism, immune system, and bacterial infectious disease. Contrastingly, more abundant genes in adenoma and cancer patients were most enriched in pathways in metabolism, environmental information processing, and cellular processes amongst other pathways. Though human genome was filtered out through pre-processing, pathways pertaining to human diseases and organismal systems were also enriched in the case of genes more abundant in adenoma and cancer patients.

Discussion

Some differences in the gut microbiome community composition between the patient groups may reflect prior findings. Members of the *Peptostreptococaceae* and *Porphyromonadaceae* family are slightly more abundant in adenoma or cancer patients, which were previously found to be associated with CRC (8). In addition, the enrichment of species included in the *Lachnospiraceae* family in the healthy patients also support prior studies indicating lower abundance of these butyrate-producing bacteria in CRC patients (30).

Moreover, higher abundance of Proteobacteria and lower abundance of Bacteroidetes in some patients with adenoma or carcinoma as compared to healthy patients in this study also demonstrate previous findings by Shen *et al.* (31). Gene abundance analysis using the metagenomes from patients indicated enrichment in genes and pathways previously shown to be associated with GI diseases including CRC. Genes K07345 and K07347 corresponding to fimbriae genes *fimA*, *fimD* that are associated with inflammatory bowel disease were found to be enriched in patients with adenoma or cancer (33). Virulence genes, specifically bacterial secretion systems, two-component systems, bacterial flagellar assembly, and biofilm formation, each part of membrane transport, signal transduction, bacterial cell motility, and prokaryotic cellular community pathways, respectively, were highly enriched in patients with adenoma or cancer. This is in line with previous studies, where virulence genes were found to be closely associated with colorectal tumor microenvironments (32).

The alpha and beta diversity of bacteria of the gut microbiome did not significantly shift in response to the presence of CRD. Standard alpha and beta diversity metrics did not adequately capture

Table 1: Genes that are more abundant in adenoma or cancer patients as compared to healthy patients.

KEGG Pathway		Number of Genes
Metabolism		
Global and overview maps	Metabolic pathways, etc.	193
Carbohydrate metabolism	Glycolysis / Gluconeogenesis, etc.	40
Energy metabolism	Oxidative phosphorylation, etc.	15
Lipid metabolism	Fatty acid biosynthesis, etc.	17
Nucleotide metabolism	Purine metabolism, etc.	4
Amino acid metabolism	Alanine, aspartate and glutamate metabolism, etc.	50
Metabolism of other amino acids	beta-Alanine metabolism, etc.	14
Glycan biosynthesis and metabolism	Lipopolysaccharide biosynthesis, etc.	8
Metabolism of cofactors and vitamins	Riboflavin metabolism, etc.	17
Metabolism of terpenoids and polyketides	Limonene and pinene degradation, etc.	7
Biosynthesis of other secondary metabolites	Isoquinoline alkaloid biosynthesis, etc.	6
Xenobiotics biodegradation and metabolism	Benzoate degradation, etc.	20
Genetic Information Processing		
Transcription	RNA polymerase	1
Folding, sorting and degradation	Protein export, etc.	2
Replication and repair	DNA replication, etc.	8
Environmental Information Processing		
Membrane transport	ABC transporters, etc.	47
Signal transduction	Two-component system, etc.	35
Cellular Processes		
Cell growth and death	Cell cycle - yeast	1
Cellular community - prokaryotes	Quorum sensing, etc.	27
Cell motility	Bacterial chemotaxis, etc.	12
Organismal Systems		
Immune system	Complement and coagulation cascades	1
Aging	Longevity regulating pathway - worm	1
Human Diseases		
Cancer: overview	Pathways in cancer, etc.	3
Cancer: specific types	Hepatocellular carcinoma, etc.	2
Cardiovascular disease	Fluid shear stress and atherosclerosis	1
Infectious disease: bacterial	Epithelial cell signaling in Helicobacter pylori infection, etc.	10
Drug resistance: antimicrobial	beta-Lactam resistance, etc.	3
Drug resistance: antineoplastic	Platinum drug resistance	1

bacterial community differences of the healthy and cancerous gut. This illuminates that examining diversity metrics alone may not be an adequate tool for identifying a healthy versus a cancerous colon. While traditional diversity metrics could not discriminate between disease states, random forest models trained on bacterial or viral taxa abundances performed modestly well at this task. This concurs with previous studies (5, 10) of this CRC microbiome dataset. An explanation for better discrimination of machine learning models is that diversity metrics characterize communities in broad strokes, but models using abundance data take into account more granular similarities, differences, and perhaps interactions within communities of interest.

Kostic *et al.* (34) have previously shown that members of the genus *Fusobacterium* are associated with CRC. The results of our study also support this finding, as an OTU containing members of this genus was the most important feature of the bacterial classification model. Interestingly, this phylum was not the most abundant in any samples. The most important members of communities may not always be the most abundant, but their metabolic activities may mediate important interactions that cannot be investigated through genetic sequence analysis alone (35).

Hannigan *et al.* achieved better performance compared to those achieved here, with their OTU and OVU-based random forest models respectively yielding AUCs of 0.809 and 0.792. Minor differences between methods of processing 16S sequence data, such as an updated OTU clustering algorithm, may account for this difference in OTU model performance. Hannigan *et al.* also implemented a more robust modeling pipeline which repeated the data splitting step for 20 iterations with different random seeds; due to time constraints we were unable to implement this level of robustness here.

Overall, the taxonomic and gene abundance analyses gave insight into enriched members that have previously established association with CRC. In addition, alpha and beta diversity metrics were not sufficient to identify CRC in patients, but the random forest model performed better in classifying disease states. Interesting analyses that were not included in this study due to time constraints are taxonomic and gene co-abundance analysis, gut microbiome network analysis to identify hubs important in cancer disease state, aligning OVUs to a virome reference for taxonomic identification, and a more thorough machine learning pipeline to select the best model for the questions at hand. A more stringent statistical analysis of taxonomic and gene abundance results would provide a measure to identify key members that are important in CRC.

Acknowledgements and project contributions

Author order was determined by alphabetizing by last name. Brittany Hicks processed 16S sequencing data, collaborated for processing virome data, and assisted in preparing the methods, results, and discussion for these analyses. Christina Kang-Yun constructed the workflow for preprocessing metagenome data, taxonomic profiling, and gene annotation and prepared the methods, results, and discussion for these analyses. Katelyn Polemi and Nicole Bowers worked together to

create figures for analyses including alpha and beta diversity plots, gene abundance MDS, taxa bar plots, OTU heatmaps (some plots not included in final paper, see GitHub), prepared the abstract and introduction for the paper, and contributed to the conclusions. Kelly Sovacool wrote the virome assembly pipeline, wrote and performed the classification modeling steps, collaborated in developing the metagenome quality control steps, integrated others' methods into a cohesive workflow, helped others troubleshoot at various steps, and participated in writing and editing the subsections of this report associated with the 16S data processing, virome pipeline, quality control, and classification modeling analyses.

References

1. **Rawla P, Sunkara T, Barsouk A.** 2019. Epidemiology of colorectal cancer: Incidence, mortality, survival, and risk factors. *Gastroenterology Rev* **14**:89–103. doi:[10.5114/pg.2018.81072](https://doi.org/10.5114/pg.2018.81072).
2. **Ewing I, Hurley JJ, Josephides E, Millar A.** 2014. The molecular genetics of colorectal cancer. *Frontline Gastroenterology* **5**:26–30. doi:[10.1136/flgastro-2013-100329](https://doi.org/10.1136/flgastro-2013-100329).
3. **Liao JB.** 2006. Viruses and human cancer. *Yale J Biol Med* **79**:115–122.
4. **Budynek P, Dąbrowska K, Skaradziński G, Górski A.** 2010. Bacteriophages and cancer. *Arch Microbiol* **192**:315–320. doi:[10.1007/s00203-010-0559-7](https://doi.org/10.1007/s00203-010-0559-7).
5. **Hannigan GD, Duhaime MB, Ruffin MT, Koumpouras CC, Schloss PD.** 2018. Diagnostic Potential and Interactive Dynamics of the Colorectal Cancer Virome. *mBio* **9**. doi:[10.1128/mBio.02248-18](https://doi.org/10.1128/mBio.02248-18).
6. **Neto AG, Hickman RA, Khan A, Nossa C, Pei Z.** 2017. The Upper Gastrointestinal Tract—Esophagus and Stomach, pp. 1–11. *In* *The Microbiota in Gastrointestinal Pathophysiology*. Elsevier.
7. **Nakatsu G, Zhou H, Wu WKK, Wong SH, Coker OO, Dai Z, Li X, Szeto C-H, Sugimura N, Lam TY-T, Yu AC-S, Wang X, Chen Z, Wong MC-S, Ng SC, Chan MTV, Chan PKS, Chan FKL, Sung JJ-Y, Yu J.** 2018. Alterations in Enteric Virome Are Associated With Colorectal Cancer and Survival Outcomes. *Gastroenterology* **155**:529–541.e5. doi:[10.1053/j.gastro.2018.04.018](https://doi.org/10.1053/j.gastro.2018.04.018).
8. **Ternes D, Karta J, Tsenkova M, Wilmes P, Haan S, Letellier E.** 2020. Microbiome in Colorectal Cancer: How to Get from Meta-omics to Mechanism? *Trends in Microbiology* **28**:401–423. doi:[10.1016/j.tim.2020.01.001](https://doi.org/10.1016/j.tim.2020.01.001).
9. **Hsu BB, Gibson TE, Yeliseyev V, Liu Q, Lyon L, Bry L, Silver PA, Gerber GK.** 2019. Dynamic Modulation of the Gut Microbiota and Metabolome by Bacteriophages in a Mouse Model. *Cell Host & Microbe* **25**:803–814.e5. doi:[10.1016/j.chom.2019.05.001](https://doi.org/10.1016/j.chom.2019.05.001).
10. **Zackular JP, Rogers MAM, Ruffin MT, Schloss PD.** 2014. The Human Gut Microbiome as a Screening Tool for Colorectal Cancer. *Cancer Prev Res* **7**:1112–1121. doi:[10.1158/1940-6207.CAPR-14-0129](https://doi.org/10.1158/1940-6207.CAPR-14-0129).
11. **Schloss PD, Westcott SL, Ryabin T, Hall JR, Hartmann M, Hollister EB, Lesniewski RA, Oakley BB, Parks DH, Robinson CJ, Sahl JW, Stres B, Thallinger GG, Van Horn DJ, Weber CF.** 2009. Introducing mothur: Open-source, platform-independent, community-supported software for describing and comparing

microbial communities. *Applied and Environmental Microbiology* **75**:7537–7541. doi:[10.1128/AEM.01541-09](https://doi.org/10.1128/AEM.01541-09).

12. **Kozich JJ, Westcott SL, Baxter NT, Highlander SK, Schloss PD.** 2013. Development of a Dual-Index Sequencing Strategy and Curation Pipeline for Analyzing Amplicon Sequence Data on the MiSeq Illumina Sequencing Platform. *Appl Environ Microbiol* **79**:5112. doi:[10.1128/AEM.01043-13](https://doi.org/10.1128/AEM.01043-13).

13. **Pruesse E, Quast C, Knittel K, Fuchs BM, Ludwig W, Peplies J, Glöckner FO.** 2007. SILVA: A comprehensive online resource for quality checked and aligned ribosomal RNA sequence data compatible with ARB. *Nucleic Acids Res* **35**:7188–7196. doi:[10.1093/nar/gkm864](https://doi.org/10.1093/nar/gkm864).

14. **Westcott SL, Schloss PD.** 2017. OptiClust, an Improved Method for Assigning Amplicon-Based Sequence Data to Operational Taxonomic Units. *mSphere* **2**:e00073–17. doi:[10.1128/mSphereDirect.00073-17](https://doi.org/10.1128/mSphereDirect.00073-17).

15. **Oksanen J.** 2018. *Vegan: Community ecology package*.

16. **Bolger AM, Lohse M, Usadel B.** 2014. Trimmomatic: A flexible trimmer for Illumina sequence data. *Bioinformatics* **30**:2114–2120. doi:[10.1093/bioinformatics/btu170](https://doi.org/10.1093/bioinformatics/btu170).

17. **Andrews S, Krueger F, Segonds-Pichon A, Biggins L, Krueger C, Wingett S.** 2010. FastQC.

18. **Bushnell B.** 2018. BBTools: A suite of fast, multithreaded bioinformatics tools designed for analysis of DNA and RNA sequence data.

19. **Li H.** 2013. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.

20. **Li D, Liu C-M, Luo R, Sadakane K, Lam T-W.** 2015. MEGAHIT: An ultra-fast single-node solution for large and complex metagenomics assembly via succinct de Bruijn graph. *Bioinformatics* **31**:1674–1676. doi:[10.1093/bioinformatics/btv033](https://doi.org/10.1093/bioinformatics/btv033).

21. **Alneberg J, Bjarnason BS, Bruijn I de, Schirmer M, Quick J, Ijaz UZ, Lahti L, Loman NJ, Andersson AF, Quince C.** 2014. Binning metagenomic contigs by coverage and composition. *Nature Methods* **11**:1144–1146. doi:[10.1038/nmeth.3103](https://doi.org/10.1038/nmeth.3103).

22. **Kuhn M.** 2013. *Caret: Classification and regression training* **1**.

23. **Truong DT, Franzosa EA, Tickle TL, Scholz M, Weingart G, Pasolli E, Tett A, Huttenhower C, Segata N.** 2015. MetaPhlAn2 for enhanced metagenomic taxonomic profiling. *Nat Methods* **12**:902–903. doi:[10.1038/nmeth.3589](https://doi.org/10.1038/nmeth.3589).

24. **Wickham H.** 2016. *Ggplot2: Elegant graphics for data analysis* Second edition. Springer, Cham.

25. **Li J, Jia H, Cai X, Zhong H, Feng Q, Sunagawa S, Arumugam M, Kultima JR, Prifti E, Nielsen T, Juncker AS, Manichanh C, Chen B, Zhang W, Levenez F, Wang J, Xu X, Xiao L, Liang S, Zhang D, Zhang Z, Chen W, Zhao H, Al-Aama JY, Edris S, Yang H, Wang J, Hansen T, Nielsen HB, Brunak S, Kristiansen K, Guarner F, Pedersen O, Doré J, Ehrlich SD, Bork P, Wang J.** 2014. An integrated catalog of reference genes in the human gut microbiome. *Nat Biotechnol* **32**:834–841. doi:[10.1038/nbt.2942](https://doi.org/10.1038/nbt.2942).

26. **Brown SM, Chen H, Hao Y, Laungani BP, Ali TA, Dong C, Lijeron C, Kim B, Wultsch C, Pei Z, Krampis K.** 2019. MGS-Fast: Metagenomic shotgun data fast annotation using microbial gene catalogs.

Gigascience **8**. doi:[10.1093/gigascience/giz020](https://doi.org/10.1093/gigascience/giz020).

27. **Robinson MD, McCarthy DJ, Smyth GK**. 2009. edgeR: A Bioconductor package for differential expression analysis of digital gene expression data. *Bioinformatics* **26**:139–140. doi:[10.1093/bioinformatics/btp616](https://doi.org/10.1093/bioinformatics/btp616).
28. **Kanehisa M, Sato Y**. 2020. KEGG Mapper for inferring cellular functions from protein sequences. *Protein Science* **29**:28–35. doi:[10.1002/pro.3711](https://doi.org/10.1002/pro.3711).
29. **Tuomisto H**. 2010. A diversity of beta diversities: Straightening up a concept gone awry. Part 1. Defining beta diversity as a function of alpha and gamma diversity. *Ecography* **33**:2–22. doi:[10.1111/j.1600-0587.2009.05880.x](https://doi.org/10.1111/j.1600-0587.2009.05880.x).
30. **Rinninella E, Raoul P, Cintoni M, Franceschi F, Miggiano GAD, Gasbarrini A, Mele MC**. 2019. What is the Healthy Gut Microbiota Composition? A Changing Ecosystem across Age, Environment, Diet, and Diseases. *Microorganisms* **7**. doi:[10.3390/microorganisms7010014](https://doi.org/10.3390/microorganisms7010014).
31. **Shen XJ, Rawls JF, Randall TA, Burcall L, Mpande C, Jenkins N, Jovov B, Abdo Z, Sandler RS, Keku TO**. 2010. Molecular characterization of mucosal adherent bacteria and associations with colorectal adenomas. *Gut Microbes* **1**:138–147. doi:[10.4161/gmic.1.3.12360](https://doi.org/10.4161/gmic.1.3.12360).
32. **Burns MB, Lynch J, Starr TK, Knights D, Blekhman R**. 2015. Virulence genes are a signature of the microbiome in the colorectal tumor microenvironment. *Genome Medicine* **7**:55. doi:[10.1186/s13073-015-0177-8](https://doi.org/10.1186/s13073-015-0177-8).
33. **Minot SS, Willis AD**. 2019. Clustering co-abundant genes identifies components of the gut microbiome that are reproducibly associated with colorectal cancer and inflammatory bowel disease. *Microbiome* **7**:110. doi:[10.1186/s40168-019-0722-6](https://doi.org/10.1186/s40168-019-0722-6).
34. **Kostic AD, Gevers D, Pedamallu CS, Michaud M, Duke F, Earl AM, Ojesina AI, Jung J, Bass AJ, Tabernero J, Baselga J, Liu C, Shivdasani RA, Ogino S, Birren BW, Huttenhower C, Garrett WS, Meyerson M**. 2012. Genomic analysis identifies association of *Fusobacterium* with colorectal carcinoma. *Genome Res* **22**:292–298. doi:[10.1101/gr.126573.111](https://doi.org/10.1101/gr.126573.111).
35. **Li Z, Quan G, Jiang X, Yang Y, Ding X, Zhang D, Wang X, Hardwidge PR, Ren W, Zhu G**. 2018. Effects of Metabolites Derived From Gut Microbiota and Hosts on Pathogens. *Front Cell Infect Microbiol* **8**. doi:[10.3389/fcimb.2018.00314](https://doi.org/10.3389/fcimb.2018.00314).