

VIBRANT: Automated recovery, annotation and curation of microbial viruses, and evaluation of viral community function from genomic sequence

Kristopher Kieft

University of Wisconsin Madison Department of Bacteriology

Zhichao Zhou

University of Wisconsin Madison Department of Bacteriology

Karthik Anantharaman (✉ karthik@bact.wisc.edu)

University of Wisconsin Madison <https://orcid.org/0000-0002-9584-2491>

Methodology

Keywords: Virome, virus, bacteriophage, metagenome, machine learning, auxiliary metabolism, software

Posted Date: March 5th, 2020

DOI: <https://doi.org/10.21203/rs.3.rs-16226/v1>

License:  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Version of Record: A version of this preprint was published at Microbiome on June 10th, 2020. See the published version at <https://doi.org/10.1186/s40168-020-00867-0>.

1 **VIBRANT: Automated recovery, annotation and curation of microbial viruses, and**
2 **evaluation of viral community function from genomic sequences**

3
4 Kristopher Kieft¹, Zhichao Zhou¹, and Karthik Anantharaman^{1*}

5
6 **Affiliations:**

7 ¹Department of Bacteriology, University of Wisconsin–Madison, Madison, WI, USA

8
9 *Corresponding author

10
11 Email: karthik@bact.wisc.edu

12
13 Address: 4550 Microbial Sciences Building, 1550 Linden Dr., Madison, WI, 53706

47 **Abstract**

48

49 **Background**

50 Viruses are central to microbial community structure in all environments. The ability to generate
51 large metagenomic assemblies of mixed microbial and viral sequences provides the opportunity to
52 tease apart complex microbiome dynamics, but these analyses are currently limited by the tools
53 available for analyses of viral genomes and assessing their metabolic impacts on microbiomes.

54

55 **Design**

56 Here we present VIBRANT, the first method to utilize a hybrid machine learning and protein
57 similarity approach that is not reliant on sequence features for automated recovery and annotation
58 of viruses, determination of genome quality and completeness, and characterization of viral
59 community function from metagenomic assemblies. VIBRANT uses neural networks of protein
60 signatures and a newly developed v-score metric that circumvents traditional boundaries to
61 maximize identification of lytic viral genomes and integrated proviruses, including highly diverse
62 viruses. VIBRANT highlights viral auxiliary metabolic genes and metabolic pathways, thereby
63 serving as a user-friendly platform for evaluating viral community function. VIBRANT was
64 trained and validated on reference virus datasets as well as microbiome and virome data.

65

66 **Results**

67 VIBRANT showed superior performance in recovering higher quality viruses and concurrently
68 reduced the false identification of non-viral genome fragments in comparison to other virus
69 identification programs, specifically VirSorter, VirFinder and MARVEL. When applied to
70 120,834 metagenomically derived viral sequences representing several human and natural
71 environments, VIBRANT recovered an average of 94% of the viruses, whereas VirFinder,
72 VirSorter and MARVEL achieved less powerful performance, averaging 48%, 87% and 71%,
73 respectively. Similarly, VIBRANT identified more total viral sequence and proteins when applied
74 to real metagenomes. When compared to PHASTER, Prophage Hunter and VirSorter for the ability
75 to extract integrated provirus regions from host scaffolds, VIBRANT performed comparably and
76 even identified proviruses that the other programs did not. To demonstrate applications of
77 VIBRANT, we studied viromes associated with Crohn's Disease to show that specific viral groups,
78 namely Enterobacteriales-like viruses, as well as putative dysbiosis associated viral proteins are
79 more abundant compared to healthy individuals, providing a possible viral link to maintenance of
80 diseased states.

81

82 **Conclusions**

83 The ability to accurately recover viruses and explore viral impacts on microbial community
84 metabolism will greatly advance our understanding of microbiomes, host-microbe interactions and
85 ecosystem dynamics.

86

87 **Keywords**

88 Virome, virus, bacteriophage, metagenome, machine learning, auxiliary metabolism, software

89

90

91

92 **Background**

93 Viruses that infect bacteria and archaea are globally abundant, and outnumber their hosts
94 in most environments [1–3]. Viruses are obligate intracellular pathogenic genetic elements capable
95 of reprogramming host cellular metabolic states during infection and can cause the lysis of 20-
96 40% of microorganisms in diverse environments every day [4,5]. Due to their abundance and
97 widespread activity, viruses are key facets in microbial communities as they contribute to cycling
98 of essential nutrients such as carbon, nitrogen, phosphorus and sulfur [6–10]. In human systems,
99 viruses have been implicated in contributing to dysbiosis that can lead to various diseases, such as
100 inflammatory bowel diseases, or even have a symbiotic role with the immune system [11–13].

101 Viruses harbor vast potential for diverse genetic content, arrangement and encoded
102 functions [14]. Recently, there has been substantial interest in “mining” viral sequences for novel
103 anti-microbial drug candidates, enzymes for biotechnological applications, and for bioremediation
104 [15–19]. Moreover, viruses have a unique capability to rapidly evolve genes via high mutation
105 rates and act as intermediate carriers to transfer these genes to the surrounding microbial
106 communities [20–22]. Our understanding of the diversity of viruses continues to expand with the
107 discovery of novel viral lineages, such as the characterizations of highly abundant crAssphage
108 within the human gut [23], megaphages that push the boundaries on the coding capacity of
109 prokaryotic viruses [24,25], and the *Autolykiviridae* family of small non-tailed bacterial viruses
110 [26]. To date, estimates of characterized viral diversity are biased towards tailed dsDNA viruses
111 and are likely underrepresenting other groups including those with ssDNA and RNA genomes
112 [27,28].

113 Recently it has been appreciated that viruses may directly link biogeochemical cycling of
114 nutrients by specifically driving metabolic processes. For example, during infection viruses can
115 acquire 40-90% of their required nutrients from the surrounding environment by taking over and
116 subsequently directing host metabolism [29–31]. To manipulate host metabolic frameworks, some
117 viruses selectively “steal” metabolic genes from their host. These host derived genes, collectively
118 termed auxiliary metabolic genes (AMGs), can be actively expressed during infection to provide
119 viruses with fitness advantages [32–35]. Viruses encoding AMGs have been found to be
120 widespread in human and natural environments and implicated in manipulating important nutrient
121 cycles [36–40]. Identifying these genes and understanding the processes underpinning their
122 function is pivotal for developing comprehensive models of the impacts of viruses on microbiomes
123 and nutrient cycling.

124 Due to the difficulty of collecting virus-only samples and the biases associated with doing
125 so [41,42], as well as the need to integrate viruses into models of ecosystem function, it has become
126 of great interest to determine which sequences within whole microbial communities are derived
127 from viruses. Even within the separated cellular fraction of a sample there can remain a large
128 number of viruses for a variety of reasons. First, these viruses can exist as active intracellular
129 infections, which may be the case for as many as 30% of all bacteria at any given time [43]. Second,
130 there may be particle-attached viruses resulting from viruses’ inherently “sticky” nature [44].
131 Lastly, many viruses exist as “proviruses”, or viral genomes either integrated into that of their host
132 or existing within the host as an episomal sequence. As such, it is crucial for the accurate evaluation
133 of microbial community characteristics, structure and functions to be able to separate these viral
134 sequences.

135 Multiple tools exist for the identification of viruses from mixed metagenomic assemblies.
136 For several years VirSorter [45], which succeeded tools such as VIROME [46] and Metavir [47],
137 has been the most widely used for its ability to identify viral metagenomic fragments (scaffolds)

138 from large metagenomic assemblies. VirSorter predominantly relies on database searches of
139 predicted proteins, using both reference homology as well as probabilistic similarity, to compile
140 metrics of enrichment of virus-like proteins and simultaneous depletion of other proteins. To do
141 this it uses a virus-specific curated database as well as Pfam [48] for non-virus annotations, though
142 it does not fully differentiate viral from non-viral Pfam annotations. It also incorporates signatures
143 of viral genomes, such as encoding short genes or having low levels of strand switching between
144 genes. VirSorter is also unique in its ability to use these annotation and sequence metrics to identify
145 and extract integrated provirus regions from host scaffolds. After prediction of viruses, VirSorter
146 labels viral sequences with one of three confidence levels: *categories* 1, 2 or 3. Categories 1 and 2
147 are generally considered trustworthy, but category 3 predictions are more likely to contain false
148 identifications. Despite its advantages, VirSorter likely underrepresents the diversity and
149 abundance of viruses within metagenomic assemblies.

150 More recent tools have been developed as alternatives or supplements of VirSorter in order
151 to expand our appreciation and understanding of viruses. VirFinder [49] was the first tool to
152 implement machine learning and be completely independent of reference databases for predicting
153 viruses which was a platform later implemented in PPR-Meta [50]. VirFinder was built with the
154 consideration that viruses tend to display distinctive patterns of 8-nucleotide frequencies
155 (otherwise known as 8-mers), which was proposed despite the knowledge that viruses can share
156 remarkably similar nucleotide patterns with their host [51]. These 8-mer patterns were used to
157 build a random forest machine learning model to quickly classify sequences as short as 500 bp
158 without the need for gene prediction. VirFinder generates model-derived scores as well as
159 probabilities of prediction accuracy, though it is up to the user to define the cutoffs which can
160 ultimately lead to uncertainties in rates of false identification of viruses. VirFinder was shown to
161 greatly improve the ability to recover viruses compared to VirSorter, but it also demonstrates
162 substantial host and source environment biases in predicting diverse viruses. For example,
163 VirFinder was able to recover viruses infecting Proteobacteria more readily than those infecting
164 Firmicutes due to reference database-associated biases while training the machine learning model.
165 Additional biases were also identified between different source environments, seen through the
166 under-recovery of viruses from certain environments compared to others [52].

167 Additional recent tools have been developed that utilize slightly different methods for
168 identifying viruses. MARVEL [53], for example, leverages annotation, sequence signatures (e.g.,
169 strand switching and gene density) and machine learning to identify viruses from metagenomic
170 bins. MARVEL differs from VirSorter in that it only utilizes a single virus-specific database for
171 annotation and also differs from VirFinder in that it does not use global nucleotide frequency
172 patterns. However, MARVEL provides no consideration for integrated proviruses and is only
173 suitable for identifying bacterial dsDNA viruses from the order *Caudovirales* which substantially
174 limits its ability to discover novel viruses. Another recently developed tool, VirMiner [54], is
175 unique in that it functions to use metagenomic reads and associated assembly data to identify
176 viruses and performs best for high abundance (i.e., high coverage when assembled) viruses.
177 VirMiner is a web-based server that utilizes a hybrid approach of employing both homology-based
178 searches to a virus-specific database as well as machine learning. VirMiner was found to have
179 improved ability to recover viruses compared to both VirSorter and VirFinder but was concurrently
180 much less accurate. Poor accuracy would lead to a skewed interpretation of viral community
181 function if the identified virome consisted of many non-viral sequences. This distinction is
182 important because VirMiner employs functional characterization as well as determination of virus-
183 host relationships.

184 Thus far, VirSorter remains the most efficient tool for identifying integrated proviruses
185 within metagenomic assemblies. Other tools, predominantly PHASTER [55] and Prophage Hunter
186 [56], are specialized in identifying integrated proviruses from whole genomes rather than scaffolds
187 generated by metagenomic assemblies. Similar to VirSorter, these two provirus predictors rely on
188 reference homology and viral sequence signatures with sliding windows to identify regions of a
189 host genome that belong to a virus. Although they are useful for whole genomes, they lack the
190 capability of identifying scaffolds belonging to lytic (i.e., non-integrated) viruses and perform
191 slower for large datasets. In addition, both PHASTER and Prophage Hunter are exclusively
192 available as web-based servers and offer no stand-alone command line tools.

193 Here we developed VIBRANT (Virus Identification By iteRative ANnoTation), a tool for
194 automated recovery, annotation, and curation of both free and integrated viruses from
195 metagenomic assemblies and genome sequences. VIBRANT is capable of identifying diverse
196 dsDNA, ssDNA and RNA viruses infecting both bacteria and archaea, and to our knowledge has
197 no evident environmental biases. VIBRANT uses neural networks of protein annotation signatures
198 from non-reference-based similarity searches with Hidden Markov Models (HMMs) as well as a
199 unique ‘v-score’ metric to maximize identification of diverse and novel viruses. After identifying
200 viruses VIBRANT implements curation steps to validate predictions. VIBRANT additionally
201 characterizes viral community function by highlighting AMGs and assesses the metabolic
202 pathways present in viral communities. All viral genomes, proteins, annotations and metabolic
203 profiles are compiled into formats for user-friendly downstream analyses and visualization. When
204 applied to reference viruses, non-reference virus datasets and various assembled metagenomes,
205 VIBRANT outperformed VirFinder, VirSorter and MARVEL in the ability to maximize virus
206 recovery and minimize false discovery. When compared to PHASTER, Prophage Hunter and
207 VirSorter for the ability to extract integrated provirus regions from host scaffolds, VIBRANT
208 performed comparably and even identified proviruses that the other programs did not. VIBRANT
209 was also used to identify differences in metabolic capabilities between viruses originating from
210 various environments. When applied to three separate cohorts of individuals with Crohn’s Disease,
211 VIBRANT was able to identify both differentially abundant viral groups compared to healthy
212 controls as well as virally encoded genes putatively influencing a diseased state. VIBRANT is
213 freely available for download at <https://github.com/AnantharamanLab/VIBRANT>. VIBRANT is
214 also available as a user-friendly, web-based application through the CyVerse Discovery
215 Environment at <https://de.cyverse.org/de/?type=apps&app-id=c2864d3c-fd03-11e9-9cf4-008cfa5ae621&system-id=de> [57].
216
217

218 **Methods**

219 **Dataset for generation and comparison of metrics**

220 To generate training and testing datasets sequences representing bacteria, archaea,
221 plasmids and viruses were downloaded from the National Center for Biotechnology Information
222 (NCBI) RefSeq and Genbank databases (accessed July 2019) (Additional File 1: Table S1). For
223 bacteria/archaea, 181 genomes were chosen by selecting from diverse phylogenetic groups.
224 Likewise, a total of 1,452 bacterial plasmids were chosen. For viruses, NCBI taxids associated
225 with viruses that infect bacteria or archaea were used to download reference virus genomes, which
226 were then limited to only sequences above 3kb. This included viruses with both DNA and RNA
227 genomes, though RNA genomes must first be converted to complementary DNA. Sequences not
228 associated with genomes, such as partial genomic regions, were identified according to sequence
229 headers and removed. This resulted in 15,238 total viral partial and complete genomes. To be

230 consistent between all sequences acquired from NCBI, proteins and genes were predicted using
231 Prodigal (-p meta, v2.6.3) [58]. All sequences were split into non-overlapping, non-redundant
232 fragments between 3kb and 15kb to simulate metagenome assembled scaffolds. These simulated
233 scaffolds are hereafter called *fragments* and were used throughout training and testing VIBRANT.
234 For RNA virus detection 33 viral (bacteriophage) genomes from NCBI RefSeq and 37 from
235 Krishnamurthy *et. al.* were used [28], and for archaeal virus detection all genomes were acquired
236 from NCBI RefSeq. The RNA and archaeal viral genomes were represented in both the training
237 and testing datasets as genomic fragments and recall evaluation was performed on whole genomes.
238 These were the only datasets in which training and evaluation datasets were semi-redundant.

239 Integrated viruses are common in both bacteria and archaea. To address this for generating
240 a dataset devoid of viruses, PHASTER (accessed July 2019) was used to predict putative integrated
241 viruses in the 181 bacteria/archaea genomes. Using BLASTn [59], any fragments that had
242 significant similarity (at least 95% identity, at least 3kb coverage and e-value < 1e-10) to the
243 PHASTER predictions were removed as contaminant virus sequence. The new bacteria/archaea
244 dataset was considered depleted of proviruses, but not entirely devoid of contamination. Next, the
245 datasets for bacteria/archaea and plasmids were annotated with KEGG, Pfam and VOG HMMs
246 (hmmsearch (v3.1), e-value < 1e-5) [60] to further remove contaminant virus sequence (see next
247 section for details of HMMs). Plasmids were included because it was noted that the dataset
248 appeared to contain virus sequences, possibly due to misclassification of episomal proviruses as
249 plasmids. Using manual inspection of the KEGG, Pfam and VOG annotations any sequence that
250 clearly belonged to a virus was removed. Manual inspection was guided first by the number of
251 KEGG, Pfam and VOG annotations, and then by the annotations themselves. For example,
252 sequences with more VOG than KEGG or Pfam annotations were inspected and removed if
253 multiple viral hallmark genes were found or if the majority of annotations represented viral-like
254 genes. The final datasets consisted of 400,291 fragments for bacteria/archaea, 14,739 for plasmids,
255 and 111,963 for viruses. Total number of fragments for all datasets used can be found in Additional
256 File 2: Table S2.

257

258 **Databases used by VIBRANT**

259 VIBRANT uses HMM profiles from three different databases: Kyoto Encyclopedia of
260 Genes and Genomes (KEGG) KoFam (March 2019 release) [61,62], Pfam (v32) [48] and Virus
261 Orthologous Groups (VOG) (release 94, vogdb.org). For Pfam all HMM profiles were used. To
262 increase speed, KEGG and VOG HMM databases were reduced in size to contain only profiles
263 likely to annotate the viruses of interest. For KEGG this was done by only retaining profiles
264 considered to be relevant to “prokaryotes” as determined by KEGG documentation. For VOG this
265 was done by only retaining profiles that had at least one significant hit to any of the 15,238 NCBI-
266 acquired viruses using BLASTp. The resulting databases consisted of 10,033 HMM profiles for
267 KEGG, 17,929 for Pfam, and 19,182 for VOG (Additional File 3: Table S3).

268

269 **V-score generation**

270 Predicted proteins from reference viral genomes from NCBI and VOG database viral
271 proteins were combined to generate v-scores, which resulted in a total of 633,194 proteins.
272 Redundancy was removed from the viral protein dataset using CD-HIT (v4.6) [63] with a identify
273 cutoff of 95%, which resulted in a total of 240,728 viral proteins. This was the final dataset used
274 to generate v-scores. All KEGG HMM profiles were used to annotate the viral proteins. A v-score
275 for each KEGG HMM profile was determined by the number of significant (e-value < 1e-5) hits

276 by hmmsearch, divided by 100, and a maximum value was set at 10 after division. The same v-
277 score generation was done for Pfam and VOG databases. Any HMM profile with no significant
278 hits to the virus dataset was given a v-score of zero. For KEGG and Pfam databases, any annotation
279 that was given a v-score above zero and contained the keyword “phage” was given a minimum v-
280 score of 1. To highlight viral hallmark genes, any annotation within all three databases with the
281 keyword *portal*, *terminase*, *spike*, *capsid*, *sheath*, *tail*, *coat*, *virion*, *lysin*, *holin*, *base plate*,
282 *lysozyme*, *head* or *structural* was given a minimum v-score of 1. Non-prokaryotic virus annotations
283 (e.g., *reovirus core-spike protein*) were not considered. Each HMM is assigned a v-score and
284 represents a metric of virus association (i.e., do not take into account virus specificity, or
285 association with non-viruses) and are manually tuned to put greater weight on viral hallmark genes
286 (Additional File 4: Table S4). Overall, annotations that are likely non-viral will have a low v-score
287 whereas annotations that are commonly associated with viruses will have a high v-score. Raw
288 HMM table outputs for v-score generation can be found in Additional Files 5, 6 and 7 for KEGG,
289 Pfam and VOG, respectively (Additional File 5: Table S5, Additional File 6: Table S6 and
290 Additional File 7: Table S7).

291

292 **Non-neural network steps and assembly of annotation metrics**

293 VIBRANT utilizes several manually curated cutoffs in order to remove the bulk of non-
294 virus input scaffolds before the neural network classifier is implemented. These steps result in the
295 generation of 27 annotation-derived metrics that are used by the neural network classifier for virus
296 identification, which is followed by additional manually set cutoffs to curate the results.

297 First, open reading frames are predicted by Prodigal (-p meta) or a user may input predicted
298 proteins. These proteins are then annotated with the 10,033 KEGG-derived HMMs. Putative
299 integrated provirus regions are extracted at this step by using sliding windows of either four or
300 nine proteins at a time (step size = 1 protein). Within these windows, scaffolds are fragmented
301 according to v-scores and total KEGG annotations. Within the 4-protein window, scaffolds can be
302 cut if (1) there are 0-1 unannotated proteins, 3-4 proteins with a v-score of 0-0.02 and a combined
303 v-score of less than 0.06, or (2) three consecutive proteins with a v-score of 0 (considered as a 3-
304 protein sub-window). Scaffolds will also be cut using a 9-protein window if nine consecutive
305 proteins are annotated. Finally, if the final two proteins on a scaffold each have a v-score of 0, the
306 scaffold will be cut. Only scaffold fragments that contain at least 8 proteins are retained. Following
307 provirus excision, several manually set cutoffs are used to remove obvious non-viral scaffolds.
308 Briefly, this is done by removing scaffolds with a high density of KEGG annotations (e.g., over
309 70% if less than 15 proteins or over 50% if greater than 15 proteins) or a high number of
310 annotations with a v-score of 0 (e.g., over 15 total). V-scores are also used such that a scaffold that
311 may be removed for having a high density of KEGG annotations will be retained if the v-score
312 meets a specific threshold (e.g., average of 0.2).

313 Scaffolds that are retained are subsequently annotated by the 17,929 Pfam HMMs. In a
314 similar manner to KEGG, scaffolds meeting set cutoffs for density and v-scores of Pfam HMMs
315 are either retained or removed. For example, scaffolds with less than 15 total or density under 60%
316 Pfam annotations are retained; a scaffold will be retained if it has greater than 60% Pfam
317 annotations as well as an average v-score of at least 0.15. For both KEGG and Pfam cutoffs, full
318 details of every cutoff can be found in Additional File 8: Table S8.

319 Following the aforementioned cutoff steps approximately 75-85% of non-viral scaffolds
320 are removed. At this point scaffolds are annotated by the 19,182 VOG HMMs. Using VOG
321 annotations and v-scores, as well as v-scores from KEGG and Pfam, putative proviruses are

322 trimmed to remove ends that may still contain host proteins. To do this, any scaffold previously
323 cut is trimmed, at both ends, to either the first instance of a VOG annotation or the first v-score of
324 at least 0.1 from KEGG or Pfam annotations.

325 Annotations from all three databases are used to assemble 27 metrics for the neural network
326 classifier. Briefly the metrics for each scaffold individually are as follows: (1) total encoded
327 proteins, (2) total KEGG annotations, (3) sum of KEGG v-scores, (4) total Pfam annotations, (5)
328 sum of Pfam v-scores, (6) total VOG annotations, (7) sum of VOG v-scores, (8) total KEGG
329 integration related annotations (e.g., integrase), (9) total KEGG annotations with a v-score of zero,
330 (10) total Pfam integration related annotations (e.g., integrase), (11) total Pfam annotations with a
331 v-score of zero, (12) total VOG redoxin (e.g., glutaredoxin) related annotations, (13) total VOG
332 non-integrase integration related annotations, (14) total VOG integrase annotations, (15) total
333 VOG ribonucleotide reductase related annotations, (16) total VOG nucleotide replication (e.g.,
334 DNA polymerase) related annotations, (17) total KEGG nuclease (e.g., restriction endonuclease)
335 related annotations, (18) total KEGG toxin/anti-toxin related annotations, (19) total VOG hallmark
336 protein (e.g., capsid) annotations, (20) total proteins annotated by KEGG, Pfam and VOG, (21)
337 total proteins annotated by Pfam and VOG only, (22) total proteins annotated by Pfam and KEGG
338 only, (23) total proteins annotated by KEGG and VOG only, (24) total proteins annotated by
339 KEGG only, (25) total proteins annotated by Pfam only, (26) total proteins annotated by VOG
340 only, and (27) total unannotated proteins. A complete list of all annotations used to generate these
341 metrics can be found in Additional File 9: Table S9. Non-annotation features such as gene density,
342 average gene length and strand switching were not used because they were found to decrease
343 performance of the neural network classifier despite being differentiating features between
344 bacteria/archaea and viruses; viruses tend to have shorter genes, less intergenic space and strand
345 switch less frequently. This decreased performance is likely due to several reasons, such as errors
346 associated with protein prediction (e.g., missed open reading frame leading to a large “intergenic”
347 gap) or that scaffolds, due to being fragmented genomes in most cases, behave differently than the
348 genome as a whole. For example, genomic regions encoding for large structural proteins will have
349 a higher average gene size, or a small window of virus proteins may have a greater average strand
350 switching level compared to the whole genome.

351

352 **Training and testing VIBRANT**

353 The bacteria/archaea genomic, plasmid and virus datasets described above were used to
354 train and test the machine learning model. Scikit-Learn (v0.21.3) [64] libraries were used to assess
355 various machine learning strategies to identify the best performing algorithm. Among support
356 vector machines, neural networks and random forests, we found that neural networks lead to the
357 most accurate and comprehensive identification of viruses. Therefore, Scikit-Learn’s supervised
358 neural network multi-layer perceptron classifier (hereafter called neural network) was used. The
359 portion of VIBRANT’s workflow up until the neural network classifier (i.e., KEGG, Pfam and
360 VOG annotation) was used to compile the 27 annotation metrics for each scaffold. To account for
361 differences in scaffold sizes all metrics are normalized (i.e., divided by) to the total number of
362 proteins encoded by the scaffold. The first metric, for total proteins, was normalized to log base
363 10 of itself. Each metric was weighted equally, though it is worth noting that the removal of several
364 metrics did not significantly impact the accuracy of model’s prediction. The normalized results
365 were randomized, and non-redundant portions of these results were taken for training or testing
366 the neural network. In total, 93,913 fragments were used for training and 9,000 different fragments

367 were used for testing the neural network specifically (Additional File 10: Table S10 and Additional
368 File 11: Table S11).

369 To test the performance of VIBRANT in its entirety, a new testing dataset was generated
370 consisting of fragments from the neural network testing set as well as additional fragments non-
371 redundant to the previous training dataset (hereafter called comprehensive test dataset). This new
372 comprehensive test dataset was comprised of 256,713 genomic fragments from bacteria/archaea,
373 29,926 from viruses and 8,968 from plasmids. Each met the minimum protein number requirement
374 of VIBRANT: at least four open reading frames.

375

376 **Calculation of evaluation metrics and benchmarking of VIBRANT**

377 For comparison of VIBRANT (v1.2.0) to VirFinder (v1.1), VirSorter (v1.0.3) and
378 MARVEL (v0.2), the comprehensive test dataset was used. Two intervals for VirFinder and
379 VirSorter were used for comparison. For VirSorter, the intervals selected were (1) category 1 and
380 2 predictions, and (2) category 1 and 2 predictions using the *virome decontamination mode*.
381 VirSorter was ran using the “Virome” database. For VirFinder, the intervals were (1) scores greater
382 than or equal to 0.90 (approximately equivalent to a p-value of 0.013), and (2) scores greater than
383 or equal to 0.75 (approximately equivalent to a p-value of 0.037). Since MARVEL was built for
384 the identification of viral bins, each scaffold was evaluated separately as a single “bin”. To ensure
385 proper identification by MARVEL and VIBRANT, different versions of Scikit-Learn were used
386 for each (v0.19.1 and v0.21.3, respectively).

387 Several metrics were used to compare performance of all four programs: recall, precision,
388 accuracy, specificity, Mathews Correlation Coefficient (MCC) and F1 score. When calculating
389 metrics, the larger bacteria/archaea and plasmid dataset was normalized to the size of the smaller
390 viral dataset in order to make accurate calculations. All equations used can be found in Additional
391 File 12: Table S12 and the results of each calculation can be found in Additional File 13: Table
392 S13. Comparison metrics were visualized using R (v3.5.2) package “ggplot2”.

393 It is worth noting that although VIBRANT was tested using sequences that were not used
394 for training, biases may still be associated with reported metrics due to the reliance of KEGG,
395 Pfam and VOG HMMs on NCBI databases. That is, NCBI databases in part were used to construct
396 the HMMs and therefore are well suited at annotating NCBI-derived sequences. This same type of
397 bias will be seen in the evaluation of VirSorter and MARVEL, both of which rely on NCBI-reliant
398 databases. Although VirFinder does not use annotation databases, the machine learning algorithm
399 it employs was trained on NCBI-derived sequences. Similarly, biases with comparisons to
400 VirFinder, VirSorter and MARVEL will arise when using NCBI databases. Sequences from NCBI
401 were used for training each of the three programs and therefore will likely contain redundancy to
402 VIBRANT’s comprehensive test dataset. This redundancy will cause artificially enhanced
403 performance. To address these biases, we further compared all four programs to non-NCBI
404 datasets (see below).

405

406 **Additional viral datasets and metagenomes**

407 The Integrated Microbial Genomes and Viruses (IMG/VR) v2.0 database (accessed July
408 2019) [65,66] was downloaded and scaffolds originating from animal-associated, aquatic
409 sediment, city, marine A (coastal, gulf, inlet, intertidal, neritic, oceanic, pelagic and strait), marine
410 B (hydrothermal vent, volcanic and oil), deep subsurface, freshwater, human-associated, plant-
411 associated, soil, wastewater and wetland environments were selected for analysis. Venn diagram
412 visualization of virus predictions with this dataset was made using Matplotlib (v3.0.0) [67].

413 For evaluation of 1kb and 3kb fragments new subset datasets were generated. For viruses,
414 all circular viruses (i.e., assumed to be complete) from IMG/VR freshwater and soil environments
415 as well as the Human Gut Virome database [68] were split into 3kb and redundant 1kb fragments.
416 Recall metrics for viruses were reported as the average from the three datasets (i.e., IMG/VR
417 freshwater, IMG/VR soil and Human Gut Virome database). For bacterial/archaeal genomic and
418 plasmid fragments, 13kb and 15kb fragments from the comprehensive test dataset were split into
419 3kb and redundant 1kb fragments.

420 For eukaryotic contamination, three likely contaminant genomic sequences were acquired
421 from NCBI: *Candida albicans* SC5314 chromosome 1 (NC_032089.1), *Naegleria gruberi* strain
422 NEG-M (ACER01000000.1), and *Ostreococcus* sp. SAG9 (VIBA01000000.1). These sequences
423 were split into fragments ranging from 1kb to 15kb.

424 Several published, assembled metagenomes from IMG/VR representing diverse
425 environments were selected for comparing VIBRANT, Virsorter and VirFinder (IMG taxon IDs:
426 3300005281, 3300017813 and 3300000439). Fifteen publicly available datasets from the human
427 gut were assembled for assessing VIBRANT and comparing the three programs [69]. Reads can
428 be found under NCBI BioProject PRJEB7774 (ERR688591, ERR688590, ERR688509,
429 ERR608507, ERR608506, ERR688584, ERR688587, ERR688519, ERR688512, ERR688508,
430 ERR688634, ERR688618, ERR688515, ERR688513, ERR688505). Reads were trimmed using
431 Sickle (v1.33) [70] and assembled using metaSPAdes (v3.12.0 65) [71] (--meta -k 21,33,55,77,99).
432 For hydrothermal vents, six publicly available hydrothermal plume samples were derived from
433 Guaymas Basin (one sample) and Eastern Lau Spreading Center (five samples). Reads can be
434 found under NCBI BioProject PRJNA314399 (SRR3577362) and PRJNA234377 (SRR1217367,
435 SRR1217459, SRR1217564, SRR1217566, SRR1217452, SRR1217567, SRR1217465,
436 SRR1217462, SRR1217460, SRR1217463, SRR1217565). Reads were trimmed using Sickle and
437 assembled using metaSPAdes (--meta -k 21,33,55,77,99). Details of assembly and processing are
438 outlined in Zhou *et al.* [72]. For analysis of Crohn's Disease metagenomes by VIBRANT, publicly
439 available metagenomes were used; the metagenomes were sequenced by He *et al.* [73], Ijaz *et al.*
440 [74] and Gevers *et al.* [75], and assembled by Pasolli *et al.* [76] (Additional File 14: Table S14).

441 The computational resource requirements and associated runtimes for VIBRANT were
442 assessed using datasets of various sizes and composition (Additional File 15: Table S15).
443 VIBRANT was able to evaluate large datasets quickly since it was built for efficient parallelization
444 across CPUs.

445

446 **AMG identification**

447 KEGG annotations were used to classify potential AMGs (Additional File 16: Table S16).
448 KEGG annotations falling under the “metabolic pathways” category as well as “sulfur relay
449 system” were considered. Manual inspection was used to remove non-AMG annotations, such as
450 *nrdAB* and *thyAX*. Other annotations not considered were associated with direct nucleotide to
451 nucleotide conversions. All AMGs were associated with a KEGG metabolic pathway map.

452

453 **Completeness estimation**

454 Scaffold completeness is determined based on four metrics: circularization of scaffold
455 sequence, VOG annotations, total VOG nucleotide replication proteins and total VOG viral
456 hallmark proteins (Additional File 9: Table S9). In order to be considered a complete genome a
457 sequence must be identified as likely circular. A kmer-based approach is used to do this.
458 Specifically, the first 20 nucleotides are compared to 20-mer sliding windows within the last 900bp

459 of the sequence. If a complete match is identified the sequence is considered a circular template.
460 Scaffolds can also be considered a low, medium or high quality draft. To benchmark completeness,
461 2466 NCBI RefSeq viruses identified as *Caudovirales*, limited to 10 kb in length, were used to
462 estimate completeness by stepwise removing 10% viral sequence at a time. VIBRANT was found
463 to identify 2465 of the 2466 viruses. This set of viruses was additionally used to assess the error
464 rate of cutting provirus regions. Viral genome diagrams to depict genome quality and
465 completeness, provirus predictions and novel virus identification, were made using Geneious
466 Prime 2019.0.3.

467

468 **Analysis of Crohn's Disease metagenomes**

469 Metagenomic reads from He *et al.* were assembled by Pasolli *et al.* and used for analysis.
470 VIBRANT (-l 5000) was used to predict viruses from 49 metagenomes originating from
471 individuals with Crohn's Disease and 53 from healthy individuals (102 total samples). A total of
472 14,121 viruses were identified. Viral sequences were dereplicated using Mash [77] and Nucmer
473 [78] to 95% nucleotide identity and 70% sequence coverage. The longest sequence was kept as the
474 representative for a total of 8,822 dereplicated viruses. A total of 96 read sets were used (59
475 Crohn's Disease and 37 healthy), trimmed using Sickle and aligned to the dereplicated viruses
476 using Bowtie2 (-N 1, v2.3.4.1) [79] and the resulting coverages were normalized to total reads.
477 The normalized relative coverage of each virus for all 96 samples were compared using DESeq2
478 [80] (Additional File 17: Table S17). Viruses that displayed significantly different abundance
479 between Crohn's Disease and control samples were determined by a p-value cutoff of 0.05. iRep
480 (default parameters) [81] was used to estimate replication activity of two highly abundant Crohn's-
481 associated viruses. EasyFig (v2.2.2) [82] was used to generate genome alignments of Escherichia
482 phage Lambda (NCBI accession number NC_001416.1) and three Crohn's-associated viruses.
483 vConTACT2 (v0.9.8) was run using default parameters on the CyVerse Discovery Environment
484 platform. Putative hosts of Crohn's-associated and healthy-associated was estimated using
485 proximity of vConTACT2 protein clustering and BLASTp identity (NCBI non-redundant protein
486 database, assessed October 2019). Two additional read sets from Gevers *et al.* [75] and Ijaz *et al.*
487 [74] were likewise assembled by Pasolli *et al.*. VIBRANT (-l 5000 -o 10) was used to predict
488 viruses from 43 metagenomes originating from individuals with Crohn's Disease and 21 from
489 healthy individuals (64 total samples). In contrast to the discovery dataset viral genomes were not
490 dereplicated and differential abundance was not determined. Instead viruses from each group were
491 directly clustered using vConTACT2. Abundances of dysbiosis associated genes in the validation
492 set were normalized to total viruses. Validation of dysbiosis associated genes' presence on viral
493 genomes, rather than microbial contamination, was done by identifying viral hallmark genes on
494 the viral scaffold (Additional File 18: Table S18). Protein networks were visualized using
495 Cytoscape (v3.7.2) [83].

496

497 **Results**

498 VIBRANT was built to extract and analyze bacterial and archaeal viruses from assembled
499 metagenomic and genome sequences, as well as provide a platform for characterizing metabolic
500 proteins and functions in a comprehensive manner. The concept behind VIBRANT's mechanism
501 of virus identification stems from the understanding that arduous manual inspection of annotated
502 genomic sequences produces the most dependable results. As such, the primary metrics used to
503 inform validated curation standards and to train VIBRANT's machine learning based neural

504 network to identify viruses reflects human-guided intuition, though in a high-throughput
 505 automated fashion.

506

507 Determination of v-score

508 We developed a unique ‘v-score’
 509 metric as an approach for providing
 510 quantitative information to VIBRANT’s
 511 algorithm in order to assess the
 512 qualitative nature of annotation
 513 information. A v-score is a value
 514 assigned to each possible protein
 515 annotation that scores its association to
 516 known viral genomes (see Methods). V-
 517 score differs from the previously used
 518 “virus quotient” metric [84,85] in that it
 519 does not take into account the
 520 annotation’s relatedness to bacteria or
 521 archaea. Not including significant
 522 similarity to non-viral genomes in the
 523 calculation of v-scores has important
 524 implications for this metric’s utility.
 525 Foremost is that annotations shared
 526 between viruses and their hosts, such as
 527 ribonucleotide reductases, will be
 528 assigned a v-score reflecting its
 529 association to viruses, not necessarily
 530 virus-specificity. Many genes are
 531 commonly associated with viruses and
 532 host organisms, but when encoded on
 533 viral genomes can be central to virus
 534 replication efficiency (e.g.,
 535 ribonucleotide reductases [86]).
 536 Therefore, a metric representing virus-
 537 association rather than virus-specificity
 538 would be more appropriate in identifying
 539 if an unknown scaffold is viral or not.
 540 Secondly, this approach takes into
 541 account widespread horizontal gene
 542 transfer of host genes by viruses as well
 543 as the presence of AMGs.

545 VIBRANT workflow

546 VIBRANT utilizes several
 547 annotation metrics in order to guide
 548 removal of non-viral scaffolds before
 549 curation of reliable viral scaffolds. The
 550 annotation metrics used are derived from
 551 HMM-based probabilistic searches of
 552 protein families from the KEGG,

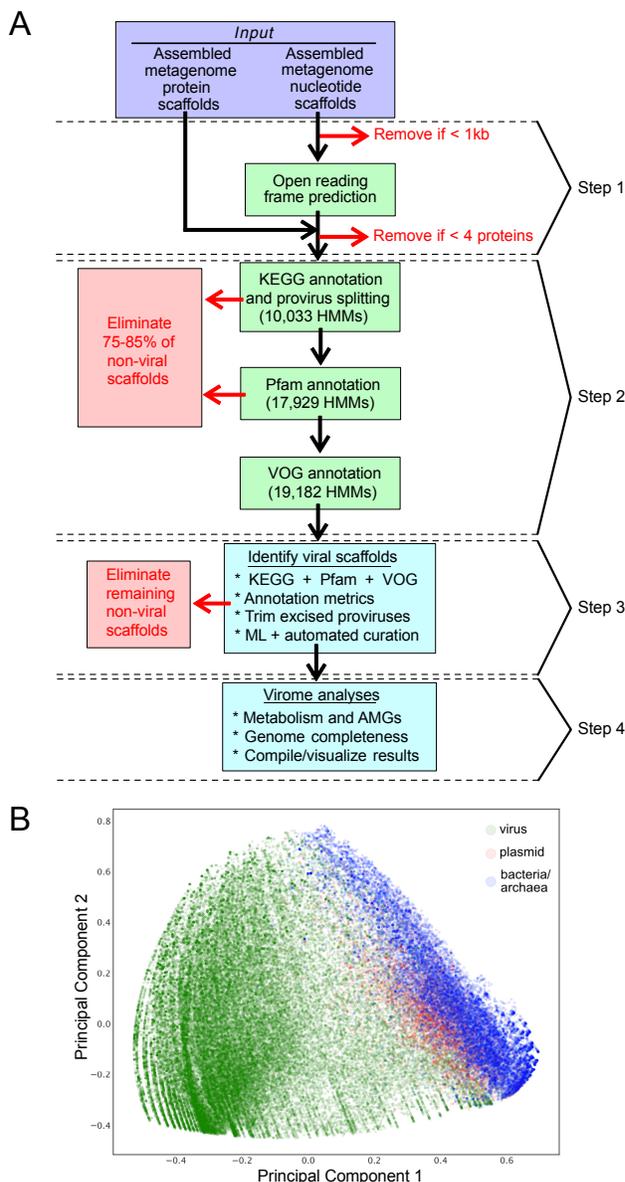


Figure 1. Representation of VIBRANT’s method for virus

identification and virome functional characterization. (A) Workflow of virome analysis. Annotations from KEGG, Pfam and VOG databases are used to construct signatures of viral and non-viral annotation signatures that are read into a neural network machine learning model. **(B)** Visual representation (PCA plot) of the metrics used by the neural network to identify viruses, depicting viral, plasmid and bacterial/archaeal genomic sequences.

The annotation metrics used are derived from HMM-based probabilistic searches of protein families from the KEGG,

550 Pfam and VOG databases. VIBRANT is not reliant on reference-based similarity and therefore
551 accounts for the large diversity of viruses on Earth and their respective proteins. Consequently,
552 widespread horizontal gene transfer, rapid mutation and the vast amount of novel sequences do
553 not hinder VIBRANT's ability to identify known and novel viruses. VIBRANT does not rely on
554 non-annotation features, such as rates of open reading frame strand switching, because these
555 features were not as well conserved in genomic scaffolds in contrast to whole genomes.

556 VIBRANT's workflow consists of four main steps (Figure 1A). Briefly, proteins (predicted
557 or user input) are used by VIBRANT to first eliminate non-viral sequences by assessing non-viral
558 annotation signatures derived from KEGG and Pfam HMM annotations. At this step potential host
559 scaffolds are fragmented using sliding windows of KEGG annotation v-scores in order to extract
560 integrated provirus sequences. Following the elimination of most non-viral scaffolds and rough
561 excision of provirus regions, proteins are annotated by VOG HMMs. Before analysis by the neural
562 network machine learning model, any extracted putative provirus is trimmed to exclude any
563 remaining non-viral sequences. Annotations from KEGG, Pfam and VOG are used to compile 27
564 metrics that are utilized by the neural network to predict viral sequences (see Methods). These 27
565 metrics were found to be adequate for the separation of viral and non-viral scaffolds (Figure 1B).

566 After prediction by the neural network a set of curation steps are used to filter the results.
567 Curation is an automated mechanism of verifying and/or altering the neural network predictions
568 in order to improve accuracy and recovery of viruses. This concept, as previously stated, originates
569 from experiences with manual inspection of viral genomes that cannot be captured even within
570 machine learning algorithms. For example, these curation steps can: (1) more accurately separate
571 plasmid sequences by discerning viral-like and plasmid-like integrase annotations, (2) remove
572 scaffolds that encode a high density of bacterial-like (i.e., v-score of zero) proteins, or (3) increase
573 true positive identifications by retaining otherwise missed scaffolds that are unique (e.g., encode
574 few but highly virus-related proteins).

575 Once viruses are identified VIBRANT automates the analysis of viral community function
576 by highlighting AMGs and assigning them to KEGG metabolic pathways. The genome quality
577 (i.e., proxy of completeness) of identified viruses is estimated using a subset of the annotation
578 metrics and viral sequences are used to identify circular templates (i.e., likely complete circular
579 viruses). These quality analyses were determined to best reflect established completeness metrics
580 for both bacteria and viruses [87,88]. Finally, VIBRANT compiles all results into a user-friendly
581 format for visualization and downstream analysis. For a detailed description of VIBRANT's
582 workflow see Methods.

583

584 **Comparison of VIBRANT to other programs**

585 VirSorter, VirFinder and MARVEL, three commonly used programs for identifying
586 bacterial and archaeal viruses from metagenomes, were selected to compare against VIBRANT
587 for the ability to accurately identify viruses. We evaluated all four programs' performance on the
588 same viral, bacterial and archaeal genomic, and plasmid datasets. Given that both VirSorter and
589 VirFinder produce various confidence ranges of virus identification, we selected certain
590 parameters for each program for comparison. For VirSorter, the parameters selected were (1)
591 category 1 and 2 predictions, and (2) category 1 and 2 predictions using the *virome*
592 *decontamination mode*. For VirFinder, the intervals were (1) scores greater than or equal to 0.90
593 (approximately equivalent to a p-value of 0.013), and (2) scores greater than or equal to 0.75
594 (approximately equivalent to a p-value of 0.037). Hereafter, we provide two statistics for each
595 VirSorter and VirFinder run that reflect results according to the two set confidence intervals,

596 respectively. Both VIBRANT and MARVEL have set output predictions and therefore will be
597 reported with a single statistic.

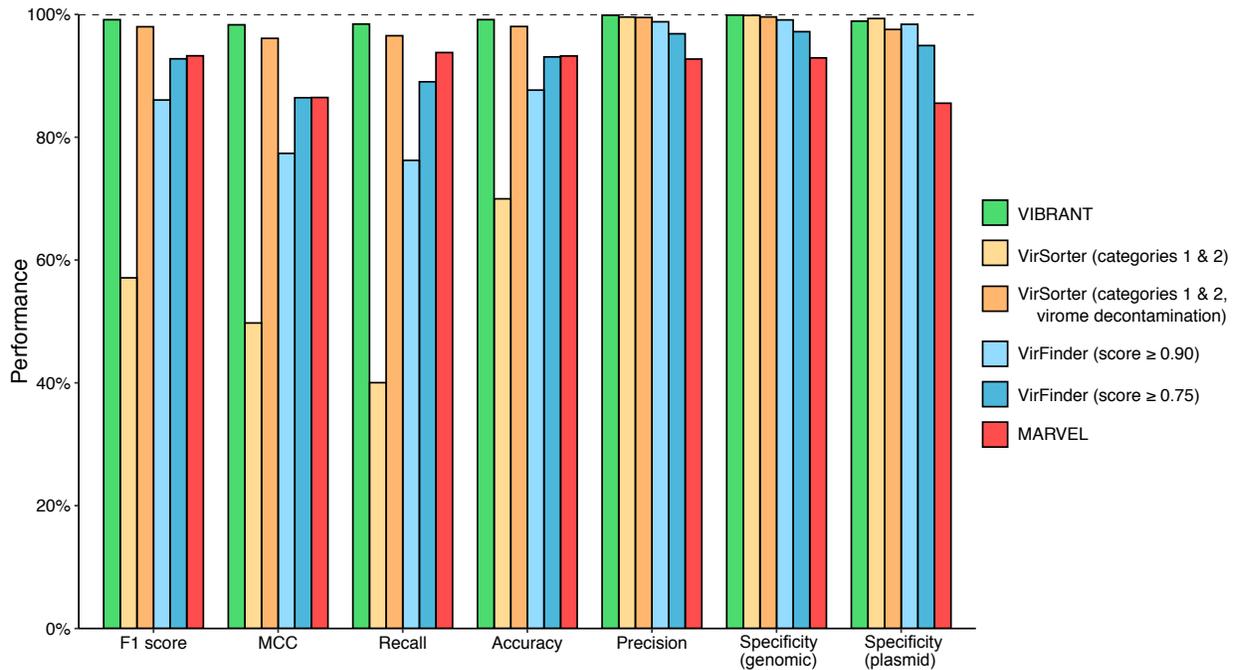


Figure 2. Performance comparison of VIBRANT, VirFinder, VirSorter and MARVEL on artificial scaffolds 3kb-15kb. Performance was evaluated using datasets of reference viruses, bacterial plasmids, and bacterial/archaeal genomes. For VirFinder and VirSorter two different confidence cutoffs were used (VirFinder: score of at least 0.90, and score of at least 0.75. VirSorter: categories 1 and 2 predictions, and categories 1 and 2 predictions using virome decontamination mode). All four programs were compared using the following statistical metrics: F1 score, MCC, recall, precision, accuracy and specificity. To ensure equal comparison all scaffolds tested encoded at least four open reading frames.

598 VIBRANT yields a single output of confident predictions and therefore does not provide
599 multiple output options. Since VIBRANT is only partially reliant on its neural network machine
600 learning model for making predictions, all comparisons are focused on VIBRANT's full workflow
601 performance. VIBRANT does not consider scaffolds shorter than 1000 bp or those that encode
602 less than four predicted open reading frames in order to maintain a low false positive rate (FPR)
603 and have sufficient annotation information for identifying viruses. Therefore, in comparison of
604 performance metrics only scaffolds meeting VIBRANT's minimum requirements were analyzed.
605 Inclusion of fragments encoding less than four open reading frames in analyses, which are
606 frequently generated by metagenomic assemblies, are discussed below. We used the following
607 statistics to compare performance: recall, precision, accuracy, specificity, MCC and F1 score
608 (Figure 2).

609 First, we evaluated the true positive rate (TPR, or recall) of viral genomic fragments as
610 well as whole viral genomes. Viral genomes were acquired from the NCBI RefSeq and GenBank
611 databases and split into various non-redundant fragments between 3 and 15 kb to simulate genomic
612 scaffolds (see Methods). VIBRANT correctly identified 98.43% of the 29,926 viral fragments,
613 which was greater than VirSorter (40.03% and 96.53%), VirFinder (76.23% and 89.03%) and
614 MARVEL (93.79%) at all scoring intervals. For VirSorter it was essential to set *virome*
615 *decontamination mode* for datasets consisting of mainly viruses, without which the TPR was
616 substantially inhibited.

617 Similar to TPR, we calculated FPR (or specificity) using two different datasets: genomic
618 fragments of bacteria and archaea (hereafter called genomic), and bacterial plasmids (plasmid).
619 Plasmids were evaluated separately because they often encode for genes similar to those on viral
620 genomes, such as those for genome replication and mobilization. Genomic and plasmid sequences
621 were acquired from NCBI RefSeq and GenBank databases and split into various non-redundant
622 fragments between 3 and 15 kb and putative proviruses were depleted from the datasets (see
623 Methods). VIBRANT had high specificity against both genomic (99.90%) and plasmid fragments
624 (98.90%). VirSorter had similar specificity against both genomic (99.84% and 99.59%) and
625 plasmid (99.33% and 97.55%) datasets, but only VirFinder set to a score cutoff of 0.90 was fully
626 comparable (genomic: 99.10%, plasmid: 98.39%). VirFinder at a score cutoff of 0.75 (genomic:
627 97.19%, plasmid: 94.93%) along with MARVEL (genomic: 92.92%, plasmid: 85.54%) were
628 slightly less specific. Although VirFinder (set to a score cutoff of 0.90) and VIBRANT had a
629 similar overall specificity, VirFinder identified 9.3 times more genomic scaffolds as viruses (false
630 discoveries) compared to VIBRANT (2,311 and 249, respectively). MARVEL was even more
631 pronounced, identifying 72.9 times more genomic scaffolds as viruses (18,164 total) compared to
632 VIBRANT.

633 We used the results from TPR of viral fragments and FPR of non-viral genomic or plasmid
634 fragments to calculate precision (i.e., proportion of true virus identifications out of all virus
635 identifications) and accuracy (i.e., proportion of correct predictions out of all predictions).
636 VIBRANT outperformed each other program at both precision (VIBRANT: 99.87%, VirFinder:
637 98.80% and 96.85%, VirSorter: 99.57% and 99.50%, and MARVEL: 92.73%) and accuracy
638 (VIBRANT: 99.15%, VirFinder: 87.67% and 93.08%, VirSorter: 69.97% and 98.03%, and
639 MARVEL: 93.23%). F1 and MCC are additional metrics (maximum values of 1) accounting for
640 both TPR and FPR, and therefore acts as a comprehensive evaluation of overall performance. Our
641 calculation of F1 indicates that VIBRANT (0.991) is able to better identify viruses while
642 subsequently reducing false identifications compared to VirFinder (0.861 and 0.928), VirSorter
643 (0.571 and 0.980) or MARVEL (0.933). MCC likewise indicated that VIBRANT (0.983) was
644 better suited at maximizing the ratio of viruses to non-viruses compared to VirSorter (0.498 and
645 0.961), VirFinder (0.774 and 0.864) and MARVEL (0.865).

646 Although VIBRANT exhibits improved performance with scaffolds at least 3kb in length,
647 it is worth noting that performance drops considerably at the set minimum length of 1kb. To
648 display this, the TPR and FPR of both 1k and 3kb scaffolds were assessed (Additional File 21:
649 Figure S1A). For this analysis, VirSorter was evaluated using virome decontamination mode and
650 VirFinder was set to a score cutoff of 0.90. MARVEL's minimum length requirement is 2kb and
651 therefore was not compared with 1kb scaffolds. For 1kb viral scaffolds, VIBRANT (1.95%) and
652 VirSorter (1.12%) recovered far fewer scaffolds compared to VirFinder (22.56%). However, at a
653 length of 3kb VIBRANT (43.54%) recovered more viral fragments than VirSorter (25.43%),
654 VirFinder (34.42%) and MARVEL (37.82%). Even at the low resolution of short scaffolds
655 VIBRANT's FPR is not impacted. For 1kb genomic and 1kb plasmid scaffolds VIBRANT
656 (<0.00% and 0.07%) and VirSorter (<0.00% and 0.10%) had fewer false positive discoveries than
657 VirFinder (2.61% and 3.70%). Similarly, for 3kb genomic and 3kb plasmid scaffolds VIBRANT
658 (0.10% and 2.69%) and VirSorter (0.11% and 2.41%) falsely identified fewer sequences than
659 VirFinder (2.26% and 5.54%) or MARVEL (6.08% and 16.30%). Overall, this suggests that
660 VirFinder is uniquely able to accurately recover short (e.g., 1kb) viral scaffolds while maintaining
661 a relatively low FPR, but this ability is not maintained with longer scaffolds. Moreover, our current
662 abilities to sequence and assemble scaffolds of lengths over 3kb will likely lead to a greater focus

663 on longer viral sequences that are more amenable to downstream analysis, such as taxonomic
664 classification and functional analyses.

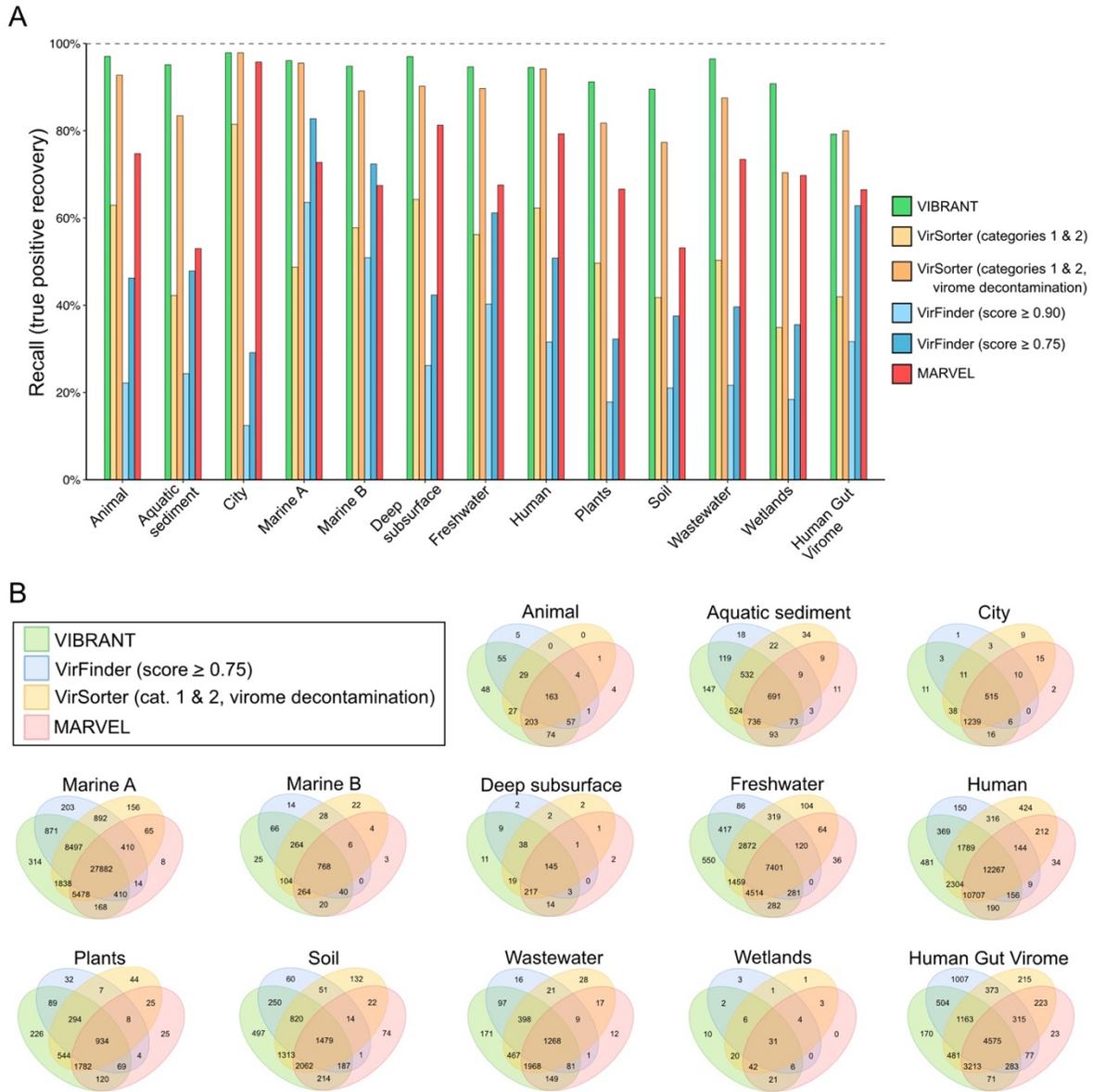
665 Next, we assessed the ability of VIBRANT to filter out eukaryotic contamination rather
666 than falsely identify these sequences as viral since eukaryotes were not represented in the training
667 or testing datasets. However, these contaminants should be sparse because the majority of
668 eukaryotic KEGG and VOG HMMs were removed from the annotation databases (see Methods).
669 Likewise, eukaryotic-like annotations should receive a low v-score. A total of 8,672 eukaryotic
670 sequences ranging from 1kb to 15kb were assessed. VIBRANT (0.62%), VirSorter (0.05% and
671 0.05%) and MARVEL (0.44%) performed well with recovering few sequences, whereas VirFinder
672 (4.92% and 15.44%) recovered contamination at a greater rate (Additional File 21: Figure S1B).

673 Finally, viruses with RNA genomes as well as those that infect archaea are rare in current
674 culture systems and sequence databases compared to bacterial dsDNA viruses. However, the true
675 abundance of RNA and archaeal viruses has yet to be explored mainly due to biases towards
676 dsDNA in genome extracting and sequencing methods [89] and the low abundance of archaea in
677 most environments. VIBRANT was built to identify all prokaryotic viruses in order to expand our
678 knowledge of understudied groups. A total of 70 RNA viral genomes and 93 archaeal viral
679 genomes were used to evaluate recall. VIBRANT was able to recover 47% of RNA viruses, or
680 84% of the those that encoded at least four predicted open reading frames. In comparison,
681 VirSorter (7% and 70%), VirFinder (33% and 57%) and MARVEL (68%) ranged from lower to
682 higher recovery (Additional File 21: Figure S1C). The high recovery of RNA viruses by MARVEL
683 is intriguing since the software was trained exclusively on dsDNA Caudovirales, but may be
684 explained by the greater rate of false positive discovery. For archaeal viruses, VIBRANT (96.77%)
685 identified significantly more viruses than VirSorter (70.97% and 93.55%), VirFinder (46.24% and
686 74.19%), and MARVEL (80.65%) (Additional File 21: Figure S1D). Taken together, VIBRANT
687 has the potential to identify RNA and archaeal viruses, though the significance of this difference
688 is hard to distinguish due to the current dearth of reference genomes with which to validate.
689

690 **Identification of viruses in diverse environments**

691 We next tested VIBRANT's ability to successfully identify viruses from a diversity of
692 environments. Using 120,834 viruses from the IMG/VR database, in which the source environment
693 of viruses is categorized, we identified that VIBRANT is more robust in identifying viruses from
694 all tested environments compared to VirFinder, VirSorter and MARVEL (Figure 3A). The 12
695 environments were: animal-associated, aquatic sediment, city, marine A (coastal, gulf, inlet,
696 intertidal, neritic, oceanic, pelagic and strait), marine B (hydrothermal vent, volcanic and oil), deep
697 subsurface, freshwater, human-associated, plant-associated, soil, wastewater and wetlands.
698 VIBRANT averaged 94.59% recall, substantially greater than VirFinder (29.19% and 48.13%),
699 VirSorter (54.37% and 87.49%) and MARVEL (71.23%). Between the 12 environments
700 VIBRANT recovered between 89.55% and 97.87% (total range of 8.33%) of the viruses.
701 Conversely, VirFinder (score cutoff of 0.75) had a range of 53.65%, VirSorter (categories 1 and
702 2, virome decontamination) had a range of 27.48% and MARVEL had a range of 42.75%. These
703 results suggest that in comparison to other software, VIBRANT has no evident environmental
704 biases and is fully capable of identifying viruses from a broad range of source environments. We
705 also used a dataset of 13,203 viruses from the Human Gut Virome database for additional
706 comparison. The vast majority of viruses (~96%) in this dataset were assumed to infect bacteria.
707 Although recall was diminished compared to IMG/VR datasets, VIBRANT (79.22%) nevertheless

708 outperformed or matched VirFinder (31.67% and 62.83%), VirSorter (41.93% and 79.97%) and
 709 MARVEL (66.49%) on this dataset.



710 Relatively few viruses from the IMG/VR dataset that were not identified by VIBRANT
 711 were identified by either VirFinder, VirSorter or MARVEL at even the most inclusive score cutoffs
 712 (Figure 3B). Furthermore, for most environments VIBRANT displayed the largest proportion of
 713 unique identifications, suggesting that VIBRANT has the propensity for discovery of viruses. The

714 differences in the overlap of identified viruses was not too distinctive in environments for which
715 many reference viruses are available, such as marine, though for more understudied environments,
716 such as plants or wastewater, VIBRANT displayed near-complete overlap with VirFinder,
717 VirSorter and MARVEL predictions. This suggests that database bias may not affect VIBRANT's
718 performance to a significant degree. Although VirFinder does not rely on an annotation database,
719 it still has been trained on a dataset of reference viral genomes which can contribute to database
720 dependency and recall bias.

721

722 **Identification of viruses in mixed metagenomes**

723 Metagenomes assembled using short read technology contain many scaffolds that do not
724 meet VIBRANT's minimum length requirements and therefore are not considered during analysis.
725 Despite this, VIBRANT's predictions contain more annotation information and greater total viral
726 sequence length than tools built to identify short sequences, such as scaffolds with less than four
727 open reading frames. VIBRANT, VirFinder (score cutoff of 0.90) and VirSorter (categories 1 and
728 2) were used to identify viruses from human gut, freshwater lake and thermophilic compost
729 metagenome sequences (Table 1). In addition, alternate program settings—VIBRANT *virome*
730 mode, VirFinder at a score cutoff of 0.75 and VirSorter *virome* decontamination mode—were used
731 to identify viruses from an estuary virome dataset. MARVEL was not considered in this analysis
732 due to the inability to achieve comparable precision. Each metagenomic assembly was limited to
733 sequences of at least 1000bp but no minimum open reading frame limit was set. For these
734 metagenomes, 31% to 40% of the scaffolds were of sufficient length (at least four open reading
735 frames) to be analyzed by VIBRANT; for the estuary virome 62% were of sufficient length. In
736 comparison, 100% of scaffolds from each dataset were long enough to be analyzed by VirFinder.
737 The ability of VirFinder to make a prediction with each scaffold is considered the major strength
738 of the tool.

739 For all six assemblies VirFinder averaged approximately 1.16 times more virus
740 identifications than VIBRANT, though for both thermophilic compost and the estuary virome
741 VIBRANT identified a greater number. Despite VirFinder averaging more total virus
742 identifications, VIBRANT averaged 2.33 times more total viral sequence length and 2.44 times
743 more total viral proteins. This is the result of VIBRANT having the capability to identify more
744 viruses of higher quality and longer sequence length. For example, among all six datasets
745 VIBRANT identified 1,320 total viruses at least 10 kb in length in comparison to VirFinder's 479.
746 VIBRANT was also able to outperform VirSorter in all metrics, averaging 2.45 times more virus
747 identifications, 1.76 times more total viral sequence length, and 1.86 times more encoded viral
748 proteins.

749 VIBRANT's method of predicting viruses provides a unique opportunity in comparison to
750 similar tools in that it yields sequences of higher quality which are more amenable for analyzing
751 protein function from virome data. It is an important distinction that the total number of viruses
752 identified may not be correlated with the total viruses identified or the total number of encoded
753 proteins. Even if VIBRANT identified fewer total viruses compared to other tools in certain
754 circumstances, more data of higher quality was generated as viral sequences of longer length were
755 identified as compared to many short fragments. This provides an important distinction that the
756 metric of total viral predictions is not necessarily an accurate representation for the quality or
757 quantity of the data generated.

Table 1. Virus recovery of VIBRANT, VirFinder and VirSorter from mixed metagenomes and a virome. Mixed community assembled metagenomes from human gut, thermophilic compost and freshwater, as well as an estuary virome, were used to compare virus prediction ability between the three programs. For each assembly the scaffolds were limited to a minimum length of 1000bp. Only a subset of each dataset contained scaffolds encoding at least four open reading frames. VIBRANT, VirFinder (score minimum of 0.90) and VirSorter (categories 1 and 2) were compared by total viral predictions, total combined length of predicted viruses, and total combined proteins of predicted viruses. Comparison columns, denoted “VIBRANT vs. VirFinder” and “VIBRANT vs. VirSorter”, display the comparison ratio of the given metric; green indicates greater performance by VIBRANT and red indicates lower performance.

Metagenome	seqs. total (≥1kb)	Seqs. ≥ 4 ORFs	Metric	VIBRANT	VirFinder (score≥0.90)	VIBRANT vs. VirFinder	VirSorter (cat. 1 & 2)	VIBRANT vs. VirSorter
human gut: adenoma	34,883	11,360	total putative viruses	527	604	0.87	284	1.86
			total virus length (bp)	5,234,242	1,696,118	3.09	3,982,292	1.31
			total virus proteins	7,661	2,134	3.59	5,484	1.40
human gut: carcinoma	53,946	18,669	total putative viruses	784	1,329	0.59	450	1.74
			total virus length (bp)	5,611,953	3,500,838	1.60	4,182,862	1.34
			total virus proteins	8,401	4,644	1.81	5,945	1.41
human gut: healthy	42,739	17,079	total putative viruses	565	672	0.84	309	1.83
			total virus length (bp)	5,623,082	2,411,049	2.33	4,512,571	1.25
			total virus proteins	8,202	3,230	2.54	6,127	1.34
thermophilic compost	68,815	21,620	total putative viruses	1,047	878	1.19	383	2.73
			total virus length (bp)	10,253,162	2,238,129	4.58	3,290,654	3.12
			total virus proteins	9,912	2,806	3.53	4,400	2.25
freshwater lake (bog)	79,862	26,832	total putative viruses	5,626	7,567	0.74	1,503	3.74
			total virus length (bp)	34,976,570	25,357,664	1.38	15,436,797	2.27
			total virus proteins	56,120	37,537	1.50	21,280	2.64
* estuary virome	5,247	3,277	total putative viruses	3,141	2,294	1.37	1,121	2.80
			total virus length (bp)	6,591,285	6,478,804	1.02	5,163,674	1.28
			total virus proteins	20,500	12,035	1.70	9,645	2.13

* VIBRANT, VirFinder and VirSorter ran with alternate settings

759 Integrated provirus prediction

760 In many environments, integrated proviruses can account for a substantial portion of the
761 active viral community [90]. Despite this, few tools exist that are capable of identifying both lytic
762 viruses from metagenomic scaffolds as well as proviruses that are integrated into host genomes.
763 To account for this important group of viruses, VIBRANT identifies provirus regions within
764 metagenomic scaffolds or whole genomes. VIBRANT is unique from most provirus prediction
765 tools in that it does not rely on sequence motifs, such as integration sites, and therefore is especially
766 useful for partial metagenomic scaffolds in which neither the provirus nor host region is complete.
767 In addition, this functionality of VIBRANT provides the ability to trim non-viral (i.e., host
768 genome) ends from viral scaffolds. This results in a more correct interpretation of genes that are
769 encoded by the virus and not those that are misidentified as being within the viral genome region.
770 Briefly, VIBRANT identifies proviruses by first identifying and isolating scaffolds and genomes
771 at regions spanning several annotations with low v-scores. These regions were found to be almost
772 exclusive to host genomes. After cutting the original sequence at these regions, a refinement step
773 trims the putative provirus fragment to the first instance of a virus-like annotation to remove

774 leftover host sequence (Figure 4A). The final scaffold fragment is then analyzed by the neural
 775 network similar to non-excised scaffolds.

776 To assess VIBRANT's ability to accurately extract provirus regions we compared its
 777 performance to PHASTER and Prophage Hunter, two programs explicitly built for this task, as

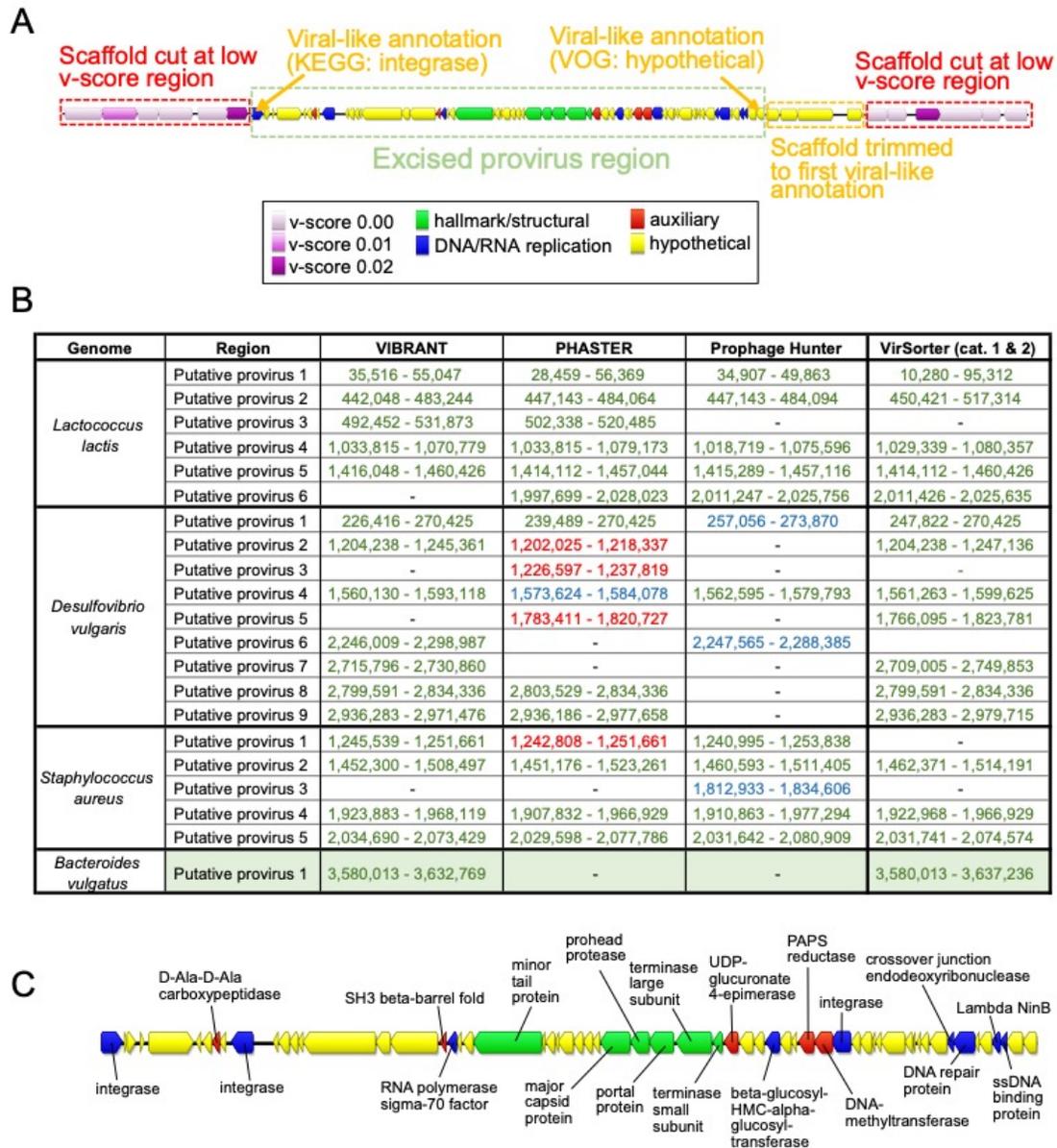


Figure 4. Prediction of integrated proviruses by VIBRANT, and comparison to PHASTER, Prophage Hunter and VirSorter. (A) Schematic representing the method used by VIBRANT to identify and extract provirus regions from host scaffolds using annotations. Briefly, v-scores are used to cut scaffolds at host-specific sites and fragments are trimmed to the nearest viral annotation. (B) Comparison of proviral predictions within four complete bacterial genomes between VIBRANT, PHASTER, Prophage Hunter and VirSorter. For PHASTER, putative proviruses are colored according to “incomplete” (red), “questionable” (blue) and “intact” (green) predictions. Prophage Hunter is colored according to “active” (green) and “ambiguous” (blue) predictions. All VirSorter predictions for categories 1 and 2 are shown in green. (C) Manual validation of the *Bacteroides vulgatus* provirus prediction made by VIBRANT. The presence of viral hallmark protein, integrase and genome replication proteins strongly suggests this is an accurate prediction.

778 well as VirSorter. We compared the performance of these programs with VIBRANT on four
779 complete bacterial genomes. VIBRANT and PHASTER predicted an equal number of proviruses,
780 17, while Prophage Hunter and VirSorter identified slightly less with 13 and 16 identifications,
781 respectively (Figure 4B). Only one putative provirus prediction (*Lactococcus lactis* putative
782 provirus 6) was shared between all programs except VIBRANT. However, VIBRANT was able to
783 identify two putative provirus regions (*Desulfovibrio vulgaris* putative provirus 7 and *Bacteroides*
784 *vulgatus* putative provirus 1) that neither PHASTER nor Prophage Hunter identified, though
785 VirSorter identified these likely due to the similar approach of extracting provirus regions. Manual
786 inspection of the putative *Bacteroides vulgatus* provirus identified a number of virus hallmark and
787 virus-like proteins suggesting that it is an accurate prediction (Figure 4C). Our results suggest
788 VIBRANT has the ability to accurately identify proviruses and, in some cases, can outperform
789 other tools in this task.

790 Both VIBRANT and VirSorter identify integrated proviruses from metagenomic
791 assemblies by cutting host scaffolds at either end of a provirus region. By employing this method
792 these programs generate a more comprehensive understanding of a virome, but errors in identified
793 cut sites may occur due to the diversity of genomic arrangements in both virus and host. This will
794 result in fragmented viral genomes that should have remained intact. We assessed the error rate of
795 VIBRANT and VirSorter (using virome decontamination mode) for cutting viral genomes. A total
796 of 2,466 *Caudovirales* complete genomes were acquired from the NCBI RefSeq database,
797 including 74 megaphages with genomes greater than 200kb. In total, VIBRANT fragmented 5
798 genomes whereas VirSorter fragmented 159 (categories 1 and 2) or 160 (categories 1, 2 and 3).
799 Although relatively comparable, VirSorter incorrectly cut 6.2% more complete viral genomes
800 compared to VIBRANT (6.4% versus 0.2%, respectively).

801

802 **Evaluating quality and completeness of predicted viral sequences**

803 Determination of quality, in relation to completeness, of a predicted viral sequence has
804 been notoriously difficult due to the absence of universally conserved viral genes. To date the most
805 reliable metric of completeness for metagenomically assembled viruses is to identify circular
806 sequences (i.e., complete circular genomes). Therefore, the remaining alternatives rely on
807 estimation based on encoded proteins that function in central viral processes: replication of
808 genomes and assembly of new viral particles.

809 VIBRANT estimates the quality of predicted viral sequences, a relative proxy for
810 completeness, and indicates sequences that are circular. To do this, VIBRANT uses annotation
811 metrics of nucleotide replication and viral hallmark proteins. Hallmark proteins are those typically
812 specific to viruses and those that are required for productive infection, such as structural (e.g.,
813 capsid, tail, baseplate), terminase or viral holin/lysin proteins. Nucleotide replication proteins are
814 a variety of proteins associated with either replication or metabolism, such as nucleases,
815 polymerases and DNA/RNA binding proteins. Viruses are categorized as low, medium or high
816 quality draft as determined by VOG annotations (Figure 5A, Additional File 19: Table S19). High
817 quality draft represents sequences that are likely to contain the majority of a virus's complete
818 genome and will contain annotations that are likely to aid in analysis of the virus, such as
819 phylogenetic relationships and true positive verification. Medium draft quality represents the
820 majority of a complete viral genome but is more likely to be a smaller portion in comparison to
821 high quality. These sequences may contain annotations useful for analysis but are under less strict
822 requirements compared to high quality. Finally, low draft quality constitutes sequences that were
823 not found to be of high or medium quality. Many metagenomic scaffolds will likely be low quality

824 genome fragments, but this quality category may still contain the higher quality genomes of some
 825 highly divergent viruses.

826 We benchmarked VIBRANT's viral genome quality estimation using a total of 2466
 827 *Caudovirales* genomes from NCBI RefSeq database. Genomes were evaluated either as complete
 828 sequences or by removing 10% of the sequence at a time stepwise between 100% and 10%
 829 completeness (Figure 5B). The results of VIBRANT's quality analysis displayed a linear trend in
 830 indicating more complete genomes as high quality and less complete genomes as lower quality.
 831 The transition from categorizing genomes as high quality to medium quality ranged from 60% and
 832 70% completeness. Although we acknowledge that VIBRANT's metrics are not perfect, we

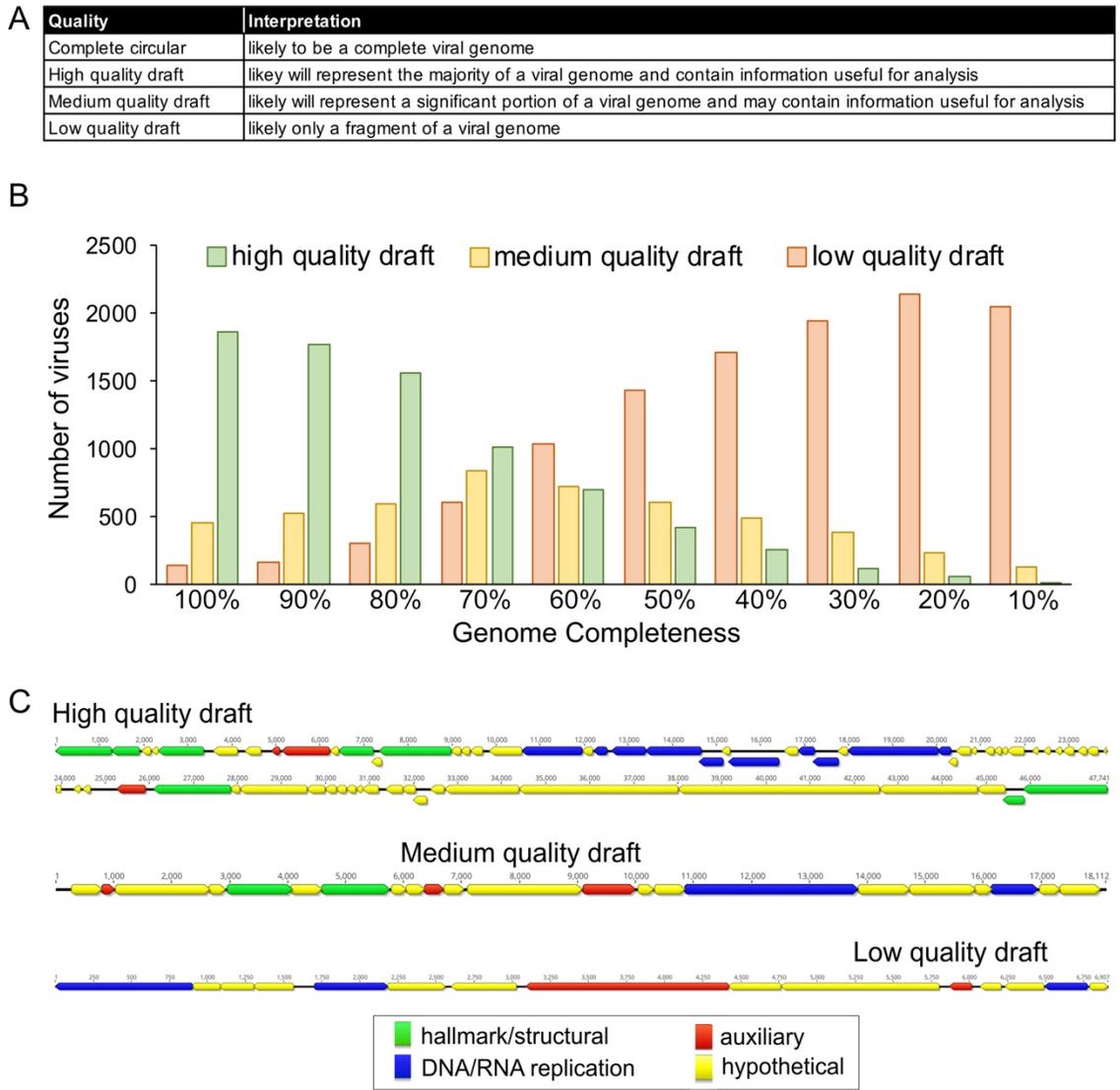


Figure 5. Estimation of genome quality of identified viral scaffolds. (A) Explanation of interpretation of quality categories: complete circular, high quality draft, medium quality draft and low quality draft. Quality generally represents total proteins, viral annotations, viral hallmark protein and nucleotide replication proteins, which are common metrics used for manual verification of viral genomes. (B) Application of quality metrics to 2466 NCBI RefSeq *Caudovirales* viruses with decreasing genome completeness from 100% to 10% completeness, respective of total sequence length. All 2466 viruses are represented within each completeness group. (C) Examples of viral scaffolds representing low, medium and high quality draft categories.

833 demonstrate the first benchmarked approach to quantify and characterize genome quality
 834 associated with completeness of viral sequences. Manual inspection and visual verification of viral
 835 genomes that were characterized into each of these genome quality categories showed that quality
 836 estimations matched annotations (Figure 5C).

837
 838 **Identifying function in viral communities:**
 839 **metabolic analysis**

840 Viruses are a dynamic and key facet in
 841 the metabolic networks of microbial
 842 communities and can reprogram the landscape
 843 of host and community metabolism during
 844 infection. This can often be achieved by
 845 modulating host metabolic networks through
 846 expression of AMGs encoded on viral genomes.
 847 Identifying these AMGs and their associated
 848 role in the function of communities is imperative
 849 for understanding complex microbiome
 850 dynamics, or in some cases can be used to
 851 predict virus-host relationships. VIBRANT is
 852 optimized for the evaluation of viral community
 853 function by identifying and classifying the
 854 metabolic capabilities encoded by a virome. To
 855 do this, VIBRANT identifies AMGs and assigns
 856 them into specific metabolic pathways and
 857 broader categories as designated by KEGG
 858 annotations.

859 To highlight the utility of this feature we
 860 compared the metabolic function of IMG/VR
 861 viruses derived from several diverse
 862 environments: freshwater, marine, soil, human-
 863 associated and city (Additional File 22: Figure
 864 S2). We found natural environments
 865 (freshwater, marine and soil) to display a
 866 different pattern of metabolic capabilities
 867 compared to human environments (human-
 868 associated and city). Viruses originating from
 869 natural environments tend to largely encode
 870 AMGs for amino acid and cofactor/vitamin
 871 metabolism with a more secondary focus on
 872 carbohydrate and glycan metabolism. On the
 873 other hand, AMGs from city and human
 874 environments are dominated by amino acid
 875 metabolism, and to some extent cofactor/
 876 vitamin and sulfur relay metabolism. In
 addition to this broad distinction, all five
 environments appear slightly different from
 each other. Despite freshwater and marine
 environments appearing similar in the ratio
 of AMGs by metabolic category, the overlap
 in specific AMGs is less extensive. The
 dissimilarity between natural and human
 environments is likewise corroborated by
 the relatively low overlap in individual
 AMGs.

877 A useful observation provided by VIBRANT's
 878 metabolic analysis is that there appears to
 be globally conserved AMGs (i.e., present
 within at least 10 of the 12 environments
 tested). These

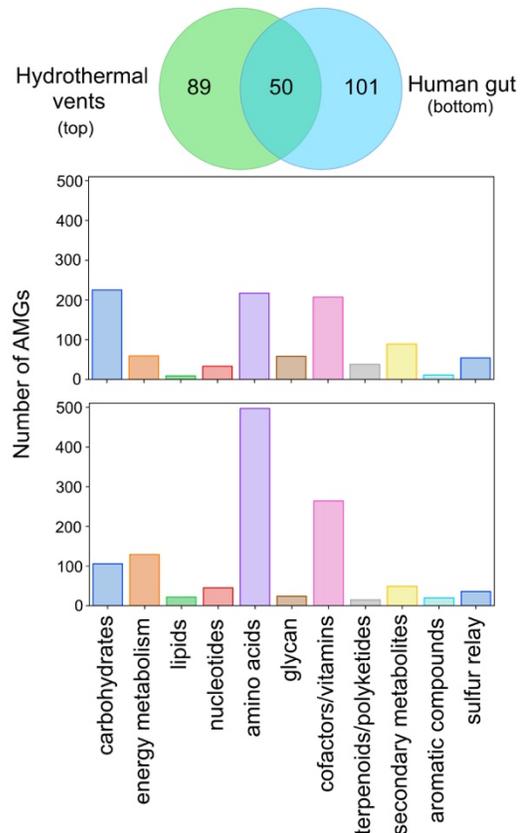


Figure 6. Comparison of AMG metabolic categories between hydrothermal vents and human gut. Venn diagram depicts the unique and shared non-redundant AMGs between 6 hydrothermal vent and 15 human gut metagenomes. The graphs depict the differential abundance of KEGG metabolic categories of respective AMGs for hydrothermal vents (top) and human gut (bottom).

879 14 genes—*dcm*, *cysH*, *folE*, *phnP*, *ubiG*, *ubiE*, *waaF*, *moeB*, *ahbD*, *cobS*, *mec*, *queE*, *queD*,
880 *queC*—likely perform functions that are central to viral replication regardless of host or
881 environment. Notably, *folE*, *queD*, *queE* and *queC* constitute the entire 7-cyano-7-deazaguanine
882 (preQ₀) biosynthesis pathway, but the remainder of queuosine biosynthesis are entirely absent with
883 the exception of *queF*. Certain AMG are unique in that they are the only common representatives
884 of a pathway amongst all AMGs identified, such as *phnP* for methylphosphonate degradation.
885 These AMGs may indicate an evolutionary advantage for manipulating a specific step of a
886 pathway, such as overcoming a reaction bottleneck, as opposed to modulating an entire pathway
887 as seen with preQ₀ biosynthesis. However, it should be noted that this list of 14 globally conserved
888 AMGs may not be entirely inclusive of the core set of AMGs in a given environment.

889 VIBRANT was evaluated for its ability to provide new insights into viral community
890 function by highlighting AMGs from mixed metagenomes. Using only data from VIBRANT's
891 direct outputs, we compared the viral metabolic profiles of 6 hydrothermal vent and 15 human gut
892 metagenomes (Figure 6). As anticipated, based on IMG/VR environment comparisons, the
893 metabolic capabilities between the two environments were different even though the number of
894 unique AMGs was relatively equal (138 for hydrothermal vents and 151 for human gut). The
895 pattern displayed by metabolic categories for each metagenome was similar to that displayed by
896 marine and human viromes. For hydrothermal vents the dominant AMGs were part of
897 carbohydrate, amino acid and cofactor/vitamin metabolism, whereas human gut AMGs were
898 mostly components of amino acid and, to some extent, cofactor/vitamin metabolism. Although the
899 observed AMGs and metabolic pathways were overall different, about a third (50 total AMGs) of
900 all AMGs from each environment were shared; between these metagenomes alone all 14 globally
901 conserved AMGs were present.

902 Observations of individual AMGs provided insights into how viruses interact within
903 different environments. For example, tryptophan 7-halogenase (*prnA*) was identified in high
904 abundance (45 total AMGs) within hydrothermal vent metagenomes but was absent from the
905 human gut. Verification using GOV2 (Global Ocean Viromes 2.0) [91] and Human Gut Virome
906 databases supported our finding that *prnA* appears to be constrained to aquatic environments,
907 which is further supported by the gene's presence on several marine cyanophages. PrnA catalyzes
908 the initial reaction for the formation of pyrrolnitrin, a strong antifungal antibiotic. Identification of
909 this AMG only within aquatic environments suggests a directed role in aquatic virus lifestyles.
910 Similarly, cysteine desulfhydrase (*iscS*) was abundant (14 total AMGs) within the human gut
911 metagenomes but not hydrothermal vents.

912 913 **Application of VIBRANT: Identification of viruses from individuals with Crohn's Disease**

914 We applied VIBRANT to identify viruses of at least 5kb in length from 102 human gut
915 metagenomes (discovery dataset): 49 from individuals with Crohn's Disease and 53 from healthy
916 individuals [73,76]. VIBRANT identified 14,121 viruses out of 511,977 total scaffolds. These viral

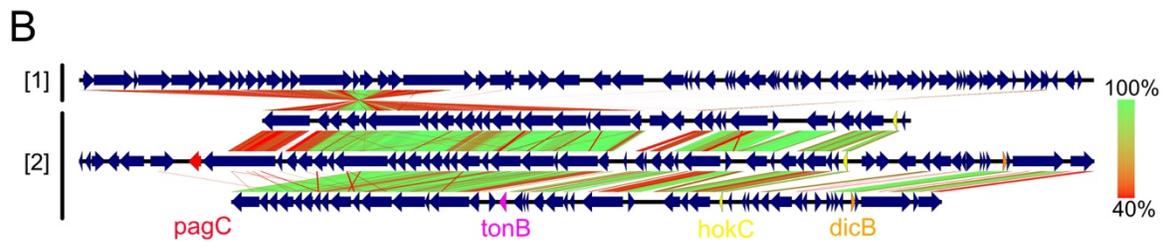
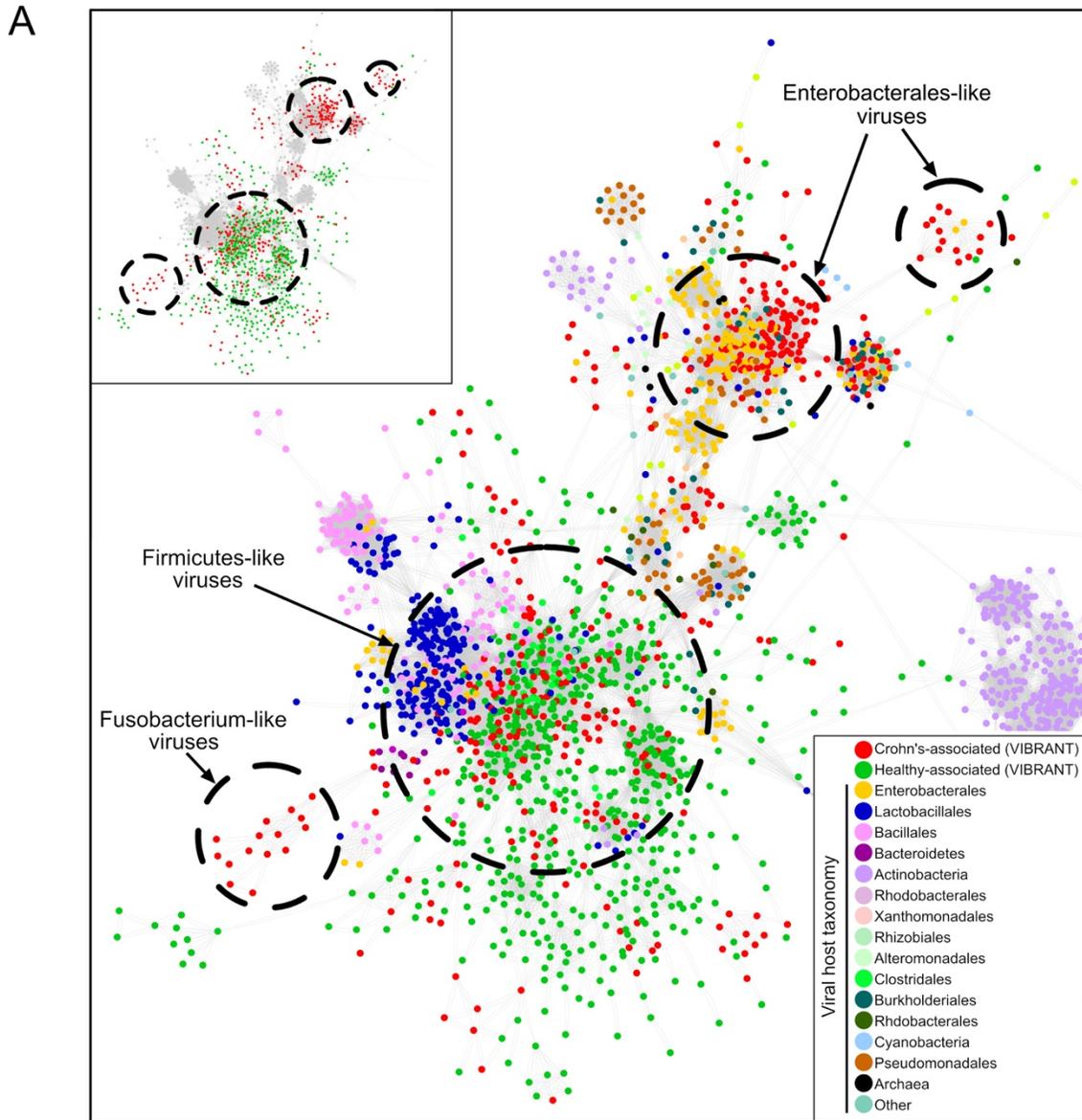


Figure 7. Viral metabolic comparison between Crohn's Disease and healthy individuals gut metagenomes. (A) Partial view of vConTACT2 protein network clustering of viruses identified by VIBRANT and reference viruses. Small clusters and clusters with no VIBRANT representatives are not shown. Each dot represents one genome and is colored according to host or dataset association. Relevant viral groups are indicated by dotted circles (circles enclose estimated boundaries). (B) tBLASTx similarity comparison between (1) Escherichia phage Lambda and (2) three Crohn's-associated viruses identified by VIBRANT. Putative virulence genes are indicated: *pagC*, *tonB*, *hokC* and *dicB*.

918 sequences were dereplicated to 8,822 non-redundant viral sequences using a cutoff of 95%
919 nucleotide identity over at least 70% of the sequence. We next used read coverage of each virus
920 sequence from all 102 metagenomes to calculate relative differential abundance across Crohn's
921 Disease and healthy individuals. In total, we found 721 viral sequences to be more abundant in the
922 gut microbiomes associated with Crohn's Disease (Crohn's-associated) and 950 to be more
923 abundant in healthy individuals (healthy-associated).

924 Using these viruses identified by VIBRANT we sought to identify taxonomic or host-
925 association relationships to differentiate the viral communities of individuals with Crohn's
926 Disease. We used vConTACT2 to cluster the 721 Crohn's- or 950 healthy-associated virus
927 sequences with reference genomes using protein similarity. The majority of virus sequences
928 (95.5%) were not clustered with any reference genome at approximately the genus level suggesting
929 VIBRANT may have identified a large pool of novel or unique viral genomes. Although fewer
930 total viruses were associated with Crohn's Disease, significantly more were clustered to at least
931 one representative at the genus level (72 for Crohn's and 4 for healthy). Interestingly, no
932 differentially abundant viruses from healthy individuals clustered with Enterobacterales-infecting
933 reference viruses (enteroviruses), yet the majority (60/76) of Crohn's-associated viruses were
934 clustered with known enteroviruses, such as Lambda- and Shigella-related viruses. The remaining
935 16 viruses mainly clustered with *Caudovirales* infecting *Lactococcus*, *Clostridium*, *Riemerella*,
936 *Klebsiella* and *Salmonella* species, though *Microviridae* and a likely complete crAssphage were
937 also identified. A significant proportion of all Crohn's-associated viruses (250/721), and the
938 majority of genus-level clustered viruses (42/76), were found to be integrated sequences within a
939 microbial genomic scaffold but were able to be identified due to VIBRANT's ability to excise
940 proviruses.

941 We also generated a protein sharing network containing all 721 Crohn's and 950 healthy-
942 associated virus sequences, which corresponded to taxonomic and host relatedness (Figure 7A).
943 This protein network identified two different clustering patterns: (1) overlapping Crohn's and
944 healthy-associated viral populations clustered with Firmicutes-like viruses which may be
945 indicative of a stable gut viral community; (2) Crohn's-associated viruses clustered with
946 Enterobacterales-like and Fusobacterium-like viruses which may be indicative of a state of
947 dysbiosis. The presence of a greater diversity and abundance of Enterobacterales and Fusobacteria
948 has previously been linked to Crohn's Disease [92,93], and therefore the presence of viruses
949 infecting these bacteria may provide similar information.

950 VIBRANT provides annotation information for all of the identified viruses which can be
951 used to infer functional characteristics in conjunction with host association. Comparison of
952 Crohn's-associated Lambda-like virus genomic content and arrangement suggested a possible role
953 of virally encoded host-persistence and virulence genes that are absent in the healthy-associated
954 virome (Figure 7B). Among all Crohn's-associated viruses, 17 total genes (*bor*, *dicB*, *dicC*, *hokC*,
955 *kilR*, *pagC*, *ydaS*, *ydaT*, *yfdN*, *yfdP*, *yfdQ*, *yfdR*, *yfdS*, *yfdT*, *ymfL*, *ymfM* and *tonB*) that have the
956 potential to impact host survival or virulence were identified. Importantly, no healthy-associated
957 viruses encoded such genes (Table 2). The presence of these putative dysbiosis-associated genes
958 (DAGs) may contribute to the manifestation and/or persistence of disease, similar to what has been
959 proposed for the bacterial microbiome [94–96]. For example, *pagC* encodes an outer membrane
960 virulence factor associated with enhanced survival of the host bacterium within the gut [97]. The
961 identification of *dicB* encoded on a putative *Escherichia* virus is unique in that it may represent a
962 'cryptic' provirus that protects the host from lytic viral infection, thus likely to enhance the ability

963 of the host to survive within the gut [98]. Finally, *hokC* may indicate mechanisms of virally
 964 encoded virulence [99].

965 To characterize the distribution and association of DAGs with Crohn's Disease, we
 966 calculated differential abundance for two highly abundant DAG-encoding viruses across all
 967 metagenome samples. The first virus encoded *pagC* and *yfdN*, and the second encoded *dicB*, *dicC*
 968 and *hokC*. Comparison of Crohn's Disease to healthy metagenomes indicates these viruses are
 969 present within the gut metagenomes of multiple individuals but more abundant in association with
 970 Crohn's Disease (Additional File 23: Figure S3A). This suggests an association of disease with
 971 not only putative DAGs, but also specific, and potentially persistent, viral groups that encode them.
 972 In order to correlate increased abundance with biological activity we calculated the index of
 973 replication (iRep) for each of the two viruses [81]. Briefly, iRep is a function of differential read
 974 coverage which is able to provide an estimate of active genome replication. Seven metagenomes
 975 containing the greatest abundance for each virus were selected for iRep analysis and indicated that
 976 each virus was likely active at the time of collection (Additional File 23: Figure S3B).

Table 2. Identification of putative DAGs encoded by Crohn's-associated viruses. The differential abundance between Crohn's Disease and healthy metagenomes of 17 putative DAGs. Abundance of each gene represents non-redundant annotations, or total gene copy number, from Crohn's-associated and healthy-associated viruses.

ID	Gene	Name	Crohn's Disease	Healthy
PF06291.11	<i>bor</i>	Bor protein	8	0
K22304	<i>dicB</i>	cell division inhibition protein	8	0
K22302	<i>dicC</i>	transcriptional repressor of cell division inhibition gene dicB	18	0
K18919	<i>hokC</i>	protein HokC/D	16	0
VOG11478	<i>kilR</i>	Killing protein	15	0
K07804	<i>pagC</i>	putative virulence related protein	13	0
PF15943.5	<i>ydaS</i>	Putative antitoxin of bacterial toxin-antitoxin system	22	0
PF06254.11	<i>ydaT</i>	Putative bacterial toxin	18	0
VOG04806	<i>yfdN</i>	Uncharacterized protein	19	0
VOG01357	<i>yfdP</i>	Uncharacterized protein	11	0
VOG11472	<i>yfdQ</i>	Uncharacterized protein	11	0
VOG01639	<i>yfdR</i>	Uncharacterized protein	17	0
VOG01103	<i>yfdS</i>	Uncharacterized protein	18	0
VOG16442	<i>yfdT</i>	Uncharacterized protein	8	0
VOG00672	<i>ymfL</i>	Uncharacterized protein	25	0
VOG21507	<i>ymfM</i>	Uncharacterized protein	9	0
K03832	<i>tonB</i>	periplasmic protein	3	0

977 To validate these aforementioned findings, we applied VIBRANT to two additional
 978 metagenomic datasets from cohorts of individuals with Crohn's disease and healthy individuals
 979 (validation dataset): 43 from individuals with Crohn's Disease and 21 from healthy individuals
 980 [74,75]. VIBRANT identified 3,759 redundant viral genomes from Crohn's-associated
 981 metagenomes and 1,444 from healthy-associated metagenomes. Determination of protein
 982 networks and visualization similarly identified clustering of Crohn's-associated viruses with
 983 reference enteroviruses (Additional File 24: Figure S4). Likewise, we were able to identify 15 out
 984 of the 17 putative DAGs to be present in higher abundance in the Crohn's Disease microbiome
 985 (Additional File 20: Table S20). This validates our findings of the presence of unique viruses and
 986 proteins associated with Crohn's Disease, and suggests Enterobacteriales-like viruses and putative

987 DAGs may act as markers of Crohn's Disease. Overall, our results suggest that VIBRANT
988 provides a platform for characterizing these relationships.
989

990 **Discussion**

991 Viruses that infect bacteria and archaea are key components in the structure, dynamics, and
992 interactions of microbial communities [2,6,10,91,100]. Tools that are capable of efficient recovery
993 of these viral genomes from mixed metagenomic samples are likely to be fundamental to the
994 growing applications of metagenomic sequencing and analyses. Importantly, such tools would
995 need to reduce bias associated with specific viral groups (e.g., *Caudovirales*) and highly
996 represented environments (e.g., marine). Moreover, viruses that exist as integrated proviruses
997 within host genomes should not be ignored as they can represent a substantial fraction of infections
998 in certain conditions and also persistent infections within a community [90].
999

1000 Here we have presented VIBRANT, a newly described method for the automated recovery
1001 of both free and integrated viral genomes from metagenomes that hybridizes neural network
1002 machine learning and protein signatures. VIBRANT utilizes metrics of non-reference based
1003 protein similarity annotation from KEGG, Pfam and VOG databases in conjunction with a unique
1004 'v-score' metric to recover viruses with little to no biases. VIBRANT was built with the
1005 consideration of the human guided intuition used to manually inspect metagenomic scaffolds for
1006 viral genomes and packages these ideas into an automated software. This platform originates from
1007 the notion that proteins generally considered as non-viral, such as ribosomal proteins [101], may
1008 be decidedly common amongst viruses and should be considered accordingly when viewing
1009 annotations. V-scores are meant to provide a quantitative metric for the level of virus-association
1010 for each annotation used by VIBRANT, especially for Pfam and KEGG HMMs. That is, v-scores
1011 provide a means for both highlighting common or hallmark viral proteins as well as differentiating
1012 viral from non-viral annotations. In addition, v-scores give a quantifiable value to viral hallmark
1013 genes instead of categorizing them in a binary fashion.

1014 VIBRANT was not only built for the recovery of viral genomes, but also to act as a platform
1015 for investigating the function of a viral community. VIBRANT supports the analysis of viromes
1016 by assembling useful annotation data and categorizing the metabolic pathways of viral AMGs.
1017 Using annotation signatures, VIBRANT furthermore is capable of estimating genome quality and
1018 distinguishing between lytic and lysogenic viruses. To our knowledge, VIBRANT is the first
1019 software that integrates virus identification, annotation and estimation of genome completeness
1020 into a stand-alone program.

1021 Benchmarking and validation of VIBRANT indicated improved performance compared to
1022 VirSorter [45], VirFinder [49] and MARVEL [53], three commonly used programs for identifying
1023 viruses from metagenomes. This included a substantial increase in the relationship between true
1024 virus identifications (recall, true positive rate) and false non-virus identifications (specificity, false
1025 positive rate). That is, VIBRANT recovered more viruses with no discernable expense to false
1026 identifications. The result was that VIBRANT was able to recover an average of 2.3 and 1.7 more
1027 viral sequence from real metagenomes than VirFinder and VirSorter, respectively. When tested on
1028 metagenome-assembled viral genomes from IMG/VR [66] representing diverse environments
1029 VIBRANT was found to have no perceivable environment bias towards identifying viruses. In
1030 comparison to provirus prediction tools, specifically PHASTER [55], Prophage Hunter [56] and
1031 VirSorter, VIBRANT was shown to be proficient in identifying viral regions within bacterial
1032 genomes. This included the identification of a putative *Bacteroides* provirus that PHASTER and
1033 Prophage Hunter were unable to identify. The importance of integrated provirus prediction was

1033 underscored in the analysis of Crohn’s Disease metagenomes since it was found that a significant
1034 proportion of disease related viruses were temperate viruses existing as host-integrated genomes.

1035 VIBRANT’s method allows for the distinction between scaffold size and coding capacity
1036 in designating the minimum length of virus identifications. Traditionally, a cutoff of 5000 bp has
1037 been used to filter for scaffolds of a sufficient length for analysis. This is under the presumption
1038 that a longer sequence will be likely to encode more proteins. For example, this cutoff has been
1039 adopted by IMG/VR. However, we suggest a total protein cutoff of four open reading frames rather
1040 than sequence length cutoff to be more suitable for comprehensive characterization of the viral
1041 community. VIBRANT’s method works as a strict function of total encoded proteins and is
1042 completely agnostic to sequence length for analysis. Therefore, the boundary of minimum encoded
1043 proteins will support a more guided cutoff for quality control of virus identifications. For example,
1044 increasing the minimum sequence length to 5000 bp will have no effect on accuracy or ability to
1045 recall viruses since VIBRANT will only be considerate of the minimum total proteins, which is
1046 set to four. The result will be the loss of all 1000 bp to 4999 bp viruses that still encode at least
1047 four proteins. To visualize this distinction, we applied VIBRANT with various length cutoffs to
1048 the previously used estuary virome (see Table 1). Input sequences were stepwise limited from
1049 1000 bp to 10000 bp (1000 bp steps) or four open reading frames to 13 open reading frames (one
1050 open reading frame steps) in length. Limiting to open reading frames indicated a reduced drop-off
1051 in total virus identifications and total viral sequence compared to a minimum sequence length limit
1052 (Additional File 25: Figure S5).

1053 The output data generated by VIBRANT—protein/gene annotation information,
1054 protein/gene sequences, HMM scores and e-values, viral sequences in FASTA and GenBank
1055 format, indication of AMGs, genome quality, etc.—provides a platform for easily replicated
1056 pipeline analyses. Application of VIBRANT to characterize the function of Crohn’s-associated
1057 viruses emphasizes this utility. VIBRANT was not only able to identify a substantial number of
1058 viral genomes, but also provided meaningful information regarding putative DAGs, viral
1059 sequences for differential abundance calculation and genome alignment, viral proteins for
1060 clustering, and AMGs for metabolic comparisons.

1061 1062 **Conclusions**

1063 Our construction of the VIBRANT platform expands the current potential for virus
1064 identification and characterization from metagenomic and genomic sequences. When compared to
1065 two widely used software programs, VirFinder and VirSorter, we show that VIBRANT improves
1066 total viral sequence and protein recovery from diverse human and natural environments. As
1067 sequencing technologies improve and metagenomic datasets contain longer sequences VIBRANT
1068 will continue to outcompete programs built for short scaffolds (e.g., 500-3000 bp) by identifying
1069 more higher quality genomes. Our workflow, through the annotation of viral genomes, aids in the
1070 capacity to discover how viruses of bacteria and archaea may shape an environment, such as
1071 driving specific metabolism during infection or dysbiosis in the human gut. Furthermore,
1072 VIBRANT is the first virus identification software to incorporate annotation information into the
1073 curation of predictions, estimation of genome quality and infection mechanism (i.e., lytic vs
1074 lysogenic). We anticipate that the incorporation of VIBRANT into microbiome analyses will
1075 provide easy interpretation of viral data, enabled by VIBRANT’s comprehensive functional
1076 analysis platform and visualization of information.

1077 1078 **Acknowledgements**

1079 We thank Upendra Devisetty for his assistance with dockerizing and integrating VIBRANT as a
1080 web-based application in the CyVerse Discovery Environment. We thank the University of
1081 Wisconsin - Office of the Vice Chancellor for Research and Graduate Education, University of
1082 Wisconsin – Department of Bacteriology, and University of Wisconsin – College of Agriculture
1083 and Life Sciences for their support.
1084

1085 **Availability of data and materials**

1086 VIBRANT is implemented in Python and all scripts and associated files are freely available
1087 at <https://github.com/AnantharamanLab/VIBRANT/>. All data and genomic sequences used for
1088 analyses are publicly available; see Supplementary Tables S1 and S14 for study and accession
1089 names. VIBRANT is also freely available for use as an application through the CyVerse Discovery
1090 Environment; to use the application visit <https://de.cyverse.org/de/>. Additional details of relevant
1091 data are available from the corresponding author on request.
1092

1093 **Author Information**

1094 **Affiliations**

1095 *Department of Bacteriology, University of Wisconsin–Madison, Madison, WI, USA*
1096 Kristopher Kieft, Zhichao Zhou, and Karthik Anantharaman
1097

1098 **Contributors**

1099 K.K and K.A designed the study, performed all analyses and interpretation of data, and wrote the
1100 manuscript. Z.Z contributed to conceptualization of study design and reviewed the manuscript. All
1101 authors have reviewed and approved the final manuscript.
1102

1103 **Corresponding Author**

1104 Correspondence to Karthik Anantharaman
1105

1106 **Ethics Declarations**

1107 **Ethics approval and consent to participate**

1108 Not applicable.
1109

1110 **Consent for publication**

1111 Not applicable.
1112

1113 **Competing interests**

1114 The authors declare that they have no competing interests.
1115

1116 **Additional Files**

1117
1118 **Additional File 1: Table S1.** List of NCBI accession numbers for bacterial and archaeal genomes,
1119 plasmids, and viral genomes used in this study.
1120

1121 **Additional File 2: Table S2.** Number and sizes of sequence fragments used to train and test
1122 VIBRANT for viruses, plasmids, and bacteria and archaea.
1123

1124 **Additional File 3: Table S3.** List of all HMM names used by VIBRANT.

1125
1126 **Additional File 4: Table S4.** List of all KEGG, Pfam and VOG annotation names and associated
1127 v-scores (if greater than zero).
1128
1129 **Additional File 5: Table S5.** Unparsed HMM table output from KEGG annotations used to
1130 generate KEGG v-scores.
1131
1132 **Additional File 6: Table S6.** Unparsed HMM table output from Pfam annotations used to generate
1133 Pfam v-scores.
1134
1135 **Additional File 7: Table S7.** Unparsed HMM table output from VOG annotations used to generate
1136 VOG v-scores.
1137
1138 **Additional File 8: Table S8.** Description of set cutoffs implemented before neural network
1139 machine learning analysis for KEGG and Pfam annotations.
1140
1141 **Additional File 9: Table S9.** Lists of KEGG, Pfam and VOG annotations used to generate
1142 annotation metrics for neural network classification. Designated by asterisks are the lists of all
1143 VOG annotations determined as nucleotide replication-associated or viral hallmark-associated,
1144 which are used during prediction and quality estimation.
1145
1146 **Additional File 10: Table S10.** Normalized data used to train the neural network machine learning
1147 classifier.
1148
1149 **Additional File 11: Table S11.** Normalized data used to test the neural network machine learning
1150 classifier.
1151
1152 **Additional File 12: Table S12.** Equations used for benchmarking analyses.
1153
1154 **Additional File 13: Table S13.** Calculations and results of benchmarking analyses for VIBRANT,
1155 VirSorter, VirFinder and MARVEL.
1156
1157 **Additional File 14: Table S14.** List of datasets used from He *et al.*, Ijaz *et al.* and Gevers *et al.*.
1158
1159 **Additional File 15: Table S15.** VIBRANT runtimes and resource requirements for datasets of
1160 various sizes and compositions.
1161
1162 **Additional File 16: Table S16.** List of all KEGG annotations determined as AMGs.
1163
1164 **Additional File 17: Table S17.** Results from DESeq2 analysis for 8,789 non-redundant viruses
1165 from the Crohn's Disease discovery dataset.
1166
1167 **Additional File 18: Table S18.** Complete VIBRANT-derived annotations of validated viral
1168 scaffolds that encode putative DAGs, respective of Table 2.
1169

1170 **Additional File 19: Table S19.** Number of *Caudovirales* genomes and genomic fragments
1171 identified per quality estimation category, exact rules used to estimate genome quality and the
1172 interpretation of quality estimations.

1173
1174 **Additional File 20: Table S20.** Validation of the greater abundance of viral DAGs in individuals
1175 with Crohn's Disease compared to healthy individuals.

1176
1177 **Additional File 21: Figure S1. Comparison of VIBRANT, VirFinder, VirSorter and**
1178 **MARVEL on additional validation datasets.** (A) The TPR and FPR of virus identifications for
1179 1kb and 3kb scaffolds, (B) the effect of eukaryotic sequence contamination, and the ability to
1180 recover complete (C) RNA and (D) archaeal viruses.

1181
1182 **Additional File 22: Figure S2. AMG and metabolic pathways between diverse environments.**
1183 VIBRANT was used to predict viruses from IMG/VR datasets and the identified metabolic
1184 pathways and AMGs were compared for freshwater, marine, soil, city and human-associated
1185 environments (graphs). The respective AMGs and their abundances were likewise compared (venn
1186 diagram).

1187
1188 **Additional File 23: Figure S3. Differential abundance and activity of two viruses associated**
1189 **with Crohn's Disease.** (A) Normalized read coverage of two abundant Crohn's-associated viruses
1190 that encode putative DAGs between Crohn's Disease and healthy gut metagenomes. Asterisks
1191 represent significant differential abundance ($p < 0.05$). (B) iRep analysis for the same two viruses
1192 as (A), representative of seven metagenomes per virus for which the virus was in high abundance.
1193 The dotted line indicates an iRep value of one, or low to no activity (i.e., genome replication).

1194
1195 **Additional File 24: Figure S4. Protein network of two Crohn's Disease validation datasets.**
1196 VIBRANT was used to predict viruses from two datasets for validation of marker virus and
1197 putative DAG discovery. The resulting viruses were used to construct a protein network indicating
1198 Crohn's-associated viruses clustering with enteroviruses more often than healthy-associated
1199 viruses.

1200
1201 **Additional File 25: Figure S5. Comparison of limiting to sequence length or open reading**
1202 **frames.** VIBRANT was used to predict viruses from an estuary virome and set to limit to either
1203 scaffold length or total encoded open reading frames. The (A) total virus identifications and (B)
1204 total viral sequence length were compared to show that limiting to open reading frames will
1205 typically yield more data.

1206 1207 **References**

- 1208
1209 1. Breitbart M, Rohwer F. Here a virus, there a virus, everywhere the same virus? Trends in
1210 Microbiology. 2005;13:278–84.
- 1211 2. Wommack KE, Colwell RR. Virioplankton: Viruses in Aquatic Ecosystems. Microbiol Mol
1212 Biol Rev. 2000;64:69–114.
- 1213 3. Danovaro R, Serresi M. Viral Density and Virus-to-Bacterium Ratio in Deep-Sea Sediments
1214 of the Eastern Mediterranean. Appl Environ Microbiol. 2000;66:1857–61.

- 1215 4. Suttle CA. Marine viruses — major players in the global ecosystem. *Nature Reviews*
1216 *Microbiology*. 2007;5:801–12.
- 1217 5. Heldal M, Bratbak G. Production and decay of viruses in aquatic environments. *Mar Ecol*
1218 *Prog Ser*. 1991;72:205–12.
- 1219 6. Gobler CJ, Hutchins DA, Fisher NS, Cosper EM, Sañudo-Wilhelmy SA. Release and
1220 bioavailability of C, N, P Se, and Fe following viral lysis of a marine chrysophyte. *Limnology*
1221 *and Oceanography*. 1997;42:1492–504.
- 1222 7. Jiao N, Herndl GJ, Hansell DA, Benner R, Kattner G, Wilhelm SW, et al. Microbial
1223 production of recalcitrant dissolved organic matter: long-term carbon storage in the global ocean.
1224 *Nature Reviews Microbiology*. 2010;8:593–9.
- 1225 8. Brussaard CPD, Wilhelm SW, Thingstad F, Weinbauer MG, Bratbak G, Heldal M, et al.
1226 Global-scale processes with a nanoscale drive: the role of marine viruses. *The ISME Journal*.
1227 2008;2:575–8.
- 1228 9. Fuhrman JA. Marine viruses and their biogeochemical and ecological effects. *Nature*.
1229 1999;399:541–8.
- 1230 10. Wilhelm SW, Suttle CA. Viruses and Nutrient Cycles in the Sea. *BioScience*. 1999;49:8.
- 1231 11. Norman JM, Handley SA, Baldridge MT, Droit L, Liu CY, Keller BC, et al. Disease-Specific
1232 Alterations in the Enteric Virome in Inflammatory Bowel Disease. *Cell*. 2015;160:447–60.
- 1233 12. Barr JJ. Missing a Phage: Unraveling Tripartite Symbioses within the Human Gut.
1234 *mSystems*. 2019;4:e00105-19.
- 1235 13. Barr JJ, Auro R, Furlan M, Whiteson KL, Erb ML, Pogliano J, et al. Bacteriophage adhering
1236 to mucus provide a non-host-derived immunity. *Proceedings of the National Academy of*
1237 *Sciences*. 2013;110:10771–6.
- 1238 14. Rohwer F. Global Phage Diversity. *Cell*. 2003;113:141.
- 1239 15. Kim B, Kim ES, Yoo Y-J, Bae H-W, Chung I-Y, Cho Y-H. Phage-Derived Antibacterials:
1240 Harnessing the Simplicity, Plasticity, and Diversity of Phages. *Viruses* [Internet]. 2019 [cited
1241 2019 Oct 24];11. Available from: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6466130/>
- 1242 16. Peng S-Y, You R-I, Lai M-J, Lin N-T, Chen L-K, Chang K-C. Highly potent antimicrobial
1243 modified peptides derived from the *Acinetobacter baumannii* phage endolysin LysAB2. *Sci Rep*.
1244 2017;7:1–12.
- 1245 17. Holt A, Cahill J, Ramsey J, O’Leary C, Moreland R, Martin C, et al. Phage-encoded cationic
1246 antimicrobial peptide used for outer membrane disruption in lysis. *bioRxiv*. 2019;515445.

- 1247 18. Harada LK, Silva EC, Campos WF, Del Fiol FS, Vila M, Dąbrowska K, et al.
1248 Biotechnological applications of bacteriophages: State of the art. *Microbiological Research*.
1249 2018;212–213:38–58.
- 1250 19. Sharma RS, Karmakar S, Kumar P, Mishra V. Application of filamentous phages in
1251 environment: A tectonic shift in the science and practice of ecorestoration. *Ecology and*
1252 *Evolution*. 2019;9:2263–304.
- 1253 20. Jiang SC, Paul JH. Gene Transfer by Transduction in the Marine Environment. *APPL*
1254 *ENVIRON MICROBIOL*. 1998;64:8.
- 1255 21. Gregory AC, Solonenko SA, Ignacio-Espinoza JC, LaButti K, Copeland A, Sudek S, et al.
1256 Genomic differentiation among wild cyanophages despite widespread horizontal gene transfer.
1257 *BMC Genomics*. 2016;17:930.
- 1258 22. Sanjuán R, Nebot MR, Chirico N, Mansky LM, Belshaw R. Viral Mutation Rates. *J Virol*.
1259 2010;84:9733–48.
- 1260 23. Dutilh BE, Cassman N, McNair K, Sanchez SE, Silva GGZ, Boling L, et al. A highly
1261 abundant bacteriophage discovered in the unknown sequences of human faecal metagenomes.
1262 *Nature Communications*. 2014;5:4498.
- 1263 24. Devoto AE, Santini JM, Olm MR, Anantharaman K, Munk P, Tung J, et al. Megaphages
1264 infect *Prevotella* and variants are widespread in gut microbiomes. *Nature Microbiology*.
1265 2019;4:693–700.
- 1266 25. Al-Shayeb B, Sachdeva R, Chen L-X, Ward F, Munk P, Devoto A, et al. Clades of huge
1267 phage from across Earth’s ecosystems. *bioRxiv*. 2019;572362.
- 1268 26. Kauffman KM, Hussain FA, Yang J, Arevalo P, Brown JM, Chang WK, et al. A major
1269 lineage of non-tailed dsDNA viruses as unrecognized killers of marine bacteria. *Nature*; London.
1270 2018;554:118-122,122A-122T.
- 1271 27. Hopkins M, Kailasan S, Cohen A, Roux S, Tucker KP, Shevenell A, et al. Diversity of
1272 environmental single-stranded DNA phages revealed by PCR amplification of the partial major
1273 capsid protein. *ISME J*. 2014;8:2093–103.
- 1274 28. Krishnamurthy SR, Janowski AB, Zhao G, Barouch D, Wang D. Hyperexpansion of RNA
1275 Bacteriophage Diversity. *PLOS Biology*. 2016;14:e1002409.
- 1276 29. Waldbauer JR, Coleman ML, Rizzo AI, Campbell KL, Lotus J, Zhang L. Nitrogen sourcing
1277 during viral infection of marine cyanobacteria. *PNAS*. 2019;116:15590–5.
- 1278 30. Stent GS, Maaløe O. Radioactive phosphorus tracer studies on the reproduction of T4
1279 bacteriophage: II. Kinetics of phosphorus assimilation. *Biochimica et Biophysica Acta*.
1280 1953;10:55–69.

- 1281 31. Kozloff LM, Knowlton K, Putnam FW, Evans EA. Biochemical Studies of Virus
1282 Reproduction V. the Origin of Bacteriophage Nitrogen. *J Biol Chem.* 1951;188:101–16.
- 1283 32. Thompson LR, Zeng Q, Kelly L, Huang KH, Singer AU, Stubbe J, et al. Phage auxiliary
1284 metabolic genes and the redirection of cyanobacterial host carbon metabolism. *PNAS.*
1285 2011;108:E757–64.
- 1286 33. Breitbart M, Thompson L, Suttle C, Sullivan M. Exploring the Vast Diversity of Marine
1287 Viruses. *Oceanography.* 2007;20:135–9.
- 1288 34. Hurwitz BL, Hallam SJ, Sullivan MB. Metabolic reprogramming by viruses in the sunlit and
1289 dark ocean. *Genome Biology.* 2013;14:R123.
- 1290 35. Roux S, Hawley AK, Beltran MT, Scofield M, Schwientek P, Stepanauskas R, et al. Ecology
1291 and evolution of viruses infecting uncultivated SUP05 bacteria as revealed by single-cell- and
1292 meta-genomics. *eLife Sciences.* 2014;3:e03125.
- 1293 36. Bragg JG, Chisholm SW. Modeling the Fitness Consequences of a Cyanophage-Encoded
1294 Photosynthesis Gene. *PLOS ONE.* 2008;3:e3550.
- 1295 37. Mann NH, Cook A, Millard A, Bailey S, Clokie M. Bacterial photosynthesis genes in a virus.
1296 *Nature.* 2003;424:741.
- 1297 38. Anantharaman K, Duhaime MB, Breier JA, Wendt KA, Toner BM, Dick GJ. Sulfur
1298 Oxidation Genes in Diverse Deep-Sea Viruses. *Science.* 2014;344:757–60.
- 1299 39. Emerson JB, Roux S, Brum JR, Bolduc B, Woodcroft BJ, Jang HB, et al. Host-linked soil
1300 viral ecology along a permafrost thaw gradient. *Nature Microbiology.* 2018;3:870.
- 1301 40. Trubl G, Jang HB, Roux S, Emerson JB, Solonenko N, Vik DR, et al. Soil Viruses Are
1302 Underexplored Players in Ecosystem Carbon Processing. *mSystems.* 2018;3:e00076-18.
- 1303 41. John SG, Mendez CB, Deng L, Poulos B, Kauffman AKM, Kern S, et al. A simple and
1304 efficient method for concentration of ocean viruses by chemical flocculation. *Environmental*
1305 *Microbiology Reports.* 2011;3:195–202.
- 1306 42. Göller PC, Haro-Moreno JM, Rodriguez-Valera F, Loessner MJ, Gómez-Sanz E. Uncovering
1307 a hidden diversity: optimized protocols for the extraction of bacteriophages from soil. *bioRxiv.*
1308 2019;733980.
- 1309 43. Labonté JM, Swan BK, Poulos B, Luo H, Koren S, Hallam SJ, et al. Single-cell genomics-
1310 based analysis of virus–host interactions in marine surface bacterioplankton. *ISME J.*
1311 2015;9:2386–99.
- 1312 44. Trubl G, Solonenko N, Chittick L, Solonenko SA, Rich VI, Sullivan MB. Optimization of
1313 viral resuspension methods for carbon-rich soils along a permafrost thaw gradient. *PeerJ.*
1314 2016;4:e1999.

- 1315 45. Roux S, Enault F, Hurwitz BL, Sullivan MB. VirSorter: mining viral signal from microbial
1316 genomic data. *PeerJ*. 2015;3:e985.
- 1317 46. Wommack KE, Bhavsar J, Polson SW, Chen J, Dumas M, Srinivasiah S, et al. VIROME: a
1318 standard operating procedure for analysis of viral metagenome sequences. *Standards in Genomic
1319 Sciences*. 2012;6:427.
- 1320 47. Roux S, Faubladiere M, Mahul A, Paulhe N, Bernard A, Debross D, et al. Metavir: a web
1321 server dedicated to virome analysis. *Bioinformatics*. 2011;27:3074–5.
- 1322 48. El-Gebali S, Mistry J, Bateman A, Eddy SR, Luciani A, Potter SC, et al. The Pfam protein
1323 families database in 2019. *Nucleic Acids Res*. 2019;47:D427–32.
- 1324 49. Ren J, Ahlgren NA, Lu YY, Fuhrman JA, Sun F. VirFinder: a novel k-mer based tool for
1325 identifying viral sequences from assembled metagenomic data. *Microbiome*. 2017;5:69.
- 1326 50. Fang Z, Tan J, Wu S, Li M, Xu C, Xie Z, et al. PPR-Meta: a tool for identifying phages and
1327 plasmids from metagenomic fragments using deep learning. *Gigascience* [Internet]. 2019 [cited
1328 2019 Aug 5];8. Available from:
1329 <https://academic.oup.com/gigascience/article/8/6/giz066/5521157>
- 1330 51. Ahlgren NA, Ren J, Lu YY, Fuhrman JA, Sun F. Alignment-free d_2 oligonucleotide
1331 frequency dissimilarity measure improves prediction of hosts from metagenomically-derived
1332 viral sequences. *Nucleic Acids Res*. 2017;45:39–53.
- 1333 52. Ponsero AJ, Hurwitz BL. The Promises and Pitfalls of Machine Learning for Detecting
1334 Viruses in Aquatic Metagenomes. *Front Microbiol* [Internet]. 2019 [cited 2019 Oct 24];10.
1335 Available from: <https://www.frontiersin.org/articles/10.3389/fmicb.2019.00806/full>
- 1336 53. Amgarten D, Braga LPP, da Silva AM, Setubal JC. MARVEL, a Tool for Prediction of
1337 Bacteriophage Sequences in Metagenomic Bins. *Front Genet* [Internet]. 2018 [cited 2019 Aug
1338 5];9. Available from: <https://www.frontiersin.org/articles/10.3389/fgene.2018.00304/full>
- 1339 54. Zheng T, Li J, Ni Y, Kang K, Misiakou M-A, Imamovic L, et al. Mining, analyzing, and
1340 integrating viral signals from metagenomic data. *Microbiome*. 2019;7:42.
- 1341 55. Arndt D, Grant JR, Marcu A, Sajed T, Pon A, Liang Y, et al. PHASTER: a better, faster
1342 version of the PHAST phage search tool. *Nucleic Acids Res*. 2016;44:W16–21.
- 1343 56. Song W, Sun H-X, Zhang C, Cheng L, Peng Y, Deng Z, et al. Prophage Hunter: an
1344 integrative hunting tool for active prophages. *Nucleic Acids Res*. 2019;47:W74–80.
- 1345 57. Merchant N, Lyons E, Goff S, Vaughn M, Ware D, Micklos D, et al. The iPlant
1346 Collaborative: Cyberinfrastructure for Enabling Data to Discovery for the Life Sciences. *PLOS
1347 Biology*. 2016;14:e1002342.
- 1348 58. Hyatt D, Chen G-L, LoCasio PF, Land ML, Larimer FW, Hauser LJ. Prodigal: prokaryotic
1349 gene recognition and translation initiation site identification. *BMC Bioinformatics*. 2010;11:119.

- 1350 59. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment search tool. *J*
1351 *Mol Biol.* 1990;215:403–10.
- 1352 60. Eddy SR. Profile hidden Markov models. *Bioinformatics.* 1998;14:755–63.
- 1353 61. Kanehisa M, Goto S. KEGG: Kyoto Encyclopedia of Genes and Genomes. *Nucleic Acids*
1354 *Res.* 2000;28:27–30.
- 1355 62. Aramaki T, Blanc-Mathieu R, Endo H, Ohkubo K, Kanehisa M, Goto S, et al.
1356 KofamKOALA: KEGG ortholog assignment based on profile HMM and adaptive score
1357 threshold. *bioRxiv.* 2019;602110.
- 1358 63. Fu L, Niu B, Zhu Z, Wu S, Li W. CD-HIT: accelerated for clustering the next-generation
1359 sequencing data. *Bioinformatics.* 2012;28:3150–2.
- 1360 64. Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, et al. Scikit-learn:
1361 Machine Learning in Python. *Journal of Machine Learning Research.* 2011;12:2825–30.
- 1362 65. Paez-Espino D, Eloie-Fadrosh EA, Pavlopoulos GA, Thomas AD, Huntemann M, Mikhailova
1363 N, et al. Uncovering Earth’s virome. *Nature.* 2016;536:425–30.
- 1364 66. Paez-Espino D, Roux S, Chen I-MA, Palaniappan K, Ratner A, Chu K, et al. IMG/VR v.2.0:
1365 an integrated data management and analysis system for cultivated and environmental viral
1366 genomes. *Nucleic Acids Res.* 2019;47:D678–86.
- 1367 67. Hunter JD. Matplotlib: A 2D graphics environment. *Computing In Science & Engineering.*
1368 2007;9:90–5.
- 1369 68. Gregory AC, Zablocki O, Howell A, Bolduc B, Sullivan MB. The human gut virome
1370 database. *bioRxiv.* 2019;655910.
- 1371 69. Feng Q, Liang S, Jia H, Stadlmayr A, Tang L, Lan Z, et al. Gut microbiome development
1372 along the colorectal adenoma–carcinoma sequence. *Nature Communications.* 2015;6:1–13.
- 1373 70. Joshi N, Fass J. Sickle: A sliding-window, adaptive, quality-based trimming tool for FastQ
1374 files [Internet]. 2011. Available from: <https://github.com/najoshi/sickle>
- 1375 71. Nurk S, Meleshko D, Korobeynikov A, Pevzner PA. metaSPAdes: a new versatile
1376 metagenomic assembler. *Genome Res.* 2017;27:824–34.
- 1377 72. Zhou Z, Tran PQ, Kieft K, Anantharaman K. Genome diversification in globally distributed
1378 novel marine Proteobacteria is linked to environmental adaptation. *bioRxiv.* 2019;814418.
- 1379 73. He Q, Gao Y, Jie Z, Yu X, Laursen JM, Xiao L, et al. Two distinct metacommunities
1380 characterize the gut microbiota in Crohn’s disease patients. *Gigascience.* 2017;6:1–11.

- 1381 74. Ijaz UZ, Quince C, Hanske L, Loman N, Calus ST, Bertz M, et al. The distinct features of
1382 microbial “dysbiosis” of Crohn’s disease do not occur to the same extent in their unaffected,
1383 genetically-linked kindred. *PLoS ONE*. 2017;12:e0172605.
- 1384 75. Gevers D, Kugathasan S, Denson LA, Vázquez-Baeza Y, Van Treuren W, Ren B, et al. The
1385 treatment-naïve microbiome in new-onset Crohn’s disease. *Cell Host Microbe*. 2014;15:382–92.
- 1386 76. Pasolli E, Asnicar F, Manara S, Zolfo M, Karcher N, Armanini F, et al. Extensive
1387 Unexplored Human Microbiome Diversity Revealed by Over 150,000 Genomes from
1388 Metagenomes Spanning Age, Geography, and Lifestyle. *Cell*. 2019;176:649-662.e20.
- 1389 77. Ondov BD, Treangen TJ, Melsted P, Mallonee AB, Bergman NH, Koren S, et al. Mash: fast
1390 genome and metagenome distance estimation using MinHash. *Genome Biology*. 2016;17:132.
- 1391 78. Delcher AL. Fast algorithms for large-scale genome alignment and comparison. *Nucleic
1392 Acids Research*. 2002;30:2478–83.
- 1393 79. Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. *Nat Methods*.
1394 2012;9:357–9.
- 1395 80. Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for RNA-
1396 seq data with DESeq2. *Genome Biology*. 2014;15:550.
- 1397 81. Brown CT, Olm MR, Thomas BC, Banfield JF. Measurement of bacterial replication rates in
1398 microbial communities. *Nature Biotechnology*. 2016;34:1256–63.
- 1399 82. Sullivan MJ, Petty NK, Beatson SA. Easyfig: a genome comparison visualizer.
1400 *Bioinformatics*. 2011;27:1009–10.
- 1401 83. Shannon P. Cytoscape: A Software Environment for Integrated Models of Biomolecular
1402 Interaction Networks. *Genome Research*. 2003;13:2498–504.
- 1403 84. Kristensen DM, Waller AS, Yamada T, Bork P, Mushegian AR, Koonin EV. Orthologous
1404 Gene Clusters and Taxon Signature Genes for Viruses of Prokaryotes. *Journal of Bacteriology*.
1405 2013;195:941–50.
- 1406 85. Graziotin AL, Koonin EV, Kristensen DM. Prokaryotic Virus Orthologous Groups
1407 (pVOGs): a resource for comparative genomics and protein family annotation. *Nucleic Acids
1408 Res*. 2017;45:D491–8.
- 1409 86. Hendricks SP, Mathews CK. Regulation of T4 Phage Aerobic Ribonucleotide Reductase:
1410 SIMULTANEOUS ASSAY OF THE FOUR ACTIVITIES. *J Biol Chem*. 1997;272:2861–5.
- 1411 87. Roux S, Adriaenssens EM, Dutilh BE, Koonin EV, Kropinski AM, Krupovic M, et al.
1412 Minimum Information about an Uncultivated Virus Genome (MIUViG). *Nature Biotechnology*.
1413 2019;37:29–37.

- 1414 88. Bowers RM, Kyrpides NC, Stepanauskas R, Harmon-Smith M, Doud D, Reddy TBK, et al.
1415 Minimum information about a single amplified genome (MISAG) and a metagenome-assembled
1416 genome (MIMAG) of bacteria and archaea. *Nature Biotechnology*. 2017;35:725–31.
- 1417 89. Tucker KP, Parsons R, Symonds EM, Breitbart M. Diversity and distribution of single-
1418 stranded DNA phages in the North Atlantic Ocean. *ISME J*. 2011;5:822–30.
- 1419 90. Payet JP, Suttle CA. To kill or not to kill: The balance between lytic and lysogenic viral
1420 infection is driven by trophic status. *Limnology and Oceanography*. 2013;58:465–74.
- 1421 91. Gregory AC, Zayed AA, Conceição-Neto N, Temperton B, Bolduc B, Alberti A, et al.
1422 Marine DNA Viral Macro- and Microdiversity from Pole to Pole. *Cell* [Internet]. 2019 [cited
1423 2019 Apr 30]; Available from:
1424 <http://www.sciencedirect.com/science/article/pii/S0092867419303411>
- 1425 92. Morgan XC, Tickle TL, Sokol H, Gevers D, Devaney KL, Ward DV, et al. Dysfunction of
1426 the intestinal microbiome in inflammatory bowel disease and treatment. *Genome Biology*.
1427 2012;13:R79.
- 1428 93. Strauss J, Kaplan GG, Beck PL, Rioux K, Panaccione R, Devinney R, et al. Invasive
1429 potential of gut mucosa-derived *Fusobacterium nucleatum* positively correlates with IBD status
1430 of the host. *Inflamm Bowel Dis*. 2011;17:1971–8.
- 1431 94. Shreiner AB, Kao JY, Young VB. The gut microbiome in health and in disease. *Curr Opin*
1432 *Gastroenterol*. 2015;31:69–75.
- 1433 95. Clemente JC, Ursell LK, Parfrey LW, Knight R. The Impact of the Gut Microbiota on
1434 Human Health: An Integrative View. *Cell*. 2012;148:1258–70.
- 1435 96. Minot SS, Willis AD. Clustering co-abundant genes identifies components of the gut
1436 microbiome that are reproducibly associated with colorectal cancer and inflammatory bowel
1437 disease. *Microbiome*. 2019;7:110.
- 1438 97. Nishio M, Okada N, Miki T, Haneda T, Danbara H. Identification of the outer-membrane
1439 protein PagC required for the serum resistance phenotype in *Salmonella enterica* serovar
1440 *Choleraesuis*. *Microbiology (Reading, Engl)*. 2005;151:863–73.
- 1441 98. Rangunathan PT, Vanderpool CK. Cryptic-Prophage-Encoded Small Protein DicB Protects
1442 *Escherichia coli* from Phage Infection by Inhibiting Inner Membrane Receptor Proteins. *Journal*
1443 *of Bacteriology* [Internet]. 2019 [cited 2019 Nov 11];201. Available from:
1444 <https://jb.asm.org/content/201/23/e00475-19>
- 1445 99. Rasko DA, Rosovitz MJ, Myers GSA, Mongodin EF, Fricke WF, Gajer P, et al. The
1446 Pangenome Structure of *Escherichia coli*: Comparative Genomic Analysis of *E. coli* Commensal
1447 and Pathogenic Isolates. *Journal of Bacteriology*. 2008;190:6881–93.

1448 100. Jover LF, Effler TC, Buchan A, Wilhelm SW, Weitz JS. The elemental composition of
1449 virus particles: implications for marine biogeochemical cycles. *Nature Reviews Microbiology*.
1450 2014;12:519–28.

1451 101. Mizuno CM, Guyomar C, Roux S, Lavigne R, Rodriguez-Valera F, Sullivan MB, et al.
1452 Numerous cultivated and uncultivated viruses encode ribosomal proteins. *Nature*
1453 *Communications*. 2019;10:752.

1454

Figures

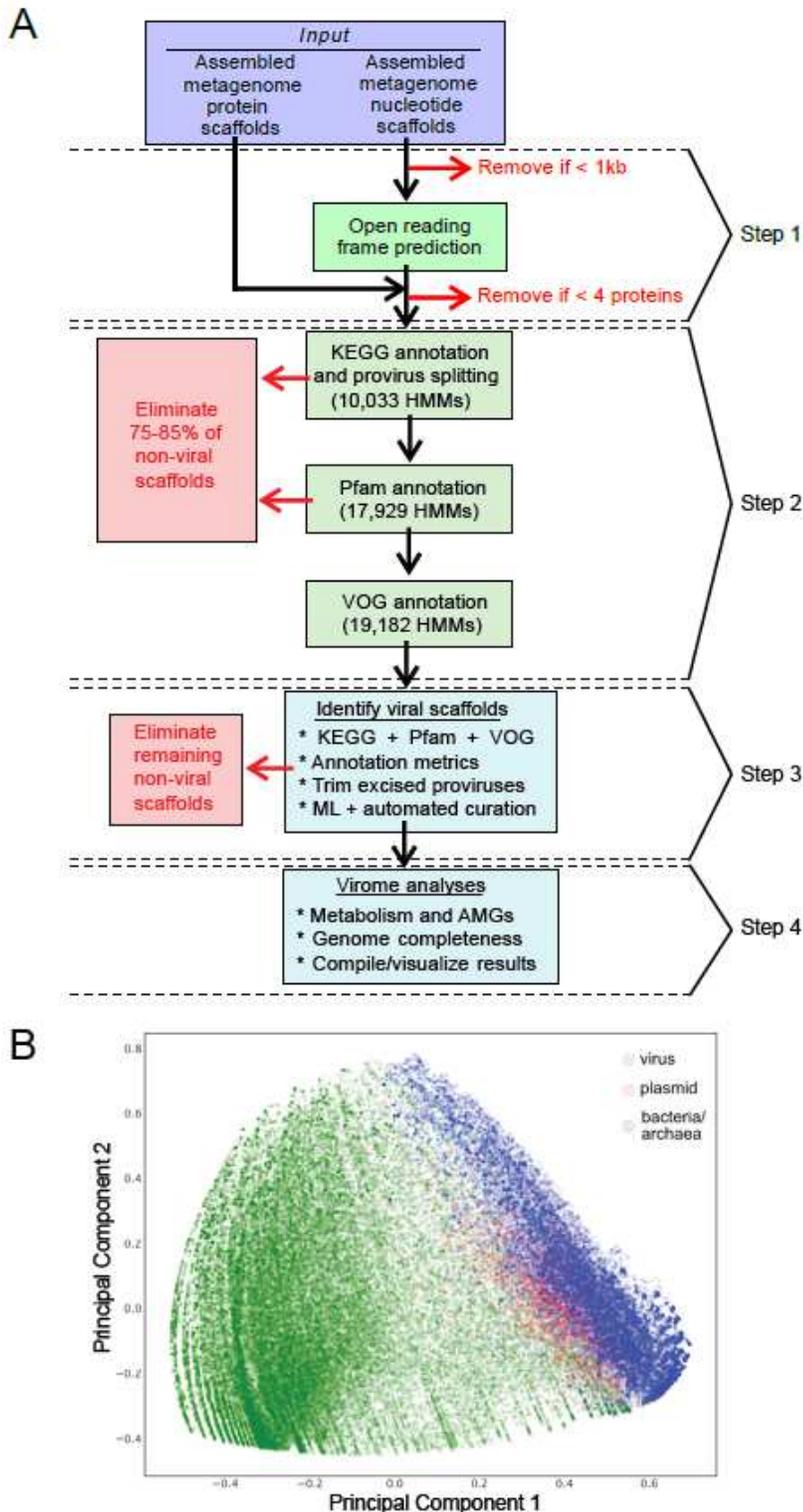


Figure 1

Representation of VIBRANT's method for virus identification and virome functional characterization. (A) Workflow of virome analysis. Annotations from KEGG, Pfam and VOG databases are used to construct signatures of viral and non-viral annotation signatures that are read into a neural network machine

learning model. (B) Visual representation (PCA plot) of the metrics used by the neural network to identify viruses, depicting viral, plasmid and bacterial/archaeal genomic sequences.

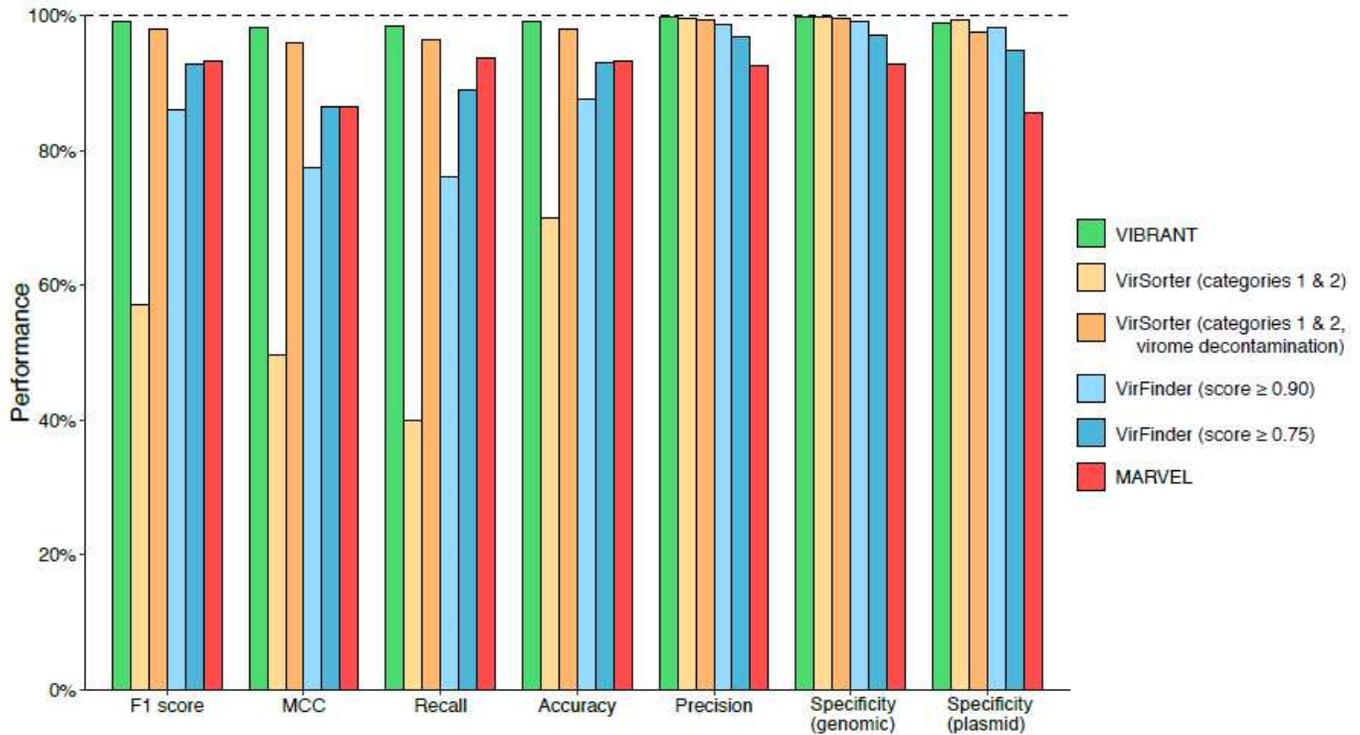


Figure 2

Performance comparison of VIBRANT, VirFinder, VirSorter and MARVEL on artificial scaffolds 3kb-15kb. Performance was evaluated using datasets of reference viruses, bacterial plasmids, and bacterial/archaeal genomes. For VirFinder and VirSorter two different confidence cutoffs were used (VirFinder: score of at least 0.90, and score of at least 0.75. VirSorter: categories 1 and 2 predictions, and categories 1 and 2 predictions using virome decontamination mode). All four programs were compared using the following statistical metrics: F1 score, MCC, recall, precision, accuracy and specificity. To ensure equal comparison all scaffolds tested encoded at least four open reading frames.

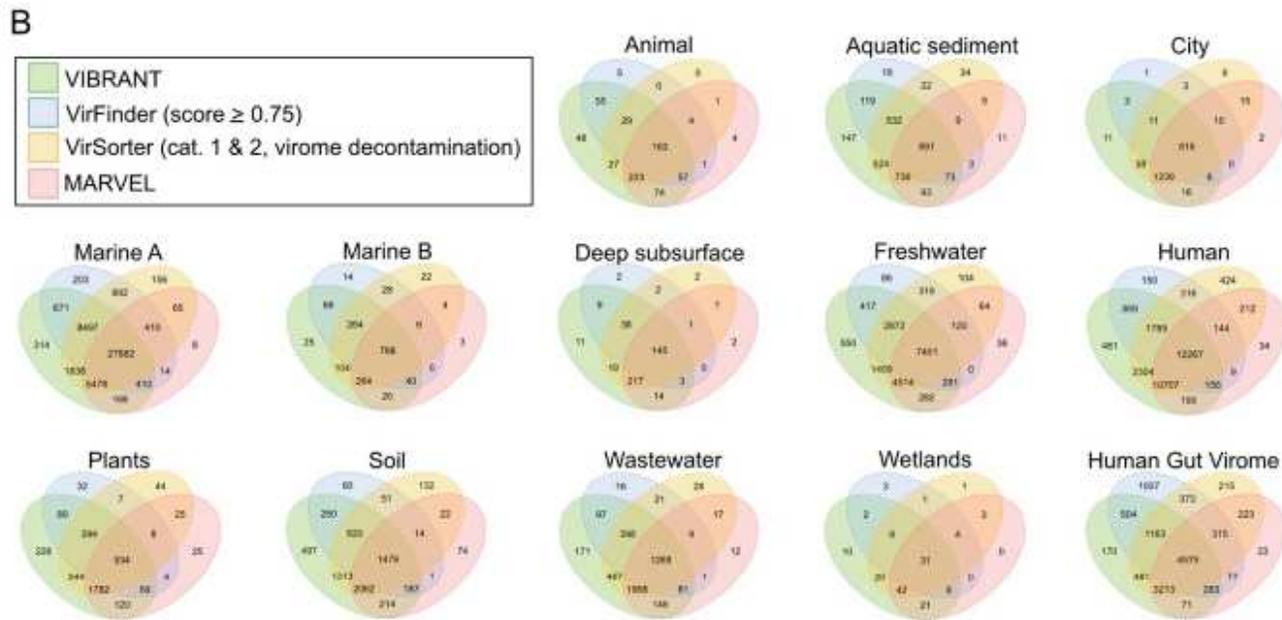
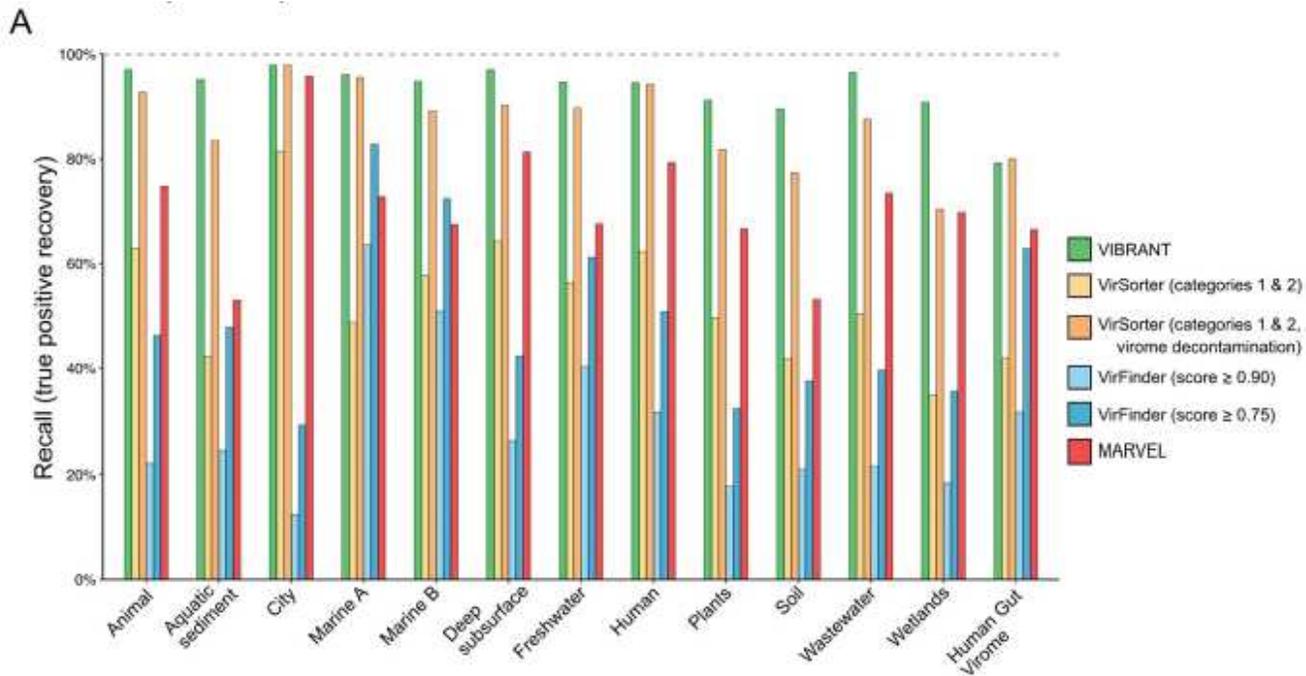
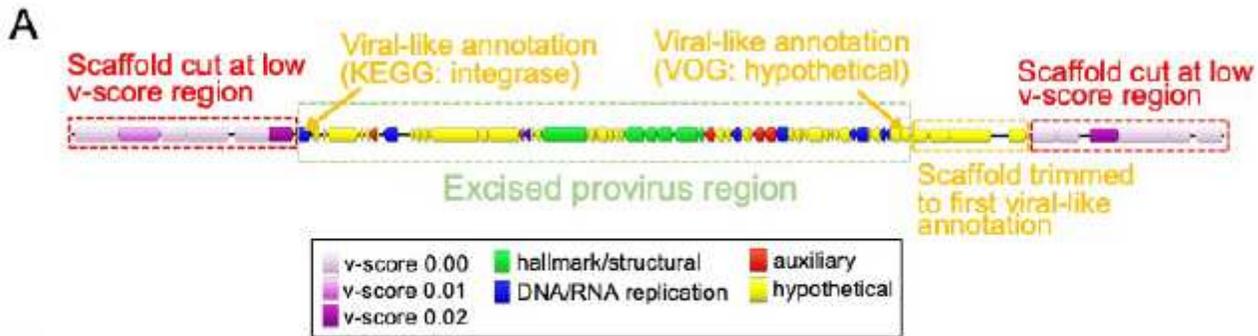


Figure 3

Effect of source environment on predictive abilities of VIBRANT, VirFinder, VirSorter and MARVEL. Viral scaffolds from IMG/VR and HGV database were used to test if VIBRANT displays biases associated with specific environments. (A) The recall (or recovery) of viral scaffolds from 12 environment groups was compared between VIBRANT and two confidence cutoffs for both VirFinder and VirSorter. Marine environments were classified into two groups: marine A (coastal, gulf, inlet, intertidal, neritic, oceanic, pelagic and strait) and marine B (hydrothermal vent, volcanic and oil). (B) Comparison of the overlap in the scaffolds identified as viruses by all three programs. Cutoffs for VirFinder (scores greater than or equal to 0.75) and VirSorter (categories 1 and 2 using virome decontamination mode) were set to display each program's ability to recover diverse viruses.



B

Genome	Region	VIBRANT	PHASTER	Prophage Hunter	VirSorter (cat. 1 & 2)
<i>Lactococcus lactis</i>	Putative provirus 1	35,516 - 55,047	26,459 - 56,369	34,907 - 49,863	10,280 - 95,312
	Putative provirus 2	442,048 - 483,244	447,143 - 484,064	447,143 - 484,094	450,421 - 517,314
	Putative provirus 3	492,452 - 531,873	502,338 - 520,485	-	-
	Putative provirus 4	1,033,815 - 1,070,779	1,033,815 - 1,079,173	1,018,719 - 1,075,596	1,029,339 - 1,080,357
	Putative provirus 5	1,415,048 - 1,460,426	1,414,112 - 1,457,044	1,415,289 - 1,457,116	1,414,112 - 1,460,426
	Putative provirus 6	-	1,997,699 - 2,028,023	2,011,247 - 2,025,756	2,011,426 - 2,025,635
<i>Desulfovibrio vulgaris</i>	Putative provirus 1	226,416 - 270,425	239,489 - 270,425	257,056 - 273,870	247,822 - 270,425
	Putative provirus 2	1,204,238 - 1,245,361	1,202,025 - 1,218,337	-	1,204,238 - 1,247,136
	Putative provirus 3	-	1,226,597 - 1,237,819	-	-
	Putative provirus 4	1,560,130 - 1,593,118	1,573,624 - 1,584,078	1,562,595 - 1,579,793	1,561,263 - 1,599,625
	Putative provirus 5	-	1,783,411 - 1,820,727	-	1,766,095 - 1,823,781
	Putative provirus 6	2,246,009 - 2,298,987	-	2,247,565 - 2,288,385	-
	Putative provirus 7	2,715,796 - 2,730,860	-	-	2,709,005 - 2,749,853
	Putative provirus 8	2,799,591 - 2,834,336	2,803,529 - 2,834,336	-	2,799,591 - 2,834,336
	Putative provirus 9	2,935,283 - 2,971,475	2,936,166 - 2,977,658	-	2,936,283 - 2,979,715
<i>Staphylococcus aureus</i>	Putative provirus 1	1,245,539 - 1,251,661	1,242,808 - 1,251,661	1,240,995 - 1,253,838	-
	Putative provirus 2	1,452,300 - 1,508,497	1,451,176 - 1,523,261	1,460,593 - 1,511,405	1,462,371 - 1,514,191
	Putative provirus 3	-	-	1,812,933 - 1,834,606	-
	Putative provirus 4	1,923,663 - 1,966,119	1,907,632 - 1,966,929	1,910,863 - 1,977,294	1,922,968 - 1,966,929
	Putative provirus 5	2,034,690 - 2,073,429	2,029,598 - 2,077,786	2,031,842 - 2,080,909	2,031,741 - 2,074,574
<i>Bacteroides vulgatus</i>	Putative provirus 1	3,580,013 - 3,632,769	-	-	3,580,013 - 3,637,236

C

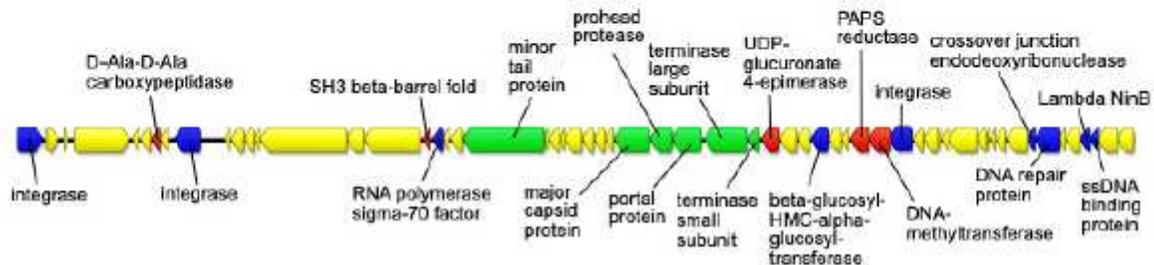


Figure 4

Prediction of integrated proviruses by VIBRANT, and comparison to PHASTER, Prophage Hunter and VirSorter. (A) Schematic representing the method used by VIBRANT to identify and extract provirus regions from host scaffolds using annotations. Briefly, v-scores are used to cut scaffolds at hostspecific sites and fragments are trimmed to the nearest viral annotation. (B) Comparison of proviral predictions within four complete bacterial genomes between VIBRANT, PHASTER, Prophage Hunter and VirSorter. For PHASTER, putative proviruses are colored according to “incomplete” (red), “questionable” (blue) and “intact” (green) predictions. Prophage Hunter is colored according to “active” (green) and “ambiguous” (blue) predictions. All VirSorter predictions for categories 1 and 2 are shown in green. (C) Manual

validation of the *Bacteroides vulgatus* provirus prediction made by VIBRANT. The presence of viral hallmark protein, integrase and genome replication proteins strongly suggests this is an accurate prediction.

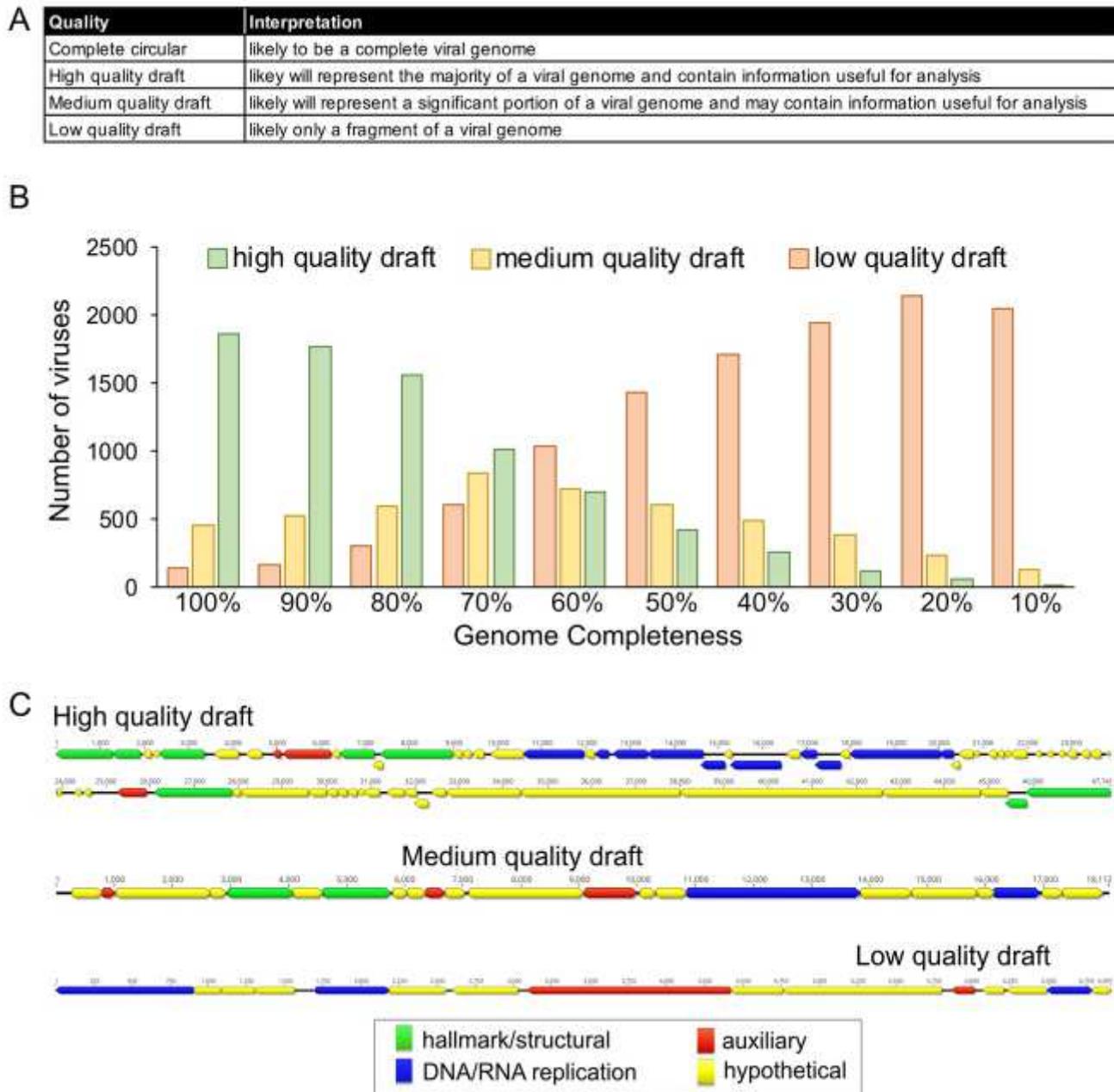


Figure 5

Estimation of genome quality of identified viral scaffolds. (A) Explanation of interpretation of quality categories: complete circular, high quality draft, medium quality draft and low quality draft. Quality generally represents total proteins, viral annotations, viral hallmark protein and nucleotide replication proteins, which are common metrics used for manual verification of viral genomes. (B) Application of quality metrics to 2466 NCBI RefSeq Caudovirales viruses with decreasing genome completeness from 100% to 10% completeness, respective of total sequence length. All 2466 viruses are represented within

each completeness group. (C) Examples of viral scaffolds representing low, medium and high quality draft categories.

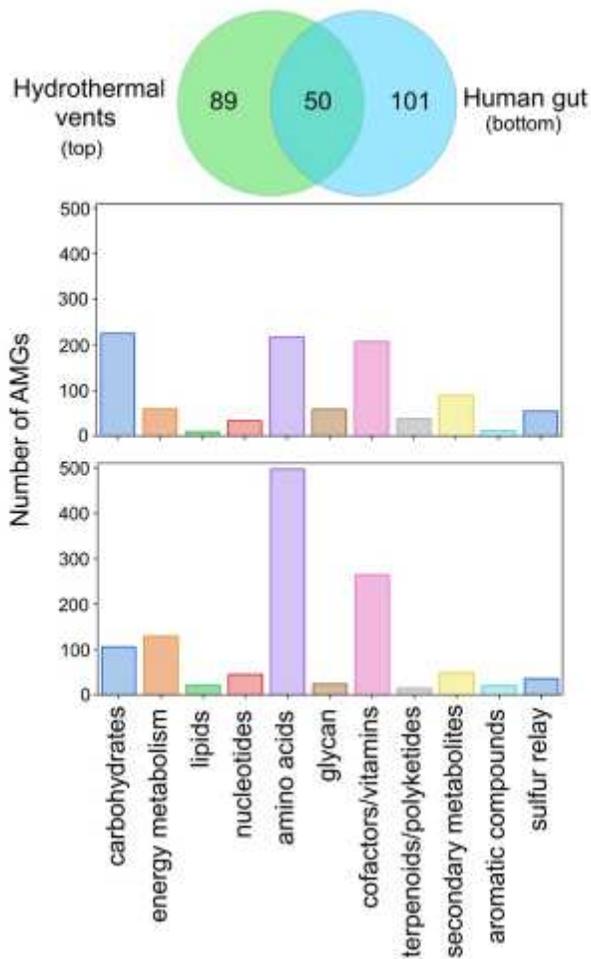


Figure 6

Comparison of AMG metabolic categories between hydrothermal vents and human gut. Venn diagram depicts the unique and shared non-redundant AMGs between 6 hydrothermal vent and 15 human gut metagenomes. The graphs depict the differential abundance of KEGG metabolic categories of respective AMGs for hydrothermal vents (top) and human gut (bottom).

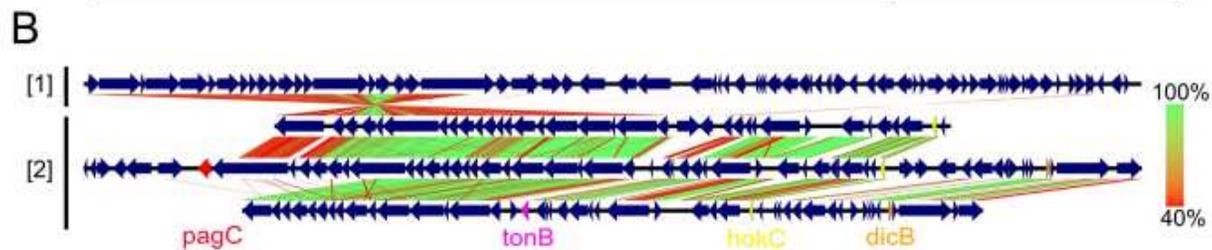
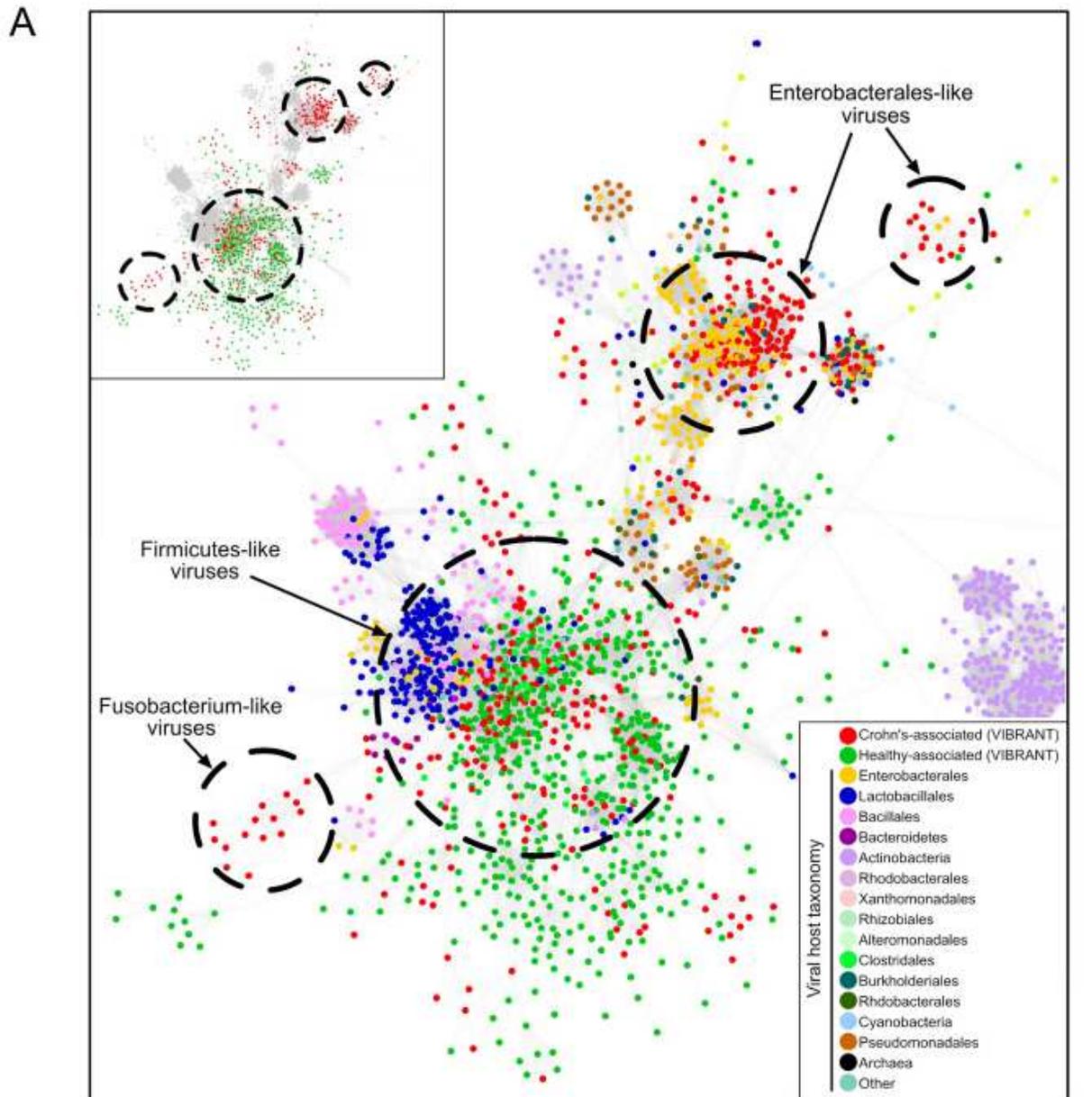


Figure 7

Viral metabolic comparison between Crohn's Disease and healthy individuals gut metagenomes. (A) Partial view of vConTACT2 protein network clustering of viruses identified by VIBRANT and reference viruses. Small clusters and clusters with no VIBRANT representatives are not shown. Each dot represents one genome and is colored according to host or dataset association. Relevant viral groups are indicated by dotted circles (circles enclose estimated boundaries). (B) tBLASTx similarity comparison between (1)

Escherichia phage Lambda and (2) three Crohn's-associated viruses identified by VIBRANT. Putative virulence genes are indicated: pagC, tonB, hokC and dicB.

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [AdditionalFile19.xlsx](#)
- [AdditionalFile17.xlsx](#)
- [AdditionalFile6.xlsx](#)
- [AdditionalFile10.xlsx](#)
- [AdditionalFile7.xlsx](#)
- [AdditionalFile18.xlsx](#)
- [AdditionalFile5.xlsx](#)
- [AdditionalFile15.xlsx](#)
- [AdditionalFile16.xlsx](#)
- [AdditionalFile22.pdf](#)
- [AdditionalFile21.pdf](#)
- [AdditionalFile24.pdf](#)
- [AdditionalFile1.xlsx](#)
- [AdditionalFile23.pdf](#)
- [AdditionalFile4.xlsx](#)
- [AdditionalFile2.xlsx](#)
- [AdditionalFile25.pdf](#)
- [AdditionalFile20.xlsx](#)
- [AdditionalFile3.xlsx](#)
- [AdditionalFile11.xlsx](#)
- [AdditionalFile14.xlsx](#)
- [AdditionalFile13.xlsx](#)
- [AdditionalFile12.xlsx](#)
- [AdditionalFile8.xlsx](#)
- [AdditionalFile9.xlsx](#)