

1                   Biogeography & environmental conditions shape  
2                   bacteriophage-bacteria networks across the human  
3                   microbiome

4       Geoffrey D Hannigan<sup>1</sup>, Melissa B Duhaime<sup>2</sup>, Danai Koutra<sup>3</sup>, and Patrick D Schloss<sup>1,\*</sup>

5       <sup>1</sup>Department of Microbiology & Immunology, University of Michigan, Ann Arbor, Michigan, 48109, USA

6       <sup>2</sup>Department of Ecology & Evolutionary Biology, University of Michigan, Ann Arbor, Michigan, 48109, USA

7       <sup>3</sup>Department of Computer Science, University of Michigan, Ann Arbor, Michigan, 48109, USA

8                   \*To whom correspondence may be addressed.

9  
10  
11  
12    **Short Title:** Network Diversity of the Healthy Human Microbiome

13    **Corresponding Author Information**

14    Patrick D Schloss, PhD

15    1150 W Medical Center Dr. 1526 MSRB I

16    Ann Arbor, Michigan 48109

17    Phone: (734) 647-5801

18    Email: pschloss@umich.edu

## 19 **Abstract**

20 Viruses and bacteria are critical components of the human microbiome and play important roles in health and  
21 disease. Most previous work has relied on studying bacteria and viruses independently, thereby reducing  
22 them to two separate communities. Such approaches are unable to capture how these microbial communities  
23 interact, such as through processes that maintain community robustness or allow phage-host populations  
24 to co-evolve. We implemented a network-based analytical approach to describe phage-bacteria network  
25 diversity throughout the human body. We built these community networks using a machine learning algorithm  
26 to predict which phages could infect which bacteria in a given microbiome. Our algorithm was applied to  
27 paired viral and bacterial metagenomic sequence sets from three previously published human cohorts. We  
28 organized the predicted interactions into networks that allowed us to evaluate phage-bacteria connectedness  
29 across the human body. We observed evidence that gut and skin network structures were person-specific  
30 and not conserved among cohabitating family members. High-fat diets appeared to be associated with  
31 less connected networks. Network structure differed between skin sites, with those exposed to the external  
32 environment being less connected and likely more susceptible to network degradation by microbial extinction  
33 events. This study quantified and contrasted the diversity of virome-microbiome networks across the human  
34 body and illustrated how environmental factors may influence phage-bacteria interactive dynamics. This  
35 work provides a baseline for future studies to better understand system perturbations, such as disease states,  
36 through ecological networks.

## 37 **Author Summary**

38 The human microbiome, the collection of microbial communities that colonize the human body, is a crucial  
39 component to health and disease. Two major components of the human microbiome are the bacterial and  
40 viral communities. These communities have primarily been studied separately using metrics of community  
41 composition and diversity. These approaches have failed to capture the complex dynamics of interacting  
42 bacteria and phage communities, which frequently share genetic information and work together to maintain  
43 ecosystem homestasis (e.g. kill-the-winner dynamics). Removal of bacteria or phage can disrupt or even  
44 collapse those ecosystems. Relationship-based network approaches allow us to capture this interaction  
45 information. Using this network-based approach with three independent human cohorts, we were able to  
46 present an initial understanding of how phage-bacteria networks differ throughout the human body, so as to  
47 provide a baseline for future studies of how and why microbiome networks differ in disease states.

## 48 Introduction

49 Viruses and bacteria are critical components of the human microbiome and play important roles in health  
50 and disease. Bacterial communities have been associated with disease states, including a range of skin  
51 conditions [1], acute and chronic wound healing conditions [2,3], and gastrointestinal diseases, such as  
52 inflammatory bowel disease [4,5], *Clostridium difficile* infections [6] and colorectal cancer [7,8]. Altered  
53 human viromes (virus communities consisting primarily of bacteriophages) also have been associated with  
54 diseases and perturbations, including inflammatory bowel disease [5,9], periodontal disease [10], spread of  
55 antibiotic resistance [11], and others [12–17]. Viruses act in concert with their microbial hosts as a single  
56 ecological community [18]. Viruses influence their living microbial host communities through processes  
57 including lysis, host gene expression modulation [19], influence on evolutionary processes such as horizontal  
58 gene transfer [22] or antagonistic co-evolution [26], and alteration of ecosystem processes and elemental  
59 stoichiometry [27].

60 Previous human microbiome work has focused on bacterial and viral communities, but have reduced them to  
61 two separate communities by studying them independently [5,9,10,12–17]. This approach fails to capture the  
62 complex dynamics of interacting bacteria and phage communities, which frequently share genetic information  
63 and work together to maintain ecosystem structure (e.g. kill-the-winner dynamics that prevent domination  
64 by a single bacterium). Removal of bacteria or phages can disrupt or even collapse those ecosystems  
65 [18,28–37]. Integrating these datasets as relationship-based networks allow us to capture this complex  
66 interaction information. Studying such bacteria-phage interactions through community-wide networks built  
67 from inferred relationships begins to provide us with insights into the drivers of human microbiome diversity  
68 across body sites and enable the study of human microbiome network dynamics overall.

69 In this study, we characterized human-associated bacterial and phage communities by their inferred  
70 relationships using three published paired virus and bacteria-dominated whole community metagenomic  
71 datasets [13,14,38,39]. We leveraged machine learning and graph theory techniques to establish  
72 and explore the human bacteria-phage network diversity therein. This approach built upon previous  
73 large-scale phage-bacteria network analyses by inferring interactions from metagenomic datasets, rather

74 than culture-dependent data [33], which is limited in the scale of possible experiments and analyses.  
75 We implemented an adapted metagenomic interaction inference model that made some improvements  
76 upon previous phage-host interaction prediction models. Previous approaches have utilized a variety of  
77 techniques, such as linear models that were used to predict bacteria-phage co-occurrence using taxonomic  
78 assignments [40], and nucleotide similarity models that were applied to both whole virus genomes [41] and  
79 clusters of whole and partial virus genomes [42]. Our approach uniquely included protein interaction data  
80 and was validated based on experimentally determined positive and negative interactions (i.e. who does  
81 and does not infect whom). We built on previous modeling work as a means to our ends, and focused on the  
82 biological insights we could gain instead of building a superior model and presenting our work as a toolkit.  
83 We therefore did not focus on extensive benchmarking against other existing models [41,43–45]. Through  
84 this approach we were able to provide an initial basic understanding of the network dynamics associated  
85 with phage and bacterial communities on and in the human body. By building and utilizing a microbiome  
86 network, we found that different people, body sites, and anatomical locations not only support distinct  
87 microbiome membership and diversity [13,14,38,39,46–48], but also support ecological communities with  
88 distinct communication structures and robustness to network degradation by extinction events. Through an  
89 improved understanding of network structures across the human body, we aim to empower future studies to  
90 investigate how these communities dynamics are influenced by disease states and the overall impact they  
91 may have on human health.

## 92 **Results**

### 93 **Cohort Curation and Sample Processing**

94 We studied the differences in virus-bacteria interaction networks across healthy human bodies by leveraging  
95 previously published shotgun sequence datasets of purified viral metagenomes (viromes) paired with  
96 bacteria-dominated whole community metagenomes. Our study contained three datasets that explored  
97 the impact of diet on the healthy human gut virome [14], the impact of anatomical location on the healthy

98 human skin virome [13], and the viromes of monozygotic twins and their mothers [38,39]. We selected  
99 these datasets because their virome samples were subjected to virus-like particle (VLP) purification,  
100 which removed contaminating DNA from human cells, bacteria, etc. To this end, the publishing authors  
101 employed combinations of filtration, chloroform/DNase treatment, and cesium chloride gradients to eliminate  
102 organismal DNA (e.g. bacteria, human, fungi, etc) and thereby allow for direct assessment of both the  
103 extracellular and fully-assembled intracellular virome (**Supplemental Figure S1 A-B**) [14,39]. Each  
104 research group reported quality control measures to ensure the purity of the virome sequence datasets,  
105 using both computational and molecular techniques (e.g. 16S rRNA gene qPCR) (**Table S1**). These reports  
106 confirmed that the virome libraries consisted of highly purified virus genomic DNA.

107 The bacterial and viral sequences from these studies were quality filtered and assembled into contigs  
108 (i.e. genomic fragments). We further grouped the related bacterial and phage contigs into operationally  
109 defined units based on their k-mer frequencies and co-abundance patterns, similar to previous reports  
110 (**Supplemental Figure S2 - S3**) [42]. This was done both for dimensionality reduction and to prevent  
111 inflation of node counts by using contigs which are expected to represent multiple fragments from the same  
112 genomes. This was also done to create genome analogs that we could use in our classification model which  
113 was built using genome sequences. We referred to these operationally defined groups of related contigs as  
114 operational genomic units (OGUs). Each OGU represented a genomically similar sub-population of either  
115 bacteria or phages. Contig lengths within clusters ranged between  $10^3$  and  $10^{5.5}$  bp (**Supplemental Figure**  
116 **S2 - S3**).

117 The original publications reported that the whole metagenomic shotgun sequence samples, which primarily  
118 consisted of bacteria, had an average viral relative abundance of 0.4% (**Table S1**) [13,14,38,39]. We  
119 confirmed these reports by finding that only 2% (6 / 280 OGUs) of bacterial OGUs had significantly strong  
120 nucleotide similarity to phage reference genomes (e-value  $< 10^{-25}$ ) [13,14,38,39]. Additionally, no OGUs  
121 were confidently identified as lytic or temperate phage OGUs in the bacterial dataset using the Virsorter  
122 algorithm [50]. We also supplemented the previous virome fraction quality control measures (**Table S1**) to  
123 find that, in light of the rigorous purification and quality control during sample collection and preparation,

124 77% (228 / 298 operational genomic units) still had some nucleotide similarity to a given bacterial reference  
125 genome (e-value <  $10^{-25}$ ). It is important to note that interpreting such alignment is complicated by the  
126 fact that most reference bacterial genomes also contain prophages (i.e. phages integrated into bacterial  
127 genomes), meaning we do not know to what extent the alignments were the result of bacterial contaminants  
128 in the virome fraction and what were true integrated prophages. As most phages in these communities  
129 have been shown to be temperate (i.e. they integrate into bacterial genomes) using methods that included  
130 nucleotide alignments of phages to bacterial reference genomes, we expected that a large fraction of  
131 those phages were temperate and therefore shared elements with bacterial reference genomes – a trend  
132 previously reported [14]. To ensure the purity of our sample sets, we supplemented the quality control  
133 measures by filtering out all OGUs that could be potential bacterial contaminants, as described previously  
134 [42]. This resulted in the removal of 143 OGUs that exhibited nucleotide similarity to bacterial genomes but  
135 no identifiable known phage elements. We were also able to identify two OGUs as representing **complete**,  
136 high confidence phages using the stringent Virsorter phage identification algorithm (class 1 confidence  
137 group) [50].

## 138 **Implementing Phage-Bacteria Interaction Prediction to Build a Community Network**

139 We predicted which phage OGUs infected which bacterial OGUs using a random forest model trained on  
140 experimentally validated infectious relationships from six previous publications [41,51–55]. Only bacteria  
141 and bacteriophages were used in the model. The training set contained 43 diverse bacterial species and 30  
142 diverse phage strains, including both broad and specific ranges of infection (**Supplemental Figure S4 A -**  
143 **B, Table S2**). While it is true that there are more known phages that infect bacteria, we were limited by the  
144 information confirming which phages do not infect certain bacteria and attempted to keep the numbers of  
145 positive and negative interactions similar. Phages with linear and circular genomes, as well as ssDNA and  
146 dsDNA genomes, were included in the analysis. Because we used DNA sequencing studies, RNA phages  
147 were not considered (**Supplemental Figure S4 C-D**). This training set included both positive relationships  
148 (i.e. a phage infects a bacterium) and negative relationships (i.e. a phage does not infect a bacterium). This

149 allowed us to validate the false positive and false negative rates associated with our candidate models,  
150 thereby building upon previous work that only considered positive relationships [41]. It is also worth noting  
151 that while a positive interaction is strong evidence that the interaction exists, we must also be conscious that  
152 negative interactions are only weak evidence for a lack of interaction because the finding could be the result  
153 of our inability to reproduce conditions in which those interactions occur. Altogether we decided to maintain a  
154 balanced dataset at the cost of under-sampling the available positive interaction information because the use  
155 of such a severely unbalanced dataset often results in over-fit and uninformative model training. However,  
156 as an additional validation measure, we used the extensive additional positive interactions as a secondary  
157 dataset to confirm that we could identify infectious interactions from a more diverse bacterial and phage  
158 dataset. Using this approach, we confirmed that 382 additional phage reference genomes, representing  
159 a diverse range of phages, were matched to at least one reference bacterial host genome of the species  
160 that they were expected to infect (**Supplemental Figure S5**). Because the model was built on full genomes  
161 and used on OGU, we also assessed whether our model was resilient to incomplete reference genomes.  
162 We found that the use of our model on random contigs representing as little as 50% length of the original  
163 reference phage and bacterial genomes resulted in minimal reduction in the ability of the model to identify  
164 relationships (**Supplemental Figure S6**).

165 Four phage and bacterial genomic features were used in our random forest model to predict infectious  
166 relationships between bacteria and phages: 1) genome nucleotide similarities, 2) gene amino acid  
167 sequence similarities, 3) bacterial Clustered Regularly Interspaced Short Palindromic Repeat (CRISPR)  
168 spacer sequences that target phages, and 4) similarity of protein families associated with experimentally  
169 identified protein-protein interactions [56]. These features were calculated using the training set described  
170 above. While the nucleotide and amino acid similarity metrics were expected to identify prophage signatures,  
171 the protein family interaction and CRISPR signatures were expected to aid in identifying lytic phages in  
172 addition to temperate phages. We chose to utilize these metrics that directly compare nucleotide sequences  
173 between sample phages and bacteria, instead of relying on alignment to reference genomes or known  
174 marker genes, because we were extrapolating our model to highly diverse communities which we expect  
175 to diverge significantly from the available reference genomes. The resulting random forest model was



176 assessed using nested cross validation, and the median area under its receiver operating characteristic  
177 (ROC) curve was 0.788, the median model sensitivity was 0.905, and median specificity was 0.538 (**Figure**  
178 **1 A**). This balance of confident true positives at the cost of fewer true negatives was ideal for this type  
179 of dataset which consisted of primarily positive connections (**Supplemental Figure S7**). Nested cross  
180 validation of the model demonstrated that the sensitivity and specificity of the model could vary but the  
181 overall model performance (AUC) remained more consistent (**Supplemental Figure S8**). This suggested  
182 that our model would perform with a similar overall accuracy despite changes in sensitivity/specificity  
183 trade-offs. The most important predictor in the model was amino acid similarity between genes, followed by  
184 nucleotide similarity of whole genomes (**Figure 1 B**). Protein family interactions were moderately important  
185 to the model, and CRISPRs were largely uninformative, due to the minimal amount of identifiable CRISPRs  
186 in the dataset and their redundancy with the nucleotide similarity methods (**Figure 1 B**). Approximately  
187 one third of the training set relationships yielded no score and therefore were unable to be assigned an  
188 interaction prediction (**Figure 1 C**).

Figure 1: **Summary of Multi-Study Network Model.** (A) Median ROC curve (dark red) used to create the microbiome-virome infection prediction model, based on nested cross validation over 25 random iterations. The maximum and minimum performance are shown in light red. (B) Importance scores associated with the metrics used in the random forest model to predict relationships between bacteria and phages. The importance score is defined as the mean decrease in accuracy of the model when a feature (e.g. Pfam) is excluded. Features include the local gene alignments between bacteria and phage genes (denoted *blastx*; the *blastx* algorithm in Diamond aligner), local genome nucleotide alignments between bacteria and phage OGUs, presence of experimentally validated protein family domains (Pfam) between phage and bacteria OGUs, and CRISPR targeting of bacteria toward phages (CRISPR). (C) Proportions of samples included (gray) and excluded (red) in the model. Samples were excluded from the model because they did not yield any scores. Those interactions without scores were automatically classified as not having interactions. (D) Network diameter (measure of graph size; the greatest number of traversed vertices required between two vertices), (E) number of vertices, and (F) number of edges (relationships) for the total network (orange) and the individual study sub-networks (diet study = red, skin study = yellow, twin study = green).

189 We used our random forest model to classify the relationships between bacteria and phage operational  
190 genomic units, which were then used to build the interactive network. The master network, analogous  
191 to the universal microbiome network concept previously described [57], contained the three studies as  
192 sub-networks, which themselves each contained sub-networks for each sample (**Supplemental Figure S9**).

193 Metadata including study, sample ID, disease, and OGU abundance within the community were stored in  
194 the master network for parsing in downstream analyses (**Supplemental Figure S9**). The phage and bacteria  
195 of the master network demonstrated both narrow broad ranges of infectious interactions (**Supplemental**  
196 **Figure S10**). Bacterial and phage relative abundance was recorded in each sample for each OGU and  
197 the weight of the edge connecting those OGUs was calculated as a function of those relative abundance  
198 values. The separate extraction of the phage and bacterial libraries ensured a more accurate measurement  
199 of the microbial communities, as has been outlined previously [58,59]. The master network was highly  
200 connected and contained 38,337 infectious relationships among 435 nodes, representing 155 phages and  
201 280 bacteria. Although the network was highly connected, not all relationships were present in all samples.  
202 Relationships were weighted by the relative abundances of their associated bacteria and phages. Like the  
203 master network, the skin network exhibited a diameter of 4 (measure of graph size; the greatest number of  
204 traversed vertices required between two vertices) and included 433 (154 phages, 279 bacteria, 99.5% total)  
205 and 38,099 (99.4%) of the master network nodes and edges, respectively (**Figure 1 E - F**). Additionally, the  
206 subnetworks demonstrated narrow ranges of eccentricity across their nodes (**Supplemental Figure S11**).  
207 The phages and bacteria in the diet and twin sample sets were more sparsely related, with the diet study  
208 consisting of 80 (32 phages, 48 bacteria) nodes and 1,290 relationships, and the twin study containing  
209 130 (29 phages, 101 bacteria) nodes and 2,457 relationships (**Figure 1 E - F**). As a validation measure,  
210 we identified five (1.7%) examples of phage OGUs which contained similar genomic elements to the four  
211 previously described, broadly infectious phages isolated from Lake Michigan (tblastx; e-value <  $10^{-25}$ ) [60].

## 212 **Role of Diet on Gut Microbiome Connectivity**

213 Diet is a major environmental factor that influences resource availability and gut microbiome composition  
214 and diversity, including bacteria and phages [14,61,62]. Previous work in isolated culture-based systems has  
215 suggested that changes in nutrient availability are associated with altered phage-bacteria network structures  
216 [30], although this has yet to be tested in humans. We therefore hypothesized that a change in diet would  
217 also be associated with a change in virome-microbiome network structure in the human gut.

218 We evaluated the diet-associated differences in gut virome-microbiome network structure by quantifying how  
219 central each sample's network was on average. We accomplished this by utilizing two common weighted  
220 centrality metrics: degree centrality and closeness centrality. Degree centrality, the simplest centrality metric,  
221 was defined as the number of connections each phage made with each bacterium. We supplemented  
222 measurements of degree centrality with measurements of closeness centrality. Closeness centrality is a  
223 metric of how close each phage or bacterium is to all of the other phages and bacteria in the network. A higher  
224 closeness centrality suggests that the effects of genetic information or altered abundance would be more  
225 impactful to all other microbes in the system. Because these are weighted metrics, the values are functions  
226 of both connectivity as well as community composition. A network with higher average closeness centrality  
227 also indicates an overall greater degree of connections, which suggests a greater resilience against network  
228 degradation by extinction events [30,63]. This is because more highly connected networks are less likely to  
229 degrade into multiple smaller networks when bacteria or phages are randomly removed [30,63]. We used  
230 this information to calculate the average connectedness per sample, which was corrected for the maximum  
231 potential degree of connectedness. Unfortunately our dataset was insufficiently powered to make strong  
232 conclusions toward this hypothesis, but this is an interesting observation that warrants further investigation.  
233 This observation also serves to illustrate the types of questions we can answer with more comprehensive  
234 microbiome sampling and integrative analyses.

235 Using our small sample set, we observed that the gut microbiome network structures associated with high-fat  
236 diets appeared less connected than those of low-fat diets, although a greater sample size will be required  
237 to more properly evaluate this trend (**Figure 2 A-B**). Five subjects were available for use, all of which had  
238 matching bacteria and virome datasets and samples from 8-10 days following the initiation of their diets.  
239 High-fat diets appeared to exhibit reduced degree centrality (**Figure 2 A**), suggesting bacteria in high-fat  
240 environments were targeted by fewer phages and that phage tropism was more restricted. High-fat diets  
241 also appeared to exhibit decreased closeness centrality (**Figure 2 B**), indicating that bacteria and phages  
242 were more distant from other bacteria and phages in the community. This would make genetic transfer and  
243 altered abundance of a given phage or bacterium less capable of impacting other bacteria and phages within  
244 the network.

**Figure 2: Impact of Diet and Obesity on Gut Network Structure.** (A) Quantification of average degree centrality (number of edges per node) and (B) closeness centrality (average distance from each node to every other node) of gut microbiome networks of subjects limited to exclusively high-fat or low-fat diets. Each point represents the centrality from a human subject stool sample that was collected 8-10 days following the beginning of their defined diet. There are five samples here, compared to the four in figure 3, because one of the was only sampled post-diet, providing us data for this analysis but not allowing us to compare to a baseline for figure 3. (C) Quantification of average degree centrality and (D) closeness centrality between obese and healthy adult women from the Twin gut study. Each point represents a stool sample taken from one of the three adult woman confirmed as obese or healthy and with matching virus and bacteria data.

245 In addition to diet, we observed a possible trend that obesity influenced network structure. This was done  
246 using the three mother samples available from the twin sample set, all of which had matching bacteria and  
247 phage samples and confirmed BMI information. The obesity-associated network appeared to have a higher  
248 degree centrality (**Figure 2 C**), but less closeness centrality than the healthy-associated networks (**Figure**  
249 **2 D**). These results suggested that the obesity-associated networks may be less connected. This again  
250 comes with the caveat that this is only an opportunistic observation using an existing sample set with too few  
251 samples to make more substantial claims. We included this observation as a point of interest, given the data  
252 was available.

## 253 **Individuality of Microbial Networks**

254 Skin and gut community membership and diversity are highly personal, with people remaining more similar  
255 to themselves than to other people over time [13,64,65]. We therefore hypothesized that this personal  
256 conservation extended to microbiome network structure. We addressed this hypothesis by calculating  
257 the degree of dissimilarity between each subject's network, based on phage and bacteria abundance  
258 and centrality. We quantified phage and bacteria centrality within each sample graph using the weighted  
259 eigenvector centrality metric. This metric defines central phages as those that are highly abundant ( $A_O$  as  
260 defined in the methods) and infect many distinct bacteria which themselves are abundant and infected by  
261 many other phages. Similarly, bacterial centrality was defined as those bacteria that were both abundant  
262 and connected to numerous phages that were themselves connected to many bacteria. We then calculated

263 the similarity of community networks using the weighted eigenvector centrality of all nodes between all  
264 samples. Samples with similar network structures were interpreted as having similar capacities for network  
265 robustness and transmitting genetic material.

266 We used this network dissimilarity metric to test whether microbiome network structures were more similar  
267 within people than between people over time. We found that gut microbiome network structures clustered by  
268 person (ANOSIM p-value = 0.008,  $R = 1$ , **Figure 3 A**). Network dissimilarity within each person over the 8-10  
269 day sampling period was less than the average dissimilarity between that person and others, although this  
270 difference was not statistically significant (p-value = 0.125, **Figure 3 B**). Four of the five available subjects  
271 were used because one of the subjects was not sampled at the initial time point. The lack of statistical  
272 confidence was likely due to the small sample size of this dataset.

273 Although there was evidence for gut network conservation among individuals, we found no evidence for  
274 conservation of gut network structures within families. The gut network structures were not more similar  
275 within families (twins and their mothers; intrafamily) compared to other families (other twins and mothers;  
276 inter-family) (p-value = 0.547, **Figure 3 C**). In addition to the gut, skin microbiome network structure was  
277 conserved within individuals (p-value < 0.001, **Figure 3 D**). This distribution was similar when separated by  
278 anatomical sites. Most sites were statistically significantly more conserved within individuals (**Supplemental**  
279 **Figure S12**).

280 As an additional validation measure, we evaluated the tolerance of these findings to inaccuracies in the  
281 underlying network. As described above, our model is not perfect and there is likely to be noise from false  
282 positive and false negative predictions. We found that additional random noise, both by creating a fully  
283 connected graph or randomly reducing the number of edges to 60% of the original, changed the statistical  
284 significance values (p-values) of our findings but not by enough to change whether they were statistically  
285 significant (p-value < 0.05). Therefore the comparisons between groups are resilient to potential noise  
286 resulting from model false positive and false negative predictions (**Supplemental Figure S13**).

**Figure 3: Intrapersonal vs Interpersonal Network Dissimilarity Across Different Human Systems.** (A) NMDS ordination illustrating network dissimilarity between subjects over time. Each sample is colored by subject, with each colored sample pair collected 8-10 days apart. Dissimilarity was calculated using the Bray-Curtis metric based on abundance weighted eigenvector centrality signatures, with a greater distance representing greater dissimilarity in bacteria and phage centrality and abundance. Only four subjects were included here, compared to the five used in figure 2, because one of the subjects was missing the initial sampling time point and therefore lacked temporal sampling. (B) Quantification of gut network dissimilarity within the same subject over time (intrapersonal) and the mean dissimilarity between the subject of interest and all other subjects (interpersonal). The *p*-value is provided near the bottom of the figure. (C) Quantification of gut network dissimilarity within subjects from the same family (intrafamily) and the mean dissimilarity between subjects within a family and those of other families (interfamily). Each point represents the inter-family and intra-family dissimilarity of a twin or mother that was sampled over time. (D) Quantification of skin network dissimilarity within the same subject and anatomical location over time (intrapersonal) and the mean dissimilarity between the subject of interest and all other subjects at the same time and the same anatomical location (interpersonal). All *p*-values were calculated using a paired Wilcoxon test.

## 287 **Network Structures Across the Human Skin Landscape**

288 Extensive work has illustrated differences in diversity and composition of the healthy human skin microbiome  
289 between anatomical sites, including bacteria, virus, and fungal communities [13,47,64]. These communities  
290 vary by degree of skin moisture, oil, and environmental exposure; features which were defined in the original  
291 publication [13]. As viruses are known to influence microbial diversity and community composition, we  
292 hypothesized that these differences would still be evident after integrating the bacterial and viral datasets  
293 and evaluating their microbe-virus network structure between anatomical sites. To test this, we evaluated  
294 the changes in network structure between anatomical sites within the skin dataset. The anatomical sites and  
295 their features (e.g. moisture & occlusion) were defined in the previous publication through consultation with  
296 dermatologists and reference to previous literature [13].

297 The average centrality of each sample was quantified using the weighted eigenvector centrality metric.  
298 Intermittently moist skin sites (dynamic sites that fluctuate between being moist and dry) were significantly  
299 more connected than the moist and sebaceous environments (*p*-value < 0.001, **Figure 4 A**). Also, skin sites  
300 that were occluded from the environment were less connected than those that were constantly exposed  
301 to the environment or only intermittently occluded (*p*-value < 0.001, **Figure 4 B**). We also confirmed that  
302 addition of noise to the underlying network, as described above, altered the values of statistical significance

303 (p-values) but not by enough to change whether they were statistically significant (**Supplemental Figure**  
304 **S14**).

**Figure 4: Impact of Skin Micro-Environment on Microbiome Network Structure.** (A) *Notched box-plot depicting differences in average eigenvector centrality between moist, intermittently moist, and sebaceous skin sites and (B) occluded, intermittently occluded, and exposed sites. Notched box-plots were created using ggplot2 and show the median (center line), the inter-quartile range (IQR; upper and lower boxes), the highest and lowest value within  $1.5 * IQR$  (whiskers), outliers (dots), and the notch which provides an approximate 95% confidence interval as defined by  $1.58 * IQR / \sqrt{n}$ . Sample sizes for each group were: Moist = 81, Sebaceous = 56, IntMoist = 56, Occluded = 106, Exposed = 61, IntOccluded = 26. (C) NMDS ordination depicting the differences in skin microbiome network structure between skin moisture levels and (D) occlusion. Samples are colored by their environment and their dissimilarity to other samples was calculated as described in figure 3. (E) The statistical differences of networks between moisture and (F) occlusion status were quantified with an anova and post hoc Tukey test. Cluster centroids are represented by dots and the extended lines represent the associated 95% confidence intervals. Significant comparisons ( $p$ -value < 0.05) are colored in red, and non-significant comparisons are gray.*

305 To supplement this analysis, we compared the network signatures using the centrality dissimilarity approach  
306 described above. The dissimilarity between samples was a function of shared relationships, degree of  
307 centrality, and bacteria/phage abundance. When using this supplementary approach, we found that network  
308 structures significantly clustered by moisture, sebaceous, and intermittently moist status (**Figure 4 C,E**).  
309 Occluded sites were significantly different from exposed and intermittently occluded sites, but there was no  
310 difference between exposed and intermittently occluded sites (**Figure 4 D,F**). These findings provide further  
311 support that skin microbiome network structure differs significantly between skin sites.

## 312 **Discussion**

313 Foundational work has provided a baseline understanding of the human microbiome by characterizing  
314 bacterial and viral diversity across the human body [13,14,46–48,66]. Here we integrated the bacterial and  
315 viral sequence sets to offer an initial understanding of how phage-bacteria networks differ throughout the  
316 human body, so as to provide a baseline for future studies of how and why microbiome networks differ in  
317 disease states. We implemented a network-based analytical model to evaluate the basic properties of the  
318 human microbiome through bacteria and phage relationships, instead of membership or diversity alone.

319 This approach enabled the application of network theory to provide a new perspective while analyzing  
320 bacterial and viral communities simultaneously. We utilized metrics of connectivity to model the extent to  
321 which communities of bacteria and phages interact through mechanisms such as horizontal gene transfer,  
322 modulated bacterial gene expression, and alterations in abundance.

323 Just as gut microbiome and virome composition and diversity are conserved in individuals [13,46,47,65], gut  
324 and skin microbiome network structures were conserved within individuals over time. Gut network structure  
325 was not conserved among family members. These findings suggested that the community properties inferred  
326 from microbiome interaction network structures, such as robustness (meaning a more highly connected  
327 network is more “robust” to network degradation because a randomly removed bacteria or phage node is less  
328 likely to divide or disintegrate [30,63] the overall network), the potential for horizontal gene transfer between  
329 members, and co-evolution of populations, were person-specific. These properties may be impacted by  
330 personal factors ranging from the body’s immune system to external environmental conditions, such as  
331 climate and diet.

332 We observed evidence supporting the ability of environmental conditions to shape gut and skin microbiome  
333 interaction network structure by observing that diet and skin location were associated with altered network  
334 structures. We observed evidence that diet was sufficient to alter gut microbiome network connectivity,  
335 although this needs to be interpreted cautiously as a case observation, due to the small sample size. Although  
336 the available sample size was small, our findings provide some preliminary evidence that high-fat diets are  
337 less connected than low-fat diets and that high-fat diets may therefore lead to less robust communities  
338 with a decreased ability for microbes to directly influence one another. We supported this finding with the  
339 observation that obesity may have been associated with decreased network connectivity. Together these  
340 findings suggest the food we eat may not only impact which microbes colonize our guts, but may also impact  
341 their interactions with infecting phages. Further work will be required to characterize these relationships with  
342 a larger cohort.

343 In addition to diet, the skin environment also influenced the microbiome interaction network structure.  
344 Network structure differed between environmentally exposed and occluded skin sites. The sites under



345 greater environmental fluctuation and exposure (the exposed and intermittently exposed sites) were more  
346 connected and therefore were predicted to have a higher resilience against network degradation when  
347 random nodes are removed from the network. Likewise, intermittently moist sites demonstrated higher  
348 connectedness than the moist and sebaceous sites. These findings agree with previous work that has shown  
349 that bacterial community networks differ by skin environment types [57]. Together these data suggested  
350 that body sites under greater degrees of fluctuation harbored more highly connected microbiomes that are  
351 potentially more robust to network disruption by extinction events. This points to a link between microbiome  
352 and environmental robustness toward network homeostasis and warrants further investigation.

353 While these findings take us an important step closer to understanding the microbiome through interspecies  
354 relationships, there are caveats and considerations to our findings. First, as with most classification models,  
355 the infection classification model developed and applied is only as good as its training set – in this case, the  
356 collection of experimentally-verified positive and negative infection data. Large-scale experimental screens  
357 for phage and bacteria infectious interactions that report high-confidence negative interactions (i.e., no  
358 infection) are desperately needed, as they would provide more robust model training and improved model  
359 performance. Furthermore, just as we have improved on previous modeling efforts, we expect that new and  
360 creative scoring metrics will improve future performance. Other creative and high performing models are  
361 currently being developed and the applications of these models to community network creation will continue  
362 to move this field forward [43–45].

363 Second, although our analyses utilized the best datasets currently available for our study, this work was done  
364 retrospectively and relied on existing data up to seven years old. These archived datasets were limited by  
365 the technology and costs of the time. For example, the diet and twin studies, relied on multiple displacement  
366 amplification (MDA) in their library preparations—an approach used to overcome the large nucleic acids  
367 requirements typical of older sequencing library generation protocols. It is now known that MDA results  
368 in biases in microbial community composition [67], as well as toward ssDNA viral genomes [68,69], thus  
369 rendering the resulting microbial and viral metagenomes largely non-quantitative. Future work that employs  
370 larger sequence datasets and that avoids the use of bias-inducing amplification steps will build on and validate

371 our findings, as well as inform the design and interpretation of further studies.

372 Although our models demonstrated satisfactory accuracy and overall performance, it was important to  
373 interpret our findings under the realization that our model was not perfect. This caveat is not new to the  
374 microbiome field, with a notable example being the use of 16S rRNA sequencing using the V4 variable  
375 region [59]. Use of the V4 variable region excluded detection of major skin bacterial members, meaning that  
376 the findings were not able to completely describe the underlying biological environment. Despite this caveat,  
377 skin microbiome studies provided valuable biological insights by focusing on the community differences  
378 between groups (e.g. disease and healthy) which were both analyzed the same way. Similarly, here we  
379 focused our conclusions on the differences between the groups which were all treated the same, so that we  
380 can minimize our dependence on a perfect predictive model. We also provided explicit evidence that the  
381 introduction of noise equally to the compared groups did not significantly impact our findings.

382 Third, the networks in this study were built using operational genomic units (OGUs), which represented  
383 groups of highly similar bacteria or phage genomes or clustered genome fragments. Similar clustering  
384 definition and validation methods, both computational and experimental, have been implemented in other  
385 metagenomic sequencing studies, as well [42,70–72]. These approaches could offer yet another level of  
386 sophistication to our network-based analyses. While this operationally defined clustering approach allows  
387 us to study whole community networks, our ability to make conclusions about interactions among specific  
388 phage or bacterial species or populations is inherently limited, compared to more focused, culture-based  
389 studies such as the work by Malki *et al* [60]. Future work must address this limitation, e.g., through improved  
390 binning methods and deeper metagenomic shotgun sequencing, but most importantly through an improved  
391 conceptual framing of what defines ecologically and evolutionarily cohesive units for both phage and bacteria  
392 [73]. Defining operational genomic units and their taxonomic underpinnings (e.g., whether OGU clusters  
393 represent genera or species) is an active area of work critical to the utility of this approach. As a first  
394 step, phylogenomic analyses have been performed to cluster cyanophage isolate genomes into informative  
395 groups using shared gene content, average nucleotide identity of shared genes, and pairwise differences  
396 between genomes [74]. Such population-genetic assessment of phage evolution, coupled with the ecological

397 implications of genome heterogeneity, will inform how to define nodes in future iterations of the ecological  
398 network developed here. Even though we are hesitant to speculate on phage host ranges at low taxonomic  
399 levels in our dataset, the data does agree with previous reports of instances of broad phage host range  
400 [60,75]. Additionally, visualization of our dataset interactions using the heat map approach previously used  
401 in other host range studies, suggests a trend toward modular and nested tropism, but we do not have the  
402 strain-level resolution that powered those previous experimental studies.

403 Finally, it is important to note that our model was built using available full genomes with known interactions,  
404 while the experimental datasets resulted in OGUs created from metagenomic shotgun sequence sets, as  
405 described above. While this is an informative approach given available data, it is not ideal. We envision  
406 future work focusing on training models using metagenomic shotgun sample sets from “mock communities”  
407 of bacteria and phages with experimentally validated infectious relationships. This would also be more  
408 informative than relying on simulated metagenomic sample sets, whose use would result in models built  
409 on simulations and more assumptions instead of empirical data. Together this way the training set can be  
410 subjected to the same pre-processing, contig assembly, and OGU binning processes as the experimental  
411 data. Furthermore, exciting advances in long read sequencing platforms such as the Oxford Nanopore  
412 Minlon system will provide more accurate genomic scaffolds than *de novo* assembled contigs, allowing for  
413 more accurate training and predictions of our models. As discussed above, it is because our current model  
414 is susceptible to this noise that we focus our conclusions on comparisons between experimental groups that  
415 were both treated the same. This is also why it was important for us to evaluate the susceptibility of our  
416 results to noise caused by the less-than-perfect prediction model.

417 Together our work takes an initial step towards defining bacteria-virus interaction profiles as a characteristic  
418 of human-associated microbial communities. This approach revealed the impacts that different human  
419 environments (e.g., the skin and gut) can have on microbiome connectivity. By focusing on relationships  
420 between bacterial and viral communities, they are studied as the interacting cohorts they are, rather than  
421 as independent entities. While our developed bacteria-phage interaction framework is a novel conceptual  
422 advance, the microbiome also consists of archaea and small eukaryotes, including fungi and *Demodex* mites

423 [1,76] – all of which can interact with human immune cells and other non-microbial community members [77].  
424 Future work will build from our approach and include these additional community members and their diverse  
425 interactions and relationships (e.g., beyond phage-bacteria). This will result in a more robust network and a  
426 more holistic understanding of the evolutionary and ecological processes that drive the assembly and function  
427 of the human-associated microbiome.

## 428 **Materials & Methods**

### 429 **Code Availability**

430 A reproducible version of this manuscript written in R markdown and all of the code used to obtain and  
431 process the sequencing data is available at the following GitHub repository:

432 [https://github.com/SchlossLab/Hannigan\\_ConjunctisViribus\\_ploscompbio\\_2017](https://github.com/SchlossLab/Hannigan_ConjunctisViribus_ploscompbio_2017)

### 433 **Data Acquisition & Quality Control**

434 Raw sequencing data and associated metadata were acquired from the NCBI sequence read archive (SRA).  
435 Supplementary metadata were acquired from the same SRA repositories and their associated manuscripts.  
436 The gut virome diet study (SRA: SRP002424), twin virome studies (SRA: SRP002523; SRP000319), and  
437 skin virome study (SRA: SRP049645) were downloaded as `.sra` files. For clarity, the sample sizes used  
438 for each study subset were described with the data in the results section. Sequencing files were converted  
439 to `fastq` format using the `fastq-dump` tool of the NCBI SRA Toolkit (v2.2.0). Sequences were quality  
440 trimmed using the `Fastx` toolkit (v0.0.14) to exclude bases with quality scores below 33 and shorter than 75  
441 bp [78]. Paired end reads were filtered to exclude sequences missing their corresponding pair using the  
442 `get_trimmed_pairs.py` script available in the source code.

## 443 **Contig Assembly**

444 Contigs were assembled using the Megahit assembly program (v1.0.6) [79]. A minimum contig length of 1  
445 kb was used. Iterative k-mer stepping began at a minimum length of 21 and progressed by 20 until 101. All  
446 other default parameters were used.

447 Contig simulations were performed by randomly extracting a string of genomic nucleotides that represented a  
448 defined percent length of that genome. This was accomplished using our `RandomContigGenerator.pl`,  
449 which was published in the associated GitHub repository.

## 450 **Contig Abundance Calculations**

451 Contigs were concatenated into two master files prior to alignment, one for bacterial contigs and one for  
452 phage contigs. Sample sequences were aligned to phage or bacterial contigs using the Bowtie2 global aligner  
453 (v2.2.1) [80]. We defined a mismatch threshold of 1 bp and seed length of 25 bp. Sequence abundance was  
454 calculated from the Bowtie2 output using the `calculate_abundance_from_sam.pl` script available in  
455 the source code.

## 456 **Operational Genomic Unit Binning**

457 Contigs often represent large fragments of genomes. In order to reduce redundancy and the resulting  
458 artificially inflated genomic richness within our dataset, it was important to bin contigs into operational  
459 units based on their similarity. This approach is conceptually similar to the clustering of related 16S rRNA  
460 sequences into operational taxonomic units (OTUs), although here we are clustering contigs into operational  
461 genomic units (OGUs) [66].

462 Contigs were clustered using the CONCOCT algorithm (v0.4.0) [81]. Because of our large dataset and limits  
463 in computational efficiency, we randomly subsampled the dataset to include 25% of all samples, and used  
464 these to inform contig abundance within the CONCOCT algorithm. CONCOCT was used with a maximum

465 of 500 clusters, a k-mer length of four, a length threshold of 1 kb, 25 iterations, and exclusion of the total  
466 coverage variable.

467 OGU abundance ( $A_O$ ) was obtained as the sum of the abundance of each contig ( $A_j$ ) associated with that  
468 OGU. The abundance values were length corrected such that:

$$A_O = \frac{10^7 \sum_{j=1}^k A_j}{\sum_{j=1}^k L_j}$$

469 Where L is the length of each contig j within the OGU.

## 470 **Operational Genomic Unit Identification**

471 To confirm a lack of phage sequences in the bacterial OGU dataset, we performed blast nucleotide alignment  
472 of the bacterial OGU representative sequences using an e-value  $< 10^{-25}$ , which was stricter than the  $10^{-10}$   
473 threshold used in the random forest model below, against all of the phage reference genomes available in  
474 the EMBL database. We used a stricter threshold because we know there are genomic similarities between  
475 bacteria and phage OGUs from the interactive model, but we were interested in contigs with high enough  
476 similarity to references that they may indeed be from phages. We also performed the converse analysis  
477 of aligning phage OGU representative sequences to EMBL bacterial reference genomes. We ran both the  
478 phage and bacteria OGU representative sequences through the Virsorter program (1.0.3) to identify phages  
479 (all default parameters were used), using only those in the high confidence identification category “class 1”  
480 [50]. Finally, we filtered out phage OGUs that had bacterial elements as described above, but also lacked  
481 known phage elements by using the tblastx algorithm and a maximum e-value of  $10^{-25}$ .

## 482 **Open Reading Frame Prediction**

483 Open reading frames (ORFs) were identified using the Prodigal program (V2.6.2) with the meta mode  
484 parameter and default settings [82].

## 485 **Classification Model Creation and Validation**

486 The classification model for predicting interactions was built using experimentally validated bacteria-phage  
487 infections or validated lack of infections from six studies [41,51–55]. No further reference databases were  
488 used in our alignment procedures. Associated reference genomes were downloaded from the European  
489 Bioinformatics Institute (see details in source code). The model was created based on the four metrics listed  
490 below.

491 The four scores were used as parameters in a random forest model to classify bacteria and bacteriophage  
492 pairs as either having infectious interactions or not. The classification model was built using the Caret  
493 R package (v6.0.73) [83]. The model was trained using five-fold cross validation with ten repeats, and  
494 the median model performance was evaluated by training the model on 80% of the dataset and testing  
495 performance on the remaining 20%. Pairs without scores were classified as not interacting. The model was  
496 optimized using the ROC value. The resulting model performance was plotted using the plotROC R package.

## 497 **Identify Bacterial CRISPRs Targeting Phages**

498 Clustered Regularly Interspaced Short Palindromic Repeats (CRISPRs) were identified from bacterial  
499 genomes using the PilerCR program (v1.06) [84]. Resulting spacer sequences were filtered to exclude  
500 spacers shorter than 20 bp and longer than 65 bp. Spacer sequences were aligned to the phage genomes  
501 using the nucleotide BLAST algorithm with default parameters (v2.4.0) [85]. The mean percent identity for  
502 each matching pair was recorded for use in our classification model.

## 503 **Detect Matching Prophages within Bacterial Genomes**

504 Temperate bacteriophages infect and integrate into their bacterial host's genome. We detected integrated  
505 phage elements within bacterial genomes by aligning phage genomes to bacterial genomes using the  
506 nucleotide BLAST algorithm and a minimum e-value of  $1e-10$ . The resulting bitscore of each alignment was  
507 recorded for use in our classification model.

## 508 **Identify Shared Genes Between Bacteria and Phages**

509 As a result of gene transfer or phage genome integration during infection, phages may share genes with  
510 their bacterial hosts, providing us with evidence of phage-host pairing. We identified shared genes between  
511 bacterial and phage genomes by assessing amino acid similarity between the genes using the Diamond  
512 protein alignment algorithm (v0.7.11.60) [86]. The mean alignment bitscores for each genome pair were  
513 recorded for use in our classification model.

## 514 **Protein - Protein Interactions**

515 The final method used for predicting infectious interactions between bacteria and phages was the detection  
516 of pairs of genes whose proteins are known to interact. We assigned bacterial and phage genes to protein  
517 families by aligning them to the Pfam database using the Diamond protein alignment algorithm. We then  
518 identified which pairs of proteins were predicted to interact using the Pfam interaction information within the  
519 Intact database [56]. The mean bitscores of the matches between each pair were recorded for use in the  
520 classification model.

## 521 **Secondary Dataset Validation**

522 The performance of our model for identifying diverse infectious relationships between bacteria and phages,  
523 beyond those that were included in the model creation step, were validated using additional bacterial and  
524 phage reference genomes, which could be linked by the records of which phage strains were isolated on  
525 which bacteria under laboratory conditions. Viral and bacterial reference genomes were downloaded from the  
526 GenBank repository on February 19, 2018 using the viral location `ftp://ftp.ncbi.nih.gov/refseq/release/viral`  
527 and the bacterial location `ftp://ftp.ncbi.nih.gov/refseq/release/bacteria/`. This resulted  
528 in the use of 539 complete phages reference genomes (with identified hosts) and 3,469 bacterial reference  
529 genomes. We used the same prediction model to predict which phages were infecting which hosts, so  
530 as to confirm that the model was capable of identifying interactions in a more diverse dataset. Bacteria  
531 interactions were identified at the species level. The random contig iteration analysis was performed using a



532 subset of bacterial reference genomes, for computational performance reasons. Only single representative  
533 genomes for each species were used.

## 534 **Interaction Network Construction**

535 The bacteria and phage operational genomic units (OGUs) were scored using the same approach as outlined  
536 above. The infectious pairings between bacteria and phage OGUs were classified using the random forest  
537 model described above. The predicted infectious pairings and all associated metadata were used to populate  
538 a graph database using Neo4j graph database software (v2.3.1) [87]. This network was used for downstream  
539 community analysis. Tolerance to false negative and false positive noise within the networks was assessed  
540 by randomly removing a defined fraction of network edges before re-running the downstream analysis work  
541 flows. This was accomplished using functionality within the igraph R package (v1.0.1) [88].

## 542 **Centrality Analysis**

543 We quantified the centrality of graph vertices using three different metrics, each of which provided different  
544 information graph structure. When calculating these values, let  $G(V, E)$  be an undirected, unweighted graph  
545 with  $|V| = n$  nodes and  $|E| = m$  edges. Also, let  $\mathbf{A}$  be its corresponding adjacency matrix with entries  
546  $a_{ij} = 1$  if nodes  $V_i$  and  $V_j$  are connected via an edge, and  $a_{ij} = 0$  otherwise.

547 Briefly, the **closeness centrality** of node  $V_i$  is calculated taking the inverse of the average length of the  
548 shortest paths ( $d$ ) between nodes  $V_i$  and all the other nodes  $V_j$ . Mathematically, the closeness centrality of  
549 node  $V_i$  is given as:

$$C_C(V_i) = \left( \sum_{j=1}^n d(V_i, V_j) \right)^{-1}$$

550 The distance between nodes ( $d$ ) was calculated as the shortest number of edges required to be traversed  
551 to move from one node to another.

552 Intuitively, the **degree centrality** of node  $V_i$  is defined as the number of edges that are incident to that node:

$$C_D(V_i) = \sum_{j=1}^n a_{ij}$$

553 where  $a_{ij}$  is the  $ij^{th}$  entry in the adjacency matrix  $\mathbf{A}$ .

554 The eigenvector centrality of node  $V_i$  is defined as the  $i^{th}$  value in the first eigenvector of the associated  
555 adjacency matrix  $\mathbf{A}$ . Conceptually, this function results in a centrality value that reflects the connections of  
556 the vertex, as well as the centrality of its neighboring vertices.

557 The **centralization** metric was used to assess the average centrality of each sample graph  $\mathbf{G}$ . Centralization  
558 was calculated by taking the sum of each vertex  $V_i$ 's centrality from the graph maximum centrality  $C_w$ , such  
559 that:

$$C(G) = \frac{\sum_{i=1}^n C_w - c(V_i)}{T}$$

560 The values were corrected for uneven graph sizes by dividing the centralization score by the maximum  
561 theoretical centralization (T) for a graph with the same number of vertices.

562 Degree and closeness centrality were calculated using the associated functions within the igraph R package  
563 (v1.0.1) [88].

## 564 **Network Relationship Dissimilarity**

565 We assessed similarity between graphs by evaluating the shared centrality of their vertices, as has been  
566 done previously. More specifically, we calculated the dissimilarity between graphs  $G_i$  and  $G_j$  using the  
567 Bray-Curtis dissimilarity metric and eigenvector centrality values such that:

$$B(G_i, G_j) = 1 - \frac{2C_{ij}}{C_i + C_j}$$

568 Where  $C_{ij}$  is the sum of the lesser centrality values for those vertices shared between graphs, and  $C_i$  and  
569  $C_j$  are the total number of vertices found in each graph. This allows us to calculate the dissimilarity between  
570 graphs based on the shared centrality values between the two graphs.

## 571 **Statistics and Comparisons**

572 Differences in intrapersonal and interpersonal network structure diversity, based on multivariate data,  
573 were calculated using an analysis of similarity (ANOSIM). Statistical significance of univariate Eigenvector  
574 centrality differences were calculated using a paired Wilcoxon test.

575 Statistical significance of differences in univariate eigenvector centrality measurements of skin virome-microbiome  
576 networks were calculated using a pairwise Wilcoxon test, corrected for multiple hypothesis tests using the  
577 Holm correction method. Multivariate eigenvector centrality was measured as the mean differences between  
578 cluster centroids, with statistical significance measured using an ANOVA and post hoc Tukey test.

## 579 **Acknowledgments**

580 We thank the members of the Schloss lab for their underlying contributions. We thank the authors of the  
581 original studies for making their data and metadata publicly available and understandable. We also thank  
582 the participants in the studies.

## 583 **Author Contributions**

584 *Conceptualization:* GDH, MBD, DK, PDS. *Data Curation:* GDH. *Formal Analysis:* GDH. *Funding Acquisition:*  
585 GDH, PDS. *Writing – Original Draft Preparation:* GDH, PDS. *Writing – Review & Editing:* GDH, MBD, DK,  
586 PDS.

## 587 **Funding Information**

588 GDH was supported in part by the Molecular Mechanisms in Microbial Pathogenesis Training Program (T32  
589 AI007528). GDH and PDS were supported in part by funding from the NIH (P30DK034933, U19AI09087,  
590 and U01AI124255).

## 591 **Competing interests**

592 The authors report no conflicts of interest.

## 593 **Supplemental Figure Captions**

Figure S1: **Sequencing Depth Summary.** *Number of sequences that aligned to (A) Phage and (B) Bacteria operational genomic units per sample and colored by study.*

Figure S2: **Contig Summary Statistics.** *Scatter plot heat map with each hexagon representing the abundance of contigs. Contigs are organized by length on the x-axis and the number of aligned sequences on the y-axis.*

**Figure S3: Operational Genomic Unit Summary Statistics.** *Scatter plot with operational genomic unit clusters organized by average contig length within the cluster on the x-axis and the number of contigs in the cluster on the y-axis. Operational genomic units of (A) bacteriophages and (B) bacteria are shown.*

**Figure S4: Summary information of validation dataset used in the interaction predictive model.** *A) Categorical heat-map highlighting the experimentally validated positive and negative interactions. Only bacteria species are shown, which represent multiple reference strains. Phages are labeled on the x-axis and bacteria are labeled on the y-axis. B) Quantification of bacterial host strains known to exist for each phage. C) Genome strandedness and D) linearity of the phage reference genomes used for the dataset.*



**Figure S5: Ability of prediction model to identify broad range of bacteria and phage interactions.** *Each complete bacteriophage labeled on the y axis. The number of complete bacterial species genomes that were correctly predicted to be infected by each phage, according to the GenBank records, is on the x-axis. True positive interactions are colored in purple, and false negative interactions are yellow. Because this dataset only had confirmed interactions and not confirmed lack of interactions, false positive and true negative values could be not determined.*

**Figure S6: Impact of incomplete genomic sequences on model performance.** *Number of correctly identified infectious interactions between bacteria and phages are presented on the y-axis. The fraction of genomic length that was used to randomly extract contigs from reference bacterial and phage sequences is presented on the x-axis (e.g. 0.9 means contig lengths were 90% of the total genome length). Overall this presents a quantification of the loss of identified infectious relationships as the percent of available genomic material is reduced.*

**Figure S7: Classification Model Performance By Nested Cross-Validation.** *Box plot illustrating the median and variance of phage-bacteria interaction prediction model. Performance was evaluated using nested cross validation, meaning that 20% of the samples were randomly withheld from model training and then used to evaluate performance. The results of 100 random iterations are shown. Metrics include area under the curve (gray), sensitivity (red), and specificity (tan).*

**Figure S8: Stable Classification Model Performance Over Random Iterations** *In addition to nested cross-validation, here we show the results from the five-fold cross validation, in which 20% of the samples were randomly withheld during the training stage for model evaluation and mtry tuning (parameter defined in the R Random Forest package, which is implemented in Caret, as “the number of variables randomly sampled as candidates at each split”). The results of 25 random iterations are shown. Metrics include area under the curve (red), sensitivity (green), and specificity (blue). Dashed line highlight the random point of 0.5.*

**Figure S9: Structure of the interactive network.** *Metadata relationships to samples (Phage Sample ID and Bacteria Sample ID) included the associated time point, the study, the subject the sample was taken from, and the associated disease. Infectious interactions were recorded between phage and bacteria operational genomic units (OGUs). Sequence count abundance for each OGU within each sample was also recorded.*

**Figure S10: Heatmap of Phage-Bacteria Interaction Relationships of Master Network.** *Heatmap illustrating the ranges of infectious interactions predicted between bacteria and bacteriophages across our three studies. Bacterial OGUs are aligned on the vertical access, and the bacteriophage OGUs are organized on the horizontal access. OGUs are organized near other OGUs with similar infectious profiles, which are further illustrated by the dendrograms. Predicted infections are tan and predicted lacks of interactions are red.*

**Figure S11: Distribution of node eccentricity across subnetworks.** *Histograms illustrating the distributions of node eccentricity values across the subnetworks, for supplementing the node, edge, and diameter values provided for the networks. Eccentricity of each node is the shortest distance of that node to the furthest other node within the graph.*

**Figure S12: Intrapersonal vs Interpersonal Dissimilarity of the Skin.** *Quantification of skin network dissimilarity within the same subject and anatomical location over time (intrapersonal) and the mean dissimilarity between the subject of interest and all other subjects at the same time and the same anatomical location (interpersonal), separated by each anatomical site (forehead [Fh], palm [Pa], toe web [Tw], umbilicus [Um], antecubital fossa [Ac], axilla [Ax], and retroauricular crease [Ra]). P-value was calculated using a paired Wilcoxon test.*



**Figure S13: P-values of interpersonal group diversity differences with graph edge noise.** *The x-axis represents the percent of errors that were randomly removed (or added) as a means to evaluate the impact of random noise on statistical significance of group differences. Resulting p-values for each graph is shown on the y-axis. The dot and bars are the mean and standard error of five iterations of group comparisons with random edge removal. Significance of A) ANOSIM p-value of diet network dissimilarity, B) p-value of interpersonal and intrapersonal diet network dissimilarity, C) p-value of interpersonal and intrapersonal skin network dissimilarity, and D) p-value of interpersonal and intrapersonal twin gut network dissimilarity. This corresponds to the findings in Figure 3.*

**Figure S14: P-values of differences in Eigen Centrality between skin site microbiome networks.** *The x-axis represents the percent of errors that were randomly removed (or added) as a means to evaluate the impact of random noise on statistical significance of group differences. Resulting p-values for each graph is shown on the y-axis. The dot and bars are the mean and standard error of five iterations of group comparisons with random edge removal. The groups compared were the degrees of site moisture (left) and occlusion (right). The findings correspond to pannels A and B in Figure 4.*

## 594 **Supplemental Table Captions**

595 Table S1: Summary of the primary quality control measures reported in the original publications of the viromes  
596 used in this current study.

597 Table S2: The positive and negative bacteria and bacteriophage interactions used to train the prediction  
598 model, as also illustrated in Figure S4. Citation sources are also included.

## 599 **References**

- 600 1. Hannigan GD, Grice EA. Microbial Ecology of the Skin in the Era of Metagenomics and Molecular  
601 Microbiology. *Cold Spring Harbor Perspectives in Medicine*. 2013;3: a015362–a015362.
- 602 2. Hannigan GD, Hodkinson BP, McGinnis K, Tyldsley AS, Anari JB, Horan AD, et al. Culture-independent  
603 pilot study of microbiota colonizing open fractures and association with severity, mechanism, location, and  
604 complication from presentation to early outpatient follow-up. *Journal of Orthopaedic Research*. 2014;32:  
605 597–605.
- 606 3. Loesche M, Gardner SE, Kalan L, Horwinski J, Zheng Q, Hodkinson BP, et al. Temporal stability in chronic  
607 wound microbiota is associated with poor healing. *Journal of Investigative Dermatology*. 2016;
- 608 4. He Q, Li X, Liu C, Su L, Xia Z, Li X, et al. Dysbiosis of the fecal microbiota in the TNBS-induced Crohn's  
609 disease mouse model. *Applied Microbiology and Biotechnology*. 2016; 1–10.
- 610 5. Norman JM, Handley SA, Baldrige MT, Droit L, Liu CY, Keller BC, et al. Disease-specific alterations in  
611 the enteric virome in inflammatory bowel disease. *Cell*. 2015;160: 447–460.
- 612 6. Seekatz AM, Rao K, Santhosh K, Young VB. Dynamics of the fecal microbiome in patients with recurrent  
613 and nonrecurrent *Clostridium difficile* infection. *Genome medicine*. 2016;8: 47.
- 614 7. Zackular JP, Rogers MAM, Ruffin MT, Schloss PD. The human gut microbiome as a screening tool for  
615 colorectal cancer. *Cancer prevention research (Philadelphia, Pa)*. 2014;7: 1112–1121.
- 616 8. Baxter NT, Zackular JP, Chen GY, Schloss PD. Structure of the gut microbiome following colonization with  
617 human feces determines colonic tumor burden. *Microbiome*. 2014;2: 20.
- 618 9. Manrique P, Bolduc B, Walk ST, Oost J van der, Vos WM de, Young MJ. Healthy human gut phageome.  
619 *Proceedings of the National Academy of Sciences of the United States of America*. 2016; 201601060.
- 620 10. Ly M, Abeles SR, Boehm TK, Robles-Sikisaka R, Naidu M, Santiago-Rodriguez T, et al. Altered Oral

- 621 Viral Ecology in Association with Periodontal Disease. *mBio*. 2014;5: e01133–14–e01133–14.
- 622 11. Modi SR, Lee HH, Spina CS, Collins JJ. Antibiotic treatment expands the resistance reservoir and  
623 ecological network of the phage metagenome. *Nature*. 2013;499: 219–222.
- 624 12. Monaco CL, Gootenberg DB, Zhao G, Handley SA, Ghebremichael MS, Lim ES, et al. Altered Virome and  
625 Bacterial Microbiome in Human Immunodeficiency Virus-Associated Acquired Immunodeficiency Syndrome.  
626 *Cell Host and Microbe*. 2016;19: 311–322.
- 627 13. Hannigan GD, Meisel JS, Tyldsley AS, Zheng Q, Hodkinson BP, SanMiguel AJ, et al. The Human Skin  
628 Double-Stranded DNA Virome: Topographical and Temporal Diversity, Genetic Enrichment, and Dynamic  
629 Associations with the Host Microbiome. *mBio*. 2015;6: e01578–15.
- 630 14. Minot S, Sinha R, Chen J, Li H, Keilbaugh SA, Wu GD, et al. The human gut virome: Inter-individual  
631 variation and dynamic response to diet. *Genome Research*. 2011;21: 1616–1625.
- 632 15. Santiago-Rodriguez TM, Ly M, Bonilla N, Pride DT. The human urine virome in association with urinary  
633 tract infections. *Frontiers in Microbiology*. 2015;6: 14.
- 634 16. Abeles SR, Ly M, Santiago-Rodriguez TM, Pride DT. Effects of Long Term Antibiotic Therapy on Human  
635 Oral and Fecal Viromes. *PLOS ONE*. 2015;10: e0134941.
- 636 17. Abeles SR, Robles-Sikisaka R, Ly M, Lum AG, Salzman J, Boehm TK, et al. Human oral viruses are  
637 personal, persistent and gender-consistent. 2014; 1–15.
- 638 18. Haerter JO, Mitarai N, Sneppen K. Phage and bacteria support mutual diversity in a narrowing staircase  
639 of coexistence. *The ISME Journal*. 2014;8: 2317–2326.
- 640 19. Lindell D, Jaffe JD, Johnson ZI, Church GM, Chisholm SW. Photosynthesis genes in marine viruses yield  
641 proteins during host infection. *Nature*. 2005;438: 86–89.
- 642 20. Tyler JS, Beerli K, Reynolds JL, Alteri CJ, Skinner KG, Friedman JH, et al. Prophage induction is  
643 enhanced and required for renal disease and lethality in an EHEC mouse model. *PLoS Pathogens*. 2013;9:

644 e1003236.

645 21. Hargreaves KR, Kropinski AM, Clokie MR. Bacteriophage behavioral ecology: How phages alter their  
646 bacterial host's habits. *Bacteriophage*. 2014;4: e29866.

647 22. Moon BY, Park JY, Hwang SY, Robinson DA, Thomas JC, Fitzgerald JR, et al. Phage-mediated horizontal  
648 transfer of a *Staphylococcus aureus* virulence-associated genomic island. *Scientific Reports*. 2015;5: 9784.

649 23. Modi SR, Lee HH, Spina CS, Collins JJ. Antibiotic treatment expands the resistance reservoir and  
650 ecological network of the phage metagenome. *Nature*. 2013;499: 219–222.

651 24. Ogg JE, Timme TL, Alemohammad MM. General Transduction in *Vibrio cholerae*. *Infection and Immunity*.  
652 1981;31: 737–741.

653 25. Frost LS, Leplae R, Summers AO, Toussaint A. Mobile genetic elements: the agents of open source  
654 evolution. *Nature Reviews Microbiology*. 2005;3: 722–732.

655 26. Koskella B, Brockhurst MA. Bacteria-phage coevolution as a driver of ecological and evolutionary  
656 processes in microbial communities. *FEMS Microbiology Reviews*. 2014;38: 916–931.

657 27. Jover LF, Effler TC, Buchan A, Wilhelm SW, Weitz JS. The elemental composition of virus particles:  
658 implications for marine biogeochemical cycles. *Nature Reviews Microbiology*. 2014;12: 519–528.

659 28. Harcombe WR, Bull JJ. Impact of phages on two-species bacterial communities. *Applied and  
660 Environmental Microbiology*. 2005;71: 5254–5259.

661 29. Middelboe M, Hagström A, Blackburn N, Sinn B, Fischer U, Borch NH, et al. Effects of Bacteriophages on  
662 the Population Dynamics of Four Strains of Pelagic Marine Bacteria. *Microbial Ecology*. 2001;42: 395–406.

663 30. Poisot T, Lepennetier G, Martinez E, Ramsayer J, Hochberg ME. Resource availability affects the  
664 structure of a natural bacteriophage community. *Biology letters*. 2011;7: 201–204.

665 31. Thompson RM, Brose U, Dunne JA, Hall RO, Hladysz S, Kitching RL, et al. Food webs: reconciling the

- 666 structure and function of biodiversity. *Trends in ecology & evolution*. 2012;27: 689–697.
- 667 32. Moebus K, Nattkemper H. Bacteriophage sensitivity patterns among bacteria isolated from marine waters.  
668 *Helgoländer Meeresuntersuchungen*. 1981;34: 375–385.
- 669 33. Flores CO, Valverde S, Weitz JS. Multi-scale structure and geographic drivers of cross-infection within  
670 marine bacteria and phages. *The ISME Journal*. 2013;7: 520–532.
- 671 34. Poisot T, Canard E, Mouillot D, Mouquet N, Gravel D. The dissimilarity of species interaction networks.  
672 *Ecology letters*. 2012;15: 1353–1361.
- 673 35. Poisot T, Stouffer D. How ecological networks evolve. *bioRxiv*. 2016;
- 674 36. Flores CO, Meyer JR, Valverde S, Farr L, Weitz JS. Statistical structure of host-phage interactions.  
675 *Proceedings of the National Academy of Sciences of the United States of America*. 2011;108: E288–97.
- 676 37. Jover LF, Flores CO, Cortez MH, Weitz JS. Multiple regimes of robust patterns between network structure  
677 and biodiversity. *Scientific Reports*. 2015;5: 17856.
- 678 38. Reyes A, Haynes M, Hanson N, Angly FE, Heath AC, Rohwer F, et al. Viruses in the faecal microbiota  
679 of monozygotic twins and their mothers. *Nature*. 2010;466: 334–338.
- 680 39. Turnbaugh PJ, Hamady M, Yatsunenko T, Cantarel BL, Duncan A, Ley RE, et al. A core gut microbiome  
681 in obese and lean twins. *Nature*. 2009;457: 480–484.
- 682 40. Lima-Mendez G, Faust K, Henry N, Decelle J, Colin S, Carcillo F, et al. Ocean plankton. Determinants  
683 of community structure in the global plankton interactome. *Science*. 2015;348: 1262073–1262073.
- 684 41. Edwards RA, McNair K, Faust K, Raes J, Dutilh BE. Computational approaches to predict bacteriophage-host  
685 relationships. *FEMS Microbiology Reviews*. 2015;40: 258–272.
- 686 42. Roux S, Brum JR, Dutilh BE, Sunagawa S, Duhaime MB, Loy A, et al. Ecogenomics and potential  
687 biogeochemical impacts of globally abundant ocean viruses. *Nature*. 2016;537: 689–693.
- 688 43. Villarroel J, Kleinheinz KA, Jurtz VI, Zschach H, Lund O, Nielsen M, et al. HostPhinder: A Phage Host

689 Prediction Tool. *Viruses*. 2016;8: 116.

690 44. Galiez C, Siebert M, Enault F, Vincent J, Söding J. WIsH: who is the host? Predicting prokaryotic hosts  
691 from metagenomic phage contigs. *Bioinformatics*. 2017;33: 3113–3114.

692 45. Ahlgren NA, Ren J, Lu YY, Fuhrman JA, Sun F. Alignment-free  $d_2^*$  oligonucleotide frequency dissimilarity  
693 measure improves prediction of hosts from metagenomically-derived viral sequences. *Nucleic Acids*  
694 *Research*. 2017;45: 39–53.

695 46. Grice EA, Kong HH, Conlan S, Deming CB, Davis J, Young AC, et al. Topographical and Temporal  
696 Diversity of the Human Skin Microbiome. *Science*. 2009;324: 1190–1192.

697 47. Findley K, Oh J, Yang J, Conlan S, Deming C, Meyer JA, et al. Topographic diversity of fungal and  
698 bacterial communities in human skin. *Nature*. 2013; 1–6.

699 48. Costello EK, Lauber CL, Hamady M, Fierer N, Gordon JI, Knight R. Bacterial community variation in  
700 human body habitats across space and time. *Science*. 2009;326: 1694–1697.

701 49. Consortium THMP. A framework for human microbiome research. *Nature*. 2012;486: 215–221.

702 50. Roux S, Enault F, Hurwitz BL, Sullivan MB. VirSorter: mining viral signal from microbial genomic data.  
703 *PeerJ*. 2015;3: e985–20.

704 51. Jensen EC, Schrader HS, Rieland B, Thompson TL, Lee KW, Nickerson KW, et al. Prevalence  
705 of broad-host-range lytic bacteriophages of *Sphaerotilus natans*, *Escherichia coli*, and *Pseudomonas*  
706 *aeruginosa*. *Applied and Environmental Microbiology*. 1998;64: 575–580.

707 52. Malki K, Kula A, Bruder K, Sible E. Bacteriophages isolated from Lake Michigan demonstrate broad  
708 host-range across several bacterial phyla. *Virology*. 2015;

709 53. Schwarzer D, Buettner FFR, Browning C, Nazarov S, Rabsch W, Bethe A, et al. A multivalent adsorption  
710 apparatus explains the broad host range of phage phi92: a comprehensive genomic and structural analysis.  
711 *Journal of Virology*. 2012;86: 10384–10398.

712 54. Kim S, Rahman M, Seol SY, Yoon SS, Kim J. *Pseudomonas aeruginosa* bacteriophage PA1Ø

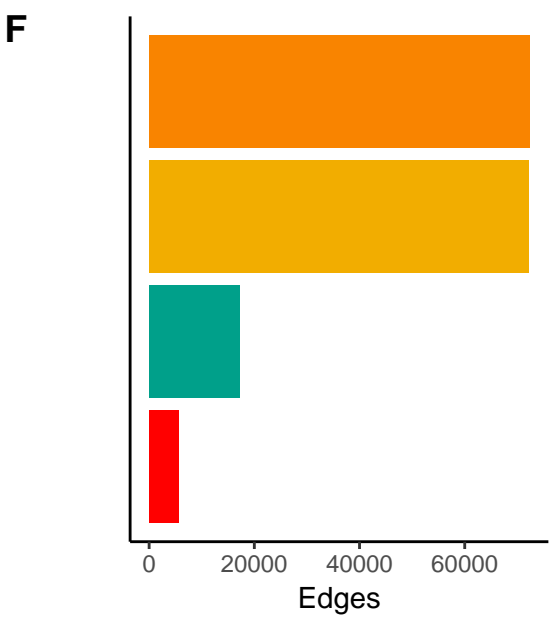
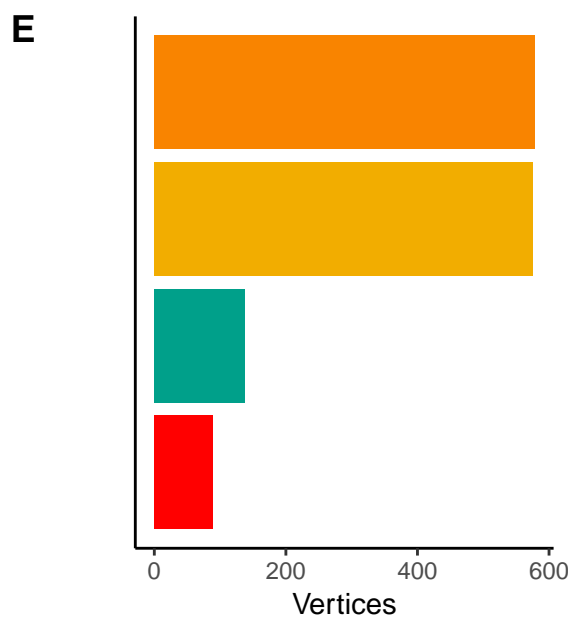
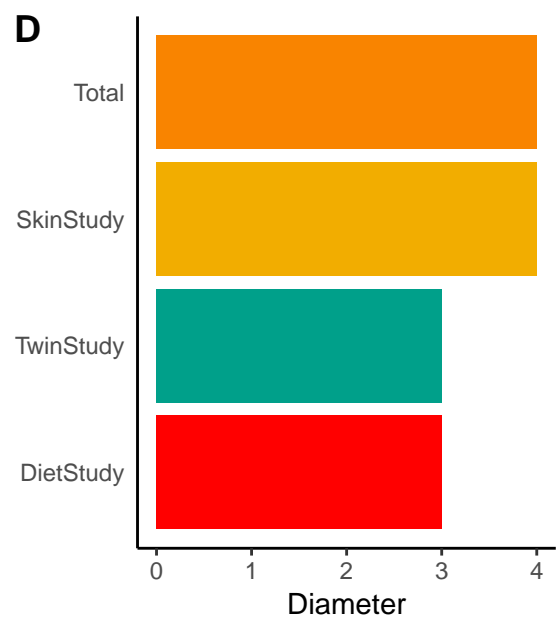
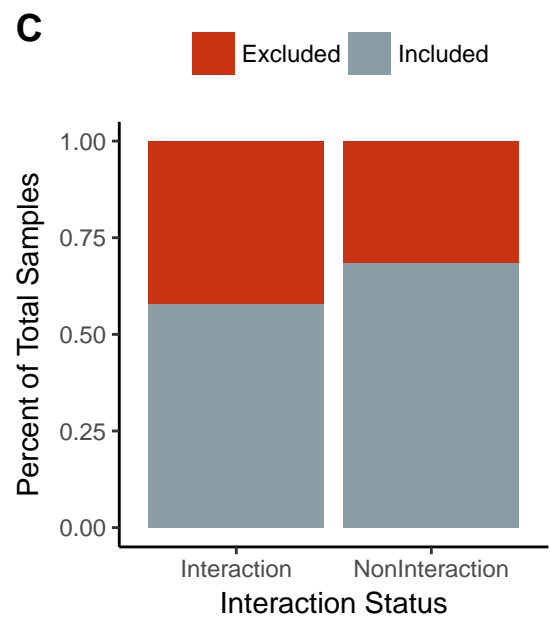
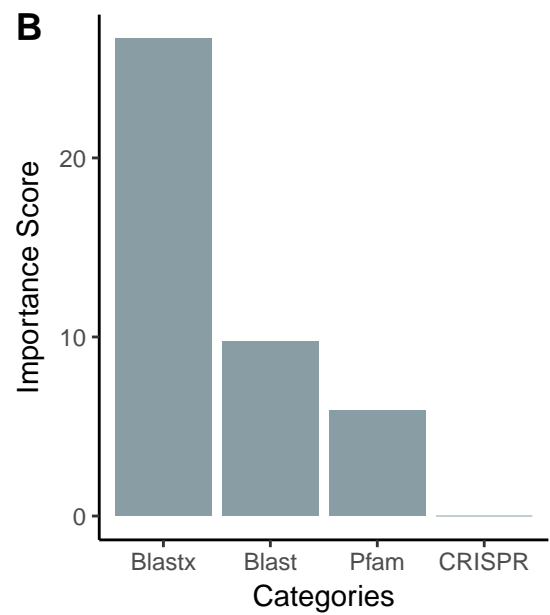
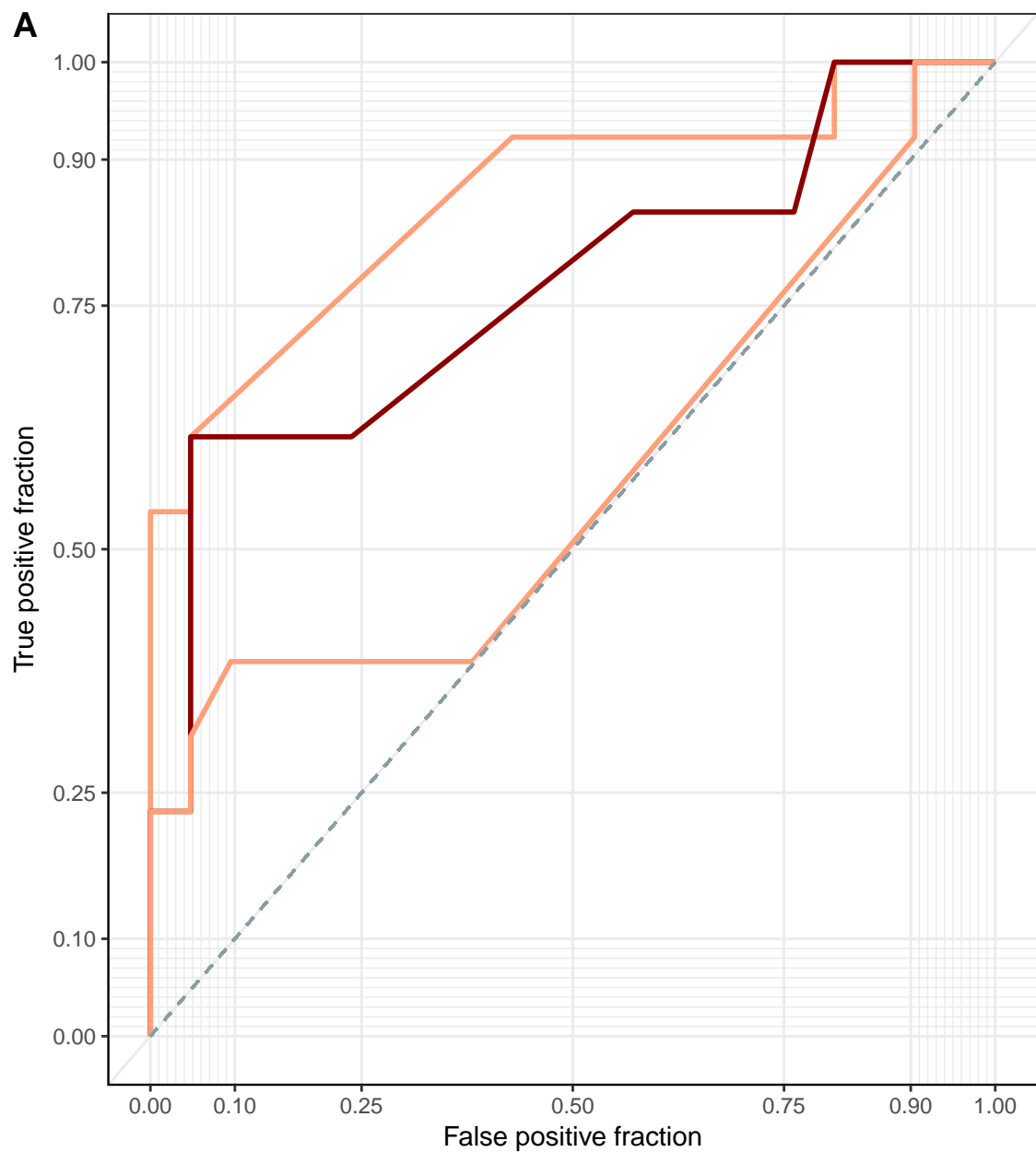


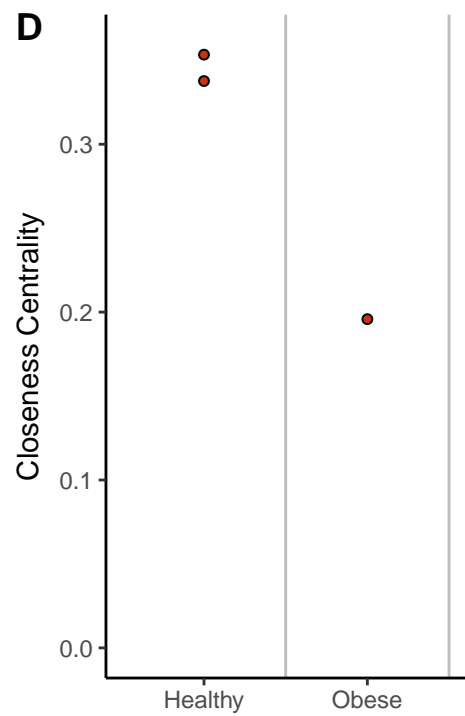
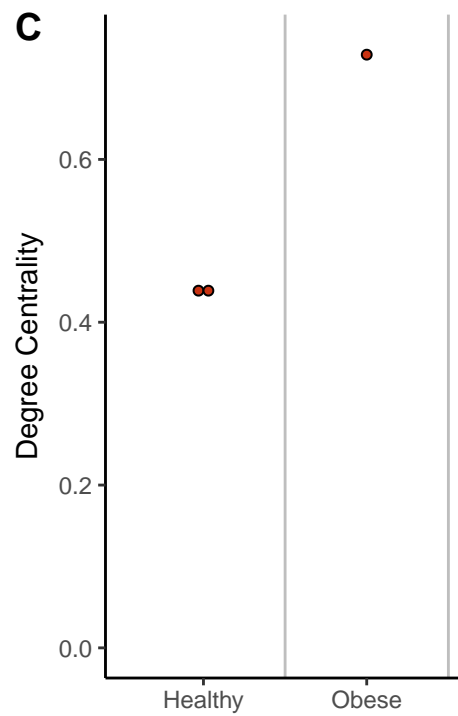
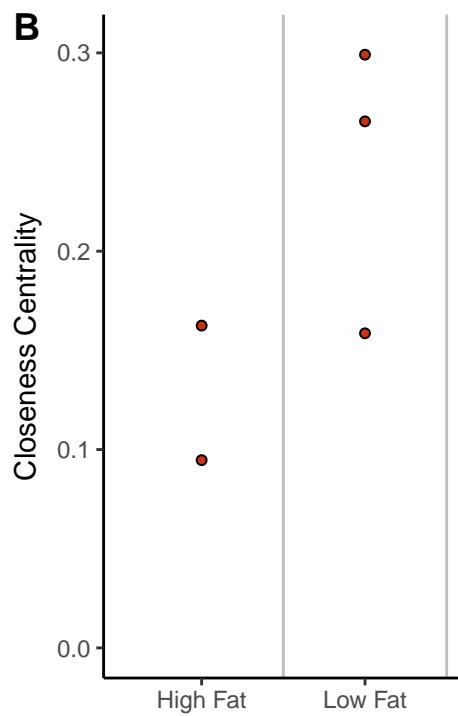
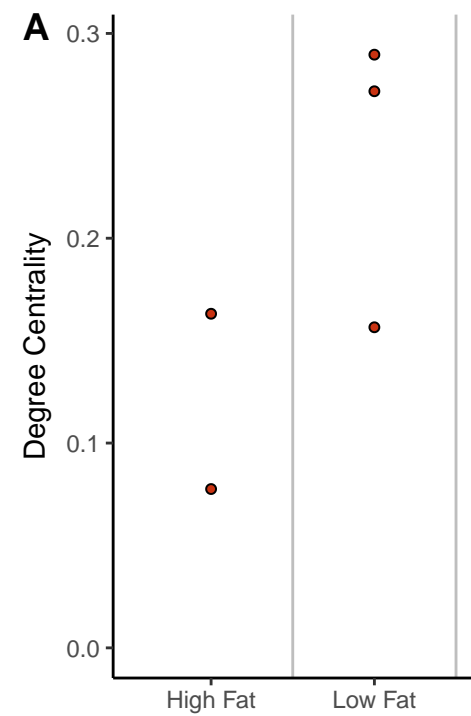
- 713 requires type IV pili for infection and shows broad bactericidal and biofilm removal activities. *Applied and*  
714 *Environmental Microbiology*. 2012;78: 6380–6385.
- 715 55. Matsuzaki S, Tanaka S, Koga T, Kawata T. A Broad-Host-Range Vibriophage, KVP40, Isolated from Sea  
716 Water. *Microbiology and Immunology*. 1992;36: 93–97.
- 717 56. Orchard S, Ammari M, Aranda B, Breuza L, Briganti L, Broackes-Carter F, et al. The MIntAct  
718 project–IntAct as a common curation platform for 11 molecular interaction databases. *Nucleic Acids*  
719 *Research*. 2014;42: D358–63.
- 720 57. Bashan A, Gibson TE, Friedman J, Carey VJ, Weiss ST, Hohmann EL, et al. Universality of human  
721 microbial dynamics. *Nature*. 2016;534: 259–262.
- 722 58. Kleiner M, Hooper LV, Duerkop BA. Evaluation of methods to purify virus-like particles for metagenomic  
723 sequencing of intestinal viromes. *BMC Genomics*. 2015;16: 7.
- 724 59. Meisel JS, Hannigan GD, Tyldsley AS, SanMiguel AJ, Hodkinson BP, Zheng Q, et al. Skin microbiome  
725 surveys are strongly influenced by experimental design. *Journal of Investigative Dermatology*. 2016;
- 726 60. Malki K, Kula A, Bruder K, Sible E, Hatzopoulos T, Steidel S, et al. Bacteriophages isolated from Lake  
727 Michigan demonstrate broad host-range across several bacterial phyla. *Virology Journal*. 2015;12: 164.
- 728 61. Turnbaugh PJ, Ridaura VK, Faith JJ, Rey FE, Knight R, Gordon JI. The effect of diet on the human  
729 gut microbiome: a metagenomic analysis in humanized gnotobiotic mice. *Science Translational Medicine*.  
730 2009;1: 6ra14–6ra14.
- 731 62. David LA, Maurice CF, Carmody RN, Gootenberg DB, Button JE, Wolfe BE, et al. Diet rapidly and  
732 reproducibly alters the human gut microbiome. *Nature*. 2014;505: 559–563.
- 733 63. Jeong H, Tombor B, Albert R, Oltvai ZN, Barabási AL. The large-scale organization of metabolic networks.  
734 *Nature*. 2000;407: 651–654.
- 735 64. Grice EA, Kong HH, Conlan S, Deming CB, Davis J, Young AC, et al. Topographical and Temporal

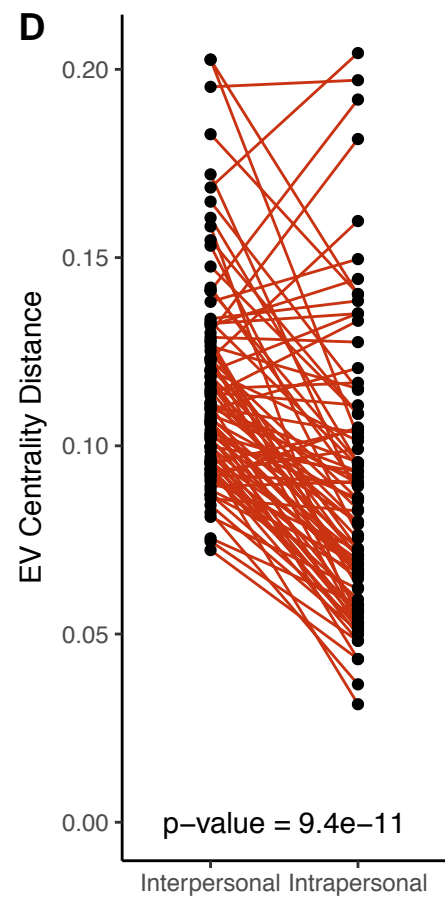
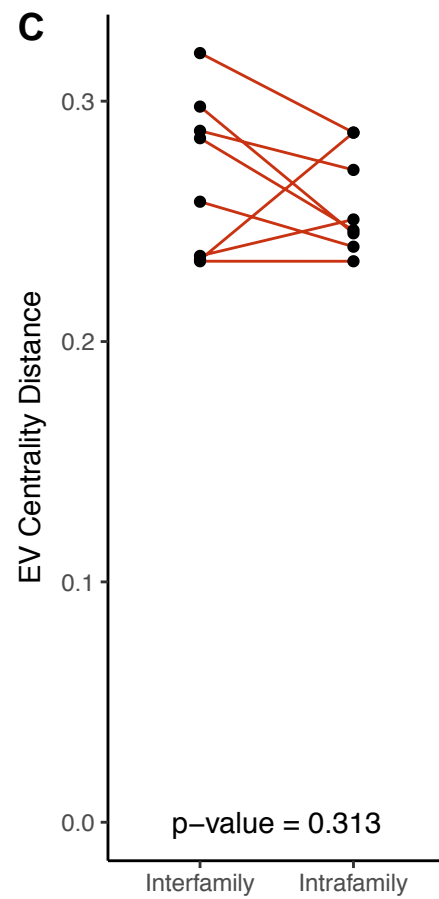
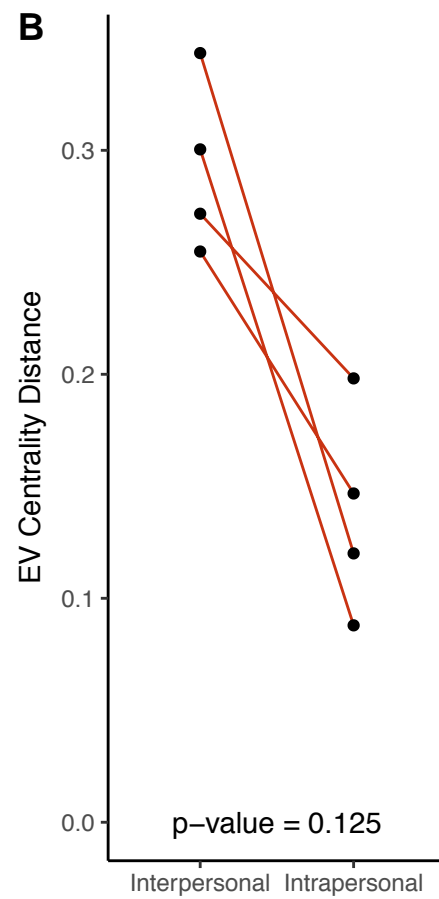
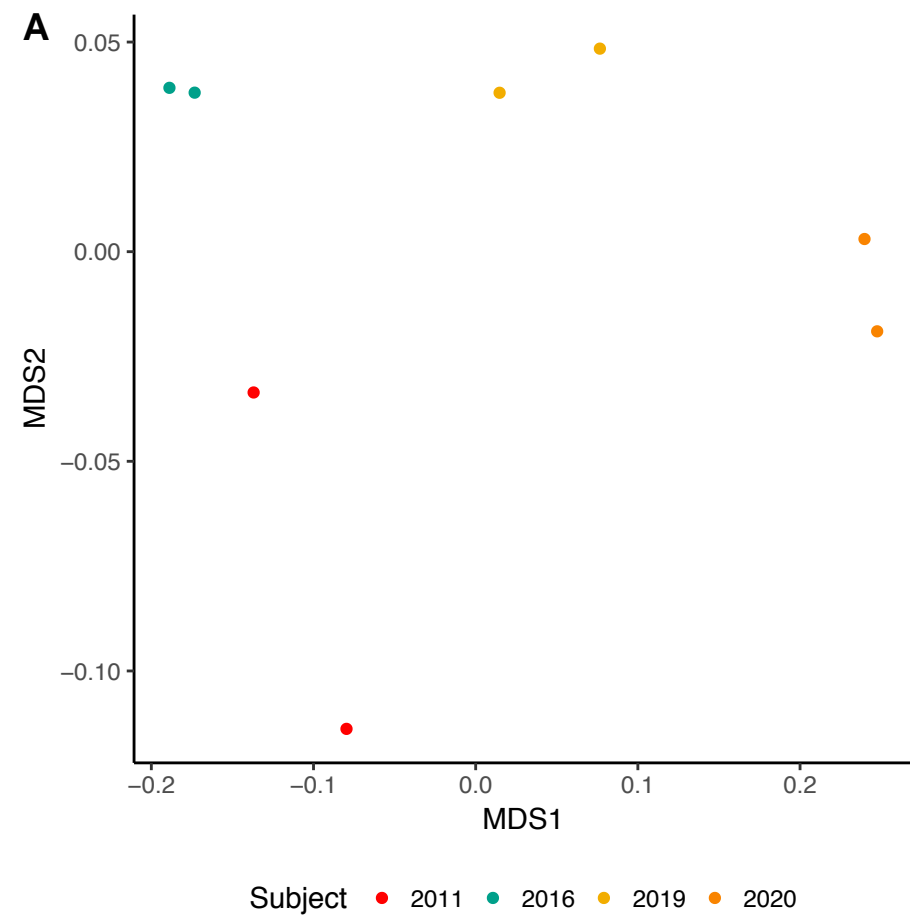
- 736 Diversity of the Human Skin Microbiome. *Science*. 2009;324: 1190–1192.
- 737 65. Minot S, Bryson A, Chehoud C, Wu GD, Lewis JD, Bushman FD. Rapid evolution of the human gut virome.  
738 *Proceedings of the National Academy of Sciences of the United States of America*. 2013;110: 12450–12455.
- 739 66. Schloss PD, Handelsman J. Introducing DOTUR, a computer program for defining operational taxonomic  
740 units and estimating species richness. *Applied and Environmental Microbiology*. 2005;71: 1501–1506.
- 741 67. Yilmaz S, Allgaier M, Hugenholtz P. Multiple displacement amplification compromises quantitative  
742 analysis of metagenomes. *Nature Methods*. 2010;7: 943–944.
- 743 68. Kim KH, Chang HW, Nam YD, Roh SW. Amplification of uncultured single-stranded DNA viruses from  
744 rice paddy soil. *Applied and ....* 2008;
- 745 69. Kim K-H, Bae J-W. Amplification methods bias metagenomic libraries of uncultured single-stranded and  
746 double-stranded DNA viruses. *Applied and Environmental Microbiology*. 2011;77: 7663–7668.
- 747 70. Minot S, Wu GD, Lewis JD, Bushman FD. Conservation of gene cassettes among diverse viruses of the  
748 human gut. *PLOS ONE*. 2012;7: e42342.
- 749 71. Deng L, Ignacio-Espinoza JC, Gregory AC, Poulos BT, Weitz JS, Hugenholtz P, et al. Viral tagging  
750 reveals discrete populations in *Synechococcus* viral genome sequence space. *Nature*. 2014;513: 242–245.
- 751 72. Brum JR, Ignacio-Espinoza JC, Roux S, Doucier G, Acinas SG, Alberti A, et al. Ocean plankton. Patterns  
752 and ecological drivers of ocean viral communities. *Science*. 2015;348: 1261498–1261498.
- 753 73. Polz MF, Hunt DE, Preheim SP, Weinreich DM. Patterns and mechanisms of genetic and phenotypic  
754 differentiation in marine microbes. *Philosophical Transactions of the Royal Society B: Biological Sciences*.  
755 2006;361: 2009–2021.
- 756 74. Gregory AC, Solonenko SA, Ignacio-Espinoza JC, LaButti K, Copeland A, Sudek S, et al. Genomic  
757 differentiation among wild cyanophages despite widespread horizontal gene transfer. *BMC Genomics*.  
758 2016;17: 930.
- 759 75. Paez-Espino D, Eloie-Fadrosch EA, Pavlopoulos GA, Thomas AD, Huntemann M, Mikhailova N, et al.

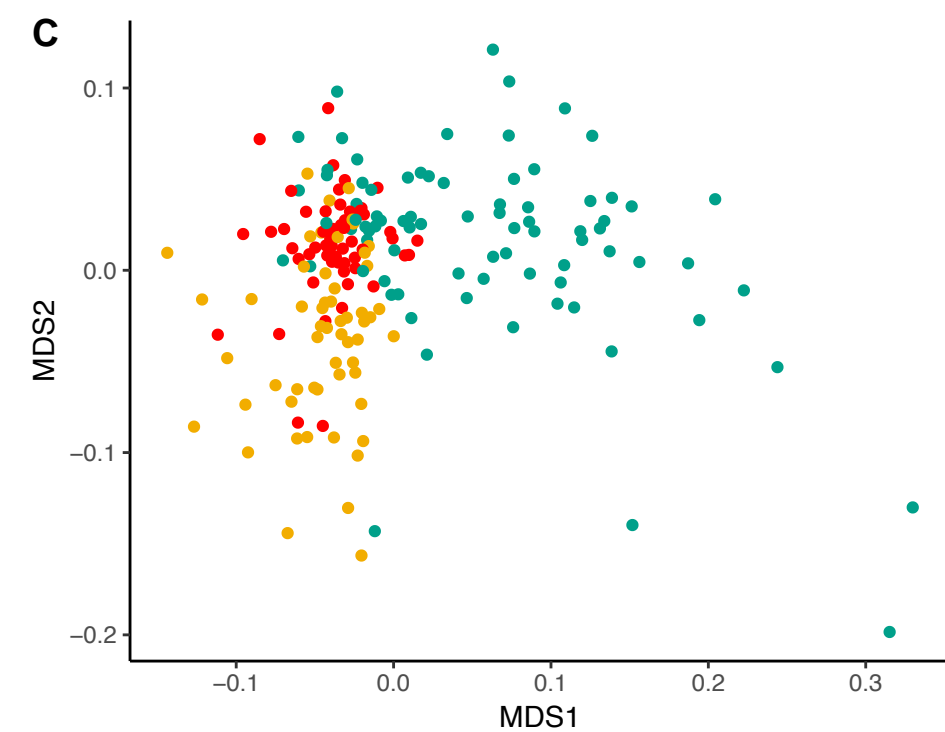
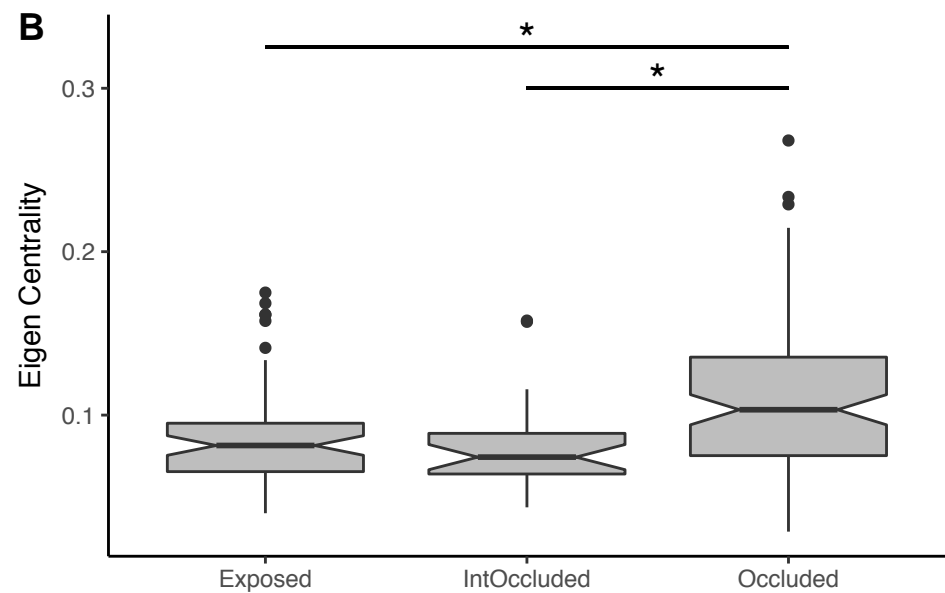
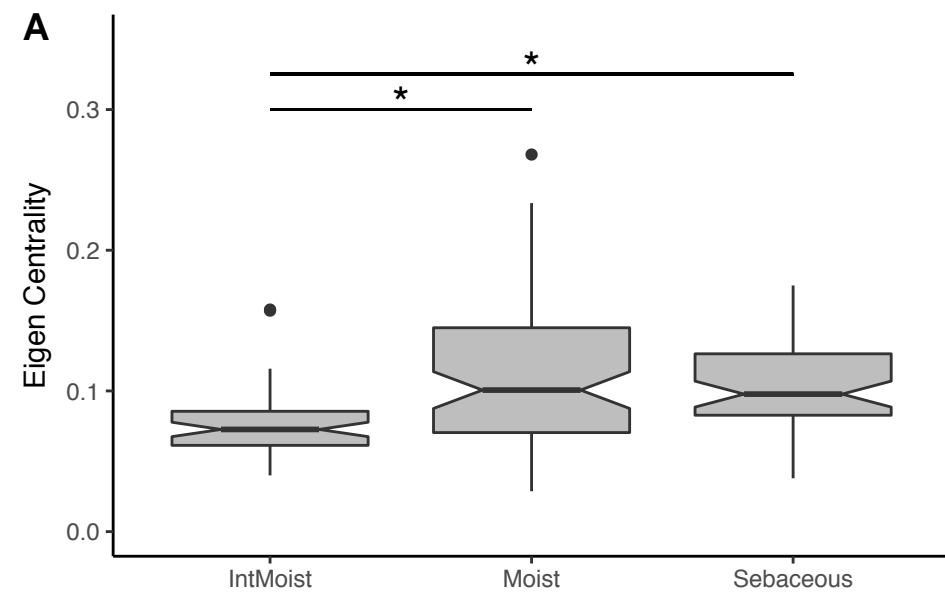
- 760 Uncovering Earth's virome. *Nature*. 2016;
- 761 76. Grice EA, Segre JA. The skin microbiome. *Nature Reviews Microbiology*. 2011;9: 244–253.
- 762 77. Round JL, Mazmanian SK. The gut microbiota shapes intestinal immune responses during health and  
763 disease. *Nature reviews Immunology*. 2009;9: 313–323.
- 764 78. Hannon GJ. FASTX-Toolkit. 2010; GNU Affero General Public License.
- 765 79. Li D, Luo R, Liu C-M, Leung C-M, Ting H-F, Sadakane K, et al. MEGAHIT v1.0: A fast and scalable  
766 metagenome assembler driven by advanced methodologies and community practices. *METHODS*.  
767 2016;102: 3–11.
- 768 80. Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. *Nature Methods*. 2012;9:  
769 357–359.
- 770 81. Alneberg J, Bjarnason BS, Ariani A, Bruijn I de, Schirmer M, Quick J, Ijaz UZ, et al. Binning metagenomic  
771 contigs by coverage and composition. *Nature Methods*. 2014; 1–7.
- 772 82. Hyatt D, LoCascio PF, Hauser LJ, Uberbacher EC. Gene and translation initiation site prediction in  
773 metagenomic sequences. *Bioinformatics*. 2012;28: 2223–2230.
- 774 83. Kuhn M. caret: Classification and Regression Training. CRAN. 2016;
- 775 84. Edgar RC. PILER-CR: fast and accurate identification of CRISPR repeats. *BMC Bioinformatics*. 2007;8:  
776 18.
- 777 85. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, et al. BLAST+: architecture and  
778 applications. *BMC Bioinformatics*. 2009;10: 1.
- 779 86. Buchfink B, Xie C, Huson DH. Fast and sensitive protein alignment using DIAMOND. *Nature Methods*.  
780 2015;12: 59–60.
- 781 87. Neo Technology, Inc. Neo4j. 2017;
- 782 88. Csardi G, Nepusz T. The igraph software package for complex network research. *InterJournal*.

783 2006;Complex Systems: 1695.

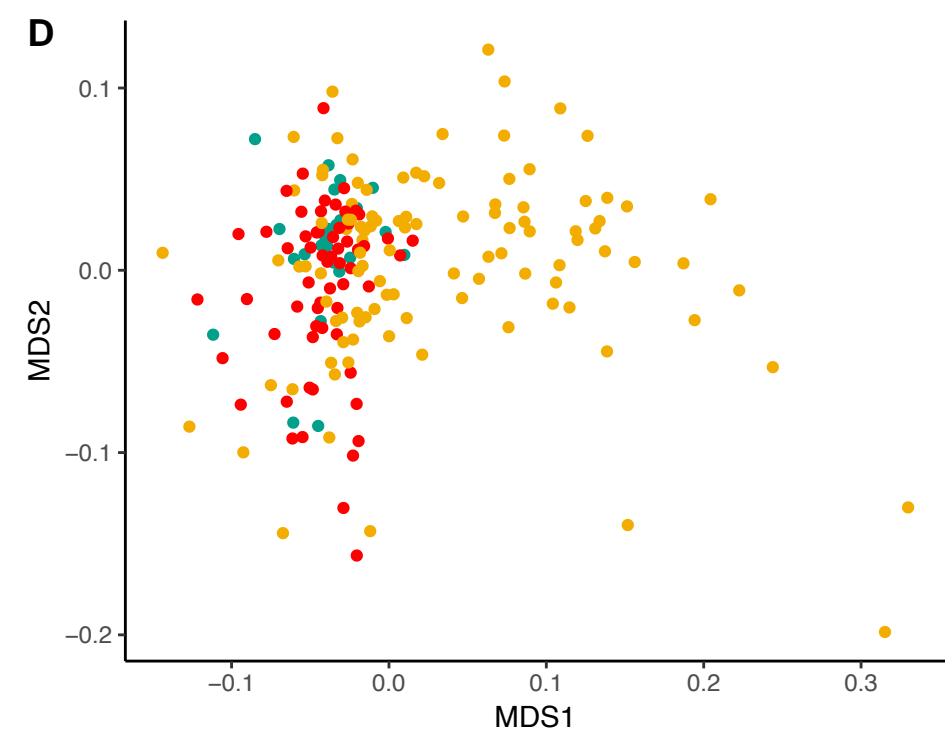








Environment ● IntMoist ● Moist ● Sebaceous



Environment ● Exposed ● IntOccluded ● Occluded

