

Viral and Bacterial Communities of Colorectal Cancer

Geoffrey D Hannigan¹, Melissa B Duhaime², Mack T Ruffin IV³, Charlie C Koumpouras¹,
and Patrick D Schloss^{1,*}

¹Department of Microbiology & Immunology, University of Michigan, Ann Arbor, Michigan, 48109

²Department of Ecology and Evolutionary Biology, University of Michigan, Ann Arbor, Michigan, 48109

³Department of Family and Community Medicine, Pennsylvania State University Hershey Medical Center,

Hershey, Pennsylvania, 17033

*To whom correspondence may be addressed.

Corresponding Author Information

Patrick D Schloss, PhD

1150 W Medical Center Dr. 1526 MSRB I

Ann Arbor, Michigan 48109

Phone: (734) 647-5801

Email: pschloss@umich.edu

Major Classification: Biological Sciences

Minor Classification: Microbiology

Keywords: Colorectal Cancer, Virome, Machine Learning, Bacteriophage, Phage, Microbiome, Microbiota

19 **Abstract**

20 Colorectal cancer is the second leading cause of cancer-related death in the United States and is a primary
21 cause of morbidity and mortality throughout the world. Colorectal cancer development has been linked to
22 differences in colonic bacterial community composition. Viruses are another important component of the
23 colonic microbial community, however they have yet to be studied in colorectal cancer despite their oncogenic
24 potential. We evaluated the colorectal cancer virome (virus community) in stool using a cohort of 90 human
25 subjects with either healthy, adenomatous (precancerous), or cancerous colons. We utilized 16S rRNA gene,
26 whole shotgun metagenomic, and purified virus metagenomic sequencing methods to compare the colorectal
27 cancer virome to the bacterial community. We identified no detectable difference in viral diversity (alpha or
28 beta) between healthy, adenomatous, or cancerous colonic samples, but more sophisticated random forest
29 models identified striking differences in the virome. The majority of the cancer-associated virome consisted
30 of temperate bacteriophages, suggesting that the community was indirectly linked to colorectal cancer by
31 modulating bacterial community structure and function. Our data suggest that the influential phages do
32 not exclusively infect influential bacteria, but rather act through the community as a whole. These results
33 provide foundational evidence that bacteriophage communities are associated with colorectal cancer and
34 likely impact cancer progression by altering the bacterial host communities.

35 **Significance Statement**

36 Colorectal cancer is a leading cause of cancer-related death in the United States and worldwide. Its risk and
37 severity have been linked to colonic bacterial community composition. Although viruses have been linked
38 to other cancers and diseases, little is known about colorectal cancer virus communities. We addressed
39 this knowledge gap by identifying differences in colonic virus communities in the stool of colorectal cancer
40 patients and how they compared to bacterial community differences. The results suggested an indirect role
41 for the virome in impacting colorectal cancer by modulating their associated bacterial community. These
42 findings both support a biological role for viruses in colorectal cancer and provide a new understanding of
43 basic colorectal cancer etiology.

44 **Introduction**

45 Due to their mutagenic abilities and propensity for functional manipulation, human viruses are strongly
46 associated with, and in many cases cause, cancer (1–4). Because bacteriophages (viruses that specifically
47 infect bacteria) are crucial for bacterial community stability and composition (5–7) and have been implicated
48 as oncogenic agents (8–11), bacteriophages have the potential to indirectly impact cancer. The gut virome
49 (the virus community of the gut) therefore has the potential to impact health and disease. Altered human
50 virome composition and diversity have been identified in diseases including periodontal disease (12), HIV
51 (13), cystic fibrosis (14), antibiotic exposure (15, 16), urinary tract infections (17), and inflammatory bowel
52 disease (18). The strong association of bacterial communities with colorectal cancer and the precedent for
53 the virome to impact other human diseases suggest that colorectal cancer may be associated with altered
54 virus communities.

55 Colorectal cancer is the second leading cause of cancer-related deaths in the United States (19). The
56 US National Cancer Institute estimates over 1.5 million Americans were diagnosed with colorectal cancer
57 in 2016 and over 500,000 Americans died from the disease (19). Growing evidence suggests that an
58 important component of colorectal cancer etiology may be perturbations in the colonic bacterial community
59 (8, 10, 11, 20, 21). Work in this area has led to a proposed disease model in which bacteria colonize the
60 colon, develop biofilms, promote inflammation, and enter an oncogenic synergy with the cancerous human
61 cells (22). This association also has allowed researchers to leverage bacterial community signatures as
62 biomarkers to provide accurate, noninvasive colorectal cancer detection from stool (8, 23, 24). While an
63 understanding of colorectal cancer bacterial communities has proven fruitful both for disease classification
64 and for identifying the underlying disease etiology, bacteria are only a subset of the colon microbiome.
65 Viruses are another important component of the colon microbial community that have yet to be studied in
66 the context of colorectal cancer. We evaluated disruptions in virus and bacterial community composition
67 in a human cohort whose stool was sampled at the three relevant stages of cancer development: healthy,
68 adenomatous, and cancerous.

69 Colorectal cancer progresses in a stepwise process that begins when healthy tissue develops into a

70 precancerous polyp (i.e., adenoma) in the large intestine (25). If not removed, the adenoma may develop
71 into a cancerous lesion that can invade and metastasize, leading to severe illness and death. Progression to
72 cancer can be prevented when precancerous adenomas are detected and removed during routine screening
73 (26, 27). Survival for colorectal cancer patients may exceed 90% when the lesions are detected early and
74 removed (26). Thus, work that aims to facilitate early detection and prevention of progression beyond early
75 cancer stages has great potential to inform therapeutic development.

76 Here we address the knowledge gap of whether virus community composition is altered in colorectal cancer
77 and, if it is, how those differences might impact cancer progression and severity. We also aimed to evaluate
78 the virome's potential for use as a diagnostic biomarker. The implications of this study are threefold. *First*,
79 this work supports a biological role for the virome in colorectal cancer development and suggests that more
80 than the bacterial members of the associated microbial communities are involved in the process. *Second*, we
81 present a supplementary, or even alternative, virus-based approach for classification modeling of colorectal
82 cancer using stool samples. *Third*, we provide initial support for the importance of studying the virome as a
83 component of the microbiome ecological network, especially in cancer.

84 **Results**

85 **Cohort Design, Sample Collection, and Processing**

86 Our study cohort consisted of 90 human subjects, 30 of whom had healthy colons, 30 of whom had
87 adenomas, and 30 of whom had carcinomas (**Figure S1**). Half of each stool sample was used to sequence
88 the bacterial communities using both 16S rRNA gene and shotgun sequencing techniques. The 16S
89 rRNA gene sequencing was performed for a previous study, and the sequences were re-analyzed using
90 contemporary methods (8). The other half of each stool sample was purified for virus like particles (VLPs)
91 before genomic DNA extraction and shotgun metagenomic sequencing. In the VLP purification, cells were
92 disrupted and extracellular DNA degraded (**Figure S1**) to allow the exclusive analysis of viral DNA within
93 virus capsids. In this manner, the *extracellular virome* of encapsulated viruses was targeted.

94 Each extraction was performed with a blank buffer control to detect contaminants from reagents or other
95 unintentional sources. Only one of the nine controls contained detectable DNA at a minimal concentration of
96 0.011 ng/ μ l, thus providing evidence of the enrichment and purification of VLP genomic DNA over potential
97 contaminants (**Figure S2 A**). As was expected, these controls yielded few sequences and were almost
98 entirely removed while rarefying the datasets to a common number of sequences (**Figure S2 B**). The high
99 quality phage and bacterial sequences were assembled into highly covered contigs longer than 1 kb (**Figure**
100 **S3**). Because contigs represent genome fragments, we further clustered related bacterial contigs into
101 operational genomic units (OGUs) and viral contigs into operational viral units (OVUs) (**Figure S3 - S4**) to
102 approximate organismal units.

103 **Unaltered Virome Diversity in Colorectal Cancer**

104 Microbiome and disease associations are often described as being of an altered diversity (i.e., “dysbiotic”).
105 We therefore first evaluated the influence of colorectal cancer on virome OVU diversity. We evaluated
106 differences in communities between disease states using the Shannon diversity, richness, and Bray-Curtis
107 metrics. We observed no significant alterations in either Shannon diversity or richness in the diseased states
108 as compared to the healthy state (**Figure S5 C-D**). There was no statistically significant clustering of the
109 disease groups (ANOSIM p-value = 0.4, **Figure S5**). Notably, there was a significant difference between
110 the few blank controls that remained after rarefying the data and the other study groups (ANOSIM p-value
111 < 0.001, **Figure S6**), further supporting the quality of the sample set. In summary, standard alpha and beta
112 diversity metrics were insufficient for capturing virus community differences between disease states (**Figure**
113 **S5**). This is consistent with what has been observed when the same metrics were applied to 16S rRNA
114 sequenced and metagenomic samples (8, 23, 24) and points to the need for alternate approaches to detect
115 the impact of colorectal cancer disease state on these communities.

116 **Altered Virome Composition in Colorectal Cancer**

117 As opposed to the diversity metrics discussed above, OTU-based relative abundance profiles generated
118 from 16S rRNA gene sequences are effective feature sets for classifying stool samples as originating from
119 individuals with healthy, adenomatous, or cancerous colons (8, 23). The exceptional performance of bacteria
120 in these classification models supports a role for bacteria in colorectal cancer. We built off of these findings by
121 evaluating the ability of virus community signatures to classify stool samples and compared their performance
122 to models built using bacterial community signatures.

123 To identify the altered virus communities associated with colorectal cancer, we built and tested random forest
124 models for classifying stool samples as belonging to individuals with either cancerous or healthy colons. We
125 confirmed that our bacterial 16S rRNA gene model replicated the performance of the original report which
126 used logit models instead of random forest models (**Figure 1 A**) (8). We then compared the bacterial OTU
127 model to a model built using OVU relative abundances. The viral model performed as well as the bacterial
128 model (corrected p-value = 0.4), with the viral and bacterial models achieving mean area under the curve
129 (AUC) values of 0.793 and 0.796, respectively (**Figure 1 A - B**). To evaluate the ability of both bacterial and
130 viral biomarkers to classify samples, we built a combined model that used both bacterial and viral community
131 data. The combined model yielded a modest but statistically significant performance improvement beyond
132 the viral (corrected p-value = 0.002) and bacterial (corrected p-value = 0.002) models, yielding an AUC of
133 0.816 (**Figure 1 A - B**). The combined features from the virus and bacterial communities improved our ability
134 to classify stool as belonging to individuals with cancerous colons.

135 To determine the advantage of viral metagenomic methods over bacterial metagenomic methods, we
136 compared the viral model to a model built using OGU relative abundance profiles from bacterial metagenomic
137 shotgun sequencing data. This model performed worse than the other models (mean AUC = 0.505) (**Figure**
138 **1 A - B**). Because the coverage provided by the metagenomic sequencing was not as deep as the equivalent
139 16S rRNA gene sequencing, we attempted to compare the approaches at a common sequencing depth.
140 This investigation revealed that the bacterial 16S rRNA gene model was strongly driven by sparse and low
141 abundance OTUs (**Figure S7**). Removal of OTUs with a median abundance of zero resulted in the removal

142 of six OTUs, and a loss of model performance down to what was observed in the metagenome-based model
143 **(Figure S7 A)**. The majority of these OTUs had a relative abundance lower than 1% across the samples
144 **(Figure S7 B)**. Although the features in the viral model also were of low abundance **(Figure S9 F)**, the
145 coverage was sufficient for high model performance, likely because viral genomes are orders of magnitude
146 smaller than bacterial genomes.

147 The association between the bacterial and viral communities and colorectal cancer was driven by a few
148 important microbes. *Fusobacterium* was the primary driver of the bacterial association with colorectal cancer,
149 which is consistent with its previously described oncogenic potential **(Figure 1 C)**(22). The virome signature
150 also was driven by a few OVUs, suggesting a role for these viruses in tumorigenesis **(Figure 1 D)**. The
151 identified viruses were bacteriophages, belonging to *Siphoviridae*, *Myoviridae*, and “unclassified” phage taxa.
152 Many of the important viruses were unidentifiable (denoted “unknown”). This is common in viromes across
153 habitats; studies have reported as much as 95% of virus sequences belonging to unknown genomic units
154 (14, 28–30). When the bacterial and viral community signatures were combined, both bacterial and viral
155 organisms drove the community association with cancer **(Figure 1 E)**.

156 **Shifted Phage Influence Between Cancer Progression Stages**

157 Because previous work has identified shifts in which bacteria were most important at different stages of
158 colorectal cancer (8, 20, 22), we explored whether shifts in the relative influence of specific phages could be
159 detected between healthy, adenomatous, and cancerous colons. We evaluated community shifts between
160 the two disease stage transitions (healthy to adenomatous and adenomatous to cancerous) by building
161 random forest models to compare only the diagnosis groups around the transitions. While bacterial OTU
162 models performed equally well for all disease class comparisons, the virome model performances differed
163 **(Figure S8 A-B)**. Like bacteria **(Figure S8 F-H)**, different virome members were important between the healthy
164 to adenomatous and adenomatous to cancerous stages **(Figure S8 C-E)**.

165 After evaluating our ability to classify samples between two disease states, we performed a three-class
166 random forest model including all disease states. The 16S rRNA gene model yielded a mean AUC of 0.771

167 and outperformed the viral community model, which yielded a mean AUC of 0.699 (p-value < 0.001, **Figure**
168 **S9 A-C**). The microbes important for the healthy versus cancer and healthy versus adenoma models were
169 also important for the three-class model (**Figure S9 D-E**). The most important bacterium in the two and three
170 class models was the same *Fusobacterium* (OTU 4) (**Figure 1 C, Figure S9 D**). The viruses most important
171 to the three-class model were identified as bacteriophages (**Figure 1 D, Figure S9 E**), but not all important
172 OVUs were of increased abundance in the diseased state (**Figure S9 F**).

173 **Bacteriophage Dominance in Colorectal Cancer Virome**

174 Differences in the colorectal cancer virome could have been driven directly by eukaryotic viruses or indirectly
175 by bacteriophages acting through their bacterial hosts. To better understand the types of viruses that were
176 important for colorectal cancer, we identified the virome OVUs as being similar to either eukaryotic viruses
177 or bacteriophages. The most important viruses to the classification model were identified as bacteriophages
178 (**Figure S9**). Overall, we were able to identify 78.8% of the OVUs as known viruses, and 93.8% of those
179 viral OVUs aligned to bacteriophage reference genomes. It is important to note that this could have been
180 influenced by our methodological biases against enveloped viruses (more common of eukaryotic viruses than
181 bacteriophage), due to chloroform and DNase treatment for purification.

182 We evaluated whether the phages in the community were primarily lytic (i.e. obligately lyse their hosts after
183 replication) or temperate (i.e. able to integrate into their host's genome to form a lysogen, and subsequently
184 transition to a lytic mode). We accomplished this by identifying three markers for temperate phages in the
185 OVU representative sequences: 1) presence of phage integrase genes, 2) presence of known prophage
186 genes, according to the ACLAME (A CLAssification of Mobile genetic Elements) database, and 3) nucleotide
187 similarity to regions of bacterial genomes (29, 31, 32). We found that the majority of the phages were
188 temperate and that the overall fraction of temperate phages remained consistent throughout the healthy,
189 adenomatous, and cancerous stages (**Figure S10 E**). These findings were consistent with previous reports
190 suggesting the gut virome is primarily composed of temperate phages (13, 18, 31, 33).

191 **Community Context of Influential Phages**

192 Because the link between colorectal cancer and the virome was driven by bacteriophages, we hypothesized
193 that the influential phages were primarily predators of the influential bacteria, and thus influenced their relative
194 abundance through predation. If this hypothesis were true, we would expect a correlation between the relative
195 abundances of influential bacteria and phages. Instead, we observed a strikingly low correlation between
196 bacterial and phage relative abundances (**Figure 2 A,C**). Overall, there was an absence of correlation
197 between the most influential OVUs and bacterial OTUs (**Figure 2 B**). This evidence supported our null
198 hypothesis that the influential phages were not primarily predators of influential bacteria.

199 Given these findings, we hypothesized that the most influential phages were acting by infecting a wide range
200 of bacteria in the overall community, instead of just the influential bacteria. In other words, we hypothesized
201 that the influential bacteriophages were community hubs (central members) within the bacteria and phage
202 interactive network. We investigated the potential host ranges of all phage OVUs using a previously
203 developed random forest model that relies on sequence features to predict which phages infected which
204 bacteria in the community (**Figure 3 A**) (34). The predicted interactions were then used to identify phage
205 community hubs. We calculated the alpha centrality (measure of importance in the ecological network)
206 of each phage OVU's connection to the rest of the network. The phages with high centrality values were
207 defined as community hubs. Next, the centrality of each OVU was compared to its importance in the
208 colorectal cancer classification model. Phage OVU centrality was significantly and positively correlated
209 with importance to the disease model (p -value = 0.02, $R = 0.14$), suggesting that phages important in
210 driving colorectal cancer also were more likely to be community hubs (**Figure 3 B**). Together these findings
211 supported our hypothesis that influential phages were hubs within their microbial communities and had
212 broad host ranges.

213 Discussion

214 Because of their propensity for mutagenesis and capacity for modulating their host functionality, many viruses
215 are oncogenic (1–4). Some bacteria also have oncogenic properties, suggesting that bacteriophages may
216 play an indirect role in promoting carcinogenesis by influencing bacterial community composition and
217 dynamics (8–10). Despite their carcinogenic potential and the strong association between bacteria and
218 colorectal cancer, a mechanistic link between virus colorectal communities and colorectal cancer has yet to
219 be evaluated. Here we show that, like colonic bacterial communities, the colon virome was altered in patients
220 with colorectal cancer relative to those with healthy colons. Our findings support a working hypothesis for
221 oncogenesis by phage-modulated bacterial community composition.

222 Here, we have begun to delineate the role the colonic virome plays in colorectal cancer (**Figure 4 A**). We
223 found that basic diversity metrics of alpha diversity (richness and Shannon diversity) and beta diversity
224 (Bray-Curtis dissimilarity) were insufficient for identifying virome community differences between healthy
225 and cancerous states. By implementing a more sophisticated machine learning approach (random forest
226 classification), we detected strong associations between the colon virus community composition and
227 colorectal cancer. The colorectal cancer virome was composed primarily of bacteriophages. These phage
228 communities were not exclusively predators of the most influential bacteria, as demonstrated by the lack
229 of correlation between the abundances of the bacterial and phage populations. Instead, we identified
230 influential phages as being community hubs, suggesting phages influence cancer by altering the greater
231 bacterial community instead of directly modulating the influential bacteria. Our previous work has shown that
232 modifying colon bacterial communities alters colorectal cancer progression and tumor burden in mice (10,
233 20). This provides a precedent for phage indirectly influencing colorectal cancer progression by altering the
234 bacterial community composition. Overall, our data support a model in which the bacteriophage community
235 modulates the bacterial community, and through those interactions indirectly influences the bacteria driving
236 colorectal cancer progression (**Figure 4 A**). Although our evidence suggested phages indirectly influenced
237 colorectal cancer development, we were not able to rule out the role of phages directly interacting with the
238 human host (35, 36).

239 In addition to modeling the potential connections between virus communities, bacteria communities, and
240 colorectal cancer, we also used our data and existing knowledge of phage biology to develop a working
241 hypothesis for the mechanisms by which this may occur. This was done by incorporating our findings into the
242 current model for colorectal cancer development (**Figure 4 B**) (22). We hypothesize that the process began
243 with broadly infectious phages in the colon lysing and thereby disrupting the existing bacterial communities.
244 This shift led to novel niche space that enabled opportunistic bacteria (such as *Fusobacterium nucleatum*)
245 to colonize. Once the initial influential founder bacteria established themselves in the epithelium, secondary
246 opportunistic bacteria were able to adhere to the founders, colonize, and begin establishing a biofilm. Phages
247 may have played a role in biofilm dispersal and growth by lysing bacteria within the biofilm, a process
248 important for effective biofilm growth (37). The oncogenic bacteria may then have been able to transform
249 the epithelial cells and disrupt tight junctions to infiltrate the epithelium, thereby initiating an inflammatory
250 immune response. As the adenomatous polyps developed and progressed towards carcinogenesis, we
251 observed a shift in the phages and bacteria whose relative abundances were most influential. As the bacteria
252 entered their oncogenic synergy with the epithelium, we conjecture that the phages continued mediating
253 biofilm dispersal. This process would thereby support the colonized oncogenic bacteria by lysing competing
254 cells and releasing nutrients to other bacteria in the form of cellular lysates. In addition to highlighting the
255 likely mechanisms by which the colorectal cancer virome is interacting with the bacterial communities, this
256 outline will guide future research investigations of the role the virome plays colorectal cancer.

257 A notable finding was the poor performance observed using bacterial metagenomic methods compared
258 to the performance of models using viral metagenomes or 16S rRNA gene sequences. We believe this
259 observation speaks to the importance of sequencing coverage in microbial community studies and the
260 advantage of the high coverage of 16S rRNA gene sequencing relative to the lower per OTU coverage
261 possible using whole metagenomic shotgun sequencing. To demonstrate this concept, consider that six
262 bacterial OTUs drove the performance of the 16S rRNA gene classification model and these OTUs were all
263 sparsely present and lowly abundant. Filtration of OTUs with a median relative abundance of zero resulted
264 in the removal of the six important OTUs and reduced model performance to being nearly random like
265 the bacterial metagenomic model. The bacterial metagenomic OGUs represented only the most abundant

266 taxa, which was not informative for this application. There has been some success in using shotgun
267 metagenomic approaches for stool colorectal cancer classification (24), but this previous approach relied on
268 lowly abundant signatures and did not utilize OGU clustering, as done here. In that former case, the models
269 only performed as well as the 16S rRNA gene model (24). Thus, the targeted 16S rRNA gene sequencing
270 approach, which represented only a fraction of the bacterial metagenomic sequencing depth, was more
271 effective for detecting colorectal cancer in stool samples. Despite a loss of enthusiasm for 16S rRNA gene
272 sequencing in favor of shotgun metagenomic techniques, 16S rRNA gene sequencing is still a superior
273 methodological approach for some important applications.

274 In addition to the diagnostic ramifications for understanding the colorectal cancer microbiome, our findings
275 suggest that viruses, while understudied and currently under-appreciated in the human microbiome, are
276 likely to be an important contributor to human disease. Viral community dynamics have the potential to
277 provide an abundance of information to supplement those of bacterial communities. Evidence has suggested
278 that the virome is a crucial component to the microbiome and that bacteriophages are important players.
279 Bacteriophage and bacterial communities cannot maintain stability and co-evolution without one another (6,
280 38). Not only is the human virome an important element to consider in human health and disease (12–18),
281 but our findings support that it is likely to have a significant impact on cancer etiology and progression.

282 **Materials & Methods**

283 This study was approved by the University of Michigan Institutional Review Board and all subjects provided
284 informed consent. A detailed description of protocols used for sample collection, processing, and analysis
285 is provided in SI Materials and Methods. All study sequences are available on the NCBI Sequence Read
286 Archive under the BioProject ID PRJNA389927. All associated source code is available at the following
287 GitHub repository:
288 https://github.com/SchlossLab/Hannigan_CRCVirome_PNAS_2017.

289 **Acknowledgments**

290 The authors thank the Schloss lab members for their underlying contributions, and the Great Lakes-New
291 England Early Detection Research Network for providing the fecal samples that were used in this study. GD
292 Hannigan was supported in part by the Molecular Mechanisms in Microbial Pathogenesis Training Program
293 (T32 AI007528). PD Schloss was supported by funding from the National Institutes of Health (P30DK034933).
294 MT Ruffin was supported by funding from the National Institutes of Health (5U01CA86400).

295 **Figures**

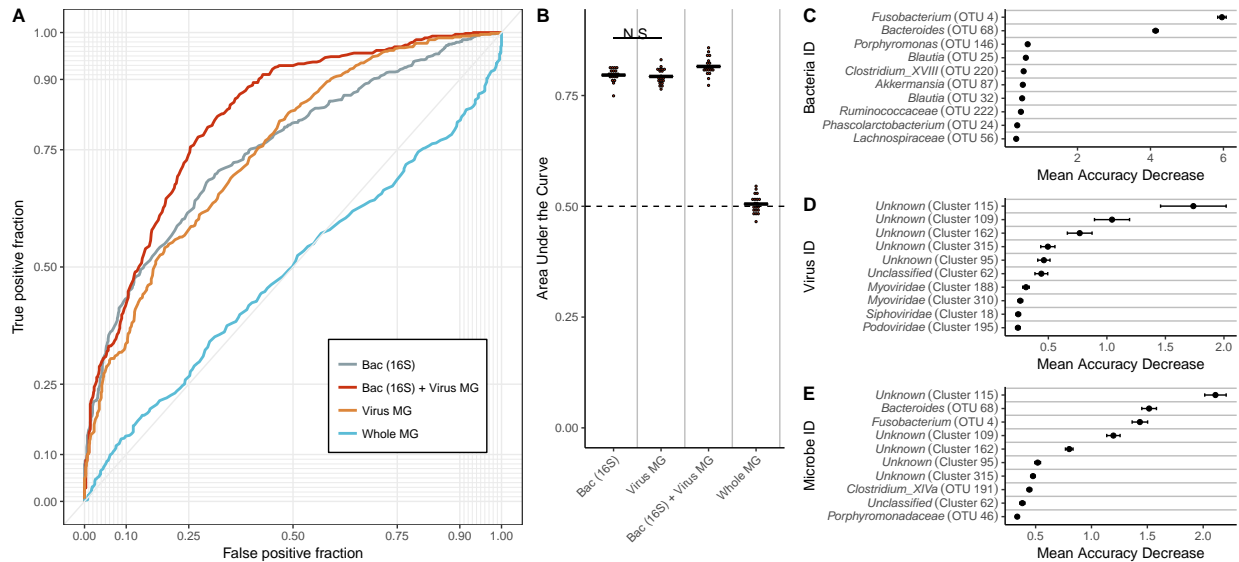


Figure 1: Results from healthy vs cancer classification models built using virome signatures, bacterial 16S rRNA gene sequence signatures, whole metagenomic signatures, and a combination of virome and 16S rRNA gene sequence signatures. A) An example ROC curve for visualizing the performance of each of the models for classifying stool as coming from either an individual with a cancerous or healthy colon. B) Quantification of the AUC variation for each model, and how it compared to each of the other models based on 15 iterations. A pairwise Wilcoxon test with a false discovery rate multiple hypothesis correction demonstrated that all models are significantly different from each other (p -value < 0.01). C) Mean decrease in accuracy (measurement of importance) of each operational taxonomic unit within the 16S rRNA gene classification model when removed from the classification model. Mean is represented by a point, and bars represent standard error. D) Mean decrease in accuracy of each operational virus unit in the virome classification model. E) Mean decrease in accuracy of each operational genomic unit and operational taxonomic unit in the model using both 16S rRNA gene and virome features.

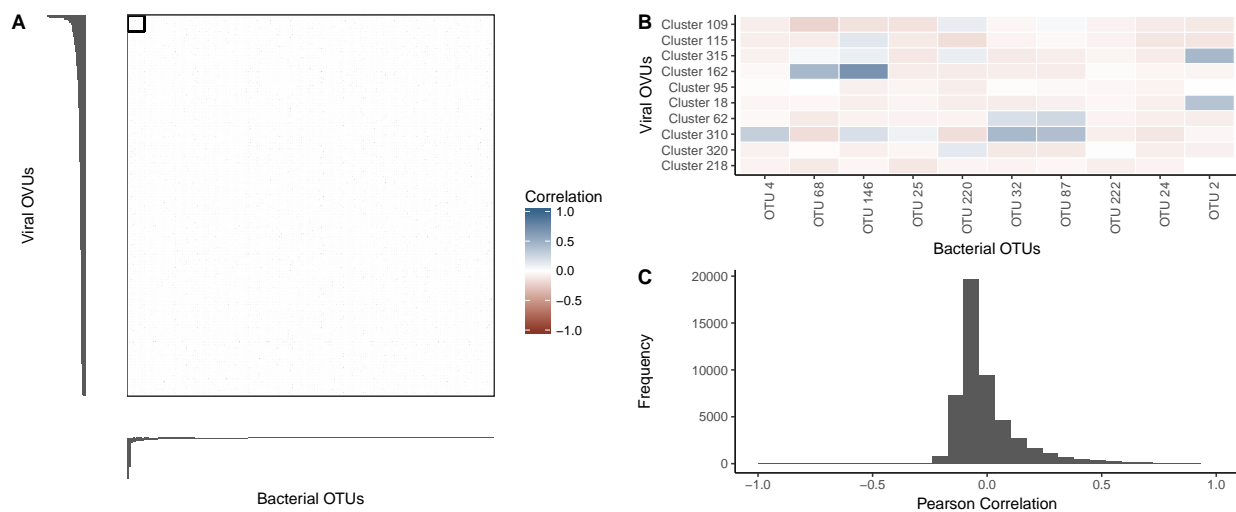


Figure 2: *Relative abundance correlations between bacterial OTUs and virome OVUs. A) Pearson correlation coefficient values between all bacterial OTUs (x-axis) and viral OVUs (y-axis) with blue being positively correlated and red being negatively correlated. Bar plots indicate the viral (left) and bacterial (bottom) operational unit importance in their colorectal cancer classification models, such that the most important units are in the top left corner. B) Magnification of the boxed region in panel (A), highlighting the correlation between the most important bacterial OTUs and virome OVUs. The most important operational units are in the top left corner of the heatmap, and the correlation scale is the same as panel (A). C) Histogram quantifying the frequencies of Pearson correlation coefficients between all bacterial OTUs and virome OVUs.*

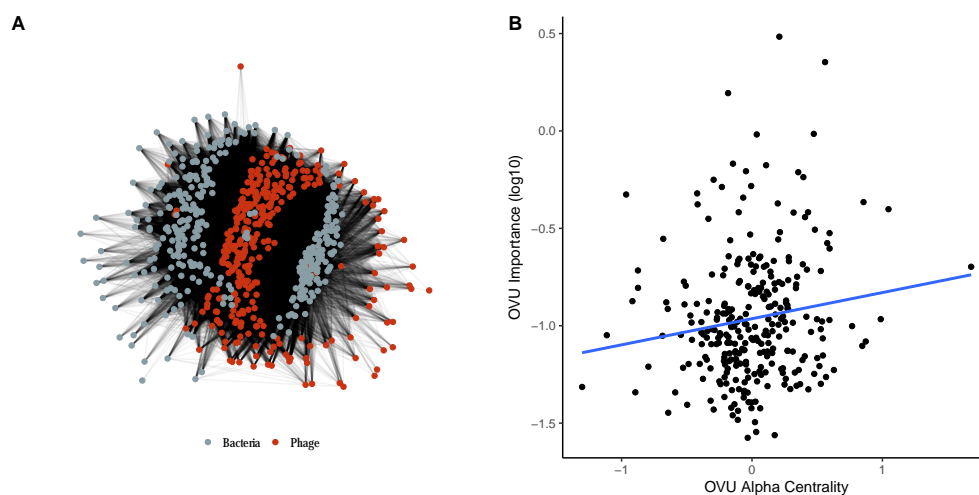


Figure 3: *Community network analysis utilizing predicted interactions between bacteria and phage operational genomic units. A) Visualization of the community network for our colorectal cancer cohort. B) Scatter plot illustrating the correlation between importance (mean decrease in accuracy) and the degree of centrality for each OVU. A linear regression line was fit to illustrate the correlation (blue) which was found to be statistically significantly and weakly correlated (p -value = 0.0173, $R = 0.14$).*

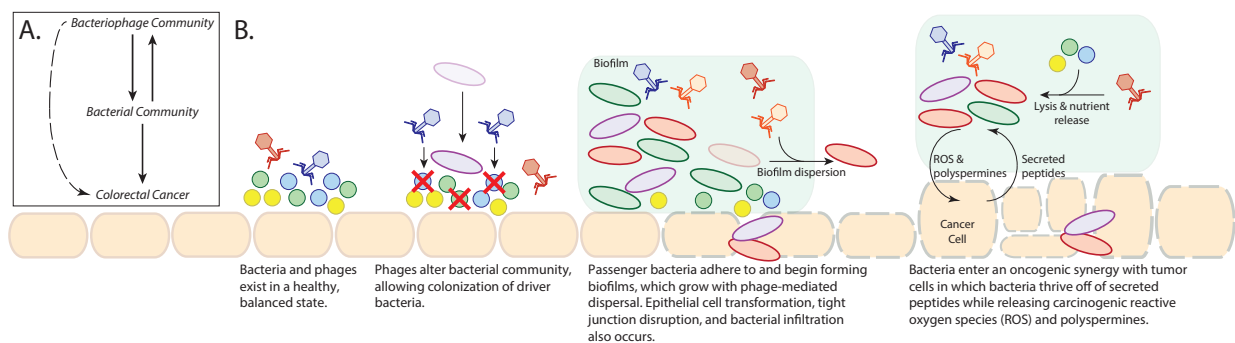


Figure 4: *Final model and working hypothesis from this study. A) Basic model illustrating the connections between the virome, bacterial communities, and colorectal cancer. B) Working hypothesis of how the bacteriophage community is associated with colorectal cancer and the associated bacterial community.*

296 **Supplemental Information**

297 **Supplemental Materials & Methods**

298 **Analysis Source Code & Data Availability**

299 All study sequences are available on the NCBI Sequence Read Archive under the BioProject ID
300 PRJNA389927.

301 All associated source code is available at the following GitHub repository:

302 https://github.com/SchlossLab/Hannigan_CRCVirome_PNAS_2017

303 **Study Design and Patient Sampling**

304 This study was approved by the University of Michigan Institutional Review Board and all subjects provided
305 informed consent. Design and sampling of this sample set have been reported previously (8). Briefly, whole
306 evacuated stool was collected from patients who were 18 years of age or older, able to provide informed
307 consent, have had colonoscopy and histologically confirmed colonic disease status, had not had surgery,
308 had not had chemotherapy or radiation, and were free of known co-morbidities including HIV, chronic viral
309 hepatitis, HNPCC, FAP, and inflammatory bowel disease. Samples were collected from four geographic
310 locations: Toronto (Ontario, Canada), Boston (Massachusetts, USA), Houston (Texas, USA), and Ann Arbor
311 (Michigan, USA). Ninety patients were recruited to the study, thirty of which were designated healthy, thirty
312 with detected adenomas, and thirty with detected carcinomas.

313 **16S rRNA Gene Sequence Data Acquisition & Processing**

314 The 16S rRNA gene sequences associated with this study were previously reported (8). Sequence (fastq)
315 and metadata files were downloaded from <http://www.mothur.org/MicrobiomeBiomarkerCRC>. The 16S rRNA

316 gene sequences were analyzed as described previously, relying on the mothur software package (v1.37.0)
317 (39, 40). Briefly, the sequences were de-replicated, aligned to the SILVA database (41), screened for
318 chimeras using UCHIME (42), and binned into operational taxonomic units (OTUs) using a 97% similarity
319 threshold. Abundances were normalized for uneven sequencing depth by randomly sub-sampling to 10,000
320 sequences, as previously reported (23).

321 **Whole Metagenomic Library Preparation & Sequencing**

322 DNA was extracted from stool samples using the PowerSoil-htp 96 Well Soil DNA Isolation Kit (Mo Bio
323 Laboratories) using an EPMotion 5075 pipetting system. Purified DNA was used to prepare a shotgun
324 sequencing library using the Illumina Nextera XT library preparation kit according to the standard kit protocol,
325 including 12 cycles of limited cycle PCR. The tagmentation time was increased from five minutes to ten
326 minutes to improve DNA fragment length distribution. The library was sequenced using one lane of the
327 Illumina HiSeq4000 platform and yielded 125 bp paired end reads.

328 **Virus Metagenomic Library Preparation & Sequencing**

329 Genomic DNA was extracted from purified virus-like particles (VLPs) from stool samples, using a modified
330 version of a previously published protocol (29, 31, 43, 44). Briefly, an aliquot of stool (~0.1 g) was
331 resuspended in SM buffer (Crystalgen; Catalog #: 221-179) and vortexed to facilitate resuspension. The
332 resuspended stool was centrifuged to remove major particulate debris then filtered through a 0.22- μ m filter
333 to remove smaller contaminants. The filtered supernatant was treated with chloroform for ten minutes
334 with gentle shaking, so as to lyse contaminating cells including bacteria, human, fungi, etc. The exposed
335 genomic DNA from the lysed cells was degraded by treating the samples with 5U of DNase for one hour
336 at 37C. DNase was deactivated by incubating the sample at 75C for ten minutes. The DNA was extracted
337 from the purified virus-like particles (VLPs) using the Wizard PCR Purification Preparation Kit (Promega).
338 Disease classes were staggered across purification runs to prevent run variation as a confounding factor.

339 As for whole community metagenomes, purified DNA was used to prepare a shotgun sequencing library
340 using the Illumina Nextera XT preparation kit according to the standard kit protocol. The tagmentation time
341 was increased from five minutes to ten minutes to improve DNA fragment length distribution. The PCR
342 cycle number was increased from twelve to eighteen cycles to address the low biomass of the samples, as
343 has been described previously (29). The library was sequenced using one lane of the Illumina HiSeq4000
344 platform and yielded 125 bp paired end reads.

345 **Metagenome Quality Control**

346 Both the viral and whole community metagenomic sample sets were subjected to the same quality control
347 procedures. The sequences were obtained as de-multiplexed fastq files and subjected to 5' and 3' adapter
348 trimming using the CutAdapt program (v1.9.1) with an error rate of 0.1 and an overlap of 10 (45). The FastX
349 toolkit (v0.0.14) was used to quality trim the reads to a minimum length of 75 bp and a minimum quality score
350 of 30 (46). Reads mapping to the human genome were removed using the DeconSeq algorithm (v0.4.3) and
351 default parameters (47).

352 **Contig Assembly & Abundance**

353 Contigs were assembled using paired end read files that were purged of sequences without a corresponding
354 pair (e.g. one read removed due to low quality). The Megahit program (v1.0.6) was used to assemble contigs
355 for each sample using a minimum contig length of 1000 bp and iterating assemblies from 21-mers to 101-mers
356 by 20 (48). Contigs from the virus and whole metagenomic sample sets were concatenated within their
357 respective groups. Abundance of the contigs within each sample was calculated by aligning sequences
358 back to the concatenated contig files using the bowtie2 global aligner (v2.2.1), with a 25 bp seed length and
359 an allowance of one mismatch (49). Abundance was corrected for contig reference length and the number of
360 contigs included in each operational genomic unit. Abundance was also corrected for uneven sampling depth
361 by randomly sub-sampling virome and whole metagenomes to 1,000,000 and 500,000 reads, respectively,

362 and by removing samples with fewer total reads than the threshold. Thresholds were set for maximizing
363 sequence information while minimizing numbers of lost samples.

364 **Operational Genomic Unit Classification**

365 Much like operational taxonomic units (OTUs) are used as an operational definition of similar 16S rRNA
366 gene sequences, we defined closely related bacterial contig sequences as operational genomic units (OGUs)
367 and virus contigs as operational viral units (OVUs) in the absence of taxonomic identity. OGUs and OVUs
368 were defined with the CONCOCT algorithm (v0.4.0) which bins related contigs by similar tetra-mer and
369 co-abundance profiles within samples using a variational Bayesian approach (50). CONCOCT was used
370 with a length threshold of 1000 bp for virus contigs and 2000 bp for bacteria.

371 **Diversity**

372 Alpha and beta diversity were calculated using the operational viral unit abundance profiles for each sample.
373 Sequences were rarefied to 100,000 sequences. Samples with less than the cutoff were removed from the
374 analysis. Alpha diversity was calculated using the Shannon diversity and richness metrics. Beta diversity
375 was calculated using the Bray-Curtis metric (mean of 25 random sub-sampling iterations), and the statistical
376 significance between the disease state clusters was assessed using an analysis of similarity (ANOSIM) with
377 a post-hoc multivariate Tukey test. All diversity calculations were performed in R using the Vegan package
378 (51).

379 **Classification Modeling**

380 Classification modeling was performed in R using the Caret package (52). OTU, OVU, and OGU abundance
381 data was preprocessed by removing features (OTUs, OVUs, and OGUs) that were present in less than thirty
382 of the samples. This served both as an effective feature reduction technique and made the calculations
383 computationally feasible. The binary random forest model was trained using the Area Under the receiver

384 operating characteristic Curve (AUC) and the three-class random forest model was trained using the mean
385 AUC. Both were validated using five-fold cross validation. Each training set was repeated five times, and the
386 model was tuned for mtry values. For consistency and accurate comparison between feature groups (e.g.,
387 bacteria, viruses), the sample model parameters were used for each group. The maximum AUC during
388 training was recorded across twenty iterations of each group model to test the significance of the differences
389 between feature set performance. Statistical significance was evaluated using a Wilcoxon test between two
390 categories, or a pairwise Wilcoxon test with Bonferroni corrected p-values when comparing more than two
391 categories.

392 **Taxonomic Identification of Operational Genomic Units**

393 Operational viral units (OVUs) were taxonomically identified using a reference database consisting of
394 all bacteriophage and eukaryotic virus genomes present in the European Nucleotide Archives. The
395 longest contiguous sequence in each operational genomic unit was used as a representative sequence for
396 classification, as described previously (53). Each representative sequence was aligned to the reference
397 genome database using the tblastx alignment algorithm (v2.2.27) and a strict similarity threshold (e-value <
398 1e-25) (54). Annotation was interpreted as phage, eukaryotic virus, or unknown.

399 **Ecological Network Analysis & Correlations**

400 The ecological network of the bacterial and phage operational genomic units was constructed and analyzed
401 as previously described (34). Briefly, a random forest model was used to predict interactions between
402 bacterial and phage genomic units, and those interactions were recorded in a graph database using *neo4j*
403 graph databasing software (v2.3.1). The degree of phage centrality was quantified using the alpha centrality
404 metric in the igraph CRAN package. A Spearman correlation was performed between model importance and
405 phage centrality scores.

406 **Phage Replication Style Identification**

407 Phage OVU replication mode was predicted using methods described previously (29, 31, 32). Briefly, we
408 identified temperate OVUs as representative contigs containing at least one of three genomic markers: 1)
409 phage integrase genes, 2) prophage genes from the ACLAME database, or 3) genomic similarity to bacterial
410 reference genomes. Integrase genes were identified in phage OVU representative contigs by aligning the
411 contigs to a reference database of all known phage integrase genes from the Uniprot database (Uniprot
412 search term: “organism:phage gene:int NOT putative”). Prophage genes were identified in the same way,
413 using the ACLAME set of reference prophage genes. In both cases, the blastx algorithm was used with an
414 e-value threshold of $10e-5$. Representative contigs were also identified as potential lysogenic phages by
415 having a high genomic similarity to bacterial genomes. To accomplish this, representative phage contigs
416 were aligned to the European Nucleotide Archive bacterial genome reference set using the blastn algorithm
417 (e-value $< 10e-25$).

418 **Supplemental Figures**

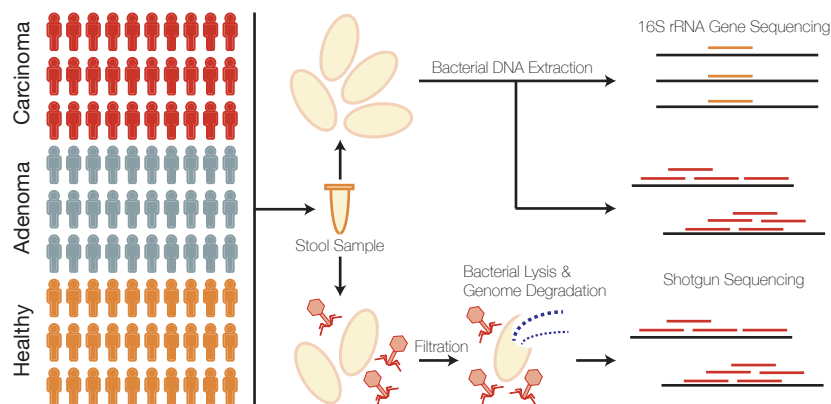


Figure S1: *Cohort and sample processing outline. Thirty subject stool samples were collected from healthy, adenoma (pre-cancer), and carcinoma (cancer) patients. Stool samples were split into two aliquots, the first of which was used for bacterial sequencing and the second which was used for virus sequencing. Bacterial sequencing was done using both 16S rRNA amplicon and whole metagenomic shotgun sequencing techniques. Virus samples were purified for viruses using filtration and a combination of chloroform (bacterial lysis) and DNase (exposed genomic DNA degradation). The resulting encapsulated virus DNA was sequenced using whole metagenomic shotgun sequencing.*

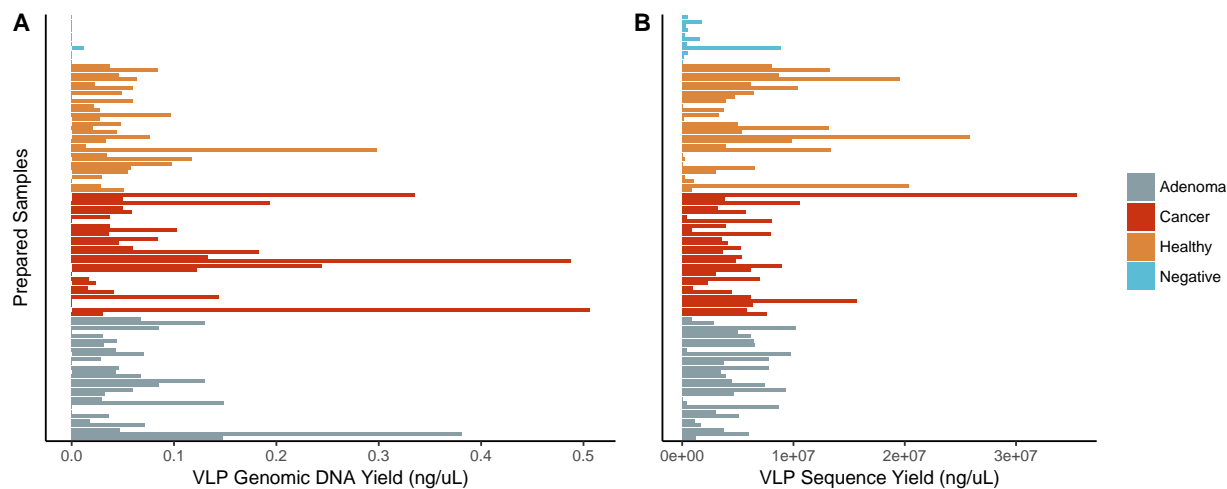


Figure S2: *Basic Quality Control Metrics.* A) VLP genomic DNA yield from all sequenced samples. Each bar represents a sample which is grouped and colored by its associated disease group. B) Sequence yield following quality control including quality score filtering and human decontamination.

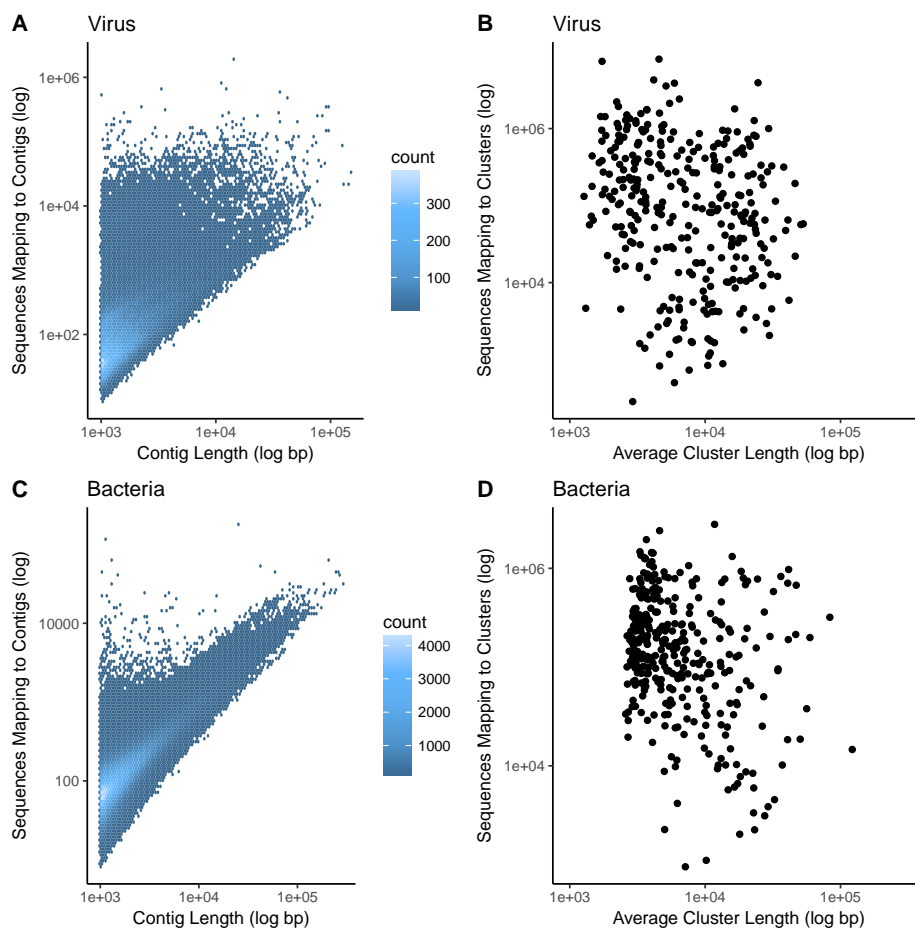


Figure S3: *Length and coverage statistics. A) Heated scatter plot demonstrating the distribution of contig coverage (number of sequences mapping to each contig) and contig length for the virus metagenomic sample set. B) Scatter plot illustrating the distribution of operational viral unit (OVU) length and sequence coverage for the virus metagenomic sample set. C) Heated scatter plot demonstrating the distribution of contig coverage and length for the whole metagenomic sample set. D) Scatter plot illustrating the distribution of operational genomic unit (OGU) length and sequence coverage for the whole metagenomic sample set.*

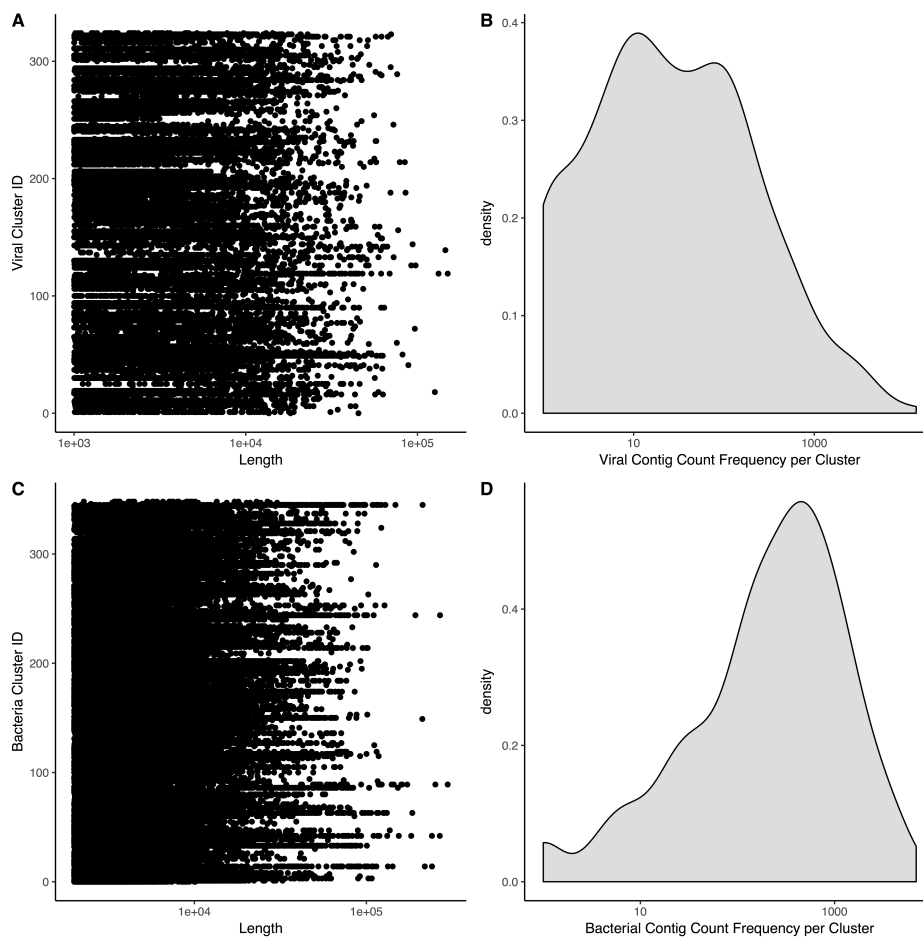


Figure S4: *Operational genomic unit composition stats. A) Strip chart demonstrating the length and frequency of contigs within each operational genomic unit of the virome sample set. The y-axis is the operational genomic unit identifier, and x-axis is the length of each contig, and each dot represents a contig found within the specified operational genomic unit. B) Density plot (analogous to histogram) of the number of virome operational genomic units containing the specific number of contigs, as indicated by the x-axis. C-D) Sample plots as panels C and D, but for the whole metagenomic sample set.*

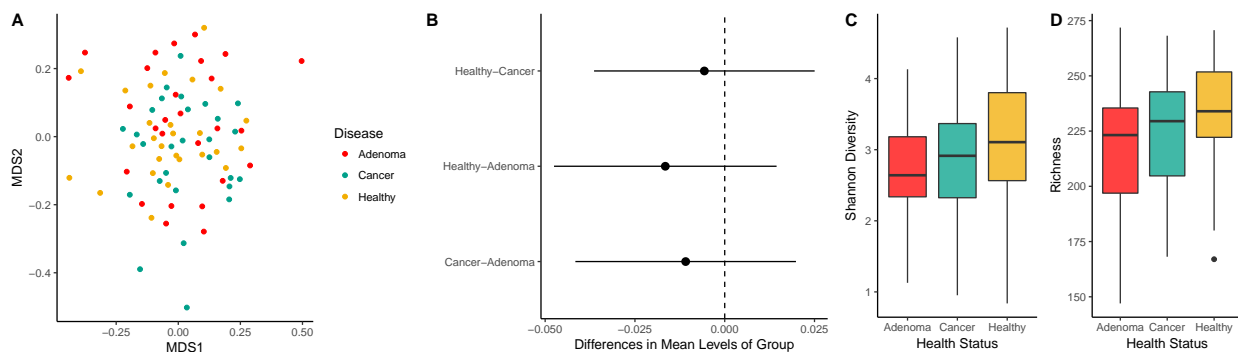


Figure S5: Diversity calculations comparing cancer states of the colorectal virome, based on relative abundance of operational genomic units in each sample. A) NMDS ordination of community samples, colored for cancerous (green), pre-cancerous (red), and healthy (yellow). B) Differences in means between disease group centroids with 95% confidence intervals based on an ANOSIM test with a post hoc multivariate Tukey test. Comparisons (indicated on y-axis) in which the intervals cross the zero mean difference line (dashed line) were not significantly different. C) Shannon diversity and D) richness alpha diversity quantification comparing pre-cancerous (grey), cancerous (red), and healthy (tan) states.

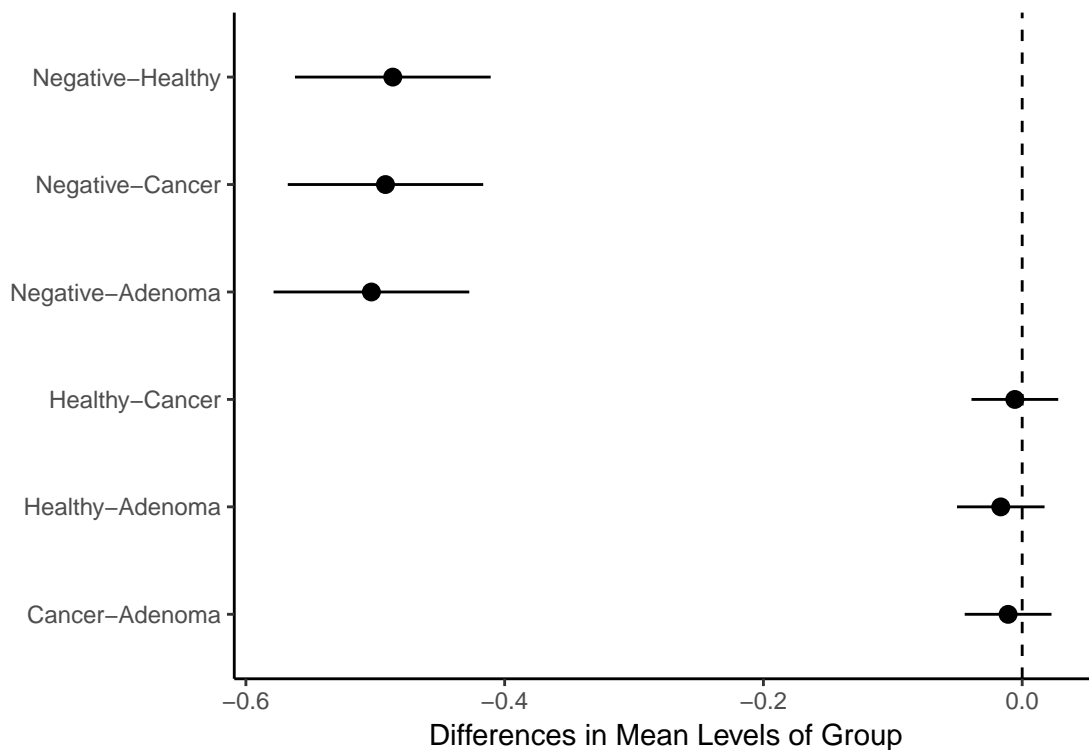


Figure S6: *Beta-diversity comparing disease states and the study negative controls. Differences in means between disease group centroids with 95% confidence intervals based on an ANOSIM test with a post hoc multivariate Tukey test. Comparisons in which the intervals cross the zero mean difference line (dashed line) were not significantly different.*

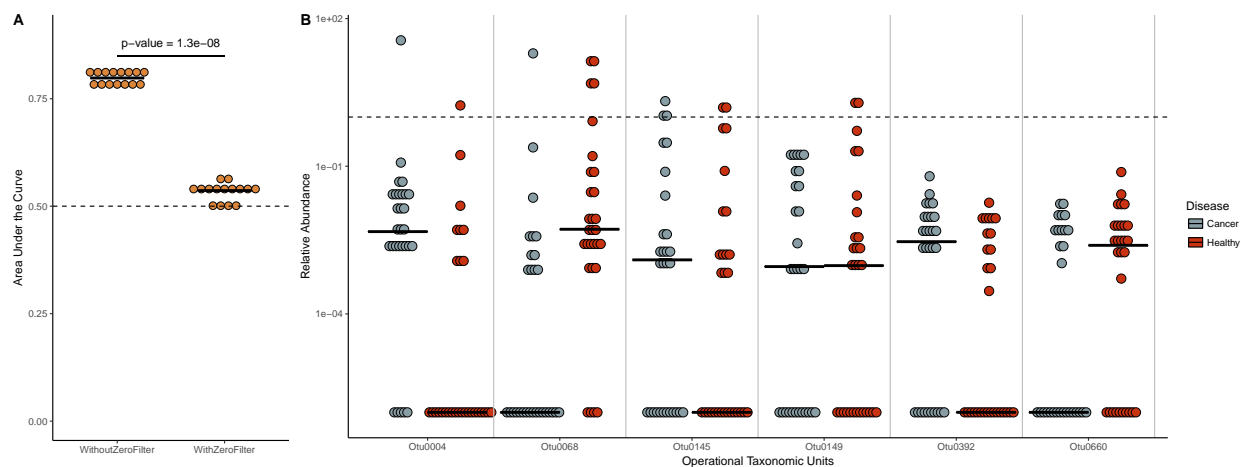


Figure S7: Comparison of bacterial 16S rRNA classification models with and without OTUs whose median relative abundance are greater than zero. A) Classification model performance (measured as area under the curve) for bacteria models using 16S rRNA data both with and without filtering of samples whose median was zero. Significance was calculated using a Wilcoxon rank sum test, and the resulting p-value is shown. The random area under the curve (0.5) is marked with a dashed line. B) Relative abundance of the six bacterial OTUs removed when filtered for OTUs with median relative abundance of zero. OTU relative abundance is separated by healthy (red) and cancerous (grey) samples. Relative abundance of 1% is marked by the dashed line.

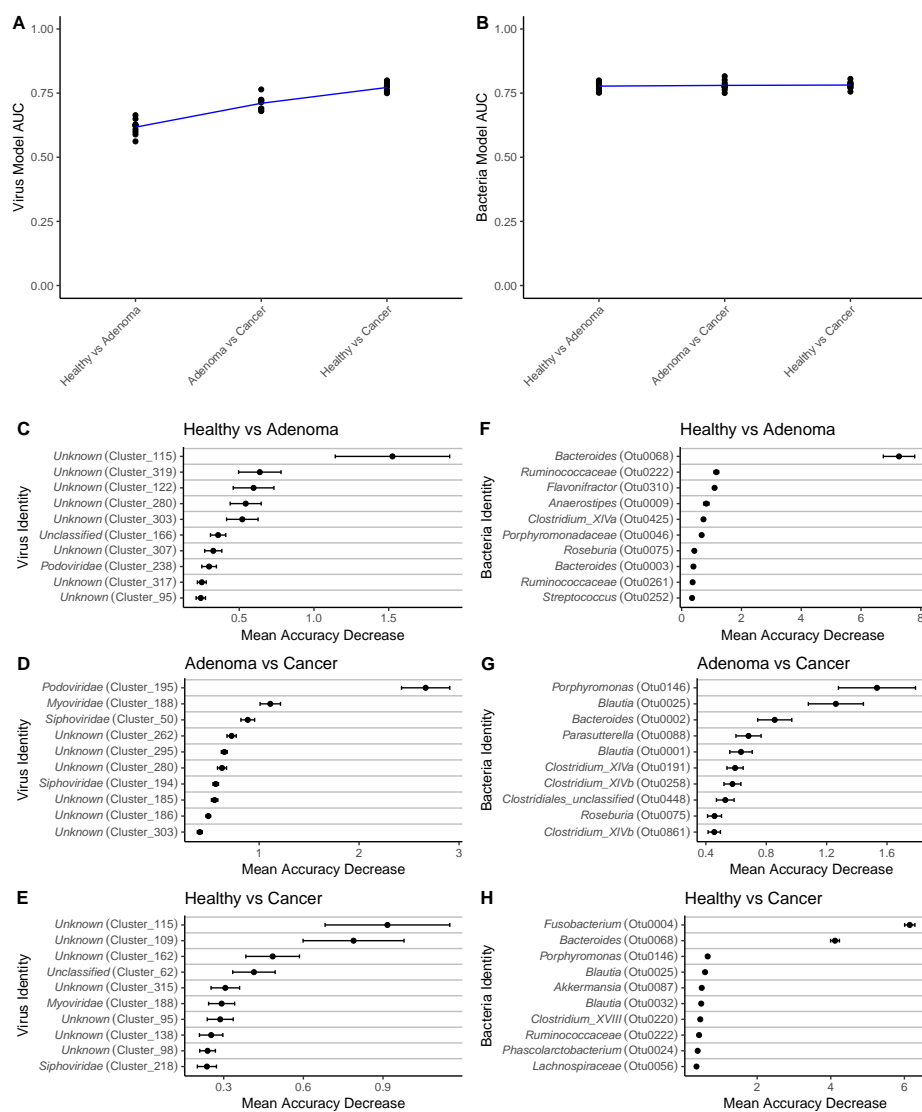


Figure S8: Transition of colorectal cancer importance through disease progression. A) Virus and B) 16S rRNA gene model performance (AUC) when discriminating all binary combinations of disease types. Blue line represents mean performance from multiple random iterations. C-E) Top ten important phage OVUs when classifying each combination of disease state, as measured by the mean decrease in accuracy metric. Mean is represented by a point, and bars represent standard error. Disease comparison is specified in the top left corner of each panel. F-H) Top ten important bacterial 16S rRNA gene OTUs for classifying each disease state combination.

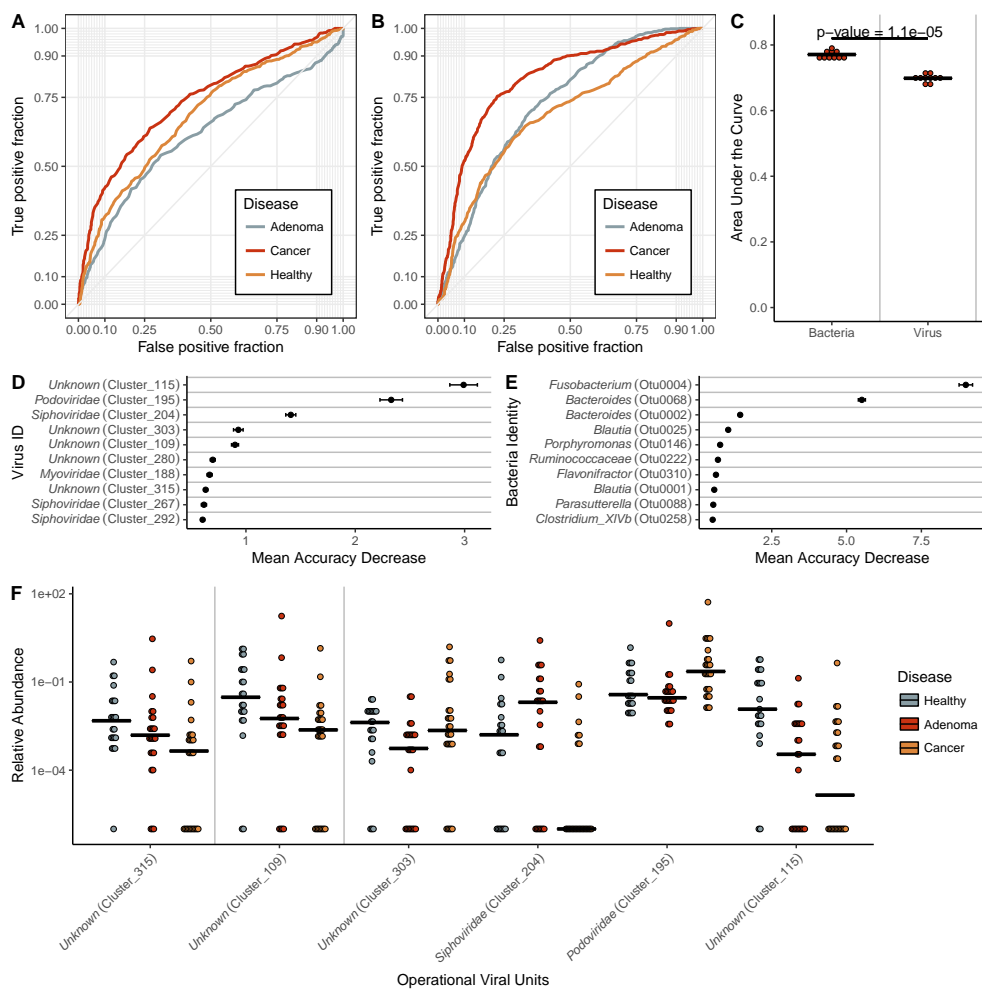


Figure S9: ROC curves from A) virome and B) bacterial 16S three-class random forest models tuned on mean AUC. Each curve represents the ability of the specified class to be classified against the other two classes. C) Quantification of the mean AUC variation for each model based on 10 model iterations. A pairwise Wilcoxon test with a Bonferroni multiple hypothesis correction demonstrated that the models are significantly different ($\alpha = 0.01$). D) Mean decrease in accuracy when virome operational genomic units and E) bacterial 16S OTUs are removed from the respective three-class classification models. Results based on 25 iterations. F) Relative abundance of the six most important virome OVUs in the model, with the most important on the right. Line indicates abundance mean.

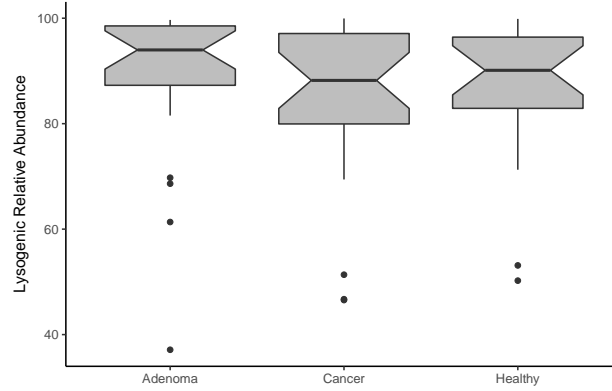


Figure S10: *Lysogenic phage relative abundance in disease states. Phage OVUs were predicted to be either lytic or lysogenic, and the relative abundance of lysogenic phages was quantified and represented as a boxplot. No disease groups were statistically significant.*

419 **References**

- 420 1. Feng H, Shuda M, Chang Y, Moore PS (2008) Clonal integration of a polyomavirus in human Merkel cell
421 carcinoma. *Science* 319(5866):1096–1100.
- 422 2. Shuda M, Kwun HJ, Feng H, Chang Y, Moore PS (2011) Human Merkel cell polyomavirus small T
423 antigen is an oncoprotein targeting the 4E-BP1 translation regulator. *Journal of Clinical Investigation*
424 121(9):3623–3634.
- 425 3. Schiller JT, Castellsagué X, Garland SM (2012) A review of clinical trials of human papillomavirus
426 prophylactic vaccines. *Vaccine* 30 Suppl 5:F123–38.
- 427 4. Chang Y, et al. (1994) Identification of herpesvirus-like DNA sequences in AIDS-associated Kaposi's
428 sarcoma. *Science* 266(5192):1865–1869.
- 429 5. Harcombe WR, Bull JJ (2005) Impact of phages on two-species bacterial communities. *Applied and*
430 *Environmental Microbiology* 71(9):5254–5259.
- 431 6. Rodriguez-Valera F, et al. (2009) Explaining microbial population genomics through phage predation.
432 *Nature Reviews Microbiology* 7(11):828–836.
- 433 7. Cortez MH, Weitz JS (2014) Coevolution can reverse predator-prey cycles. *Proceedings of the National*
434 *Academy of Sciences of the United States of America* 111(20):7486–7491.
- 435 8. Zackular JP, Rogers MAM, Ruffin MT, Schloss PD (2014) The human gut microbiome as a screening tool
436 for colorectal cancer. *Cancer prevention research (Philadelphia, Pa)* 7(11):1112–1121.
- 437 9. Garrett WS (2015) Cancer and the microbiota. *Science* 348(6230):80–86.
- 438 10. Baxter NT, Zackular JP, Chen GY, Schloss PD (2014) Structure of the gut microbiome following
439 colonization with human feces determines colonic tumor burden. *Microbiome* 2(1):20.
- 440 11. Arthur JC, et al. (2012) Intestinal inflammation targets cancer-inducing activity of the microbiota. *Science*

441 338(6103):120–123.

442 12. Ly M, et al. (2014) Altered Oral Viral Ecology in Association with Periodontal Disease. *mBio*
443 5(3):e01133–14–e01133–14.

444 13. Monaco CL, et al. (2016) Altered Virome and Bacterial Microbiome in Human Immunodeficiency
445 Virus-Associated Acquired Immunodeficiency Syndrome. *Cell Host and Microbe* 19(3):311–322.

446 14. Willner D, et al. (2009) Metagenomic analysis of respiratory tract DNA viral communities in cystic fibrosis
447 and non-cystic fibrosis individuals. *PLOS ONE* 4(10):e7370.

448 15. Abeles SR, Ly M, Santiago-Rodriguez TM, Pride DT (2015) Effects of Long Term Antibiotic Therapy on
449 Human Oral and Fecal Viromes. *PLOS ONE* 10(8):e0134941.

450 16. Modi SR, Lee HH, Spina CS, Collins JJ (2013) Antibiotic treatment expands the resistance reservoir and
451 ecological network of the phage metagenome. *Nature* 499(7457):219–222.

452 17. Santiago-Rodriguez TM, Ly M, Bonilla N, Pride DT (2015) The human urine virome in association with
453 urinary tract infections. *Frontiers in Microbiology* 6:14.

454 18. Norman JM, et al. (2015) Disease-specific alterations in the enteric virome in inflammatory bowel disease.
455 *Cell* 160(3):447–460.

456 19. Siegel R, Desantis C, Jemal A (2014) Colorectal cancer statistics, 2014. *CA: a cancer journal for clinicians*
457 64(2):104–117.

458 20. Zackular JP, Baxter NT, Chen GY, Schloss PD (2016) Manipulation of the Gut Microbiota Reveals Role
459 in Colon Tumorigenesis. *mSphere* 1(1):e00001–15.

460 21. Dejea CM, et al. (2014) Microbiota organization is a distinct feature of proximal colorectal cancers.
461 *Proceedings of the National Academy of Sciences of the United States of America* 111(51):18321–18326.

462 22. Flynn KJ, Baxter NT, Schloss PD (2016) Metabolic and Community Synergy of Oral Bacteria in Colorectal

- 463 Cancer. *mSphere* 1(3):e00102–16.
- 464 23. Baxter NT, Ruffin MT, Rogers MAM, Schloss PD (2016) Microbiota-based model improves the sensitivity
465 of fecal immunochemical test for detecting colonic lesions. *Genome medicine* 8(1):37.
- 466 24. Zeller G, et al. (2014) Potential of fecal microbiota for early-stage detection of colorectal cancer.
467 *Molecular systems biology* 10(11):766–766.
- 468 25. Fearon ER (2011) Molecular genetics of colorectal cancer. *Annual review of pathology* 6(1):479–507.
- 469 26. Levin B, et al. (2008) Screening and surveillance for the early detection of colorectal cancer and
470 adenomatous polyps, 2008: a joint guideline from the American Cancer Society, the US Multi-Society Task
471 Force on Colorectal Cancer, and the American College of Radiology. *CA: A Cancer Journal for Clinicians*
472 (The University of Texas MD Anderson Cancer Center, Houston, TX, USA. John Wiley & Sons, Ltd.), pp
473 130–160.
- 474 27. Zauber AG (2015) The impact of screening on colorectal cancer mortality and incidence: has it really
475 made a difference? *Digestive diseases and sciences* 60(3):681–691.
- 476 28. Pedulla ML, et al. (2003) Origins of highly mosaic mycobacteriophage genomes. *Cell* 113(2):171–182.
- 477 29. Hannigan GD, et al. (2015) The Human Skin Double-Stranded DNA Virome: Topographical and Temporal
478 Diversity, Genetic Enrichment, and Dynamic Associations with the Host Microbiome. *mBio* 6(5):e01578–15.
- 479 30. Brum JR, et al. (2015) Ocean plankton. Patterns and ecological drivers of ocean viral communities.
480 *Science* 348(6237):1261498–1261498.
- 481 31. Minot S, et al. (2011) The human gut virome: Inter-individual variation and dynamic response to diet.
482 *Genome Research* 21(10):1616–1625.
- 483 32. Hannigan GD, et al. (2017) Evolutionary and functional implications of hypervariable loci within the skin
484 virome. *PeerJ* 5(4):e2959.
- 485 33. Reyes A, et al. (2010) Viruses in the faecal microbiota of monozygotic twins and their mothers. *Nature*

486 466(7304):334–338.

487 34. Hannigan GD, Duhaime MB, Koutra D, Schloss PD (2017) Biogeography & Environmental Conditions
488 Shape Phage & Bacteria Interaction Networks Across The Human Microbiome. *bioRxiv*:1–40.

489 35. Lengeling A, Mahajan A, Gally DL (2013) Bacteriophages as Pathogens and Immune Modulators? *mBio*
490 4(6):e00868–13–e00868–13.

491 36. Gorski A, et al. (2012) Phage as a Modulator of Immune Responses. *Bacteriophages, Part B*
492 (Bacteriophage Laboratory, Ludwik Hirsfeld Institute of Immunology; Experimental Therapy, Polish
493 Academy of Sciences, Wrocław, Poland. agorski@ikp.pl; Elsevier), pp 41–71.

494 37. Rossmann FS, et al. (2015) Phage-mediated Dispersal of Biofilm and Distribution of Bacterial Virulence
495 Genes Is Induced by Quorum Sensing. *PLoS Pathogens* 11(2):e1004653–17.

496 38. Brockhurst MA, Koskella B (2013) Experimental coevolution of species interactions. *Trends in ecology*
497 *& evolution* 28(6):367–375.

498 39. Kozich JJ, Westcott SL, Baxter NT, Highlander SK, Schloss PD (2013) Development of a dual-index
499 sequencing strategy and curation pipeline for analyzing amplicon sequence data on the MiSeq Illumina
500 sequencing platform. *Applied and Environmental Microbiology* 79(17):5112–5120.

501 40. Schloss PD, et al. (2009) Introducing mothur: open-source, platform-independent, community-supported
502 software for describing and comparing microbial communities. *Applied and Environmental Microbiology*
503 75(23):7537–7541.

504 41. Pruesse E, et al. (2007) SILVA: a comprehensive online resource for quality checked and aligned
505 ribosomal RNA sequence data compatible with ARB. *Nucleic Acids Research* 35(21):7188–7196.

506 42. Edgar RC, Haas BJ, Clemente JC, Quince C, Knight R (2011) UCHIME improves sensitivity and speed
507 of chimera detection. *Bioinformatics* 27(16):2194–2200.

508 43. Thurber RV, Haynes M, Breitbart M, Wegley L, Rohwer F (2009) Laboratory procedures to generate viral

509 metagenomes. *Nature protocols* 4(4):470–483.

510 44. Kleiner M, Hooper LV, Duerkop BA (2015) Evaluation of methods to purify virus-like particles for
511 metagenomic sequencing of intestinal viromes. *BMC Genomics* 16(1):7.

512 45. Martin M (2011) Cutadapt removes adapter sequences from high-throughput sequencing reads.
513 *EMBnetjournal* 17(1):10.

514 46. Hannon GJ FASTX-Toolkit. GNU Affero General Public License.

515 47. Schmieder R, Edwards R (2011) Fast identification and removal of sequence contamination from genomic
516 and metagenomic datasets. *PLOS ONE* 6(3):e17288.

517 48. Li D, et al. (2016) MEGAHIT v1.0: A fast and scalable metagenome assembler driven by advanced
518 methodologies and community practices. *METHODS* 102:3–11.

519 49. Langmead B, Salzberg SL (2012) Fast gapped-read alignment with Bowtie 2. *Nature Methods*
520 9(4):357–359.

521 50. Alneberg J, et al. (2014) Binning metagenomic contigs by coverage and composition. *Nature*
522 *Methods*:1–7.

523 51. Oksanen J, et al. vegan: Community Ecology Package.

524 52. Kuhn M caret: Classification and Regression Training.

525 53. Guidi L, et al. (2016) Plankton networks driving carbon export in the oligotrophic ocean. *Nature*
526 532(7600):465–470.

527 54. Camacho C, et al. (2009) BLAST+: architecture and applications. *BMC Bioinformatics* 10(1):1.