# Human Phageprints: A high-resolution exploration of oral phages reveals globally-distributed phage families with individual-specific and temporally-stable community compositions

Gita Mahmoudabadi, Kelsey Homyk, Adam Catching, Helen Foley, Arbel Tadmor, Ana Mahmoudabadi, Allison Cheung, Rob Phillips

# Abstract

Metagenomic studies have revolutionized the study of novel phages. However these studies trade the depth of coverage for breadth. In this study we show that the targeted sequencing of a phage genomic region as small as 200-300 base pairs, can provide sufficient sequence diversity to serve as an individual-specific barcode or "Phageprint". The targeted approach reveals a high-resolution view of phage communities that is not available through metagenomic datasets. By creating instructional videos and collection kits, we enabled citizen scientists to gather ~700 oral samples spanning ~100 individuals residing in different parts of the world. In examining phage communities at 6 different oral sites, and by comparing phage communities of individuals living across the globe, we were able to study the effect of spatial separation, ranging from several millimeters to thousands of kilometers. We found that the spatial separation of just a few centimeters (the distance between two oral sites) can already result in highly distinct phage community compositions. For larger distances, spanning the phage communities of different individuals living in different parts of the world, we did not observe any correlation between spatial distance and phage community composition as individuals residing in the same city did not have any more similar phage communities than individuals living on different continents. Additionally, we found that neither genetics nor cohabitation seem to play a role in the relatedness of phage community compositions across individuals. Cohabitating siblings and even identical twins did not have phage community compositions that were any more similar than those of unrelated individuals. The primary factor contributing to phage community composition relatedness is direct contact between two habitats, as is demonstrated by the similarity between oral phage community compositions of partners. Furthermore, by exploring phage communities across the span of a month, and in some cases several years, we observed highly stable community compositions. These

24    studies consistently point to the existence of remarkably diverse and personal phage families that are

25    stable in time and apparently present in people around the world.

26

## Introduction

28         The study of bacteriophages, or viruses of bacteria, has traditionally relied on the culturing of

29    the bacterial hosts. Because the vast majority of bacteria remain unculturable, we have only recently

30    begun to recognize the overwhelming presence of phages through culture-independent techniques

31    (1, 2). These advances collectively paint phages not only as the most numerous and diverse

32    biological entities on our planet, but also as regulators of microbial ecosystems through rapid

33    infection cycles and gene transfer events (3-7). Yet, compared to their bacterial hosts, and despite

34    their proven potential to transform fields such as medicine (8-10), agriculture (11, 12), and

35    biotechnology (13-15), phages are, in general, poorly characterized (16-20).

36         Even across familiar microbial habitats such as those within the human body, the identity of

37    phages and their corresponding bacterial hosts, their community compositions, their modes of

38    transfer between habitats, their co-evolutionary history with bacterial and human hosts, their role in

39    health and disease, among other important topics remain highly unexplored. We thus chose to study

40    the human oral cavity, not only because it represents a multifaceted and medically important

41    ecosystem, but also because there are very few studies focused on oral phage communities (21-25).

42         Several intriguing studies have revealed phages as the most abundant members of the human

43    oral cavity ($10^8$ virus like particles per mL of saliva) (22, 26), with distinct communities at sites of

44    disease (27, 28), capable of augmenting the bacterial arsenal of pathogenic genes (29, 30). These

45    studies have relied on the shotgun metagenomic approach, in part because one of the defining

46    features of viral genomes is the lack of a universally conserved sequence. Given that the ribosomal

47    RNA sequence can be used as a universal marker for cellular genomes, its sequence variation is used

48  to draw conclusions about cellular evolution and taxonomic classification (31-33). This marker-

49  based approach to microbiology is additionally indispensible to microbial ecology as it allows a high

50  coverage depth of the 16S region, which in turn, enables precise and reproducible depictions of

51  bacterial community compositions (34-38).

52      Using the current sequencing platforms, the trade-off for coverage depth is typically the

53  coverage breadth (Figure 1). In comparison to the marker-based approach, shotgun metagenomics

54  provides a much greater breadth in coverage and offers several advantages. However, it suffers from

55  several disadvantages. The coverage depth is often heterogeneous and remains comparatively low in

56  these studies, a manifestation of which is that the *de novo* assembly of genomes from complex

57  environments remains a significant challenge (39), even for abundant members with relatively short

58  genome lengths (40). Moreover, the genomes assembled through shotgun metagenomics are often

59  consensus genomes or an average representation of similar genomes within an environment (41). It

60  is typical to see genomic segments with ~100x coverage that are islands in the sea of lower coverage

61  depth regions (42-44). Even across regions with high coverage depth, a notable limitation surfaces

62  when there are variants that occur with a frequency below the detection limit.

63      Due to these technical challenges, the marker-based approach, which allows orders of

64  magnitude greater coverage depth by focusing the reads on a small genomic segment, provides a

65  higher resolution view of community compositions. The targeted approach is therefore widely used

66  to complement shotgun metagenomic depictions of bacterial communities (45-47).
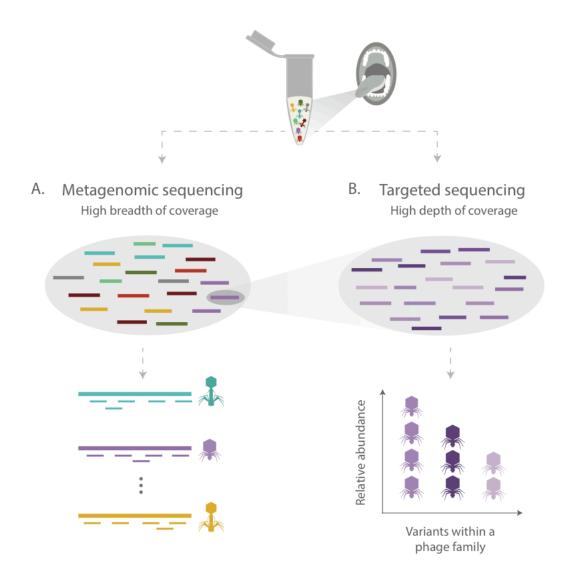
Figure 1. Comparison of A) shotgun metagenomic sequencing and B) targeted sequencing approaches. A) Shotgun metagenomic sequencing offers high breadth of coverage, spanning genomes from many different organisms, however it suffers from low depth of coverage (shown here by the incomplete assembly of phage genomes). B) Targeted sequencing approaches, which use PCR to amplify a specific genomic region, exchange breadth of coverage for depth. Targeted sequencing studies, due to their greater depth of coverage, provide much higher resolution for constructing the community composition by equating coverage depth with relative abundance of species or strains.

Due to the immense sequence diversity of phage genomes (48-50), an in-depth view of their communities through a targeted sequencing approach could provide novel insights. As such, the overarching aim of this study was to explore oral phage communities and their inter-and intra-personal diversity, their spatial patterns of distribution, as well as temporal dynamics in a large-scale and high-resolution fashion. Towards this aim, we first had to choose regions within phage genomes on which to perform targeted sequencing. Because of the vast diversity of uncultured phages, we relied on oral metagenomic datasets to identify candidate marker sequences from such phages.

We have described the methods for phage marker discovery and validation in our recent manuscript (Tadmor *et al.*, in preparation). Briefly, we arrived at seven phage markers, each corresponding to a distinct lineage of the terminase large subunit gene (TerL), which is involved in the packaging of DNA inside the phage capsid. The terminase was chosen as a target because it is considered to be a uniquely phage gene encoded by many double-stranded DNA phages (51, 52). In the absence of a genomic taxonomy for viruses, we will refer to those phages (or prophages) that share similar TerL sequences as members of a phage family. This assumption is predicated on previous studies that have shown no significant sequence similarity between TerL sequences of unrelated phages (53-55). In this study we will explore three of the seven markers, and will refer to those phages that contain these TerL markers as members of the HA, HB1 or PCA2 phage families. Refer to Materials and Methods for further information on marker discovery.

By designing primers to target these phage families, we were able to obtain at least several thousand sequences per marker, per individual (see Materials and Methods). As a direct comparison to our exploration of the same markers from hundreds of shutgun metagenomic samples (Tadmor *et al.*, in preparation), this study increases the coverage depth of a marker per subject by at least three orders of magnitude (from a few sequences to a few thousand sequences). We will demonstrate that at high sequencing depth, the phage community composition derived from members of just a single
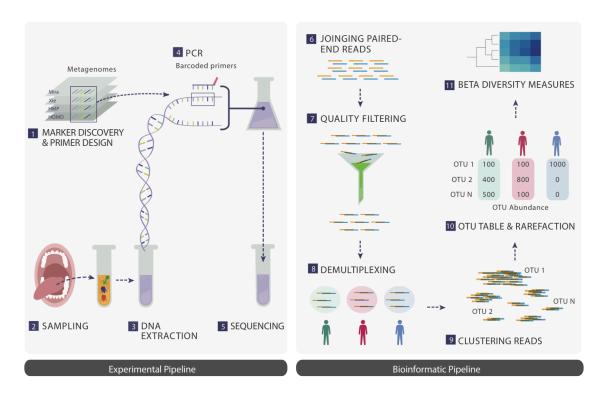
102    phage family can already serve as a fingerprint, or a "phageprint" – highly unique to a microbial

103    habitat and stable over time.

104        By creating instructional videos and collection kits, we enabled citizen scientists to gather

105    ~700 oral samples spanning ~100 individuals residing in different parts of the world. As a point of

106    comparison, one of the largest studies of the human microbiome recently reported on data from 265

107    individuals (56). By examining phage communities at 6 different oral sites, and by comparing phage

108    communities of individuals living across the globe, we were able to study the effect of spatial

109    separation, ranging from several millimeters to thousands of kilometers. We found that the spatial

110    separation of just a few centimeters (the distance between an individual's gingival sites and the hard

111    palate, for example) can already result in highly distinct phage community compositions. For larger

112    distances, spanning the phage communities of different individuals, we did not observe any

113    correlation between spatial distance and phage community composition. In other words, individuals

114    residing in the same city did not have any more similar phage communities than individuals living on

115    different continents.

116        Additionally, we found that neither genetics nor cohabitation seem to play a role in the

117    relatedness of phage community compositions across individuals. Cohabitating siblings and even

118    identical twins did not have phage communities that were any more similar than those of unrelated

119    individuals. The only factor we observed that contributes to phage community relatedness is direct

120    contact between two habitats, as is demonstrated by the similarity between oral phage community

121    compositions of partners. Furthermore, by exploring phage communities across the span of a

122    month, and in some cases several years, we observed highly stable communities. These studies

123    consistently point to the existence of remarkably diverse and personal phage families that are stable

124    in time and apparently present in individuals living in different parts of the world.

125

## Methods Summary

From a methodological standpoint, targeted sequencing of these phage markers is very similar to 16S sequencing (35, 57). Using barcoded primers, we employ PCR and next generation sequencing to attain millions of paired-end reads (Figure 2). After several quality control filters, the reads are demultiplexed based on their barcoded primer sequence and linked back to the sample and the marker from which they originated. All reads derived from the same primer sets (i.e. reads corresponding to the same TerL marker) are then pooled and clustered based on their sequence similarity into Operational Taxonomic Units or OTUs. An OTU table is constructed wherein the number of reads belonging to each OTU across each sample is denoted. Using the OTU table, we can plot the relative abundance of each OTU within a sample. We refer to this plot, which represents relative abundance profile of members within a phage family as a "community composition" plot. Finally, we will use various diversity metrics to further explore phage communities within and between individuals. This procedure is repeated for each of the three phage families. Detailed description of these protocols can be found in the Materials and Methods.

Figure 2. A schematic summary of the main experimental and bioinformatic methods: 1) Discovery of ubiquitous phage families by examining large terminase sequences that occur across different metagenomic datasets, 2) experimental sampling, 3) DNA extraction from oral biofilm samples, 4) PCR using barcoded primers followed by PCR clean-up, 5) paired-end Next Generation Sequencing, 6) joining paired-end reads to eliminate sequencing errors, 7) additional quality control steps to further eliminate errors, 8) demultiplexing of reads based on their barcode sequence and linking sequences to the sample they originate from, 9) gathering reads from all samples and clustering them based on sequence similarity (OTUs), 10) counting the number of sequences belonging to each OTU from each sample (i.e. constructing an OTU table), and rarefying the table so that each sample is represented by the same total number of sequences, and 11) performing various downstream diversity analysis (e.g. community composition plots or phageprints) using the constructed OTU table as the basis. Note that these steps are performed separately for each of the three phage families (three separate OTU tables are constructed).

9

## Results

158    **Results**

159    **An exploration of phage families reveals the presence of highly diverse, personal**

160    **phage community compositions that are stable in time.** As previously described in the

161    introduction, due to the high depth of coverage afforded through targeted sequencing, we are able

162    to explore phage sequence diversity in extraordinary detail. With bacterial 16S data, sequences are

163    generally clustered at 97% sequence similarity into operational taxonomic units (OTUs), primarily to

164    manage the large volume of data. At this threshold, each OTU is conventionally referred to as a

165    bacterial species. In fact, it is based on OTU counts that the number of bacterial species in a habitat

166    is estimated. In the absence of any convention for handling viral targeted sequencing data, we have

167    used here various sequence similarity thresholds for clustering (including 100% sequence similarity

168    threshold). We found the results to be largely robust to variations in the sequence similarity

169    threshold (see Materials and Methods).

170    Figure 3 demonstrates the HA phage community composition from a subject's tongue

171    dorsum (top surface) at two time points. The x-axis is a list of OTUs that were generated when HA

172    phage family sequences from all subjects were clustered based on sequence similarity. The y-axis

173    counts the relative abundance of the subject's HA sequences that fall into each OTU. As shown in

174    this representative figure, and across all other community composition plots we have seen for all

175    three phage families, the community composition is highly skewed towards a small number of

176    dominant OTUs (typically one or two OTUs). In addition to these OTUs, there are many other

177    OTUs with abundance values that are fairly stable in time. Generally, the dominant OTUs are not

178    the same across different individuals, and the presence of numerous other OTUs with temporally

179    stable relative abundances, gives rise to phage community compositions that are highly personal.

180    Therefore, we coined "phageprint" as shorthand to refer to a community composition plot.

181

Figure 3. HA Phage community compositions (phageprints) from subject 37 at two different time points. Samples were collected from the tongue dorsum. A) Subject 37's phageprint at $0^{th}$ time point, collected right after brushing tongue dorsal and teeth surfaces. B) Subject 37's phageprint 24 hours after the initial time point (no brushing in between time points). Each phageprint is derived from the analysis of 4000 sequences. OTUs are defined at 98% sequence similarity.

We have so far demonstrated the highly personal nature of phage communities residing in the human mouth. To better explore the temporal dynamics of these phage communities, 10 subjects

192     collected biofilm from the tongue dorsum every 24 hours for 30 days. The HB1 community

193     composition as it evolved over 30 days on subject 1's tongue dorsum is depicted in Figure 4. Here,

194     to provide a more detailed view of this community, we cluster the HB1 TerL sequences into OTUs

195     based on 100% similarity.

196



197

198          Figure 4. A 3D surface plot depicting the HB1 phage community composition as it
199          evolves over 30 days on subject 1's tongue dorsum. The x-axis contains ~10,000
200          OTUs ordered according to the depicted phylogenetic tree of the OTU sequences
201          (the phylogenetic tree is provided largely to serve as a schematic since it is hard to
202          visualize the details of this tree). Each OTU is composed of identical sequences (i.e.
203          100% sequence similarity threshold). The y-axis depicts the relative abundance of
204          each OTU, and the z-axis shows the fluctuations in relative abundance of each OTU
205          in time.

206

207        Surprisingly, over 30 days, the main features of each phage community composition is

208    preserved, though there are also interesting fluctuations that are well above the experimental error

209    and detection threshold (see SI). Figure 5 demonstrates different degrees of temporal stability and

210    phylogenetic diversity across individuals. However, a global trend is that the dominant OTU(s)

211    remain dominant over the span of 30 days in all subjects. This observation is especially interesting in

212    light of the inter-and intra-personal differences in diet and oral hygiene practices over time (Figure

213    6).



214

215        Figure 5. Depictions of HB1 phage community composition evolution in different

216        subjects over 30 days. The format of the plots is the same as that of Figure 4, and the

217        order of OTUs is based on their phylogenetic distance and identical across all plots.

218        All samples are collected from the tongue dorsum. Note that subject 2 and 4 are a

219        couple, and their phage community compositions share some main features. The

220        metadata associated with these subjects is provided in Figure 6.

221

222    To make quantitative pairwise comparisons between community compositions we employed

223    several commonly used metrics such as the Bray-Curtis and Unifrac (see Materials and Methods),

224    and in doing so, we distill the comparison of thousands of sequences from any two samples to a

225    single score. We will therefore present heatmaps of pairwise comparison scores for each phage

226    family.

227    All distance metrics explored paint similar pictures of the HB1 phage communities, depicting

228    them as highly personal and stable over time (Figure 7, Figure 8). Because phage communities in

229    different individuals have such distinct compositions, abundance-based metrics are especially

230    suitable for describing them. However, even the binary Jaccard and weighted Unifrac distance

231    metrics demonstrate a similar message. Figure 8 further demonstrates the intra-and -interpersonal

232    distances as measured through these various distance metrics. As is expected from the heatmaps

233    shown in Figure 7, the intra-personal distances are markedly different from the inter-personal, with

234    the notable exception being subject 2 and 4, who are partners.

235

236

Figure 6. Subject daily metadata during 30 consecutive days. Top panel for each subject represents the Caloric intake from fats, carbohydrates, and protein. Mean Caloric intake or MCI reports the Caloric intake averaged over 30 days (the x-axis for all plots is number of days). Pie charts demonstrate the diet over 30 days based on

242    median fat, carbohydrate and protein consumption. The second panel depicts the

243    change in Calorie intake from the previous day. The third and fourth panel

244    correspond to the number of times that the subject brushed his or her teeth and

245    tongue, respectively, during the 24 hour sampling interval. We have used asterisk to

246    denote days for which we did not receive data from subject, and to distinguish them

247    from zero values in third and fourth panel, they have been given "-.1" value. The

248    subjects are the same as those shown in the previous two figures, however two of

249    the subjects did not report dietary information so they are not included in this figure.

250

Subject IDs
01 02 04 06 07 08 09 10 11 12

Pearson

1.0
0.8
0.6
0.4
0.2
0.0

B.

Binary Jaccard    Abundance Jaccard    Bray Curtis    Unifrac (unweighted)

251

252    Figure 7. HB1 phage community temporal dynamics (previously shown graphically in
253    Figure 5) depicted here by pairwise distance metrics: A) Peason, B) Binary Jaccard,
254    Abundance Jaccard, Bray Curtis and unweighted Unifrac. The heatmap scale applies
255    to all heatmaps shown. Subjects 02 and 04 are a couple. Samples from each subject
256    are chronologically ordered.

Figure 8. Intra-and inter-personal distances between HB1 phage communities from 10 subjects, over the span of 30 days (further quantifying the heatmaps from Figure 7). Box-plots depict distances from pairwise comparisons made using the following metrics: A) Binary Jaccard, B) Abundance Jaccard, C) Bray-Curtis, D) Pearson, and E) unweighted Unifrac. The outliers defined as those outside of the 1.5 x IQR (inter-quartile range) are denoted by "+". The box-plots corresponding to the comparisons between the couple in this study are highlighted.

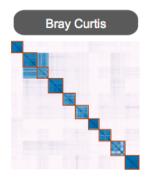**Phage community comparisons across siblings, couples, and non-related individuals residing across the globe.** Given the ubiquitous presence of the phage families across subjects residing in the U.S. we wondered whether phage families (HA and HB1, specifically) are globally distributed, and whether subjects residing in the same country would have more similar phage communities. We discovered that phage families were in fact found in individuals from various ethnicities, nationalities, and ages. Surprisingly, neither from abundance-based nor phylogenetic distance comparisons did we find an indication that people residing in the same country share more similar phage communities (Figure 9). Instead, we continued to find that individuals typically have highly unique phage communities.

Even siblings who were either living in the same household or had previously, do not have any more similar phage communities than unrelated individuals. In fact, one of the four sibling groups with uncorrelated phageprints are identical twins (Figure 9). However, 3 out of 4 couples in this study exhibited highly similar phage communities. The dissimilar couple may be due to celiac disease diagnosed in one of the partners, which is known to alter oral ecology (58). These results suggest that genetics and cohabitation do not significantly impact a person's oral phage community. The more impactful factor appears to be direct oral contact with another person. To further test these trends, larger studies encompassing a greater number of individuals and regions in the world are required.

A.



B.

285    Figure 9. HB1 phage community across 61 individuals residing across different parts
286    of the globe. Samples are obtained from the tongue dorsum. A) Pearson distance (1
287    – Pearson correlation) is shown as a heatmap. A subset of individuals residing in the
288    U.S. are either couples or siblings. Green and red boxes are drawn around samples
289    from each sibling group and couple, respectively. B) Intra- and inter-country
290    distances from pairwise comparisons made using Bray-Curtis and unweighted
291    Unifrac distance metrics. The outliers are denoted as points outside of the 1.5 x IQR
292    (inter-quartile range). Siblings and couples are excluded from this analysis.
293

294    **Different oral sites.** Thus far, all phage communities shown are those sampled from the tongue

295    dorsum. In order to examine the spatial patterns of phage communities we obtained additional oral

296    samples spanning 9 individuals and 6 oral sites (Courtesy of Bik *et al.*).  Figure 10 shows the HB1

297    phage community compositions of a subject at four oral sites where HB1 phage family was found.

298    Clearly, different oral sites in this subject have very similar HB1 phageprints. When examining all

299    HB1 positive samples, an immediately recognizable pattern is that the HB1 phage community

300    compositions of an individual are highly correlated. In stark contrast are the correlations between

301    the phage community compositions of different individuals.

302

Figure 10. HB1 phage community compositions (phageprints) across 4 different oral sites in subject 16. Each phageprint is derived from the analysis of 4000 sequences (see SI). OTUs are defined at 98% sequence similarity and OTUs with less than or equal to 0.1% relative abundance across all phageprints were filtered out (see SI).

Figure 11. Pearson correlation coefficient matrix of HB1 phage community compositions spanning 9 subjects and four oral sites. Each community composition is derived from the analysis of 4000 sequences associated with an individual and a particular oral site. OTUs are defined at 98% sequence similarity and OTUs with less than or equal to 0.1% relative abundance across all phageprints were filtered out (see SI). Phageprints are color-coded based on the individual they originate from. Community compositions that have been replicated at least twice and averaged have an asterisk next to them (see SI).

As in the case of the HB1 phage family, there is low to non-existing correlation between the HA phage community compositions of different individuals at the same oral site (Figure 12), reinforcing the notion of highly personal phage communities. However, unlike HB1, not all oral sites within the same subject are highly or even moderately correlated (see subjects 3, 12, and 17). In subject 12 for

24

322    example, the tongue dorsum has a correlation close to zero with supra-gingiva and sub-gingiva sites,

323    which are nearly perfectly correlated. Similarly, in subject 3, the hard palate and the tongue ventral

324    surface have nearly identical phage community compositions while they have a very low correlation

325    with the community at the tongue dorsum. However, unlike subject 12, the tongue dorsum in

326    subject 3 seems to be an intermediate community, having a moderate correlation with all other sites

327    that are distinct from each other. In subject 17 as well, buccal mucosa serves as the intermediate

328    community, having a moderate correlation with the disparate communities of sub-gingiva and the

329    hard palate. Phage-host network representations for HB1 (SI Figure 1) and HA (SI Figure 2) phage

330    families across this cohort demonstrates in extensive detail the cause of weak or strong correlations

331    between different oral sites.

332

Figure 12. Pearson correlation coefficient matrix of HA community compositions encompassing 11 subjects and six oral sites. Each community composition is derived from the analysis of 4000 sequences associated with an individual and a particular oral site. Samples are color-coded based on the individual they originate from. Oral sites shown are the tongue dorsum (TD), buccal mucosa (BM), supra-gingiva (SP), sub-gingiva (SB), hard palate (HP), and ventral surface of the tongue (TV). Samples whose community composition has been replicated at least twice and averaged have an asterisk next to them (see SI).

**A bioinformatic search for the bacterial hosts.** Because we aimed to study previously uncharacterized phages, the bacterial hosts for the phage families in this study have not yet been cultured or identified. However, using homology-based searches we can identify candidate host

26

346   species. For each phage family, the most abundant sequence in each OTU served as its

347   representative sequence and was used as a query for BLASTx homology search. With the exception

348   of a few sequences tagged as "putative proteins", all resulting homologs were terminase sequences

349   (SI Table 1, SI Table 3, SI Table 5). Additionally, the bacterial species with the highest chance of

350   harboring members of these three phage families were determined based on the results of the

351   BLASTx homology search (SI Table 2, SI Table 4, SI Table 6, SI Figure 3).

352   The HA phage family infects only a single genus of Firmicutes (*Streptococcus*), but appears in

353   the genomes of many different species within this genus (SI Table 4). The majority of HB1

354   homologs belonged to ReqiPepy6 phage isolated from *Rhodococcus equi* (phylum Actinobacteria).

355   Other OTU homologs were matched to ReqiPoco6, another *R. equi* phage, and six species spread

356   across two different families of the Firmicutes phylum.

357

## Discussion

359   Our method for finding ubiquitous human oral phages relied on a relatively small

360   metagenomic dataset, which contained sequences from 6 individuals residing in Spain (59). Yet, on

361   the basis of markers designed from this small dataset we were able to identify the same phage

362   families in at least 10 times as many individuals from across the globe. This finding seems to suggest

363   that there exist certain phage families that are a stable feature of the human oral microbiome. Studies

364   of phages from various natural environments (e.g. marine, soil, lakes) also report the finding of

365   phage families that are distributed across similar types of habitats despite vast geographical distances

366   and barriers that exist between these habitats (17, 60, 61). The discovery of core bacterial members

367   within the human microbiome (21, 62, 63) that are present globally (23) further support our

368   discovery of globally distributed phage families. Similar to our findings for phages, the oral bacteria

369    of individuals from the same part of the world was as different from each other as they were to

370    individuals from other parts of the world (23).

371         The ubiquitous presence of the identified phage families in individuals, together with their

372    temporal stability, seems to suggest that they likely play important roles in this environment. The

373    observed temporal stability of these phage community compositions over the span of a month is

374    supported through metagenomic studies of oral phages (64, 65) as well as 16S sequencing of

375    bacterial communities inhabiting various sites in the human body (62, 66, 67). Our study represents

376    one of the largest studies of human oral phages. As a comparison, the most recent version of the

377    Human Microbiome Project, contains samples from 265 individuals (56). However, future studies

378    are required to expand our dataset to include many more individuals and many more parts of the

379    world.

380         A particularly important aspect of our study is that it combined the advantages of

381    metagenomics with targeted sequencing to not only identify core phage families inhabiting the

382    human oral cavity, but to also characterize their communities with a resolution that is unavailable

383    through metagenomic studies of phages. This detailed view allowed us to clearly observe the highly

384    complex, and personal nature of phage community compositions. Moreover, the emergence of

385    phageprints is directly the result of the remarkable phage sequence diversity that we were able to

386    capture via targeted sequencing. For example, we observed a few hundred HB1 OTUs (defined at

387    97% sequence similarity). Even though the HB1 phage family is only one of many oral phage

388    families, it contains the same level of sequence diversity as the entire bacterial population in the

389    human mouth. This perhaps explains the need to use highly elaborate algorithms applied onto 16S

390    and metagenomic sequences from all bacterial strains to be able to identify a person based on

391    his/her microbiome (68) whereas the personal nature of phage community composition plots is a

392    clearly apparent feature of the datasets. Although it is widely known from shotgun metagenomic

393  studies that viruses are highly diverse, targeted studies employing different types of bacterial and

394  viral markers could perhaps enable more quantitative comparisons of sequence diversity across these

395  organisms.

396  Because of the great diversity of sequences associated with just one phage family, in our

397  study of about 60 individuals we found one case where unrelated individuals whose phageprints had

398  similar correlation coefficients. This may be due to experimental error or due to the course-graining

399  associated with Pearson correlation matrices. However, if we conservatively assume that each phage

400  family can provide just 50 unique patterns, then the combination of phageprints from just 6

401  independent phage families would already provide a greater number of possible patterns than the

402  size of the current human population. Identifying individuals based on characteristic patterns of

403  anonymized personal microbiomes is a legitimate ethical concern (69) and should be considered with

404  the potential uses of phageprints as well. Future studies are needed to test the long-term stability of

405  the human phageprints especially with regard to perturbations such as exposure to antibiotics. To

406  our knowledge, this is the first study that demonstrates the potential application of phage sequences

407  for human identification.



408

409  Figure 13. An estimate for the number of additional globally-distributed phage
410  families needed to achieve the number of possible phageprint patterns that surpass
411  the current human population. Assuming that phageprints from each phage family
412  can provide 50 unique patterns, there would only be 3 additional globally distributed
413  phage families needed.

## Materials and Methods

**Phage marker discovery**. In our search for candidate phage markers, we combined the advantages of metagenomic and marker-based approaches. To study previously unexplored phage families, we first used existing metagenomic datasets to identify candidate phage families, and then by targeting these families using PCR, we were able to explore them with high depth of coverage. We limited our bioinformatic and later experimental search for ubiquitous phage families to those inhabiting the human oral cavity.

We have described our method for finding candidate phage markers (Tadmor *et al.*, in preparation), but here we will provide a brief summary. Depending on the question of interest, there are potentially numerous ways by which phage marker sequences can be selected. In our method, we imposed several search criteria: 1) candidate markers should be unique to phages; 2) candidate markers should not share any significant sequence similarity so that they are more likely to represent distinct phage groups, and 3) candidate markers should be present across different oral metagenomic datasets so that they are more likely to represent core (though not necessarily abundant) members of the human oral phageome. Just as one gut phage genome was detected and assembled by an analysis that took into account its presence across multiple metagenomic datasets (40), we suspected that presence of a phage marker across multiple metagenomes could suggest that it belongs to a ubiquitous phage family. In the absence of a taxonomic convention for viral genomic data, we use the term "phage family" to refer to phages that share a given marker.

To meet the first criterion, we focused our search on the terminase large subunit (TerL), which is a powerful motor used to package DNA into the phage capsid. We have previously used terminases as phage markers to study phage-host interactions within the termite gut (70). Unlike many other viral genes such as integrases and lysins, terminases lack bacterial homologs, and thus, are considered to be unique to viruses (51, 52). In meeting the second criterion, pairwise sequence

30

438    similarity analysis was performed to exclude candidate markers that shared any significant similarity.

439    In exploring thousands of double-stranded DNA phage genomes, we have found that terminases

440    from different phages typically do not share any significant sequence similarity (55). In cases where

441    we detected sequence similarity between two terminase sequences, they belonged to highly similar

442    phages infecting the same host species (55). As such, by imposing the second criterion, we likely

443    arrived at distinct TerL lineages, with each lineage representing a distinct phage family.

444    Moreover, to meet the third criterion, we used two small oral metagenomic datasets (59, 71)

445    and chose only TerL sequences that appeared in more than one individual. Together, these criteria

446    led us to 7 non-homologous TerL lineages. We then searched for these 7 markers across larger,

447    publically available metagenomic datasets (24, 72), and found the markers to be present across many

448    individuals (Tadmor *et al.*, in preparation). Now that we had used shotgun metagenomic datasets to

449    identify several ubiquitous phage markers, we aimed to use the benefits of targeted sequencing to

450    develop a high-resolution survey of phage communities across space and time.

451    **The sample collection kit and measures against sampling contamination.** To obtain

452    samples, we developed a sample collection kit and prepared kit contents within the PCR flowhood.

453    Before and after every kit preparation session, the flowhood surfaces and pipettes were wiped using

454    sterile wipes, DNA AWAY™, and 95% ethanol. At the end of each session the surfaces were also

455    UV-sterilized (60 minutes). Each kit contains plastic tongue scrapers (Yellow CeraSpoon Safe Ear

456    Curettes, Bionix) that were first autoclaved and then UV-sterilized for 60 minutes, 1.5 mL gamma-

457    sterilized and pre-packaged collection tubes certified as pyrogen- RNase- DNA- and ATP-free

458    (VWR), each containing 200 uL sterile 1X PBS buffer (VWR), along with pre-packaged sterile gloves

459    (VWR). Each collection tube and tongue scraper pair was placed inside a sterile bag and the bags

460    were placed in another bag. The next steps were performed outside of the flowhood. Each

461    collection bag was put inside a Styrofoam box along with ice gel packs. Ice gel packs and Styrofoam

462    boxes were not reused to prevent cross contamination between individuals in case of a spill, which

463    would already be highly unlikely due to multiple layers of packaging. Upon arrival of samples,

464    collection tubes were taken out of their original bags, wiped with 95% ethanol and DNA AWAY™

465    using sterile wipes and placed into a new sterile bag. Gloves were frequently exchanged both during

466    this step and before proceeding to the next collection tube to prevent cross contamination. In

467    addition to standard lab attire such as gloves and lab coat, a facemask was worn to prevent

468    contamination during kit preparation and sample storage.

469        **Subject recruitment and sample collection.** For the bulk of our sample collection, we

470    relied heavily on citizen scientists. We made an educational video to introduce a diverse audience to

471    the fascinating world of phages, explain our study and to recruit volunteers. We also created an

472    instructional video for prospective volunteers on subject disqualifying criteria and subject rights, and

473    to provide a step-by-step demonstration of sample collection, storage, and shipment. Among other

474    exclusion criteria, subjects could not have taken antibiotics for the preceding 3 months and subjects

475    could not have active cavities or gum disease. Qualified subjects were sent a kit and were asked not

476    to brush their teeth or tongue for a minimum of 8 hours prior to sample collection to allow for a

477    substantial build up of plaque on the tongue dorsum. Put simply, subjects were instructed to 1) wear

478    gloves, 2) scrape their tongue (dorsal surface) several times using the tongue scraper, 3) deposit their

479    sample into the collection tube, 4) place the tube back into the bag, and 5) store the bag in their

480    freezer along with ice gel packs prior to an over-night shipment of their samples. They were also

481    instructed to report any sources of error that occurred at any step, and to send their samples along

482    with their signed consent form and questionnaire. Our sample collection and processing protocols

483    were approved by Caltech Institutional Review Board (IRB protocol 14-0430) and Institutional

484    Biosafety Committee (IBC protocol 13-198).

485    Nine subjects included in this study are those included in a previous study of oral microbial

486    diversity by Bik *et al.* (21). Briefly, samples were collected from individuals by a dentist who

487    examined subjects for their oral health, thereby excluding subjects with active cavities, gingivitis, or

488    periodontal disease. For each subject, samples from different oral sites were collected using sterile

489    curettes and deposited separately in 1.5 mL collection tubes containing PBS buffer. The 6 oral sites

490    sampled include plaque from tongue dorsum, tongue ventral, buccal mucosa, hard palate, supra-

491    gingiva, and sub-gingiva.

492    **Measures against contamination**. A common source of contamination in PCR originates

493    from previously amplified template sequences that enter new PCR reactions. To prevent

494    contamination this type of contamination, four physically separated workstations were developed

495    for DNA extraction (station A1), PCR preparation (station A2), PCR and gel electrophoresis

496    (station B1), and PCR cleanup (station B2). A and B specify two different buildings at Caltech while

497    1 and 2 refer to two different rooms within the same building. The flow of materials was from

498    building A to B and never the vice-versa. Every station had its own set of lab equipment, materials,

499    and storage space. Disposable lab coats (Sigma-Aldrich®) were worn and disposed of at the end of

500    every procedure to ensure that DNA was not carried between stations via clothing. Facemasks

501    (Fisher Scientific) were also worn at all times to prevent any oral or nasal droplets from entering

502    reactions. Prior to the start of every DNA extraction, lab equipment and bench tops were cleaned

503    using sterile wipes and DNA AWAY™ (Thermo Scientific), a surface decontaminant that eliminates

504    DNA and DNAses. PCR preparations and aliquoting of reagents were carried out in a PCR

505    flowhood (AirClean® Systems) equipped with a UV light and laminar airflow capabilities. Lab

506    equipment required for PCR preparation was designated to the PCR preparation flowhood. At the

507    end of every experimental session and when introducing new equipment into the flowhood, all

508    surfaces were first wiped with DNA AWAY™ solution and then exposed to UV radiation for 60

509     minutes. Prepackaged, sterile gloves were used for PCR preparation. To prevent sample-to-sample

510     contamination during DNA extraction, PCR preparation, and PCR cleanup, gloves were frequently

511     exchanged. Most importantly, 5 No Template Control (NTC) reactions accompanied every PCR

512     run. Similarly, to test the presence of contaminants in extraction reagents, for every extraction

513     experiment, 3 reactions were carried out without the addition of any sample. PCR using phage

514     primers was performed on these extraction control reactions.

515         **DNA Extraction (Station A1).** DNA extraction of human oral samples was done

516     according to the manual from MoBio PowerBiofilm® DNA Isolation Kit. The advantage to using

517     this kit for DNA extraction and purification is that it combines the use of chemical and mechanical

518     (bead-beating) treatments for an increased efficiency in biofilm disruption, lysis, and removal of

519     inhibitors such as humic acid. The final concentrations of DNA were measured using Nanodrop.

520     The concentration range of the total extracted genomic DNA was typically between 5 to 50 ng/μL.

521         **PCR preparation (Station A2) and PCR (Station B1).** Each PCR reaction contained 12.5

522     μL of PerfeCTa® qPCR SuperMix, ROX™ (Quanta Biosciences), a premix containing AccuStart™

523     Taq DNA polymerase, $MgCl_2$, dNTPs, and ROX reference dye for qPCR applications. Additionally,

524     each reaction contained 10.5 μL of RT-PCR Grade Water (Ambion®) which is free of nucleic acids

525     and nucleases, 1 μL of extracted DNA at 1 ng/μL, 0.5 μL of forward and 0.5 μL of reverse primers,

526     each at 50 ng/μL (synthesized by IDT). A higher than recommended primer concentration was used

527     because the phage primers used are 32-64 fold degenerate. The thermocycling protocol was made

528     according to PerfeCTa qPCR SuperMix recommendations: 1) a 10-minute activation of AccuStart™

529     Taq DNA polymerase at 95°C, 2) 10 seconds of DNA denaturation at 95°C, 3) 20 seconds of

530     annealing at 60°C, and 4) 30 seconds of extension at 72°C, 40 cycles repeating steps 2 to 4, followed

531     by 5 minutes of final extension at 72°C.

532       **Gel electrophoresis (Station B1) and PCR cleanup (Station B2).** Phage PCR products

533       were visualized using 2% agarose in TAE buffer. After gels were cast, 5 μL of each PCR product

534       was mixed with 1 μL of 6X loading dye and loaded into a well. 5 μL of 100 base-pair ladder was

535       used, and the gel electrophoresis instrument was set to run for 30 minutes at 100V. Phage PCR

536       positive hits were purified using the QIAquick PCR Purification Kit (QIAGEN). 20μL of PCR

537       products were used and purified according to the QIAquick PCR Purification manual.

538       **Illumina sequencing.** Upon PCR cleanup, double stranded DNA concentration in each

539       sample was measured using Qubit instrument. Qubit measurements were performed in Building C

540       due to practical considerations rather than a necessary treatment for preventing contamination.

541       Samples were combined into one reaction (~2 μg dsDNA) and submitted to GENEWIZ, Inc for

542       library preparation and MiSeq 2x300bp Paired-End sequencing.

543       **DNA barcodes for multiplexed sequencing.** To enable multiplexing, unique DNA

544       barcodes (Table 1) were appended onto the forward primer sequences (Table 3) used to amplify

545       each phage marker. These barcoded primer sequences were synthesized by IDT. Using this scheme,

546       ~100 samples were submitted per MiSeq sequencing run (Table 1) and by matching the barcode

547       sequence to the sample ID, information about who and where the sample came from was accessible.

548       More specifically, Hamady error-correcting 8-letter barcodes (73) were used. Hamady DNA

549       barcodes are an example of Hamming code wherein the addition of parity bits allow for detection

550       and correction of errors within the barcode sequence. In the case of Hamady barcodes, up to 2

551       errors in the barcode sequence can be detected and one error can be corrected.

552       **Quality control steps to eliminate sequencing errors.** We used Illumina MiSeq's

553       2x300bp paired-end configuration (GENEWIZ, Inc). Each sequencing run produced about 20-25

554       Million paired-end reads. Paired-end reads were joined using *join_paired_ends.py* script from QIIME

(Quantitative Insights Into Microbial Ecology) package, and unless noted otherwise scripts used in this chapter are part of QIIME (74). When a base is confirmed by both reads, higher Phred score is increased by up to 3 points. If paired reads had any mismatches across their overlapping bases, the paired reads was eliminated from any further analysis (QC step #1). For markers HB1, PCA2, and HA the overlap between the paired reads entirely covers the marker sequence, hence eliminating many sequencing errors.

Upon joining reads and eliminating those with mismatches in the region of overlap *seqQualityFilters.py*, an in-house script, was used to preform QC step #2: taking joined reads from QC step #1, and eliminating any sequences that have one or more bases marked by a Phred score below 30. Excluded from QC step #2 were the first two bases in the beginning and end of each sequence, which for majority of reads have much lower quality scores.

Using *seqQualityFilters.py,* sequences were placed in 3 different bins according to their primer sequences, and any sequence that did not have the correct barcode length, or the correct primer sequences at the expected positions, was eliminated (QC step #3). Additionally, nearly half of remaining sequences had to be reverse complemented so that all sequences were oriented in the 5' to 3' direction. Using the same script, primer and barcode sequences were removed, and barcode sequences were written to a separate file (to be used as input to *split_libraries_fastq.py*). At this point sequences that did not have the correct length were filtered out (QC step #3). Sequences were demultiplexd using *split_libraries_fastq.py* and reads with errors in the barcode sequence were eliminated (QC step #4).

**Phage community composition plots ("Phageprints").** After demultiplexing quality-controlled reads, sequences were clustered according to a specified sequence similarity threshold using UCLUST *de novo* clustering algorithm (75) used in *pick_otus.py* script. Using *make_otu_table.py*, OTU tables were generated. An OTU table summarizes counts of sequences assigned to each OTU

579    across each sample. We refer to this per-sample sequence count as the OTU size. As long as an

580    OTU of size 1 or greater exists in at least one sample, it is included in the OTU table. In this way,

581    the counts of OTUs for samples containing the same marker remains the same, though their size

582    could vary widely across different samples. Later we will demonstrate the effects of noise filters

583    applied to the OTU table. The relative abundance of each OTU within each sample was calculated

584    via *processOtuTable.py,* another in-house script. In plotting the relative OTU abundance values for

585    different samples, we arrived at complex, individual-specific patterns. We dubbed these phage

586    community composition plots as "phageprints".

587        **Metrics for quantitative comparison of phageprints.** The first metric explored is binary

588    Jaccard distance, which is equal to one minus the ratio of the intersection to the union of two

589    samples' OTUs: $1 - \frac{|A \cap B|}{|A|+|B|-|A \cap B|}$. Here, $A$ and $B$ represent the OTUs that are present in sample 1

590    and 2, respectively. This is a binary method of comparing samples simply based on the

591    presence/absence of the OTUs. In addition to the Pearson distance (1- Pearson correlation), we

592    chose two other abundance-based distance metrics, namely abundance-weighted Jaccard and Bray-

593    Curtis. Abundance-weighted Jaccard, which is equal to $1 - \frac{UV}{U+V-UV}$ (76), is similar to Jaccard but

594    here $U$ and $V$ represent the sum of relative abundances of OTUs shared between samples 1 and 2,

595    respectively. Bray-Curtis dissimilarity (77) is defined as $\frac{\sum |x_{ik}-x_{jk}|}{\sum x_{ik}+x_{jk}}$, where $x_{ik}$ and $x_{jk}$ correspond to

596    the relative abundance of OTU $k$ in samples $i$ and $j$.

597        Lastly, we explored unweighted Unifrac, a phylogenetic distance metric (78). The Unifrac

598    algorithm operates on a phylogenetic tree containing sequences from all samples. It proceeds to

599    create pairwise comparisons between samples by identifying the branch lengths that are shared

600    between two samples, as well as the branch lengths that are unique to each sample. The Unifrac

601    distance is then defined as the unshared branch lengths divided by the total branch lengths, where

602      total branch lengths is the sum of shared and unshared branch lengths. If two samples are identical,

603      the fraction of the tree's branch lengths that is unique to one sample or the other will be zero, and

604      thus, the Unifrac distance will be zero.

605          **Examining the effect of OTU sequence similarity threshold.** In analyzing 16S

606      sequences, clusters or Operational Taxonomic Units (OTUs) are conventionally defined at 97%

607      sequence similarity threshold. To examine the effect of sequence similarity threshold for phage OTU

608      formation, we tested OTU sequence similarity thresholds of 98%, 97%, 95%, 90%, and 80%. Figure

609      14 is a matrix of Pearson correlation coefficients calculated during the pairwise comparison of HB1

610      community compositions using different sequence similarity thresholds for defining OTUs. Very

611      similar Pearson correlation matrices are obtained as the sequence similarity threshold is lowered

612      from 98% to 80%. However, because the number of cluster is reduced as we reduce the sequence

613      similarity threshold, with lower sequence similarity thresholds, the chance that individual-specific

614      variations are lumped into the same cluster is increased. If noise-induced sequence variations are

615      effectively accounted for, higher sequence similarity thresholds for defining OTUs can enable a

616      more accurate and detailed depiction of a person's phage community composition. For this reason,

617      we used a sequence similarity threshold of 98% for the study of different oral sites, and later we used

618      a 100% sequence similarity threshold for the temporal and the global study.

| A | | 3.3 | 3.5 | 6.3 | 6.5 | 6.6 | 10.1 | 10.3 | 12.1 | 12.2 | 12.3 | 12.6 | 16.1 | 16.3 | 17.2 | 17.3 | 17.5 | 17.6 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 3.3 | 1.00 | 0.95 | 0.12 | 0.12 | 0.12 | 0.02 | 0.04 | 0.12 | 0.11 | 0.11 | 0.12 | 0.05 | 0.02 | 0.00 | 0.02 | 0.00 | 0.00 |
| | 3.5 | 0.95 | 1.00 | 0.08 | 0.08 | 0.08 | 0.01 | 0.03 | 0.10 | 0.10 | 0.10 | 0.11 | 0.04 | 0.01 | 0.07 | 0.08 | 0.07 | 0.07 |
| | 6.3 | 0.12 | 0.08 | 1.00 | 1.00 | 1.00 | 0.02 | 0.05 | 0.12 | 0.12 | 0.14 | 0.16 | 0.18 | 0.06 | 0.00 | 0.04 | 0.01 | 0.00 |
| | 6.5 | 0.12 | 0.08 | 1.00 | 1.00 | 1.00 | 0.02 | 0.06 | 0.13 | 0.14 | 0.15 | 0.17 | 0.18 | 0.06 | 0.00 | 0.04 | 0.01 | 0.00 |
| | 6.6 | 0.12 | 0.08 | 1.00 | 1.00 | 1.00 | 0.02 | 0.05 | 0.12 | 0.12 | 0.14 | 0.16 | 0.19 | 0.06 | 0.00 | 0.04 | 0.01 | 0.00 |
| | 10.1 | 0.02 | 0.01 | 0.02 | 0.02 | 0.02 | 1.00 | 0.47 | 0.02 | 0.02 | 0.02 | 0.02 | 0.00 | -0.01 | -0.01 | 0.00 | 0.00 | -0.01 |
| | 10.3 | 0.04 | 0.03 | 0.05 | 0.06 | 0.05 | 0.47 | 1.00 | 0.04 | 0.04 | 0.04 | 0.04 | 0.01 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| | 12.1 | 0.12 | 0.10 | 0.12 | 0.13 | 0.12 | 0.02 | 0.04 | 1.00 | 1.00 | 1.00 | 0.95 | 0.06 | 0.04 | 0.00 | 0.03 | 0.01 | 0.00 |
| | 12.2 | 0.11 | 0.10 | 0.12 | 0.14 | 0.12 | 0.02 | 0.04 | 1.00 | 1.00 | 1.00 | 0.95 | 0.06 | 0.04 | 0.00 | 0.03 | 0.01 | 0.00 |
| | 12.3 | 0.11 | 0.10 | 0.14 | 0.15 | 0.14 | 0.02 | 0.04 | 1.00 | 1.00 | 1.00 | 0.95 | 0.06 | 0.04 | 0.01 | 0.04 | 0.02 | 0.01 |
| | 12.6 | 0.12 | 0.11 | 0.16 | 0.17 | 0.16 | 0.02 | 0.04 | 0.95 | 0.95 | 0.95 | 1.00 | 0.06 | 0.03 | 0.10 | 0.12 | 0.11 | 0.10 |
| | 16.1 | 0.05 | 0.04 | 0.18 | 0.18 | 0.19 | 0.00 | 0.01 | 0.06 | 0.06 | 0.06 | 0.06 | 1.00 | 0.26 | 0.00 | 0.01 | 0.00 | 0.00 |
| | 16.3 | 0.02 | 0.01 | 0.06 | 0.06 | 0.06 | -0.01 | 0.00 | 0.04 | 0.04 | 0.04 | 0.03 | 0.26 | 1.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| | 17.2 | 0.00 | 0.07 | 0.00 | 0.00 | 0.00 | -0.01 | 0.00 | 0.00 | 0.00 | 0.01 | 0.10 | 0.00 | 0.00 | 1.00 | 0.92 | 1.00 | 1.00 |
| | 17.3 | 0.02 | 0.08 | 0.04 | 0.04 | 0.04 | 0.00 | 0.00 | 0.03 | 0.03 | 0.04 | 0.12 | 0.01 | 0.00 | 0.92 | 1.00 | 0.95 | 0.92 |
| | 17.5 | 0.00 | 0.07 | 0.01 | 0.01 | 0.01 | 0.00 | 0.00 | 0.01 | 0.01 | 0.02 | 0.11 | 0.00 | 0.00 | 1.00 | 0.95 | 1.00 | 1.00 |
| | 17.6 | 0.00 | 0.07 | 0.00 | 0.00 | 0.00 | -0.01 | 0.00 | 0.00 | 0.00 | 0.01 | 0.10 | 0.00 | 0.00 | 1.00 | 0.92 | 1.00 | 1.00 |



Heatmap scale

| -0.01 | 0.00 | 0.10 | 0.20 | 0.30 | 0.40 | 0.50 | 0.60 | 0.70 | 0.80 | 0.90 | 1.00 |

Figure 14. Pairwise Pearson correlation coefficient values calculated for HB1 phage community compositions as a function of A) 98%, B) 97%, C) 95%, D) 90%, and E) 80% sequence similarity thresholds for OTU formation. Sample IDs can be decoded as before: subject ID precedes oral site ID. Oral sites 1-6 correspond to tongue dorsum, hard palate, buccal mucosa, ventral tongue, supra-gingiva, and sub-gingiva respectively (e.g. 3.3 corresponds to subject 3 community composition derived from the buccal mucosa, and 3.5 is subject 3 supra-gingiva community composition). The number of OTUs generated at 98%, 97%, 95%, 90%, and 80% sequence similarity thresholds are 210, 181, 172, 170, and 80, respectively.

**Detecting experimental noise.** How reproducible is a phage community composition plot? Figure 15 summarizes the sources of noise from all experimental processes performed during

632　this study. First, it's important to capture sampling variation. How consistently can we capture a

633　phage community from an individual's oral site given that we are sampling different parts of the

634　biofilm each time? Another factor that could contribute to sampling variation are the personal

635　differences in the rate of biofilm mass accumulation on the tongue dorsum. Secondly, we need to

636　ask whether processes of lysis and DNA extraction allow for the availability of the same template

637　DNA sequences in the same relative abundances across different extraction runs.

638　　　　　Third, we need to evaluate the OTU abundance variations that could result in PCR due to

639　both errors as well as other stochastic events. For example, it's possible that very rare template

640　sequences are left out of the initial cycles of PCR and their relative abundance at the end of PCR is

641　lower than their relative abundance prior to PCR. In this hypothetical scenario PCR could serve as a

642　biased amplifier. PCR purification is similar to extraction and sampling in that it does not introduce

643　sequence errors; however it is unlike these processes because after PCR billions of template copies

644　are created and it's unlikely that the loss of a fraction of templates during PCR purification will

645　dramatically change OTU relative abundances. Finally, Illumina MiSeq sequencing is another error-

646　prone process not only at the level of base-calling, but at the level of bridge amplification which like

647　PCR could introduce errors that propagate exponentially. Refer to Figure 15 for a summary of

648　processes that could result in irreproducible OTUs or variation in OTU relative abundances.

Identifying sources of noise

**A. Sampling**

OTU abundance variation in sampling the same oral site

**B. DNA extraction**

OTU abundance variation due to lysis and DNA purification processes

**C. PCR**

- OTU abundance variation due to PCR bias
- PCR errors

**D. Sequencing**

- Base-call errors
- Amplification errors
- OTU abundance variation due to coverage difference across different samples

649

Figure 15. Sources of error and variation in experimental processes used in this study. A) Sampling of the same oral site in the same individual could result in collection of different microbial communities, which could introduce new OTUs or change relative abundance of existing OTUs. B) DNA extraction is not 100% efficient and the fraction of DNA extracted from an environment could serve as a source of variation across different samples. C) PCR introduces errors that could present themselves as novel OTUs or cause variation in abundance of genuine OTUs. D) Sequencing also introduces errors both at the level of base-calling and bridge amplification.

To quantify how reproducible a given phage community composition is, we obtained 3 different samples from subject 37 tongue dorsum. We then performed DNA extraction and PCR separately on each sample and sent samples for sequencing (sequencing run #2). The logic behind this experiment was to capture a lumped measure of noise arising from various processes depicted in Figure 15. After performing quality control steps 1-4, demultiplexing reads based on their barcode

41

665     sequences, clustering reads based on 98% sequence similarity threshold for OTU formation,

666     rarefying the OTU table to 4000 reads per sample, and calculating the relative abundances of OTUs,

667     we measured the standard deviation in the relative abundance of each OTU across these three

668     samples (Figure 16). Remarkably, relative abundance values across these three samples were highly

669     consistent, with the majority of OTUs having standard deviations below 0.2% and the maximum

670     standard deviation observed was less than 0.7% relative abundance.



**OTU relative abundance standard deviation across 3 samples from subject 37 tongue dorsum (HB1 marker)**

671

672     Figure 16. Standard deviations of OTU relative abundances calculated for all
673     experimental processes. Three data points per OTU are used for standard deviation
674     calculations. These three data points correspond to measurements of OTU relative
675     abundances obtained for three different samples obtained from subject 37 tongue
676     dorsum (HB1 marker) which underwent separate sampling, DNA extraction, PCR
677     and PCR cleanup procedures. The maximum standard deviation observed is less than
678     0.007 relative abundance, and majority are close to 0.

679

680     **Identifying non-reproducible OTUs.** To identify OTUs that were non-reproducible

681     across the three samples from subject 37's tongue dorsum (HB1 marker), we flagged OTUs that had

682     appeared in only one or two samples out of three. We then plotted the histogram of non-

683     reproducible OTUs as a function of their relative abundance (for those OTUs appearing in 2 out 3

684     samples, the higher relative abundance value was used). The thresholds defining each bin, $b$, were

685    selected to be the following: $0 > b_1 \geq 0.00025$ (OTU of size 1 sequence since the total number of

686    sequences per sample is 4000), $0.00025 > b_2 \geq 0.0005$ (2 sequences), $0.0005 > b_3 \geq 0.00075$ (3

687    sequences), $0.00075 > b_4 \geq 0.001$ (4 sequences), and $0.001 < b_5$ (5 or more sequences).

688    Figure 17 demonstrates the number of non-reproducible OTUs drops as a function of OTU

689    relative abundance, and all OTUs with more than 4 sequences (0.001 relative abundance) are

690    reproducible. To conclude, we arrived at 0.001 relative abundance as the detection threshold for

691    OTUs.

692



693    Figure 17. Number of non-reproducible OTUs across three samples obtained from
694    subject 37 tongue dorsum (HB1 marker), presented as a function of OTU relative
695    abundance. A total of 30 OTUs appear in one or two samples out of three, and
696    therefore are considered non-reproducible. 21 out of 30 OTUs are defined by a
697    single sequence which translates into 0.00025 relative abundance since samples are
698    rarefied to 4000 sequences. The number of non-reproducible OTUs drops as a
699    function of OTU relative abundance, and all OTUs with more than 4 sequences
700    (0.001 relative abundance) are reproducible across three samples.

701

702    In addition to capturing a lumped sum of noise across all experimental processes for subject

703    37 tongue dorsum sample (Figure 16,Figure 17), for samples from subjects 3, 6, 10, 16, and 17, we

704    performed a second set of PCR on previously extracted DNA samples, and submitted those samples

705    for sequencing (Figure 18). In addition to these replicates, we acquired new samples from the tongue

706    dorsum for subjects 31, 35, 37, and 38, and submitted these samples for the second sequencing run.

707    In obtaining replicate phageprints, we were able to demonstrate that with proper quality filtration

708    steps phageprints are highly reproducible even when they are generated from two separate PCR and

709    sequencing steps (Figure 18).



710

Figure 18. Panel A is the Pearson correlation matrix of all HB1 phageprints. Each
phageprint is derived from the analysis of 4000 sequences associated with an
individual and a particular oral site. OTUs are defined at 98% sequence similarity and
OTUs with less than or equal to 0.1% relative abundance across all phageprints were
filtered out. Phageprints are color-coded based on the individual they originate from.
Oral sites shown to be positive for the HB1 marker are the tongue dorsum (TD),
buccal mucosa (BM), supra-gingiva (SP), and sub-gingiva (SB). Phageprints that were
acquired from sequencing run #1, are those marked as replicate #1. Panel B shows
that to confirm reproducibility of phageprints, a second set of PCR was performed
on previously extracted DNA from all samples included in sequencing run #1 and

721      those PCR products were included in sequencing run #2. Phageprints derived from

722      the second sequencing run are marked as replicate #2.

723

724      **Identifying phage marker homologs.** The most abundant sequence from each OTU was

725  retrieved using *pick_rep_set.py* to serve as a representative sequence. BLASTx function was used to

726  detect the closest homolog to each OTU's representative sequence from within the NCBI's non-

727  redundant protein database. HB1 representative sequences were aligned using Geneious (79), using a

728  gap open penalty of 30 and gap extension penalty of 15 and a 65% similarity cost matrix. No gaps

729  were introduced. The alignment is shown in SI Figure 4.

730      **Phage-host networks.** OTU tables were input to *createNetwork.py,* an in-house script that

731  creates node and edge tables. The nodes represent samples and phage OTUs, and a directed edge

732  connects samples to the OTUs that they host. The weight of this connection is based on the relative

733  abundance of the OTU in that sample. Gephi software (80) was used to visualize the resulting

734  networks, and to obtain the degree distribution.

735

736

737

738

739

740

741

742

# Supplementary Information

**Subject node ID and color-code**

| 17 | | | 16 | | 12 | | | 6 | | | 3 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| BM | SP | SB | BM | TD | SB | TD | BM | SB | SP | BM | BM | SP |

Subject nodes

**TD**: Tongue Dorsum (top surface)
**BM**: Buccal Mucosa (cheek lining)
**SP**: Supra-gingiva (above gum surface)
**SB**: Sub-gingiva (below gum surface)

745     SI Figure 1. HB1 phage family network. Purple nodes are the OTU nodes and all
746     other nodes represent samples. Sample nodes and edges are color-coded based on
747     the individual they originate from. The oral site associated with each sample is
748     abbreviated next to the sample's node. Each edge connects an OTU a sample it
749     exists in, and the edge weight is proportional to the relative abundance of the OTU
750     in that sample. Node IDs are displayed. For OTU nodes, the node ID is the OTU
751     ID which can be matched to IDs in SI Table 1 for identifying taxonomic
752     information regarding each OTU. For sample nodes, the nodes IDs are simply the
753     subjects' IDs.

754

755

SI Figure 2. HA phage-host network. Purple nodes are the OTU nodes and all other nodes represent samples. Sample nodes and edges are color-coded based on the individual they originate from. Subject node color code, ID, and the oral sites are displayed above sample nodes. Each edge connects an OTU a sample it exists in, and the edge weight is proportional to the relative abundance of the OTU in that sample. Node IDs are displayed. For OTU nodes, the node ID is the OTU ID which can be

763      matched to IDs in SI Table 3 for identifying taxonomic information regarding each

764      OTU. For sample nodes, the nodes IDs are simply the subjects' IDs.

765

766      SI Table 1. Closest homolog to each OTU's representative sequence (HB1 phage

767      family). Each OTU's representative sequence was used as a query for NCBI's

768      BLASTx homology search against the non-redundant protein database. The table

769      summarizes the E-value and the percent amino acid identity across the query

770      sequence and the closest homolog, as well as the closest homolog's name, sequence

771      ID, and taxon ID. The taxon ID is color coded, and the taxonomic classification

772      corresponding to each taxon ID can be retrieved from the following table. Note with

773      the exception of a few "putative uncharacterized" homolog names that most are

774      identified as terminases or TerLs (terminase large subunits).

775

| Query Sequence ID (OTU ID) | Percent Identity | E value | Closest Homolog | Closest Homolog Sequence ID | Closest Homolog Taxon ID |
|---|---|---|---|---|---|
| 0 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 1 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 10 | 67.9 | 3.00E-30 | terminase [[Clostridium] symbiosum] | gi\|489596073\|ref\|WP_003500516.1\| | 1512 |
| 103 | 69.14 | 5.00E-30 | terminase [[Clostridium] scindens] | gi\|639772655\|ref\|WP_024738760.1\| | 29347 |
| 104 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 106 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 109 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 11 | 72.84 | 1.00E-34 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 112 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 117 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 118 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 12 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 122 | 70.37 | 2.00E-19 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 123 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 128 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 13 | 70.37 | 6.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 132 | 70.37 | 4.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 134 | 67.9 | 3.00E-30 | terminase [[Clostridium] symbiosum] | gi\|489596073\|ref\|WP_003500516.1\| | 1512 |
| 14 | 67.9 | 3.00E-30 | terminase [[Clostridium] symbiosum] | gi\|489596073\|ref\|WP_003500516.1\| | 1512 |
| 140 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 142 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 149 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 15 | 71.6 | 5.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 161 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 162 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 164 | 72.84 | 2.00E-32 | TerL [Rhodococcus phage ReqiPoco6] | gi\|593774729\|ref\|YP_009012597.1\| | 691964 |
| 165 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 166 | 71.6 | 5.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 17 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 170 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 176 | 71.6 | 3.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 178 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 18 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 183 | 66.67 | 7.00E-30 | terminase [[Clostridium] symbiosum] | gi\|489596073\|ref\|WP_003500516.1\| | 1512 |
| 184 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 189 | 71.6 | 4.00E-33 | TerL [Rhodococcus phage ReqiPoco6] | gi\|593774729\|ref\|YP_009012597.1\| | 691964 |
| 19 | 72.84 | 1.00E-34 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 195 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 196 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 197 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 2 | 66.67 | 7.00E-30 | terminase [[Clostridium] symbiosum] | gi\|489596073\|ref\|WP_003500516.1\| | 1512 |
| 202 | 67.9 | 3.00E-30 | terminase [[Clostridium] symbiosum] | gi\|489596073\|ref\|WP_003500516.1\| | 1512 |
| 203 | 65.43 | 2.00E-28 | TerL [Rhodococcus phage ReqiPoco6] | gi\|593774729\|ref\|YP_009012597.1\| | 691964 |
| 206 | 70.37 | 2.00E-32 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 207 | 71.6 | 5.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 208 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 21 | 67.9 | 3.00E-30 | terminase [[Clostridium] symbiosum] | gi\|489596073\|ref\|WP_003500516.1\| | 1512 |
| 210 | 71.6 | 2.00E-31 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 213 | 69.14 | 2.00E-29 | TerL [Rhodococcus phage ReqiPoco6] | gi\|593774729\|ref\|YP_009012597.1\| | 691964 |
| 216 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 218 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 22 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 220 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 221 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 222 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 225 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 227 | 66.67 | 7.00E-30 | terminase [[Clostridium] symbiosum] | gi\|489596073\|ref\|WP_003500516.1\| | 1512 |
| 229 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 231 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 233 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |

776

| | | | | | |
|---|---|---|---|---|---|
| 236 | 70.37 | 2.00E-32 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 237 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 238 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 24 | 71.6 | 5.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 241 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 242 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 245 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 249 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 25 | 71.6 | 3.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 250 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 255 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 256 | 72.84 | 4.00E-34 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 257 | 66.67 | 7.00E-30 | terminase [[Clostridium] symbiosum] | gi\|489596073\|ref\|WP_003500516.1\| | 1512 |
| 259 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 26 | 71.6 | 5.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 260 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 261 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 262 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 264 | 71.6 | 3.00E-32 | TerL [Rhodococcus phage ReqiPoco6] | gi\|593774729\|ref\|YP_009012597.1\| | 691964 |
| 265 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 27 | 62.96 | 2.00E-26 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 272 | 70.37 | 3.00E-32 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 273 | 67.9 | 2.00E-30 | TerL [Rhodococcus phage ReqiPoco6] | gi\|593774729\|ref\|YP_009012597.1\| | 691964 |
| 274 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 29 | 74.07 | 2.00E-35 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 3 | 67.9 | 3.00E-29 | terminase [Clostridiales bacterium VE202-03] | gi\|639707411\|ref\|WP_024723669.1\| | 1232439 |
| 30 | 64.2 | 2.00E-28 | putative uncharacterized protein [Ruminococcus sp. CAG:17] | gi\|547240587\|ref\|WP_021976510.1\| | 1262951 |
| 32 | 67.9 | 3.00E-30 | terminase [[Clostridium] symbiosum] | gi\|489596073\|ref\|WP_003500516.1\| | 1512 |
| 33 | 67.9 | 3.00E-31 | TerL [Rhodococcus phage ReqiPoco6] | gi\|593774729\|ref\|YP_009012597.1\| | 691964 |
| 34 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 36 | 76.54 | 2.00E-35 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 37 | 67.9 | 3.00E-30 | terminase [[Clostridium] symbiosum] | gi\|489596073\|ref\|WP_003500516.1\| | 1512 |
| 38 | 66.67 | 2.00E-29 | TerL [Rhodococcus phage ReqiPoco6] | gi\|593774729\|ref\|YP_009012597.1\| | 691964 |
| 4 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 40 | 66.67 | 3.00E-32 | putative uncharacterized protein [Ruminococcus sp. CAG:17] | gi\|547240587\|ref\|WP_021976510.1\| | 1262951 |
| 42 | 64.2 | 3.00E-26 | terminase [[Clostridium] symbiosum] | gi\|489596073\|ref\|WP_003500516.1\| | 1512 |
| 44 | 66.67 | 5.00E-19 | terminase [[Ruminococcus] torques] | gi\|490985259\|ref\|WP_004846995.1\| | 33039 |
| 46 | 67.9 | 3.00E-30 | terminase [[Clostridium] symbiosum] | gi\|489596073\|ref\|WP_003500516.1\| | 1512 |
| 48 | 71.6 | 9.00E-34 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 5 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 50 | 86.42 | 3.00E-44 | TerL [Rhodococcus phage ReqiPoco6] | gi\|593774729\|ref\|YP_009012597.1\| | 691964 |
| 51 | 67.9 | 3.00E-30 | terminase [[Clostridium] symbiosum] | gi\|489596073\|ref\|WP_003500516.1\| | 1512 |
| 54 | 65.43 | 6.00E-31 | terminase [[Clostridium] symbiosum] | gi\|489596073\|ref\|WP_003500516.1\| | 1512 |
| 59 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 6 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 60 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 62 | 69.14 | 2.00E-30 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 67 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 68 | 74.07 | 1.00E-34 | TerL [Rhodococcus phage ReqiPoco6] | gi\|593774729\|ref\|YP_009012597.1\| | 691964 |
| 7 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 71 | 70.37 | 2.00E-31 | TerL [Rhodococcus phage ReqiPoco6] | gi\|593774729\|ref\|YP_009012597.1\| | 691964 |
| 72 | 62.96 | 4.00E-27 | TerL [Rhodococcus phage ReqiPoco6] | gi\|593774729\|ref\|YP_009012597.1\| | 691964 |
| 75 | 70.37 | 8.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 77 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 8 | 56.79 | 9.00E-26 | TerL [Rhodococcus phage ReqiPoco6] | gi\|593774729\|ref\|YP_009012597.1\| | 691964 |
| 80 | 69.14 | 3.00E-31 | TerL [Rhodococcus phage ReqiPoco6] | gi\|593774729\|ref\|YP_009012597.1\| | 691964 |
| 81 | 74.07 | 7.00E-33 | putative uncharacterized protein [Ruminococcus sp. CAG:17] | gi\|547240587\|ref\|WP_021976510.1\| | 1262951 |
| 82 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 86 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 87 | 67.9 | 1.00E-27 | terminase [[Clostridium] hathewayi] | gi\|493833739\|ref\|WP_006781000.1\| | 154046 |
| 89 | 74.07 | 1.00E-34 | TerL [Rhodococcus phage ReqiPoco6] | gi\|593774729\|ref\|YP_009012597.1\| | 691964 |
| 9 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |
| 90 | 71.6 | 2.00E-33 | TerL [Rhodococcus phage ReqiPepy6] | gi\|593779801\|ref\|YP_009017628.1\| | 691965 |

777

778

779         SI Table 2. Taxonomic classification of closest homologs (HB1 phage family).

780         Majority of OTUs (86 out of 123) have the closest match to ReqiPoco6 terminase

781         large subunit, whereas 15 OTUs have closest homologs belonging to ReqiPepy6.

| Closest Homolog Taxon ID | Kingdom | Phylum | Class | Order | Family | Genus | Species |
|---|---|---|---|---|---|---|---|
| 1262951 | Bacteria | Firmicutes | Clostridia | Clostridiales | Ruminococcaceae | Ruminococcus | Ruminococcus sp. CAG:17 |
| 691964 | Viruses | dsDNA viruses, no RNA stage | Caudovirales | Siphoviridae | unclassified Siphoviridae | Rhodococcus phage | ReqiPoco6 |
| 691965 | Viruses | dsDNA viruses, no RNA stage | Caudovirales | Siphoviridae | unclassified Siphoviridae | Rhodococcus phage | ReqiPepy6 |
| 1512 | Bacteria | Firmicutes | Clostridia | Clostridiales | Lachnospiraceae | Lachnoclostridium | Clostridium symbiosum |
| 29347 | Bacteria | Firmicutes | Clostridia | Clostridiales | Lachnospiraceae | Lachnoclostridium | Clostridium scindens |
| 33039 | Bacteria | Firmicutes | Clostridia | Clostridiales | Lachnospiraceae | Blautia | Ruminococcus torques |
| 1232439 | Bacteria | Firmicutes | Clostridia | Clostridiales | unclassified Clostridiales | unclassified Clostridiales | Clostridiales bacterium VE202-03 |
| 154046 | Bacteria | Firmicutes | Clostridia | Clostridiales | Lachnospiraceae | Lachnoclostridium | Clostridium hathewayi |

782

783

784         SI Table 3. Closest homolog to each OTU's representative sequence (HA phage

785         family). Each OTU's representative sequence was used as a query for NCBI's

786         BLASTx homology search against the non-redundant protein database. The table

787         summarizes the E-value and the percent amino acid identity across the query

788         sequence and the closest homolog, as well as the closest homolog's name, sequence

789         ID, and taxon ID.  The taxon ID is color coded, and the taxonomic classification

790         corresponding to each taxon ID can be retrieved from the following table. Note with

791         the exception of a few "putative uncharacterized" homolog names that most are

792         identified as terminases or TerLs (terminase large subunits).

| Query Sequence ID (OTU ID) | Percent Identity | E value | Closest Homolog | Closest Homolog Sequence ID | Closest Homolog Taxon ID |
|---|---|---|---|---|---|
| 0 | 97.56 | 4.00E-49 | hypothetical protein [Streptococcus sp. F0442] | gi\|497418421\|ref\|WP_009732619.1\| | 999425 |
| 1 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 10 | 100 | 9.00E-50 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 100 | 98.78 | 4.00E-49 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 101 | 100 | 9.00E-50 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 102 | 100 | 2.00E-49 | terminase [Streptococcus infantis] | gi\|493136448\|ref\|WP_006154887.1\| | 68892 |
| 103 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 104 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 106 | 100 | 2.00E-49 | terminase [Streptococcus infantis] | gi\|493136448\|ref\|WP_006154887.1\| | 68892 |
| 107 | 100 | 2.00E-49 | terminase [Streptococcus infantis] | gi\|493136448\|ref\|WP_006154887.1\| | 68892 |
| 108 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 109 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 11 | 98.78 | 1.00E-48 | terminase [Streptococcus oralis] | gi\|446545997\|ref\|WP_000623343.1\| | 1303 |
| 110 | 100 | 9.00E-50 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 111 | 100 | 9.00E-50 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 112 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 113 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 114 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 117 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 118 | 100 | 9.00E-50 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 119 | 98.78 | 3.00E-49 | terminase [Streptococcus oralis] | gi\|446545997\|ref\|WP_000623343.1\| | 1303 |
| 12 | 100 | 2.00E-49 | terminase [Streptococcus infantis] | gi\|493136448\|ref\|WP_006154887.1\| | 68892 |
| 120 | 100 | 9.00E-50 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 121 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 122 | 100 | 9.00E-50 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 123 | 98.78 | 3.00E-49 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 125 | 100 | 9.00E-50 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 126 | 98.78 | 3.00E-49 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 127 | 98.78 | 3.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 128 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 129 | 98.78 | 1.00E-48 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 13 | 98.78 | 1.00E-49 | hypothetical protein [Streptococcus sp. F0442] | gi\|497418421\|ref\|WP_009732619.1\| | 999425 |
| 130 | 98.78 | 2.00E-49 | terminase [Streptococcus oralis] | gi\|446545997\|ref\|WP_000623343.1\| | 1303 |
| 131 | 98.78 | 6.00E-48 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 132 | 98.78 | 1.00E-48 | terminase [Streptococcus oralis] | gi\|446545997\|ref\|WP_000623343.1\| | 1303 |
| 133 | 98.78 | 7.00E-49 | terminase [Streptococcus infantis] | gi\|493136448\|ref\|WP_006154887.1\| | 68892 |
| 134 | 98.78 | 6.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 135 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 136 | 98.78 | 5.00E-49 | terminase [Streptococcus oralis] | gi\|446545997\|ref\|WP_000623343.1\| | 1303 |
| 137 | 100 | 9.00E-50 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 138 | 98.78 | 3.00E-49 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 139 | 98.78 | 1.00E-48 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 140 | 100 | 9.00E-50 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 141 | 98.78 | 1.00E-48 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 142 | 97.56 | 4.00E-48 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 143 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 145 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 146 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 147 | 100 | 9.00E-50 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 148 | 98.78 | 2.00E-48 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 15 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 150 | 98.78 | 9.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 151 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 152 | 98.78 | 3.00E-49 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 153 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 154 | 97.56 | 1.00E-48 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 155 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 156 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 157 | 98.78 | 3.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 158 | 97.56 | 1.00E-48 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 159 | 98.78 | 9.00E-49 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 16 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 160 | 98.78 | 4.00E-49 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 161 | 100 | 9.00E-50 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 162 | 98.78 | 3.00E-49 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 164 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 165 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 166 | 97.56 | 1.00E-48 | terminase [Streptococcus oralis] | gi\|446545997\|ref\|WP_000623343.1\| | 1303 |
| 167 | 100 | 2.00E-49 | terminase [Streptococcus infantis] | gi\|493136448\|ref\|WP_006154887.1\| | 68892 |
| 168 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 169 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 17 | 100 | 9.00E-50 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 170 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 171 | 98.78 | 3.00E-49 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 172 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 173 | 100 | 2.00E-49 | terminase [Streptococcus infantis] | gi\|493136448\|ref\|WP_006154887.1\| | 68892 |
| 174 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 175 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 176 | 100 | 9.00E-50 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 177 | 97.56 | 1.00E-48 | terminase [Streptococcus oralis] | gi\|446545997\|ref\|WP_000623343.1\| | 1303 |
| 178 | 98.78 | 4.00E-49 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 179 | 100 | 9.00E-50 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 18 | 100 | 2.00E-49 | terminase [Streptococcus infantis] | gi\|493136448\|ref\|WP_006154887.1\| | 68892 |
| 180 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 181 | 100 | 9.00E-50 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 182 | 100 | 9.00E-50 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |

| | | | | | |
|---|---|---|---|---|---|
| 183 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 184 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 185 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 186 | 100 | 9.00E-50 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 19 | 98.78 | 3.00E-49 | terminase [Streptococcus parasanguinis] | gi\|671602035\|ref\|WP_031575397.1\| | 1318 |
| 2 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 20 | 96.34 | 2.00E-48 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 21 | 100 | 9.00E-50 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 22 | 91.46 | 8.00E-46 | hypothetical protein [Streptococcus sp. F0442] | gi\|497418421\|ref\|WP_009732619.1\| | 999425 |
| 23 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 24 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 25 | 100 | 9.00E-50 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 26 | 98.78 | 3.00E-49 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 27 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 28 | 98.78 | 7.00E-50 | terminase [Streptococcus parasanguinis] | gi\|671602035\|ref\|WP_031575397.1\| | 1318 |
| 29 | 98.78 | 7.00E-49 | terminase [Streptococcus infantis] | gi\|493136448\|ref\|WP_006154887.1\| | 68892 |
| 3 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 30 | 100 | 2.00E-49 | terminase [Streptococcus infantis] | gi\|493136448\|ref\|WP_006154887.1\| | 68892 |
| 31 | 98.78 | 6.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 32 | 98.78 | 4.00E-49 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 33 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 34 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 35 | 98.78 | 1.00E-49 | hypothetical protein [Streptococcus sp. F0442] | gi\|497418421\|ref\|WP_009732619.1\| | 999425 |
| 36 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|565851306\|ref\|WP_023933954.1\| | 257758 |
| 37 | 100 | 9.00E-50 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 38 | 98.78 | 7.00E-50 | terminase [Streptococcus parasanguinis] | gi\|671602035\|ref\|WP_031575397.1\| | 1318 |
| 39 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 4 | 98.78 | 4.00E-49 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 40 | 98.78 | 2.00E-49 | hypothetical protein [Streptococcus sp. F0442] | gi\|497418421\|ref\|WP_009732619.1\| | 999425 |
| 42 | 100 | 9.00E-50 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 43 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 44 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 45 | 98.78 | 9.00E-49 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 47 | 98.78 | 7.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 48 | 98.78 | 3.00E-49 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 49 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 5 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 50 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 51 | 100 | 9.00E-50 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 52 | 98.78 | 3.00E-49 | terminase [Streptococcus oralis] | gi\|446545997\|ref\|WP_000623343.1\| | 1303 |
| 53 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 54 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 56 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 57 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 59 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 6 | 100 | 9.00E-50 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 60 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 61 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 62 | 98.78 | 4.00E-49 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 63 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 64 | 98.78 | 3.00E-49 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 65 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 66 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 68 | 100 | 9.00E-50 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 69 | 98.78 | 4.00E-49 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 7 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 71 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 72 | 100 | 1.00E-49 | ative phage terminase, large subunit [Streptococcus tigurin | gi\|494783687\|ref\|WP_007519095.1\| | 1077464 |
| 73 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 74 | 98.78 | 3.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 76 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 77 | 98.78 | 3.00E-49 | terminase [Streptococcus oralis] | gi\|446545997\|ref\|WP_000623343.1\| | 1303 |
| 78 | 98.78 | 3.00E-49 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 79 | 100 | 9.00E-50 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 8 | 100 | 9.00E-50 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 80 | 98.78 | 3.00E-49 | terminase [Streptococcus oralis] | gi\|446545997\|ref\|WP_000623343.1\| | 1303 |
| 82 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 83 | 97.56 | 1.00E-48 | terminase [Streptococcus oralis] | gi\|446545997\|ref\|WP_000623343.1\| | 1303 |
| 84 | 100 | 9.00E-50 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 85 | 100 | 9.00E-50 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 86 | 95.12 | 3.00E-48 | terminase [Streptococcus oralis] | gi\|446545997\|ref\|WP_000623343.1\| | 1303 |
| 87 | 98.78 | 3.00E-49 | terminase [Streptococcus oralis] | gi\|446545997\|ref\|WP_000623343.1\| | 1303 |
| 88 | 100 | 9.00E-50 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 89 | 98.78 | 2.00E-48 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 9 | 98.78 | 3.00E-49 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 90 | 98.78 | 4.00E-49 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 93 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 95 | 100 | 9.00E-50 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 96 | 98.78 | 4.00E-49 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 97 | 100 | 1.00E-49 | terminase [Streptococcus pseudopneumoniae] | gi\|446545996\|ref\|WP_000623342.1\| | 257758 |
| 98 | 100 | 9.00E-50 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |
| 99 | 97.56 | 1.00E-48 | terminase [Streptococcus sp. SR1] | gi\|726981126\|ref\|WP_033585808.1\| | 1161416 |

794

795

796     SI Table 4. Taxonomic classification of closest homologs to each OTU's
797     representative sequence (HA phage family).

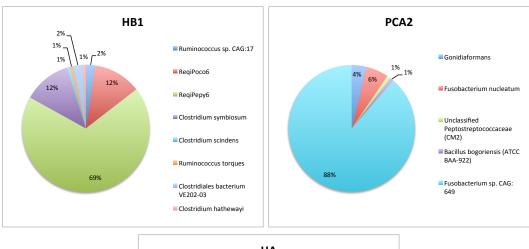| Closest Homolog Taxon ID | Kingdom | Phylum | Class | Order | Family | Genus | Species |
|---|---|---|---|---|---|---|---|
| 1318 | Bacteria | Firmicutes | Bacilli | Lactobacillales | Streptococcaceae | Streptococcus | Streptococcus parasanguinis |
| 1077464 | Bacteria | Firmicutes | Bacilli | Lactobacillales | Streptococcaceae | Streptococcus | Streptococcus tigurinus |
| 257758 | Bacteria | Firmicutes | Bacilli | Lactobacillales | Streptococcaceae | Streptococcus | Streptococcus pseudopneumoniae |
| 999425 | Bacteria | Firmicutes | Bacilli | Lactobacillales | Streptococcaceae | Streptococcus | Streptococcus sp. F0442 |
| 1161416 | Bacteria | Firmicutes | Bacilli | Lactobacillales | Streptococcaceae | Streptococcus | Streptococcus sp. SR1 |
| 1303 | Bacteria | Firmicutes | Bacilli | Lactobacillales | Streptococcaceae | Streptococcus | Streptococcus oralis |
| 68892 | Bacteria | Firmicutes | Bacilli | Lactobacillales | Streptococcaceae | Streptococcus | Streptococcus infantis |

798

799

800     SI Table 5. Closest homolog to each OTU's representative sequence (PCA2 phage
801     family). Each OTU's representative sequence was used as a query for NCBI's
802     BLASTx homology search against the non-redundant protein database. The table
803     summarizes the E-value and the percent amino acid identity across the query
804     sequence and the closest homolog, as well as the closest homolog's name, sequence
805     ID, and taxon ID.  The taxon ID is color coded, and the taxonomic classification
806     corresponding to each taxon ID can be retrieved from the following table. Note with
807     the exception of a few "putative uncharacterized" homolog names, most are
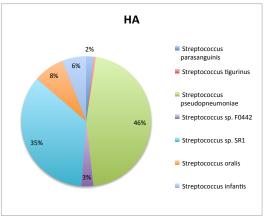808     identified as terminases or TerLs (terminase large subunits).

| Query Sequence ID (OTU ID) | Percent Identity | E value | Closest Homolog | Closest Homolog Sequence ID | Closest Homolog Taxon ID |
|---|---|---|---|---|---|
| 0 | 100 | 7.00E-32 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 10 | 90 | 2.00E-28 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 100 | 76.27 | 4.00E-21 | terminase [Peptostreptococcaceae bacterium CM2] | gi\|497213446\|ref\|WP_009527708.1\| | 796939 |
| 101 | 98.33 | 1.00E-30 | terminase [Fusobacterium nucleatum] | gi\|495968206\|ref\|WP_008692785.1\| | 851 |
| 103 | 98.33 | 2.00E-30 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 104 | 75 | 1.00E-21 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 107 | 98.33 | 1.00E-31 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 108 | 100 | 7.00E-32 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 109 | 75 | 1.00E-21 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 11 | 93.33 | 3.00E-30 | terminase [Fusobacterium periodonticum] | gi\|496096975\|ref\|WP_008821482.1\| | 860 |
| 110 | 75 | 1.00E-21 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 111 | 98.33 | 3.00E-31 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 112 | 100 | 7.00E-32 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 113 | 98.33 | 6.00E-31 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 115 | 73.33 | 4.00E-21 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 116 | 73.33 | 5.00E-21 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 117 | 100 | 7.00E-32 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 12 | 100 | 7.00E-32 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 120 | 100 | 7.00E-32 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 123 | 83.33 | 6.00E-25 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 124 | 100 | 7.00E-32 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 126 | 98.33 | 1.00E-30 | terminase [Fusobacterium nucleatum] | gi\|495968206\|ref\|WP_008692785.1\| | 851 |
| 127 | 100 | 7.00E-32 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 128 | 100 | 7.00E-32 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 129 | 100 | 7.00E-32 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 13 | 35.59 | 4.00E-04 | terminase [Bacillus bogoriensis] | gi\|651939129\|ref\|WP_026673624.1\| | 246272 |
| 132 | 98.33 | 1.00E-30 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 135 | 100 | 7.00E-32 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 136 | 100 | 7.00E-32 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 137 | 100 | 7.00E-32 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 138 | 75 | 1.00E-21 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 139 | 96.67 | 4.00E-31 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 14 | 96.67 | 5.00E-30 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 140 | 100 | 7.00E-32 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 141 | 100 | 7.00E-32 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 142 | 100 | 7.00E-32 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 143 | 91.67 | 1.00E-29 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 144 | 98.33 | 4.00E-31 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 146 | 98.33 | 2.00E-31 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 147 | 98.33 | 6.00E-31 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 149 | 100 | 7.00E-32 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 15 | 73.33 | 2.00E-21 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 150 | 98.33 | 2.00E-31 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 151 | 98.33 | 3.00E-31 | terminase [Fusobacterium periodonticum] | gi\|496096975\|ref\|WP_008821482.1\| | 860 |
| 152 | 98.33 | 1.00E-30 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 153 | 98.33 | 2.00E-31 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 154 | 100 | 7.00E-32 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 155 | 98.33 | 6.00E-31 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 156 | 98.33 | 5.00E-31 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 16 | 98.33 | 1.00E-31 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 2 | 91.67 | 1.00E-29 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 21 | 73.33 | 7.00E-21 | terminase [Fusobacterium nucleatum] | gi\|495968206\|ref\|WP_008692785.1\| | 851 |
| 22 | 75 | 1.00E-21 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 23 | 73.33 | 5.00E-21 | terminase [Fusobacterium nucleatum] | gi\|495968206\|ref\|WP_008692785.1\| | 851 |
| 24 | 100 | 7.00E-32 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 25 | 98.33 | 4.00E-31 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 26 | 98.33 | 1.00E-30 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 27 | 75 | 1.00E-21 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 29 | 98.33 | 2.00E-31 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 3 | 100 | 7.00E-32 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 32 | 98.33 | 3.00E-31 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 33 | 100 | 7.00E-32 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 34 | 98.33 | 4.00E-31 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 36 | 96.67 | 1.00E-29 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 4 | 96.67 | 1.00E-30 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 40 | 75 | 1.00E-21 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 41 | 98.33 | 3.00E-31 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 45 | 100 | 7.00E-32 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 47 | 73.33 | 1.00E-21 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 5 | 75 | 1.00E-21 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 50 | 100 | 7.00E-32 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 51 | 98.33 | 2.00E-31 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 54 | 100 | 7.00E-32 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 56 | 100 | 7.00E-32 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 60 | 100 | 7.00E-32 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 62 | 100 | 7.00E-32 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 63 | 100 | 7.00E-32 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 64 | 93.33 | 3.00E-29 | terminase [Fusobacterium periodonticum] | gi\|496096975\|ref\|WP_008821482.1\| | 860 |
| 65 | 100 | 7.00E-32 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 66 | 98.33 | 3.00E-31 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 67 | 100 | 7.00E-32 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |

809

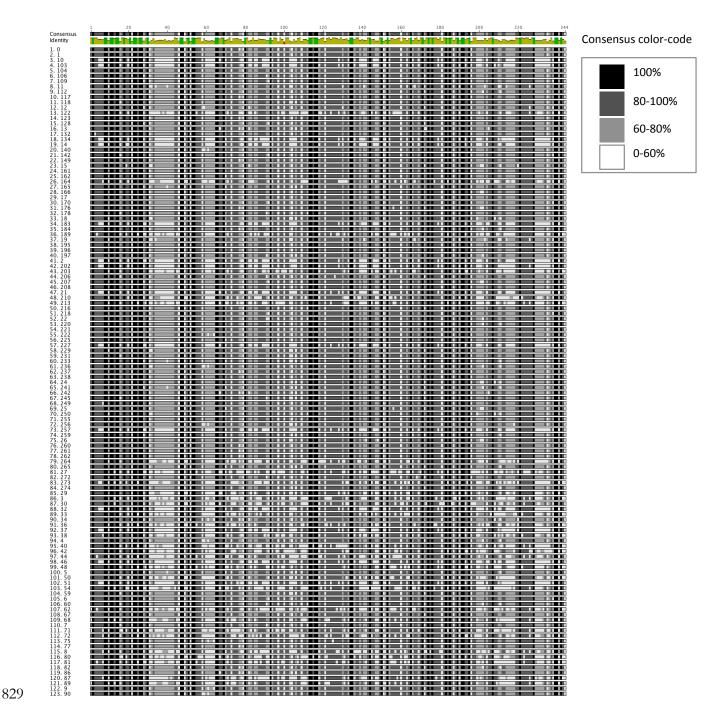| | | | | | | |
|---|---|---|---|---|---|---|
| 70 | 100 | 7.00E-32 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 72 | 98.33 | 3.00E-31 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 73 | 100 | 7.00E-32 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 75 | 96.67 | 2.00E-30 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 77 | 100 | 7.00E-32 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 78 | 90 | 5.00E-29 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 79 | 100 | 7.00E-32 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 8 | 95 | 5.00E-30 | terminase [Fusobacterium periodonticum] | gi\|496096975\|ref\|WP_008821482.1\| | 860 |
| 80 | 100 | 7.00E-32 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 81 | 100 | 7.00E-32 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 82 | 100 | 7.00E-32 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 83 | 100 | 7.00E-32 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 84 | 98.33 | 1.00E-30 | terminase [Fusobacterium nucleatum] | gi\|495968206\|ref\|WP_008692785.1\| | 851 |
| 86 | 100 | 7.00E-32 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 89 | 90 | 5.00E-29 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 90 | 98.33 | 3.00E-31 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 91 | 98.33 | 3.00E-31 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 92 | 98.33 | 8.00E-31 | terminase [Fusobacterium nucleatum] | gi\|495968206\|ref\|WP_008692785.1\| | 851 |
| 93 | 96.67 | 9.00E-31 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 94 | 96.67 | 4.00E-31 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 95 | 100 | 7.00E-32 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 96 | 100 | 7.00E-32 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 97 | 98.33 | 3.00E-31 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |
| 99 | 98.33 | 3.00E-31 | putative phage terminase large subunit [Fusobacterium sp. CAG:649] | gi\|547450305\|ref\|WP_022069933.1\| | 1262900 |

SI Table 6. Taxonomic classification of closest homologs to each OTU's

representative sequence (PCA2 phage family).

| Closest Homolog Taxon ID | Kingdom | Phylum | Class | Order | Family | Genus | Species |
|---|---|---|---|---|---|---|---|
| 860 | Bacteria | Fusobacteria | Fusobacteria | Fusobacteriales | Fusobacteriaceae | Fusobacterium | gonidiaformans |
| 851 | Bacteria | Fusobacteria | Fusobacteria | Fusobacteriales | Fusobacteriaceae | Fusobacterium | Fusobacterium nucleatum |
| 796939 | Bacteria | Firmicutes | Clostridia | Clostridiales | Peptostreptococcaceae | unclassified Peptostreptococcaceae | unclassified Peptostreptococcaceae (CM2) |
| 246272 | Bacteria | Firmicutes | Bacilli | Bacillales | Bacillaceae | Bacillus | Bacillus bogoriensis (ATCC BAA-922) |
| 1262900 | Bacteria | Fusobacteria | Fusobacteriia | Fusobacteriales | Fusobacteriaceae | Fusobacterium | Fusobacterium sp. CAG:649 |

SI Figure 3. Percentage of HB1, PCA2 and HA phage family OTUs belonging to each taxonomic group identified in SI Figure 2, SI Figure 1, and SI Table 3.

829

830    SI Figure 4. The nucleotide alignment of HB1 phage family OTU representative

831    sequences. Sequences were aligned using Geneious (79). No gaps were introduced.

832    Each base is color-coded according to its relative abundance within a column in the

833    alignment. Conserved bases are black and highly variable sites are denoted in white

## References

1. Suttle CA, Chan AM, & Cottrell MT (1990) Infection of phytoplankton by viruses and reduction of primary productivity. *Nature* 347(6292):467-469.
2. Bergh Ø, BØrsheim KY, Bratbak G, & Heldal M (1989) High abundance of viruses found in aquatic environments. *Nature* 340(6233):467-468.
3. Suttle CA (2007) Marine viruses—major players in the global ecosystem. *Nature Reviews Microbiology* 5(10):801-812.
4. Rohwer F & Thurber RV (2009) Viruses manipulate the marine environment. *Nature* 459(7244):207-212.
5. Roux S*, et al.* (2016) Ecogenomics and potential biogeochemical impacts of globally abundant ocean viruses. *Nature.*
6. Mokili JL, Rohwer F, & Dutilh BE (2012) Metagenomics and future perspectives in virus discovery. *Current opinion in virology* 2(1):63-77.
7. Navarro F & Muniesa M (2017) Phages in the Human Body. *Frontiers in microbiology* 8.
8. Szafrański SP, Winkel A, & Stiesch M (2017) The use of bacteriophages to biocontrol oral biofilms. *Journal of biotechnology* 250:29-44.
9. Abedon ST, Kuhl SJ, Blasdel BG, & Kutter EM (2011) Phage treatment of human infections. *Bacteriophage* 1(2):66-85.
10. Kutter E*, et al.* (2010) Phage therapy in clinical practice: treatment of human infections. *Current pharmaceutical biotechnology* 11(1):69-86.
11. Endersen L*, et al.* (2014) Phage therapy in the food industry. *Annual review of food science and technology* 5:327-349.
12. Nagy JK, Király L, & Schwarczinger I (2012) Phage therapy for plant disease control with a focus on fire blight. *Central European Journal of Biology* 7(1):1-12.
13. Haq IU, Chaudhry WN, Akhtar MN, Andleeb S, & Qadri I (2012) Bacteriophages and their implications on future biotechnology: a review. *Virology journal* 9(1):9.
14. Schoenfeld T*, et al.* (2010) Functional viral metagenomics and the next generation of molecular tools. *Trends in microbiology* 18(1):20-29.

870    15.    Hamzeh-Mivehroud M, Alizadeh AA, Morris MB, Church WB, & Dastmalchi S (2013)
871            Phage display as a technology delivering on the promise of peptide drug discovery. *Drug*
872            *discovery today* 18(23-24):1144-1157.

873    16.    Youle M, Haynes M, & Rohwer F (2012) Scratching the surface of biology's dark matter.
874            *Viruses: Essential agents of life*, (Springer), pp 61-81.

875    17.    Paez-Espino D*, et al.* (2016) Uncovering Earth's virome. *Nature* 536(7617):425-430.

876    18.    Reyes A, Semenkovich NP, Whiteson K, Rohwer F, & Gordon JI (2012) Going viral: next
877            generation sequencing applied to human gut phage populations. *Nature Reviews. Microbiology*
878            10(9):607.

879    19.    Penadés JR, Chen J, Quiles-Puchalt N, Carpena N, & Novick RP (2015) Bacteriophage-
880            mediated spread of bacterial virulence genes. *Current opinion in microbiology* 23:171-178.

881    20.    Stone R (2002) Food and agriculture: testing grounds for phage therapy. *Science*
882            298(5594):730-730.

883    21.    Bik EM*, et al.* (2010) Bacterial diversity in the oral cavity of 10 healthy individuals. *The ISME*
884            *journal* 4(8):962-974.

885    22.    Edlund A, Santiago-Rodriguez TM, Boehm TK, & Pride DT (2015) Bacteriophage and their
886            potential roles in the human oral cavity. *Journal of oral microbiology* 7.

887    23.    Nasidze I, Li J, Quinque D, Tang K, & Stoneking M (2009) Global diversity in the human
888            salivary microbiome. *Genome research* 19(4):636-643.

889    24.    Chen T*, et al.* (2010) The Human Oral Microbiome Database: a web accessible resource for
890            investigating oral microbe taxonomic and genomic information. *Database* 2010.

891    25.    Abeles SR & Pride DT (2014) Molecular bases and role of viruses in the human microbiome.
892            *Journal of molecular biology* 426(23):3892-3906.

893    26.    Pride DT*, et al.* (2012) Evidence of a robust resident bacteriophage population revealed
894            through analysis of the human salivary virome. *The ISME journal* 6(5):915.

895    27.    Ly M*, et al.* (2014) Altered oral viral ecology in association with periodontal disease. *MBio*
896            5(3):e01133-01114.

897    28.    Santiago-Rodriguez TM*, et al.* (2015) Transcriptome analysis of bacteriophage communities
898            in periodontal health and disease. *BMC genomics* 16(1):549.

899    29.    Willner D*, et al.* (2011) Metagenomic detection of phage-encoded platelet-binding factors in
900            the human oral cavity. *Proceedings of the National Academy of Sciences* 108(Supplement 1):4547-
901            4553.

902    30.    Roberts AP & Mullany P (2010) Oral biofilms: a reservoir of transferable, bacterial,
903            antimicrobial resistance. *Expert review of anti-infective therapy* 8(12):1441-1450.

904    31.    Hug LA*, et al.* (2016) A new view of the tree of life. *Nature Microbiology* 1:16048.

905    32.    Yarza P*, et al.* (2014) Uniting the classification of cultured and uncultured bacteria and
906            archaea using 16S rRNA gene sequences. *Nature Reviews. Microbiology* 12(9):635.

907    33.    Woese CR, Kandler O, & Wheelis ML (1990) Towards a natural system of organisms:
908            proposal for the domains Archaea, Bacteria, and Eucarya. *Proceedings of the National Academy of*
909            *Sciences* 87(12):4576-4579.

910    34.    Wu GD*, et al.* (2011) Linking long-term dietary patterns with gut microbial enterotypes.
911            *Science* 334(6052):105-108.

912    35.    Caporaso JG*, et al.* (2011) Global patterns of 16S rRNA diversity at a depth of millions of
913            sequences per sample. *Proceedings of the National Academy of Sciences* 108(Supplement 1):4516-
914            4522.

915    36.    Yatsunenko T*, et al.* (2012) Human gut microbiome viewed across age and geography. *nature*
916            486(7402):222.

917    37.    Shanks OC, *et al.* (2013) Comparison of the microbial community structures of untreated
918          wastewaters from different geographic locales. *Applied and environmental microbiology*
919          79(9):2906-2913.

920    38.    Barberán A, *et al.* (2015) Continental-scale distributions of dust-associated bacteria and fungi.
921          *Proceedings of the National Academy of Sciences* 112(18):5756-5761.

922    39.    Yu F, *et al.* (2017) Microfluidic-based mini-metagenomics enables discovery of novel
923          microbial lineages from complex environmental samples. *bioRxiv*:114496.

924    40.    Dutilh BE, *et al.* (2014) A highly abundant bacteriophage discovered in the unknown
925          sequences of human faecal metagenomes. *Nature communications* 5.

926    41.    Yang X, *et al.* (2012) De novo assembly of highly diverse viral populations. *BMC genomics*
927          13(1):475.

928    42.    Li Y, *et al.* (2016) VIP: an integrated pipeline for metagenomics of virus identification and
929          discovery. *Scientific reports* 6.

930    43.    Wüthrich D, *et al.* (2016) Exploring the virome of cattle with non-suppurative encephalitis of
931          unknown etiology by metagenomics. *Virology* 493:22-30.

932    44.    Manso CF, Bibby DF, & Mbisa JL (2017) Efficient and unbiased metagenomic recovery of
933          RNA virus genomes from human plasma samples. *Scientific Reports* 7.

934    45.    Sanschagrin S & Yergeau E (2014) Next-generation sequencing of 16S ribosomal RNA gene
935          amplicons. *Journal of visualized experiments: JoVE* (90).

936    46.    Steven B, Gallegos‐Graves LV, Starkenburg SR, Chain PS, & Kuske CR (2012) Targeted
937          and shotgun metagenomic approaches provide different descriptions of dryland soil
938          microbial communities in a manipulated field study. *Environmental microbiology reports* 4(2):248-
939          256.

940    47.    Tessler M, *et al.* (2017) Large-scale differences in microbial biodiversity discovery between
941          16S amplicon and shotgun sequencing. *Scientific Reports* 7.

942    48.    Ng TFF, *et al.* (2011) Broad surveys of DNA viral diversity obtained through viral
943          metagenomics of mosquitoes. *PloS one* 6(6):e20579.

944    49.    Grose JH & Casjens SR (2014) Understanding the enormous diversity of bacteriophages: the
945          tailed phages that infect the bacterial family Enterobacteriaceae. *Virology* 468:421-443.

946    50.    Suttle CA (2005) Viruses in the sea. *Nature* 437(7057):356-361.

947    51.    Casjens S (2003) Prophages and bacterial genomics: what have we learned so far? *Molecular*
948          *microbiology* 49(2):277-300.

949    52.    Tadmor AD & Phillips R (2015) Host-Virus Interaction: From Metagenomics to Single-Cell
950          Genomics. *Encyclopedia of Metagenomics*, (Springer), pp 257-265.

951    53.    Moore SD & Prevelige Jr PE (2002) DNA packaging: a new class of molecular motors.
952          *Current biology* 12(3):R96-R98.

953    54.    Rao VB & Feiss M (2008) The bacteriophage DNA packaging motor. *Annual review of genetics*
954          42:647-681.

955    55.    Mahmoudabadi G & Phillips R (2018) A comprehensive and quantitative exploration of
956          thousands of viral genomes. *eLife* 7:e31955.

957    56.    Lloyd-Price J, *et al.* (2017) Strains, functions and dynamics in the expanded Human
958          Microbiome Project. *Nature* 550(7674):61-66.

959    57.    Huttenhower C, *et al.* (2012) Structure, function and diversity of the healthy human
960          microbiome. *Nature* 486(7402):207.

961    58.    Tian N, *et al.* (2017) Salivary Gluten Degradation and Oral Microbial Profiles in Health and
962          Celiac Disease. *Applied and environmental microbiology*:AEM. 03330-03316.

963  59.  Belda-Ferre P*, et al.* (2012) The oral metagenome in health and disease. *The ISME journal*
964      6(1):46.
965  60.  Breitbart M & Rohwer F (2005) Here a virus, there a virus, everywhere the same virus?
966      *Trends in microbiology* 13(6):278-284.
967  61.  Holmfeldt K*, et al.* (2013) Twelve previously unknown phage genera are ubiquitous in global
968      oceans. *Proceedings of the National Academy of Sciences* 110(31):12798-12803.
969  62.  Hall MW*, et al.* (2017) Inter-personal diversity and temporal dynamics of dental, tongue, and
970      salivary microbiota in the healthy oral cavity. *npj Biofilms and Microbiomes* 3(1):2.
971  63.  Costello EK*, et al.* (2009) Bacterial community variation in human body habitats across space
972      and time. *Science* 326(5960):1694-1697.
973  64.  Abeles SR*, et al.* (2014) Human oral viruses are personal, persistent and gender-consistent.
974      *The ISME journal* 8(9):1753-1767.
975  65.  Pride DT*, et al.* (2012) Evidence of a robust resident bacteriophage population revealed
976      through analysis of the human salivary virome. *The ISME journal* 6(5):915-926.
977  66.  Oh J*, et al.* (2016) Temporal stability of the human skin microbiome. *Cell* 165(4):854-866.
978  67.  Belstrøm D*, et al.* (2016) Temporal stability of the salivary microbiota in oral health. *PLoS*
979      *One* 11(1):e0147472.
980  68.  Franzosa EA*, et al.* (2015) Identifying personal microbiomes using metagenomic codes.
981      *Proceedings of the National Academy of Sciences* 112(22):E2930-E2938.
982  69.  Franzosa EA*, et al.* (2015) Identifying personal microbiomes using metagenomic codes.
983      *Proceedings of the National Academy of Sciences*:201423854.
984  70.  Tadmor AD, Ottesen EA, Leadbetter JR, & Phillips R (2011) Probing individual
985      environmental bacteria for viruses by using microfluidic digital PCR. *Science* 333(6038):58-62.
986  71.  Xie G*, et al.* (2010) Community and gene composition of a human dental plaque microbiota
987      obtained by metagenomic sequencing. *Molecular oral microbiology* 25(6):391-405.
988  72.  Turnbaugh PJ*, et al.* (2007) The human microbiome project: exploring the microbial part of
989      ourselves in a changing world. *Nature* 449(7164):804.
990  73.  Hamady M, Walker JJ, Harris JK, Gold NJ, & Knight R (2008) Error-correcting barcoded
991      primers for pyrosequencing hundreds of samples in multiplex. *Nature methods* 5(3):235-237.
992  74.  Caporaso JG*, et al.* (2010) QIIME allows analysis of high-throughput community sequencing
993      data. *Nature methods* 7(5):335-336.
994  75.  Edgar RC (2010) Search and clustering orders of magnitude faster than BLAST. *Bioinformatics*
995      26(19):2460-2461.
996  76.  Chao A, Chazdon RL, Colwell RK, & Shen TJ (2005) A new statistical approach for
997      assessing similarity of species composition with incidence and abundance data. *Ecology letters*
998      8(2):148-159.
999  77.  Bray JR & Curtis JT (1957) An ordination of the upland forest communities of southern
1000      Wisconsin. *Ecological monographs* 27(4):325-349.
1001  78.  Lozupone C & Knight R (2005) UniFrac: a new phylogenetic method for comparing
1002      microbial communities. *Applied and environmental microbiology* 71(12):8228-8235.
1003  79.  Kearse M*, et al.* (2012) Geneious Basic: an integrated and extendable desktop software
1004      platform for the organization and analysis of sequence data. *Bioinformatics* 28(12):1647-1649.
1005  80.  Bastian M, Heymann S, & Jacomy M (2009) Gephi: an open source software for exploring
1006      and manipulating networks. *Icwsm* 8:361-362.
1007