1 **Phanta: Phage-inclusive profiling of human gut metagenomes**

2

3 Yishay Pinto[1,2,*], Meenakshi Chakraborty[1,*], Navami Jain[1,2], Ami S Bhatt[1,2]

4

5 1. Department of Genetics, Stanford University, Stanford, CA, USA.

6 2. Department of Medicine, Divisions of Hematology and Blood & Marrow

7 Transplantation, Stanford University, Stanford, CA, USA.

8 * These authors contributed equally

9
10
11
12 ## **Abstract**

13

14 The human gut microbiome is a diverse ecosystem that encompasses multiple domains of life
15 and plays a vital role in human health. Due to technical limitations, most microbiome studies have
16 focused on gut prokaryotes, overlooking bacteriophages and other gut viruses. The most common
17 method to profile viruses is to assemble shotgun metagenomic reads - often from virus-enriched
18 samples - and identify viral genomes *de novo*. While valuable, this resource-intensive and
19 reference-independent method has limited sensitivity. To overcome these drawbacks, we
20 developed Phanta, which profiles human gut metagenomes in a virus-inclusive manner directly
21 from short reads utilizing recently published catalogs of gut viral genomes. Phanta incorporates
22 *k*-mer based classification tools and was developed with virus-specific properties in mind.
23 Specifically, it includes optimizations considering viruses' small genome size, sequence
24 homology with prokaryotes, and interactions with other members of the gut microbial community.
25 Based on simulations, the workflow is fast and accurate with respect to both prokaryotes and
26 viruses, minimizing false positive species identification using a novel genome coverage-based
27 strategy. When applied to metagenomes from healthy adults, Phanta identified ~200 viral species
28 per sample, ~5x more than the standard assembly-based methods. Notably, we observed a 2:1
29 ratio between gut viruses and bacteria, with higher interindividual variability of the gut virome
30 compared to the gut bacteriome. Phanta performs equally well on bulk vs. virus-enriched
31 metagenomes, making it possible to study prokaryotes and viruses in a single experiment, with a
32 single analysis. Phanta can tandemly profile gut viruses and prokaryotes in existing and novel
33 datasets, and can therefore identify cross-domain interactions with likely relevance to human
34 health. We expect that Phanta will reduce the barrier to virus-inclusive studies of the human gut
35 microbiome, thus making it standard practice.
36
37
38
39
40
41
42

## **Introduction**

The human gut microbiome is an ecosystem of diverse microorganisms including archaea, bacteria, viruses, and fungi. It plays a vital role in human health by interacting with our immune, digestive, and nervous systems[1–4]. Since the 1970s, tools such as 16S rRNA sequencing have enabled us to identify prokaryotic taxa present in the gut[5], and therefore to determine crucial relationships between these taxa and human health, age, lifestyle, environment, geography, and demographics[6–9]. However, these fundamental techniques overlook the viral fraction of the microbiome, preventing us from evaluating the impact of the human gut virome on human health.

Shotgun metagenomics is a popular and affordable method to sequence metagenomic samples[10–13]. This method captures genomic DNA from all gut organisms, not only prokaryotes, making it an optimal tool to study DNA viruses of the virome[14–16]. In the past decade, thousands of human microbiome samples have been analyzed using this "domain-inclusive" method[17–20]. Human gut prokaryotes can be well-quantified from shotgun metagenomes through direct read classification by comparing sequencing reads to reference genomes[18,19,21–24]. However, in the absence of comprehensive catalogs of viral genomes, the most common method for profiling the virome from shotgun metagenomes has been to assemble sequencing reads into contigs and identify viral genomes *de novo*[25,26]. Assembly-based approaches overcome the fundamental limitation that, until recently, a majority of phages had no reference genome[27]. However, despite their strengths at *de novo* phage discovery, assembly-based approaches have limited ability to detect low-abundance phages, due to the relative difficulty of assembling the genomes of low abundance taxa[28–30].

With increases in shotgun metagenomes from human gut samples, more comprehensive databases of gut viral genomes have recently been created[27,31–37]. By using these new compendiums, it is now possible to profile gut viruses and their prokaryotic hosts simultaneously through read-based, reference-dependent methods. This approach can address the sensitivity limitation of assembly-based methods to profile the virome, resulting in much more complete profiles of the gut microbiome with both prokaryotes and viruses accurately represented.

In this paper, we present Phanta, a fast and accurate virus-inclusive profiler of human gut metagenomes based on classification of short reads to our newly constructed, comprehensive database of human gut microbes. The provided database contains the latest genome catalogs from multiple domains of life, including more than 190,000 phage genomes and the entire HumGut collection of prokaryotic genomes[19,27]. Phanta incorporates the state-of-the-art tools Kraken2[22] and Bracken[38], and complements them with additional filtering steps and optimizations specifically tailored to the challenges of gut viral quantification. Phanta accurately quantifies both bacteria and phage abundances in simulated mixed communities. In metagenomes from healthy human adults, Phanta identifies >100-fold more viral reads and minimizes unclassified reads when compared to the default Kraken2/Bracken databases and workflow. In addition, due to its high sensitivity, Phanta identifies 5-fold more viral species than a common workflow of contig assembly and viral sequence identification. Finally, Phanta quantifies just as many viruses when applied to bulk shotgun metagenomes vs. matched metagenomes enriched for virus-like particles. This

87  demonstrates that it is possible to profile multiple domains of life from a single metagenomic
88  sequencing experiment, as opposed to needing an additional sequencing experiment after
89  enrichment for virus-like particles. Taken together, we anticipate that Phanta, which is freely
90  available at https://github.com/bhattlab/phanta, will facilitate improved profiling of cross-domain
91  interactions in gut microbiomes.
92
93  **Results**
94
95  **Phanta: A workflow for phage-inclusive profiling of human gut metagenomes**
96  Phanta was developed to generate accurate and complete profiles of human gut metagenomes,
97  with the goal of deepening our understanding of cross-domain interactions in the gut. To achieve
98  this objective, we first constructed a comprehensive database of gut microbial genomes found in
99  humans. To minimize false mapping, it was important to curate comprehensive collections of
100  genomes from all groups of taxa residing in the human gut - not only phages and other viruses,
101  but also prokaryotes, eukaryotes, and possible contaminants. For this purpose, we used the
102  HumGut collection as a reference for both human gut bacteria and archaea[19]. HumGut includes
103  dereplicated genomes from both UHGG and RefSeq. For viruses, we used the Metagenomic Gut
104  Virus catalog (MGV; dominated by human gut phages)[21,27] and RefSeq. For gut eukaryotes, we
105  also used RefSeq, and for contaminants, we used the human genome (hg38) and the Core
106  UniVec database from NCBI[22]. To create an informative viral taxonomy, MGV genomes were first
107  clustered to species-level operational taxonomic units (vOTUs). MGV vOTUs with high similarity
108  to a RefSeq viral species were labeled with the NCBI-assigned taxonomy of that species. For the
109  remaining MGV vOTUs, higher levels of taxonomy were assigned iteratively (see Methods).
110
111  The first step of Phanta is read classification to a database of reference genomes, such as that
112  described above (Fig. 1A). As viruses have relatively low abundance in a typical metagenomic
113  sample, we chose to use whole genome classification, which is typically more sensitive in the low-
114  coverage regime than methods relying on clade-specific marker genes[22,39]. Specifically, Phanta
115  classifies reads to the lowest possible taxonomic rank by Kraken2[22,24], a *k*-mer-based method that
116  has been shown to be both fast and accurate given the correct database and optimized
117  parameters[39]. Second, Phanta reduces false positive species by filtering out species based on a
118  calculated proxy for genome coverage (see Methods), a known issue in taxonomic classification[40].
119  Third, Phanta quantifies species-level relative abundances by executing Bracken, a tool
120  complementary to Kraken2 that redistributes all classified reads to the species level using a
121  Bayesian inference approach[38]. By default, Bracken calculates the "relative read abundance" -
122  the proportion of reads assigned to a species out of all reads. However, since viral genomes can
123  be orders of magnitude smaller than prokaryotic genomes, read abundance approaches inflate
124  the relative signal from prokaryotes within a community. Therefore, we additionally calculate
125  "relative taxonomic abundance", which instead estimates the relative proportion of different
126  organisms (not proportion of DNA sequence) within a given sample [41]. Briefly, we adjust the
127  relative read abundance of each species using the median length of the species' genomes. This
128  provides a comparable abundance estimation to amplicon sequencing or marker gene-based
129  approaches (Fig. 1B). Lastly, Phanta allows users to determine cross-domain relationships by

130 summing viral abundances by predicted host, providing information about the predicted virulence
131 of the viral community, and correlating the abundances of phages and bacteria.
132
133 **Phanta accurately classifies short reads from simulated mixed microbial communities**
134 To evaluate the performance of Phanta, we simulated 10 mixed communities, each containing a
135 total of ~6.5M 150 base pair (bp) paired-end reads from a combination of 300 prokaryotic
136 genomes and 50 viral genomes (see Methods). The relative read abundance of prokaryotes and
137 viruses in the resulting simulated samples was 0.95 and 0.05, respectively (Figure 2A). Phanta
138 accurately assigned reads to the right domain with average read abundance of 0.951±0.004 mean
139 read abundance for prokaryotes, and 0.048±0.004 mean read abundance for viruses (Figure 2B;
140 Supplementary Data File 1).
141
142 We next used the simulated communities to test the accuracy of classification of reads by
143 Kraken2. Reads were classified with high precision to all taxonomic ranks, with 63% of reads
144 classified to the species level or lower (median across simulated communities; see Figure 2C).
145 Next, we tested the accuracy of Phanta in estimating the abundance of each simulated species.
146 Phanta's species-level estimates for relative read abundance were highly correlated with the true
147 simulated values - Pearson's R=0.997 for all species (including bacteria and archaea), R=0.998
148 for bacterial species (Figure 2D), and R=0.925 for viral species (Figure 2E).
149
150 **Phanta's filtering step significantly reduces false positive species identification**
151 While developing Phanta, we observed that even a small fraction of mis-classified reads can lead
152 to a non-negligible number of falsely identified species. Therefore, to increase the signal-to-noise
153 of the identified species, we made the following modifications to the default Kraken2-Bracken
154 workflow. First, we introduced a filtering step between Kraken2 and Bracken that estimates the
155 breadth of genome coverage for species detected by Kraken2 and filters out likely false positive
156 species based on a user-adjustable coverage threshold. In addition, for a read to be classified by
157 Kraken2, we required that a certain fraction of a read's *k*-mers be mapped to a given taxon, in
158 order for the read to receive that classification. To achieve this, we adjusted Kraken2's confidence
159 threshold. By default, Phanta uses a confidence threshold of 0.1 (vs. 0 for default Kraken2; also
160 recommended by [39]), and this can be further adjusted by the user. These steps reduced false
161 positive species by 50-fold with minimal reduction of true positive species relative to a consecutive
162 run of Kraken2 and Bracken using default parameters (Figures 2F-G). Overall, we demonstrated
163 that Phanta performs with high accuracy in both classifying reads and estimating abundance while
164 substantially reducing false species identification.
165
166 **Masking prophages in prokaryotic genomes further increases sensitivity to viral reads**
167 Due to genetic flow between viruses and their hosts, phage genomes share a relatively high
168 proportion of their genome with their bacterial hosts (Supplementary Fig. 1A). This can limit
169 detection of viral sequences in metagenomes, because portions of the viral sequences will also
170 be present in bacterial genomes. Therefore, we decided to construct an alternative version of
171 Phanta's default database, in which prophage sequences, which are phage sequences that are
172 integrated into the bacterial chromosome, are "masked". This is accomplished by replacing the
173 prophage sequences with Ns in all bacterial genomes where they appear. Prophage sequences

174  were predicted using VIBRANT[42]. We anticipated that masking would further increase Phanta's
175  sensitivity to viral reads in simulated communities. Indeed, using the masked database reduced
176  the number of "ambiguous" read classifications - i.e., the number of reads that Kraken2 classified
177  to the "root" of the taxonomy tree. The vast majority of reads that were classified to the root using
178  the default database, but received a new classification after masking, were reclassified to the viral
179  domain (Supplementary Fig. 1B). This result demonstrates that: (1) shared sequences between
180  bacteria and viruses can indeed result in ambiguous read classification, and (2) this ambiguity
181  can be partially resolved by masking prophages in bacterial genomes. Importantly, masking does
182  not lead to over-detection of viruses; Phanta's final read abundance estimate for viruses remained
183  highly accurate (Supplementary Fig. 1C).

184

185  **Phanta improves the overall proportion of reads classified in shotgun metagenomes from**
186  **healthy adults**
187  Given the good performance of Phanta on simulated samples, we wished to assess whether
188  Phanta could improve viral identification in samples from healthy adults. We applied Phanta with
189  the default (no prophage masking) database on human gut metagenomes sampled from 245
190  healthy adults (age range 21-79, from Yachida *et al.*)[43]. In total, across 245 samples, the workflow
191  took ~60 minutes to run using 1 core, 16 threads, and 32GB memory. Given that Phanta
192  incorporates Kraken2 and Bracken, we were easily able to compare the workflow's performance
193  using Phanta's default database, compared to existing Kraken2/Bracken-compatible databases.
194  In particular, we compared against four existing databases (Table 1): the standard Kraken2
195  database[22,44] (May 2021), the Unified Human Gastrointestinal Genome (UHGG) collection[18] (July
196  2021), RefSeq Complete[39] (April 2022), and HumGut[19] (July 2021). Phanta's default database
197  was able to minimize the number of unclassified reads to 2% (Fig. 3A), and notably, it requires
198  ~97% less disk space than the most comprehensive database tested, RefSeq Complete (32GB
199  for Phanta, 1.2 TB for RefSeq Complete[39]).

200

201  **Phanta substantially increases viral identification in shotgun metagenomes**
202  In addition to maximizing classified reads, Phanta's default database led to the highest level of
203  viral identification, detecting 25-fold and 188-fold more viral sequences compared to RefSeq
204  Complete and the standard Kraken2 database, respectively (Fig. 3B; Supplementary Data File
205  2). Using Phanta, we now estimate that viral DNA constitutes 3-5% of the DNA in the human gut.
206  Taken together, Phanta improves read classification both by enabling the classification of
207  previously unclassified reads and by improving the recognition of viral sequences.

208

209  **Phanta outperforms standard assembly-based methods in identifying viruses in shotgun**
210  **metagenomes**
211  A current gold-standard workflow commonly used to identify viruses in shotgun metagenomes
212  involves assembling reads into contigs and labeling the likely viral contigs[42,45–49]. To compare
213  Phanta to this gold standard, we randomly selected 50 metagenomes from the healthy adult
214  cohort and ran a standard assembly workflow. In short, reads were assembled to contigs using
215  metaSPAdes[50], short/low-quality contigs were filtered using CheckV[51], and viral contigs were
216  identified using both VIBRANT[42] and VirSorter[45]. For each sample, the total set of viral contigs
217  from both methods was de-replicated to 95% ANI to calculate a number of viral species. Phanta

218    was able to identify a higher number of viral species than the assembly workflow in all samples,
219    with a median of 190 (IQR: 149-252) viral species per sample relative to 35 (IQR: 25-42) (Fig. 3C)
220    identified with assembly-based approaches. Of note, the vast majority of viral contigs predicted
221    by assembly were highly similar to genomes in the viral portion of Phanta's database
222    (Supplementary Fig. 2).

223

224    **There are twice as many viral as bacterial genomes in the human gut**
225    By default, Bracken calculates relative read abundance for each identified taxon - i.e., the fraction
226    of reads classified to it. This measurement serves as an estimation of the fraction of genomic
227    DNA belonging to each taxon, out of the total DNA in a sample. While this measurement is highly
228    valuable, an ecological perspective of a community requires understanding the proportions of
229    "individuals" in the community - i.e., relative taxonomic abundance[41]. Relative read abundance is
230    typically similar to taxonomic abundance in communities with similar genome lengths. However,
231    in mixed communities containing taxa with orders of magnitude differences in genome length, like
232    bacteria and viruses, relative read abundance is biased towards taxa with longer genomes (as
233    illustrated in Fig. 1B). Hence, Phanta calculates an estimation of relative taxonomic abundance
234    by correcting the relative read abundance by genome length. Using our relative taxonomic
235    abundance calculation, we estimate the ratio between copies of viral genomes to bacterial
236    genomes in the human gut to be ~2:1 (Figs. 4A-B; Supplementary Data File 2). Phanta also
237    reports several other normalizations - reads per million base pairs, reads per million reads, reads
238    per million base pairs per million reads (analogous to RPKM in transcriptomics) and genome
239    copies per million reads.

240

241    **High individuality of the human gut virome**
242    We further used the viral and bacterial profiles reported by Phanta to describe core differences
243    between the virome and bacteriome of healthy adults. We observed a higher between-sample
244    dissimilarity of the virome relative to the bacteriome in healthy adults (Fig. 4C). The high
245    dissimilarity of the virome between individuals points to a highly personalized virome, as has been
246    suggested previously[31,52–54]. Consistent with this result, individual viral species are skewed
247    towards lower prevalence than bacterial species (Fig. 4D). However, a number of lowly prevalent
248    viruses show high mean abundance across individuals, indicating that they are highly abundant
249    when present. As previously suggested, the prototypical crAssphage[55,56] (RefSeq ID 1211417)
250    was one of the most abundant viral species, although it was not among the most prevalent (Fig.
251    4D and Supplementary Table 1). Two of the most prevalent and abundant species were OTU-
252    66229 and OTU-72541. These phages are highly similar to the recently described Bacteroides
253    phages LoVEphage[37] and Hankyphage (p00)[57], respectively (Supplementary Fig. 3). The most
254    abundant and prevalent phage detected was Caudovirales OTU-21255, a temperate phage likely
255    of family Siphoviridae whose presumed host is *Bacteroides uniformis*. This species was found in
256    232/245 (95%) of individuals in this cohort of healthy adults, and comprises 1512 genomes in
257    Phanta's default database.

258

259    **Prevalent phages infect Bacteroides**
260    We next examined relationships between viral species prevalence and predicted host.
261    Bacteroides is the most commonly predicted host genus for viral species detected in the healthy

262  adult cohort. Specifically, it was the predicted host for 6.5% of detected viral species, compared
263  with 3% of species in Phanta's database, more than twice than expected. The dominance of
264  Bacteroides as a predicted host further increases among the more prevalent viral species (Fig.
265  4E).
266
267  **Temperate phages dominate the human gut phageome**
268  Phanta's default database includes estimates of virulence per species (see Methods), which we
269  used to determine the ratio between different phage lifestyles (virulent vs. temperate) in the
270  human gut. We observed that in the vast majority of samples temperate phages are dominant
271  with a median of 0.54 for the ratio of virulent/temperate species identified, and 0.55 for the
272  corresponding abundance ratio. Notably, more prevalent phages are skewed towards the
273  temperate lifestyle (Supplementary Fig. 4), potentially reflecting the ability of some temperate
274  phages to remain dormant in their hosts. Interestingly, the abundance of virulent phages in the
275  community, relative to temperate phages, is positively correlated with overall phage abundance
276  in the microbiome (Fig. 4F). This is consistent with the nature of virulent phages, whose active
277  replication increases their ratio relative to their bacterial host.
278
279  **Phanta performs well on virus-enriched metagenomes**
280  Viral enrichment, either through filtration or other approaches to achieve viral particle isolation, is
281  commonly used in viromics studies to enhance the detection of viral DNA in metagenomes[58].
282  Therefore, we wanted to test Phanta's performance in metagenomes originating from virus-
283  enriched samples. We applied Phanta to paired bulk and virus-enriched shotgun metagenomes
284  from infants (Supplementary Data Files 3-6; source data: Liang *et al.*[59]). We first tested the
285  performance of Phanta on the virus-enriched samples by correlating the viral-like particle counts
286  (from Supplementary Table 2 in [59]) to the number of viral species identified (i.e., viral species
287  richness) by various assembly or classification methods. Phanta-based richness was the most
288  strongly correlated with VLP counts (Fig. 5A).
289
290  **Viral profiles from bulk and virus-enriched metagenomes overlap, but complement each**
291  **other**
292  Given its high sensitivity, we hypothesized that Phanta would detect a comparable number of viral
293  species in bulk metagenomes as in virus-enriched metagenomes. Indeed, the number of species
294  detected was similar in paired bulk and virus-enriched metagenomes (Fig. 5B). We further tested
295  whether the bulk and viral-enriched metagenomes provide a similar profile of the viral community
296  by examining pairs of bulk and viral-enriched metagenomes from the same sample. First, we
297  examined 10 pairs of metagenomes with relatively deep sequencing of the bulk metagenomes
298  (range of 150bp paired-end reads 8.6M - 13.3M; median = 9.3M). Species present in bulk
299  metagenomes captured a median of 94% of the viral abundance in virus-enriched metagenomes
300  (Fig. 5C). The variance in this quantity is mostly explained by the sequencing depth of the bulk
301  metagenomes (Supplementary Fig. 5). To complement this analysis, we examined 10 pairs of
302  metagenomes with highly successful viral enrichment (see Methods). Species present in virus-
303  enriched metagenomes captured a median of 69% of the viral abundance in bulk metagenomes
304  (Fig. 5D). Those differences are expected as shotgun metagenomes can capture viruses that did
305  not enrich in the VLP enrichment process for a variety of reasons, technical or biological[31]. For

306     example, prophages lack viral-like particles, and are therefore more likely to be captured by bulk
307     metagenomes. Given the inclusive nature of bulk metagenomes, they capture more viral species
308     per total number of viral reads (Supplementary Fig. 6A), whereas viral-enriched metagenomes
309     capture more viral species per total number of metagenome reads (Supplementary Fig. 6B). With
310     the ability to identify prophages in bulk metagenomes, we hypothesized that the fraction of
311     temperate phages would be higher in virus-enriched metagenomes. Indeed, we observed a 3-fold
312     higher virulent/temperate abundance ratio in virus-enriched metagenomes relative to bulk (Fig.
313     5E; Supplementary Fig. 6C).

314

315     **Phanta is highly effective for simultaneous quantification of phages and their hosts from**
316     **a single metagenomics experiment**
317     One advantage of using Phanta to profile bulk metagenomes, as opposed to virus-enriched
318     metagenomes, is the ability to examine phages and their hosts simultaneously and from a single
319     dataset, instead of two separately generated datasets. Using a Phanta-based analysis of the bulk
320     metagenomic dataset from Liang *et al.*[59] investigating the impact of diet on the infant gut, we found
321     that Bifidobacterium and its phages are ~2-fold more abundant in breastfed infants relative to
322     formula-fed or infants that were fed with a mixed breast milk and formula diet (Fig. 5F). This
323     observation, although expected, demonstrates the power of Phanta to simultaneously identify
324     phages and their bacterial hosts and to associate them with known traits.

325

326     **Phanta accurately identifies and quantifies human-infecting viruses**
327     Lastly, we wished to test the ability of Phanta to accurately identify human-infecting viruses in
328     metagenomes. Liang *et al.* were able to identify viruses in the family of Adenoviridae using qPCR
329     from their infant stool samples[59]. Phanta identified 5 samples with the mastadenovirus C species,
330     with almost perfect correlation between the estimation of genome copies per uL using qPCR and
331     Phanta's estimation of genome copies per million reads (Fig. 5G). Phanta was able to identify
332     Adenoviruses in bulk shotgun metagenomic samples with as low as 88 copies/uL in qPCR and
333     successfully identified all samples with >550 copies/uL. Phanta demonstrated higher sensitivity
334     in identifying Adenoviruses relative to using assembly-based methods (Fig. 5H), which only
335     detected Adenoviruses in samples that had >20,000 copies/uL by qPCR. Of note, we used the
336     assembled contigs to confirm that Phanta successfully identified the right Adenovirus species, by
337     aligning the contigs to all Adenovirus genomes from RefSeq.

338

339     **Discussion**

340

341     A major goal of microbiome studies is to identify microbial features associated with traits of
342     interest, such as phenotypes, lifestyle factors, and health status. In an ideal world, organisms
343     from all domains could be accurately quantified in a single experiment. The first step in achieving
344     this goal is to profile microbial communities - i.e., to determine their composition from sequencing
345     data. Although shotgun metagenomes capture both prokaryotes and viruses, profiling the viral
346     fraction of microbial communities has historically presented a greater challenge and has required
347     specially tailored methods. For example, popular reference-based methods have allowed
348     accurate profiling of prokaryotes from metagenomes[24] without being able to accurately capture
349     viruses due to the historical lack of comprehensive reference databases of viral genomes[27].

350 Because of these limitations, profiling viruses has required additional orthogonal analyses, based
351 on assembling metagenomic reads and identifying viral genomes *de novo*[25,26]. In addition, due to
352 the relatively low abundance of viral sequences in bulk metagenomes, it has been common to
353 conduct an entirely separate experiment to profile the virome by enriching for viral sequences
354 prior to making sequencing libraries[58].

356 With the recent development of much more comprehensive databases of viral genomes[27,31–34],
357 deeply sequenced bulk metagenomes, and fast and accurate read classifiers[22,38], technical and
358 experimental advances have converged to make it possible to integrate prokaryotic and viral
359 profiling. By harnessing the latest developments, Phanta enables simultaneous profiling of
360 bacteriophages and their prokaryotic hosts, in a single experiment and with a single analysis. This
361 simultaneous profiling has several advantages. First, it reduces the need to sequence both viral-
362 enriched and bulk metagenomes, thus reducing research time and costs, in addition to eliminating
363 technical differences between two separate experiments. Second, it bypasses the need to use
364 separate computational workflows to profile prokaryotes and viruses. Lastly, and most
365 importantly, it allows the study of cross-domain interactions between phages and their hosts,
366 either in novel datasets, or in the wealth of metagenomic datasets that are already publicly
367 available.

369 Although Phanta can be applied with different databases, Phanta's default database was
370 constructed with the human gut in mind. For decades, the viral portion of the human gut was
371 mostly unknown, and considered as "dark matter"[60,61]. There is still much to learn, with some basic
372 discoveries occurring only in the past few years. For example, the first representative of one of
373 the most abundant bacteriophage clades - crAss-like viruses - was discovered only in 2014[56].
374 Similar examples, such as the highly prevalent Hankyphage (p00)[57] and LoVEphage[37], were
375 discovered only in 2018 and 2021, respectively. We anticipate that Phanta, when applied with its
376 default database, will allow similar key discoveries to be made. In this study alone, we were able
377 to estimate a ~2:1 ratio of viruses to bacteria in the human gut, determine that temperate and
378 Bacteroides-infecting phages dominate the gut phageome, and demonstrate a high interindividual
379 variability of the gut virome, as compared to the bacteriome. These and other core principles can
380 serve as a springboard for more extensive discovery, such that "gut microbiome" will no longer
381 be publicly synonymous with "gut bacteria," but rather understood as a complex community with
382 many types of interacting members.

384 Importantly, Phanta was developed with careful attention to the risk of spurious discovery, as
385 read-based classifiers are frequently known to make mistakes, and thus to identify false positive
386 taxa[40]. As described, to mitigate false classification we increase classification confidence and filter
387 out species with low genome coverage, an idea that was previously described in the
388 implementation of KrakenUniq[40]. Of course, these decisions come with potential costs. For
389 example, increasing the required confidence of classification may lead reads from some species
390 to all classify at higher taxonomic ranks during the Kraken2 step of the workflow. In such a
391 scenario, the sensitivity of viral identification would be decreased, since during Bracken, classified
392 reads are only redistributed to species that initially received some direct classifications. Similarly,
393 requiring a certain genome coverage reduces the probability of identifying lowly abundant species

394    with long genomes. However, all the relevant parameters of Phanta are user-adjustable, and
395    using our simulations we were able to show that a combination of minor increments in both
396    thresholds is sufficient to reduce most of the noise with a very small cost to signal (Figs. 2D-E).
397

398    More broadly, Phanta offers a flexible setup that can be modified according to the user's analysis
399    goals and main concerns. If a user aims to minimize false negatives, i.e. to increase the probability
400    of identifying all species while allowing a substantial increase in false positives, the user can
401    decrease (1) the confidence cutoff, (2) the coverage requirement, and (3) the minimal number of
402    reads directly classified to a species for it to receive an abundance estimate. On the other hand,
403    if a user wishes to minimize false positives while taking the risk of decreasing true positives, the
404    user can increase these three parameters. Phanta also provides an alternative database to the
405    default, in which predicted prophages in the HumGut genomes were masked. This masked
406    database can be used to increase the likelihood of identifying prophages. In addition to
407    parameters and database choice, the characteristics of a sequencing experiment can impact the
408    power of identification by Phanta. Although we did find high agreement between viral-enriched
409    and bulk shotgun metagenomes (Fig. 5C), enriching the library for viral particles would be
410    recommended if a researcher (i) prioritizes identification of viruses that are particles over
411    prophages, (ii) is not focused on determining cross-domain interactions, and (iii) is limited by the
412    possible depth of metagenomic sequencing. Conversely, bulk metagenomic analysis allows users
413    to: (i) profile prophages in addition to virulent phages, (ii) avoid potential biases introduced by the
414    process of isolating viral particles, and (iii) identify cross-domain interactions, both within and
415    across samples. Given the low and rapidly decreasing costs of shotgun sequencing, and our
416    findings that bulk metagenomes of fairly standard depth allow for comparable virus identification
417    to viral particle-enriched fractions, we anticipate that many researchers may opt to enhance their
418    standard analyses of bulk metagenomes by applying Phanta.
419

420    While Phanta enhances the knowledge that can be gained about viruses from bulk metagenomes,
421    it has several limitations. First, while Phanta has high sensitivity, using a reference-based method
422    restricts identifications to the genomes in the database, and thus limits resolution. For example,
423    Phanta's default database is biased toward dsDNA viruses identified in the human gut. Similarly,
424    while Phanta does include some eukaryotes in its default database, our knowledge of this domain
425    in the gut is still limited; this is, in part, due to limitations in reference databases for protists,
426    amoeba, helminths and fungi. Improvements in eukaryotic reference databases should enhance
427    eukaryote classification in the coming years. Second, extending Phanta to characterize the virome
428    in other human microbiomes, such as the skin or vaginal microbiome, may require curation of
429    additional metagenome-derived virome databases generated from these niches. Furthermore,
430    classifying short reads to reference genomes is challenging when reads originate from genomic
431    regions that are conserved between species. Moreover, the usage of $k$-mer-based methods,
432    although fast and computationally efficient, does not provide information required for aligning
433    reads to a specific region in the genome, and thus does not allow investigation of genome
434    variation. Finally, viral taxonomy is not as well-defined as prokaryotic taxonomy, and thus Phanta
435    cannot currently provide specific named designations to many viral species, beyond family- or
436    order-level assignments. We anticipate that as knowledge of the virome increases, this challenge
437    will begin to be addressed.

438

439    Despite these limitations, Phanta is benchmarked, easy to use, carefully tuned to limit false
440    positives, and able to provide simultaneous profiling of various domains from a single experiment.
441    These advantages suggest that Phanta will help accelerate the study of the virome in human gut
442    microbiomes, as well as illuminate cross-domain interactions in this niche. Phanta enables much
443    higher resolution of the viral portion in a human gut sample when analyzing a bulk metagenome
444    relative to current approaches or databases, and thus it promises to provide exciting insights when
445    applied to the tens of thousands of human gut metagenomes that have already been sequenced,
446    to date. We expect that Phanta will be both: (1) used to re-analyze publicly available data, and
447    (2) taken into account when planning new experiments. Overall, Phanta lowers the barrier to virus-
448    inclusive studies of the gut microbiome, and we expect that its application will confidently identify
449    numerous novel associations between viruses, prokaryotes, and human traits.
450
451

452 **Online Methods**

453

454 **Constructing a comprehensive, taxonomy-aware, domain-inclusive database of human gut**
455 **microbes**

456

457 Phanta's default database was constructed to be compatible with the Kraken2/Bracken tools[22,38].
458 Therefore, its construction required curating: 1) a large collection of genomes, and 2) taxonomy
459 files placing each genome within a tree of named nodes.

460

461 The viral genomes within the database were sourced from: 1) the recently published human gut-
462 focused MGV catalog (available at (https://portal.nersc.gov/MGV/)[27] and 2) RefSeq[21], the
463 database of reference genomes maintained by NCBI (MM/YY of download: 02/22).

464

465 After downloading the viral genomes, the viral taxonomy tree was constructed. The first step was
466 to download the complete NCBI taxonomy using the kraken2-build --download-taxonomy utility.
467 Next, branches of the taxonomy were pruned so that only the branches leading to the RefSeq
468 viral genomes remained.

469

470 After providing taxonomic assignments to RefSeq genomes, assignments were provided to the
471 MGV genomes. The first step in doing so was to group the MGV genomes into the 54,118 ANI-
472 based species specified by the MGV paper[27]. Each of these species came with a designated
473 "species representative genome" that was chosen based on features such as circularity and
474 length. Code on the MGV GitHub page
475 (https://github.com/snayfach/MGV/tree/master/aai_cluster) was then used to cluster species into
476 genera based on amino acid identity (AAI) and gene sharing between the representative
477 genomes.

478

479 To avoid species duplications between MGV and RefSeq viruses, and to provide a full NCBI
480 taxonomy for MGV genomes where available, average nucleotide identity was calculated between
481 all of the 54,118 species representative genomes in MGV and all the RefSeq viral genomes using
482 fastANI[62]. In cases where an MGV species representative genome had > 95% ANI to a RefSeq
483 viral genome, all of the genomes in the relevant MGV species were re-assigned to RefSeq, i.e.,
484 designated as strains of the RefSeq viral genome.

485

486 To determine where each AAI-based MGV genus fit into the NCBI taxonomy, we utilized a file
487 from the MGV website (https://portal.nersc.gov/MGV/) that provides - when possible - NCBI-
488 recognized taxonomic annotations for each genome at the genus, family, and/or order levels,
489 based on amino acid alignments to a protein database[27]. We used this information to remove
490 some of the AAI-based genera and re-assign their contained species to the relevant NCBI-
491 recognized genus. Specifically, for each AAI-based genus, we calculated the percentage of
492 species representative genomes within the genus that had an NCBI genus annotation provided.
493 If this percentage was greater than 50%, and the NCBI genus annotation was consistent for >
494 90% of the species representative genomes with annotations, the AAI-based genus was removed
495 and all of its species were re-assigned to the NCBI genus.

496
497     The remaining AAI-based genera were then assigned as direct descendants of the lowest
498     possible NCBI-recognized taxonomic level, by iterating a variant of the strategy described above.
499     More specifically, starting with family: if > 50% of the species representative genomes within a
500     given AAI-based genus had an NCBI family annotation, and the NCBI family annotation was
501     consistent for > 90% of the species representative genomes with annotations, the AAI-based
502     genus was assigned as a direct descendant of the relevant NCBI family. The remaining AAI-
503     based genera (i.e., those without a family assignment) were then assigned to an order - when
504     possible - in the same manner. All of the AAI-based genera without an order assignment were
505     assigned as direct descendants of the superkingdom of Viruses.
506
507     The prokaryotic genomes within Phanta's database were sourced from HumGut, a recently
508     published human gut-focused catalog of prevalent bacterial and archaeal genomes[19]. The
509     HumGut catalog was in turn sourced from both the Unified Human Gastrointestinal Genome
510     (UHGG) collection[18] and RefSeq[21]. An NCBI-compatible taxonomy file for the HumGut genomes
511     was downloaded directly from the HumGut website (http://arken.nmbu.no/~larssn/humgut/). The
512     branches of the NCBI taxonomy leading to the human genome were also included in this
513     taxonomy file and thus we also included the human genome (hg38) in our database.
514
515     We sourced fungal genomes from RefSeq and common contaminant sequences from the Core
516     UniVec database using the kraken2-build download-library command provided by the Kraken2
517     developers (MM/YY of download: 02/22). The relevant branches of the NCBI taxonomy were then
518     obtained in the same way that they were obtained for the RefSeq viral genomes (i.e., by "pruning"
519     the full NCBI taxonomy, please see above).
520
521     Finally, the constructed taxonomy files for each portion of the database were concatenated, and
522     a Kraken2/Bracken-compatible database was built using the commands provided on the Github
523     sites (https://github.com/DerrickWood/kraken2; https://github.com/jenniferlu717/Bracken).
524
525     **Masking prophages in prokaryotic genomes**
526
527     An alternative version of Phanta's default database was also created, in which predicted
528     prophages were masked (i.e., replaced with Ns) within all the prokaryotic genomes from HumGut.
529     VIBRANT (v1.2.1)[42] was used to predict prophages within the HumGut genomes. Prophage
530     coordinates were extracted and masking was conducted using the bedtools utility
531     MaskFastaFromBed[63]. All of the analyses in this paper were conducted using the unmasked
532     version of the database, except where explicitly noted otherwise.
533
534     **Workflow implementation**
535
536     Phanta was implemented using the workflow management system Snakemake. Core scripts are
537     written in Python, bash, and R. A step-by-step tutorial detailing workflow installation and usage is
538     provided on the main page of the Phanta GitHub (https://github.com/bhattlab/phanta). Briefly,
539     after cloning the GitHub repository to their system, users should: 1) download the desired

540  database - default (unmasked) or masked - via the command line, 2) make slight edits to a
541  configuration file, and 3) execute the provided Snakemake command on the command line, within
542  the appropriate conda environment that is fully specified in a provided yaml file. As detailed in the
543  GitHub tutorial, the repository also provides a test data set that can be used to verify that the
544  workflow was installed correctly.
545
546  **Classification of metagenomic reads to taxa**
547
548  The first step of the Phanta workflow is classification of metagenomic reads in each sample
549  against the desired database of genomes (default/unmasked or masked, see above).
550  Classification is conducted using the Kraken2 tool (currently v2.1.2)[22], which classifies reads
551  using a *k*-mer-based approach. More specifically, to classify each read, Kraken2 slides along the
552  read length, computes a "minimizer" (i.e., compact version) of each *k*-mer, and looks up where
553  the minimizer maps in the genome database. After all the minimizers in the read have been looked
554  up, Kraken2 classifies the read to the lowest taxonomic level possible, considering the user's
555  preference for the confidence in the assignment (supplied via the --confidence parameter to
556  Kraken2). By default, Phanta supplies a confidence of 0.1 to Kraken2, but this value can be
557  adjusted by the user in the Snakemake configuration file. This parameter ranges from 0 to 1 and
558  essentially specifies a certain fraction of a read's *k*-mers to be mapped to a given taxon, in order
559  for Kraken2 to make that classification. E.g., 0.1 = 10%.
560
561  Phanta also makes use of the --report-minimizer-data parameter available in Kraken2 v2.1.2, that
562  is based on ideas from KrakenUniq[40]. Providing this parameter modifies the standard Kraken2
563  output to report an additional data point for each taxon, specifically: how many unique minimizers
564  in the genomes of this taxon are covered by read sequences?
565
566  **Filtering of false positive species after classification**
567
568  Phanta filters likely false positive species from each sample after the initial classification step and
569  before species-level abundance estimates are calculated (Figure 1). This filtering step makes use
570  of the minimizer data reported by Kraken2 during the classification step (described above, in the
571  section "Classification of metagenomic reads to taxa").
572
573  Specifically, a proxy for genome coverage is calculated for each genome of each species
574  identified during classification. This proxy is calculated by dividing: 1) the reported number of
575  unique minimizers in the genome that are covered by read sequences, by 2) the total number of
576  unique minimizers contained in the genome. The denominator of this fraction is not reported in
577  the Kraken2 output, but is obtained by Phanta from an "inspect.out" file contained within the
578  genome database (originally generated using the kraken2-inspect functionality).
579
580  Bacterial and viral species are marked as false positives and filtered out if none of their strain-
581  level genomes have a calculated coverage above a user-specified threshold. Suggested
582  thresholds are provided in the Snakemake configuration file (0.01 for bacterial species; 0.1 for
583  viral species). These suggested thresholds were chosen because they yielded a high signal-to-

584  noise ratio in identified species when tested on the mixed simulated metagenomes described
585  below.

587  Users can also require that the numerator of the fraction above (i.e., the number of unique
588  minimizers covered by reads) be higher than a specified threshold for at least one strain-level
589  genome of each "true positive" species. In other words, it is possible to specify that a high
590  calculated genome coverage will not "count" unless the number of unique minimizers is higher
591  than a specific value (e.g., > 300 unique minimizers) for at least one strain-level genome. By
592  default, this option is not utilized by Phanta but can be implemented by the user by making use
593  of the minimizer_thresh_viral and minimizer_thresh_bacterial parameters in the Snakemake
594  configuration file.

596  **Species abundance estimation and correction for genome length**

598  After species are filtered from the Kraken2 output, abundances of the remaining species are
599  estimated using the Kraken2-compatible tool Bracken (currently v2.7)[38]. Bracken estimates
600  species-level abundances by redistributing all classified reads to the species level.

602  Of note, Bracken accepts a threshold parameter that specifies one last filter for false positive
603  species - how many sample reads must have been classified to a species during Kraken2
604  classification for Bracken to estimate its abundance? By default, Phanta specifies this threshold
605  as 10 reads - the accepted standard for running Bracken - but this number can be adjusted by
606  the user through the filter_thresh argument in the Snakemake configuration file.

608  We also utilize Bracken output to calculate relative taxonomic abundance estimates for each
609  species by considering genome length. Specifically, the abundance estimate for each species is
610  scaled by the median length of the genomes under the species. Additional normalizations are
611  also provided in this corrected output file, such as reads per million reads per million base pairs
612  (analogous to RPKM in transcriptomics), copies per million reads, and more.

614  **Provided post-processing scripts**

616  There are three main post-processing scripts in the Phanta GitHub.

618  The first calculates "lifestyle statistics" for the viral community in each metagenome (e.g., ratio of
619  virulent:temperate viruses), based on lifestyle predictions for viral species that are provided in
620  Phanta's default database. Lifestyle predictions for species from MGV were obtained from the
621  mgv_contig_info file provided in the MGV database[27]. These predictions were calculated using
622  BACPHLIP[64] and we used the same tool (v0.9.6) to make lifestyle predictions for viral species
623  from RefSeq. Throughout the manuscript, viruses with a BACPHLIP-predicted virulence score
624  above 0.5 were considered virulent; others were considered temperate.

626   The second collapses viral abundances in each sample by predicted host, based on provided
627   host predictions for viral species in Phanta's default database. Host predictions were made using
628   iPHoP[65].

630   The third correlates the abundances of bacterial and viral species in each sample. This cross-
631   kingdom correlation is done by fastspar[66,67]- a method designed to correlate compositional data.

633   Also provided are post-processing scripts to filter or sum abundance tables (counts, relative read
634   abundances, or relative taxonomic abundances) to a desired taxonomic rank (e.g., species or
635   genus).

637   **Simulating mixed metagenomes**

639   10 mixed metagenomes (each containing ~6.5M paired-end 150bp reads) were simulated using
640   CAMISIM (v1.3)[68]. These simulated metagenomes were used to generate the data in Figure 2.
641   Each simulated metagenome consisted of: 1) 95% prokaryotic reads from 300 randomly chosen
642   genomes from the HumGut catalog, and 2) 5% viral reads from 50 randomly chosen genomes
643   from the MGV catalog.

645   **Download and processing of publicly available, short-read human gut metagenomes**

647   245 shotgun gut metagenomes from healthy human adults in a Japanese cohort were
648   downloaded from SRA (accession DRP004793 - Yachida *et al.*[43]). Shotgun gut metagenomes
649   from infants were also downloaded from SRA (accession PRJNA524703 - Liang *et al.*[59]). The full
650   list of downloaded samples, along with accession numbers, is available within Supplementary
651   Table 2.

653   Following download, each metagenome was preprocessed as follows. First, reads that exactly
654   matched each other (PCR duplicates) were removed using hts_SuperDeduper (v1.2.0). Next,
655   TrimGalore (v0.6.5 healthy adults; v0.6.7 infants) was used to: 1) trim low-quality bases (Phred
656   score < 30) from the ends of reads, and 2) discard reads with a final length of < 60bp. Human
657   reads were then removed using BWA alignment against the human genome (GRCh37). Initial
658   and final quality checks were performed using MultiQC (v1.7 healthy adults; v1.11 infants).

660   All results from applying Phanta to these metagenomes were obtained using Phanta's default
661   database and parameters, except where explicitly noted otherwise (i.e., varied databases were
662   tested in Figures 3A and 3B). Note also that for the infant cohort, the database file required for
663   running Bracken was slightly modified from default (adjusted for 120bp reads rather than 150bp,
664   following the instructions on the Bracken GitHub).

666   Separate from running Phanta, a subset of these metagenomes was assembled into contigs and
667   scaffolds using metaSPADES[50] version 3.15. Specifically, the following metagenomes were
668   assembled: 1) fifty randomly selected metagenomes from the healthy adult cohort, and 2) all bulk

669　metagenomes from the "four-month" subgroup of the infant cohort. The specific metagenomes
670　that were successfully assembled are indicated in Supplementary Table 2.
671
672　**Assembly-based method for identifying viral species in healthy adult gut metagenomes**
673
674　To generate the results in Fig. 3C, the 50 assembled healthy adult gut metagenomes were run
675　through two standard methods for phage identification from metagenomic assemblies. The first
676　method, VIBRANT, uses a hybrid machine learning and protein similarity approach to identify viral
677　signatures[42]. The second method, VirSorter, predicts protein-coding genes in assembled DNA
678　sequences and assesses their similarity to known viral proteins[45].
679
680　VIBRANT (v.1.2.0) was run on assembled scaffolds and the quality and completeness of identified
681　phages were estimated by CheckV (v.0.7.0)[51] using database v0.6. Low-quality phage scaffolds
682　were filtered out.
683
684　A similar procedure was performed using VirSorter (v1.0.6, downloaded in February 2018), where
685　phage contigs were classified as category 1, 2, or 3 depending on confidence level. Category 3
686　predictions were filtered out before running CheckV.
687
688　Finally, dRep (v3.2.2)[69] was applied to the combined set of quality-filtered phage contigs predicted
689　by VIBRANT+VirSorter in each sample to extract a unique set of phage genomes based on an
690　ANI threshold of 0.95 and coverage threshold of 0.5. fastANI was applied for secondary clustering
691　and genome filters included a minimum length of 1000 bp, an N50 weight of 0, and a size weight
692　of 1.
693
694　The full list of parameters utilized with the "drep dereplicate" utility was: *-sa 0.95 --S_algorithm*
695　*fastANI -nc .5 -l 1000 -N50W 0 -sizeW 1 --ignoreGenomeQuality --clusterAlg single*
696
697　**Assembly-based method for identifying Adenoviruses in stool samples**
698　The assembled bulk metagenomes from the "four-month" subgroup of the infant cohort were used
699　to calculate the column labeled "Assembly" in Fig. 5H. FastANI was used to calculate average
700　nucleotide identity between all assembled contigs in each sample and 1801 *Adenoviridae*
701　genomes available in NCBI (retrieved by *datasets download genome taxon Adenoviridae*). Each
702　contig with ANI score >=95% to at least one *Adenoviridae* genome was counted as an Adenovirus.
703
704　**Calculation of dissimilarities between metagenomes**
705
706　Bray-Curtis and Jaccard distances were calculated using the R package vegan, version 2.5-7.
707
708　**Choosing pairs of metagenomes for overlap analysis**
709
710　For the analyses in Figs. 5C-5D, we wanted to determine how well each type of metagenome
711　could represent the information in the other, excluding samples with low sequencing depth of the
712　bulk metagenomes, or low enrichment of virus-enriched metagenomes.

713

714 For the analysis in Fig. 5C, we chose pairs of metagenomes with decent viral enrichment and
715 deeply sequenced bulk metagenomes. Specifically: (1) We identified the top 50% of samples
716 based on the percent of reads that Phanta assigned to viruses in the virus-enriched
717 metagenomes; (2) Of these, we selected 10 samples whose paired bulk metagenomes were the
718 most deeply sequenced..

719

720 For the analysis in Fig. 5D, we chose pairs of metagenomes with decent bulk sequencing depth
721 and highly successful viral enrichment. Specifically: (1) We identified the top 50% of samples
722 based on the sequencing depth of the bulk metagenomes; (2) Of these, we selected 10 samples
723 with the highest percent of reads that Phanta assigned to viruses in the virus-enriched
724 metagenomes.

725

726 **Determination of size and number of genomes in each Kraken2/Bracken-compatible**
727 **database tested**
728 To determine the size of each Kraken2/Bracken-compatible database tested (listed in Table 1),
729 we summed the sizes of the following files and rounded to the nearest GB: hash.k2d, opts.k2d,
730 taxo.k2d, seqid2taxid.map, database150mers.kmer_distrib. These are the files necessary for
731 running Kraken2 and Bracken. We obtained the number of prokaryotic and viral genomes in each
732 database that we did not construct directly from the relevant publications: Wright *et al.*, 2022 (for
733 Standard Kraken2 and RefSeq Complete)[39]; Almeida *et al.*, 2021 (for UHGG)[18]; Hiseni *et al.*, 2021
734 (for HumGut)[19].

735

**Data and code availability**

Phanta is publicly available at https://github.com/bhattlab/phanta with a detailed tutorial describing installation and usage. Accession numbers of all publicly available metagenomes used for analysis are provided in Supplementary Table 2. Workflows used for preprocessing and assembly were used in this manuscript and are available at: https://github.com/bhattlab/bhattlab_workflows.

756

| Database | Size (GB) | Prokaryotic genomes* | Viral genomes* | Median classification time (sec)** | Median % classified reads** |
|---|---|---|---|---|---|
| Standard Kraken2 | 51 | 21,920 | 10,489 | 491 | 56.31 |
| RefSeq Complete | 1,192 | 215,725 | 10,863 | 4,889 | 93.61 |
| UHGG | 16 | 4,644 | 0 | 548 | 88.75 |
| HumGut | 26 | 30,691 | 0 | 549 | 97.57 |
| Phanta | 32 | 30,691 | 201,305 | 544 | 98.07 |

757 *Numbers were obtained from the original papers, see Methods.
758 **Classification times and percentages of classified reads were determined by conducting
759 Kraken2 classification of five random samples from the healthy human cohort from Figure 3,
760 and calculating median classification times and percentages across the samples.
761
762 **Table 1.** Characteristics of the different Kraken2/Bracken-compatible databases tested in this
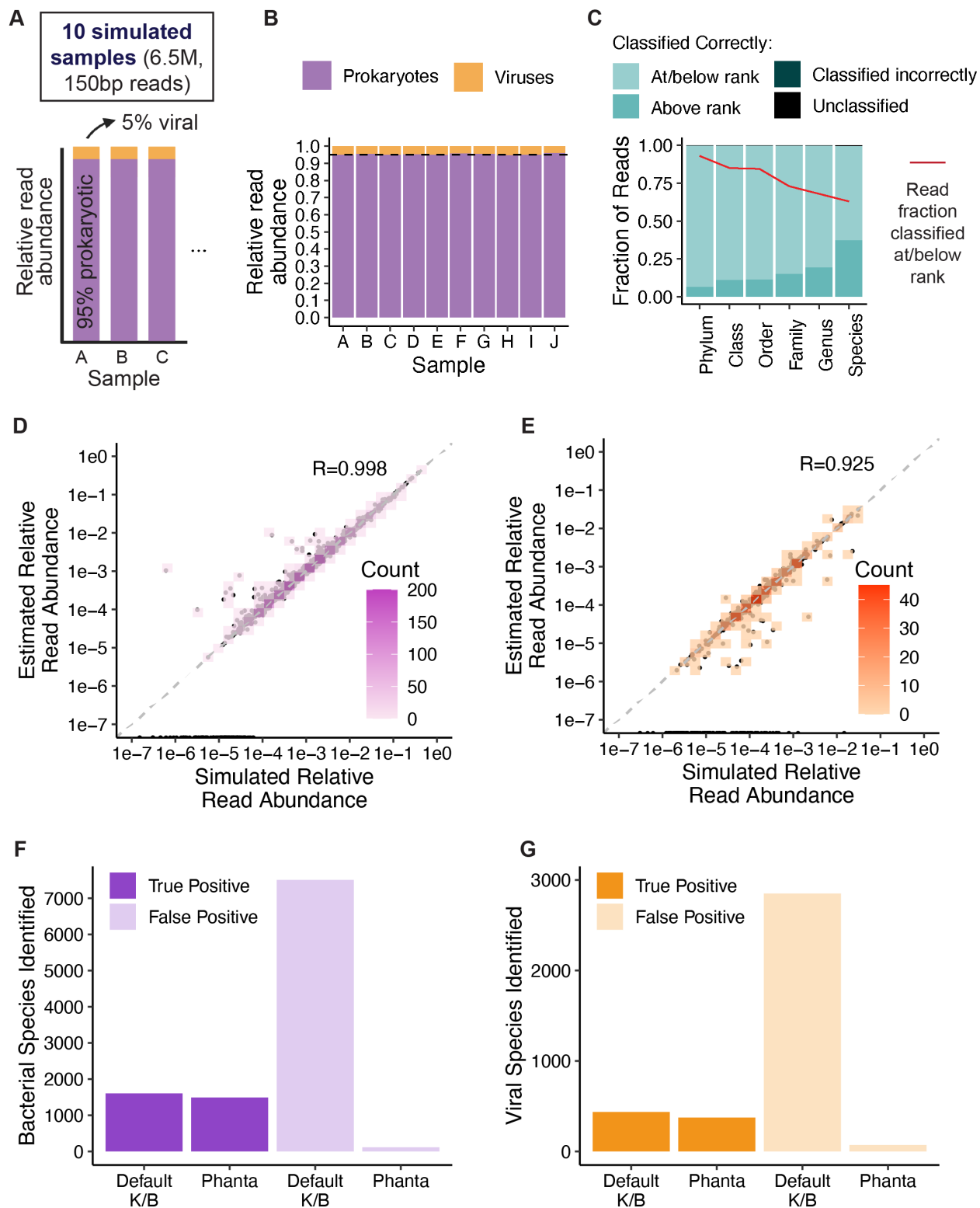763 study.
764

765

766 **Figure 1. Overview of Phanta's comprehensive, virus-inclusive metagenomic**
767 **annotation workflow.**
768 **(A)** Phanta's workflow. First, reads from each sample are classified against a
769 comprehensive, virus-inclusive database of genomes from the human gut. Reads are
770 classified to the lowest possible taxonomic level. After classification, genome coverage is
771 estimated for each detected species in each sample. Species with low estimated genome
772 coverage are filtered out to prevent false positive identifications. Next, Phanta quantifies
773 the abundances of the remaining species in each sample. Reads originally classified
774 above the species level (for example to the genus or family level) are redistributed
775 downwards. Then, two types of abundance are calculated: (1) relative read abundance,
776 which normalizes species-level read counts by read depth, and (2) relative taxonomic
777 abundance (see panel B). Post-processing scripts are provided to determine cross-
778 domain relationships.
779 **(B)** Motivation behind Phanta's provided correction of relative read abundance of relative
780 taxonomic abundance. Shown here is a simple gut microbial community with a 1:1 ratio
781 between bacteria and viruses (one bacterium of species *E. coli*; one virus of species T4).
782 Even if *E. coli* and T4 genomes are equally covered by reads in a shotgun metagenome,
783 the dramatic difference between their genome lengths will inflate the ratio of bacteria to
784 viruses, if relative read abundance is used as the metric. By contrast, relative taxonomic
785 abundance, which corrects for genome length, accurately captures the 1:1 ratio of these
786 species.
787

788

789

790

791

792 **Figure 2. Evaluation of Phanta's performance using simulated metagenomes.**
793 **(A)** Composition of simulated metagenomes. Results in (B) - (G) were obtained by
794 applying Phanta to these simulated metagenomes while using Phanta's default database
795 and parameters.
796 **(B)** Accuracy of Phanta's final estimates of relative read abundance at the domain level.
797 The dashed line indicates the true relative read abundance of prokaryotes.
798 **(C)** Accuracy of Phanta's classification step at each taxonomic rank. For each rank, the
799 two shades of blue represent reads that were classified to a lineage that included the
800 correct value of the rank. Specifically, light blue shading indicates the median fraction of
801 reads (across simulated samples) that were classified correctly at or below the rank - e.g.
802 for family, they were classified either to the correct family, or to the correct genus/species,
803 which is even more specific than the correct family. Darker blue shading indicates the
804 median fraction of reads that were classified correctly above the rank - e.g. for family,
805 they were classified to the correct order or phylum, which is less specific than the correct
806 family but still accurate. The dark green and black portions of each bar represent reads
807 that were either: (i) classified to a lineage that did not include the correct value of the rank,
808 or (ii) unclassified, respectively. The red line indicates the median fraction of classified
809 reads that were classified at or below each rank (e.g., what fraction of reads were
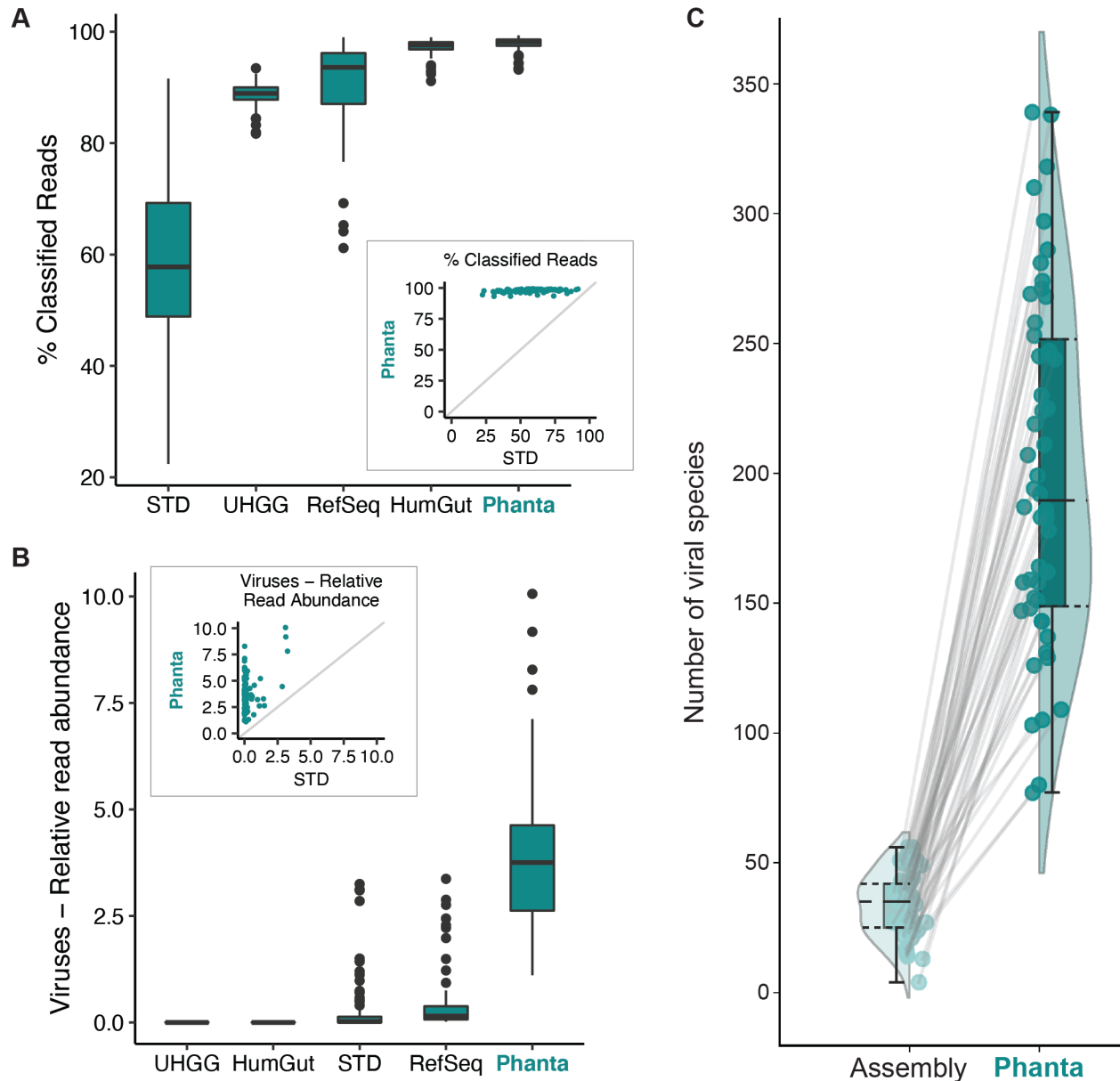810 classified at or below the family level).
811 **(D)** Accuracy of Phanta's final estimates of relative read abundance for 1,606 bacterial
812 species, across all simulated samples. The dashed line is the x=y diagonal. Each dot
813 represents one bacterial species in one simulated sample. The x-axis is the simulated
814 abundance in the sample, and the y-axis is the abundance estimated by Phanta. The *R*
815 value indicates Pearson's correlation coefficient, considering all the dots, i.e. all bacterial
816 species in all simulated samples. Colors of the overlaid boxes represent numbers (counts)
817 of dots.
818 **(E)** Same as (D), for 500 viral species.
819 **(F)** Signal-to-noise ratio of bacterial species identified by Phanta vs. the Kraken2/Bracken
820 workflow, using default parameters for both workflows and using Phanta's default
821 database as the reference database.
822 **(G)** Same as (F), for viral species.
823

**Figure 3. Evaluation of Phanta's performance using shotgun gut metagenomes from 245 healthy human adults.** Metagenomes sourced from Yachida et al.[43]

**(A)** Percentage of sample reads that could be classified during Phanta's initial classification step using Phanta's default parameters and a variety of Kraken2/Bracken-compatible databases. Boxplots display the percentage distribution across the set of metagenomes. Database abbreviations: STD = standard Kraken2[44], UHGG = Unified Human Gastrointestinal Genome Collection[18], RefSeq = RefSeq Complete v205[39], HumGut = HumGut[19], Phanta = Phanta's default database. The insert shows the same information as the boxplots for STD and Phanta.

**(B)** Similar to (A) but comparing the relative read abundance of viruses after Phanta's filtering and abundance estimation steps.

836   **(C)** Comparing the number of distinct viral species identified by Phanta using the default
837   database and parameters vs. a standard, assembly-based workflow to identify viral
838   species in shotgun metagenomes. Dots represent individual metagenomes and lines
839   are drawn between dots representing the same metagenome. Distributions of dots are
840   shown using both boxplots and violin plots.
841   **Note:** in all boxplots, boxes represent the interquartile range (IQR), the horizontal line
842   indicates the median, and whiskers extend between (25th percentile - 1.5*IQR) and
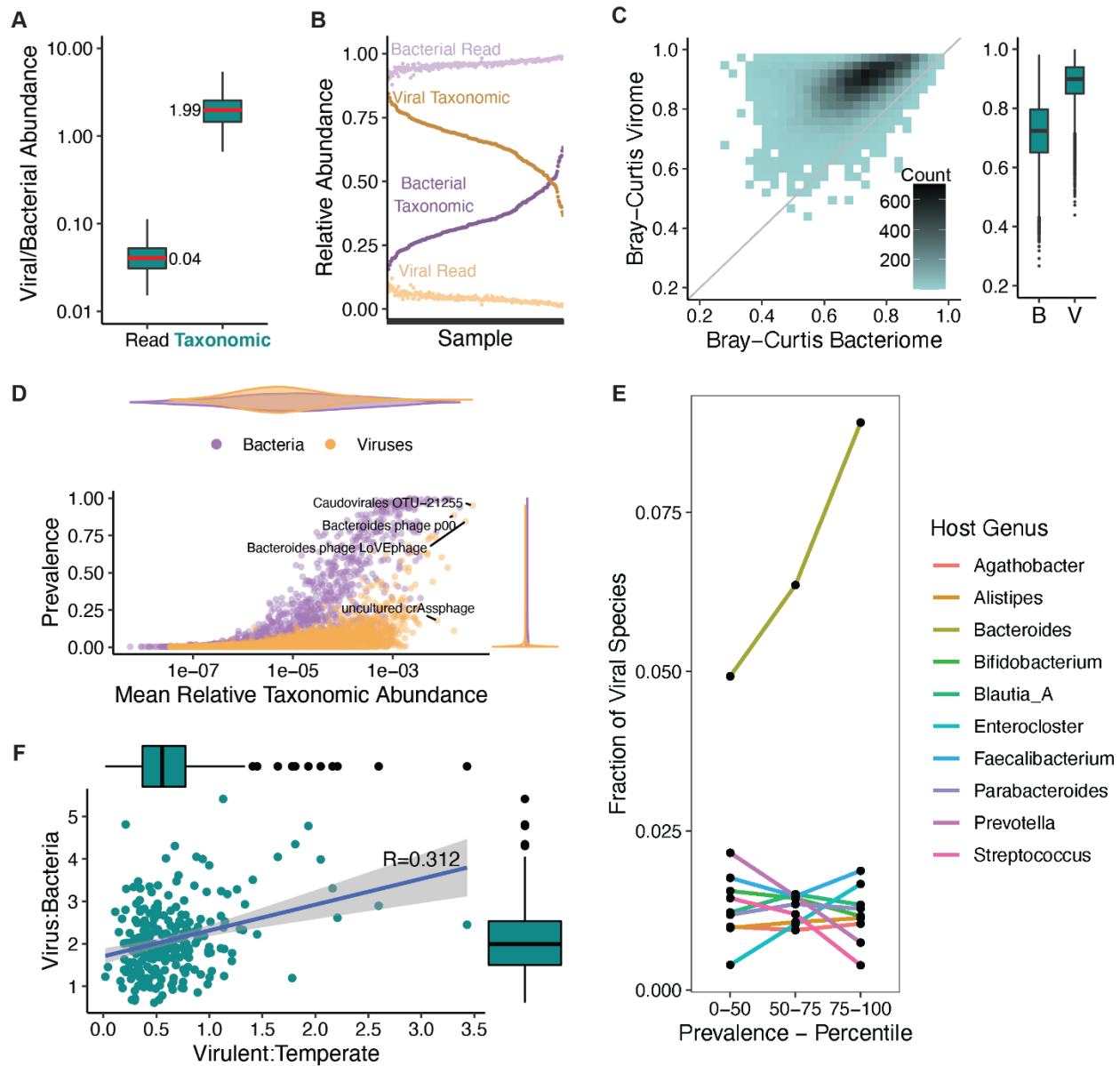843   (75th percentile + 1.5*IQR).
844
845
846
847
848
849

850



851
852  **Figure 4. Core properties of the healthy adult virome.** Metagenomes sourced from
853  Yachida *et al.*[43] (same as Figure 3).
854  **(A)** Ratio of viral to bacterial abundance in the gut, using relative read abundance vs.
855  relative taxonomic abundance. Boxplots display the distribution of this ratio across the
856  set of 245 healthy adult metagenomes.
857  **(B)** Abundance values used to calculate the ratios in (A).
858  **(C)** Comparing the variability of the gut phageome and bacteriome across
859  metagenomes. Bray-Curtis dissimilarities were calculated twice between all
860  metagenome pairs, once using relative taxonomic abundances of bacterial species
861  (horizontal axis of scatterplot) and once using relative taxonomic abundances of viral
862  species (vertical axis of scatterplot). The boxplots display the same data - B =
863  bacteriome, V = virome. The gray line on the scatterplot is the x=y diagonal.
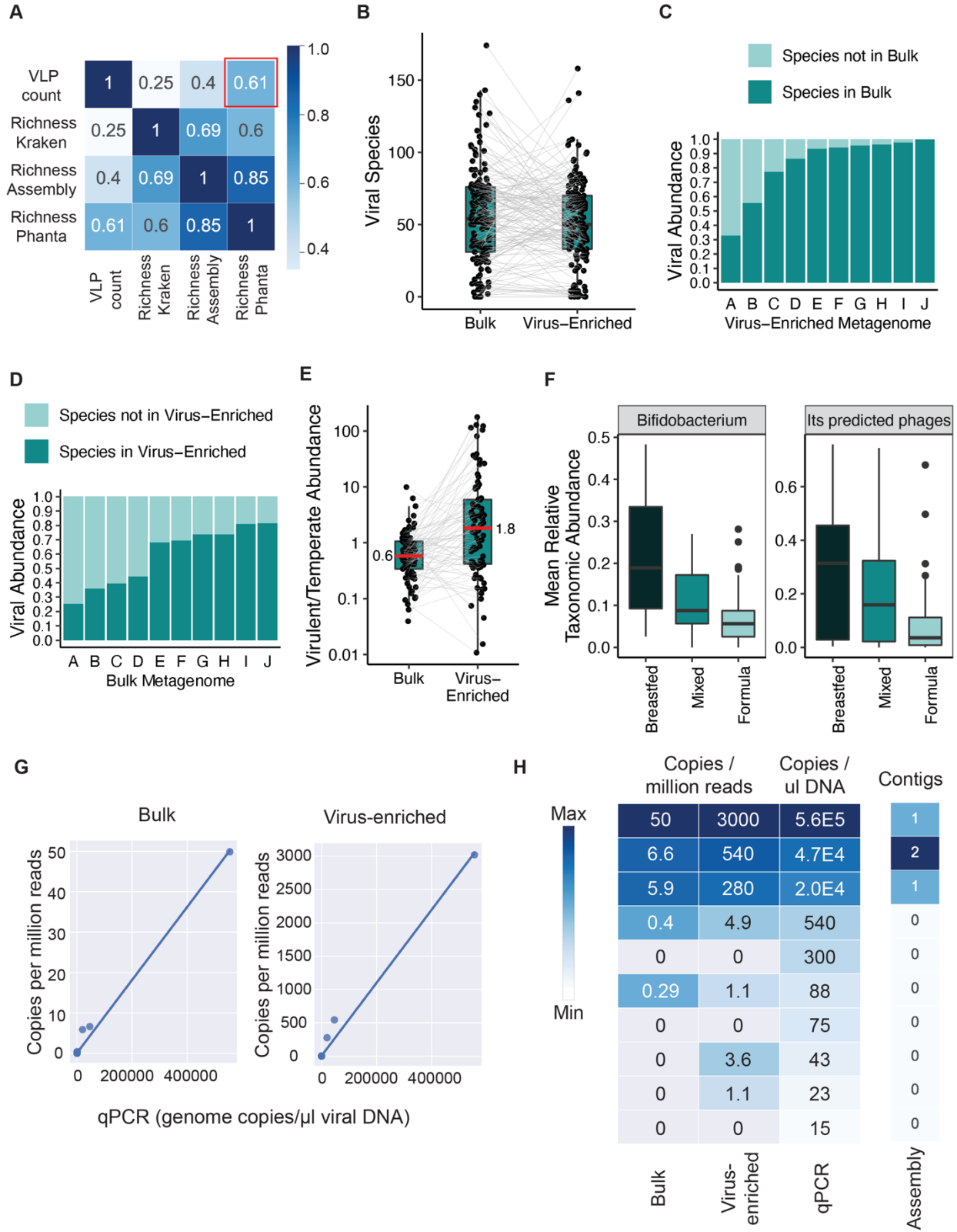
864   **(D)** Abundance and prevalence of bacterial and viral species. Abundance is the mean
865   relative taxonomic abundance across metagenomes and prevalence is the number of
866   positive individuals divided by the cohort size (245). Violin plots aligned with the x- and
867   y-axes represent distributions of abundance and prevalence, respectively.
868   **(E)** Distribution of predicted host genera for viral species in various prevalence
869   categories (e.g., category 75-100 represents the top 25% of viruses in terms of
870   prevalence). These results are based on host genus predictions that were made using
871   iPHoP[65] and are provided in Phanta's default database.
872   **(F)** Relationship between abundance ratio of viruses and bacteria and abundance ratio
873   of virulent and temperate phages. Boxplots aligned with the x- and y-axes display the
874   distributions of each ratio. Results are based on viral lifestyle predictions made by
875   BACPHLIP[64] (provided in Phanta's default database). Displayed $R$ is Pearson's
876   correlation coefficient. Relative taxonomic abundance was used as the abundance
877   metric. To prevent low quality samples from affecting the analysis, 11 outliers for
878   sequencing depth - i.e., >1.5*IQR above or below the median depth - were removed
879   (n=234).
880   **Note:** in all boxplots, boxes represent the interquartile range (IQR), the horizontal line
881   indicates the median, and whiskers extend between (25th percentile - 1.5*IQR) and
882   (75th percentile + 1.5*IQR).

**Figure 5. Application of Phanta to paired virus-enriched and bulk metagenomes from the infant gut.** Metagenomes sourced from Liang *et al.*[59] Longitudinal cohort = 20 infants sampled at months 0, 1, and 4 (60 samples total). Four-month cohort = 83 infants sampled at month 4.

**(A)** All-by-all Spearman's correlations between statistics related to viral content, for all virus-enriched metagenomes from infants in the longitudinal cohort (n=60). Specifically, four statistics were correlated: (1) VLP Count: number of viral-like particles per gram feces, (2) Richness Kraken: viral species richness based on applying Kraken2 with a RefSeq-based database, (3) Richness Assembly: viral species richness based on applying an assembly-based method, and (4) Richness Phanta: viral species richness based on applying Phanta. Phanta's richness estimation has the highest correlation with VLP count (red box). VLP Count, Richness Kraken and Richness Assembly were originally reported by Liang *et al.*.

**(B)** Number of viral species identified by Phanta in all metagenome pairs, from both infant cohorts. Each dot represents a metagenome and lines connect metagenome pairs.

**(C)** Overlap between viral species identified by Phanta in 10 pairs of metagenomes (see Methods) from infants in the four-month cohort. Each bar represents the total relative taxonomic abundance of viruses identified in a virus-enriched metagenome. Colors depict the proportion of this abundance from species also found in the paired bulk metagenome.

**(D)** Complementary analysis to (C), showing the proportion of relative taxonomic abundance in bulk metagenomes from species also found in virus-enriched metagenomes.

**(E)** Abundance ratio of virulent to temperate species detected by Phanta in virus-enriched and bulk metagenomes from the four-month cohort. Ratios were obtained using one of Phanta's provided post-processing scripts, along with viral lifestyle predictions that were made by BACPHLIP and are provided in Phanta's default database.

**(F)** Phanta's abundance estimates for *Bifidobacterium* and predicted *Bifidobacterium* phages in bulk metagenomes from infants in the four-month cohort (who had a range of diets). This analysis was facilitated by one of Phanta's provided post-processing scripts, along with host genus predictions that were made by iPHoP[65] and are provided in Phanta's default database.

**(G)** Relationship between the originally reported abundance of Adenovirus in infant stool samples (based on qPCR), vs. the newly determined abundance, based on applying Phanta to the corresponding metagenomes. This analysis considered all stool samples from the four-month cohort; most were negative or weakly positive by both methods (i.e. plotted close to (0, 0)).

**(H)** Heatmap of Adenovirus abundance in stool samples from infants in the four-month cohort, as determined by four complementary methods. Shown are stool samples originally reported to be positive for Adenovirus by qPCR. Method abbreviations: qPCR = qPCR for Adenovirus from DNA extracted from virus-like particles; quantified by genome copies / μl DNA. Assembly = alignment of assembled contigs from bulk metagenomes to Adenovirus genomes; quantified by number of contigs identified as Adenovirus. Bulk/virus-enriched = application of Phanta to bulk or virus-enriched metagenomes, using the default Phanta database; quantified by genome copies per million reads.

939 **Note:** in all boxplots, boxes represent the interquartile range (IQR), the horizontal line
940 indicates the median, and whiskers extend between (25th percentile - 1.5*IQR) and (75th
941 percentile + 1.5*IQR).
942
943
944 **References**
945

946 1. Shreiner, A. B., Kao, J. Y. & Young, V. B. The gut microbiome in health and in disease.

947 *Curr. Opin. Gastroenterol.* **31**, 69–75 (2015).

948 2. Pflughoeft, K. J. & Versalovic, J. Human microbiome in health and disease. *Annu. Rev.*

949 *Pathol.* **7**, 99–122 (2012).

950 3. Cryan, J. F. *et al.* The Microbiota-Gut-Brain Axis. *Physiol. Rev.* **99**, 1877–2013 (2019).

951 4. Kau, A. L., Ahern, P. P., Griffin, N. W., Goodman, A. L. & Gordon, J. I. Human nutrition, the

952 gut microbiome and the immune system. *Nature* **474**, 327–336 (2011).

953 5. Tringe, S. G. & Hugenholtz, P. A renaissance for the pioneering 16S rRNA gene. *Curr.*

954 *Opin. Microbiol.* **11**, 442–446 (2008).

955 6. Drewes, J. L. *et al.* High-resolution bacterial 16S rRNA gene profile meta-analysis and

956 biofilm status reveal common colorectal cancer consortia. *NPJ Biofilms Microbiomes* **3**, 34

957 (2017).

958 7. Romano, S. *et al.* Meta-analysis of the Parkinson's disease gut microbiome suggests

959 alterations linked to intestinal inflammation. *NPJ Parkinsons Dis* **7**, 27 (2021).

960 8. Yatsunenko, T. *et al.* Human gut microbiome viewed across age and geography. *Nature*

961 **486**, 222–227 (2012).

962 9. Davis-Richardson, A. G. *et al.* Bacteroides dorei dominates gut microbiome prior to

963 autoimmunity in Finnish children at high risk for type 1 diabetes. *Front. Microbiol.* **5**, 678

964 (2014).

965 10. Xu, W. *et al.* Characterization of Shallow Whole-Metagenome Shotgun Sequencing as a

966 High-Accuracy and Low-Cost Method by Complicated Mock Microbiomes. *Front. Microbiol.*

967     **12**, 678319 (2021).

968     11. Sharpton, T. J. An introduction to the analysis of shotgun metagenomic data. *Front. Plant*

969         *Sci.* **5**, 209 (2014).

970     12. Quince, C., Walker, A. W., Simpson, J. T., Loman, N. J. & Segata, N. Shotgun

971         metagenomics, from sampling to analysis. *Nat. Biotechnol.* **35**, 833–844 (2017).

972     13. Kunin, V., Copeland, A., Lapidus, A., Mavromatis, K. & Hugenholtz, P. A bioinformatician's

973         guide to metagenomics. *Microbiol. Mol. Biol. Rev.* **72**, 557–78, Table of Contents (2008).

974     14. Gregory, A. C. *et al.* MetaPop: a pipeline for macro- and microdiversity analyses and

975         visualization of microbial and viral metagenome-derived populations. *Microbiome* **10**, 49

976         (2022).

977     15. Pandolfo, M., Telatin, A., Lazzari, G., Adriaenssens, E. M. & Vitulo, N. MetaPhage: an

978         automated pipeline for analyzing, annotating, and classifying bacteriophages in

979         metagenomics sequencing data. doi:10.1101/2022.04.17.488583.

980     16. Shen, W. *et al.* KMCP: accurate metagenomic profiling of both prokaryotic and viral

981         populations by pseudo-mapping. doi:10.1101/2022.03.07.482835.

982     17. Lopera-Maya, E. A. *et al.* Effect of host genetics on the gut microbiome in 7,738

983         participants of the Dutch Microbiome Project. *Nat. Genet.* **54**, 143–151 (2022).

984     18. Almeida, A. *et al.* A unified catalog of 204,938 reference genomes from the human gut

985         microbiome. *Nat. Biotechnol.* **39**, 105–114 (2021).

986     19. Hiseni, P., Rudi, K., Wilson, R. C., Hegge, F. T. & Snipen, L. HumGut: a comprehensive

987         human gut prokaryotic genomes collection filtered by metagenome data. *Microbiome* **9**, 1–

988         12 (2021).

989     20. Zhernakova, A. *et al.* Population-based metagenomics analysis reveals markers for gut

990         microbiome composition and diversity. *Science* **352**, 565–569 (2016).

991     21. O'Leary, N. A. *et al.* Reference sequence (RefSeq) database at NCBI: current status,

992         taxonomic expansion, and functional annotation. *Nucleic Acids Res.* **44**, D733–45 (2016).

993   22. Wood, D. E., Lu, J. & Langmead, B. Improved metagenomic analysis with Kraken 2.

994        *Genome Biol.* **20**, 257 (2019).

995   23. Segata, N. *et al.* Metagenomic microbial community profiling using unique clade-specific

996        marker genes. *Nat. Methods* **9**, 811–814 (2012).

997   24. Wood, D. E. & Salzberg, S. L. Kraken: ultrafast metagenomic sequence classification using

998        exact alignments. *Genome Biol.* **15**, R46 (2014).

999   25. Khan Mirzaei, M. *et al.* Challenges of Studying the Human Virome - Relevant Emerging

1000       Technologies. *Trends Microbiol.* **29**, 171–181 (2021).

1001  26. Sutton, T. D. S., Clooney, A. G., Ryan, F. J., Ross, R. P. & Hill, C. Choice of assembly

1002       software has a critical impact on virome characterisation. *Microbiome* **7**, 12 (2019).

1003  27. Nayfach, S. *et al.* Metagenomic compendium of 189,680 DNA viruses from the human gut

1004       microbiome. *Nature Microbiology* **6**, 960–970 (2021).

1005  28. Bharti, R. & Grimm, D. G. Current challenges and best-practice protocols for microbiome

1006       analysis. *Brief. Bioinform.* **22**, 178–193 (2021).

1007  29. Ghurye, J. S., Cepeda-Espinoza, V. & Pop, M. Metagenomic Assembly: Overview,

1008       Challenges and Applications. *Yale J. Biol. Med.* **89**, 353–362 (2016).

1009  30. Rose, R., Constantinides, B., Tapinos, A., Robertson, D. L. & Prosperi, M. Challenges in

1010       the analysis of viral metagenomes. *Virus Evol* **2**, vew022 (2016).

1011  31. Gregory, A. C. *et al.* The Gut Virome Database Reveals Age-Dependent Patterns of Virome

1012       Diversity in the Human Gut. *Cell Host Microbe* **28**, 724–740.e8 (2020).

1013  32. Soto-Perez, P. *et al.* CRISPR-Cas System of a Prevalent Human Gut Bacterium Reveals

1014       Hyper-targeting against Phages in a Human Virome Catalog. *Cell Host Microbe* **26**, 325–

1015       335.e5 (2019).

1016  33. Paez-Espino, D. *et al.* IMG/VR v.2.0: an integrated data management and analysis system

1017       for cultivated and environmental viral genomes. *Nucleic Acids Res.* **47**, D678–D686 (2019).

1018  34. Tisza, M. J. & Buck, C. B. A catalog of tens of thousands of viruses from human

1019      metagenomes reveals hidden associations with chronic diseases. *Proc. Natl. Acad. Sci. U.*

1020      *S. A.* **118**, (2021).

1021   35. Camarillo-Guerrero, L. F., Almeida, A., Rangel-Pineros, G., Finn, R. D. & Lawley, T. D.

1022      Massive expansion of human gut bacteriophage diversity. *Cell* **184**, 1098–1109.e9 (2021).

1023   36. Benler, S. *et al.* Thousands of previously unknown phages discovered in whole-community

1024      human gut metagenomes. *Microbiome* **9**, 1–17 (2021).

1025   37. Van Espen, L. *et al.* A Previously Undescribed Highly Prevalent Phage Identified in a

1026      Danish Enteric Virome Catalog. *mSystems* **6**, e0038221 (2021).

1027   38. Lu J, Breitwieser FP, Thielen P, Salzberg SL. Bracken: estimating species abundance in

1028      metagenomics data. *PeerJ Computer Science* **3**, (2017).

1029   39. Wright, R. J., Comeau, A. M. & Langille, M. G. I. From defaults to databases: parameter

1030      and database choice dramatically impact the performance of metagenomic taxonomic

1031      classification tools. *bioRxiv* 2022.04.27.489753 (2022) doi:10.1101/2022.04.27.489753.

1032   40. Breitwieser, F. P., Baker, D. N. & Salzberg, S. L. KrakenUniq: confident and fast

1033      metagenomics classification using unique k-mer counts. *Genome Biol.* **19**, 198 (2018).

1034   41. Sun, Z. *et al.* Challenges in benchmarking metagenomic profilers. *Nat. Methods* **18**, 618–

1035      626 (2021).

1036   42. Kieft, K., Zhou, Z. & Anantharaman, K. VIBRANT: automated recovery, annotation and

1037      curation of microbial viruses, and evaluation of viral community function from genomic

1038      sequences. *Microbiome* **8**, 90 (2020).

1039   43. Yachida, S. *et al.* Metagenomic and metabolomic analyses reveal distinct stage-specific

1040      phenotypes of the gut microbiota in colorectal cancer. *Nat. Med.* **25**, 968–976 (2019).

1041   44. Index zone by BenLangmead. https://benlangmead.github.io/aws-indexes/k2.

1042   45. Roux, S., Enault, F., Hurwitz, B. L. & Sullivan, M. B. VirSorter: mining viral signal from

1043      microbial genomic data. *PeerJ* **3**, e985 (2015).

1044   46. Guo, J. *et al.* VirSorter2: a multi-classifier, expert-guided approach to detect diverse DNA

1045    and RNA viruses. *Microbiome* **9**, 37 (2021).

1046    47. Ren, J. *et al.* Identifying viruses from metagenomic data using deep learning. *Quant Biol* **8**,

1047        64–77 (2020).

1048    48. Ren, J., Ahlgren, N. A., Lu, Y. Y., Fuhrman, J. A. & Sun, F. VirFinder: a novel k-mer based

1049        tool for identifying viral sequences from assembled metagenomic data. *Microbiome* **5**, 69

1050        (2017).

1051    49. Amgarten, D., Braga, L. P. P., da Silva, A. M. & Setubal, J. C. MARVEL, a Tool for

1052        Prediction of Bacteriophage Sequences in Metagenomic Bins. *Front. Genet.* **9**, 304 (2018).

1053    50. Nurk, S., Meleshko, D., Korobeynikov, A. & Pevzner, P. A. metaSPAdes: a new versatile

1054        metagenomic assembler. *Genome Res.* **27**, 824–834 (2017).

1055    51. Nayfach, S. *et al.* CheckV assesses the quality and completeness of metagenome-

1056        assembled viral genomes. *Nat. Biotechnol.* **39**, 578–585 (2021).

1057    52. Moreno-Gallego, J. L. *et al.* Virome Diversity Correlates with Intestinal Microbiome Diversity

1058        in Adult Monozygotic Twins. *Cell Host Microbe* **25**, 261–272.e5 (2019).

1059    53. Chen, W. *et al.* Vast human gut virus diversity uncovered by combined short- and long-read

1060        sequencing. doi:10.1101/2022.07.03.498593.

1061    54. Shkoporov, A. N. *et al.* The Human Gut Virome Is Highly Diverse, Stable, and Individual

1062        Specific. *Cell Host Microbe* **26**, 527–541.e5 (2019).

1063    55. Stachler, E. & Bibby, K. Metagenomic Evaluation of the Highly Abundant Human Gut

1064        Bacteriophage CrAssphage for Source Tracking of Human Fecal Pollution. *Environ. Sci.*

1065        *Technol. Lett.* **1**, 405–409 (2014).

1066    56. Dutilh, B. E. *et al.* A highly abundant bacteriophage discovered in the unknown sequences

1067        of human faecal metagenomes. *Nature Communications* vol. 5 (2014).

1068    57. Benler, S. *et al.* A diversity-generating retroelement encoded by a globally ubiquitous

1069        Bacteroides phage. *Microbiome* **6**, 191 (2018).

1070    58. Kleiner, M., Hooper, L. V. & Duerkop, B. A. Evaluation of methods to purify virus-like

particles for metagenomic sequencing of intestinal viromes. *BMC Genomics* **16**, 7 (2015).

59. Liang, G. *et al.* The stepwise assembly of the neonatal virome is modulated by breastfeeding. *Nature* **581**, 470–474 (2020).

60. Krishnamurthy, S. R. & Wang, D. Origins and challenges of viral dark matter. *Virus Res.* **239**, 136–142 (2017).

61. Roux, S., Hallam, S. J., Woyke, T. & Sullivan, M. B. Viral dark matter and virus–host interactions resolved from publicly available microbial genomes. *eLife* vol. 4 (2015).

62. Jain, C., Rodriguez-R, L. M., Phillippy, A. M., Konstantinidis, K. T. & Aluru, S. High throughput ANI analysis of 90K prokaryotic genomes reveals clear species boundaries. *Nat. Commun.* **9**, 1–8 (2018).

63. Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics* vol. 26 841–842 (2010).

64. Hockenberry, A. J. & Wilke, C. O. BACPHLIP: predicting bacteriophage lifestyle from conserved protein domains. *PeerJ* **9**, e11396 (2021).

65. Roux, S. *et al.* iPHoP: an integrated machine-learning framework to maximize host prediction for metagenome-assembled virus genomes. (2022) doi:10.1101/2022.07.28.501908.

66. Watts, S. C., Ritchie, S. C., Inouye, M. & Holt, K. E. FastSpar: rapid and scalable correlation estimation for compositional data. *Bioinformatics* **35**, 1064–1066 (2019).

67. Friedman, J. & Alm, E. J. Inferring Correlation Networks from Genomic Survey Data. *PLoS Comput. Biol.* **8**, e1002687 (2012).

68. Fritz, A. *et al.* CAMISIM: simulating metagenomes and microbial communities. *Microbiome* **7**, 17 (2019).

69. Olm, M. R., Brown, C. T., Brooks, B. & Banfield, J. F. dRep: a tool for fast and accurate genomic comparisons that enables improved genome recovery from metagenomes through de-replication. *ISME J.* **11**, 2864–2868 (2017).