

1 **Improving gut virome comparisons using predicted phage host information**

2

3 Michael Shamash ^{‡ a}, Anshul Sinha ^{‡ a}, and Corinne F. Maurice ^{a, b *}

4

5 [‡] These authors contributed equally to this work

6 ^a Department of Microbiology & Immunology, McGill University, Montreal, QC, Canada

7 ^b McGill Centre for Microbiome Research, Montreal, QC, Canada

8 * Corresponding author: corinne.maurice@mcgill.ca

9

10 **Competing interests statement**

11 The authors have no conflicts of interest to disclose.

12

13 **Abstract**

14 The human gut virome is predominantly made up of bacteriophages (phages), viruses that infect
15 bacteria. Metagenomic studies have revealed that phages in the gut are highly individual specific
16 and dynamic. These features make it challenging to perform meaningful cross-study comparisons.
17 While several taxonomy frameworks exist to group phages and improve these comparisons, these
18 strategies provide little insight into the potential effects phages have on their bacterial hosts. Here,
19 we propose the use of predicted phage host families (PHFs) as a functionally relevant, higher rank
20 unit of phage taxonomy to improve these cross-study analyses. We first show that bioinformatic
21 predictions of phage hosts are accurate at the host family level by measuring their concordance to
22 Hi-C sequencing-based predictions in human and mouse fecal samples. Next, using phage host
23 family predictions, we determined that PHFs reduce intra- and interindividual ecological distances
24 compared to viral contigs in a previously published cohort of 10 healthy individuals, while
25 simultaneously improving longitudinal virome stability. Lastly, by reanalyzing a previously
26 published metagenomics dataset with > 1,000 samples, we determined that PHFs are prevalent
27 across individuals and can aid in the detection of inflammatory bowel disease-specific virome
28 signatures. Overall, our analyses support the use of predicted phage hosts in reducing between-
29 sample distances and providing a biologically relevant framework for making between-sample
30 virome comparisons.

31

32 **Introduction**

33 The human gut virome, the collection of viruses in the human gastrointestinal tract, is dominated
34 by bacteria-infecting bacteriophages (phages). This community is highly diverse and individual-
35 specific at the nucleotide level (1–5). This vast diversity makes it challenging to perform cross-
36 individual or cross-cohort comparisons, as it is rare for all individuals in a cohort group to share a
37 single viral OTU (vOTU). Most recently, a longitudinal analysis of 59 individuals further
38 demonstrated that the individuality of the gut virome confounded disease signal detection in the
39 context of inflammatory bowel diseases (IBDs), in part due to intraindividual fluctuations of
40 viruses over time (6).

41 Viral clusters, such as those generated by vConTACT 2 (7), have been proposed as a potential
42 solution to this, by grouping together viruses based on their shared protein content. While this
43 approach is useful for complete genome sequences, it may not always be reliable in the context of
44 virome datasets. Indeed, current short-read virome studies often generate several contigs for a
45 single viral genome, raising the risk that each contig from a given virus would be placed into a
46 different viral cluster, confounding ecological conclusions (7).

47 An important limitation of previous virome analyses was the inability to confidently link
48 uncultured phages with their hosts. Recent experimental and bioinformatic advances aim to
49 address this issue. For instance, proximity ligation sequencing is being used to assign phages to
50 their hosts *in situ*. With this approach, phage DNA inside of host cells at the time of sampling is
51 covalently crosslinked to the bacterial host DNA, leading to generation of chimeric reads during
52 the sequencing process (8, 9). On the computational side, tools such as iPHoP enable the high-
53 throughput prediction of hosts using phage sequence data alone (10). Using a combination of
54 existing tools and machine learning models, iPHoP can consistently predict hosts down to the

55 genus level. These two approaches, proximity ligation and iPHoP, have yet to be formally
56 compared with each other for assigning hosts to gut virome-derived sequences.
57 Here, we propose using predicted phage host range to allow for ecologically relevant comparisons
58 of viromes across individuals, regardless of nucleotide-level diversity. These comparisons could
59 provide broad insight on ecosystem function, as phages have the ability to alter bacterial
60 abundances and metabolism (11, 12). We introduce the term Phage Host Family (PHF) as a term
61 to describe the predicted bacterial host of a phage sequence, at the family level. This family-level
62 cut-off was determined based on comparisons of predicted phage host range from iPHoP with
63 experimental assignments via proximity ligation sequencing of human and mouse fecal samples,
64 where high concordance was seen down to the family, but not genus level. Using this metric, we
65 then re-evaluate two previously published large datasets. First, we apply PHF analysis to viromes
66 from a cohort of 10 healthy individuals (1), sampled longitudinally for approximately 1 year, and
67 conclude that incorporating PHFs reduced interindividual variation, while also increasing within-
68 individual virome stability over time. Second, we analyze the phageome of a large cohort of
69 individuals with IBDs (13), where we determine that aggregating vOTUs using PHFs allows for
70 the detection of greater disease-specific differences in the virome, in addition to reducing
71 interindividual variability. We propose that the use of PHFs as an ecologically informed unit of
72 phage taxonomy is useful in allowing for cross-sample comparisons in gut virome studies.
73

74 **Methods**

75 **Preparing fecal samples for proximity ligation sequencing**

76 Human fecal samples were collected with the approval of protocol A04-M27-15B from the McGill
77 University Institutional Review Board. Participants provided informed written consent for the
78 utilization of their samples. Fresh fecal samples were collected, aliquoted in an anaerobic chamber,
79 and kept at -70 °C until processing.

80 Adult female germ-free C57BL/6 mice were maintained in Tecniplast IsoCages at McGill
81 University. Mice had unlimited access to irradiated diet (Research Diets, New Brunswick, NJ) and
82 autoclaved water. Germ-free mice were humanized by oral gavage of 200 uL of resuspended
83 human donor feces. Mouse fecal samples were collected with the approval of McGill University
84 animal use protocol MCGL-7999.

85 A total of 10 mouse fecal samples (from 6 mice), and 2 human fecal samples, were collected for
86 proximity ligation and bulk metagenome sequencing. Fresh fecal samples were collected and
87 stored at -70 °C until processing. All mice had unlimited access to standard chow and water. Fecal
88 samples were resuspended in 1 mL PBS (0.02 um filter-sterilized). After an initial centrifugation
89 at 1,000xg for 1 minute to pellet large debris, the bacterial cell-containing supernatant was
90 centrifuged again at 10,000xg for 10 minutes. Pelleted bacterial cells were resuspended in 1 mL
91 of a PBS-formaldehyde solution (1% formaldehyde) and incubated at room temperature for 20
92 minutes to cross-link DNA. Glycine was added in excess to quench unused formaldehyde and
93 incubated for 15 minutes at room temperature. The fixed bacterial cells were pelleted (10,000xg
94 for 10 minutes) and washed twice with PBS. The final resuspended bacterial pellet was transferred
95 to a BeadBug tube with 0.1mm silica glass beads (Benchmark Scientific, Sayreville, United States),
96 and vortexed at maximum speed for 5 minutes. The sample was transferred to a DNA LoBind tube

97 and sent to Phase Genomics (Seattle, United States) for library preparation and sequencing using
98 their ProxiMeta kit and analysed using the corresponding bioinformatic pipeline (14). In addition,
99 a bulk metagenome was sequenced for each sample: immediately after resuspending the original
100 fecal sample in PBS, 250 uL of sample was used for DNA extraction using the QIAGEN
101 PowerFecal Pro DNA kit (QIAGEN, Hilden, Germany) following the manufacturer's instructions.
102 All libraries (proximity ligation and bulk metagenome) were sequenced using the Illumina
103 NovaSeq platform with 2x 150bp reads.

104

105 **Comparing host predictions between iPHoP and proximity-ligation sequencing**

106 Viral contigs identified by the ProxiMeta pipeline were used as input for iPHoP (v. 1.3.3) (10) to
107 computationally predict hosts, the output was imported into R (v. 4.2.2) for analysis. The
108 proximity-ligation linkage data was imported into R and filtered to keep only viral contigs who
109 also had hosts predicted with iPHoP. These two datasets were then compared for concordance at
110 the following taxonomic ranks: phylum, class, order, family, and genus. Concordance was
111 calculated at each taxonomic rank using three distinct approaches to account for phages which
112 have multiple assigned/predicted hosts: (1) consider the pairing concordant if the most confident
113 iPHoP prediction matches the top Hi-C hit (most stringent); (2) consider the pairing concordant if
114 the most confident iPHoP prediction matches any of the Hi-C hits; and (3) consider the pairing
115 concordant if any of the iPHoP predictions match any of the Hi-C hits (least stringent). The percent
116 concordance was calculated as the number of viral contigs with concordant hosts, divided by the
117 total number of viral contigs, multiplied by 100%.

118

119

120 **Re-analysis of Shkoporov *et al.* dataset**

121 The published phyloseq object (1) was downloaded and imported into R (v. 4.2.2) using phyloseq
122 (v. 1.42) (15). Published virome contigs were also downloaded and filtered to keep those > 1 kb
123 using seqkit (v. 2.5.1), resulting in 57,721 contigs. iPHoP (v. 1.3.3) (10) was used to predict the
124 bacterial hosts of these contigs, and the output was imported into R for analysis. The GTDB tree
125 used by iPHoP (bac120_r202.tree) was also imported into R using phyloseq's read_tree command
126 and combined with the downloaded phyloseq object. Viral contig relative abundance was
127 calculated and added to the phyloseq object, replacing the existing otu_table object. Predicted host
128 information was added to the phyloseq object as a tax_table object. In cases where a viral contig
129 had more than 1 predicted host by iPHoP, the most confident host prediction was selected. Taxa
130 bar plots were generated using microshades (v. 1.10) (16). Samples with less than 30% of the
131 community consisting of contigs with unknown hosts were retained for subsequent analysis.
132 Phyloseq's tax_glom function was used to agglomerate viral contigs that have the same predicted
133 host at the family level. Vegan (v. 2.6-6.1) (17) was used to calculate distances between samples
134 using the bray, wunifrac, and unifrac metrics. Distances were evaluated by comparison type, either
135 inter- or intraindividual sample comparisons, and the Friedman test with post-hoc Wilcoxon
136 signed-rank test (using Bonferroni correction for multiple comparisons) were used test for
137 significance.

138 We define virome stability as the similarity between two sequential samples and it is calculated as
139 follows: $\text{stability} = (1 - \text{distance from previous sample})$. We calculated virome stability using
140 distances between consecutively collected samples from the same individual using Bray-Curtis
141 distances at the contig and PHF levels and tested for significance using the Wilcoxon signed-rank
142 test.

143 **Re-analysis of the HMP2 dataset**

144 In the human IBD cohort, originally analyzed by Lloyd-Price *et al.*, bulk metagenome reads were
145 obtained from 1,595 samples belonging to 130 individuals (27 non-IBD, 65 CD, and 38 UC)
146 sampled longitudinally over one year (13). Data was downloaded from: <https://ibdmdb.org/results>.
147 Paired-end sequencing reads (101 bp) were generated using Illumina HiSeq 2000 or 2500. Raw
148 reads were trimmed based on sequence quality using Trimmomatic (v 0.33) (18). Quality-
149 controlled sequences that aligned to human and mouse genomes were removed by Bowtie2 (v.
150 2.2) (19). These steps were performed using the kneaddata workflow (20). Quality-controlled reads
151 were then grouped by individual and co-assembled into 3,249,501 contigs > 1kb using MEGAHIT
152 (v. 1.2.9) (21). The contigs within each co-assembly were classified as phages by VIBRANT (v.
153 1.2.1) (22). In total, there were 81,422 predicted phages across all co-assemblies. These contigs
154 were then filtered for completeness using CheckV (v. 1.0.3) (23), keeping only the 6,741 contigs
155 that were over 50% complete. A 50% completeness cutoff was used to balance the trade-offs of
156 overestimating viral richness due to fragmentation during assembly and maintaining viral richness.
157 These remaining contigs were then dereplicated using blastn, keeping contigs with an average
158 nucleotide identity of 95% over 80% alignment fraction relative to the shorter sequence (24).
159 Similar to the analyses of the Shkoporov *et al.* dataset, iPHoP (v. 1.3.3) (10) was used to predict
160 the bacterial host of each phage contig, keeping the most confident host prediction if there were
161 multiple predictions. Quality-controlled reads were mapped to the phage contigs library using
162 Bowtie2 (19). Contigs were considered present in a given sample using mapping thresholds
163 defined by Stockdale *et al.* (6), where a contig was present if Bowtie2 mapped reads covered 50%
164 of contigs <5 kb, 30% of contigs \geq 5 kb and <20 kb, or 10% of contigs \geq 20 kb. After calculating
165 Good's coverage and generating rarefaction curves for each sample, 502 samples were removed

166 which had below 1,500 length-normalized read counts (25). PCoAs and Bray-Curtis distances on
167 the remaining 1,093 samples were generated using MicroViz (v. 0.12.1) (26), which uses Vegan
168 as a wrapper. DESeq2 (v. 1.44) (27) was used to calculate differentially abundant PHFs based on
169 dysbiosis status, using the simple formula: design = ~ Participant.ID + dysbiosis_binary. PHFs
170 with an adjusted p value ≤ 0.05 and with a \log_2 fold-change ≥ 1 or with a \log_2 fold-change ≤ -1
171 were considered differentially abundant. For differential abundance analyses, only individuals
172 which had both a dysbiotic and non-dysbiotic sample were included so that a paired analysis could
173 be conducted. Only the 18 PHFs that were more than 50% prevalent across individuals were
174 considered for analysed PHFs. This arbitrary threshold was used to consider only the features that
175 were widely distributed and abundant across samples. Auxiliary metabolic genes (AMGs) were
176 predicted from viral contigs using VIBRANT. Using KEGG annotations, VIBRANT categorizes
177 these AMGs into metabolic categories (22). In some cases, a single AMG belonged to multiple
178 metabolic categories.

179

180 **Code and data availability**

181 Code used for data analysis is available at: https://github.com/mshamash/PHF_manuscript. Whole
182 genome and Hi-C sequencing reads are available on the NCBI SRA using accession number
183 PRJNA1145458.

184

185 **Results**

186 *Computational prediction of phage hosts is concordant with proximity ligation sequencing* 187 *assignments to the family level*

188 We conducted proximity ligation (Hi-C) sequencing on 10 fecal samples from human
189 microbiota-associated mice, and 2 fecal samples from healthy human donors. After Hi-C host
190 assignment, we identified 1,577 phage-host pairings consisting of 1,547 phages targeting 77
191 unique hosts at the genus level, with some phages being linked to more than one host. Using iPHoP,
192 we then predicted hosts for the 1,547 phages with Hi-C-assigned hosts, yielding 1,587 phage-host
193 pairings, comprising 1,243 phages targeting 108 unique hosts at the genus level. These 1,243
194 phages, which had hosts assigned by both Hi-C and iPHoP, were used for subsequent comparisons
195 between approaches.

196 Concordance between the two approaches was calculated at each taxonomic rank from
197 phylum to genus using three distinct approaches to account for some phages having multiple
198 assigned/predicted hosts: (1) pairing is concordant if the most confident iPHoP prediction matches
199 the top Hi-C hit (most stringent); (2) pairing is concordant if the most confident iPHoP prediction
200 matches any of the Hi-C hits; and (3) pairing is concordant if any of the iPHoP predictions match
201 any of the Hi-C hits (least stringent). Regardless of the comparison approach, there was
202 consistently over 98% concordance at the phylum level, over 97% concordance at the class level,
203 over 96% concordance at the order level, and over 92% concordance at the family level (**Figure**
204 **1**). Genus-level concordance was lower, with approximately 67% concordance using comparison
205 metrics (1) and (2), and 73% concordance using metric (3) (**Figure 1**).

206 Based on these findings, we decided to use iPHoP predictions at the family-level for our
207 subsequent analyses as this was the lowest taxonomic rank which still had high concordance

208 between the tool's predictions and Hi-C experimental assignments. These family-level host
209 predictions will now be referred to as PHF.

210

211 *Using predicted hosts as a functional measure of virome diversity reduces interindividual*
212 *variation and increases intraindividual stability*

213 We next wanted to use PHFs to evaluate functional virome diversity within and across
214 individuals. Using a previously published dataset consisting of 140 total samples from 10 healthy
215 individuals (1), we predicted hosts for the provided assembled phage contigs (n=57,721 contigs).
216 In total, iPHoP yielded 197,994 host predictions for 49,852 (86.3%) of the viral contigs (an average
217 of 3.97 bacterial hosts predicted per viral contig). The remaining 7,869 (13.7%) contigs had no
218 host predicted. The most confident iPHoP prediction for each contig was retained and used in
219 downstream analysis. Overall, there was large variation in the proportion of the virome made up
220 of contigs with known hosts (**Figure 2A**). To ensure that the subsequent comparisons between
221 samples are fair (i.e., by comparing samples with similar proportions of the community represented
222 by contigs having assigned hosts), we filtered the dataset to retain only samples which had less
223 than 30% of the community consisting of contigs with unknown hosts, resulting in a new dataset
224 composed of 63 samples from 10 individuals. Keeping all samples, including those with a high
225 proportion of contigs with unknown hosts, would introduce bias into our analyses by over-
226 representing incomplete or ambiguous community structures, potentially leading to inaccurate
227 conclusions about the relationships between samples. Using CheckV to filter for contigs >50%
228 complete did not significantly change the proportion of the community consisting of contigs with
229 unknown hosts (data not shown).

230 Phages with the same family-level host predictions were agglomerated into PHF groups
231 and the resulting abundance matrix was used for subsequent analyses. Pairwise distances were
232 calculated between all samples using traditional contig-level Bray-Curtis, PHF-level Bray-Curtis,
233 and PHF-level weighted UniFrac metrics, as described in Methods. Intraindividual sample
234 distances were consistently lower than interindividual sample distances (**Figure 2B**). Regardless
235 of inter- or intraindividual sample comparison, PHF-level weighted UniFrac distances were
236 significantly lower than PHF-level Bray-Curtis distances, which were themselves significantly
237 lower than contig-level Bray-Curtis distances (**Figure 2B**). Finally, we evaluated the effects of
238 using PHF-level distances on longitudinal intraindividual virome stability, calculated as the
239 pairwise distance between all pairs of consecutive samples from each individual. Virome stability
240 was consistently higher when incorporating PHF-level distances using the Bray-Curtis metric
241 (**Figure 2C**).

242

243 *PHFs are prevalent and can provide biological insight into the IBD virome*

244 We next wanted to characterize PHFs in a larger dataset and determine whether
245 agglomerating at the PHF-level could improve detection of disease-specific signatures. To do so,
246 we re-analyzed the human microbiome project 2 (HMP2) dataset containing longitudinal bulk
247 metagenome samples from IBD and non-IBD controls. After removing samples which contained
248 low read counts (see methods), 1,093 samples from 57 CD, 31 UC, 27 non-IBD controls remained
249 for further analyses (68.5% of total samples) (13). From these samples, we co-assembled contigs,
250 dereplicated these contigs, predicted phages using VIBRANT, and filtered for phage completeness
251 > 50% using CheckV. Using this approach, we obtained a total of 3,862 distinct virus operational
252 taxonomic units (vOTUs) across the samples within the dataset. Of these 3,862 vOTUs, 87.1%

253 (3,365/3,862) had an iPHoP predicted host. In total, these vOTUs belonged to 75 distinct PHFs.
254 Interestingly, the amount of vOTUs comprising each PHF varied greatly, with some PHFs
255 comprised of hundreds of distinct vOTUs, whereas some rare PHFs were only comprised of a
256 single vOTU (**Supplementary Figure 1**).

257 To further characterize PHFs, and to link their host associations with metabolic
258 functionality, we searched for viral-encoded AMGs. These genes, which are expressed throughout
259 the process of viral infection, are thought to provide phages with increased fitness via modulation
260 of host metabolism (22, 28). In total, 45/75 PHFs carried at least one AMG and 12/75 PHFs carried
261 at least 10 AMGs. In general, PHFs were enriched in amino acid metabolism, energy metabolism,
262 and cofactor and vitamin metabolism genes (**Supplementary Figure 2A**), in line with previous
263 surveys of AMGs in human microbiomes (22). Notably, compared to other PHFs,
264 *Bifidobacteriaceae*-infecting phages were enriched in carbohydrate metabolism genes
265 (**Supplementary Figure 2B**). *Enterobacteriaceae*-infecting phages on the other hand were
266 enriched in protein folding, sorting, and degradation genes (**Supplementary Figure 2C**), and in
267 particular *cysO*, which encodes a sulfur-carrier protein important in cysteine biosynthesis and
268 resistance to oxidative stress (29). Fourteen distinct *Enterobacteriaceae*-infecting phage vOTUs
269 carried *cysO* (**Supplementary Figure 2D**). Only 3/75 other PHFs (*Pasteurellaceae*,
270 *Pseudomonadaceae*, *Burkholderiaceae*) carried *cysO* on 7 distinct contigs (**Supplementary**
271 **Figure 2D**). Interestingly, all of these PHFs infect bacteria from the phylum Proteobacteria,
272 potentially reflecting host-specific adaptation through the carriage of this AMG.

273 Consistent with the high levels of interindividuality at the vOTU level observed in the
274 Shkoporov *et al.* dataset, we found that only 236/3,862 (6.11%) vOTUs were found in more than
275 50% of individuals in the HMP2 dataset (**Figure 3A**). In contrast, a higher proportion of PHFs

276 (18/75; 24%) were found in more than 50% of individuals (**Figure 3B**). Importantly, these
277 prevalent features made up a significantly higher mean relative abundance in samples at the PHF
278 level compared to the vOTU level (**Figure 3C**). Thus, prevalent PHFs represent a larger fraction
279 of the total community in comparison to prevalent vOTUs. In line with these observations,
280 intraindividual and interindividual Bray-Curtis distance between samples was significantly lower
281 at the PHF level in comparison to the vOTU level (**Figure 3D**).

282 Given that PHFs reduced ecological distance between samples, we hypothesized that this
283 would also allow for more biologically relevant comparisons between individuals, and ultimately
284 a greater ability to detect disease-specific signatures in the human virome. We first generated
285 PCoA plots using Bray-Curtis distance and found that the first two principal components explained
286 more cumulative variance when agglomerating the virome at the PHF level in comparison to the
287 vOTU level (**Figure 4A**; 39.3% vs. 11%). Importantly, the proportion of variance explained by
288 diagnosis (non-IBD, CD, UC) was higher using PHFs than using vOTUs (**Supplementary Figure**
289 **3**; $R^2 = 0.0261$ vs. $R^2 = 0.0185$). Lloyd Price *et al.* defined dysbiotic samples within this HMP2
290 dataset as those with high microbiota divergence from non-IBD controls (13). Using this
291 designation, we also found that dysbiosis status explained a higher proportion of variance using
292 PHFs when compared to vOTUs (**Figure 4A**; $R^2 = 0.0394$ vs. $R^2 = 0.0157$). We also performed
293 differential abundance analyses to determine whether certain PHFs were enriched or depleted
294 depending on dysbiosis status. Including only prevalent PHFs (found in > 50% of individuals), we
295 identified a single PHF enriched in dysbiotic samples (*Enterobacteriaceae*) and 4 significantly
296 depleted PHFs (*CAG-74*, *Ruminococcaceae*, *Acidaminococcaceae*, *Acutalibacteraceae*) (**Figure**
297 **4B**). These observations suggest that predicted phage hosts can be used to identify certain IBD-
298 specific virome signatures.

299 **Discussion**

300 In the past decade, the development of phage-specific bioinformatic tools, alongside large
301 cohort viral metagenomic studies, has revealed key characteristics of the human gut virome.
302 Notably, gut viromes exhibit high levels of interindividuality (1, 6) and temporal variation (6).
303 While we can now appreciate the sheer genomic phage diversity that our collective guts harbor, it
304 remains a challenge to understand how similar our viromes are over time and from one another.
305 Here, we demonstrate that the use of predicted phage host families (PHFs) can improve virome
306 comparisons between and within individuals, resulting in valuable functional information typically
307 lost with current approaches.

308 Using PHFs as a unit of taxonomy in two independent published datasets, we showed that
309 in comparison to vOTUs, intra- and inter-personal ecological distance is reduced, indicating that
310 despite phages differing between samples at the contig/vOTU level, their functionality remains
311 similar. These findings are reminiscent of the functional redundancy characteristic of gut bacterial
312 communities, whereby phylogenetic differences between individuals exist despite conserved
313 functional profiles (30, 31). The conserved functionality of both phage and bacterial communities
314 over time likely contribute to the stability and resilience of both subsets.

315 The advantages of working with reduced between-sample virome distance were evident as
316 we showed that the first two principal coordinate axes of PCoA plots explained more variance
317 when using PHFs as the unit of taxonomy. We also showed that the proportion of variance
318 explained by disease and dysbiosis status were greater when using PHFs. These findings are in
319 line with those from Clooney *et al.* (2), who analyzed human IBD viromes. They showed that
320 gene-sharing-based genus-level taxonomy, compared to contig-based analyses, better identified
321 disease-associated compositional changes and increased the variance explained by the first two

322 principal coordinate axes. These observations together highlight the importance of using a higher
323 taxonomic rank when making cross-individual comparisons of gut viromes.

324 A key additional benefit of using predicted hosts lies in the biologically relevant
325 information they provide. This contrasts with existing gene-sharing and phage morphology-based
326 taxonomy approaches, where taxonomic groups are not necessarily informative of how phages
327 interact with their bacterial hosts or the ecosystem at-large (32). Phages in several ecosystems,
328 including the gut, have been shown to be strong regulators of bacterial abundance, diversity, and
329 metabolism (11, 12, 33). Therefore, grouping phages by their predicted hosts provides context for
330 the effects that they may have on the bacterial community and beyond. We showed that in the
331 context of IBD, dysbiotic samples were enriched in *Enterobacteriaceae* PHFs, and depleted in
332 *CAG-74*, *Ruminococcaceae*, *Acidaminococcaceae* and *Acutalibacteraceae* PHFs. Our analyses
333 provide a framework to identify interactions relevant to disease, although follow-up studies are
334 needed to understand the importance of these phage-host interactions. For instance, phage
335 enrichment in tandem with host depletion could be relevant to several diseases (34–36). While the
336 increased abundance of *Enterobacteriaceae* PHFs we observed in dysbiotic samples is likely a
337 consequence of increased host abundance, it is interesting to note that these phages were enriched
338 in *cysO*, a gene involved in cysteine biosynthesis. Notably, *cysO* has been directly tied to defence
339 against oxidative stress (37). As *Enterobacteriaceae* are known to proliferate in the inflamed gut
340 in the face of oxidative stress (38, 39), our data imply that phage-encoded AMGs could be a source
341 of this resistance. More broadly, the observation that different PHFs carry distinct AMGs suggests
342 that grouping at the phage host family level provides an additional layer of functional insight
343 beyond phage-host relationships. However, an important consideration is that potential

344 inaccuracies in defining prophage borders (40) could lead to an overestimation of AMGs (41).
345 Thus, while these findings merit further investigation, they should be interpreted with caution.

346 While we do propose the use of PHFs for between-sample virome comparisons, it should
347 be noted that this method is reliant on the sensitivity of iPHoP (or any other phage-host matching
348 bioinformatic tool used). In our analyses, between 12.9% (Lloyd-Price *et al.* dataset) and 13.7%
349 (Shkoporov *et al.* dataset) vOTUs did not have an iPHoP-assigned host. This may become an even
350 larger issue if this approach is applied to non-human associated microbiomes where iPHoP
351 performs with less sensitivity (10). Regardless, it is reasonable to assume that the sensitivity of
352 bioinformatic phage-host prediction tools will improve alongside recent improvements in phage
353 genome reconstruction approaches such as contig extension (42) and viral binning (43).

354 To assess the accuracy of bioinformatic phage-host predictions, we measured the
355 concordance at different taxonomic ranks between iPHoP, a bioinformatic tool, and Hi-C
356 sequencing, which relies on physical linkage between phage and host. Due to the prohibitive costs
357 associated with Hi-C sequencing, especially when applied to large volumes of sample, we suggest
358 using iPHoP for family-level host predictions as a suitable alternative. Still, this approach should
359 be interpreted with caution as the concordance between iPHoP and Hi-C sequencing was only
360 assessed using fecal samples. These trends could feasibly differ depending on the environment
361 sampled.

362 Lastly, as phage-host range is often not beyond the species and strain level, by grouping
363 phages at the host family level, this method lacks the sensitivity to detect trends in specific phage-
364 host pairs. Despite these limitations, as bioinformatics methods to detect phage-host pairs improve
365 their resolution, similar approaches to PHFs could be used at lower taxonomic ranks.

366

367 **References**

368

369 1. Shkoporov AN, Clooney AG, Sutton TDS, Ryan FJ, Daly KM, Nolan JA, McDonnell SA,
370 Khokhlova EV, Draper LA, Forde A, Guerin E, Velayudhan V, Ross RP, Hill C. 2019. The
371 human gut virome is highly diverse, stable, and individual specific. *Cell Host Microbe*
372 26:527–541.

373 2. Clooney AG, Sutton TDS, Shkoporov AN, Holohan RK, Daly KM, O’Regan O, Ryan FJ,
374 Draper LA, Plevy SE, Ross RP, Hill C. 2019. Whole-virome analysis sheds light on viral
375 dark matter in inflammatory bowel disease. *Cell Host Microbe* 26:764–778.

376 3. Nishijima S, Nagata N, Kiguchi Y, Kojima Y, Miyoshi-Akiyama T, Kimura M, Ohsugi M,
377 Ueki K, Oka S, Mizokami M, Itoi T, Kawai T, Uemura N, Hattori M. 2022. Extensive gut
378 virome variation and its associations with host and environmental factors in a population-
379 level cohort. *Nat Commun* 13:5252.

380 4. Beller L, Deboutte W, Vieira-Silva S, Falony G, Tito RY, Rymenans L, Yinda CK,
381 Vanmechelen B, Van Espen L, Jansen D, Shi C, Zeller M, Maes P, Faust K, Van Ranst M,
382 Raes J, Matthijnssens J. 2022. The virota and its transkingdom interactions in the healthy
383 infant gut. *Proc Natl Acad Sci* 119:e2114619119.

384 5. Zuo T, Sun Y, Wan Y, Chan FKL, Miao Y, Ng SC. 2020. Human-Gut-DNA Virome
385 Variations across Geography, Ethnicity, and Urbanization. *Cell Host Microbe* 1–11.

386 6. Stockdale SR, Shkoporov AN, Khokhlova EV, Daly KM, McDonnell SA, O’Regan O, Nolan
387 JA, Sutton TDS, Clooney AG, Ryan FJ, Sheehan D, Lavelle A, Draper LA, Shanahan F, Ross

- 388 RP, Hill C. 2023. Interpersonal variability of the human gut virome confounds disease signal
389 detection in IBD. *Commun Biol* 6:221.
- 390 7. Bin Jang H, Bolduc B, Zablocki O, Kuhn JH, Roux S, Adriaenssens EM, Brister JR,
391 Kropinski AM, Krupovic M, Lavigne R, Turner D, Sullivan MB. 2019. Taxonomic
392 assignment of uncultivated prokaryotic virus genomes is enabled by gene-sharing networks.
393 *Nat Biotechnol* 37:632–639.
- 394 8. Marbouty M, Thierry A, Millot GA, Koszul R. 2021. MetaHiC phage-bacteria infection
395 network reveals active cycling phages of the healthy human gut. *eLife* 10:e60608.
- 396 9. Shamash M, Maurice CF. 2021. Phages in the infant gut: a framework for virome
397 development during early life. *ISME J* <https://doi.org/10.1038/s41396-021-01090-x>.
- 398 10. Roux S, Camargo AP, Coutinho FH, Dabdoub SM, Dutilh BE, Nayfach S, Tritt A. 2023.
399 iPHoP: An integrated machine learning framework to maximize host prediction for
400 metagenome-derived viruses of archaea and bacteria. *PLOS Biol* 21:e3002083.
- 401 11. Reyes A, Wu M, McNulty NP, Rohwer FL, Gordon JI. 2013. Gnotobiotic mouse model of
402 phage-bacterial host dynamics in the human gut. *Proc Natl Acad Sci U S A*, 2013/11/20 ed.
403 110:20236–20241.
- 404 12. Hsu BB, Gibson TE, Yeliseyev V, Liu Q, Lyon L, Bry L, Silver PA, Gerber GK. 2019.
405 Dynamic Modulation of the Gut Microbiota and Metabolome by Bacteriophages in a Mouse
406 Model. *Cell Host Microbe* 25:803-814.e5.

- 407 13. Lloyd-Price J, Arze C, Ananthakrishnan AN, Schirmer M, Avila-Pacheco J, Poon TW,
408 Andrews E, Ajami NJ, Bonham KS, Brislawn CJ, Casero D, Courtney H, Gonzalez A,
409 Graeber TG, Hall AB, Lake K, Landers CJ, Mallick H, Plichta DR, Prasad M, Rahnavard G,
410 Sauk J, Shungin D, Vázquez-Baeza Y, White RA, Bishai J, Bullock K, Deik A, Dennis C,
411 Kaplan JL, Khalili H, McIver LJ, Moran CJ, Nguyen L, Pierce KA, Schwager R, Sirota-Madi
412 A, Stevens BW, Tan W, ten Hove JJ, Weingart G, Wilson RG, Yajnik V, Braun J, Denson
413 LA, Jansson JK, Knight R, Kugathasan S, McGovern DPB, Petrosino JF, Stappenbeck TS,
414 Winter HS, Clish CB, Franzosa EA, Vlamakis H, Xavier RJ, Huttenhower C. 2019. Multi-
415 omics of the gut microbial ecosystem in inflammatory bowel diseases. *Nature* 569:655–662.
- 416 14. Uritskiy G, Press M, Sun C, Huerta GD, Zayed AA, Wisner A, Grove J, Auch B, Eacker SM,
417 Sullivan S, Bickhart DM, Smith TPL, Sullivan MB, Liachko I. 2021. Accurate viral genome
418 reconstruction and host assignment with proximity-ligation sequencing. *bioRxiv*
419 <https://doi.org/10.1101/2021.06.14.448389>.
- 420 15. McMurdie PJ, Holmes S. 2013. phyloseq: An R Package for Reproducible Interactive
421 Analysis and Graphics of Microbiome Census Data. *PLoS ONE* 8:e61217.
- 422 16. Dahl EM, Neer E, Bowie KR, Leung ET, Karstens L. 2022. *microshades* : An R Package for
423 Improving Color Accessibility and Organization of Microbiome Data. *Microbiol Resour*
424 *Announc* 11:e00795-22.
- 425 17. Oksanen J, Blanchet FG, Friendly M, Kindt R, Legendre P, McGlinn D, Minchin PR, O’Hara
426 RB, Simpson GL, Solymos P, Stevens MHH, Szoecs E, Wagner H. 2019. *vegan*: Community
427 Ecology Package.

- 428 18. Bolger AM, Lohse M, Usadel B. 2014. Trimmomatic: A flexible trimmer for Illumina
429 sequence data. *Bioinformatics*, 2014/04/01 ed. 30:2114–2120.
- 430 19. Langmead B, Salzberg SL. 2012. Fast gapped-read alignment with Bowtie 2. *Nat Methods*
431 9:357–359.
- 432 20. Beghini F, McIver LJ, Blanco-Míguez A, Dubois L, Asnicar F, Maharjan S, Mailyan A,
433 Manghi P, Scholz M, Thomas AM, Valles-Colomer M, Weingart G, Zhang Y, Zolfo M,
434 Huttenhower C, Franzosa EA, Segata N. 2021. Integrating taxonomic, functional, and strain-
435 level profiling of diverse microbial communities with bioBakery 3. *eLife* 10:e65088.
- 436 21. Li D, Liu C-M, Luo R, Sadakane K, Lam T-W. 2015. MEGAHIT: an ultra-fast single-node
437 solution for large and complex metagenomics assembly via succinct *de Bruijn* graph.
438 *Bioinformatics* 31:1674–1676.
- 439 22. Kieft K, Zhou Z, Anantharaman K. 2020. VIBRANT: automated recovery, annotation and
440 curation of microbial viruses, and evaluation of viral community function from genomic
441 sequences. *Microbiome* 8:90.
- 442 23. Nayfach S, Pedro Camargo A, Eloë-Fadrosch E, Roux S. 2020. CheckV: assessing the quality
443 of metagenome-assembled viral genomes. *bioRxiv* 1–20.
- 444 24. Roux S, Adriaenssens EM, Dutilh BE, Koonin EV, Kropinski AM, Krupovic M, Kuhn JH,
445 Lavigne R, Brister JR, Varsani A, Amid C, Aziz RK, Bordenstein SR, Bork P, Breitbart M,
446 Cochrane GR, Daly RA, Desnues C, Duhaime MB, Emerson JB, Enault F, Fuhrman JA,
447 Hingamp P, Hugenholtz P, Hurwitz BL, Ivanova NN, Labonté JM, Lee KB, Malmstrom RR,
448 Martinez-Garcia M, Mizrachi IK, Ogata H, Páez-Espino D, Petit MA, Putonti C, Rattei T,

- 449 Reyes A, Rodriguez-Valera F, Rosario K, Schriml L, Schulz F, Steward GF, Sullivan MB,
450 Sunagawa S, Suttle CA, Temperton B, Tringe SG, Thurber RV, Webster NS, Whiteson KL,
451 Wilhelm SW, Wommack KE, Woyke T, Wrighton KC, Yilmaz P, Yoshida T, Young MJ,
452 Yutin N, Allen LZ, Kyrpides NC, Eloe-Fadrosh EA. 2019. Minimum information about an
453 uncultivated virus genome (MIUVIG). *Nat Biotechnol* 37:29–37.
- 454 25. Zhou Z, Yu M, Ding G, Gao G, He Y. 2020. Diversity and structural differences of bacterial
455 microbial communities in rhizocompartments of desert leguminous plants. *PLOS ONE*
456 15:e0241057.
- 457 26. Barnett D, Arts I, Penders J. 2021. microViz: an R package for microbiome data visualization
458 and statistics. *J Open Source Softw* 6:3201.
- 459 27. Love MI, Huber W, Anders S. 2014. Moderated estimation of fold change and dispersion for
460 RNA-seq data with DESeq2. *Genome Biol* 15:550.
- 461 28. Hurwitz BL, U'Ren JM. 2016. Viral metabolic reprogramming in marine ecosystems. *Curr*
462 *Opin Microbiol* 31:161–168.
- 463 29. Burns-Huang K, Mundhra S. 2019. Mycobacterium tuberculosis cysteine biosynthesis genes
464 *mec+*-*cysO*-*cysM* confer resistance to clofazimine. *Tuberculosis* 115:63–66.
- 465 30. Allison SD, Martiny JBH. 2008. Resistance, resilience, and redundancy in microbial
466 communities. *Proc Natl Acad Sci* 105:11512–11519.

- 467 31. Tian L, Wang X-W, Wu A-K, Fan Y, Friedman J, Dahlin A, Waldor MK, Weinstock GM,
468 Weiss ST, Liu Y-Y. 2020. Deciphering functional redundancy in the human microbiome. *Nat*
469 *Commun* 11:6217.
- 470 32. Turner D, Kropinski AM, Adriaenssens EM. 2021. A roadmap for genome-based phage
471 taxonomy. *Viruses* 13:506.
- 472 33. Rasmussen TS, Mentzel CMJ, Kot W, Castro-Mejía JL, Zuffa S, Swann JR, Hansen LH,
473 Vogensen FK, Hansen AK, Nielsen DS. 2020. Faecal virome transplantation decreases
474 symptoms of type 2 diabetes and obesity in a murine model. *Gut* 69:2122–2130.
- 475 34. Cornuault JK, Petit M-A, Mariadassou M, Benevides L, Moncaut E, Langella P, Sokol H, De
476 Paepe M. 2018. Phages infecting *Faecalibacterium prausnitzii* belong to novel viral genera
477 that help to decipher intestinal viromes. *Microbiome* 6:65.
- 478 35. Tetz G, Brown SM, Hao Y, Tetz V. 2019. Type 1 Diabetes: an Association Between
479 Autoimmunity, the Dynamics of Gut Amyloid-producing *E. coli* and Their Phages. *Sci Rep*
480 9:9685.
- 481 36. Madi N, Cato ET, Abu Sayeed Md, Creasy-Marrazzo A, Cuénod A, Islam K, Khabir MdIU,
482 Bhuiyan MdTR, Begum YA, Freeman E, Vustepalli A, Brinkley L, Kamat M, Bailey LS,
483 Basso KB, Qadri F, Khan AI, Shapiro BJ, Nelson EJ. 2024. Phage predation, disease severity,
484 and pathogen genetic diversity in cholera patients. *Science* 384:eadj3166.
- 485 37. Tikhomirova A, Rahman MM, Kidd SP, Ferrero RL, Roujeinikova A. 2024. Cysteine and
486 resistance to oxidative stress: implications for virulence and antibiotic resistance. *Trends*
487 *Microbiol* 32:93–104.

- 488 38. Winter SE, Winter MG, Xavier MN, Thiennimitr P, Poon V, Keestra AM, Laughlin RC,
489 Gomez G, Wu J, Lawhon SD, Popova IE, Parikh SJ, Adams LG, Tsolis RM, Stewart VJ,
490 Bäumlér AJ. 2013. Host-Derived Nitrate Boosts Growth of *E. coli* in the Inflamed Gut.
491 Science 339:708–711.
- 492 39. Rivera-Chávez F, Lopez CA, Bäumlér AJ. 2017. Oxygen as a driver of gut dysbiosis. Free
493 Radic Biol Med 105:93–101.
- 494 40. Zünd M, Ruscheweyh H-J, Field CM, Meyer N, Cuenca M, Hoces D, Hardt W-D, Sunagawa
495 S. 2021. High throughput sequencing provides exact genomic locations of inducible
496 prophages and accurate phage-to-host ratios in gut microbial strains. Microbiome 9:77.
- 497 41. Enault F, Briet A, Bouteille L, Roux S, Sullivan MB, Petit MA. 2017. Phages rarely encode
498 antibiotic resistance genes: A cautionary tale for virome analyses. ISME J 11:237–247.
- 499 42. Chen L, Banfield JF. 2024. COBRA improves the completeness and contiguity of viral
500 genomes assembled from metagenomes. Nat Microbiol 9:737–750.
- 501 43. Kieft K, Adams A, Salamzade R, Kalan L, Anantharaman K. 2022. vRhyme enables binning
502 of viral genomes from metagenomes. Nucleic Acids Res 50:e83–e83.
- 503
- 504

505 **Figure Legends**

506

507 **Figure 1. Computationally-predicted bacterial hosts for vOTUs are concordant with *in situ***
508 **associations to the bacterial family level.** Agreement between iPHoP predicted host range and
509 Hi-C assigned host range at various taxonomic ranks for 1,243 vOTUs. Additional comparisons
510 were made when iPHoP predicted multiple hosts for a vOTU (see main text for details on the three
511 comparisons).

512

513

514 **Figure 2. PHFs reduce interindividual variation and increase intraindividual virome**
515 **stability in a cohort of 10 healthy individuals.** Data were analyzed from a previously published
516 study of 10 healthy individuals (1). **(A)** Taxonomic bar plots of virome composition at the PHF
517 level for each individual over time. Facet labels above the bar plots correspond to the subject IDs
518 from the original study. **(B)** Ecological distances between samples with Bray-Curtis at the contig
519 level, Bray-Curtis at the PHF level, and Weighted UniFrac at the PHF level. Interindividual and
520 intraindividual comparisons are both shown. Significance was assessed using the Friedman test
521 with the post-hoc Wilcoxon signed-rank test, using Bonferroni correction for multiple comparisons
522 ($p < 0.001$, ***; $p < 0.0001$, ****). **(C)** Virome stability, defined here as (1 - ecological distance
523 from previous sample), was calculated for each individual using the Bray-Curtis distance metrics
524 at the contig level and at the PHF level. Significance was assessed using the Wilcoxon signed-rank
525 test ($p < 0.01$, **; $p < 0.001$, ***).

526

527 **Figure 3. PHFs are prevalent and reduce intra- and interindividuality in a large human IBD**
528 **cohort.** Data were analyzed from the previously published HMP2 dataset (13). Samples with low
529 read counts (< 1,500) were removed from analyses. In total, bulk metagenomes from 1,093
530 samples from 115 individuals (57 CD, 31 UC, 27 non-IBD controls) were included for downstream
531 analyses. **(A, B)** Rank prevalence distributions of vOTUs **(A)** and PHFs **(B)** across individuals. In
532 total, there were 3,886 distinct vOTUs and 75 distinct PHFs. The dotted red line indicates the rank
533 at which features are more, or less than, 50% prevalent. **(C)** Mean relative abundance of features
534 (PHFs vs. vOTUs) that were present in more than 50% of individuals in the dataset. **(D)** Bray-
535 Curtis distance between samples according to interindividual or intraindividual comparisons.
536 Significance was assessed using the Wilcoxon signed-rank test ($p \leq 0.0001$, ****).

537

538

539 **Figure 4. PHFs reveal disease-specific signatures of IBD.** Data were analyzed from the
540 previously published HMP2 dataset (13). Samples with low read counts (< 1,500) were removed
541 from analyses. **(A)** PCoA plots generated from Bray-Curtis distances matrices using vOTUs (left)
542 and PHFs (right). Samples are color-coded according to the dysbiotic status identified in (13). **(B)**
543 Differentially abundant PHFs based on dysbiosis status. Only individuals which had both a
544 dysbiotic and non-dysbiotic sample were included. Only PHFs that were more than 50% prevalent
545 across individuals were considered for these analyses. PHFs with an adjusted p value ≤ 0.05 and
546 with a \log_2 fold-change ≥ 1 or with a \log_2 fold-change ≤ -1 were considered differentially abundant.

547

548 **Supplementary Figure 1. vOTU membership of PHFs.** Data were analyzed from the previously
549 published HMP2 dataset (13). Samples with low read counts (< 1,500) were removed from
550 analyses. The distribution of PHFs is based on the number of distinct vOTUs that comprises them.
551 The 5 PHFs comprised of the most vOTUs are indicated on the plot.

552

553

554 **Supplementary Figure 2. AMG distribution across PHFs.** Data were analyzed from the
555 previously published HMP2 dataset (13). Only the 12 PHFs that contained > 10 AMGs are shown
556 here. AMGs were detected using VIBRANT, which uses KEGG annotations to assign metabolic
557 categories. **(A)** Distribution of the AMGs found within each PHF. We then determined the number
558 of AMGs per Mb of assembled vOTUs for each PHF, broken down by **(B)** carbohydrate
559 metabolism genes, and **(C)** folding, sorting and degradation genes. **(D)** Number of vOTUs
560 containing *cysO* across PHFs.

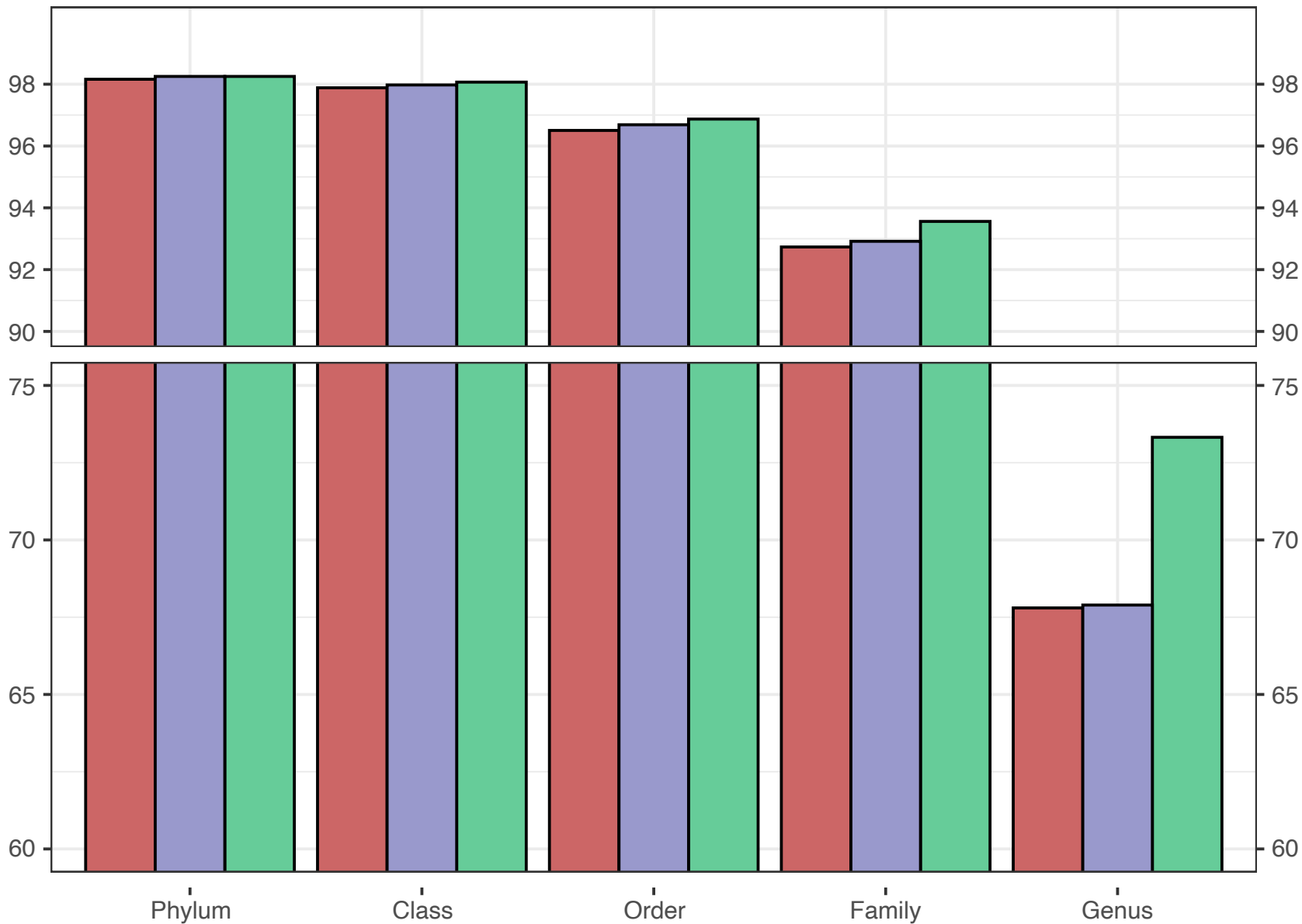
561

562

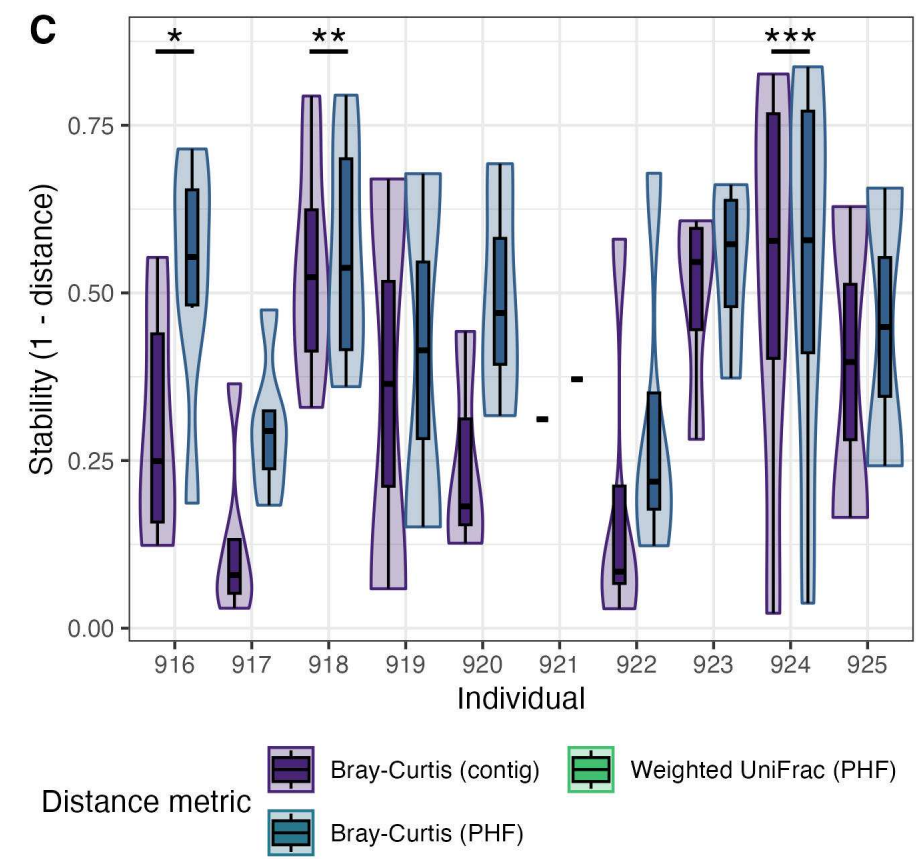
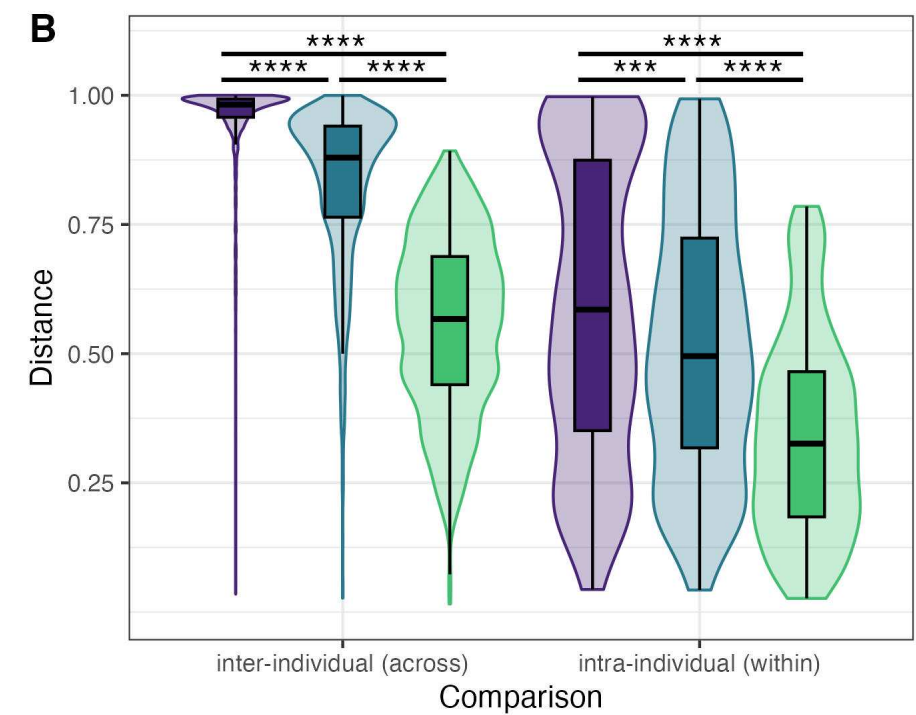
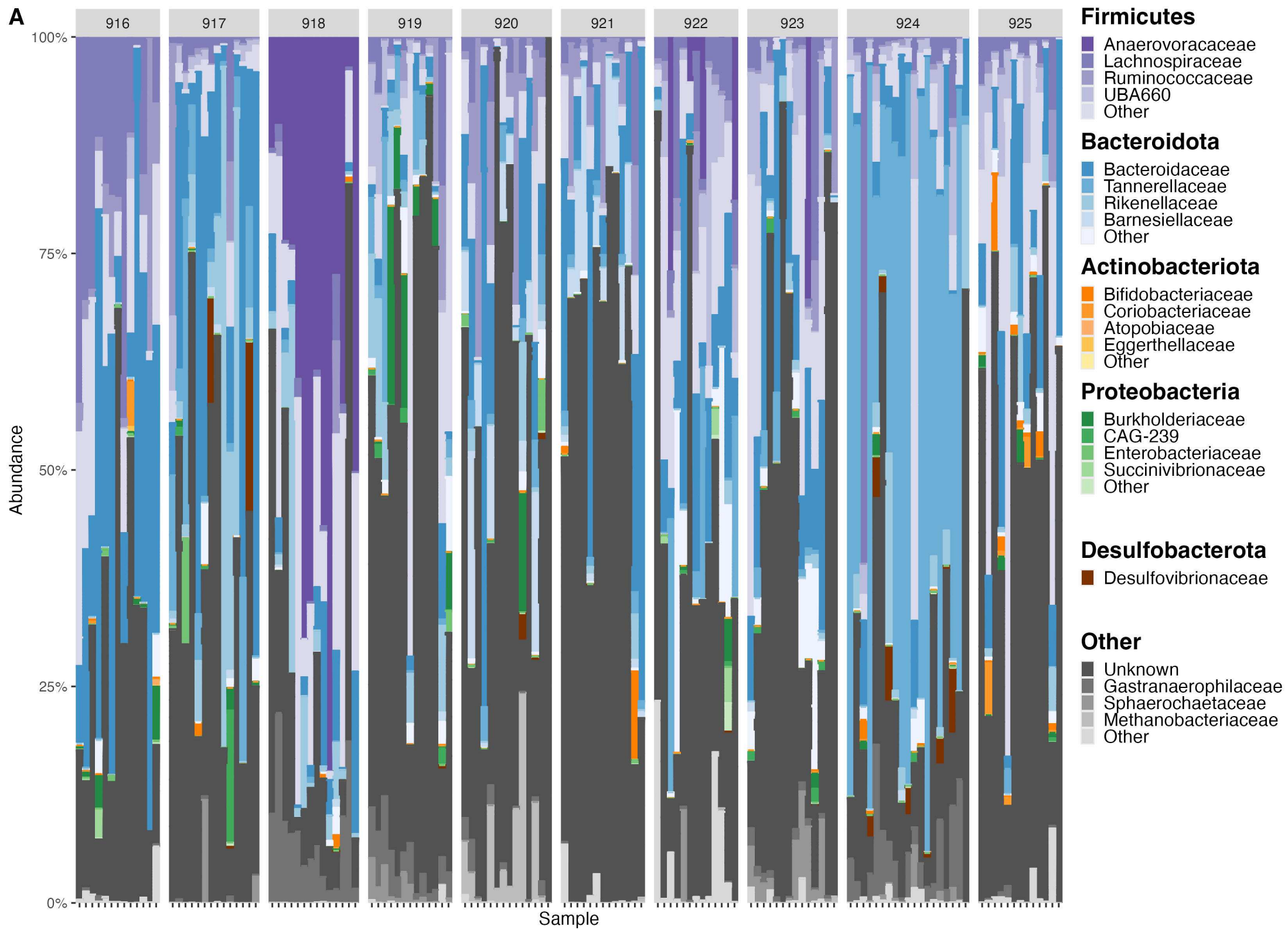
563 **Supplementary Figure 3. Ordination of samples based on patient diagnosis.** Data were
564 analyzed from the previously published HMP2 dataset (13). Samples with low read counts (<
565 1,500) were removed from analyses. PCoA plots were generated from Bray-Curtis distances
566 matrices using vOTUs (left) and PHFs (right). Samples are color-coded according to the diagnosis
567 status identified in (13).

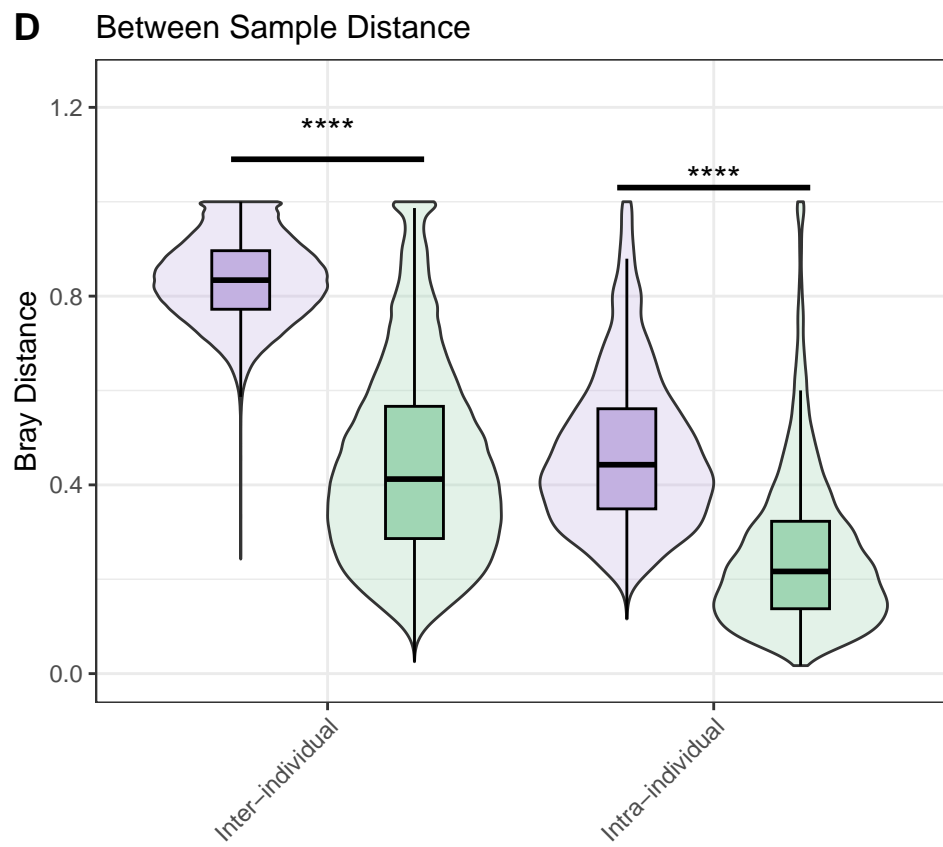
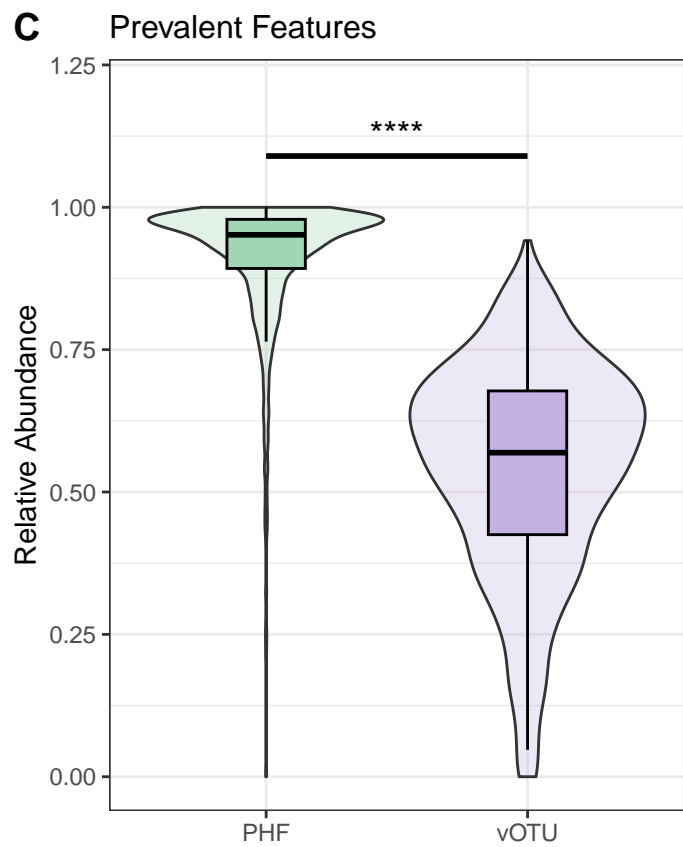
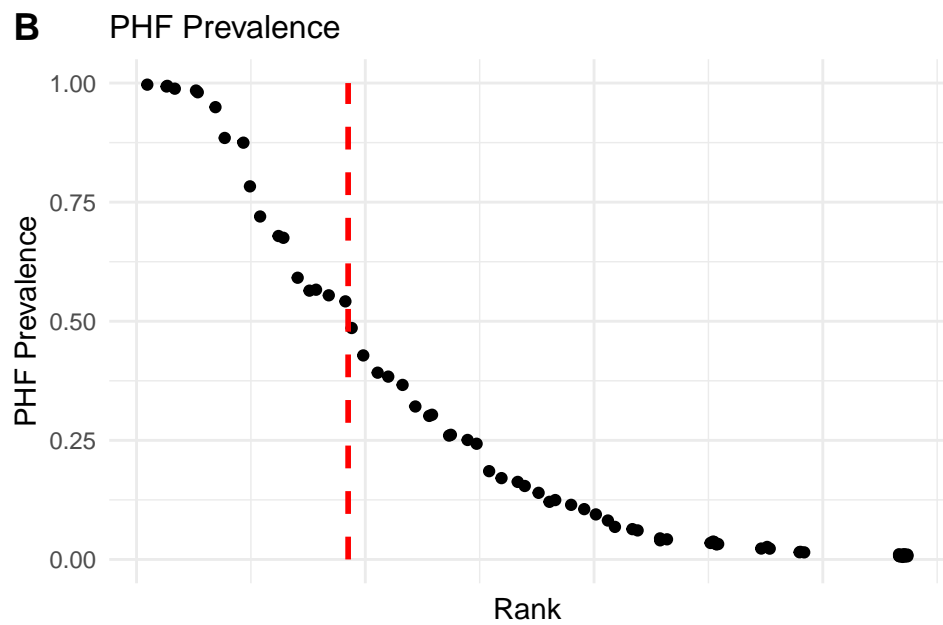
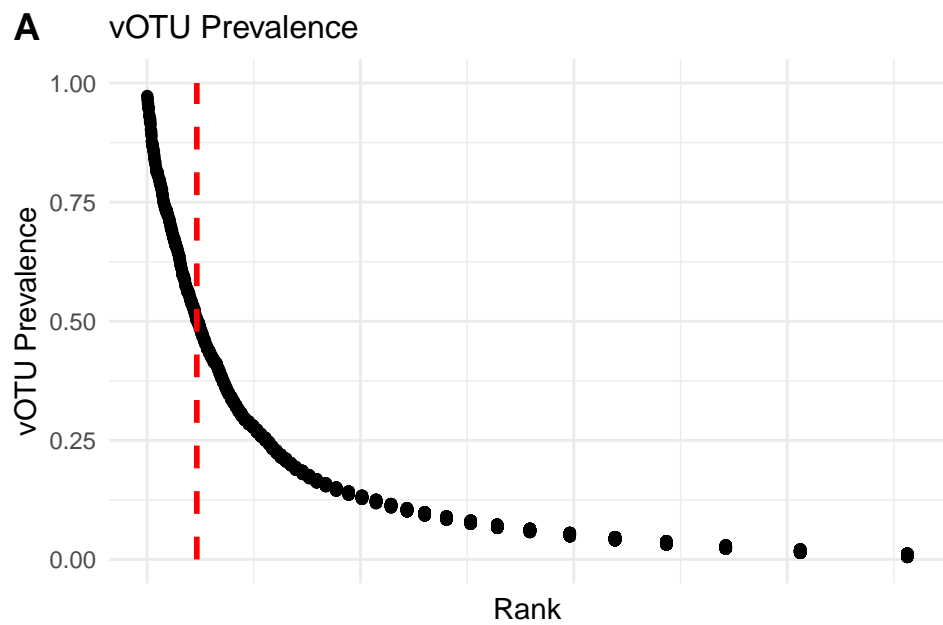
568

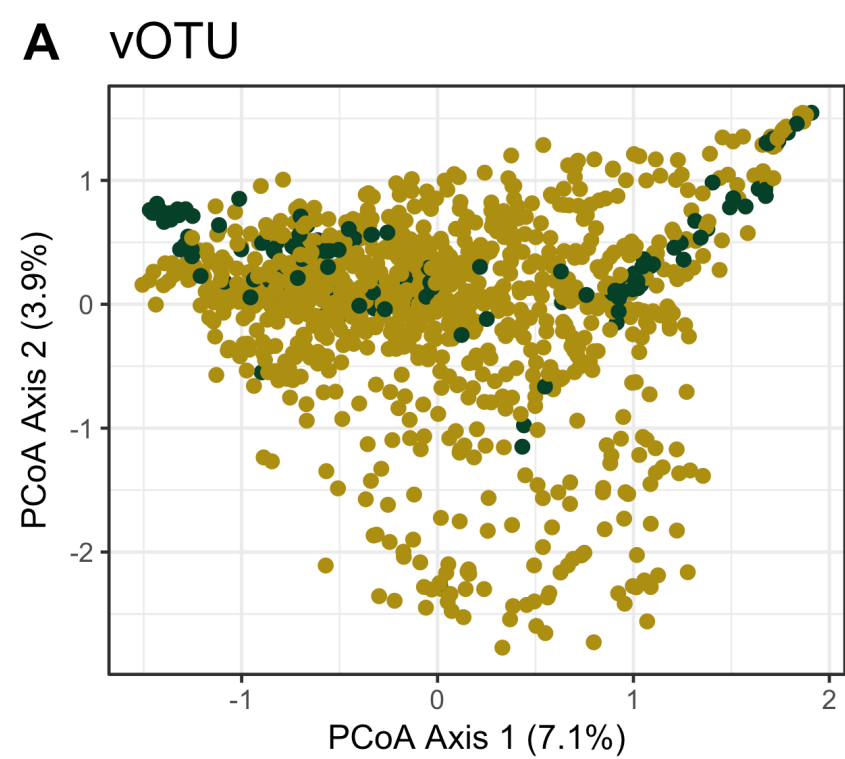
% concordance between iPHoP and Hi-C



Top hit (iPHoP) vs top hit (Hi-C) Top hit (iPHoP) vs all hits (Hi-C) All hits (iPHoP) vs all hits (Hi-C)

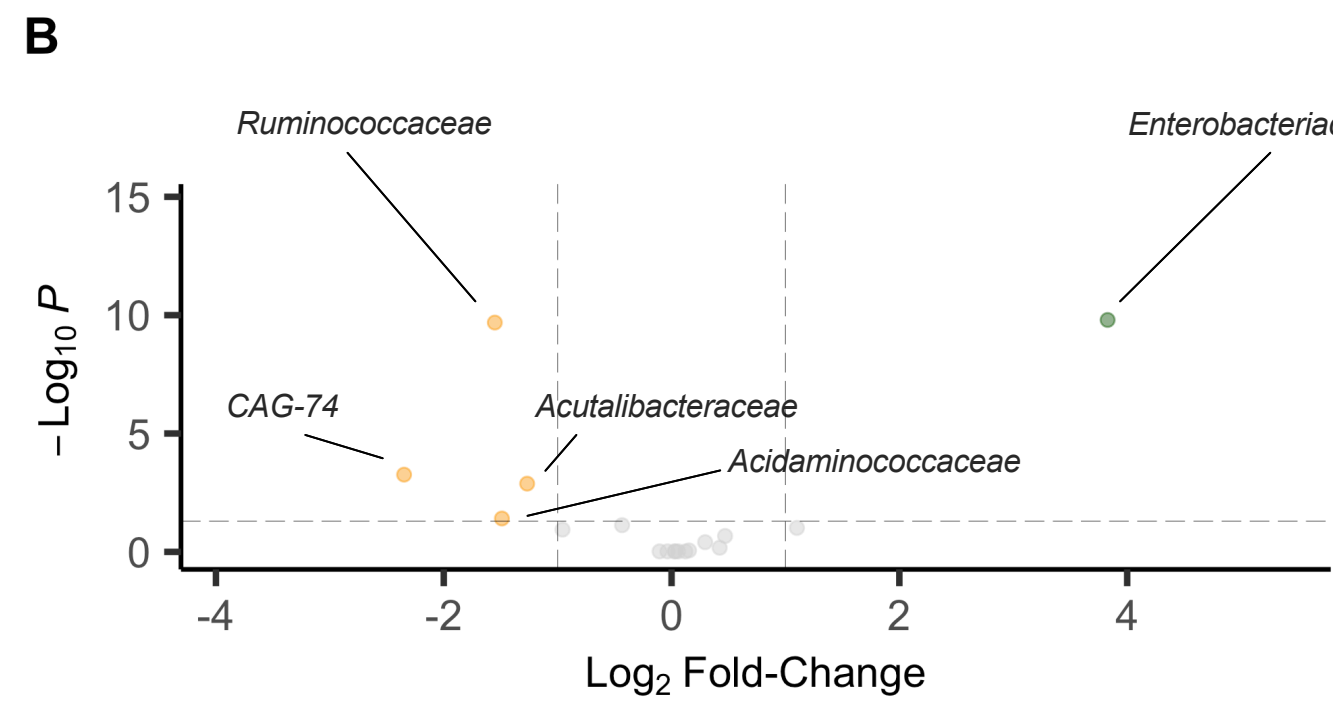
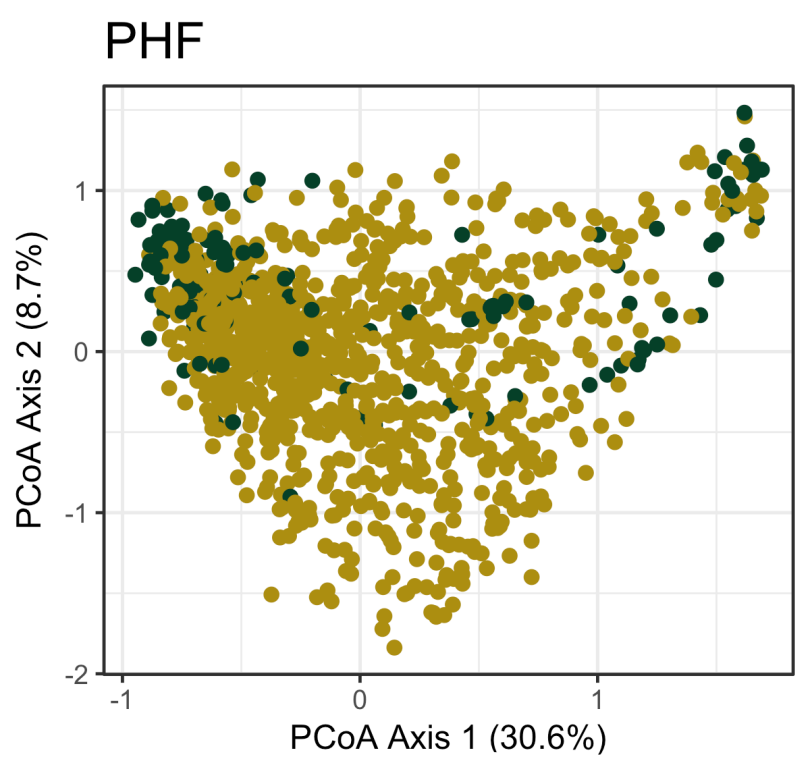






Dysbiosis Status

- Non-Dysbiotic
- Dysbiotic



- Depleted in Dysbiotic Samples
- Enriched in Dysbiotic Samples
- Non-significant