

VIBRANT: Automated recovery, annotation and curation of microbial viruses, and evaluation of virome function from genomic sequences

Kristopher Kieft

University of Wisconsin Madison <https://orcid.org/0000-0002-9120-1961>

Zhichao Zhou

University of Wisconsin Madison

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

<https://orcid.org/0000-0002-9584-2491>

Methodology

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

Posted Date: December 13th, 2019

DOI: <https://doi.org/10.21203/rs.2.18878/v1>

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

[Read Full License](#)

1 **VIBRANT: Automated recovery, annotation and curation of microbial viruses, and**
2 **evaluation of virome 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 virome
59 function from metagenomic assemblies. VIBRANT uses neural networks of protein signatures and
60 a novel v-score metric that circumvents traditional boundaries to maximize identification of lytic
61 viral genomes and integrated proviruses, including highly diverse viruses. VIBRANT highlights
62 viral auxiliary metabolic genes and metabolic pathways, thereby serving as a user-friendly
63 platform for evaluating virome function. VIBRANT was trained and validated on reference virus
64 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 and VirFinder. When applied to 120,834
70 metagenomically derived viral sequences representing several human and natural environments,
71 VIBRANT recovered an average of 94.5% of the viruses, whereas VirFinder and VirSorter
72 achieved less powerful performance, averaging 48.1% and 56.0%, respectively. Similarly,
73 VIBRANT identified more total viral sequence and proteins when applied to real metagenomes.
74 When compared to PHASTER and Prophage Hunter for the ability to extract integrated provirus
75 regions from host scaffolds, VIBRANT performed comparably and even identified proviruses that
76 the other programs did not. To demonstrate applications of VIBRANT, we studied viromes
77 associated with Crohn's Disease to show that specific viral groups, namely Enterobacteriales-like
78 viruses, as well as putative dysbiosis associated viral proteins are more abundant compared to
79 healthy individuals, providing a possible viral link to maintenance of diseased states.

80

81 **Conclusions**

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

85

86 **Keywords**

87 Virome, virus, bacteriophage, metagenome, machine learning, auxiliary metabolism

88

89

90

91 **Background**

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

101 It is estimated that viral diversity exceeds that of living organisms, and therefore harbors
102 enormous potential for diverse genomic content, arrangement and encoded functions [14].
103 Accordingly, there is substantial interest in “mining” viral sequences for novel anti-microbial drug
104 candidates, enzymes for biotechnological applications, and for bioremediation efforts [15–19].
105 Moreover, viruses have a unique capability to rapidly evolve genes via high mutation rates and act
106 as intermediate carriers to transfer these genes to their hosts and subsequently to the surrounding
107 communities [20–22].

108 Our understanding of the diversity of viruses continues to expand with the discovery of
109 novel viral lineages within the last decade. One of the most striking examples is the
110 characterization of crAssphage, an extraordinarily abundant virus infecting *Bacteroides* species
111 within the human gut that went unnoticed for years due to its lack of homology with known viral
112 sequences [23]. Moreover, the discovery of megaphages, the largest known bacterial viruses
113 infecting the human gut bacteria *Prevotella*, has pushed the boundaries on the coding capacity of
114 viruses [24,25]. In the oceans, a newly discovered lineage of *Vibrio*-infecting non-tailed viruses,
115 generally considered unconventional since most known bacterial viruses are tailed, fueled the
116 notion that current viral recovery methods are skewing our understanding of viruses in the
117 environment [26]. Taken together, this highlights that estimates of viral diversity are biased
118 towards tailed dsDNA viruses and are likely underrepresenting other families of viruses including
119 those with ssDNA and RNA genomes [27,28].

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

131 Due to the difficulty of collecting virus-only samples as well as the need to integrate viruses
132 into models of ecosystem function, it has become of great interest to determine which sequences
133 within microbial communities are derived from viruses. Within the cellular fraction of a sample
134 there can remain a large number of viruses for a variety of reasons. First, these viruses can exist as
135 active intracellular infections, which may be the case for as many as 30% of all bacteria at any
136 given time [41]. Second, there may be particle-attached viruses resulting from viruses’ inherently

137 “sticky” nature [42]. Lastly, many viruses exist as “proviruses”, or viral genomes either integrated
138 into that of their host or existing within the host as an episomal sequence. As such, it is crucial for
139 the accurate evaluation of microbial community characteristics, structure and functions to be able
140 to separate these viral sequences.

141 Multiple tools exist for the identification of viral sequences from mixed metagenomic
142 assemblies. For several years VirSorter [43], which succeeded tools such as VIROME [44] and
143 Metavir [45], has been the most widely used for its ability to accurately identify viral metagenomic
144 fragments (scaffolds) from large metagenomic assemblies. VirSorter predominantly relies on
145 database searches of predicted proteins, using both reference homology as well as probabilistic
146 similarity, to compile metrics of enrichment of virus-like proteins and simultaneous depletion of
147 other proteins. To do this it uses a virus-specific curated database as well as Pfam [46] for non-
148 virus annotations, though it does not fully differentiate viral from non-viral Pfam annotations. It
149 also incorporates signatures of viral genomes, such as encoding short genes or having low levels
150 of strand switching between genes. VirSorter is also unique in its ability to use these annotation
151 and sequence metrics to identify and extract integrated provirus regions from host scaffolds. After
152 prediction of viral sequences, VirSorter labels viral scaffolds with one of three confidence levels:
153 *categories* 1, 2 or 3. Categories 1 and 2 are generally considered accurate, but category 3
154 predictions are more likely to contain false identifications. While VirSorter is quite accurate, it
155 likely underrepresents the diversity and abundance of viruses within metagenomic assemblies.

156 More recent tools have been developed to compete with the performance of VirSorter in
157 order to expand our appreciation and understanding of viruses. VirFinder [47] was the first tool to
158 implement machine learning and be completely independent of reference databases for predicting
159 viral sequences which was a platform later implemented in PPR-Meta [48]. VirFinder was built
160 with the consideration that viruses tend to display distinctive patterns of 8-nucleotide frequencies
161 (otherwise known as 8-mers), which was proposed despite the knowledge that viruses can share
162 remarkably similar nucleotide patterns with their host [49]. These 8-mer patterns were used to
163 build a random forest machine learning model to quickly classify sequences as short as 500 bp
164 without the need for gene prediction. VirFinder generates model-derived scores as well as
165 probabilities of prediction accuracy, though it is up to the user to define the cutoffs which can
166 ultimately lead to uncertainties in rates of false identification of viral sequences. VirFinder was
167 shown to greatly improve the ability to recover viral sequences compared to VirSorter, but it also
168 demonstrates substantial host and source environment biases in predicting diverse viruses. For
169 example, VirFinder was able to recover viruses infecting Proteobacteria more readily than those
170 infecting Firmicutes due to reference database-associated biases while training the machine
171 learning model. Additional biases were also identified between different source environments,
172 seen through the under-recovery of viruses from certain environments compared to others [50].

173 Additional recent tools have been developed that utilize slightly different methods for
174 identifying viral scaffolds. MARVEL [51], for example, leverages annotation, sequence signatures
175 (e.g., strand switching and gene density) and machine learning to identify viruses from
176 metagenomic bins. MARVEL differs from VirSorter in that it only utilizes a single virus-specific
177 database for annotation and also differs from VirFinder in that it does not use global nucleotide
178 frequency patterns. However, MARVEL provides no consideration for integrated proviruses and
179 is only suitable for identifying bacterial viruses from the order *Caudovirales* which substantially
180 limits its ability to discover novel viruses. Another recently developed tool, VirMiner [52], is
181 unique in that it functions to use metagenomic reads and associated assembly data to identify
182 viruses and performs best for high abundance (i.e., high coverage when assembled) viruses.

183 VirMiner is a web-based server that utilizes a hybrid approach of employing both homology-based
184 searches to a virus-specific database as well as machine learning. VirMiner was found to have
185 improved ability to recover viral scaffolds compared to both VirSorter and VirFinder but was
186 concurrently much less accurate. Poor accuracy would lead to a skewed interpretation of virome
187 function if the identified virome consisted of many non-viral sequences. This distinction is
188 important because VirMiner employs functional characterization as well as determination of virus-
189 host relationships.

190 Thus far, VirSorter remains the most efficient tool for identifying integrated proviruses
191 within metagenomic assemblies. Other tools, predominantly PHASTER [53] and Prophage Hunter
192 [54], are specialized in identifying integrated proviruses from whole genomes rather than scaffolds
193 generated by metagenomic assemblies. Similar to VirSorter, these two provirus predictors rely on
194 reference homology and viral sequence signatures with sliding windows to identify regions of a
195 host genome that belong to a virus. Although they are useful for whole genomes, they lack the
196 capability of identifying scaffolds belonging to lytic (i.e., non-integrated) viruses and perform
197 slower for large datasets. In addition, both PHASTER and Prophage Hunter are exclusively
198 available as web-based servers and offer no stand-alone command line tools.

199 Here we developed VIBRANT (Virus Identification By iterative ANnotation), a tool for
200 automated recovery, annotation, and curation of both free and integrated viruses from
201 metagenomic assemblies and genome sequences. VIBRANT is capable of identifying diverse
202 dsDNA, ssDNA and RNA viruses infecting both bacteria and archaea, and to our knowledge has
203 no evident environmental biases. VIBRANT uses neural networks of protein annotation signatures
204 from non-reference-based similarity searches with Hidden Markov Models (HMMs) as well as a
205 unique ‘v-score’ metric to maximize identification of diverse and novel viruses. After identifying
206 viral scaffolds VIBRANT implements curation steps to validate predictions. VIBRANT
207 additionally characterizes virome function by highlighting AMGs and assesses the metabolic
208 pathways present in viral communities. All viral genomes, proteins, annotations and metabolic
209 profiles are compiled into formats for user-friendly downstream analyses and visualization. When
210 applied to reference viruses, non-reference virus datasets and various assembled metagenomes,
211 VIBRANT outperformed both VirFinder and VirSorter in the ability to maximize virus recovery
212 and minimize false discovery. When compared to PHASTER and Prophage Hunter for the ability
213 to extract integrated provirus regions from host scaffolds, VIBRANT performed comparably and
214 even identified proviruses that the other programs did not. VIBRANT was also used to identify
215 differences in metabolic capabilities between viruses originating from various environments.
216 When applied to three separate cohorts of individuals with Crohn’s Disease, VIBRANT was able
217 to identify both differentially abundant viral groups compared to healthy controls as well as virally
218 encoded genes putatively influencing a diseased state. VIBRANT is freely available for download
219 at <https://github.com/AnantharamanLab/VIBRANT>. VIBRANT is also available as a user-
220 friendly, web-based application through the CyVerse Discovery Environment at
221 [https://de.cyverse.org/de/?type=apps&app-id=c2864d3c-fd03-11e9-9cf4-008cfa5ae621&system-
222 id=de](https://de.cyverse.org/de/?type=apps&app-id=c2864d3c-fd03-11e9-9cf4-008cfa5ae621&system-id=de) [55].
223

224 Results

225 VIBRANT was built to extract and analyze bacterial and archaeal viruses from assembled
226 metagenomic and genome sequences, as well as provide a platform for characterizing metabolic
227 proteins and functions in a comprehensive manner. The concept behind VIBRANT’s mechanism
228 of virus identification stems from the understanding that arduous manual inspection of annotated

229 genomic sequences produces the most dependable
 230 results. As such, the primary metrics used to inform
 231 validated curation standards and to train VIBRANT’s
 232 machine learning based neural network to identify
 233 viral sequences reflects human-guided intuition,
 234 though in a high-throughput automated fashion.

235
 236 **Determination of v-score**

237 We developed a unique ‘v-score’ metric as an
 238 approach for providing quantitative information to
 239 VIBRANT’s algorithm in order to assess the
 240 qualitative nature of annotation information. A v-
 241 score is a value assigned to each possible protein
 242 annotation that scores its association to viral
 243 genomes. V-score differs from the previously used
 244 “virus quotient” metric [56,57] in that it does not take
 245 into account the annotation’s relatedness to bacteria
 246 or archaea. Not including significant similarity to
 247 non-viral genomes in the calculation of v-scores has
 248 important implications for this metric’s utility.
 249 Foremost is that annotations shared between viruses
 250 and their hosts, such as ribonucleotide reductases,
 251 will be assigned a v-score reflecting its association to
 252 viruses, not necessarily virus-specificity. Many genes
 253 are commonly associated with viruses and host
 254 organisms, but when encoded on viral genomes can
 255 be central to virus replication efficiency (e.g.,
 256 ribonucleotide reductases [58]). Therefore, a metric
 257 representing virus-association rather than virus-
 258 specificity would be more appropriate in identifying
 259 if an unknown scaffold is viral or not. Secondly, this
 260 approach takes into account widespread horizontal
 261 gene transfer of host genes by viruses as well as the
 262 presence of AMGs.

263
 264 **VIBRANT workflow**

265 VIBRANT utilizes several annotation metrics
 266 in order to guide removal of non-viral sequences
 267 before curation of reliable viral scaffolds. The
 268 annotation metrics used are derived from HMM-
 269 based probabilistic searches of protein families from
 270 the Kyoto Encyclopedia of Genes and Genomes
 271 (KEGG) KoFam [59,60], Pfam [46] and Virus
 272 Orthologous Group (VOG) (vogdb.org) databases.
 273 VIBRANT is not reliant on reference-based similarity
 274 and therefore accounts for the large diversity of

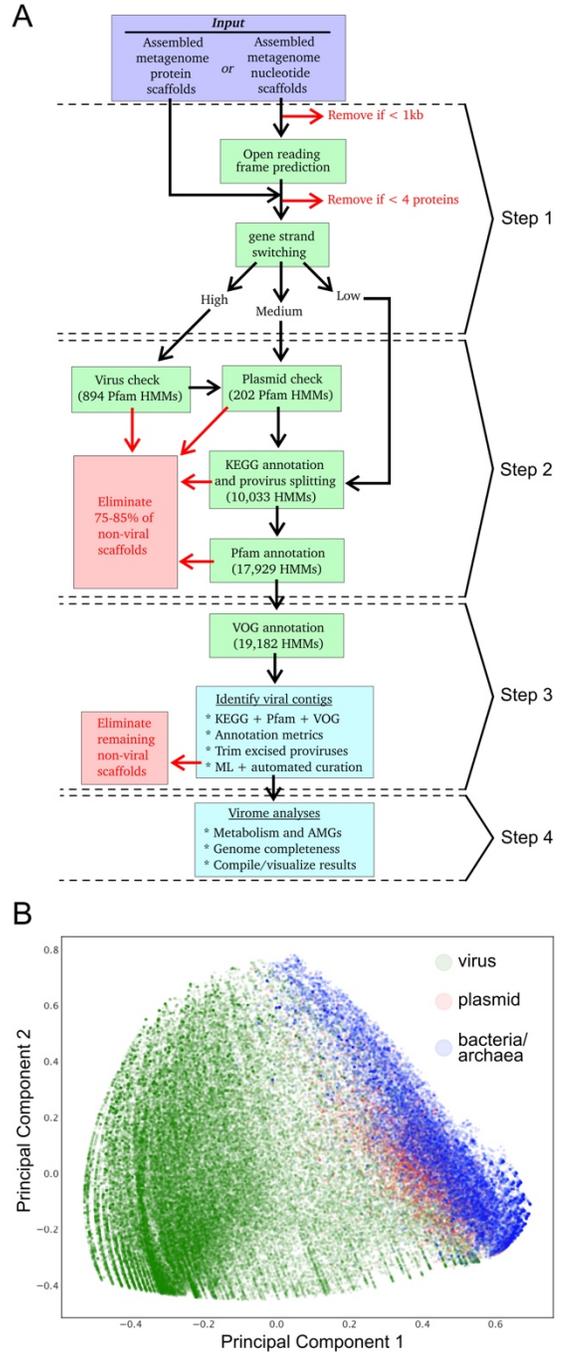


Figure 1. Representation of VIBRANT’s method for virus identification and virome functional characterization. (A) Workflow of virome analysis. Protein HMMs for analysis from KEGG, Pfam and VOG databases were used to construct signatures of viral and non-viral annotation metrics. (B) Visual representation (PCA plot) of the metrics used by the neural network to identify viruses, depicting viral, plasmid and bacterial/archaeal genomic sequences.

275 viruses on Earth and their respective proteins. Consequently, widespread horizontal gene transfer,
276 rapid mutation and the vast amount of novel sequences do not hinder VIBRANT's ability to
277 identify known and novel viruses. VIBRANT relies minimally on non-annotation features, such
278 as rates of open reading frame strand switching, because these features were not as well conserved
279 in genomic scaffolds in contrast to whole genomes.

280 VIBRANT's workflow consists of four main steps (Figure 1A). Briefly, proteins (predicted
281 or user input) are used by VIBRANT to first eliminate non-viral sequences by assessing non-viral
282 annotation signatures derived from KEGG and Pfam HMM annotations. At this step potential host
283 scaffolds are fragmented using sliding windows of KEGG v-scores in order to extract integrated
284 provirus sequences. Following the elimination of most non-viral scaffolds and rough excision of
285 provirus regions, proteins are annotated by VOG HMMs. Before analysis by the neural network
286 machine learning model, any extracted putative provirus is trimmed to exclude any remaining non-
287 viral sequences. Annotations from KEGG, Pfam and VOG are used to compile 27 metrics that are
288 utilized by the neural network to predict viral sequences. These 27 metrics were found to be
289 adequate for the separation of viral and non-viral scaffolds (Figure 1B). After prediction by the
290 neural network a set of curation steps are implemented to filter the results in order to improve
291 accuracy as well as recovery of viruses. Once viruses are identified VIBRANT automates the
292 analysis of virome function by highlighting AMGs and assigning them to KEGG metabolic
293 pathways. The genome quality (i.e., proxy of completeness) of identified viral scaffolds is
294 estimated using a subset of the annotation metrics and viral sequences are used to identify circular
295 templates (i.e., likely complete circular viruses). These quality analyses were determined to best
296 reflect established completeness metrics for both bacteria and viruses [61,62]. Finally, VIBRANT
297 compiles all results into a user-friendly format for visualization and downstream analysis. For a
298 detailed description of VIBRANT's workflow see Methods.

299

300 **Comparison of VIBRANT to other programs**

301 VirSorter and VirFinder, two commonly used programs for identifying bacterial and
302 archaeal viruses from metagenomes, were selected to compare against VIBRANT for the ability
303 to accurately identify viruses. We evaluated all three programs' performance on the same viral,
304 bacterial and archaeal genomic, and plasmid datasets. Given that both VirSorter and VirFinder
305 produce various confidence ranges of virus identification, we selected certain parameters for each
306 program for comparison. For VirSorter, the parameters selected were [1] category 1 and 2
307 predictions, and [2] categories 1, 2 and 3 (i.e., all) predictions. For VirFinder, the intervals were
308 [1] scores greater than or equal to 0.90 (approximately equivalent to a p-value of 0.013), and [2]
309 scores greater than or equal to 0.75 (approximately equivalent to a p-value of 0.037). Hereafter,
310 we provide two statistics for each VirSorter and VirFinder run that reflect results according to the
311 two set confidence intervals, respectively.

312 VIBRANT yields a single output of confident predictions and therefore does not provide
313 multiple output options. Since VIBRANT is only partially reliant on its neural network machine
314 learning model for making predictions, all comparisons are focused on VIBRANT's full workflow
315 performance. VIBRANT does not consider scaffolds shorter than 1000 bp or those that encode
316 less than four predicted open reading frames in order to maintain a low false positive rate (FPR)
317 and have sufficient annotation information for identifying viruses. Therefore, in comparison of
318 performance metrics only scaffolds meeting VIBRANT's minimum requirements were analyzed.
319 Inclusion of fragments encoding less than four open reading frames in analyses, which are
320 frequently generated by metagenomic assemblies, are discussed below. We used the following

321 calculations to compare performance: recall, precision, accuracy, specificity and F1 score (Figure
 322 2).

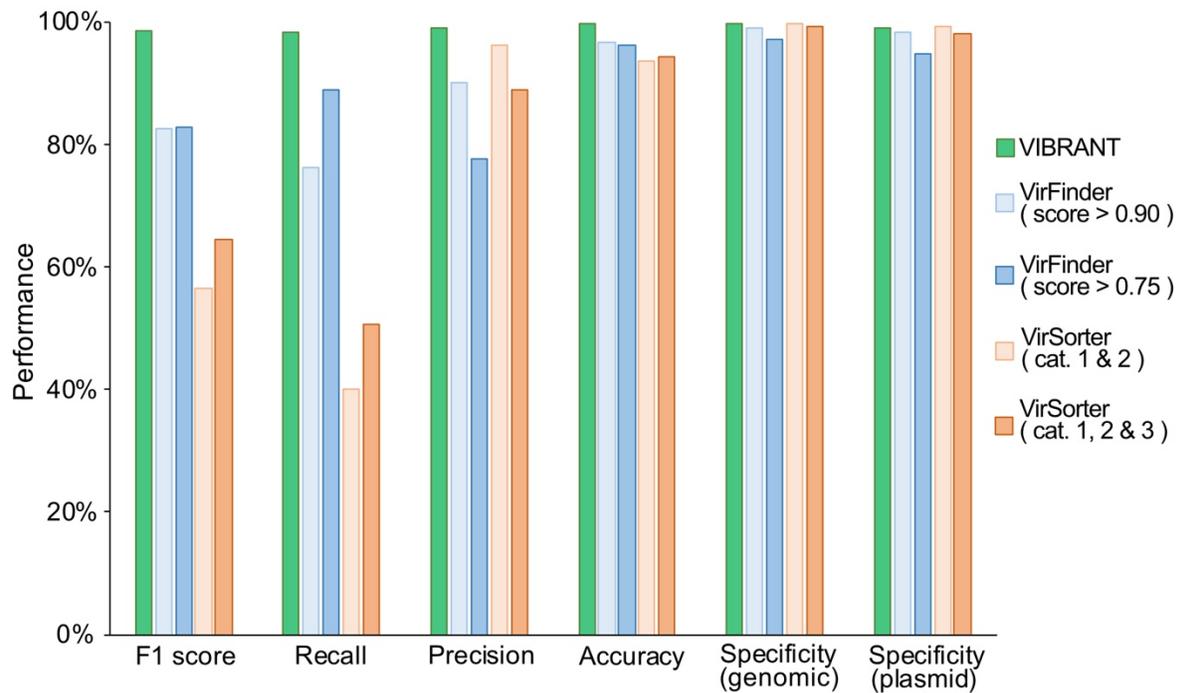


Figure 2. Performance comparison of VIBRANT, VirFinder and VirSorter on artificial scaffolds 3kb-15kb. Performance was evaluated using datasets of reference viruses, bacterial plasmids, and bacterial/archaeal genomes. For VirFinder and VirSorter two different confidence cutoffs were used (VirFinder: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

323 First, we evaluated the true positive rate (TPR, or recall) of viral genomic fragments as
 324 well as whole viral genomes. Viral genomes were acquired from the National Center for
 325 Biotechnology Information (NCBI) RefSeq and GenBank databases and split into various non-
 326 redundant fragments between 3 and 15 kb to simulate genomic scaffolds (Additional File 1: Table
 327 S1). VIBRANT correctly identified 98.38% of the 29,926 viral fragments, which was substantially
 328 greater than either VirSorter (40.00% and 50.67%) and VirFinder (76.23% and 89.02%).

329 Similar to TPR, we calculated FPR (or specificity) using two different datasets: genomic
 330 fragments of bacteria and archaea (hereafter genomic), and bacterial plasmids (plasmid). Plasmids
 331 were evaluated separately because they often encode for genes similar to those on viral genomes,
 332 such as those for genome replication and mobilization. Genomic and plasmid sequences were
 333 acquired from NCBI RefSeq and GenBank databases and split into various non-redundant
 334 fragments between 3 and 15 kb (Additional File 1: Table S1). Before analysis, putative proviruses
 335 were depleted from the datasets (see Methods). With the exception of VirSorter set at categories 1
 336 and 2 for the plasmid dataset, VIBRANT had the highest specificity for both genomic (99.92%)
 337 and plasmid fragments (99.04%). VirSorter had similar specificity for both genomic (99.84% and
 338 99.33%) and plasmid (99.33% and 98.10%) datasets, but only VirFinder set to a score cutoff of
 339 0.90 was fully comparable (genomic: 99.10%, plasmid: 98.40%). At a score cutoff of 0.75,
 340 VirFinder was slightly less specific (genomic: 97.19%, plasmid: 94.92%). Although VirFinder (set

341 to a score cutoff of 0.90) and VIBRANT had a similar overall specificity, VirFinder identified 11.8
342 times more bacterial/archaeal scaffolds as viruses (false discoveries) compared to VIBRANT
343 (2,311 and 196, respectively).

344 We used the results from TPR of viral fragments and FPR of non-viral genomic or plasmid
345 fragments to calculate precision, accuracy and F1 score. VIBRANT outperformed VirFinder and
346 VirSorter at either score criteria in both precision (99.01%) and accuracy (99.74%). F1 is a metric
347 (maximum value of 1) accounting for both TPR and FPR, and therefore acts as a comprehensive
348 evaluation of overall performance. Our calculation of F1 indicates that VIBRANT (0.99) is able
349 to better identify viruses while subsequently reducing false identifications compared to VirSorter
350 (0.57 and 0.65) or VirFinder (0.83 and 0.83).

351

352 **Identification of viruses in diverse environments**

353 We next tested VIBRANT's ability to successfully identify viruses from a diversity of
354 environments. Using 120,834 viruses from the Integrated Microbial Genomes and Viruses
355 (IMG/VR v2.0) database [63,64], in which the source environment of viruses is categorized, we
356 identified that VIBRANT is more robust in identifying viruses from all tested environments
357 compared to VirFinder and VirSorter (Figure 3A). Excluding air, in which there were only 62
358 representative viruses, VIBRANT averaged 94.5% recall, substantially greater than VirFinder
359 (29.2% and 48.1%) and VirSorter (54.4% and 56.0%). These results suggest that in comparison to
360 other software, VIBRANT has no evident database or environmental biases and is fully capable of
361 identifying viruses from a broad range of source environments. We also used a dataset of 13,203
362 viruses from the Human Gut Virome database [65] for additional comparison. The vast majority
363 of viruses (~96%) in this dataset were assumed to infect bacteria. Although recall was diminished
364 compared to IMG/VR datasets, VIBRANT (78.7%) nevertheless outperformed both VirFinder
365 (31.7% and 62.8%) and VirSorter (41.9% and 46.5%) on this dataset.

366 Many viruses from the IMG/VR dataset that were identified by VIBRANT were not
367 identified by either VirFinder or VirSorter, indicating that VIBRANT has the propensity for
368 discovery of novel viruses (Figure 3B). For most environments, the majority of viruses identified
369 by VirFinder were already identified by either VIBRANT or VirSorter. The differences in the
370 overlap of identified viruses was not too distinctive in environments for which many reference
371 viruses are available, such as marine. For more understudied environments, such as plants or
372 wastewater, VIBRANT displayed near-complete overlap with VirFinder and VirSorter predictions
373 in conjunction with identifying over 40% more viruses.

374

375 **Identification of viruses in mixed metagenomes**

376 Metagenomes assembled using short read technology contain many scaffolds that do not
377 meet VIBRANT's minimum length requirements and therefore are not considered during analysis.
378 Despite this, VIBRANT's predictions contain more annotation information and greater total viral
379 sequence length than tools built to identify short sequences, such as scaffolds with less than four
380 open reading frames. VIBRANT, VirFinder (score cutoff of 0.90) and VirSorter (categories 1 and
381 2) were used to identify viruses from human gut, freshwater lake and thermophilic compost
382 metagenome sequences (Table 1). In addition, alternate program settings—VIBRANT “virome”
383 mode, VirFinder score cutoff of 0.75 and VirSorter “virome decontamination” mode—were used
384 to identify viruses from an estuary virome dataset. Each metagenomic assembly was limited to

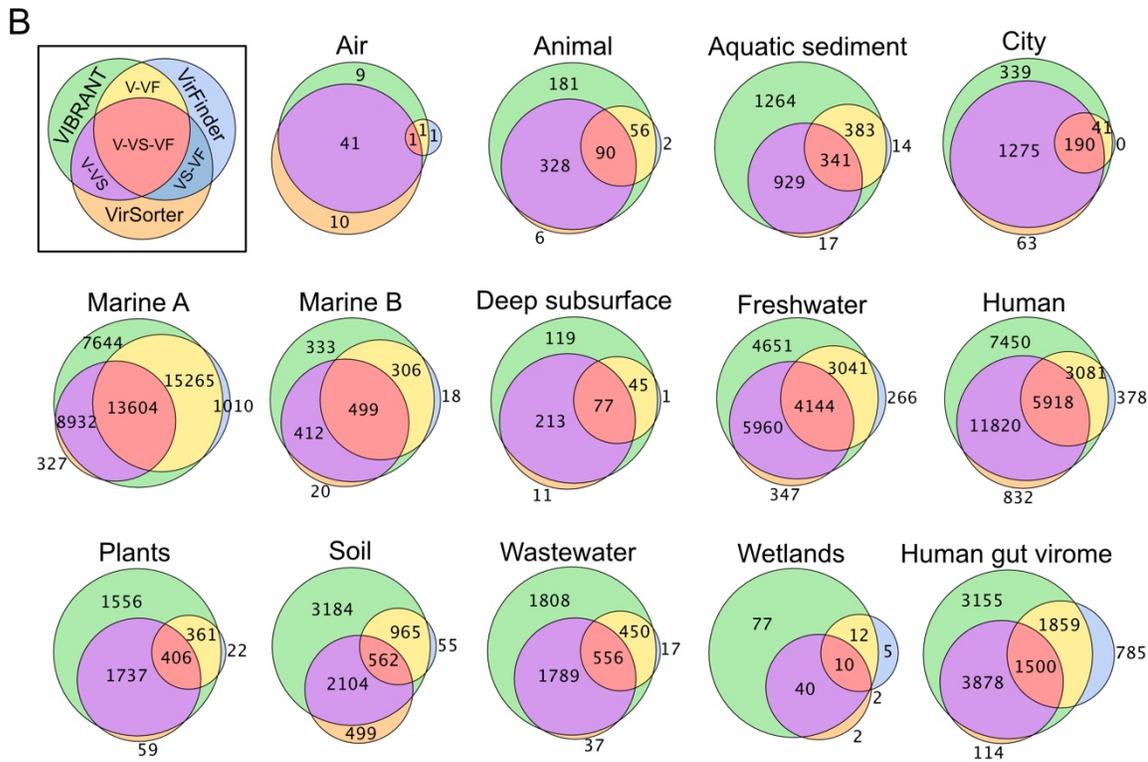
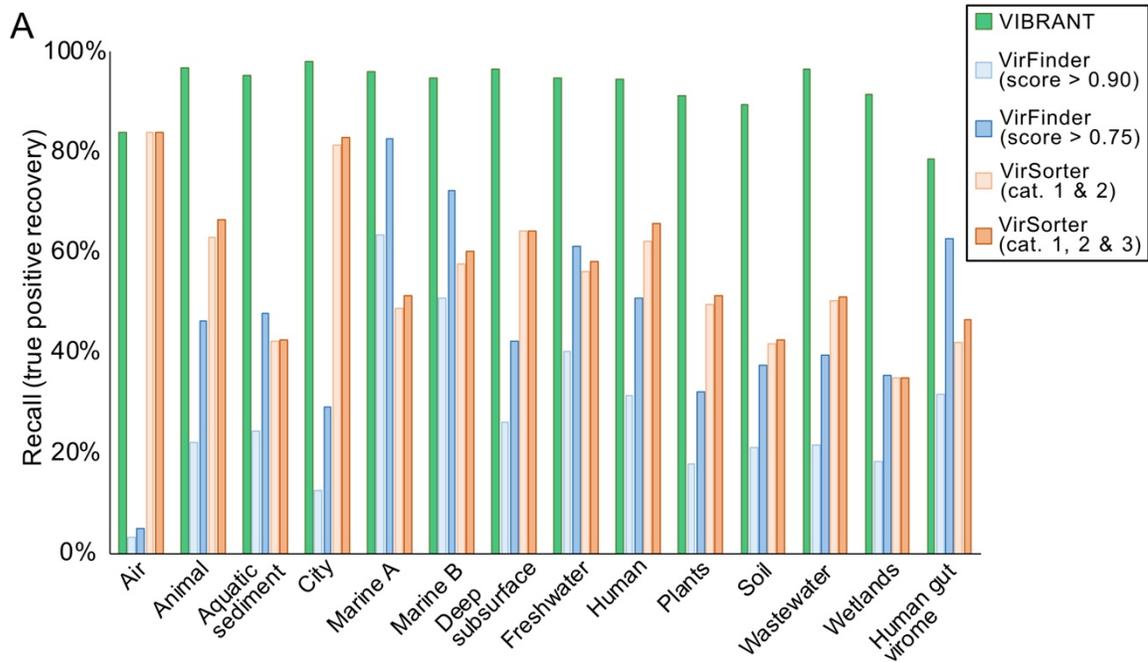


Figure 3. Effect of source environment on predictive abilities of VIBRANT, VirFinder and VirSorter. Viral scaffolds from IMG/VR and HGV database were used to test if VIBRANT displays biases associated with specific environments. (A) The recall (or recovery) of viral scaffolds from 14 environment groups 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. Unique scaffolds identified by each program are in green (VIBRANT), orange (VirSorter) and light blue (VirFinder). The size of the circles represents the relative size of the group.

Table 1. Virus recovery of VIBRANT, VirFinder and VirSorter from mixed metagenomes and a virome. Mixed community assembled metagenomes from the human gut, thermophilic compost and a freshwater lake, 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.

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	505	604	0.84	284	1.78
			total virus length (bp)	5,159,390	1,696,118	3.04	3,982,292	1.30
			total virus proteins	7,534	2,134	3.53	5,484	1.37
human gut: carcinoma	53,946	18,669	total putative viruses	744	1,329	0.56	450	1.65
			total virus length (bp)	5,415,994	3,500,838	1.55	4,182,862	1.29
			total virus proteins	8,108	4,644	1.75	5,945	1.36
human gut: healthy	42,739	17,079	total putative viruses	548	672	0.82	309	1.77
			total virus length (bp)	5,468,452	2,411,049	2.27	4,512,571	1.21
			total virus proteins	7,998	3,230	2.48	6,127	1.31
thermophilic compost	68,815	21,620	total putative viruses	1,057	878	1.20	383	2.76
			total virus length (bp)	6,577,000	2,238,129	2.94	3,290,654	2.00
			total virus proteins	9,908	2,806	3.53	4,400	2.25
freshwater lake (bog)	79,862	26,832	total putative viruses	5,600	7,567	0.74	1,503	3.73
			total virus length (bp)	34,861,470	25,357,664	1.37	15,436,797	2.26
			total virus proteins	55,976	37,537	1.49	21,280	2.63
* estuary virome	5,247	3,277	total putative viruses	3,135	2,294	1.37	1,121	2.80
			total virus length (bp)	10,241,625	6,478,804	1.58	5,163,674	1.98
			total virus proteins	20,475	12,035	1.70	9,645	2.12

* VIBRANT, VirFinder and VirSorter ran with alternate settings

386 sequences of at least 1000bp but no minimum open reading frame limit was set. For these
387 metagenomes, 31% to 40% of the scaffolds were of sufficient length (at least four open reading
388 frames) to be analyzed by VIBRANT; for the estuary virome 62% were of sufficient length. In
389 comparison, 100% of scaffolds from each dataset were long enough to be analyzed by VirFinder.
390 The ability of VirFinder to make a prediction with each scaffold is considered the major strength
391 of the tool.

392 For all six assemblies VirFinder averaged approximately 1.2 times more virus
393 identifications than VIBRANT, though for both thermophilic compost and the estuary virome
394 VIBRANT identified a greater number. Despite VirFinder averaging more total virus
395 identifications, VIBRANT averaged just over 2.1 times more total viral sequence length and 2.4
396 times more total viral proteins. This is the result of VIBRANT having the capability to identify
397 more viruses of higher quality and longer
398 sequence length. For example, among all six datasets VIBRANT identified 1,309 total viruses at
399 least 10 kb in length in comparison to VirFinder's 479. VIBRANT was also able to outperform
400 VirSorter in all metrics, averaging 2.4 times more virus identifications, nearly 1.7 times more total
401 viral sequence length, and 1.8 times more encoded viral proteins.

402 VIBRANT's method of predicting viral scaffolds provides a unique opportunity in
403 comparison to similar tools in that it yields scaffolds of higher quality which are more amenable
404 for analyzing protein function in viromes. It is an important distinction that the total number of
405 viruses identified may not be correlated with the total viral sequence identified or the total number

406 of encoded proteins. Even if VIBRANT identified fewer total viral sequences compared to other
 407 tools in certain circumstances, more data of higher quality was generated as viral sequences of
 408 longer length were identified as compared to many short fragments. This provides an important
 409 distinction that the metric of total viral predictions is not necessarily an accurate representation for
 410 the quality or quantity of the data generated.
 411

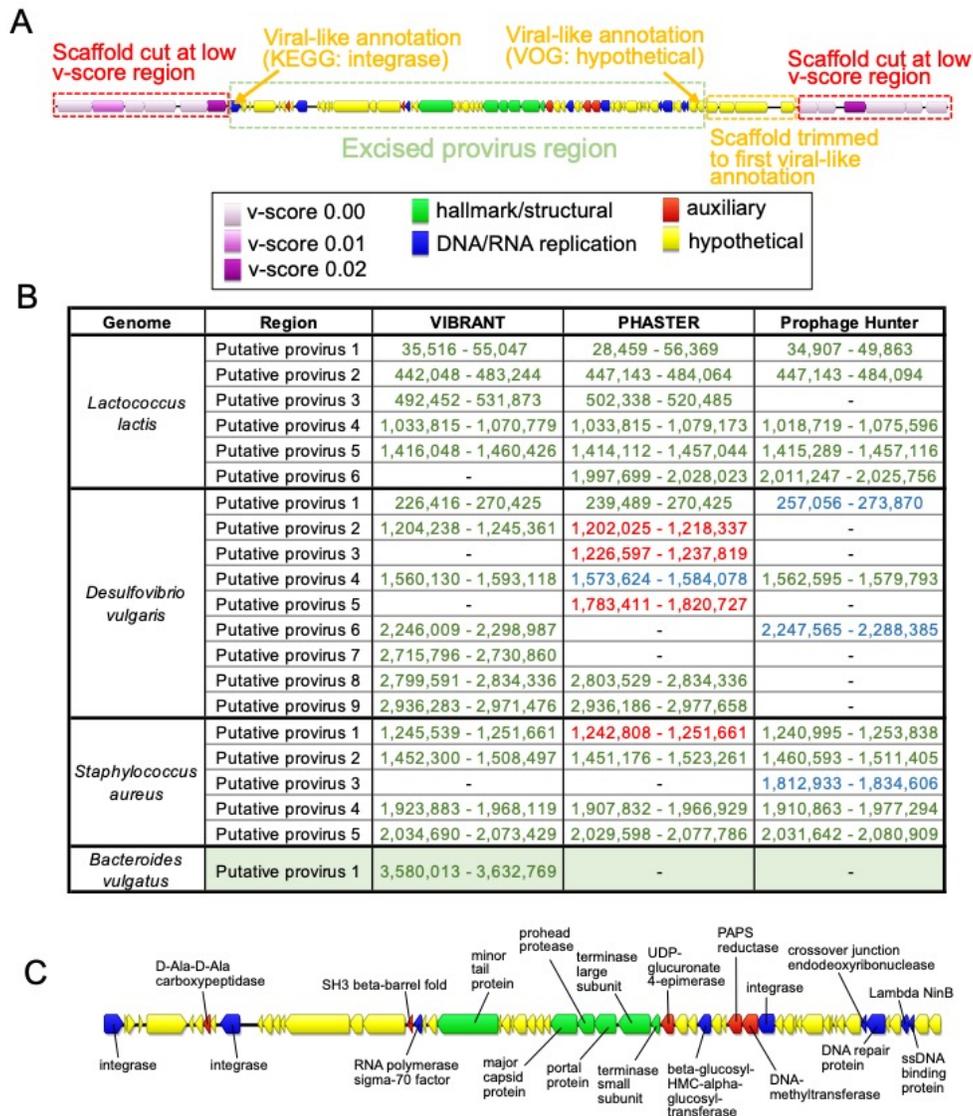


Figure 4. Prediction of integrated proviruses by VIBRANT, and comparison to PHASTER and Prophage Hunter. (A) Schematic representing the method used by VIBRANT to identify and extract provirus regions from host scaffolds using annotations. Briefly, v-scores are used to cut scaffolds at host-specific sites and fragments are trimmed to the nearest viral annotation. (B) Comparison of proviral predictions within four complete bacterial genomes between VIBRANT, PHASTER and Prophage Hunter. 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. (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.

412 **Integrated provirus prediction**

413 In many environments, integrated proviruses can account for a substantial portion of the
414 active viral community [66]. Despite this, few tools exist that are capable of identifying both lytic
415 viruses from metagenomic scaffolds as well as proviruses that are integrated into host genomes.
416 To account for this important group of viruses, VIBRANT identifies provirus regions within
417 metagenomic scaffolds or whole genomes. VIBRANT is unique from most provirus prediction
418 tools in that it does not rely on sequence motifs, such as integration sites, and therefore is especially
419 useful for partial metagenomic scaffolds in which neither the provirus nor host region is complete.
420 In addition, this functionality of VIBRANT provides the ability to trim non-viral (i.e., host
421 genome) ends from viral scaffolds. This results in a more correct interpretation of genes that are
422 encoded by the virus and not those that are misidentified as being within the viral genome region.
423 Briefly, VIBRANT identifies proviruses by first identifying and isolating scaffolds and genomes
424 at regions spanning several annotations with low v-scores. These regions were found to be almost
425 exclusive to host genomes. After cutting the original sequence at these regions, a refinement step
426 trims the putative provirus fragment to the first instance of a virus-like annotation to remove
427 leftover host sequence (Figure 4A). The final scaffold fragment is then analyzed by the neural
428 network similar to non-excised scaffolds.

429 To assess VIBRANT's ability to accurately extract provirus regions we compared its
430 performance to PHASTER and Prophage Hunter, two programs explicitly built for this task. We
431 compared the performance of these programs with VIBRANT on four bacterial genomes.
432 VIBRANT and PHASTER predicted an equal number of proviruses, 17, while Prophage Hunter
433 identified less, 13 (Figure 4B). Only one putative provirus prediction (*Lactococcus lactis* putative
434 provirus 6) was shared between PHASTER and Prophage Hunter but not VIBRANT. However,
435 VIBRANT was able to identify two putative provirus regions (*Desulfovibrio vulgaris* putative
436 provirus 7 and *Bacteroides vulgatus* putative provirus 1) that neither PHASTER nor Prophage
437 Hunter identified. Manual inspection of the putative *Bacteroides vulgatus* provirus identified a
438 number of *bona fide* virus hallmark and virus-like proteins suggesting that it is an accurate
439 prediction (Figure 4C). Our results suggest VIBRANT has the ability to accurately identify
440 proviruses and, in some cases, can outperform other tools in this task.

441 **Evaluating quality of viral scaffolds and genomes**

442 Determination of quality, in relation to completeness, of a viral scaffold has been
443 notoriously difficult due to the absence of universally conserved viral genes. To date the most
444 reliable metric of completeness for metagenomically assembled viruses is to identify circular
445 sequences (i.e., complete circular genomes). Therefore, the remaining alternatives rely on
446 estimation based on encoded proteins that function in central viral processes: replication of
447 genomes and assembly of new viral particles.

448 VIBRANT estimates the quality of predicted viral scaffolds, a relative proxy for
449 completeness, and indicates scaffolds that are circular. To do this, VIBRANT uses annotation
450 metrics of nucleotide replication and viral hallmark proteins. Hallmark proteins are those typically
451 specific to viruses and those that are required for productive infection, such as structural (e.g.,
452 capsid, tail, baseplate), terminase or viral holin/lysin proteins. Nucleotide replication proteins are
453 a variety of proteins associated with either replication or metabolism, such as nucleases,
454 polymerases and DNA/RNA binding proteins. Genomic scaffolds are categorized as low, medium
455 or high quality draft as determined by VOG annotations (Figure 5A, Additional File 2: Table S2).
456 High quality draft represents scaffolds that are likely to contain the majority of a virus's complete
457

458 genome and will contain annotations that are likely to aid in analysis of the virus, such as
 459 phylogenetic relationships and true positive verification. Medium draft quality represents the
 460 majority of a complete viral genome but is more likely to be a smaller portion in comparison to
 461 high quality. These scaffolds may contain annotations useful for analysis but are under less strict
 462 requirements compared to high quality. Finally, low draft quality constitutes scaffolds that were
 463 not found to be of high or medium quality. Many metagenomic scaffolds will likely be low quality
 464 genome fragments, but this quality category may still contain the higher quality genomes of some
 465 highly divergent viruses.

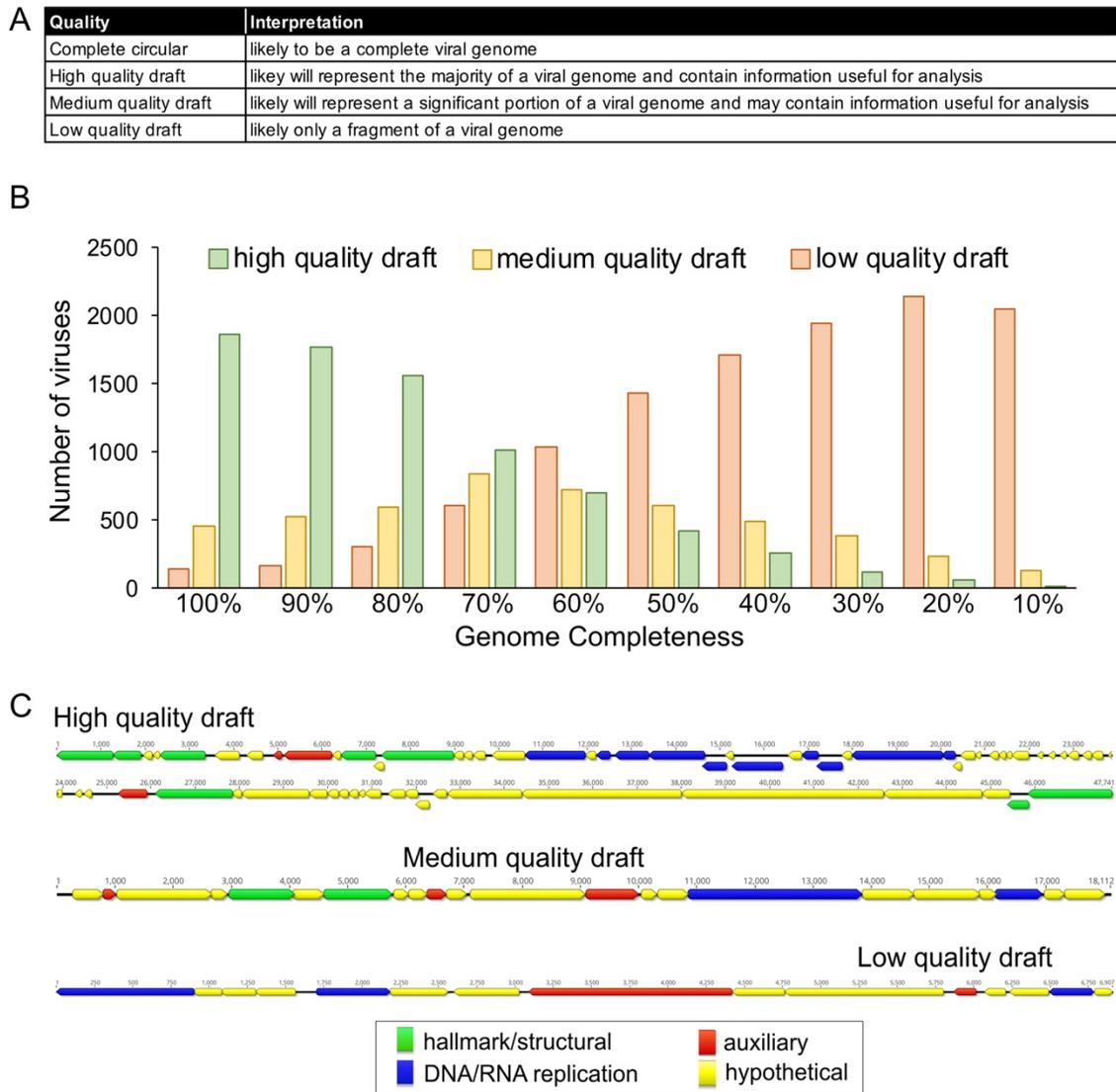


Figure 5. Estimation of genome quality of identified viral scaffolds. (A) Explanation of interpretation of quality categories: complete circular, high quality draft, medium quality draft and low quality draft. Quality generally represents total proteins, viral annotations, viral hallmark protein and 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.

466 We benchmarked VIBRANT's viral genome quality estimation using a total of 2466
467 *Caudovirales* genomes from NCBI RefSeq database. Genomes were evaluated either as complete
468 sequences or by removing 10% of the sequence at a time stepwise between 100% and 10%
469 completeness (Figure 5B). The results of VIBRANT's quality analysis displayed a linear trend in
470 indicating more complete genomes as high quality and less complete genomes as lower quality.
471 The transition from categorizing genomes as high quality to medium quality ranged from 60% and
472 70% completeness. Although we acknowledge that VIBRANT's metrics are not perfect, we
473 demonstrate the first benchmarked approach to quantify and characterize genome quality
474 associated with completeness of viral scaffolds. Manual inspection and visual verification of viral
475 genomes that were characterized into each of these genome quality categories showed that quality
476 estimations matched annotations (Figure 5C).

477

478 **Identifying function in virome: metabolic analysis**

479 Viruses are a dynamic and key facet in the metabolic networks of microbial communities
480 and can reprogram the landscape of host metabolism during infection. This can often be achieved
481 by modulating host metabolic networks through expression of AMGs encoded on viral genomes.
482 Identifying these AMGs and their associated role in the function of communities is imperative for
483 understanding complex microbiome dynamics, or in some cases can be used to predict virus-host
484 relationships. VIBRANT is optimized for the evaluation of function in viromes by identifying and
485 classifying the metabolic capabilities of the viral community. To do this, VIBRANT identifies
486 AMGs and assigns them into specific metabolic pathways and broader categories as designated by
487 KEGG annotations.

488 To highlight the utility of this feature we compared the metabolic function of viruses
489 derived from several diverse environments: freshwater, marine, soil, human-associated and city
490 (Additional File 21: Figure S1). We found natural environments (freshwater, marine and soil) to
491 display a different pattern of metabolic capabilities compared to human environments (human-
492 associated and city). Viruses originating from natural environments tend to largely encode AMGs
493 for amino acid and cofactor/vitamin metabolism with a more secondary focus on carbohydrate and
494 glycan metabolism. On the other hand, AMGs from city and human environments are dominated
495 by amino acid metabolism, and to some extent cofactor/vitamin and sulfur relay metabolism. In
496 addition to this broad distinction, all five environments appear slightly different from each other.
497 Despite freshwater and marine environments appearing similar in the ratio of AMGs by metabolic
498 category, the overlap in specific AMGs is less extensive. The dissimilarity between natural and
499 human environments is likewise corroborated by the relatively low overlap in individual AMGs.

500 A useful observation provided by VIBRANT's metabolic analysis is that there appears to
501 be globally conserved AMGs (i.e., present within at least 10 of the 13 environments tested). These
502 14 genes—*dcm*, *cysH*, *folE*, *phnP*, *ubiG*, *ubiE*, *waaF*, *moeB*, *ahbD*, *cobS*, *mec*, *queE*, *queD*,
503 *queC*—likely perform functions that are central to viral replication regardless of host or
504 environment. Notably, *folE*, *queD*, *queE* and *queC* constitute the entire 7-cyano-7-deazaguanine
505 (preQ₀) biosynthesis pathway, but the remainder of queuosine biosynthesis are entirely absent with
506 the exception of *queF*. Certain AMGs are unique in that they are the only common representatives
507 of a pathway amongst all AMGs identified, such as *phnP* for methylphosphonate degradation.
508 These AMGs may indicate an evolutionary advantage for manipulating a specific step of a
509 pathway, such as overcoming a reaction bottleneck, as opposed to modulating an entire pathway
510 as seen with preQ₀ biosynthesis. However, it should be noted that this list of 14 globally conserved
511 AMGs may not be entirely inclusive of the core set of AMGs in a given environment.

512 VIBRANT was evaluated for its
 513 ability to provide new insights into virome
 514 function by highlighting AMGs from mixed
 515 metagenomes. Using only data from
 516 VIBRANT's direct outputs, we compared the
 517 viral metabolic profiles of 6 hydrothermal
 518 vent and 15 human gut metagenomes (Figure
 519 6). As anticipated, based on IMG/VR
 520 environment comparisons, the metabolic
 521 capabilities between the two environments
 522 were different even though the number of
 523 unique AMGs was relatively equal (138 for
 524 hydrothermal vents and 151 for human gut).
 525 The pattern displayed by metabolic
 526 categories for each metagenome was similar
 527 to that displayed by marine and human
 528 viromes. For hydrothermal vents the
 529 dominant AMGs were part of carbohydrate,
 530 amino acid and cofactor/vitamin metabolism,
 531 whereas human gut AMGs were mostly
 532 components of amino acid and, to some
 533 extent, cofactor/vitamin metabolism.
 534 Although the observed AMGs and metabolic
 535 pathways were overall different, about a third
 536 (50 total AMGs) of all AMGs from each
 537 environment were shared; between these
 538 metagenomes alone all 14 globally conserved
 539 AMGs were present.

540 Observations of individual AMGs
 541 provided insights into how viruses interact
 542 within different environments. For example,
 543 tryptophan 7-halogenase (*prnA*) was
 544 identified in high abundance (45 total AMGs)
 545 within hydrothermal vent metagenomes but
 546 was absent from the human gut. Verification using GOV2 (Global Ocean Viromes 2.0) [67] and
 547 Human Gut Virome databases supported our finding that *prnA* appears to be constrained to aquatic
 548 environments, which is further supported by the gene's presence on several marine cyanophages.
 549 PrnA catalyzes the initial reaction for the formation of pyrrolnitrin, a strong antifungal antibiotic.
 550 Identification of this AMG only within aquatic environments suggests a directed role in aquatic
 551 virus lifestyles. Similarly, cysteine desulphydrase (*iscS*) was abundant (14 total AMGs) within the
 552 human gut metagenomes but not hydrothermal vents.

553

554 Application of VIBRANT: Identification of viruses from individuals with Crohn's Disease

555 We applied VIBRANT to identify viruses of at least 5kb in length from 102 human gut
 556 metagenomes (discovery dataset): 49 from individuals with Crohn's Disease and 53 from healthy
 557 individuals [68,69]. VIBRANT identified 14,121 viruses out of 511,977 total scaffolds. These viral

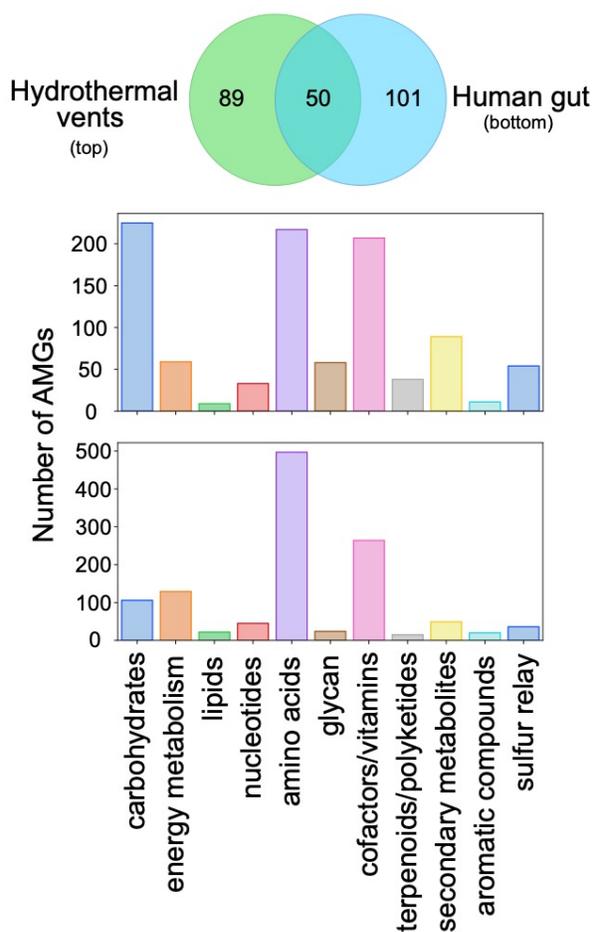


Figure 6. Comparison of AMG metabolic categories between hydrothermal vents and human gut. The Venn diagram depicts the unique and shared non-redundant AMGs between 6 hydrothermal vent and 15 human gut metagenomes. The graphs depict the differential abundance of KEGG metabolic categories of respective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

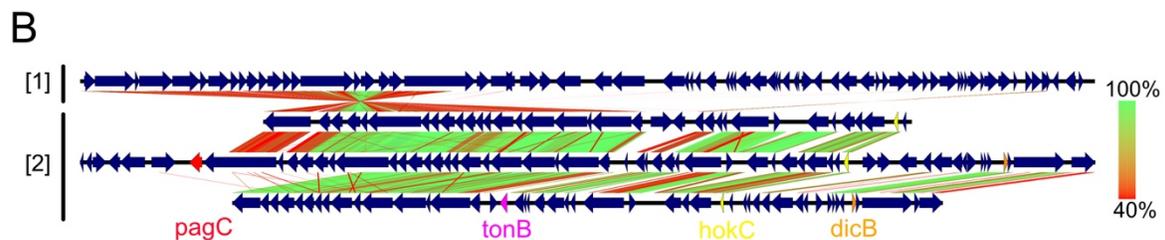
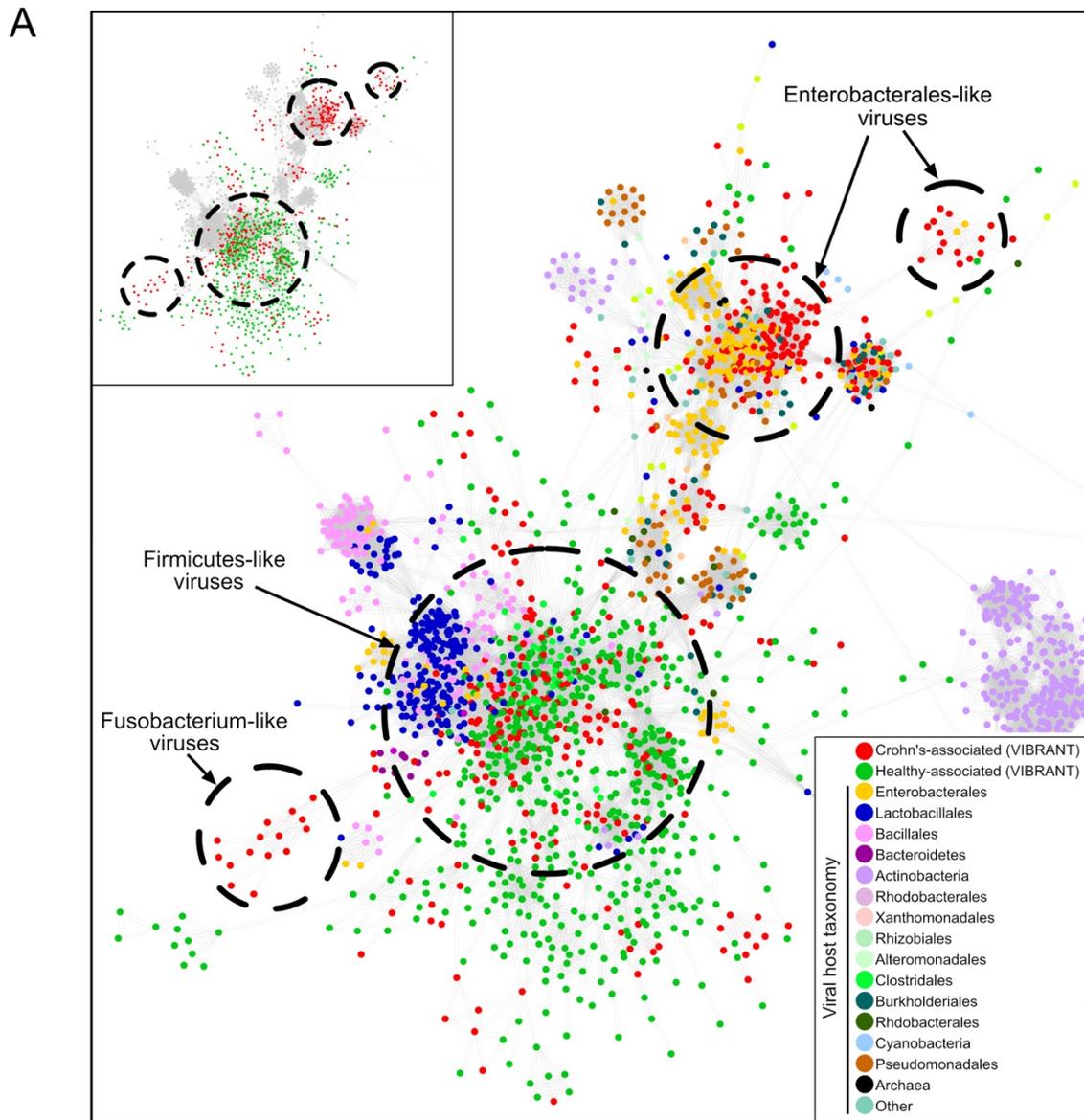


Figure 7. Viral metabolic comparison between Crohn's Disease and healthy individuals gut metagenomes.

(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*.

559 scaffolds were dereplicated to 8,822 non-redundant viral genomes using a cutoff of 95% nucleotide
560 identity over at least 70% of the scaffold. We next used read coverage of each virus from all 102
561 metagenomes to calculate relative differential abundance across Crohn's Disease and healthy
562 individuals. In total, we found 721 viruses to be more abundant in the gut microbiomes associated
563 with Crohn's Disease (Crohn's-associated) and 950 to be more abundant in healthy individuals
564 (healthy-associated).

565 Using these viruses identified by VIBRANT we sought to identify taxonomic or host-
566 association relationships to differentiate the virome of individuals with Crohn's Disease. We used
567 vConTACT2 to cluster the 721 Crohn's- or 950 healthy-associated viruses with reference genomes
568 using protein similarity. The majority of viruses (95.5%) were not clustered with any reference
569 genome at approximately the genus level suggesting VIBRANT may have identified a large pool
570 of novel or unique viral genomes. Although fewer total viruses were associated with Crohn's
571 Disease, significantly more were clustered to at least one representative at the genus level (72 for
572 Crohn's and 4 for healthy). Interestingly, no differentially abundant viruses from healthy
573 individuals clustered with Enterobacterales-infecting reference viruses (enteroviruses), yet the
574 majority (60/76) of Crohn's-associated viruses were clustered with known enteroviruses, such as
575 Lambda- and Shigella-related viruses. The remaining 16 viruses mainly clustered with
576 *Caudovirales* infecting *Lactococcus*, *Clostridium*, *Riemerella*, *Klebsiella* and *Salmonella* species,
577 though *Microviridae* and a likely complete crAssphage were also identified. A significant
578 proportion of all Crohn's-associated viruses (250/721), and the majority of genus-level clustered
579 viruses (42/76), were found to be integrated sequences within a microbial genomic scaffold but
580 were able to be identified due to VIBRANT's ability to excise proviruses.

581 We also generated a protein sharing network containing all 721 Crohn's and 950 healthy-
582 associated viruses, which corresponded to taxonomic and host relatedness (Figure 7A). This
583 protein network identified two different clustering patterns: [1] overlapping Crohn's and healthy-
584 associated viral populations clustered with Firmicutes-like viruses which may be indicative of a
585 stable gut virome; [2] Crohn's-associated viruses clustered with Enterobacterales-like and
586 Fusobacterium-like viruses which may be indicative of a state of dysbiosis. The presence of a
587 greater diversity and abundance of Enterobacterales and Fusobacteria has previously been linked
588 to Crohn's Disease [70,71], and therefore the presence of viruses infecting these bacteria may
589 provide similar information.

590 VIBRANT provides annotation information for all of the identified viruses which can be
591 used to infer functional characteristics in conjunction with host association. Comparison of
592 Crohn's-associated Lambda-like virus genomic content and arrangement suggested a possible role
593 of virally encoded host-persistence and virulence genes that are absent in the healthy-associated
594 virome (Figure 7B). Among all Crohn's-associated viruses, 17 total genes (*bor*, *dicB*, *dicC*, *hokC*,
595 *kilR*, *pagC*, *ydaS*, *ydaT*, *yfdN*, *yfdP*, *yfdQ*, *yfdR*, *yfdS*, *yfdT*, *ymfL*, *ymfM* and *tonB*) that have the
596 potential to impact host survival or virulence were identified. Importantly, no healthy-associated
597 viruses encoded such genes (Table 2). The presence of these putative dysbiosis-associated genes
598 (DAGs) may contribute to the manifestation and/or persistence of disease, similar to what has been
599 proposed for the bacterial microbiome [72–74]. For example, *pagC* encodes an outer membrane
600 virulence factor associated with enhanced survival of the host bacterium within the gut [75]. The
601 identification of *dicB* encoded on a putative *Escherichia* virus is unique in that it may represent a
602 'cryptic' provirus that protects the host from lytic viral infection, thus likely to enhance the ability
603 of the host to survive within the gut [76]. Finally, *hokC* may indicate mechanisms of virally
604 encoded virulence [77].

605 To characterize the distribution and association of DAGs with Crohn's Disease, we
 606 calculated differential abundance for two DAG-encoding viruses across all metagenome samples.
 607 The first virus encoded *pagC* and *yfdN*, and the second encoded *dicB*, *dicC* and *hokC*. Comparison
 608 of Crohn's Disease to healthy metagenomes indicates these viruses are present within the gut
 609 metagenomes of multiple individuals but more abundant in association with Crohn's Disease
 610 (Figure 8A). This suggests an association of disease with not only putative DAGs, but also specific,
 611 and potentially persistent, viral groups that encode them. In order to correlate increased abundance
 612 with biological activity we calculated the index of replication (iRep) for each of the two viruses
 613 [78]. Briefly, iRep is a function of differential read coverage which is able to provide an estimate
 614 of active genome replication. Seven metagenomes containing the greatest abundance for each virus
 615 were selected for iRep analysis and indicated that each virus was likely active at the time of
 616 collection (Figure 8B).

617 To validate these aforementioned findings, we applied VIBRANT to two additional
 618 metagenomic datasets from cohorts of individuals with Crohn's disease and healthy individuals
 619 (validation dataset): 43 from individuals with Crohn's Disease and 21 from healthy individuals
 620 [79,80]. VIBRANT identified 3,759 redundant viral genomes from Crohn's-associated
 621 metagenomes and 1,444 from healthy-associated metagenomes. Determination of protein
 622 networks and visualization similarly identified clustering of Crohn's-associated viruses with
 623 reference enteroviruses (Additional File 2: Figure S2). Likewise, we were able to identify 15 out
 624 of the 17 putative DAGs to be present in higher abundance in the Crohn's Disease microbiome.
 625 This validates our findings of the presence of unique viruses and proteins associated with Crohn's
 626 Disease, and suggests Enterobacterales-like viruses and putative DAGs may act as markers of
 627 Crohn's Disease. Overall, our results suggest that VIBRANT provides a platform for
 628 characterizing these relationships.

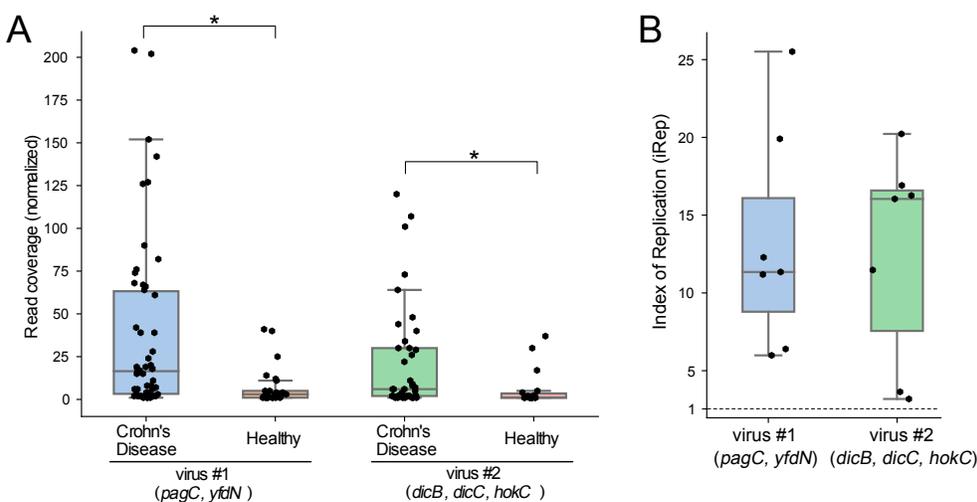


Figure 8. Differential abundance and activity of two viruses associated with Crohn's Disease. (A) Normalized read coverage of two Crohn's-associated viruses that encode putative DAGs between Crohn's Disease and healthy gut metagenomes. Asterisks represent significant differential abundance ($p < 0.05$). (B) iRep analysis for the same two viruses as (A), representative of seven metagenomes per virus. The dotted line indicates an iRep value of one, or low to no activity.

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 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 <i>dicB</i>	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

630

631 Discussion

632 Viruses that infect bacteria and archaea are key components in the structure, dynamics, and
633 interactions of microbial communities. Tools that are capable of efficient recovery of these viral
634 genomes from mixed metagenomic samples are likely to be fundamental to the growing
635 applications of metagenomic sequencing and analyses. Importantly, such tools would need to
636 reduce bias associated with specific viral groups (e.g., *Caudovirales*) and highly represented
637 environments (e.g., marine). Moreover, viruses that exist as integrated proviruses within host
638 genomes should not be ignored as they can represent a substantial fraction of infections in certain
639 conditions and also persistent infections within a community.

640 Here we have presented VIBRANT, a novel method for the automated recovery of both
641 free and integrated viral genomes from metagenomes that hybridizes neural network machine
642 learning and protein signatures. VIBRANT utilizes metrics of non-reference based protein
643 similarity annotation from KEGG, Pfam and VOG databases in conjunction with a novel ‘v-score’
644 metric to recover viruses with little to no biases. VIBRANT was built with the consideration of
645 the human guided intuition used to manually inspect metagenomic scaffolds for viral genomes and
646 packages these ideas into an automated software. This platform originates from the notion that
647 proteins generally considered as non-viral, such as ribosomal proteins [81], may be decidedly
648 common amongst viruses and should be considered accordingly when viewing annotations. V-
649 scores are meant to provide a quantitative metric for the level of virus-association for each
650 annotation used by VIBRANT, especially for Pfam and KEGG HMMs. That is, v-scores provide
651 a means for both highlighting common or hallmark viral proteins as well as differentiating viral
652 from non-viral annotations. In addition, v-scores give a quantifiable value to viral hallmark genes
653 instead of categorizing them in a binary fashion.

654 VIBRANT was not only built for the recovery of viral genomes, but also to act as a platform
655 for investigating the function of a virome. VIBRANT supports the analysis of virome function by
656 assembling useful annotation data and categorizing the metabolic pathways of viral AMGs. Using
657 annotation signatures, VIBRANT furthermore is capable of estimating genome quality and
658 distinguishing between lytic and lysogenic viruses. To our knowledge, VIBRANT is the first
659 software that integrates virus identification, annotation and estimation of genome completeness
660 into a stand-alone program.

661 Benchmarking and validation of VIBRANT indicated improved performance compared to
662 VirSorter and VirFinder, two commonly used programs for identifying viruses from metagenomes.
663 This included a substantial increase in the relationship between true virus identifications (recall,
664 true positive rate) and false non-virus identifications (specificity, false positive rate). That is,
665 VIBRANT recovered more viruses with no discernable expense to false identifications. The result
666 was that VIBRANT was able to recover an average of 2.4 and 1.7 more viral sequence from real
667 metagenomes than VirFinder and VirSorter, respectively. When tested on metagenome-assembled
668 viral genomes from IMG/VR representing diverse environments VIBRANT was found to have no
669 perceivable environment bias towards identifying viruses. In comparison to provirus prediction
670 tools, specifically PHASTER and Prophage Hunter, VIBRANT was shown to be proficient in
671 identifying viral regions within bacterial genomes. This included the identification of a putative
672 *Bacteroides* provirus that the other two programs were unable to identify. The importance of
673 integrated provirus prediction was underscored in the analysis of Crohn's Disease metagenomes
674 since it was found that a significant proportion of disease related viruses were temperate viruses
675 existing as host-integrated genomes.

676 VIBRANT's method allows for the distinction between scaffold size and coding capacity
677 in designating the minimum length of virus identifications. Traditionally, a cutoff of 5000 bp has
678 been used to filter for scaffolds of a sufficient length for analysis. This is under the presumption
679 that a longer sequence will be likely to encode more proteins. For example, this cutoff has been
680 adopted by IMG/VR. However, we suggest a total protein cutoff of four open reading frames rather
681 than sequence length cutoff to be more suitable for comprehensive characterization of the viral
682 community. VIBRANT's method works as a strict function of total encoded proteins and is
683 completely agnostic to sequence length for analysis. Therefore, the boundary of minimum encoded
684 proteins will support a more guided cutoff for quality control of virus identifications. For example,
685 increasing the minimum sequence length to 5000 bp will have no effect on accuracy or ability to
686 recall viruses since VIBRANT will only be considerate of the minimum total proteins, which is
687 set to four. The result will be the loss of all 1000 bp to 4999 bp viruses that still encode at least
688 four proteins. To visualize this distinction, we applied VIBRANT with various length cutoffs to
689 the previously used estuary virome (see Table 1). Input sequences were stepwise limited from
690 1000 bp to 10000 bp (1000 bp steps) or four open reading frames to 13 open reading frames (one
691 open reading frame steps) in length. Limiting to open reading frames indicated a reduced drop-off
692 in total virus identifications and total viral sequence compared to a minimum sequence length limit
693 (Additional File 23: Figure S3).

694 The output data generated by VIBRANT—protein/gene annotation information,
695 protein/gene sequences, HMM scores and e-values, viral sequences in FASTA and GenBank
696 format, indication of AMGs, genome quality, etc.—provides a platform for easily replicated
697 pipeline analyses. Application of VIBRANT to characterize the function of Crohn's-associated
698 viruses emphasizes this utility. VIBRANT was not only able to identify a substantial number of
699 viral genomes, but also provided meaningful information regarding putative DAGs, viral

700 sequences for differential abundance calculation and genome alignment, viral proteins for
701 clustering, and AMGs for metabolic comparisons.

702

703 **Conclusions**

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

718

719 **Methods**

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

721 To generate training and testing datasets sequences representing bacteria, archaea,
722 plasmids and viruses were downloaded from NCBI databases (accessed July 2019) (Additional
723 File 3: Table S3). For bacteria/archaea, 181 genomes from diverse phylogenetic groups were
724 randomly chosen. Likewise, a total of 1,452 bacterial plasmids were chosen. For viruses, NCBI
725 taxids associated with viruses that infect bacteria or archaea were used to download reference virus
726 genomes, which were then limited to only sequences above 3kb. Sequences not associated with
727 genomes, such as partial genomic regions, were manually removed. This resulted in 15,238 total
728 viral genomes. All sequences were split into non-overlapping fragments between 3kb and 15kb to
729 simulate metagenome assembled scaffolds (hereafter called *fragments*).

730 Integrated viruses are common in both bacteria and archaea. To address this for generating
731 a dataset devoid of viruses, PHASTER (accessed July 2019) was used to predict putative integrated
732 viruses in the 181 bacteria/archaea genomes. Using BLASTn [82], any fragments that had
733 significant similarity (at least 95% identity, at least 3kb coverage and e-value < 1e-10) to the
734 PHASTER predictions were removed as contaminant virus sequence. The new bacteria/archaea
735 dataset was considered depleted of prophages, but not entirely devoid of contamination. Next, the
736 datasets for bacteria/archaea and plasmids were annotated with KEGG, Pfam and VOG
737 (hmmsearch (v3.1), e-value < 1e-5) [83] to further remove contaminant virus sequence. Plasmids
738 were included because it was noted that the dataset appeared to contain virus sequences, possibly
739 due to misclassification of episomal proviruses as plasmids. Using manual inspection of the
740 KEGG, Pfam and VOG annotations any sequence that clearly belonged to a virus was removed.
741 The final datasets consisted of 400,291 fragments for bacteria/archaea, 14,739 for plasmids, and
742 111,963 for viruses.

743

744 **V-score generation**

745 Reference and database viral proteins were used to generate v-scores. To be consistent
746 between all 15,238 viruses acquired from NCBI, proteins were predicted for all genomes using
747 Prodigal (-p meta, v2.6.3) [84]. All VOG proteins were added to this dataset, which resulted in a
748 total of 633,194 proteins. Redundancy was removed from the generated viral protein dataset using
749 cdhit (v4.6) [85] with a identify cutoff of 95%, which resulted in a total of 240,728 viral proteins
750 (Additional File 4: Table S4). This was the final dataset used to generate v-scores. All KEGG
751 HMM profiles to be used by VIBRANT (method described below) were used to annotate the viral
752 proteins. A v-score for each KEGG HMM profile was determined by the number of significant (e-
753 value < 1e-5) hits by hmmsearch, divided by 100, and a maximum value was set at 10 after
754 division. The same v-score generation was done for Pfam and VOG databases. Any HMM profile
755 with no significant hits to the virus dataset was given a v-score of zero. For KEGG and Pfam
756 databases, any annotation that was given a v-score above zero and contained the keyword “phage”
757 was given a minimum v-score of 1. To highlight viral hallmark genes, any annotation within all
758 three databases with the keyword *portal*, *terminase*, *spike*, *capsid*, *sheath*, *tail*, *coat*, *virion*, *lysin*,
759 *holin*, *base plate*, *lysozyme*, *head* or *structural* was given a minimum v-score of 1. Non-phage
760 annotations (e.g., *phage shock protein*, *reovirus core-spike protein*) were not considered. The
761 resulting v-scores are a metric of virus association (i.e., do not take into account virus specificity,
762 or association with non-viruses) and are manually tuned to put greater weight on viral hallmark
763 genes (Additional File 5: Table S5). Raw HMM table outputs can be found in Additional Files 6,
764 7 and 8 for KEGG, Pfam and VOG, respectively (Additional File 6: Table S6, Additional File 7:
765 Table S7, and Additional File 8: Table S8).

766

767 **Databases used by VIBRANT**

768 VIBRANT uses HMM profiles from three different databases: KEGG, Pfam and VOG
769 (Additional File 9: Table S9). For Pfam all HMM profiles were used. To increase speed, KEGG
770 and VOG HMM databases were reduced in size to contain only profiles likely to annotate the
771 viruses of interest. For KEGG this was done by only retaining profiles considered to be relevant
772 to “prokaryotes” as determined by KEGG. For VOG this was done by only retaining profiles that
773 had at least one significant hit to an NCBI-acquired viral protein database using BLASTp. That is,
774 any VOG HMM profile given a v-score of zero was removed. The resulting databases consisted
775 of 10,033 HMM profiles for KEGG, 17,929 for Pfam, and 19,182 for VOG.

776 Two additional databases consisting of redundant Pfam HMM profiles were also generated.
777 The first database consisted of virus annotations which were determined by a text search of
778 “bacteriophage” to the Pfam database. Only HMM profiles with v-scores above zero were
779 considered and those common to bacteria/archaea (e.g., glutaredoxin) were manually removed.
780 This resulted in 894 virus specific HMMs. The second database consisted of common plasmid
781 annotations. Proteins were predicted for the plasmid dataset using Prodigal (-p meta) and all Pfam
782 HMMs with a v-score of zero were used to annotate the plasmid proteins (e-value < 1e-5). Any
783 annotation with at least 50 hits was retained as a common plasmid HMM profile, which resulted
784 in 202 common plasmid HMMs.

785

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

787 VIBRANT utilizes several manually curated cutoffs in order to remove the bulk of non-
788 virus input scaffolds before the neural network classifier is implemented. These steps will result
789 in the assembly of 27 annotation metrics that are used by the neural network classifier for virus
790 identification, which is followed by additional manually set cutoffs to curate the results.

791 First, open reading frames predicted by Prodigal (-p meta) or user input proteins are used
792 to calculate the fraction of strand switching per scaffold (strand switches divided by total genes).
793 Scaffolds are then classified as having either a *low* (5%), *medium* (5-35%) or *high* (>35%) level
794 of strand switching. Scaffolds with a high level are annotated with the 894 virus-specific Pfam
795 HMMs and only retained if there is at least one significant hit (score > 50). Throughout, scaffolds
796 that are not retained are eliminated from further analysis. Scaffolds with a medium-level, and those
797 with a high-level that passed the previous cutoff, are annotated with the 202 common plasmid
798 Pfam HMMs and only retained if there are three or less significant hits (score > 50). Scaffolds with
799 a low level are combined with those from high/medium that passed the previous cutoff(s).

800 Scaffolds are then annotated with the 10,033 KEGG-derived HMMs. Putative integrated
801 provirus regions are extracted at this step by using sliding windows of either four or nine proteins
802 at a time (step size = 1 protein). Within these windows scaffolds are fragmented according to v-
803 scores and total KEGG annotations. Within the 4-protein window, scaffolds can be cut if [1] there
804 are 0-1 unannotated proteins, 3-4 proteins with a v-score of 0-0.02 and a combined v-score of less
805 than 0.06, or [2] three consecutive proteins with a v-score of 0 (considered as a 3-protein window).
806 Scaffolds will also be cut using a 9-protein window if nine consecutive proteins are annotated.
807 Finally, if the final two proteins on a scaffold each have a v-score of 0, the scaffold will be cut.
808 Only scaffold fragments that contain at least 8 proteins are retained. Following provirus excision,
809 several manual cutoffs are used to remove obvious non-viral scaffolds. Briefly, this is done by
810 removing scaffolds with a high density of KEGG annotations (e.g., over 70% if less than 15
811 proteins or over 50% if greater than 15 proteins) or a high number of annotations with a v-score
812 of 0 (e.g., over 15). V-scores are also used such that a scaffold that may be removed for having a
813 high density of KEGG annotations will be retained if the v-score meets a specific threshold (e.g.,
814 average of 0.2).

815 Scaffolds that are retained are annotated by the 17,929 Pfam HMMs. In a similar manner
816 to KEGG, scaffolds meeting set cutoffs for density and v-scores of Pfam HMMs are either retained
817 or removed. For example, scaffolds with less than 15 total or density under 60% Pfam annotations
818 are retained; a scaffold will be retained if it has greater than 60% Pfam annotations as well as an
819 average v-score of at least 0.15. For both KEGG and Pfam cutoffs full details of every cutoff see
820 Additional File 10: Table S10.

821 Following the aforementioned cutoff steps approximately 75-85% of non-viral scaffolds
822 are removed. At this point scaffolds are annotated by the 19,182 VOG HMMs. Using VOG
823 annotations and v-scores from KEGG and Pfam, putative proviruses that were cut during KEGG
824 annotation are now trimmed to remove ends that may still contain host proteins. To do this, any
825 scaffold previously cut is trimmed, at both ends, to either the first instance of a VOG annotation
826 or the first v-score of at least 0.1.

827 Annotations from all three databases are used to assemble 27 metrics for the neural network
828 classifier. Briefly the metrics are as follows: [1] total proteins, [2] total KEGG annotations, [3]
829 sum of KEGG v-scores, [4] total Pfam annotations, [5] sum of Pfam v-scores, [6] total VOG
830 annotations, [7] sum of VOG v-scores, [8] total KEGG integration related annotations (e.g.,
831 integrase), [9] total KEGG annotations with a v-score of zero, [10] total KEGG integration related
832 annotations (e.g., integrase), [11] total Pfam annotations with a v-score of zero, [12] total VOG
833 redoxin (e.g., glutaredoxin) related annotations, [13] total VOG non-integrase integration related
834 annotations, [14] total VOG integrase annotations, [15] total VOG ribonucleotide reductase related
835 annotations, [16] total VOG nucleotide replication (e.g., DNA polymerase) related annotations,
836 [17] total KEGG nuclease (e.g., restriction endonuclease) related annotations, [18] total KEGG

837 toxin/anti-toxin related annotations, [19] total VOG hallmark protein (e.g., capsid) annotations,
838 [20] total proteins annotated by KEGG, Pfam and VOG, [21] total proteins annotated by Pfam and
839 VOG only, [22] total proteins annotated by Pfam and KEGG only, [23] total proteins annotated by
840 KEGG and VOG only, [24] total proteins annotated by KEGG only, [25] total proteins annotated
841 by Pfam only, [26] total proteins annotated by VOG only, and [27] total unannotated proteins.
842 Non-annotation features such as gene density, average gene length and strand switching were not
843 used because they were found to decrease performance of the neural network classifier despite
844 being differentiating features between bacteria/archaea and viruses; viruses tend to have shorter
845 genes, less intergenic space and strand switch less frequently. This decreased performance is likely
846 due to several reasons, such as errors associated with protein prediction (e.g., missed open reading
847 frame leading to a large “intergenic” gap) or that scaffolds, due to being fragmented genomes in
848 most cases, behave differently than the genome as a whole. For example, genomic regions
849 encoding for large structural proteins will have a higher average gene size or a small window of
850 virus proteins may have a greater average strand switching level compared to the whole genome.

851

852 **Training and testing VIBRANT**

853 The bacteria/archaea genomic, plasmid and virus datasets described above were used to
854 train and test the machine learning model. Scikit Learn libraries were used to assess various
855 machine learning strategies to identify the best performing algorithm. Among support vector
856 machines, neural networks and random forests, we found that neural networks lead to the most
857 accurate and comprehensive identification of viruses. Therefore, Scikit Learn’s [86] supervised
858 neural network multi-layer perceptron classifier (hereafter neural network) was used. The portion
859 of VIBRANT up until the neural network classifier (i.e., KEGG, Pfam and VOG annotation) was
860 used to compile the 27 annotation metrics for each of the three datasets. To account for differences
861 in scaffold sizes all metrics were normalized (i.e., divided by) to the total number of proteins
862 encoded by the scaffold. The first metric, for total proteins, was normalized to log base 10 of itself.
863 Each metric was weighted equally, though it is worth noting that the removal of several metrics,
864 mainly metrics 8-18, did not significantly impact the accuracy of model’s prediction. The
865 normalized results were randomized and non-redundant portions of these results were taken for
866 training or testing the neural network. It is important to note that the testing set here was not used
867 as the comprehensive testing set for the entire workflow. In total, 93,913 fragments were used for
868 training and 9,000 were used for testing the neural network (Additional File 11: Table S11 and
869 Additional File 12: Table S12).

870 To comprehensively test the performance of VIBRANT in its entirety a new testing dataset
871 was generated consisting of fragments from the neural network testing set as well as additional
872 fragments non-redundant to the previous training dataset. This new testing dataset was comprised
873 of 256,713 fragments from bacteria/archaea, 29,926 from viruses and 8,968 from plasmids. Each
874 met the minimum size requirement of VIBRANT: at least four open reading frames. For
875 comparison to VirFinder (v1.1) and VirSorter (v1.0.3), the latter testing dataset was used. Two
876 intervals for VirFinder and VirSorter were used for comparison. For VirSorter, the intervals
877 selected were [1] category 1 and 2 predictions, and [2] categories 1, 2 and 3 (i.e., all) predictions.
878 VirSorter was ran using the “Virome” database. For VirFinder, the intervals were [1] scores greater
879 than or equal to 0.90 (approximately equivalent to a p-value of 0.013), and [2] scores greater than
880 or equal to 0.75 (approximately equivalent to a p-value of 0.037). All equations used can be found
881 in Additional File 13: Table S13 and results used for the generation of Figure 1 can be found in
882 Additional File 14: Table S14.

883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928

AMG identification

KEGG annotations were used to classify potential AMGs (Additional File 15: Table S15). KEGG annotations falling under the “metabolic pathways” category as well as “sulfur relay system” were considered. Manual inspection was used to remove non-AMG annotations, such as *nrdAB* and *thyAX*. Other annotations not considered dealt with direct nucleotide to nucleotide conversions. All AMGs were associated with a KEGG metabolic pathway map.

Completeness estimation

Scaffold completeness is determined based on four metrics: circularization of scaffold sequence, VOG annotations, total VOG nucleotide replication proteins and total VOG viral hallmark proteins (Additional File 16: Table S16). In order to be considered a complete genome a sequence must be identified as likely circular. A kmer-based approach is used to do this. Specifically, the first 20 nucleotides are compared to 20-mer sliding windows within the last 900bp of the sequence. If a complete match is identified the sequence is considered a circular template. Scaffolds can also be considered a low, medium or high quality draft. To benchmark completeness, NCBI RefSeq viruses identified as *Caudovirales*, limited to 10 kb in length, were used to estimate completeness by stepwise removing 10% viral sequence at a time (Additional File 22: Table S2). Viral genome diagrams to depict genome quality and completeness, as well as provirus predictions, were made using Geneious Prime 2019.0.3.

Additional viral datasets and metagenomes

IMG/VR v2.0 (accessed July 2019) was downloaded and scaffolds originating from air, animal, aquatic sediment, city, marine A (coastal, gulf, inlet, intertidal, neritic, oceanic, pelagic and strait), marine B (hydrothermal vent, volcanic and oil), deep subsurface, freshwater, human, plants, soil wastewater and wetland environments were selected for analysis. Venn diagram visualization of virus predictions with this dataset was made using Matplotlib (v3.0.0) [87]. Several published, assembled metagenomes from IMG/VR representing diverse environments were selected for comparing VIBRANT, Virsorter and VirFinder (IMG taxon IDs: 3300005281, 3300017813 and 3300000439). Fifteen publicly available datasets from the human gut were assembled for assessing VIBRANT and comparing the three programs [88]. Reads can be found under NCBI BioProject PRJEB7774 (ERR688591, ERR688590, ERR688509, ERR608507, ERR608506, ERR688584, ERR688587, ERR688519, ERR688512, ERR688508, ERR688634, ERR688618, ERR688515, ERR688513, ERR688505). Reads were trimmed using Sickle (v1.33) [89] and assembled using metaSPAdes (v3.12.0 65) [90] (--meta -k 21,33,55,77,99). For hydrothermal vents, six publicly available hydrothermal plume samples were derived from Guaymas Basin (one sample) and Eastern Lau Spreading Center (five samples). Reads can be found under NCBI BioProject PRJNA314399 (SRR3577362) and PRJNA234377 (SRR1217367, SRR1217459, SRR1217564, SRR1217566, SRR1217452, SRR1217567, SRR1217465, SRR1217462, SRR1217460, SRR1217463, SRR1217565). Reads were trimmed using Sickle and assembled using metaSPAdes (--meta -k 21,33,55,77,99). Details of assembly and processing are outlined in Zhou *et al.* [91]. For analysis of Crohn’s Disease metagenomes by VIBRANT, publicly available metagenomes were used; the metagenomes were sequenced by He *et al.*, Ijaz *et al.* and Gevers *et al.*, and assembled by Pasolli *et al.* (Supplementary Tables 17, 18 and 19).

Analysis of Crohn’s Disease metagenomes

929 Metagenomic reads from He *et al.* were assembled by Pasolli *et al.* and used for analysis.
930 VIBRANT (-l 5000) was used to predict viruses from 49 metagenomes originating from
931 individuals with Crohn's Disease and 53 from healthy individuals (102 total samples). A total of
932 14,121 viruses were identified. Viral scaffolds were dereplicated using Mash [92] and Nucmer
933 [93] to 95% nucleotide identity and 70% scaffold coverage. The longest scaffold was kept as the
934 representative for a total of 8,822 dereplicated viral scaffolds. A total of 96 read sets were used
935 (59 Crohn's Disease and 37 healthy), trimmed using Sickle and aligned to the dereplicated
936 scaffolds using Bowtie2 (-N 1, v2.3.4.1) [94] and the resulting coverages were normalized to total
937 reads. The normalized relative coverage of each scaffold for all 96 samples were compared using
938 DESeq2 [95] (Additional File 20: Table S20). Scaffolds in significantly different abundance
939 between Crohn's Disease and control samples were determined by a p-value cutoff of 0.05. iRep
940 (default parameters) [78] was used to estimate replication activity of two Crohn's-associated
941 viruses. EasyFig (v2.2.2) [96] was used to generate genome alignments of Escherichia phage
942 Lambda (NCBI accession number NC_001416.1) and three Crohn's-associated viruses.
943 vConTACT2 (v0.9.8) was ran using default parameters on the CyVerse Discovery Environment
944 platform. Putative hosts of Crohn's-associated and healthy-associated was estimated using
945 proximity of vConTACT2 protein clustering and BLASTp identity (NCBI non-redundant protein
946 database, assessed October 2019). Two additional read sets from Gevers *et al.* [80] and Ijaz *et al.*
947 [79] were likewise assembled by Pasolli *et al.*. VIBRANT (-l 5000 -o 10) was used to predict
948 viruses from 43 metagenomes originating from individuals with Crohn's Disease and 21 from
949 healthy individuals (64 total samples). In contrast to the discovery dataset viral genomes were not
950 dereplicated and differential abundance was not determined. Instead viruses from each group were
951 directly clustered using vConTACT2. Abundances of DAGs in the validation set were normalized
952 to total viruses. Protein networks were visualized using Cytoscape (v3.7.2) [98].

953

954 **Acknowledgements**

955 We thank Upendra Devisetty for his assistance with dockerizing and integrating VIBRANT as a
956 web-based application in the CyVerse Discovery Environment. We thank the University of
957 Wisconsin - Office of the Vice Chancellor for Research and Graduate Education, University of
958 Wisconsin – Department of Bacteriology, and University of Wisconsin – College of Agriculture
959 and Life Sciences for their support.

960

961 **Availability of data and materials**

962 VIBRANT is implemented in Python and all scripts and associated files are freely available
963 at <https://github.com/AnantharamanLab/VIBRANT/>. All data and genomic sequences used for
964 analyses are publicly available; see Supplementary Tables 3, 17, 18 and 19 for study and accession
965 names. VIBRANT is also freely available for use as an application through the CyVerse Discovery
966 Environment. To use the application visit [https://de.cyverse.org/de/?type=apps&app-](https://de.cyverse.org/de/?type=apps&app-id=c2864d3c-fd03-11e9-9cf4-008cfa5ae621&system-id=de)
967 [id=c2864d3c-fd03-11e9-9cf4-008cfa5ae621&system-id=de](https://de.cyverse.org/de/?type=apps&app-id=c2864d3c-fd03-11e9-9cf4-008cfa5ae621&system-id=de), and for more details see
968 <https://wiki.cyverse.org/wiki/display/DEapps/VIBRANT-1.0.1>. Additional details of relevant data
969 are available from the corresponding author on request.

970

971 **Author Information**

972 **Affiliations**

973 *Department of Bacteriology, University of Wisconsin–Madison, Madison, WI, USA*

974 Kristopher Kieft, Zhichao Zhou, and Karthik Anantharaman

975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020

Contributors

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

Corresponding Author

Correspondence to Karthik Anantharaman

Ethics Declarations

Ethics approval and consent to participate

Not applicable.

Consent for publication

Not applicable.

Competing interests

The authors declare that they have no competing interests.

Additional Files

Additional File 1: Table S1. Number and sizes of sequence fragments used to train and test VIBRANT for viruses, plasmids, and bacteria and archaea.

Additional File 2: Table S2. Number of *Caudovirales* genomes and genomic fragments identified per quality estimation category, exact rules used to estimate genome quality and the interpretation of quality estimations.

Additional File 3: Table S3. List of NCBI accession numbers for bacterial and archaeal genomes, plasmids, and viral genomes used in this study.

Additional File 4: Table S4. Protein prediction coordinates of dereplicated proteins from NCBI viruses used to generate v-scores.

Additional File 5: Table S5. List of all KEGG, Pfam and VOG annotation names and associated v-scores (if greater than zero).

Additional File 6: Table S6. Unparsed HMM table output from KEGG annotations used to generate KEGG v-scores.

Additional File 7: Table S7. Unparsed HMM table output from Pfam annotations used to generate Pfam v-scores.

Additional File 8: Table S8. Unparsed HMM table output from VOG annotations used to generate VOG v-scores.

1021 **Additional File 9: Table S9.** List of all HMM names used by VIBRANT.
1022
1023 **Additional File 10: Table S10.** Description of set cutoffs implemented before neural network
1024 machine learning analysis for KEGG and Pfam annotations.
1025
1026 **Additional File 11: Table S11.** Normalized data used to train the neural network machine learning
1027 classifier.
1028
1029 **Additional File 12: Table S12.** Normalized data used to test the neural network machine learning
1030 classifier.
1031
1032 **Additional File 13: Table S13.** Equations used for benchmarking analyses.
1033
1034 **Additional File 14: Table S14.** Calculations and results of benchmarking analyses.
1035
1036 **Additional File 15: Table S15.** List of all KEGG annotations determined as AMGs.
1037
1038 **Additional File 16: Table S16.** List of all VOG annotations determined as nucleotide replication-
1039 associated or viral hallmark-associated, which are used during prediction and quality estimation.
1040
1041 **Additional File 17: Table S17.** List of datasets used from He *et al.*
1042
1043 **Additional File 18: Table S18.** List of datasets used from Ijaz *et al.*
1044
1045 **Additional File 19: Table S19.** List of datasets used from Gevers *et al.*
1046
1047 **Additional File 20: Table S20.** Results from DESeq2 analysis for 8,789 non-redundant viruses
1048 from the Crohn's Disease discovery dataset.
1049
1050 **Additional File 21: Figure S1. AMG and metabolic pathways between diverse environments.**
1051 VIBRANT was used to predict viruses from IMG/VR datasets and the identified metabolic
1052 pathways and AMGs were compared for freshwater, marine, soil, city and human-associated
1053 environments (graphs). The respective AMGs and their abundances were likewise compared (venn
1054 diagram).
1055
1056 **Additional File 22: Figure S2. Protein network of two Crohn's Disease validation datasets.**
1057 VIBRANT was used to predict viruses from two datasets for validation of marker virus and
1058 putative DAG discovery. The resulting viruses were used to construct a protein network indicating
1059 Crohn's-associated viruses clustering with enteroviruses more often than healthy-associated
1060 viruses.
1061
1062 **Additional File 23: Figure S3. Comparison of limiting to sequence length or open reading
1063 frames.** VIBRANT was used to predict viruses from an estuary virome and set to limit to either
1064 scaffold length or total encoded open reading frames. The (A) total virus identifications and (B)
1065 total viral sequence length were compared to show that limiting to open reading frames will
1066 typically yield more data.

1067
1068
1069
1070
1071

1072
1073

1074
1075

1076
1077

1078
1079

1080
1081
1082

1083
1084
1085

1086
1087
1088

1089
1090

1091

1092
1093

1094
1095

1096
1097
1098

1099

References

1. Breitbart M, Rohwer F. Here a virus, there a virus, everywhere the same virus? Trends in Microbiology. 2005;13:278–84.
2. Wommack KE, Colwell RR. Virioplankton: Viruses in Aquatic Ecosystems. Microbiol Mol Biol Rev. 2000;64:69–114.
3. Danovaro R, Serresi M. Viral Density and Virus-to-Bacterium Ratio in Deep-Sea Sediments of the Eastern Mediterranean. Appl Environ Microbiol. 2000;66:1857–61.
4. Suttle CA. Marine viruses — major players in the global ecosystem. Nature Reviews Microbiology. 2007;5:801–12.
5. Heldal M, Bratbak G. Production and decay of viruses in aquatic environments. Mar Ecol Prog Ser. 1991;72:205–12.
6. Gobler CJ, Hutchins DA, Fisher NS, Coper EM, Sañudo-Wilhelmy SA. Release and bioavailability of C, N, P Se, and Fe following viral lysis of a marine chrysophyte. Limnology and Oceanography. 1997;42:1492–504.
7. Jiao N, Herndl GJ, Hansell DA, Benner R, Kattner G, Wilhelm SW, et al. Microbial production of recalcitrant dissolved organic matter: long-term carbon storage in the global ocean. Nature Reviews Microbiology. 2010;8:593–9.
8. Brussaard CPD, Wilhelm SW, Thingstad F, Weinbauer MG, Bratbak G, Heldal M, et al. Global-scale processes with a nanoscale drive: the role of marine viruses. The ISME Journal. 2008;2:575–8.
9. Fuhrman JA. Marine viruses and their biogeochemical and ecological effects. Nature. 1999;399:541–8.
10. Wilhelm SW, Suttle CA. Viruses and Nutrient Cycles in the Sea. BioScience. 1999;49:8.
11. Norman JM, Handley SA, Baldrige MT, Droit L, Liu CY, Keller BC, et al. Disease-Specific Alterations in the Enteric Virome in Inflammatory Bowel Disease. Cell. 2015;160:447–60.
12. Barr JJ. Missing a Phage: Unraveling Tripartite Symbioses within the Human Gut. mSystems. 2019;4:e00105-19.
13. Barr JJ, Auro R, Furlan M, Whiteson KL, Erb ML, Pogliano J, et al. Bacteriophage adhering to mucus provide a non-host-derived immunity. Proceedings of the National Academy of Sciences. 2013;110:10771–6.
14. Rohwer F. Global Phage Diversity. Cell. 2003;113:141.

- 1100 15. Kim B, Kim ES, Yoo Y-J, Bae H-W, Chung I-Y, Cho Y-H. Phage-Derived Antibacterials:
1101 Harnessing the Simplicity, Plasticity, and Diversity of Phages. *Viruses* [Internet]. 2019 [cited
1102 2019 Oct 24];11. Available from: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6466130/>
- 1103 16. Peng S-Y, You R-I, Lai M-J, Lin N-T, Chen L-K, Chang K-C. Highly potent antimicrobial
1104 modified peptides derived from the *Acinetobacter baumannii* phage endolysin LysAB2. *Sci Rep.*
1105 2017;7:1–12.
- 1106 17. Holt A, Cahill J, Ramsey J, O’Leary C, Moreland R, Martin C, et al. Phage-encoded cationic
1107 antimicrobial peptide used for outer membrane disruption in lysis. *bioRxiv.* 2019;515445.
- 1108 18. Harada LK, Silva EC, Campos WF, Del Fiol FS, Vila M, Dąbrowska K, et al.
1109 Biotechnological applications of bacteriophages: State of the art. *Microbiological Research.*
1110 2018;212–213:38–58.
- 1111 19. Sharma RS, Karmakar S, Kumar P, Mishra V. Application of filamentous phages in
1112 environment: A tectonic shift in the science and practice of ecorestoration. *Ecology and*
1113 *Evolution.* 2019;9:2263–304.
- 1114 20. Jiang SC, Paul JH. Gene Transfer by Transduction in the Marine Environment. *APPL*
1115 *ENVIRON MICROBIOL.* 1998;64:8.
- 1116 21. Gregory AC, Solonenko SA, Ignacio-Espinoza JC, LaButti K, Copeland A, Sudek S, et al.
1117 Genomic differentiation among wild cyanophages despite widespread horizontal gene transfer.
1118 *BMC Genomics.* 2016;17:930.
- 1119 22. Sanjuán R, Nebot MR, Chirico N, Mansky LM, Belshaw R. Viral Mutation Rates. *J Virol.*
1120 2010;84:9733–48.
- 1121 23. Dutilh BE, Cassman N, McNair K, Sanchez SE, Silva GGZ, Boling L, et al. A highly
1122 abundant bacteriophage discovered in the unknown sequences of human faecal metagenomes.
1123 *Nature Communications.* 2014;5:4498.
- 1124 24. Devoto AE, Santini JM, Olm MR, Anantharaman K, Munk P, Tung J, et al. Megaphages
1125 infect *Prevotella* and variants are widespread in gut microbiomes. *Nature Microbiology.*
1126 2019;4:693–700.
- 1127 25. Al-Shayeb B, Sachdeva R, Chen L-X, Ward F, Munk P, Devoto A, et al. Clades of huge
1128 phage from across Earth’s ecosystems. *bioRxiv.* 2019;572362.
- 1129 26. Kauffman KM, Hussain FA, Yang J, Arevalo P, Brown JM, Chang WK, et al. A major
1130 lineage of non-tailed dsDNA viruses as unrecognized killers of marine bacteria. *Nature; London.*
1131 2018;554:118-122,122A-122T.
- 1132 27. Hopkins M, Kailasan S, Cohen A, Roux S, Tucker KP, Shevenell A, et al. Diversity of
1133 environmental single-stranded DNA phages revealed by PCR amplification of the partial major
1134 capsid protein. *ISME J.* 2014;8:2093–103.

- 1135 28. Krishnamurthy SR, Janowski AB, Zhao G, Barouch D, Wang D. Hyperexpansion of RNA
1136 Bacteriophage Diversity. *PLOS Biology*. 2016;14:e1002409.
- 1137 29. Waldbauer JR, Coleman ML, Rizzo AI, Campbell KL, Lotus J, Zhang L. Nitrogen sourcing
1138 during viral infection of marine cyanobacteria. *PNAS*. 2019;116:15590–5.
- 1139 30. Stent GS, Maaløe O. Radioactive phosphorus tracer studies on the reproduction of T4
1140 bacteriophage: II. Kinetics of phosphorus assimilation. *Biochimica et Biophysica Acta*.
1141 1953;10:55–69.
- 1142 31. Kozloff LM, Knowlton K, Putnam FW, Evans EA. Biochemical Studies of Virus
1143 Reproduction V. the Origin of Bacteriophage Nitrogen. *J Biol Chem*. 1951;188:101–16.
- 1144 32. Thompson LR, Zeng Q, Kelly L, Huang KH, Singer AU, Stubbe J, et al. Phage auxiliary
1145 metabolic genes and the redirection of cyanobacterial host carbon metabolism. *PNAS*.
1146 2011;108:E757–64.
- 1147 33. Breitbart M, Thompson L, Suttle C, Sullivan M. Exploring the Vast Diversity of Marine
1148 Viruses. *Oceanography*. 2007;20:135–9.
- 1149 34. Hurwitz BL, Hallam SJ, Sullivan MB. Metabolic reprogramming by viruses in the sunlit and
1150 dark ocean. *Genome Biology*. 2013;14:R123.
- 1151 35. Roux S, Hawley AK, Beltran MT, Scofield M, Schwientek P, Stepanauskas R, et al. Ecology
1152 and evolution of viruses infecting uncultivated SUP05 bacteria as revealed by single-cell- and
1153 meta-genomics. *eLife Sciences*. 2014;3:e03125.
- 1154 36. Bragg JG, Chisholm SW. Modeling the Fitness Consequences of a Cyanophage-Encoded
1155 Photosynthesis Gene. *PLOS ONE*. 2008;3:e3550.
- 1156 37. Mann NH, Cook A, Millard A, Bailey S, Clokie M. Bacterial photosynthesis genes in a virus.
1157 *Nature*. 2003;424:741.
- 1158 38. Anantharaman K, Duhaime MB, Breier JA, Wendt KA, Toner BM, Dick GJ. Sulfur
1159 Oxidation Genes in Diverse Deep-Sea Viruses. *Science*. 2014;344:757–60.
- 1160 39. Emerson JB, Roux S, Brum JR, Bolduc B, Woodcroft BJ, Jang HB, et al. Host-linked soil
1161 viral ecology along a permafrost thaw gradient. *Nature Microbiology*. 2018;3:870.
- 1162 40. Trubl G, Jang HB, Roux S, Emerson JB, Solonenko N, Vik DR, et al. Soil Viruses Are
1163 Underexplored Players in Ecosystem Carbon Processing. *mSystems*. 2018;3:e00076-18.
- 1164 41. Labonté JM, Swan BK, Poulos B, Luo H, Koren S, Hallam SJ, et al. Single-cell genomics-
1165 based analysis of virus–host interactions in marine surface bacterioplankton. *ISME J*.
1166 2015;9:2386–99.

- 1167 42. Trubl G, Solonenko N, Chittick L, Solonenko SA, Rich VI, Sullivan MB. Optimization of
1168 viral resuspension methods for carbon-rich soils along a permafrost thaw gradient. *PeerJ*.
1169 2016;4:e1999.
- 1170 43. Roux S, Enault F, Hurwitz BL, Sullivan MB. VirSorter: mining viral signal from microbial
1171 genomic data. *PeerJ*. 2015;3:e985.
- 1172 44. Wommack KE, Bhavsar J, Polson SW, Chen J, Dumas M, Srinivasiah S, et al. VIROME: a
1173 standard operating procedure for analysis of viral metagenome sequences. *Standards in Genomic*
1174 *Sciences*. 2012;6:427.
- 1175 45. Roux S, Faubladiere M, Mahul A, Paulhe N, Bernard A, Debroas D, et al. Metavir: a web
1176 server dedicated to virome analysis. *Bioinformatics*. 2011;27:3074–5.
- 1177 46. El-Gebali S, Mistry J, Bateman A, Eddy SR, Luciani A, Potter SC, et al. The Pfam protein
1178 families database in 2019. *Nucleic Acids Res*. 2019;47:D427–32.
- 1179 47. Ren J, Ahlgren NA, Lu YY, Fuhrman JA, Sun F. VirFinder: a novel k-mer based tool for
1180 identifying viral sequences from assembled metagenomic data. *Microbiome*. 2017;5:69.
- 1181 48. Fang Z, Tan J, Wu S, Li M, Xu C, Xie Z, et al. PPR-Meta: a tool for identifying phages and
1182 plasmids from metagenomic fragments using deep learning. *Gigascience* [Internet]. 2019 [cited
1183 2019 Aug 5];8. Available from:
1184 <https://academic.oup.com/gigascience/article/8/6/giz066/5521157>
- 1185 49. Ahlgren NA, Ren J, Lu YY, Fuhrman JA, Sun F. Alignment-free d_2 oligonucleotide
1186 frequency dissimilarity measure improves prediction of hosts from metagenomically-derived
1187 viral sequences. *Nucleic Acids Res*. 2017;45:39–53.
- 1188 50. Ponso AJ, Hurwitz BL. The Promises and Pitfalls of Machine Learning for Detecting
1189 Viruses in Aquatic Metagenomes. *Front Microbiol* [Internet]. 2019 [cited 2019 Oct 24];10.
1190 Available from: <https://www.frontiersin.org/articles/10.3389/fmicb.2019.00806/full>
- 1191 51. Amgarten D, Braga LPP, da Silva AM, Setubal JC. MARVEL, a Tool for Prediction of
1192 Bacteriophage Sequences in Metagenomic Bins. *Front Genet* [Internet]. 2018 [cited 2019 Aug
1193 5];9. Available from: <https://www.frontiersin.org/articles/10.3389/fgene.2018.00304/full>
- 1194 52. Zheng T, Li J, Ni Y, Kang K, Misiakou M-A, Imamovic L, et al. Mining, analyzing, and
1195 integrating viral signals from metagenomic data. *Microbiome*. 2019;7:42.
- 1196 53. Arndt D, Grant JR, Marcu A, Sajed T, Pon A, Liang Y, et al. PHASTER: a better, faster
1197 version of the PHAST phage search tool. *Nucleic Acids Res*. 2016;44:W16–21.
- 1198 54. Song W, Sun H-X, Zhang C, Cheng L, Peng Y, Deng Z, et al. Prophage Hunter: an
1199 integrative hunting tool for active prophages. *Nucleic Acids Res*. 2019;47:W74–80.

- 1200 55. Merchant N, Lyons E, Goff S, Vaughn M, Ware D, Micklos D, et al. The iPlant
1201 Collaborative: Cyberinfrastructure for Enabling Data to Discovery for the Life Sciences. *PLOS*
1202 *Biology*. 2016;14:e1002342.
- 1203 56. Kristensen DM, Waller AS, Yamada T, Bork P, Mushegian AR, Koonin EV. Orthologous
1204 Gene Clusters and Taxon Signature Genes for Viruses of Prokaryotes. *Journal of Bacteriology*.
1205 2013;195:941–50.
- 1206 57. Graziotin AL, Koonin EV, Kristensen DM. Prokaryotic Virus Orthologous Groups
1207 (pVOGs): a resource for comparative genomics and protein family annotation. *Nucleic Acids*
1208 *Res*. 2017;45:D491–8.
- 1209 58. Hendricks SP, Mathews CK. Regulation of T4 Phage Aerobic Ribonucleotide Reductase:
1210 SIMULTANEOUS ASSAY OF THE FOUR ACTIVITIES. *J Biol Chem*. 1997;272:2861–5.
- 1211 59. Kanehisa M, Goto S. KEGG: Kyoto Encyclopedia of Genes and Genomes. *Nucleic Acids*
1212 *Res*. 2000;28:27–30.
- 1213 60. Aramaki T, Blanc-Mathieu R, Endo H, Ohkubo K, Kanehisa M, Goto S, et al.
1214 KofamKOALA: KEGG ortholog assignment based on profile HMM and adaptive score
1215 threshold. *bioRxiv*. 2019;602110.
- 1216 61. Roux S, Adriaenssens EM, Dutilh BE, Koonin EV, Kropinski AM, Krupovic M, et al.
1217 Minimum Information about an Uncultivated Virus Genome (MIUViG). *Nature Biotechnology*.
1218 2019;37:29–37.
- 1219 62. Bowers RM, Kyrpides NC, Stepanauskas R, Harmon-Smith M, Doud D, Reddy TBK, et al.
1220 Minimum information about a single amplified genome (MISAG) and a metagenome-assembled
1221 genome (MIMAG) of bacteria and archaea. *Nature Biotechnology*. 2017;35:725–31.
- 1222 63. Paez-Espino D, Chen I-MA, Palaniappan K, Ratner A, Chu K, Szeto E, et al. IMG/VR: a
1223 database of cultured and uncultured DNA Viruses and retroviruses. *Nucleic Acids Res*.
1224 2017;45:D457–65.
- 1225 64. Paez-Espino D, Roux S, Chen I-MA, Palaniappan K, Ratner A, Chu K, et al. IMG/VR v.2.0:
1226 an integrated data management and analysis system for cultivated and environmental viral
1227 genomes. *Nucleic Acids Res*. 2019;47:D678–86.
- 1228 65. Gregory AC, Zablocki O, Howell A, Bolduc B, Sullivan MB. The human gut virome
1229 database. *bioRxiv*. 2019;655910.
- 1230 66. Payet JP, Suttle CA. To kill or not to kill: The balance between lytic and lysogenic viral
1231 infection is driven by trophic status. *Limnology and Oceanography*. 2013;58:465–74.
- 1232 67. Gregory AC, Zayed AA, Conceição-Neto N, Temperton B, Bolduc B, Alberti A, et al.
1233 Marine DNA Viral Macro- and Microdiversity from Pole to Pole. *Cell [Internet]*. 2019 [cited
1234 2019 Apr 30]; Available from:
1235 <http://www.sciencedirect.com/science/article/pii/S0092867419303411>

- 1236 68. Pasolli E, Asnicar F, Manara S, Zolfo M, Karcher N, Armanini F, et al. Extensive
1237 Unexplored Human Microbiome Diversity Revealed by Over 150,000 Genomes from
1238 Metagenomes Spanning Age, Geography, and Lifestyle. *Cell*. 2019;176:649-662.e20.
- 1239 69. He Q, Gao Y, Jie Z, Yu X, Laursen JM, Xiao L, et al. Two distinct metacommunities
1240 characterize the gut microbiota in Crohn's disease patients. *Gigascience*. 2017;6:1-11.
- 1241 70. Morgan XC, Tickle TL, Sokol H, Gevers D, Devaney KL, Ward DV, et al. Dysfunction of
1242 the intestinal microbiome in inflammatory bowel disease and treatment. *Genome Biology*.
1243 2012;13:R79.
- 1244 71. Strauss J, Kaplan GG, Beck PL, Rioux K, Panaccione R, Devinney R, et al. Invasive
1245 potential of gut mucosa-derived *Fusobacterium nucleatum* positively correlates with IBD status
1246 of the host. *Inflamm Bowel Dis*. 2011;17:1971-8.
- 1247 72. Shreiner AB, Kao JY, Young VB. The gut microbiome in health and in disease. *Curr Opin*
1248 *Gastroenterol*. 2015;31:69-75.
- 1249 73. Clemente JC, Ursell LK, Parfrey LW, Knight R. The Impact of the Gut Microbiota on
1250 Human Health: An Integrative View. *Cell*. 2012;148:1258-70.
- 1251 74. Minot SS, Willis AD. Clustering co-abundant genes identifies components of the gut
1252 microbiome that are reproducibly associated with colorectal cancer and inflammatory bowel
1253 disease. *Microbiome*. 2019;7:110.
- 1254 75. Nishio M, Okada N, Miki T, Haneda T, Danbara H. Identification of the outer-membrane
1255 protein PagC required for the serum resistance phenotype in *Salmonella enterica* serovar
1256 *Choleraesuis*. *Microbiology (Reading, Engl)*. 2005;151:863-73.
- 1257 76. Rangunathan PT, Vanderpool CK. Cryptic-Prophage-Encoded Small Protein DicB Protects
1258 *Escherichia coli* from Phage Infection by Inhibiting Inner Membrane Receptor Proteins. *Journal*
1259 *of Bacteriology* [Internet]. 2019 [cited 2019 Nov 11];201. Available from:
1260 <https://j.b.asm.org/content/201/23/e00475-19>
- 1261 77. Rasko DA, Rosovitz MJ, Myers GSA, Mongodin EF, Fricke WF, Gajer P, et al. The
1262 Pangenome Structure of *Escherichia coli*: Comparative Genomic Analysis of *E. coli* Commensal
1263 and Pathogenic Isolates. *Journal of Bacteriology*. 2008;190:6881-93.
- 1264 78. Brown CT, Olm MR, Thomas BC, Banfield JF. Measurement of bacterial replication rates in
1265 microbial communities. *Nature Biotechnology*. 2016;34:1256-63.
- 1266 79. Ijaz UZ, Quince C, Hanske L, Loman N, Calus ST, Bertz M, et al. The distinct features of
1267 microbial "dysbiosis" of Crohn's disease do not occur to the same extent in their unaffected,
1268 genetically-linked kindred. *PLoS ONE*. 2017;12:e0172605.
- 1269 80. Gevers D, Kugathasan S, Denson LA, Vázquez-Baeza Y, Van Treuren W, Ren B, et al. The
1270 treatment-naive microbiome in new-onset Crohn's disease. *Cell Host Microbe*. 2014;15:382-92.

- 1271 81. Mizuno CM, Guyomar C, Roux S, Lavigne R, Rodriguez-Valera F, Sullivan MB, et al.
1272 Numerous cultivated and uncultivated viruses encode ribosomal proteins. *Nature*
1273 *Communications*. 2019;10:752.
- 1274 82. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment search tool. *J*
1275 *Mol Biol*. 1990;215:403–10.
- 1276 83. Eddy SR. Profile hidden Markov models. *Bioinformatics*. 1998;14:755–63.
- 1277 84. Hyatt D, Chen G-L, LoCascio PF, Land ML, Larimer FW, Hauser LJ. Prodigal: prokaryotic
1278 gene recognition and translation initiation site identification. *BMC Bioinformatics*. 2010;11:119.
- 1279 85. Fu L, Niu B, Zhu Z, Wu S, Li W. CD-HIT: accelerated for clustering the next-generation
1280 sequencing data. *Bioinformatics*. 2012;28:3150–2.
- 1281 86. Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, et al. Scikit-learn:
1282 Machine Learning in Python. *Journal of Machine Learning Research*. 2011;12:2825–30.
- 1283 87. Hunter JD. Matplotlib: A 2D graphics environment. *Computing In Science & Engineering*.
1284 2007;9:90–5.
- 1285 88. Feng Q, Liang S, Jia H, Stadlmayr A, Tang L, Lan Z, et al. Gut microbiome development
1286 along the colorectal adenoma–carcinoma sequence. *Nature Communications*. 2015;6:1–13.
- 1287 89. Joshi N, Fass J. Sickle: A sliding-window, adaptive, quality-based trimming tool for FastQ
1288 files [Internet]. 2011. Available from: <https://github.com/najoshi/sickle>
- 1289 90. Nurk S, Meleshko D, Korobeynikov A, Pevzner PA. metaSPAdes: a new versatile
1290 metagenomic assembler. *Genome Res*. 2017;27:824–34.
- 1291 91. Zhou Z, Tran PQ, Kieft K, Anantharaman K. Genome diversification in globally distributed
1292 novel marine Proteobacteria is linked to environmental adaptation. *bioRxiv*. 2019;814418.
- 1293 92. Ondov BD, Treangen TJ, Melsted P, Mallonee AB, Bergman NH, Koren S, et al. Mash: fast
1294 genome and metagenome distance estimation using MinHash. *Genome Biology*. 2016;17:132.
- 1295 93. Delcher AL. Fast algorithms for large-scale genome alignment and comparison. *Nucleic*
1296 *Acids Research*. 2002;30:2478–83.
- 1297 94. Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. *Nat Methods*.
1298 2012;9:357–9.
- 1299 95. Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for RNA-
1300 seq data with DESeq2. *Genome Biology*. 2014;15:550.
- 1301 96. Sullivan MJ, Petty NK, Beatson SA. Easyfig: a genome comparison visualizer.
1302 *Bioinformatics*. 2011;27:1009–10.

- 1303 97. Jang HB, Bolduc B, Zablocki O, Kuhn J, Roux S, Adriaenssens E, et al. Gene sharing
1304 networks to automate genome-based prokaryotic viral taxonomy. *bioRxiv*. 2019;533240.
- 1305 98. Shannon P. Cytoscape: A Software Environment for Integrated Models of Biomolecular
1306 Interaction Networks. *Genome Research*. 2003;13:2498–504.
- 1307

Supplementary Files

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

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