

1 *Bacteroides thetaiotaomicron*-infecting bacteriophage isolates inform sequence-  
2 based host range predictions

3

4

5

6 \*Andrew J. Hryckowian<sup>1,5</sup>, \*Bryan D. Merrill<sup>1</sup>, Nathan T. Porter<sup>2</sup>, William Van  
7 Treuren<sup>1</sup>, Eric J. Nelson<sup>3</sup>, Rebecca A. Garlena<sup>4</sup>, Daniel A. Russell<sup>4</sup>, Eric C.  
8 Martens<sup>2</sup>, Justin L. Sonnenburg<sup>1,5</sup>

9

10

11

12 <sup>1</sup>Department of Microbiology & Immunology, Stanford University School of  
13 Medicine, Stanford, CA

14 <sup>2</sup>Department of Microbiology & Immunology, University of Michigan Medical  
15 School, Ann Arbor, MI

16 <sup>3</sup>Emerging Pathogens Institute, University of Florida, Gainesville, FL

17 <sup>4</sup>Department of Biological Sciences, University of Pittsburgh, Pittsburgh, PA

18 <sup>5</sup>Address correspondence to Andrew Hryckowian

19 ([Andrew.hryckowian@gmail.com](mailto:Andrew.hryckowian@gmail.com)) and Justin Sonnenburg

20 ([jsonnenburg@stanford.edu](mailto:jsonnenburg@stanford.edu)).

21 \*Equal contribution

22

23

## 24 Summary

25

26 Our emerging view of the gut microbiome largely focuses on bacteria and less is  
27 known about other microbial components such as of bacteriophages (phages).

28 Though phages are abundant in the gut, very few phages have been isolated  
29 from this ecosystem. Here, we report the genomes of 27 phages from the United  
30 States and Bangladesh that infect the prevalent human gut bacterium

31 *Bacteroides thetaiotaomicron*. These phages are mostly distinct from previously  
32 sequenced phages with the exception of two, which are crAss-like phages. We  
33 compare these isolates to existing human gut metagenomes, revealing

34 similarities to previously inferred phages and additional unexplored phage  
35 diversity. Finally, we use host tropisms of these phages to identify alleles of  
36 phage structural genes associated with infectivity. This work provides a detailed  
37 view of the gut's "viral dark matter" and a framework for future efforts to further  
38 integrate isolation- and sequencing-focused efforts to understand gut-resident  
39 phages.

40

## 41 Introduction

42

43 Bacteriophages (phages) are highly abundant constituents of free-living  
44 and host-associated microbial communities (microbiomes) (Brussow and  
45 Hendrix, 2002, Barr et al., 2013). Like other microbiome members (e.g. bacteria,  
46 fungi), the diversity and abundance of phages differ between healthy and

47 diseased individuals. While some gut resident phages appear to be unique to  
48 individual humans and stable across long time scales (Shkoporov et al., 2019),  
49 others correlate with host disease status (Manrique et al., 2016, Duerkop et al.,  
50 2018). These observations highlight the possibility that phages of host-  
51 associated microbiomes play central roles in the structure and function of these  
52 communities and may therefore impact human health. Taken together with the  
53 burgeoning antibiotic resistance crisis, this possibility amplifies the importance of  
54 phage therapy as an alternative or supplement to existing paradigms of  
55 microbiome management (e.g. widespread antibiotic use).

56 In contrast to the enthusiasm for phages and phage-based therapeutics is  
57 an underlying reality that gut-resident phages, as a microcosm of the global  
58 phage population, are poorly understood. This deficiency is highlighted when a  
59 typical path for generating and addressing hypotheses in the microbiome field is  
60 considered. For example, sequencing-based approaches are often used in the  
61 microbiome field to generate hypotheses about complex microbial ecosystems.  
62 These hypotheses are subsequently addressed experimentally using a varied  
63 and expanding toolkit of molecular, genetic, immunological, and microbiological  
64 tools. This sequencing-to-mechanism approach is not straightforward for phages  
65 for several reasons. First, unlike bacteria, phages do not have conserved marker  
66 genes (e.g. the 16S rRNA marker gene) that enable phylogenetic classification  
67 and analysis. Instead, phage genomes must be inferred from metagenomic  
68 studies, either based on conservation of phage-like genes (e.g., terminase, DNA  
69 polymerase) (Grazziotin et al., 2017), sequence identity relative to known phage

70 isolates (Roux et al., 2015), or by database-independent approaches (Ren et al.,  
71 2017). While powerful for general characterization of changes in the composition  
72 of phage communities, inter-study methodological variation (e.g. sample  
73 preparation, contig assembly, reference databases used) can impact a study's  
74 conclusions to a greater extent than the treatment effects (e.g. health or disease  
75 status) (Gregory et al., 2019).

76 Furthermore, metagenomic approaches fail to provide definitive  
77 information on the bacterial hosts of these phages. To address this deficiency,  
78 many methods have been developed to predict the bacterial hosts of phages  
79 inferred from metagenomes. For example, homology searches, identification of  
80 CRISPR spacers, and co-occurrence analysis were used to make the prediction  
81 that the highly prevalent and abundant crAssphage infects bacteria in the phylum  
82 *Bacteroidetes* (Dutilh et al., 2014). This prediction was validated in part when a  
83 crAss-like phage (CrAss001) was isolated on *Bacteroides intestinalis* (Shkoporov  
84 et al., 2018). However, based on the divergence of CrAss001 from the  
85 prototypical crAssphage combined with the astounding diversity of crAss-like  
86 phages across host associated and environmental microbiomes (Yutin et al.,  
87 2018, Guerin et al., 2018), it is likely that other crAss-like phages infect other  
88 bacterial strains within the *Bacteroidetes* phylum. Furthermore, crAss-like phages  
89 can simultaneously be biomarkers of healthy and diseased states. For example,  
90 one crAss-like phage, IAS virus, is enriched in HIV<sup>+</sup> individuals with low CD4  
91 counts (Oude Munnink et al., 2014) while some crAss-like phages are stable over  
92 a 12 month period in healthy humans (Shkoporov et al., 2019).



93 CrAss001 is one of four isolated and sequenced phages confirmed to  
94 infect *Bacteroides*, the most abundant bacterial genus in the human gut  
95 microbiome. The other three phages are B40-8 and B124-14 (which infect *B.*  
96 *fragilis*) and Hankyphage, which is present as a prophage in many *Bacteroides*  
97 strains (Benler et al., 2018, Ogilvie et al., 2012, Hawkins et al., 2008). Despite  
98 the prevalence of Hankyphage lysogens in the public databases, Hankyphage  
99 induced from a Hankyphage-containing *B. dorei* lysogen was unable to form  
100 plaques on Hankyphage-naïve *B. dorei* or additional *Bacteroides* species (Benler  
101 et al., 2018). Additionally, CrAss001 does not form robust plaques on *B.*  
102 *intestinalis* despite persisting at high levels in co-culture with its host for nearly  
103 one month (Shkoporov et al., 2018). These observations suggest that unexplored  
104 factors influence *Bacteroides*-phage interactions.

105 Additional *Bacteroides* phages were previously isolated but not  
106 sequenced, representing limitations in connecting experiment-focused studies  
107 with bioinformatics-focused studies. Our recent work incorporated 71 phage  
108 isolates to show that multiple phase-variable mechanisms, including capsular  
109 polysaccharides (CPS), modify bacteriophage susceptibility in *B.*  
110 *thetaitaomicron*. These findings mirror work in other isolated phages where cps  
111 is a major determinant of phage-host tropism (Porter et al., 2019). However, in  
112 the absence of genome sequences, phage-encoded determinants of host tropism  
113 were not previously explored.

114 Here, we report the genomes of 27 phages that infect *B. thetaitaomicron*  
115 (18 previously described isolates (Porter et al., 2019) and 9 new isolates; **Table**

116 **S1)**. By comparing these genomes with those of existing *Bacteroides* phage  
117 isolates and with phage genomes identified from publicly available metagenomic  
118 studies, we simultaneously reveal similarities to previously inferred phage  
119 genomes and additional unexplored phage diversity. Finally, analysis of these  
120 genomes in the context of their cps-specific host ranges reveals targets for future  
121 study aimed at understanding the structure-function relationship of phage host  
122 range and phylogeny. We suggest the utility and feasibility of future efforts that  
123 integrate both isolation- and computational-based methods. Such an approach  
124 would enrich databases of known phage-host pairs by providing additional  
125 reference genomes and definitive host information. Furthermore, isolated phages  
126 will enable investigators to build experimental systems to test hypotheses and  
127 theoretical predictions and contribute to a growing collection of phages that may  
128 be used in the future for therapeutic or biotechnological applications.

129

## 130 Results

131

132 *Isolation and comparative analysis of 27 phages infecting B. thetaiotaomicron*

133

134 Our study centers of phages isolated from four geographic locations, three  
135 within the US and one in Bangladesh (**Figs. 1A, Table S1**). Using a previously  
136 reported protocol for phage isolation (Porter et al., 2019), we isolated 9  
137 bacteriophages from primary wastewater effluent from the Sand Island  
138 Wastewater Treatment Plant (Honolulu, Hawaii) or from sewer-adjacent pond

139 water at two locations in Dhaka, Bangladesh. High titer stocks were prepared of  
140 these 9 phage isolates and of a subset of 17 phages from an existing collection  
141 of 71 *B. thetaiotaomicron*-infecting phages isolated from Ann Arbor, Michigan  
142 and San Jose, California (Porter et al., 2019). Phage genomes were sequenced  
143 and assembled. The phages are grouped into three genomically related clusters  
144 ( $\alpha$ ,  $\beta$ ,  $\gamma$ ) and have genomes that are on average 38kb +/- 0.4kb, 99kb +/- 0.3kb,  
145 and 177kb +/- 4.5kb, respectively, and exhibit extensive genomic mosaicism  
146 (**Figs. S1-S3; Table S1**). Transmission electron microscopy of one  
147 representative from each cluster reveals distinct virion morphologies. Based on  
148 these representatives, cluster  $\alpha$  phages are siphoviruses, cluster  $\beta$  phages are  
149 podoviruses, cluster  $\gamma$  phages are myoviruses, and the capsid sizes of these  
150 phages scale with genome size (**Figure 1B-D**).

151 Phage genomes were annotated and compared on the basis of shared  
152 gene content (pham membership) (Cresawn et al., 2011). Phams are built and  
153 expanded when a candidate protein shares  $\geq 32.5\%$  identity or blastp e-value  
154  $\leq 1e-50$  with one or more existing members of the pham. A dendrogram was built  
155 based on the presence or absence of each pham in each phage, which  
156 confirmed the three distinct genome clusters (**Figs. 1E, S1-S3**). Genome maps of  
157 representatives of each of these clusters are shown in **Figs. 1F-H**. tRNAs were  
158 detected in cluster  $\beta$  and  $\gamma$  phages (n=12-13 and n=2-3, respectively) but not in  
159 cluster  $\alpha$  phages (**Tables S1, S2**).

160 While there is a high degree of intra-cluster sequence identity, there are  
161 only two phams shared between clusters: pham 150 (encoding a putative

162 thymidylate synthase) and pham 23 (encoding collagen triple helix repeats),  
163 which are shared between all cluster  $\beta$  and  $\gamma$  representatives. Consistent with  
164 observations from previously isolated phages (Hatfull and Hendrix, 2011), the  
165 majority (roughly 80%) of phams in these *B. thetaiotaomicron*-infecting phages  
166 have no detectable conserved domains or known functions (**Tables S3-S5**).

167

168 *Comparative analysis of B. thetaiotaomicron phages with existing Bacteroides*  
169 *phage isolates.*

170

171 We compared these 27 *B. thetaiotaomicron*-infecting phages to 4 other  
172 previously sequenced *Bacteroides*-infecting phages (Benler et al., 2018, Ogilvie  
173 et al., 2012, Hawkins et al., 2008, Shkoporov et al., 2018) (**Fig. 2**). We noted  
174 extensive shared phams (n=53) and genome organization between the cluster  $\beta$   
175 phages (DAC15 and DAC17) and CrAss001 (**Fig. 2, Fig. S4, Table S6**),  
176 reinforcing previous predictions that at least a subset of crAss-like phages prey  
177 on *Bacteroides*. A small number of phams are shared between the other isolated  
178 *B. thetaiotaomicron*-infecting phages and the previously isolated *Bacteroides*-  
179 infecting phages. Pham 22 is present in CrAss001 and the cluster  $\gamma$  phages.  
180 Pham 400 is present in cluster  $\gamma$  phages and in B124-14/B40-8, pham 887 is  
181 shared between B124-14 and the cluster  $\alpha$  phages, pham 423 is shared between  
182 cluster  $\alpha$  phages and B40-8, and pham 394 is shared between Hankyphage and  
183 the cluster  $\alpha$  phages (**Table S6**). Based on this lack of relatedness, our cluster  $\alpha$   
184 and cluster  $\gamma$  phages represent the first isolates of two novel clades of

185 *Bacteroides*-infecting phages. Furthermore, B40-8 and B124-14 are members of  
186 a separate cluster (cluster  $\delta$ ) and Hankyphage is a singleton with no isolated  
187 relatives (**Fig. 2**). These cluster assignments are validated by vConTACT2 (Bin  
188 Jang et al., 2019) (see **Methods**). No RefSeq phage genomes from the  
189 ProkaryoticViralRefSeq94-Merged database were grouped into clusters with  
190 these 31 isolated phages.

191

192 *Identification of phages related to isolated B. thetaiotaomicron phages in existing*  
193 *metagenomes.*

194 Because the majority of phage-focused work in the gut microbiome field is  
195 based on metagenomic sequencing, we wondered if relatives of the sequenced  
196 *B. thetaiotaomicron*-infecting phage isolates could be found in existing  
197 metagenomes. To identify relatives of these phages, we used the protein search  
198 feature of SearchSRA (Torres et al., 2017, Levi et al., 2018, Towns et al., 2014,  
199 Stewart et al., 2015, Buchfink et al., 2015b, Langmead and Salzberg, 2012) to  
200 map 100,000 subsampled reads from each of the ~100,000 metagenomes in the  
201 Sequence Read Archive (SRA) onto representatives of clusters  $\alpha$ ,  $\beta$ , and  $\gamma$   
202 (SJC01, DAC15, and DAC20, respectively). We identified 812 candidate  
203 metagenomes in the SRA where at least one of the representative phage  
204 genomes was covered by reads at an estimated real depth of >15% (given the  
205 true sequencing depth of the sample) and the percent of the genome detected  
206 was >20% for SJC01, DAC15, or DAC20 (**Fig. 3A-C**). We subsequently focused  
207 on human gut-derived metagenomes possessing sequences that are SJC01-like

208 (>50% detected, >30x estimated coverage), DAC15-like (>40% detected, >15x  
209 estimated coverage), or DAC20-like (>20% detected) genomes for further  
210 analysis (**Table S7**). These metagenomes were downloaded from NCBI and  
211 assembled. Contigs containing significant hits (blastp e-value <1e-3) for >25% of  
212 the genes in SJC01, DAC15, or DAC20 were compared to the genomes of the  
213 isolated *Bacteroides*-infecting phages described above. See Methods for a more  
214 detailed description of this method of identifying Phage in SearchSRA (PhiSh).  
215 Several PhiSh genomes were identified which are related to SJC01. These  
216 genomes include previously uncharacterized contigs from prior studies (PhiSh01  
217 – PhiSh03, PhiSh05 – PhiSh07)(Monaco et al., 2016, He et al., 2017, Liu et al.,  
218 2016, Zheng et al., 2017, Guthrie et al., 2017) and a genome previously identified  
219 in a study examining the rapid evolution of the human gut virome (PhiSh04)  
220 (Minot et al., 2013). Serendipitously, we noticed that HSC01, a genome of a  
221 phage predicted to infect *Bacteroides caccae* (Reyes et al., 2013) is related to  
222 these cluster  $\alpha$  isolates and cluster  $\alpha$ -like PhiSh genomes (**Figs. 3DE; Table S8**).  
223 vConTACT2 also places all of these SJC01-like PhiSh genomes within cluster  $\alpha$ .  
224 Six DAC15-like genomes were identified with this method (PhiSh08 –  
225 PhiSh13) (**Table S8**). Five of these genomes (PhiSh08 – PhiSh12) were  
226 previously identified in a study aimed at identifying crAss-like phages in human  
227 fecal metagenomes (Guerin et al., 2018) while PhiSh13 represents a novel  
228 crAss-like phage genome (He et al., 2017). Importantly, these DAC15-like PhiSh  
229 genomes are diverse (they can be classified into the previously described  
230 candidate crAss-like genera 6, 7, and 10; **Table S8**) and are differentially

231 clustered by vConTACT2 (clusters  $\beta$  and  $\epsilon$ ), demonstrating that the PhiSh  
232 identification approach is capable of detecting genomes that are closely and  
233 distantly related to the PhiSh bait genome used (**Fig. 3D**).

234 Despite identifying a diverse collection of cluster  $\alpha$ -like and cluster  $\beta$ -like  
235 PhiSh genomes, only partial  $\gamma$ -like PhiSh genomes were identified (**Fig. 3C**). The  
236 lack of full-length  $\gamma$ -like PhiSh genomes may be due to insufficient sequencing  
237 depth of the original studies or the presence of highly divergent phages which  
238 share subsets of genes with cluster  $\gamma$  phages.

239

#### 240 *Identification of infection-associated phams*

241 Our previous work demonstrated that multiple phase-variable  
242 mechanisms, including capsular polysaccharides (*cps*), modify bacteriophage  
243 susceptibility in *B. thetaiotaomicron* (Porter et al., 2019). However, phage-  
244 encoded determinants of host tropism in these phages were previously  
245 unexplored. When the *cps* specificities of these phages are compared with  
246 genome cluster membership (**Fig. 1B, Table S1**), relationships between genome  
247 cluster membership and host range become evident (**Fig. 4A**). For example,  
248 cluster  $\gamma$  phages tend to be most restrictive in their host range, primarily infecting  
249 *cps7*, *cps8*, and acapsular *B. thetaiotaomicron*. Cluster  $\beta$  phages are similarly  
250 restricted in their host range but are unique in their ability to efficiently infect *B.*  
251 *thetaiotaomicron cps3*. Some cluster  $\alpha$  phages have promiscuous host ranges  
252 while other cluster  $\alpha$  phages have restrictive host ranges (more similar to those of  
253 the cluster  $\beta$  and cluster  $\gamma$  phages). This variation in host range among cluster  $\alpha$

254 phages prompted us to search for phams that are associated with the different  
255 infection patterns in the cluster  $\alpha$  phages.

256 We noted two major themes driving variation among the cluster  $\alpha$  phages:  
257 variation between shared predicted structural components in these phages, such  
258 as gene products (gps) 4, 5, and 8; and mosaicism in genes at the 3' end of the  
259 genomes, representing genes encoding small hypothetical proteins and genes  
260 encoding predicted DNA methylases. (**Fig. S1**). Therefore, we considered the  
261 possibility that allelic variation and presence/absence of phams could contribute  
262 to differences in host range among the phages. To account for each of these  
263 possibilities, we used an algorithm to identify infection-associated phams (IAPs).  
264 Specifically, we computed phams at alternative cutoffs such that membership  
265 was dictated by only by varying levels of amino acid (AA) identity between (25  
266 and 100%; see **Methods**). As the threshold value increases, the total number of  
267 phams increases, with a concomitant decrease in mean pham membership (**Fig.**  
268 **4B**), as previously observed (Cresawn et al., 2011). With the possibility that  
269 different thresholds may reveal allelic variants that correspond to infectivity, we  
270 compared these alternative pham tables with infection thresholds. This approach  
271 identified 662 total phams across all 64 infectivity/pham threshold comparisons in  
272 the 19 cluster  $\alpha$  phages. Of these, 135 were identified as IAPs.

273 Among the IAPs is cluster  $\alpha$  gp8, which is present in all cluster  $\alpha$  phages,  
274 exhibits substantial sequence variation among these phages, and is predicted to  
275 encode a tail protein (**Figs. 1F, 4C, S1**). At the 85% AA identity cutoff, gp8 is  
276 grouped into two distinct phams and phages that have the SJC01-like variant



277 infect *cps1*, *cps5*, *cps6*, and  $\Delta$ *cps* *B. thetaiotaomicron* more efficiently than those  
278 that do not (**Fig. 4D**). Analysis of this IAP in the context of metagenome-derived  
279 cluster  $\alpha$ -like phages reveals additional variation not represented in our isolates  
280 (**Fig. 4E**). Interestingly, the variants of this IAP in PhiSh02 and HSC01 contain  
281 Bacteroides-Associated Carbohydrate Binding Often N-terminal (BACON)  
282 domains (See (Reyes et al., 2013) and **Supplementary Data 1**). These  
283 combined observations suggest a role for this IAP in differential recognition of  
284 complex polysaccharides (e.g. capsular polysaccharides) across confirmed  
285 *Bacteroides*-infecting phages and related genomes.

286

## 287 Discussion

288

289 In this work, we integrate phenotypic and genomic analysis of isolated  
290 phages with metagenomic analysis to highlight several opportunities for future  
291 study of gut-resident phages. In particular, though metagenome-focused studies  
292 of phages continue to generate tremendous insights into the composition and  
293 dynamics of viromes in the gut and other ecosystems, they are limited in scope  
294 due to a lack of definitive connections between predicted phages and their  
295 bacterial hosts. Several approaches have been developed to predict phage host  
296 range (Edwards et al., 2016). These approaches have been validated in part,  
297 notably for CrAss001, which was isolated on *B. intestinalis* after predictions that  
298 crAss-like phages infect members of the phylum Bacteroidetes (Dutilh et al.,  
299 2014, Yutin et al., 2018). We further validate these predictions with two more

300 crAss-like phages, DAC15 and DAC17, which infect *B. thetaiotaomicron*. This  
301 brings the total number of published isolates of this highly abundant viral family to  
302 three (**Figs. 2, S4**). As more crAss-like phages are isolated, we anticipate that  
303 existing discrepancies relating to the roles of these phages in the gut (e.g. some  
304 crAss-like phages are associated with disease (Oude Munnink et al., 2014) while  
305 others are stably maintained in healthy individuals (Shkoporov et al., 2019)) can  
306 be disentangled with controlled experimental approaches.

307         Similarly, based on our genomic and metagenomic analysis of cluster  $\alpha$   
308 phages, we show that two previously reported phage genomes (PhiSh04 and  
309 HSC01) are related despite differences in temporal dynamics and predicted host  
310 range (**Fig. 3D**) (Reyes et al., 2013, Minot et al., 2013). This raises questions  
311 regarding what determinants encoded by phage or bacterial host are responsible  
312 for observed differences in host range and phage population dynamics. Some  
313 insights come from experiments using the cluster  $\alpha$  phage ARB25 in gnotobiotic  
314 mice. ARB25 is stably maintained in a bi-colonization with its host for months and  
315 the phase-variable mechanisms used by *B. thetaiotaomicron* to evade ARB25  
316 are dependent on the presence of CPS among other cell surface features (Porter  
317 et al., 2019). HSC01, unlike PhiSh04 or ARB25, does not stably co-exist with its  
318 predicted host (*B. caccae*) in gnotobiotic mice (Reyes et al., 2013), suggesting  
319 that although HSC01 is related to phages that are stably maintained, it may have  
320 distinct ecological impacts in the gut. Alternatively, it is possible that other  
321 members of the gut microbiome affect the relationship between a phage and its  
322 host bacterium.

323 By combining work that involves phage isolation, sequencing, and  
324 phenotypic characterization, with metagenomic analyses, we hope to reciprocally  
325 inform these studies (e.g., by adding phages and information on IAPs to  
326 publically available databases) and to provide the reagents necessary to  
327 experimentally test hypotheses using the broad toolkit available in the gut  
328 microbiome field (e.g., by probing phage-host interactions using gnotobiotics and  
329 molecular genetics). Future isolation efforts can be further optimized with high  
330 throughput approaches (e.g. robotics and automated liquid handling) or as part of  
331 educational efforts like those pioneered by the SEA-PHAGES program (Hanauer  
332 et al., 2017), which would simultaneously crowd source the effort while providing  
333 training opportunities for the next generation of microbiome scientists. Together,  
334 this integration will allow for a more comprehensive consideration of the  
335 interactions that occur between phages and their hosts at the population,  
336 individual, and molecular scales.

337

### 338 Methods

#### 339 *Bacterial strains and culture conditions.*

340 The bacterial strains used in this study are listed in **Table S9**. Frozen  
341 stocks of these strains were maintained in 25% glycerol at -80°C and were  
342 routinely cultured in an anaerobic chamber (Coy) under 5% H<sub>2</sub>, 10% CO<sub>2</sub>, 85%  
343 N<sub>2</sub> at 37°C in *Bacteroides* Phage Recovery Medium (BPRM), as described  
344 previously (Porter et al., 2019): per 1 liter of broth, 10 g meat peptone, 10 g  
345 casein peptone, 2 g yeast extract, 5 g NaCl, 0.5 g L-cysteine monohydrate, 1.8 g

346 glucose, and 0.12 g MgSO<sub>4</sub> heptahydrate were added; after autoclaving and  
347 cooling to approximately 55 °C, 10 ml of 0.22 µm-filtered hemin solution (0.1%  
348 w/v in 0.02% NaOH), 1 ml of 0.22 µm-filtered 0.05 g/ml CaCl<sub>2</sub> solution, and 25 ml  
349 of 0.22µm-filtered 1 M Na<sub>2</sub>CO<sub>3</sub> solution were added. For BPRM agar plates, 15  
350 g/L agar was added prior to autoclaving and hemin and Na<sub>2</sub>CO<sub>3</sub> were added as  
351 above prior to pouring the plates. For BPRM top agar used in soft agar overlays,  
352 3.5 g/L agar was added prior to autoclaving. Hemin, CaCl<sub>2</sub>, and Na<sub>2</sub>CO<sub>3</sub> were  
353 added to the top agar as above immediately before conducting experiments.  
354 Bacterial strains were routinely struck from the freezer stocks onto BPRM agar  
355 and grown anaerobically for up to 2 days. A single colony was picked for each  
356 bacterial strain, inoculated into 5 mL BPRM, and grown anaerobically overnight  
357 to provide the starting culture for experiments.

358

359 *Bacteriophage isolation from primary wastewater effluent and sewer-adjacent*  
360 *pond water*

361 The bacteriophages described in this study were isolated from primary  
362 wastewater effluent from the Ann Arbor, Michigan Wastewater Treatment Plant  
363 and from the San Jose-Santa Clara Regional Wastewater Treatment Facility, as  
364 described previously (Porter et al., 2019). For the current study, phages were  
365 isolated from primary wastewater effluent from the Sand Island Wastewater  
366 Treatment Plant (Honolulu, Hawaii) or from sewer-adjacent pond water in Dhaka,  
367 Bangladesh (**Table S1**). Water samples were centrifuged at 5,500 rcf for 10  
368 minutes at room temperature to remove any remaining solids. The supernatant

369 was then sequentially filtered through 0.45  $\mu\text{m}$  and 0.22  $\mu\text{m}$  polyvinylidene  
370 fluoride (PVDF) filters. This processed primary effluent was concentrated up to  
371 500-fold via 100 kDa PVDF size exclusion columns.

372 Initial screening for plaques was done using a soft agar overlay method  
373 where 50  $\mu\text{L}$  of the concentrated primary effluent was combined with 0.5 mL  
374 overnight culture and 4.5 mL BPRM top agar and poured onto a standard circular  
375 petri dish [100 mm x 15 mm]. Soft agar overlays were incubated anaerobically at  
376 37 °C overnight. To promote a diverse collection of phages, no more than 5  
377 plaques from the same plate were plaque purified and a diversity of plaque  
378 morphologies were selected as applicable.

379 Single, isolated plaques were picked into 100  $\mu\text{L}$  phage buffer (prepared  
380 as an autoclaved solution of 5 ml of 1 M Tris pH 7.5, 5 ml of 1 M  $\text{MgSO}_4$ , 2 g  
381 NaCl in 500 ml with ddH<sub>2</sub>O). Phages were plaque purified using a 96-well plate-  
382 based method, where serial dilutions were prepared in 96-well plates and 1  $\mu\text{L}$  of  
383 each dilution was spotted onto a solidified top agar overlay. This procedure was  
384 repeated at least 3 times to plaque purify each phage.

385 High titer phage stocks were generated by flooding a soft agar overlay on  
386 a plate that yielded a “lacey” pattern of bacterial growth (near confluent lysis).  
387 Following overnight incubation of each plate, 5 ml of sterile phage buffer was  
388 added to the plate to re-suspend the phage. After at least 2 hours of incubation at  
389 room temperature, the lysate was spun at 5,500 rcf for 10 minutes to clear debris  
390 and then filter sterilized through a 0.22  $\mu\text{m}$  PVDF filter. For more details on  
391 phages used in this work, see **Table S1**.

392

393 *Phage genome sequencing and assembly*

394 DNA was extracted from high-titer phage lysates and sequencing libraries  
395 were prepared using the Ultra II FS Kit (New England Biolabs) or for ARB14 and  
396 ARB25, the TruSeq Nano DNA LT Kit (Illumina). Libraries were quantified using a  
397 BioAnalyzer (Agilent) and subsequently sequenced using 150-base single-end  
398 reads (Illumina MiSeq), or for ARB14 and ARB25, 250-base paired-end reads  
399 (Illumina MiSeq). Phage genomes were assembled using Geneious version 9.1.5  
400 with default options after trimming reads with an error probability limit of 0.05. All  
401 genomes published here circularized during assembly. Phage genomes  
402 belonging to the same cluster were rearranged to have identical 5' ends.  
403 Coverage for each assembly was calculated by mapping reads onto each  
404 assembled genome using bowtie2 (Langmead and Salzberg, 2012) (--very-  
405 sensitive) and then using jgi\_summarize\_ban\_contig\_depths from the MetaBAT2  
406 tool (Kang et al., 2019) to calculate mean coverage depth.

407

408 *Annotation and comparative analyses of B. thetaiotaomicron infecting phages*

409 Protein-coding genes and tRNAs were predicted and annotated using  
410 DNA-Master default parameters (<http://cobamide2.pitt.edu/>), which incorporates  
411 Genemark (Besemer and Borodovsky, 2005), Glimmer (Delcher et al., 1999),  
412 and tRNAscan-SE (Lowe and Eddy, 1997). To classify genes into related groups  
413 (phams) and identify conserved domains, Phamerator was used with default  
414 parameters (Cresawn et al., 2011). Phage genome ends and packaging

415 strategies for cluster  $\beta$  phages were inferred using PhageTerm (Garneau et al.,  
416 2017) which identified clear direct terminal repeats (DTRs). PhageTerm was  
417 unable to identify DTRs or cohesive ends in the cluster  $\alpha$  or  $\gamma$  phages, possibly  
418 indicating a headful packaging strategy. The large terminase proteins share  
419 significant similarity (BLASTP e-value  $<1e-3$ ) with the PBSX-family of large  
420 terminases, which also use a headful packaging strategy (**Table S1**) (Anderson  
421 and Bott, 1985). To predict virion structural genes, iVireons was used with default  
422 parameters (Seguritan et al., 2012). Protein-coding genes were classified as  
423 “predicted structural genes” (e.g. general structural, tail, or capsid, annotated in  
424 **Fig. 1**) for genes with score 0.7 and above. To visualize genome-level  
425 relationships among phages, pham tables were processed with Janus  
426 (<http://cobamide2.pitt.edu/>) and Splitstree using default parameters. Phage  
427 genomes were clustered together using vConTACT2 and the  
428 ProkaryoticViralRefSeq94-Merged database with default parameters (Bin Jang et  
429 al., 2019). CRISPR protospacers were identified and used as the basis for host  
430 prediction of the isolated *B. thetaiotaomicron* phages and PhiSh genomes with  
431 CRISPRdb (Grissa et al., 2007) and the JGI IMG/VR Spacer Database (Paez-  
432 Espino et al., 2019) with an E-value cutoff of 1. Matches with the highest percent  
433 sequence identity are shown in **Tables S1 and S8**. Genomes of the phage  
434 isolates described in **Table S1** are uploaded to NCBI (BioProject ID  
435 PRJNA606391).

436

437 *Quantitative host range analysis*

438 Host range analysis was carried out as previously described (Porter et al.,  
439 2019). Briefly, high titer phage stocks were prepared on their “preferred host  
440 strain,” which is the strain yielding the highest titer of phages in a pre-screen of  
441 phage host range (**Table S1**). Lysates were then diluted to approximately  $10^6$   
442 PFU/mL, were added to the wells of a 96-well plate, then further diluted to  $10^5$ ,  
443  $10^4$ , and  $10^3$  PFU/mL. One microliter of each dilution was plated onto solidified  
444 top agar overlays containing wildtype *B. thetaiotaomicron*, acapsular *B.*  
445 *thetaiotaomicron*, or *B. thetaiotaomicron* expressing a single capsule (**Table S9**).  
446 After spots dried, plates were incubated anaerobically for 15-24 hours prior to  
447 counting plaques. Phage titers were normalized to the “preferred host strain.”

448

#### 449 *Infection associated pham identification*

450 We defined an infection-associated pham (IAP) as a pham that (1) was  
451 found in every phage of a given cluster ( $\alpha$ ,  $\beta$ , and  $\gamma$ ; see **Fig. 1**) that infected the  
452 *B. thetaiotaomicron* isolate in question, but (2) was not found in every phage of  
453 the same cluster. Criterion (1) is a stringent threshold. For example, if 10 different  
454 phages infected a given bacterial strain, but only 9 shared a particular pham, it  
455 would fail criterion (1). Criterion (2) was included to eliminate core genes.

456 We employed two important thresholds when identifying IAPs.

457 The first of these is an infection threshold - the normalized percentage of  
458 infectivity a given phage on a given isolate as described in the methods section  
459 'Quantitative host range analysis'. Here, a stringent threshold is 100%, which  
460 considers “infection” to be a case where the phage generates as many plaques



461 on a given *B. thetaiotaomicron* strain as it does on its preferred host strain. A  
462 permissive threshold is 1% - here a phage would have to cause 1/100th as many  
463 plaques as it did on its preferred host. The second of these is the pham identity  
464 threshold - the percentage sequence identity that two genes must share to be  
465 counted as in the same pham. This clustering is described in methods section  
466 'Annotation and comparative analyses of *B. thetaiotaomicron* infecting phages.'  
467 Here, a stringent clustering threshold is 100%, where genes sharing 100%  
468 sequence identity are grouped in the same pham. A permissive threshold would  
469 be 1%. The lower this threshold, the more disparate the sequences that are  
470 grouped together.

471 We computed our IAP identification algorithm using as thresholds each  
472 member of the product set of [1%, 5%, 10%, 50%] X [25%, 27.5%, 30%, 35%,  
473 40%, 45%, 50%, 60%, 70%, 75%, 80%, 85%, 90%, 95%, 100%] (infection  
474 threshold and pham identity threshold, respectively). Code and data are available  
475 as **Supplementary Data 2**, which provides a simple python script and the  
476 accompanying data allowing exact reproduction of the method.

477 Comparisons of cluster  $\alpha$  gp8 and homologs from metagenome-derived  
478 cluster  $\alpha$  genomes (PhiSh01-PhiSh07, HSC01) were conducted using Clustal  
479 Omega (Sievers et al., 2011) and visualized using The Interactive Tree of Life  
480 (Letunic and Bork, 2019) with default parameters.

481

482 *Transmission Electron Microscopy*

483 High titer phage lysates of representatives from each genome cluster  
484 (SJC01, DAC15, DAC20) were precipitated overnight at 4°C with gentle rocking  
485 in a solution of 1M NaCl and 10% w/v PEG8000. Phages were then precipitated  
486 via centrifugation (5500xg for 10 minutes at 4°C). Six milliliters of phage buffer  
487 was added to the pellet and broken with gentle agitation and swirling and the  
488 mixture was incubated overnight at 4°C with gentle rocking. The following day,  
489 the sample was centrifuged at 5500xg for 10 minutes at 4°C. CsCl was slowly  
490 added to the supernatant and gently dissolved via gentle swirling (final  
491 concentration 75% w/v solution). Samples were centrifuged at 26,000 RPM for  
492 24 hours at 5°C. Phage bands were extracted and stored at 4°C.

493 CsCl-banded lysates were applied directly to glow discharged Carbon  
494 Type-B 200 mesh copper grids. Samples were allowed to adsorb to the grids for  
495 3 minutes and were subsequently washed with 2 drops of ultrapure water. Three  
496 drops of uranyl acetate (1% w/v in water) were applied to the grid and the third  
497 drop was maintained on the grid for 1 minute. Filter paper was used to remove  
498 the majority of the uranyl acetate and allowed to dry at room temperature.  
499 Samples were then viewed at 120 kV on a JEOL JEM-1400 transmission  
500 electron microscope and images were collected using a Gatan Orius digital  
501 camera.

502

503 *Comparative genomic analyses between isolated B. thetaiotaomicron infecting*  
504 *phages, other isolated Bacteroides-infecting phages, and PhiSh genomes.*

505 Genomes of representatives of each genome cluster (SJC01, DAC15,

506 DAC20) were queried against the entire SRA using SearchSRA (Torres et al.,  
507 2017, Levi et al., 2018, Towns et al., 2014, Stewart et al., 2015, Buchfink et al.,  
508 2015b, Langmead and Salzberg, 2012). To determine whether these genome  
509 clusters are found in human gut metagenomes, one representative from each  
510 cluster (SJC01, DAC15, DAC20) was queried using SearchSRA using the  
511 “protein search” option. SearchSRA uses DIAMOND blastx to query 100,000  
512 reads from each of ~100,000 metagenomes publicly available in NCBI SRA  
513 against a single query amino acid sequence. The input data for each  
514 representative phage genome consisted of a single amino acid sequence  
515 consisting of every translated gene in order of appearance in the genome,  
516 separated by “XXX”. This input format was required when the analysis was  
517 conducted (July 24, 2019).

518 Data were retrieved from SearchSRA in the typical BLAST M8 format (one  
519 file per NCBI metagenome aligned to the reference phage) and parsed into BED  
520 format. BEDTools (Quinlan, 2014) coverage was used to calculate the coverage  
521 depth of each base pair along the genome. These tables were read into R 3.6.2.  
522 For each sequence run (SRR) that had  $\geq 1$  read aligning to a query amino acid  
523 sequence, SRADB (Zhu et al., 2013) was used to get the associated sample  
524 accession number (SRS) and other related sample metadata. Coverage data  
525 from sequencing runs belonging to the same sample were combined, and then  
526 average coverage depth and detection (% of bases with  $\geq 1x$  coverage) was  
527 calculated for each metagenome sample mapped.

528 For each metagenome sample mapped where the number of reads

529 sequenced was >10000, the estimated true coverage depth of the reference  
530 phage in that metagenome sample was calculated as # spots  
531 sequenced\*SearchSRA average coverage / 100000. To determine whether to  
532 assemble a given metagenome and search for a relative of a given  
533 representative phage, we filtered the list of metagenome samples based on  
534 whether the estimated real coverage was >15% and the percent of the genome  
535 detected was >20%. This list was filtered further by selecting only human gut  
536 metagenomes and by selecting samples where coverage and detection were the  
537 highest (**Table S7**).

538 Metagenomes were downloaded from NCBI SRA using parallel-fastq-  
539 dump 0.6.5 (<https://github.com/rvalieris/parallel-fastq-dump>). For each  
540 metagenome assembled, reads were trimmed using BBDuk  
541 (<https://sourceforge.net/projects/bbmap/>) 38.69 (parameters ref=adapters,phix  
542 threads=\$(( \$coreNum - 2 )) ktrim=r k=23 mink=11 hdist=2 tpe tbo qtrim=rl  
543 trimq=20 minlen=55) and assembled using MEGAHIT v1.2.9 (--mem-flag 2 -k-list  
544 21,29,39,49,59,69,79,89,99) for all samples, or -k-list  
545 21,29,39,49,59,69,79,89,99,109,119,129,139,149 if read length was >=2x250bp.

546 To identify contigs in the metagenome assemblies that might be putative  
547 relatives of the representative phages, we used DIAMOND 0.9.24(Buchfink et al.,  
548 2015a) to build a blastx database containing all individual amino acid sequences  
549 from all three representative genomes. DIAMOND blastx queries consisted all  
550 contigs from a single metagenome assembly. Individual contigs containing  
551 significant (e <= 0.001) hits for >25% of the genes from a given representative

552 phage genome were reoriented to align the 5' ends with isolated phage genomes  
553 and then included in subsequent Phamerator analysis (**Table S8**). See  
554 **Supplementary Data 1** for Genbank and fasta files of the PhiSh genomes. A  
555 tutorial for performing this analysis can be found as **Supplementary Data 3**.

556

#### 557 **Supplementary Data**

558 **Supplementary Data 1-3** are accessible at

559 <https://drive.google.com/open?id=100fpin0IDy-6iGDDe17-OEzs0yrNHIBu>.

560

#### 561 **Author contributions**

562 AJH, BDM, NTP, WVT, DAR, and RAG performed experiments/computational  
563 analyses, and analyzed the data. AJH, BDM, NTP, and WVT prepared the  
564 display items. EJN, ECM, and JLS provided key insights, tools, and reagents.  
565 AJH wrote the paper. All authors edited the manuscript prior to submission.

566

#### 567 **Acknowledgements**

568 We thank Jackson Gardner for assistance with host range analyses; Gayatri  
569 Vithanage, Lyle Shizumura, Greig Steward, and Ned Ruby for logistical  
570 assistance in phage isolation from Sand Island Wastewater Treatment Plant; and  
571 John Perrino for transmission electron microscopy expertise. This work was  
572 funded by NIH grants (GM099513 and DK096023 to ECM; DP5OD019893 to  
573 EJN, DK085025 and AT00989203 to JLS), an NIH postdoctoral NRSA  
574 (5T32AI007328 to AJH), a Stanford University School of Medicine Dean's

575 Postdoctoral Fellowship (AJH), the NIH Cellular Biotechnology Training Program  
576 (T32GM008353 to NTP), by a NCCR ARRA Award (1S10RR026780-01 to  
577 Stanford University Cell Sciences Imaging Facility), and by a National Science  
578 Foundation Graduate Research Fellowship (DGE-114747 to BDM). JLS is a  
579 Chan Zuckerberg Biohub Investigator.

580

### 581 **Declaration of Interests**

582 The authors declare no competing interests.

583

### 584 **Figure legends**

585

586 **Figure 1. Isolation and characterization of 27 *Bacteroides thetaiotaomicron*-**  
587 **infecting phages.** (A) Phages were isolated from wastewater samples collected  
588 from 3 locations in the United States and from 2 locations in Dhaka, Bangladesh.  
589 (B) Network phylogeny analysis of phage genomes, compared according to  
590 shared gene content using Phamerator, as described in **Methods**, reveals 3  
591 genomically distinct phage clusters. Colored circles indicate groups of phages  
592 according to cluster assignment, assigned by vConTACT2. (C-E) Annotated  
593 genome maps of representative members of each cluster (SJC01, DAC15, and  
594 DAC20). Genes are represented as colored boxes and conserved domains are  
595 inlaid yellow boxes within genes. If a gene has a conserved domain, it is  
596 annotated in black text. iVireons was used to predict structural genes as  
597 described in **Methods** and are annotated in red as predicted tail, major capsid, or

598 general structural (tail, MCP, +S, respectively). **(F-H)** Transmission electron  
599 micrographs of SJC01, DAC15, and DAC20 show morphological differences  
600 between these representatives of the phage clusters. See **Table S1** for additional  
601 details on the isolation locations and genotypic/phenotypic characterization of  
602 these phages.

603

604 **Figure 2. Network phylogeny of 31 *Bacteroides*-infecting phages based on**  
605 **gene content.** The genomes of 31 *Bacteroides* infecting phages were compared  
606 according to shared gene content using Splitstree and cluster assignments were  
607 made using vConTACT2, as described in **Methods**.

608

609 **Figure 3. Identification of Phage in SearchSRA (PhiSh) related to isolated**  
610 ***B. thetaiotaomicron*-infecting phages.** Representatives of each genome  
611 cluster (SJC01, DAC15, DAC20) were used to query the entire NCBI SRA using  
612 SearchSRA as described in Methods. **(A-C)** Coverage depth (log<sub>10</sub>-transformed)  
613 of SJC01, DAC15, and DAC20 genomes, respectively in the 100 best hits to  
614 SJC01 identified via Search SRA (tDNA mode). The percentage of SJC01,  
615 DAC15, and DAC20 genomes detected ( $\geq 1$  read) in each metagenome is  
616 indicated by the gray shaded column on the right of each panel. **(D)** Network  
617 phylogeny of *Bacteroides*-infecting phage genomes described in **Figure 2** and  
618 related genomes identified in publicly available metagenomes. Genomes were  
619 clustered using vConTACT2. The subset of cluster alpha phages enclosed in a  
620 rectangle is shown in greater detail at the right-hand side of the panel. Phages

621 highlighted in panel E are in bold. **(E)** Genome maps of 4 cluster  $\alpha$  phages  
622 (SJC01, ARB25, PhiSh04, and HSC01). The genes are color-coded according to  
623 pham membership and are numbered. Pairwise nucleotide identity is represented  
624 as shading between genomes. The color of this shading represents the degree of  
625 sequence similarity with violet being the most similar, progressing through the  
626 color spectrum to red, which is the least similar. Regions with no shading indicate  
627 no similarity with a BLASTN score of  $10^{-4}$  or greater.

628

629 **Figure 4. Prediction of infection-associated phams (IAPs) in *Bacteroides-***  
630 **infecting phages. (A)** Host range of *B. thetaiotaomicron* phages on strains  
631 expressing a variety of CPS (WT, wild type), a single CPS (cps1-cps8 strains) or  
632 no CPS ( $\Delta$ cps, acapsular). Tenfold serial dilutions of phage lysates ranging from  
633 approximately  $10^6$  to  $10^3$  plaque-forming units (PFU) / mL were spotted onto top  
634 agar plates containing each of the 10 bacterial strains. Plates were then  
635 incubated overnight, and plaques on each host were counted. Phage titers  
636 (PFU/ml) were calculated for each host and normalized to the titer on the  
637 “preferred host strain” for each replicate (individual replicates are shown, n=3 per  
638 phage). The phages were then clustered based on their plaquing efficiencies on  
639 the different strains (see **Methods**). Each row in the heatmap corresponds to one  
640 of three individual experimental replicates with a phage, whereas each column  
641 corresponds to one of the 10 host strains. **(B)** Changes in the total number of  
642 phams and average pham size as a function of percent amino acid identity. **(C)**  
643 Partial genome maps of 4 cluster  $\alpha$  phages (SJC01, SJC10, HNL05, and ARB25)



644 highlighting variation in gp4, gp5, and gp8. The genes are color coded according  
645 to pham membership at standard cutoffs and are numbered. Pairwise nucleotide  
646 identity is represented as shading between genomes. The color of this shading  
647 represents the degree of sequence similarity with violet being the most similar,  
648 progressing through the color spectrum to red, which is the least similar. Regions  
649 with no shading indicate no similarity with a BLASTN score of  $10^{-4}$  or greater. The  
650 red asterisk highlights gp8 from these phages. Data corresponding to these 4  
651 phages in panels A and E are in bold. **(D)** Phages containing SJC01-like gp8  
652 were compared against phages containing the alternative allele of gp8 (85% AA  
653 identity threshold) in terms of infectivity on bacterial strains highlighted in panel  
654 A. SJC01 gp8 is associated with higher infectivity of *B. thetaiotaomicron cps1*,  
655 *cps5*, *cps6*, and  $\Delta cps$  as assessed by Mann-Whitney U Test ( $p < 0.05 = *$ ,  $p < 0.01$   
656  $= **$ ,  $p < 0.001 = ***$ ). (E) gp8 from cluster  $\alpha$  isolates and the gene in the same  
657 position in cluster  $\alpha$  genomes identified from metagenomes (PhiSh01-07, and  
658 HSC01) were aligned using ClustalW and a dendrogram of these alleles was  
659 created using The Interactive Tree of Life (See **Methods**).

660

661 **Figure S1. Genome maps of *B. thetaiotaomicron*-infecting cluster  $\alpha$  phages.**

662 The genes are color-coded according to pham membership and are numbered.  
663 Pairwise nucleotide identity is represented as shading between genomes. The  
664 color of this shading represents the degree of sequence similarity with violet  
665 being the most similar, progressing through the color spectrum to red, which is

666 the least similar. Regions with no shading indicate no similarity with a BLASTN  
667 score of  $10^{-4}$  or greater.

668

669 **Figure S2. Genome maps of *B. thetaiotaomicron*-infecting cluster  $\beta$  phages.**

670 The genes are color-coded according to pham membership and are numbered.  
671 Pairwise nucleotide identity is represented as shading between genomes. The  
672 color of this shading represents the degree of sequence similarity with violet  
673 being the most similar, progressing through the color spectrum to red, which is  
674 the least similar. Regions with no shading indicate no similarity with a BLASTN  
675 score of  $10^{-4}$  or greater.

676

677 **Figure S3. Genome maps of *B. thetaiotaomicron*-infecting cluster  $\gamma$  phages.**

678 The genes are color-coded according to pham membership and are numbered.  
679 Pairwise nucleotide identity is represented as shading between genomes. The  
680 color of this shading represents the degree of sequence similarity with violet  
681 being the most similar, progressing through the color spectrum to red, which is  
682 the least similar. Regions with no shading indicate no similarity with a BLASTN  
683 score of  $10^{-4}$  or greater.

684

685 **Figure S4. Genome maps of DAC15, DAC17, and CrAss001.** The genes are  
686 color coded according to pham membership and are numbered. Pairwise  
687 nucleotide identity is represented as shading between genomes. The color of this  
688 shading represents the degree of sequence similarity with violet being the most

689 similar, progressing through the color spectrum to red, which is the least similar.  
690 Regions with no shading indicate no similarity with a BLASTN score of  $10^{-4}$  or  
691 greater.

692

## 693 **References**

- 694 ANDERSON, L. M. & BOTT, K. F. 1985. DNA packaging by the *Bacillus subtilis*  
695 defective bacteriophage PBSX. *J Virol*, 54, 773-80.
- 696 BARR, J. J., AURO, R., FURLAN, M., WHITESON, K. L., ERB, M. L.,  
697 POGLIANO, J., STOTLAND, A., WOLKOWICZ, R., CUTTING, A. S.,  
698 DORAN, K. S., SALAMON, P., YOULE, M. & ROHWER, F. 2013.  
699 Bacteriophage adhering to mucus provide a non-host-derived immunity.  
700 *Proc Natl Acad Sci U S A*, 110, 10771-6.
- 701 BENLER, S., COBIAN-GUEMES, A. G., MCNAIR, K., HUNG, S. H., LEVI, K.,  
702 EDWARDS, R. & ROHWER, F. 2018. A diversity-generating retroelement  
703 encoded by a globally ubiquitous Bacteroides phage. *Microbiome*, 6, 191.
- 704 BESEMER, J. & BORODOVSKY, M. 2005. GeneMark: web software for gene  
705 finding in prokaryotes, eukaryotes and viruses. *Nucleic Acids Res*, 33,  
706 W451-4.
- 707 BIN JANG, H., BOLDUC, B., ZABLOCKI, O., KUHN, J. H., ROUX, S.,  
708 ADRIAENSSENS, E. M., BRISTER, J. R., KROPINSKI, A. M.,  
709 KRUPOVIC, M., LAVIGNE, R., TURNER, D. & SULLIVAN, M. B. 2019.  
710 Taxonomic assignment of uncultivated prokaryotic virus genomes is  
711 enabled by gene-sharing networks. *Nat Biotechnol*, 37, 632-639.
- 712 BRUSSOW, H. & HENDRIX, R. W. 2002. Phage genomics: small is beautiful.  
713 *Cell*, 108, 13-6.
- 714 BUCHFINK, B., XIE, C. & HUSON, D. H. 2015a. Fast and sensitive protein  
715 alignment using DIAMOND. *Nat Methods*, 12, 59-60.
- 716 BUCHFINK, B., XIE, C. & HUSON, D. H. 2015b. Fast and sensitive protein  
717 alignment using DIAMOND. *Nat Methods*.
- 718 CRESAWN, S. G., BOGEL, M., DAY, N., JACOBS-SERA, D., HENDRIX, R. W.  
719 & HATFULL, G. F. 2011. Phamerator: a bioinformatic tool for comparative  
720 bacteriophage genomics. *BMC Bioinformatics*, 12, 395.
- 721 DELCHER, A. L., HARMON, D., KASIF, S., WHITE, O. & SALZBERG, S. L.  
722 1999. Improved microbial gene identification with GLIMMER. *Nucleic  
723 Acids Res*, 27, 4636-41.
- 724 DUERKOP, B. A., KLEINER, M., PAEZ-ESPINO, D., ZHU, W., BUSHNELL, B.,  
725 HASSELL, B., WINTER, S. E., KYRPIDES, N. C. & HOOPER, L. V. 2018.  
726 Murine colitis reveals a disease-associated bacteriophage community. *Nat  
727 Microbiol*.
- 728 DUTILH, B. E., CASSMAN, N., MCNAIR, K., SANCHEZ, S. E., SILVA, G. G.,  
729 BOLING, L., BARR, J. J., SPETH, D. R., SEGURITAN, V., AZIZ, R. K.,

- 730 FELTS, B., DINSDALE, E. A., MOKILI, J. L. & EDWARDS, R. A. 2014. A  
731 highly abundant bacteriophage discovered in the unknown sequences of  
732 human faecal metagenomes. *Nat Commun*, 5, 4498.
- 733 EDWARDS, R. A., MCNAIR, K., FAUST, K., RAES, J. & DUTILH, B. E. 2016.  
734 Computational approaches to predict bacteriophage-host relationships.  
735 *FEMS Microbiol Rev*, 40, 258-72.
- 736 GARNEAU, J. R., DEPARDIEU, F., FORTIER, L. C., BIKARD, D. & MONOT, M.  
737 2017. PhageTerm: a tool for fast and accurate determination of phage  
738 termini and packaging mechanism using next-generation sequencing data.  
739 *Sci Rep*, 7, 8292.
- 740 GRAZZIOTIN, A. L., KOONIN, E. V. & KRISTENSEN, D. M. 2017. Prokaryotic  
741 Virus Orthologous Groups (pVOGs): a resource for comparative genomics  
742 and protein family annotation. *Nucleic Acids Res*, 45, D491-d498.
- 743 GREGORY, A. C., ZABLOCKI, O., HOWELL, A., BOLDUC, B. & SULLIVAN, M.  
744 B. 2019. The human gut virome database. *bioRxiv*.
- 745 GRISSA, I., VERGNAUD, G. & POURCEL, C. 2007. The CRISPRdb database  
746 and tools to display CRISPRs and to generate dictionaries of spacers and  
747 repeats. *BMC Bioinformatics*, 8, 172.
- 748 GUERIN, E., SHKOPOROV, A., STOCKDALE, S. R., CLOONEY, A. G., RYAN,  
749 F. J., SUTTON, T. D. S., DRAPER, L. A., GONZALEZ-TORTUERO, E.,  
750 ROSS, R. P. & HILL, C. 2018. Biology and Taxonomy of crAss-like  
751 Bacteriophages, the Most Abundant Virus in the Human Gut. *Cell Host*  
752 *Microbe*, 24, 653-664.e6.
- 753 GUTHRIE, L., GUPTA, S., DAILY, J. & KELLY, L. 2017. Human microbiome  
754 signatures of differential colorectal cancer drug metabolism. *NPJ Biofilms*  
755 *Microbiomes*, 3, 27.
- 756 HANAUER, D. I., GRAHAM, M. J., BETANCUR, L., BOBROWNICKI, A.,  
757 CRESAWN, S. G., GARLENA, R. A., JACOBS-SERA, D., KAUFMANN,  
758 N., POPE, W. H., RUSSELL, D. A., JACOBS, W. R., JR., SIVANATHAN,  
759 V., ASAI, D. J. & HATFULL, G. F. 2017. An inclusive Research Education  
760 Community (iREC): Impact of the SEA-PHAGES program on research  
761 outcomes and student learning. *Proc Natl Acad Sci U S A*, 114, 13531-  
762 13536.
- 763 HATFULL, G. F. & HENDRIX, R. W. 2011. Bacteriophages and their genomes.  
764 *Curr Opin Virol*, 1, 298-303.
- 765 HAWKINS, S. A., LAYTON, A. C., RIPP, S., WILLIAMS, D. & SAYLER, G. S.  
766 2008. Genome sequence of the *Bacteroides fragilis* phage ATCC 51477-  
767 B1. *Virol J*, 5, 97.
- 768 HE, Q., GAO, Y., JIE, Z., YU, X., LAURSEN, J. M., XIAO, L., LI, Y., LI, L.,  
769 ZHANG, F., FENG, Q., LI, X., YU, J., LIU, C., LAN, P., YAN, T., LIU, X.,  
770 XU, X., YANG, H., WANG, J., MADSEN, L., BRIX, S., WANG, J.,  
771 KRISTIANSEN, K. & JIA, H. 2017. Two distinct metacommunities  
772 characterize the gut microbiota in Crohn's disease patients. *Gigascience*,  
773 6, 1-11.

- 774 KANG, D. D., LI, F., KIRTON, E., THOMAS, A., EGAN, R., AN, H. & WANG, Z.  
775 2019. MetaBAT 2: an adaptive binning algorithm for robust and efficient  
776 genome reconstruction from metagenome assemblies. *PeerJ*, 7, e7359.
- 777 LANGMEAD, B. & SALZBERG, S. L. 2012. Fast gapped-read alignment with  
778 Bowtie 2. *Nat Methods*
- 779 LETUNIC, I. & BORK, P. 2019. Interactive Tree Of Life (iTOL) v4: recent updates  
780 and new developments. *Nucleic Acids Res*, 47, W256-w259.
- 781 LEVI, K., RYNGE, M., EROMA, A. & EDWARDS, R. A. 2018. Searching the  
782 Sequence Read Archive using Jetstream and Wrangler. *Proceedings of*  
783 *the Practice and Experience on Advanced Research Computing*.
- 784 LIU, W., ZHANG, J., WU, C., CAI, S., HUANG, W., CHEN, J., XI, X., LIANG, Z.,  
785 HOU, Q., ZHOU, B., QIN, N. & ZHANG, H. 2016. Unique Features of  
786 Ethnic Mongolian Gut Microbiome revealed by metagenomic analysis. *Sci*  
787 *Rep*, 6, 34826.
- 788 LOWE, T. M. & EDDY, S. R. 1997. tRNAscan-SE: a program for improved  
789 detection of transfer RNA genes in genomic sequence. *Nucleic Acids Res*,  
790 25, 955-64.
- 791 MANRIQUE, P., BOLDUC, B., WALK, S. T., VAN DER OOST, J., DE VOS, W.  
792 M. & YOUNG, M. J. 2016. Healthy human gut phageome. *Proc Natl Acad*  
793 *Sci U S A*, 113, 10400-5.
- 794 MINOT, S., BRYSON, A., CHEHOUD, C., WU, G. D., LEWIS, J. D. &  
795 BUSHMAN, F. D. 2013. Rapid evolution of the human gut virome. *Proc*  
796 *Natl Acad Sci U S A*, 110, 12450-5.
- 797 MONACO, C. L., GOOTENBERG, D. B., ZHAO, G., HANDLEY, S. A.,  
798 GHEBREMICHAEL, M. S., LIM, E. S., LANKOWSKI, A., BALDRIDGE, M.  
799 T., WILEN, C. B., FLAGG, M., NORMAN, J. M., KELLER, B. C.,  
800 LUEVANO, J. M., WANG, D., BOUM, Y., MARTIN, J. N., HUNT, P. W.,  
801 BANGSBERG, D. R., SIEDNER, M. J., KWON, D. S. & VIRGIN, H. W.  
802 2016. Altered Virome and Bacterial Microbiome in Human  
803 Immunodeficiency Virus-Associated Acquired Immunodeficiency  
804 Syndrome. *Cell Host Microbe*, 19, 311-22.
- 805 OGILVIE, L. A., CAPLIN, J., DEDI, C., DISTON, D., CHEEK, E., BOWLER, L.,  
806 TAYLOR, H., EBDON, J. & JONES, B. V. 2012. Comparative  
807 (meta)genomic analysis and ecological profiling of human gut-specific  
808 bacteriophage phiB124-14. *PLoS One*, 7, e35053.
- 809 OUDE MUNNINK, B. B., CANUTI, M., DEIJS, M., DE VRIES, M., JEBBINK, M.  
810 F., REBERS, S., MOLENKAMP, R., VAN HEMERT, F. J., CHUNG, K.,  
811 COTTEN, M., SNIJDERS, F., SOL, C. J. & VAN DER HOEK, L. 2014.  
812 Unexplained diarrhoea in HIV-1 infected individuals. *BMC Infect Dis*, 14,  
813 22.
- 814 PAEZ-ESPINO, D., ROUX, S., CHEN, I. A., PALANIAPPAN, K., RATNER, A.,  
815 CHU, K., HUNTEMANN, M., REDDY, T. B. K., PONS, J. C., LLABRES,  
816 M., ELOE-FADROSH, E. A., IVANOVA, N. N. & KYRPIDES, N. C. 2019.  
817 IMG/VR v.2.0: an integrated data management and analysis system for  
818 cultivated and environmental viral genomes. *Nucleic Acids Res*, 47, D678-  
819 d686.

- 820 PORTER, N. T., HRYCKOWIAN, A. J., MERRILL, B. D., GARDNER, J. O.,  
821 SINGH, S., SONNENBURG, J. L. & MARTENS, E. C. 2019. Multiple  
822 phase-variable mechanisms, including capsular polysaccharides, modify  
823 bacteriophage susceptibility in *Bacteroides thetaiotaomicron*. *bioRxiv*.
- 824 QUINLAN, A. R. 2014. BEDTools: The Swiss-Army Tool for Genome Feature  
825 Analysis. *Curr Protoc Bioinformatics*, 47, 11.12.1-34.
- 826 REN, J., AHLGREN, N. A., LU, Y. Y., FUHRMAN, J. A. & SUN, F. 2017.  
827 VirFinder: a novel k-mer based tool for identifying viral sequences from  
828 assembled metagenomic data. *Microbiome*, 5, 69.
- 829 REYES, A., WU, M., MCNULTY, N. P., ROHWER, F. L. & GORDON, J. I. 2013.  
830 Gnotobiotic mouse model of phage-bacterial host dynamics in the human  
831 gut. *Proc Natl Acad Sci U S A*, 110, 20236-41.
- 832 ROUX, S., ENAULT, F., HURWITZ, B. L. & SULLIVAN, M. B. 2015. VirSorter:  
833 mining viral signal from microbial genomic data. *PeerJ*, 3, e985.
- 834 SEGURITAN, V., ALVES, N., JR., ARNOULT, M., RAYMOND, A., LORIMER, D.,  
835 BURGIN, A. B., JR., SALAMON, P. & SEGALL, A. M. 2012. Artificial  
836 neural networks trained to detect viral and phage structural proteins. *PLoS*  
837 *Comput Biol*, 8, e1002657.
- 838 SHKOPOROV, A. N., CLOONEY, A. G., SUTTON, T. D. S., RYAN, F. J., DALY,  
839 K. M., NOLAN, J. A., MCDONNELL, S. A., KHOKHLOVA, E. V., DRAPER,  
840 L. A., FORDE, A., GUERIN, E., VELAYUDHAN, V., ROSS, R. P. & HILL,  
841 C. 2019. The Human Gut Virome Is Highly Diverse, Stable, and Individual  
842 Specific. *Cell Host Microbe*, 26, 527-541.e5.
- 843 SHKOPOROV, A. N., KHOKHLOVA, E. V., FITZGERALD, C. B., STOCKDALE,  
844 S. R., DRAPER, L. A., ROSS, R. P. & HILL, C. 2018. PhiCrAss001  
845 represents the most abundant bacteriophage family in the human gut and  
846 infects *Bacteroides intestinalis*. *Nat Commun*, 9, 4781.
- 847 SIEVERS, F., WILM, A., DINEEN, D., GIBSON, T. J., KARPLUS, K., LI, W.,  
848 LOPEZ, R., MCWILLIAM, H., REMMERT, M., SODING, J., THOMPSON,  
849 J. D. & HIGGINS, D. G. 2011. Fast, scalable generation of high-quality  
850 protein multiple sequence alignments using Clustal Omega. *Mol Syst Biol*,  
851 7, 539.
- 852 STEWART, C. A., COCKERILL, T. M., FOSTER, I., HANCOCK, D.,  
853 MERCHANT, N., SKIDMORE, E., STANZIONE, D., TAYLOR, J.,  
854 TUECKE, S., TURNER, G., VAUGHN, M. & GAFFNEY, N. I. 2015.  
855 Jetstream: a self-provisioned, scalable science and engineering cloud  
856 environment. *Proceedings of the 2015 XSEDE Conference: Scientific*  
857 *Advancements Enabled by Enhanced Cyberinfrastructure*.
- 858 TORRES, P. J., EDWARDS, R. A. & MCNAIR, K. A. 2017. {PARTIE}: a partition  
859 engine to separate metagenomic and amplicon projects in the Sequence  
860 Read Archive. *Bioinformatics*.
- 861 TOWNS, J., COCKERILL, T., DAHAN, M., FOSTER, I., GAITHER, K.,  
862 GRIMSHAW, A., HAZLEWOOD, V., LATHROP, S., LIFKA, D.,  
863 PETERSON, G. D., ROSKIES, R., SCOTT, J. R. & WILKINS-DIEHR, N.  
864 2014. XSEDE: Accelerating Scientific Discovery. *Computing in Science*  
865 *Engineering*.



- 866 YUTIN, N., MAKAROVA, K. S., GUSSOW, A. B., KRUPOVIC, M., SEGALL, A.,  
867 EDWARDS, R. A. & KOONIN, E. V. 2018. Discovery of an expansive  
868 bacteriophage family that includes the most abundant viruses from the  
869 human gut. *Nat Microbiol*, 3, 38-46.
- 870 ZHENG, S., SHAO, S., QIAO, Z., CHEN, X., PIAO, C., YU, Y., GAO, F., ZHANG,  
871 J. & DU, J. 2017. Clinical Parameters and Gut Microbiome Changes  
872 Before and After Surgery in Thoracic Aortic Dissection in Patients with  
873 Gastrointestinal Complications. *Sci Rep*, 7, 15228.
- 874 ZHU, Y., STEPHENS, R. M., MELTZER, P. S. & DAVIS, S. R. 2013. SRadb:  
875 query and use public next-generation sequencing data from within R. *BMC*  
876 *Bioinformatics*, 14, 19.  
877









