

Simulation Study and Comparative Evaluation of Viral Contiguous Sequence Identification Tools

Cody Glickman (✉ cody.glickman@cuanschutz.edu)

University of Colorado Anschutz Medical Campus

Jo Hendrix

University of Colorado Anschutz Medical Campus

Michael Strong

National Jewish Health

Research Article

Keywords: virus, bacteriophage, prophage, metagenomics, tool comparison

Posted Date: March 13th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-287089/v1>

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

[Read Full License](#)

RESEARCH

Simulation Study and Comparative Evaluation of Viral Contiguous Sequence Identification Tools

Cody Glickman^{1,2*}, Jo Hendrix^{1,2} and Michael Strong^{1,2}

Abstract

Background:

Viruses, including bacteriophage, are important components of environmental and human associated microbial communities. Viruses can act as extracellular reservoirs of bacterial genes, can mediate microbiome dynamics, and can influence the virulence of clinical pathogens. It is essential, therefore, to have robust sequence analysis methods in place to detect and annotate viral elements within microbial communities. Various targeted metagenomic analysis techniques detect viral sequences, but these methods often exclude large and genome integrated viruses. In this study, we evaluate and compare the ability of nine state-of-the-art bioinformatic tools, including Vibrant, VirSorter, VirSorter2, VirFinder, DeepVirFinder, MetaPhinder, JGI Earth Virome Pipeline, Kraken 2, and VirBrant, to identify viral contiguous sequences (contigs) across simulated metagenomes with different read distributions, taxonomic compositions, and complexities.

Results:

Of the tools tested in this study, VirSorter achieved the best F1 score while Vibrant had the highest average F1 score at predicting integrated prophages. Though less balanced in its precision and recall, Kraken2 had the highest average precision by a substantial margin.

We introduced the machine learning tool, VirBrant, which demonstrated an improvement in average F1 score over tools such as MetaPhinder. The tool utilizes machine learning with both protein compositional and nucleotide features. The addition of nucleotide features improves the precision and recall compared to the protein compositional features alone.

Viral identification by all tools was not impacted by underlying read distribution but did improve with contig length. Tool performance was inversely related to taxonomic complexity and varied by the phage host. *Rhizobium* and *Enterococcus* phage were identified consistently by the tools; whereas, *Neisseria* phage were commonly missed in this study.

Conclusion:

This study benchmarked the performance of nine state-of-the-art bioinformatic tools to identify viral contigs across different simulation conditions. This study explored the ability of the tools to identify integrated prophage elements traditionally excluded from targeted sequencing approaches. Our comprehensive analysis of viral identification tools to assess their performance in a variety of situations provides valuable insights to viral researchers looking to mine viral elements from publicly available metagenomic data.

Keywords: virus; bacteriophage; prophage; metagenomics; tool comparison

*Correspondence:

cody.glickman@cuanschutz.edu

¹Center for Genes, Environment, and Health, National Jewish Health, 1400 Jackson Street, 80206 Denver, Colorado, USA
Full list of author information is available at the end of the article

1 Background

2 Viruses are the most abundant biological entities on Earth [1]. However, the
3 collective knowledge of environmental viral sequences, including bacteriophages, re-
4 mains underrepresented relative to the amount of genetic information for eukaryotic
5 viruses and bacteria. Bacteriophages are viruses that infect bacteria and are com-
6 monly referred to as phages. In addition to bacteria, phages are capable of infecting
7 archaea. Phages are obligate parasites that play an important role in the genomic
8 composition and evolution of their bacterial hosts. Phages directly contribute to

9 bacterial infections in humans by acting as a genetic reservoir for virulent genes
10 in bacteria such as *Escherichia coli*, *Salmonella enterica*, *Pseudomonas aeruginosa*,
11 *Vibrio cholerae*, *Corynebacterium diphtheriae*, and *Streptococcus pyogenes* [2, 3].

12 In addition, some phages utilize Ig-like domains to attach to mucosal layers
13 in humans to lie in wait for bacterial prey. This bacteriophage adherence to mucus
14 (BAM) model suggests that phages may act as a non-host derived innate immunity
15 system to modulate the bacterial microbiome [4]. A longitudinal study of the human
16 virome revealed composition conservation that mimicked the stability of healthy
17 bacterial microbiomes [5, 6]. Dysbiosis in the virome has been observed in disease
18 states such as inflammatory bowel disease (IBD), Crohn’s disease, and asthma [7–9].

19 The study of viruses has traditionally relied on the ability to cultivate vi-
20 ral particles from a cultured host; however, many bacteria cannot be cultured in
21 a laboratory setting [10]. The limited number of culturable hosts, in combination
22 with the additional complexities of viral isolation limit the study of viruses. The
23 advancements in next generation sequencing technologies created an opportunity
24 to study viruses with culture independent methods. However, because viruses do
25 not share a common universal marker gene, like the bacterial small subunit RNA,
26 sequencing techniques such as metagenomics are a necessity [11]. Metagenomics is a
27 non-targeted sequencing approach to elucidate the totality of genetic material within
28 a sample, either DNA or RNA. However, in part due to small genomes, viruses are
29 traditionally underrepresented in metagenomic studies from a read abundance per-
30 spective. It is common for viral reads to comprise less than 5% of metagenomic
31 sequences [12]. A way to enrich viral reads in metagenomic studies is to filter or
32 directly select viral like particles (VLPs). However, these techniques tend to remove
33 large viruses and viruses integrated into bacterial genomes called prophages, before
34 sequencing. Therefore, the ability to identify viral elements directly from metage-
35 nomic sequencing studies is also important for understanding the composition of
36 the virome. The advent of computational tools dedicated to the identification of
37 viral sequences in metagenomics has improved our ability to identify known, novel,
38 and integrated viruses.

39 MetaPhinder is an approach that uses BLASTn and average nucleotide iden-
40 tity thresholds to classify viral contigs in metagenomics [13]. Methods that use
41 sequence similarity suffers worsening performance with smaller contig lengths. Do-
42 main recognition is utilized by more tools to counter the limitations of contig length
43 on traditional sequence homology approaches, but these tools are often reliant on
44 specialized viral domains like those from pVOGs (prokaryotic virus orthologous
45 groups) [14]. Unlike prophage identification methods that use viral domain enrich-
46 ment or presence/absence to calculate a score, a new method, called Vibrant, uses
47 domain abundances in a neural network framework to classify contigs having more
48 than 4 proteins [15]. VirSorter2 follows a similar methodology, using domain per-
49 centages, protein compositionality features, and key homology genes in a tree-based
50 machine learning framework to classify viral reads [16].

51 Homology of viral protein domains is limited to known viruses, which are
52 thought to represent only a small slice of the vast viral dark matter [17]. Another
53 homology approach sought to expand known viral hidden Markov models (HMMs)
54 through a semi-supervised expansion of existing viral protein families. Paez-Espino

55 et al. (JGI Earth Virome Pipeline) collected viral coding regions from NCBI servers
56 and known viral metagenomic contigs; then clustered those peptides into protein
57 families to create new viral HMMs [18]. This initial set was used as bait to identify
58 potential viral contigs in thousands of metagenomic data sets. Predicted proteins
59 from these captured viral contigs were added to the original set of peptides and re-
60 clustered to create thousands of new viral protein families and HMMs. Even with
61 the expansion of viral families, both VirSorter and the JGI Earth Virome Pipeline
62 are at least partially reliant on domain homology. A reference-free viral identifica-
63 tion tool was developed using machine learning to address limitations of homology
64 searching. VirFinder is a logistic regression classification using nucleotide sequence
65 8-mers as features [19]. The authors of VirFinder expanded the concept of using
66 machine learning to identify viral contiguous sequences with DeepVirFinder, a con-
67 volutional neural network that takes raw sequences as inputs and learns features
68 that are useful for viral contig prediction [20]. VirFinder relies solely on sequence-
69 based features, which is analogous to another k-mer approach, Kraken2. Kraken2
70 uses discriminatory 35-mers to uniquely identify sequences to the species and even
71 subspecies level [21]. In order to use Kraken2 in a viral identification context, we
72 created the tool, VirKraken, that parses the Kraken2 classification output to iden-
73 tify viral contigs in metagenomic reads. VirKraken is available on PyPI and at
74 <https://github.com/Strong-Lab/VirKraken>. VirKraken references the Kraken2 as-
75 signed taxonomy identification number against an edited NCBI Taxonomy database
76 to assign kingdom and to filter sequences when requested [22].

77 Another approach to identify viral elements in metagenomics involves nega-
78 tion of known bacteria contigs. VirMine uses a homology search against a bacterial
79 protein database; if hits of bacterial genes outnumber the number of unknown hits
80 the contig is removed, thus leaving viral contigs [23]. All previously described tools
81 identify viral elements from assembled sequences. MARVEL is a machine learning
82 method that classifies binned reads as viral clusters using a random forest approach
83 with three features (gene density, strand shifts, and fraction of homology hits to
84 a viral protein database) [24]. We developed a machine learning model called Vir-
85 Brant that uses both protein compositional features such as gene density and strand
86 shift frequency, in addition to sequence-based features to classify viral contigs. The
87 addition of protein compositional features is hypothesized to offset the dip in perfor-
88 mance of sequence-based machine learning models compared to homology methods
89 on longer contigs [19].

90 Many approaches exist to identify viral elements in metagenomics. However, a
91 systemic evaluation among many of these tools has not been performed. This study
92 is meant to provide information and guidance to researchers regarding when to use
93 a specific viral identification tool to further study viral elements or to remove them
94 for downstream analyses. The characterization of more viral elements in the public
95 domain could lead to the discovery of novel viruses [25] and provide insight into the
96 functional potential residing in an extracellular genetic reservoir [2].

97 **Methods**

98 VirBrant, a hybrid protein composition and nucleotide feature set for viral classification

99 To build VirBrant, 1,849 complete phage and 2,327 complete archaeal/bacterial
100 genomes were compiled from RefSeq (Accessed on January 8th, 2020). Prophage

101 elements in the archaeal and bacterial genomes were identified using VirSorter [22].
102 Category 4 prophages were selected, and the predicted nucleotide sequences were
103 added to the complete viral genomes. Custom scripts were used to identify and
104 remove the predicted prophage sequences from the host genome. The total number
105 of prophages predicted was 730 in 339 bacterial genomes (14.57% of genomes con-
106 tained at least 1 prophage) culminating in an average prophage per genome ratio
107 of 0.314.

108 After removing integrated prophages, the complete genomes were fragmented
109 into k-mers of 4 sizes using an n-step kmerization method. The n-step method
110 removes contig end-overlap and ensures that the maximum number of k-mers is
111 the length of the base sequence over the length k. The complete genomes were
112 fragmented into sizes of 1KB, 3KB, 5KB, and 10KB sequences. Due to the size
113 of bacterial and archaeal genomes relative to phage genomes, the fragments from
114 non-phage sampling were down sampled to evenly distribute the classes. The four
115 different fragment lengths were used to train four separate models.

116 The n-step fragments were subjected to a sliding window kmerization of size
117 8 using a k-mer counting program written in C [26]. A sliding window kmerization
118 calculates k-mer abundance with significant overlap and the maximum number of
119 k-mers is the length of the base sequences minus 1. The program stores all 8-mer
120 values (65,536 possible 8-mers) in a hash table. In real world metagenomic sampling,
121 the directionality of a sequence fragment may be ambiguous. Therefore, similar to
122 VirFinder [19], we developed custom scripts to sum complement, reverse, and reverse
123 complement sequences thus reducing our feature space from 65,536 possible k-mers
124 to 16,384 possible k-mers. The nucleotide feature space is further reduced to 888
125 k-mers using Gini importance or total decrease in node impurity above 0.001, which
126 is a weighted probability of reaching a feature averaged over all trees in a random
127 forest [27].

128 Protein composition feature set creation

129 The use of protein compositional features is built into tools such as MARVEL
130 and VirSorter [24, 28]. MARVEL and VirSorter both utilize gene density as a marker
131 of viral elements. In this study, four protein features associated with viral genomes
132 were included as part of the feature set in VirBrant; gene density, operon length,
133 average peptide length, and percentage of overlapping peptides. Due to the physical
134 restraints of some viral capsids, viral genomes are commonly tightly packed and
135 translate shorter proteins than bacterial genomes [29]. In addition, viral genomes
136 often have overlapped genes for different life cycles and have long stretches of genes
137 located on the same strand [30]. Custom scripts were used to calculate the four
138 protein characteristics from the output of the Prodigal gene prediction software
139 [31]. **Figure 1** shows the observed distribution of protein features in the training
140 data of the 10KB model.

141 **[Figure 1]**

142 Model and hyperparameter selection

143 After combining the complementary nucleotide features and the protein com-
144 positional features, the total feature space of VirBrant was 892 features. During

145 training, the performance of a random forest, multi-layer perceptron, and an ad-
146 ditive boosting model were compared using 5-fold cross validation [32]. At every
147 fragment size, the additive boosting model performed the best. We selected XG-
148 Boost (version 0.81) and performed a RandomSearchCV (version 0.20.1) analysis
149 to determine hyperparameters [32, 33]. The pre-trained models were added to the
150 tool repository for use classifying metagenomic fasta sequences. VirBrant generates
151 outputs as a header file containing the header sequences of viral elements and a
152 fasta file containing the nucleotide information of the predicted viral elements.

153 Building simulated Illumina metagenomes

154 To build a simulated test set, all complete genomes were downloaded from
155 NCBI RefSeq (accessed on 12/15/2020). The genomes deposited since May 1st,
156 2020, were selected to test the viral contig identification tools because many of the
157 tools were trained or relied on databases last updated prior to this date. Bacterial
158 hosts of phages were collected using a dataset from Virus-Host DB [34] (Accessed
159 on December 17, 2020). Phage were assigned bacteria genera values by their host
160 organism. Using information from the Earth Microbiome Project (EMP) and from
161 Qiita, the recently submitted genomes were further filtered by 53 genera commonly
162 found in soils (37 genera) and in clinical samples (26 genera) with 8 genera in both
163 niches [35, 36]. This resulted in 297 unique bacterial genomes being used for the
164 simulations with 82 genomes found in both clinical sampling (160 genomes) and soil
165 sampling (219 genomes). The reliance on recently submitted genomes to produce
166 the testing set did not produce traditional bacterial distributions seen in clinical
167 and soil microbiomes. For example, while the genera *Bacteroides* are commonly
168 present in the clinical microbe samples, the amount in this study does not repre-
169 sent a substantial portion of the community as seen in other clinical microbiome
170 studies [37]. The distribution of bacterial genera was used as a confounder for viral
171 classification in this study. The goal of this study was to observe the performance
172 of phage identification in the presence of genetically similar bacteria.

173 Phage genomes were also filtered by their host bacterial genera and randomly
174 down sampled to match the number of bacterial genomes in the simulations. While
175 phages are thought to outnumber bacteria ten to one in the environment [38], we
176 matched the complexity of phage and bacteria in our simulations across taxonomic
177 levels due to limitations in the number of available phage genomes for the full
178 datasets. In order to test the impact of taxonomic complexity on viral identification
179 tool performance, we subsampled phage and bacterial genomes into medium (50
180 bacterial genomes and 50 phage genomes) and low (10 bacterial genomes and 10
181 phage genomes) complexity subsets. **Supplementary Table 1** (clinical) and **Sup-**
182 **plementary Table 2** (soil) detail the taxonomic abundance of the top 6 genera
183 and phage host genera in the testing set across taxonomic complexity levels. While
184 both lower complexities draw from the full distribution of genomes, there is no over-
185 lap in the selected genomes between the medium and low taxonomic levels. This
186 was accomplished through setting a random seed in the subsampling procedure and
187 using set operations to confirm no overlap of genomes.

188 Simulated metagenomes were created using InSilicoSeq (version 1.2.0). InSili-
189 coSeq and another popular metagenomic simulator, CAMISIM use a lognormal read

190 distribution by default, however, four additional read distributions are provided as
191 a part of the InSilicoSeq software suite: uniform, exponential, zero inflated lognor-
192 mal, and halfnormal [39, 40]. Due to the enormous diversity of naturally occurring
193 communities, read distribution profiles are likely to fluctuate. To understand the
194 impact of read distribution and taxonomic complexities on the performance of viral
195 identification, we created 30 MiSeq simulations with 12 million 2x300 reads. The
196 30 simulations were composed of two environmental conditions (clinical and soil
197 microbes) with five read distributions across three taxonomic levels (full, medium,
198 low). Bacterial reads represented 93.75% of the total composition in each simulation
199 and phages represented 6.25%. Prior studies suggest phages commonly represent less
200 than 5% of metagenomic sequencing reads [12] due to genomes that are orders of
201 magnitude smaller than prokaryotic genomes. Our decision to exceed the 5% of
202 viral reads in metagenomics was driven by the need to identify an expanded set
203 of phages from taxonomically diverse testing sets. After assembly and filtration of
204 contigs less than 1KB in length, phages comprise an average of 1.54% of total contig
205 abundance.

206 After simulating, the reads were perfectly binned by sequence origin to limit
207 the creation of chimeric contigs. Chimeric contigs are assembly errors when reads
208 from different organisms are assembled together resulting in a shorter fragmented
209 assembly or taxonomic misclassification downstream. The decision to bin prior to
210 assembly was to allow for genera labeled contigs in order to explore false positive
211 and recall rates of bacteria and phage, respectively. The perfect bins were assembled
212 using metaSpades (version 3.11.1) and only contigs of length 1KB or greater were
213 retained [41]. The relative abundance of bacteria genera in the simulations are shown
214 in **Figure 2**.

215 [Figure 2]

216 Integrated prophage identification

217 Integrated prophage elements were identified in complete bacterial genomes us-
218 ing VirSorter prior to read simulations [28]. Integrated prophages were selected for
219 downstream processing if assigned as category 4, the highest confidence category for
220 prophages within VirSorter [28]. A nucleotide BLAST database was created with the
221 identified prophage elements. After read simulation and assembly, bacterial contigs
222 were identified as prophages using a BLASTn search against the prophage database
223 with a bitscore greater than 1000 and a percent identity greater than 95%. **Sup-**
224 **plementary Figure 1** shows the genera distribution of the identified prophage
225 elements separated by read distribution and sampling site.

226 Tools used in simulation study

227 [Table 1]

228 The tools used in the study shown in Table 1 were tested on their performance
229 to identify viral elements from assembled contigs in the simulations. The tools
230 used in this study to identify viral contigs were Vibrant (Version 1.2.0), VirSorter,

231 VirSorter2, VirFinder, DeepVirFinder, MetaPhinder, JGI Earth Virome Pipeline,
232 Kraken 2, and VirBrant [15, 16, 19, 28].

233 Any VirSorter predictions that were classified to the lowest confidence category
234 were removed via evidence by the tool developers [28]. VirFinder and DeepVirFinder
235 assign a probability value and any contigs that had a value less than 0.01 were clas-
236 sified as viral. A diamond blast database was created with the viral proteins from
237 the JGI Earth Virome Pipeline [18, 42]. Proteins from the simulation contigs were
238 predicted using Prodigal and searched for viral homology using diamond BLAST
239 against proteins from the JGI Earth Virome Pipeline with matches retained that
240 had a bit score greater than 100 and an e-value less than 1e-05 [31]. Contigs with
241 more than one hit were classified as viral. MetaPhinder, VirBrant, VirBrant Pro-
242 teins, and Vibrant were run with default parameters [13, 15]. VirSorter2 uses the
243 include groups flag to capture both double-stranded DNA phage and single stranded
244 DNA viruses specific to phage as described by the authors [16]. Kraken 2 was run
245 with default parameters using the minikraken database from March 2020 [21]. The
246 resulting Kraken 2 report was parsed for viral reads using VirKraken (Version 0.0.5).

247 Tool performance scoring

248 The structure of the simulation allowed for each contig to possess a true origin
249 label. These labels were used to identify the performance of the tools to identify viral
250 elements in the simulations. The performance was measured by precision, recall, and
251 F1 score. Prophages were considered viral in this study and an additional analysis
252 of tool performance on prophage identification was performed. The performance
253 measures were used in a simulation performance ranking system to determine the
254 best performing tool across different scenarios. The performance of each tool was
255 ranked within each condition with 1 representing the best performing tool. The
256 highest-ranking value (worst performing tool) changes as some tools were unable to
257 properly calculate a score. This occurs when a tool did not predict any viral element
258 in a simulation.

259 In addition to overall performance, tool performance is evaluated at four dis-
260 cretized contig lengths: 1KB-2.5KB, 2.5KB-5KB, 5KB-10KB, 10KB+. The recall of
261 the tools to identify viral elements by genera was used to determine any systematic
262 biases for or against specific viral groups. Visualizations of scoring metrics were per-
263 formed in Python using a combination of Matplotlib (version 2.2.3) and Seaborn
264 (version 0.9.0) plotting software [43, 44]. Kruskal-Wallis nonparametric testing
265 was performed to determine if the scoring values arose from the same distributions.

266 Results

267 Overall tool performance

268 The F1 performance across different read simulation conditions was not sig-
269 nificantly different ($H = 4.02$, $p = 0.404$, Kruskal-Wallis). The F1 performance
270 was significantly different by taxonomic complexity with tool performance in lower
271 complexity simulations enriched relative to both medium and full complexity sim-
272 ulations ($H = 47.65$, $p = 4.50e-11$, Kruskal-Wallis). The F1 performance, as well as
273 precision and recall, of longer contigs specially the 10KB+ bin was enriched rela-
274 tive to other contig length bins ($H = 275.7$, $p = 1.82e-59$, Kruskal-Wallis). **Table**

275 **2** contains the mean performance of the tools and the average ranking across the
276 30 simulations. The F1 performance of the tools in the simulation discretized by
277 taxonomic complexity is shown in **Figure 3**.

278 [Table 2]

279 [Figure 3]

280 Kraken2 led both average precision and precision rank. In this study, Earth
281 Virome led in recall and recall rank. The tool with the highest average F1 score and
282 best F1 rank was VirSorter. VirSorter was also the tool used to perform prophage
283 identification. This may provide VirSorter with an advantage over other tools in
284 prophage identification.

285 Prophage identification performance

286 The prophage performance of the low complexity simulations are removed due
287 to the presence of only a single prophage contig in all 10 simulations. The F1
288 performance of the tools to identify prophage in 20 medium and high complexity
289 simulations is shown in **Table 3**.

290 [Table 3]

291 Tool performance by contig length

292 As the length of the contigs increase, the performance of the tools improved.
293 **Figure 4** demonstrates the F1 performance of each tool within defined contig length
294 bins. If the F1 score of a tool was 0, the record was removed as some simulations
295 lack shorter contiguous sequences.

296 [Figure 4]

297 Viral recall by host genera

298 Recall scores of viral elements from the medium and full distributions were
299 calculated across 30 host genera. Recall was only retained if greater than 0 to
300 prevent the absence of a phage host genera by niche. **Figure 5** shows the recall of
301 viral contigs by host genera across all tools. The viral host genera with the best
302 recall was *Xanthomonas*, however, phage with *Xanthomonas* as a host were not well
303 represented in the data set. Phage known to infect *Enterococcus* achieved an average
304 recall over 0.83 across all tools. DeepVirFinder performed the best at identifying
305 phage known to infect *Enterococcus* with an average recall rate of 0.97. *Neisseria*
306 phage had the lowest average recall performance across all tools (0.23), with only
307 7 tools correctly predicting at least one *Neisseria* phage contig. The Earth Virome
308 performed the best at identifying this elusive phage (0.68) and the next best tool
309 was MetaPhinder with a recall rate of 0.24.

310 [Figure 5]

311 False positive genera

312 In addition to the recall rate of viral elements by host genera, the percent
313 of genera associated with bacterial false positives was calculated for each tool in
314 medium and full complexity simulations. Bacterial genera that represent more than
315 one third of false positives of a tool in a simulation were retained. Eleven genera were
316 represented with *Streptomyces* present in 9 of 10 tools. Additionally, *Citrobacter* and
317 *Pseudomonas* were major false positive genera in more than 5 tools. **Supplemental**
318 **Figure 2** shows the genera of false positives that represent more than 33% by tool.

319 Discussion

320 This study benchmarked and evaluated the ability of nine viral classification
321 tools to identify viral and prophage elements within shotgun metagenomics. The
322 study consisted of 30 Illumina MiSeq simulations across two communities, five read
323 abundance distributions, and three taxonomic levels. The performance of the tools
324 was consistent across read distributions ($H = 4.02$, $p = 0.404$, Kruskal-Wallis),
325 whereas, the average performance increased with a reduction in taxonomic com-
326 plexity ($H = 47.65$, $p = 4.50e-11$, Kruskal-Wallis). Lower taxonomic complexity
327 was associated with longer contig lengths in the assemblies and longer contigs were
328 associated with improved overall performance.

329 The differences between performance scores suggests the selection of a tool
330 may depend upon the desired application. VirSorter scored the highest average F1
331 score and had the best F1 ranking across all the simulations. Kraken2 may be the
332 ideal tool when minimizing the number of false positives. The Earth Virome pipeline
333 had the best recall; however, the application of this tool is not meant for traditional
334 viral identification due to the large false positive rate. The Earth Virome protein
335 set was derived from an iterative viral protein domain search and may include many
336 unknown proteins that may not truly be derived from viral sources [18]. Even so,
337 the broad homology search space still failed to capture all viral derived contigs
338 demonstrating the difficulty of viral identification within metagenomic samples.

339 Prophage identification in metagenomics is a difficult problem as many in-
340 tegrated viral elements are degraded in bacterial hosts to drive evolution [45]. As
341 such, remnants of prophage particles are scattered across bacterial genomes and vi-
342 ral genes can be mistakenly attributed as bacterial in origin. Many tools to identify
343 prophages in whole genome experiments fail to generalize to metagenomics due to
344 fragmentation that breaks down traditional viral enrichment measurements. The
345 decision to select the highest confidence prophage predictions using VirSorter from
346 the complete genomes prior to simulation may have provided VirSorter with an
347 added performance boost. Vibrant had the highest average F1 score and best F1
348 ranking at identifying prophages across all 20 simulations. Kraken 2 had the highest
349 average precision and VirSorter had the best precision ranking. The Earth Virome
350 proteins excelled at recall; however, the next best tools were VirFinder and Deep-
351 VirFinder. VirFinder and DeepVirFinder like many other tools that perform well
352 with prophage recall have a high false positive rate.

353 The performance of all tools would increase with an additional step of remov-
354 ing known bacterial contigs. One approach is to search for genes unique to bacteria
355 and archaea, the 16S rRNA. 16S rRNA profiles from RFAM can be applied to the

356 RNA domain search tool, Infernal, to remove contigs with known bacterial genes
357 [46, 47]. This approach may impact the recovery of prophage contigs if the integra-
358 tion site of the virus was near a 16S rRNA.

359 Viral identification tools performed well at identifying phages known to infect
360 genera such as *Enterococcus* (0.83), *Mycobacterium* (0.77), and *Salmonella* (0.81).
361 The performance of the tools to identify phages that infect genera such as *Neisseria*
362 (0.23), *Brevibacterium* (0.30), and *Mesorhizobium* (0.33) dropped substantially. De-
363 tecting the presence of *Neisseria* phage may be important for a diagnostic of invasive
364 meningococcal disease as prophage-like elements are commonly found throughout
365 the *Neisseria* genera [48].

366 The performance of VirBrant including the nucleotide features showed im-
367 provement over the protein compositional features alone. The precision of VirBrant
368 dramatically improved with contigs over 10KB, however, smaller bins were plagued
369 with many false positives. Integrated prophages added to the viral class in the
370 training data represented 28.3% of the total viral genomes. Prophages are com-
371 monly degraded in bacterial hosts to drive evolution [45], therefore degraded viral
372 elements in bacterial contigs with similar nucleotide structures as the complete
373 prophages may be misclassified. In addition, the use of k-mer profiles for smaller
374 contig classification created sparse data sets, which may have led to overfitting.

375 The performance of the tools presented needs to be weighted with the com-
376 putational cost to run each tool. This study was performed on a shared high per-
377 formance computing cluster and individual tool performance and memory require-
378 ments were not captured on an isolated node. However, the mechanism of viral
379 identification can be used infer the relative time and memory consumption of the
380 tools. The fastest tool in this study was Kraken2, which uses discriminatory k-mers
381 to compare against a pre-computed hash table. The amount of memory needed to
382 build the full hash table may be a drawback against using Kraken2 on a personal
383 machine. The machine learning tool, Vibrant, uses protein features derived from
384 multiple HMM searches. As a result of a large domain space, this tool ran for a
385 significantly longer amount of time (1 week for full complexity simulations) relative
386 to the other tools on the shared compute cluster.

387 This study benchmarked and compared the performance of viral identification
388 tools in metagenomics. The viral identification performance measures, in conjunc-
389 tion with the genera and prophage recall, highlights the advantages and challenges
390 of using specific viral identification tools, and can be used as a guide to assist the
391 selection of tools for subsequent research.

392 Conclusion

393 In summary, we tested the performance of nine viral identification tools on
394 30 simulated metagenomes. The underlying read distribution has little impact on
395 average tool performance. Increasing contig length and decreasing taxonomic com-
396 plexity improved the average performance of the tools. Vibrant performed the best
397 at the identification of prophages in metagenomics. Overall, the tool that averaged
398 the best F1 score was VirSorter, while Kraken2 lead all other tools in precision. The
399 results of these simulations should provide researchers with a guide to selecting the
400 appropriate tool for their own viral identification research.

401 **List of Abbreviations**

402 contig – contiguous sequence
 403 BAM – bacteriophage adherence to mucus
 404 IBD – inflammatory bowel disease
 405 VLP – viral like particles
 406 pVOG – prokaryotic virus orthologous group
 407 HMM – hidden Markov model
 408 EMP – Earth Microbiome Project
 409 KB – kilo-basepairs
 410

411 **Declarations**

412 **Ethics approval and consent to participate**
 413 Not applicable

414 **Consent for publication**

415 Not applicable

416 **Availability of data and materials**

417 All scripts used to derive the figures and additional preprocessing workflows are available on the Strong Lab GitHub
 418 at <https://github.com/Strong-Lab/Viral-Classification.in.Metagenomics>. VirBrant is available on the Strong Lab
 419 GitHub at <https://github.com/Strong-Lab/VirBrant>. VirKraken, the Kraken2 extension used in this study is
 420 available on the Strong Lab GitHub at <https://github.com/Strong-Lab/VirKraken> and on PyPI. The fasta files of
 421 the simulations are found at <https://tinyurl.com/fastavm>.

422 **Competing interests**

423 The authors have declared no competing interests in this work.

424 **Funding**

425 CG is supported by NLM 5 T15 LM009451-12.

426 **Author's contributions**

427 CG and MS conceived the study. CG and JH performed data analysis. CG and MS wrote the manuscript. All authors
 428 read, revised, and approved the final draft.

429 **Acknowledgements**

430 The authors would like to thank Chris Miller, James Costello, Catherine Lozupone, and Kirk Harris for their helpful
 431 comments.

432 **Author details**

433 ¹Center for Genes, Environment, and Health, National Jewish Health, 1400 Jackson Street, 80206 Denver,
 434 Colorado, USA. ²Computational Bioscience, University of Colorado Anschutz, 12801 E 17th Avenue, 80045 Aurora,
 435 Colorado, USA.

436 **References**

- 437 1. Ackermann, H.-W.: 5500 phages examined in the electron microscope. *Archives of virology* **152**(2), 227–243
 438 (2007)
- 439 2. Modi, S.R., Lee, H.H., Spina, C.S., Collins, J.J.: Antibiotic treatment expands the resistance reservoir and
 440 ecological network of the phage metagenome. *Nature* **499**(7457), 219–222 (2013)
- 441 3. Brüßow, H., Chanchaya, C., Hardt, W.-D.: Phages and the evolution of bacterial pathogens: from genomic
 442 rearrangements to lysogenic conversion. *Microbiology and molecular biology reviews* **68**(3), 560–602 (2004)
- 443 4. Barr, J.J., Auro, R., Furlan, M., Whiteson, K.L., Erb, M.L., Pogliano, J., Stotland, A., Wolkowicz, R., Cutting,
 444 A.S., Doran, K.S., *et al.*: Bacteriophage adhering to mucus provide a non-host-derived immunity. *Proceedings*
 445 *of the National Academy of Sciences* **110**(26), 10771–10776 (2013)
- 446 5. Martínez, I., Muller, C.E., Walter, J.: Long-term temporal analysis of the human fecal microbiota revealed a
 447 stable core of dominant bacterial species. *PLoS one* **8**(7), 69621 (2013)
- 448 6. Minot, S., Bryson, A., Chehoud, C., Wu, G.D., Lewis, J.D., Bushman, F.D.: Rapid evolution of the human gut
 449 virome. *Proceedings of the National Academy of Sciences* **110**(30), 12450–12455 (2013)
- 450 7. Gogokhia, L., Buhrke, K., Bell, R., Hoffman, B., Brown, D.G., Hanke-Gogokhia, C., Ajami, N.J., Wong, M.C.,
 451 Ghazaryan, A., Valentine, J.F., *et al.*: Expansion of bacteriophages is linked to aggravated intestinal
 452 inflammation and colitis. *Cell host & microbe* **25**(2), 285–299 (2019)
- 453 8. Wagner, J., Maksimovic, J., Farries, G., Sim, W.H., Bishop, R.F., Cameron, D.J., Catto-Smith, A.G.,
 454 Kirkwood, C.D.: Bacteriophages in gut samples from pediatric crohn's disease patients: metagenomic analysis
 455 using 454 pyrosequencing. *Inflammatory bowel diseases* **19**(8), 1598–1608 (2013)
- 456 9. Megremis, S., Constantinides, B., Xepapadaki, P., Bachert, C., Neurath-Finotto, S., Jartti, T., Kowalski, M.L.,
 457 Sotiropoulos, A.G., Tapinos, A., Vuorinen, T., *et al.*: Bacteriophage deficiency characterizes respiratory virome
 458 dysbiosis in childhood asthma. *bioRxiv* (2020)
- 459 10. Vartoukian, S.R., Palmer, R.M., Wade, W.G.: Strategies for culture of 'unculturable' bacteria. *FEMS*
 460 *microbiology letters* **309**(1), 1–7 (2010)
- 461 11. Rohwer, F., Edwards, R.: The phage proteomic tree: a genome-based taxonomy for phage. *Journal of*
 462 *bacteriology* **184**(16), 4529–4535 (2002)
- 463 12. Edwards, R.A., Rohwer, F.: Viral metagenomics. *Nature Reviews Microbiology* **3**(6), 504–510 (2005)
- 464 13. Jurtz, V.I., Villarreal, J., Lund, O., Voldby Larsen, M., Nielsen, M.: Metaphinder—identifying bacteriophage
 465 sequences in metagenomic data sets. *PLoS One* **11**(9), 0163111 (2016)

- 466 14. Graziotin, A.L., Koonin, E.V., Kristensen, D.M.: Prokaryotic virus orthologous groups (pvogs): a resource for
467 comparative genomics and protein family annotation. *Nucleic acids research*, 975 (2016)
- 468 15. Kieft, K., Zhou, Z., Anantharaman, K.: Vibrant: automated recovery, annotation and curation of microbial
469 viruses, and evaluation of viral community function from genomic sequences. *Microbiome* **8**(1), 1–23 (2020)
- 470 16. Guo, J., Bolduc, B., Zayed, A.A., Varsani, A., Dominguez-Huerta, G., Delmont, T.O., Pratama, A.A., Gazitúa,
471 M.C., Vik, D., Sullivan, M.B., *et al.*: Virsorter2: a multi-classifier, expert-guided approach to detect diverse dna
472 and rna viruses. *Microbiome* **9**(1), 1–13 (2021)
- 473 17. Martínez-García, M., Santos, F., Moreno-Paz, M., Parro, V., Antón, J.: Unveiling viral–host interactions within
474 the ‘microbial dark matter’. *Nature communications* **5**(1), 1–8 (2014)
- 475 18. Paez-Espino, D., Eloie-Fadrosch, E.A., Pavlopoulos, G.A., Thomas, A.D., Huntemann, M., Mikhailova, N.,
476 Rubin, E., Ivanova, N.N., Kyrpides, N.C.: Uncovering earth’s virome. *Nature* **536**(7617), 425–430 (2016)
- 477 19. Ren, J., Ahlgren, N.A., Lu, Y.Y., Fuhrman, J.A., Sun, F.: Virfinder: a novel k-mer based tool for identifying
478 viral sequences from assembled metagenomic data. *Microbiome* **5**(1), 69 (2017)
- 479 20. Ren, J., Song, K., Deng, C., Ahlgren, N.A., Fuhrman, J.A., Li, Y., Xie, X., Poplin, R., Sun, F.: Identifying
480 viruses from metagenomic data using deep learning. *Quantitative Biology*, 1–14 (2020)
- 481 21. Wood, D.E., Lu, J., Langmead, B.: Improved metagenomic analysis with kraken 2. *Genome biology* **20**(1), 1–13
482 (2019)
- 483 22. Federhen, S.: The ncbi taxonomy database. *Nucleic acids research* **40**(D1), 136–143 (2012)
- 484 23. Garretto, A., Hatzopoulos, T., Putonti, C.: virmine: automated detection of viral sequences from complex
485 metagenomic samples. *PeerJ* **7**, 6695 (2019)
- 486 24. Amgarten, D., Braga, L.P., da Silva, A.M., Setubal, J.C.: Marvel, a tool for prediction of bacteriophage
487 sequences in metagenomic bins. *Frontiers in genetics* **9**, 304 (2018)
- 488 25. Dutilh, B.E., Cassman, N., McNair, K., Sanchez, S.E., Silva, G.G., Boling, L., Barr, J.J., Speth, D.R.,
489 Seguritan, V., Aziz, R.K., *et al.*: A highly abundant bacteriophage discovered in the unknown sequences of
490 human faecal metagenomes. *Nature communications* **5**(1), 1–11 (2014)
- 491 26. Alex Reynolds: Kmer-counter. <https://github.com/alexpreynolds/kmer-counter>
- 492 27. Breiman, L., Friedman, J., Stone, C.J., Olshen, R.A.: Classification and Regression Trees. CRC press, ???
493 (1984)
- 494 28. Roux, S., Enault, F., Hurwitz, B.L., Sullivan, M.B.: Virsorter: mining viral signal from microbial genomic data.
495 *PeerJ* **3**, 985 (2015)
- 496 29. Hatfull, G.F., Jacobs-Sera, D., Lawrence, J.G., Pope, W.H., Russell, D.A., Ko, C.-C., Weber, R.J., Patel, M.C.,
497 Germane, K.L., Edgar, R.H., *et al.*: Comparative genomic analysis of 60 mycobacteriophage genomes: genome
498 clustering, gene acquisition, and gene size. *Journal of molecular biology* **397**(1), 119–143 (2010)
- 499 30. Hatfull, G.F., Cresawn, S.G., Hendrix, R.W.: Comparative genomics of the mycobacteriophages: insights into
500 bacteriophage evolution. *Research in microbiology* **159**(5), 332–339 (2008)
- 501 31. Hyatt, D., Chen, G.-L., LoCascio, P.F., Land, M.L., Larimer, F.W., Hauser, L.J.: Prodigal: prokaryotic gene
502 recognition and translation initiation site identification. *BMC bioinformatics* **11**(1), 1–11 (2010)
- 503 32. Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P.,
504 Weiss, R., Dubourg, V., *et al.*: Scikit-learn: Machine learning in python. *the Journal of machine Learning*
505 *research* **12**, 2825–2830 (2011)
- 506 33. Chen, T., Guestrin, C.: Xgboost: A scalable tree boosting system. In: Proceedings of the 22nd Acm Sigkdd
507 International Conference on Knowledge Discovery and Data Mining, pp. 785–794 (2016)
- 508 34. Mihara, T., Nishimura, Y., Shimizu, Y., Nishiyama, H., Yoshikawa, G., Uehara, H., Hingamp, P., Goto, S.,
509 Ogata, H.: Linking virus genomes with host taxonomy. *Viruses* **8**(3), 66 (2016)
- 510 35. Thompson, L.R., Sanders, J.G., McDonald, D., Amir, A., Ladau, J., Locoy, K.J., Prill, R.J., Tripathi, A.,
511 Gibbons, S.M., Ackermann, G., *et al.*: A communal catalogue reveals earth’s multiscale microbial diversity.
512 *Nature* **551**(7681), 457–463 (2017)
- 513 36. Gonzalez, A., Navas-Molina, J.A., Kosciolk, T., McDonald, D., Vázquez-Baeza, Y., Ackermann, G., DeReus,
514 J., Janssen, S., Swafford, A.D., Orchanian, S.B., *et al.*: Qiita: rapid, web-enabled microbiome meta-analysis.
515 *Nature methods* **15**(10), 796–798 (2018)
- 516 37. Wexler, H.M.: Bacteroides: the good, the bad, and the nitty-gritty. *Clinical microbiology reviews* **20**(4),
517 593–621 (2007)
- 518 38. Labrie, S.J., Samson, J.E., Moineau, S.: Bacteriophage resistance mechanisms. *Nature Reviews Microbiology*
519 **8**(5), 317–327 (2010)
- 520 39. Gourel, H., Karlsson-Lindsjö, O., Hayer, J., Bongcam-Rudloff, E.: Simulating illumina metagenomic data with
521 insilicoseq. *Bioinformatics* **35**(3), 521–522 (2019)
- 522 40. Fritz, A., Hofmann, P., Majda, S., Dahms, E., Dröge, J., Fiedler, J., Lesker, T.R., Belmann, P., DeMaere,
523 M.Z., Darling, A.E., *et al.*: Camisim: simulating metagenomes and microbial communities. *Microbiome* **7**(1),
524 1–12 (2019)
- 525 41. Nurk, S., Meleshko, D., Korobeynikov, A., Pevzner, P.A.: metaspades: a new versatile metagenomic assembler.
526 *Genome research* **27**(5), 824–834 (2017)
- 527 42. Buchfink, B., Xie, C., Huson, D.H.: Fast and sensitive protein alignment using diamond. *Nature methods*
528 **12**(1), 59–60 (2015)
- 529 43. Hunter, J.D.: Matplotlib: A 2d graphics environment. *IEEE Annals of the History of Computing* **9**(03), 90–95
530 (2007)
- 531 44. Waskom, M., the seaborn development team: Mwaskom/seaborn. doi:10.5281/zenodo.592845.
532 <https://doi.org/10.5281/zenodo.592845>
- 533 45. Bobay, L.-M., Touchon, M., Rocha, E.P.: Pervasive domestication of defective prophages by bacteria.
534 Proceedings of the National Academy of Sciences **111**(33), 12127–12132 (2014)
- 535 46. Griffiths-Jones, S., Bateman, A., Marshall, M., Khanna, A., Eddy, S.R.: Rfam: an rna family database. *Nucleic*
536 *acids research* **31**(1), 439–441 (2003)
- 537 47. Nawrocki, E.P., Eddy, S.R.: Infernal 1.1: 100-fold faster rna homology searches. *Bioinformatics* **29**(22),

- 538 2933–2935 (2013)
- 539 48. Al Suwayyid, B.A., Rankine-Wilson, L., Speers, D.J., Wise, M.J., Coombs, G.W., Kahler, C.M.: Meningococcal
- 540 disease-associated prophage-like elements are present in neisseria gonorrhoeae and some commensal neisseria
- 541 species. *Genome biology and evolution* 12(2), 3938–3950 (2020)

542 **Figures**

543 **Figure 1: Protein compositional feature performance** The performance of the 4 protein compositional features in
544 the 10KB training dataset. **A)** Gene density represented by number of genes per 1KB. **B)** Median operon length is a
545 representative measure of strand switching frequency. An operon is defined as a set of closely linked genes on the
546 same strand. **C)** Percentage of overlapping peptides measured as a percentage of all predicted genes. Viruses that
547 have a lysogeny phase are known overlap genes for different life cycles. **D)** Median amino acid length as viral
548 peptides are commonly shorter than bacterial peptides.

549 **Figure 2: Relative abundance of genera in simulations** These figures highlight the relative abundance of the contigs
550 greater than 1KB. Bacterial contigs represent 98.46% of contigs, while phages and prophages combine for the
551 remaining 1.54%. **A)** The contig distribution within 15 soil simulations. **B)** The contig distribution within 15 clinical
552 simulations.

553 **Figure 3: F1 scores of tools by taxonomic conditions**

554 Dodge boxplots by taxonomic complexity arranged by average F1 performance with the best performing tools on the
555 right side of the x axis.

556 **Figure 4: F1 scores of tools across contig length bins in all simulations**

557 The average F1 performance of all tools increases as the bin representing contig lengths increases. All thirty
558 simulations are included as part of this figure, however in some simulations, predicted viral contigs of a specific
559 length are absent. This may cause some tools to have more data points than others.

560 **Figure 5: Viral recall by host genera in medium and full complexity simulations**

561 The 30 host genera of phage are listed in order of mean recall along the x-axis. The dotted line in the figure
562 demarcates 0.5 recall.

563 Tables

Table 1: Tools used in viral identification benchmarking study

| Tool | Last updated | Target | Viral homology matching | Compositional protein features | Machine learning classification | Programming skills required |
|--------------------------|--------------|--------|-------------------------|--------------------------------|---------------------------------|-----------------------------|
| <i>VirSorter</i> | 2015 | Virus | Yes | Yes | No | No |
| <i>VirSorter2</i> | 2020 | Virus | Yes | Yes | Yes | No |
| <i>VirFinder</i> | 2017 | Virus | No | No | Yes | Yes |
| <i>DeepVirFinder</i> | 2020 | Virus | No | No | Yes | Yes |
| <i>Vibrant</i> | 2020 | Virus | Yes | No | Yes | Yes |
| <i>MetaPhinder</i> | 2016 | Phage | Yes | No | No | Yes |
| <i>Earth Virome</i> | 2020 | Virus | Yes | No | No | Yes |
| <i>VirBrant</i> | 2020 | Phage | No | Yes | Yes | Yes |
| <i>Kraken2+VirKraken</i> | 2020 | Virus | Yes | No | No | Yes |

Table 2: Average performance and simulation rankings of tools at identifying phage

| Tool | F1 Rank | Precision Rank | Recall Rank | Average F1 | Average Precision | Average Recall |
|--------------------------|-------------|----------------|-------------|--------------|-------------------|----------------|
| <i>VirSorter</i> | 2.10 | 3.10 | 6.40 | 0.636 | 0.640 | 0.658 |
| <i>Kraken2</i> | 2.93 | 1.07 | 7.80 | 0.609 | 0.962 | 0.467 |
| <i>Vibrant</i> | 3.52 | 4.10 | 7.26 | 0.560 | 0.573 | 0.598 |
| <i>VirFinder</i> | 3.93 | 2.30 | 9.32 | 0.548 | 0.717 | 0.450 |
| <i>DeepVirFinder</i> | 5.04 | 5.38 | 7.90 | 0.432 | 0.392 | 0.496 |
| <i>VirSorter2</i> | 5.27 | 5.93 | 3.10 | 0.463 | 0.341 | 0.797 |
| <i>VirBrant</i> | 5.40 | 6.00 | 3.60 | 0.413 | 0.317 | 0.755 |
| <i>VirBrant Proteins</i> | 7.47 | 7.70 | 4.63 | 0.142 | 0.213 | 0.717 |
| <i>MetaPhinder</i> | 8.73 | 8.67 | 2.83 | 0.082 | 0.138 | 0.842 |
| <i>Earth Virome</i> | 9.73 | 9.73 | 1.78 | 0.023 | 0.044 | 0.872 |

Table 3: Average performance and simulation rankings of tools at identifying prophage

| Tool | F1 Rank | Precision Rank | Recall Rank | Prophage F1 | Prophage Precision | Prophage Recall |
|--------------------------|-------------|----------------|-------------|--------------|--------------------|-----------------|
| <i>Vibrant</i> | 1.15 | 1.95 | 7.45 | 0.169 | 0.146 | 0.231 |
| <i>VirSorter</i> | 2.11 | 1.66 | 8.76 | 0.147 | 0.144 | 0.164 |
| <i>VirSorter2</i> | 2.70 | 3.40 | 4.70 | 0.117 | 0.0685 | 0.453 |
| <i>VirBrant</i> | 3.95 | 4.45 | 7.13 | 0.0446 | 0.0269 | 0.252 |
| <i>VirBrant Proteins</i> | 5.20 | 5.50 | 6.05 | 0.0188 | 0.00978 | 0.342 |
| <i>Kraken2</i> | 6.70 | 1.85 | 9.90 | 0.0169 | 0.172 | 0.00896 |
| <i>MetaPhinder</i> | 6.65 | 6.85 | 2.90 | 0.0152 | 0.00776 | 0.588 |
| <i>Earth Virome</i> | 6.85 | 7.00 | 1.60 | 0.0117 | 0.00588 | 0.728 |
| <i>VirFinder</i> | 8.88 | 8.88 | 2.73 | 0.00725 | 0.00365 | 0.705 |
| <i>DeepVirFinder</i> | 8.79 | 9.03 | 2.79 | 0.00647 | 0.00325 | 0.637 |

Figures

Figure 1

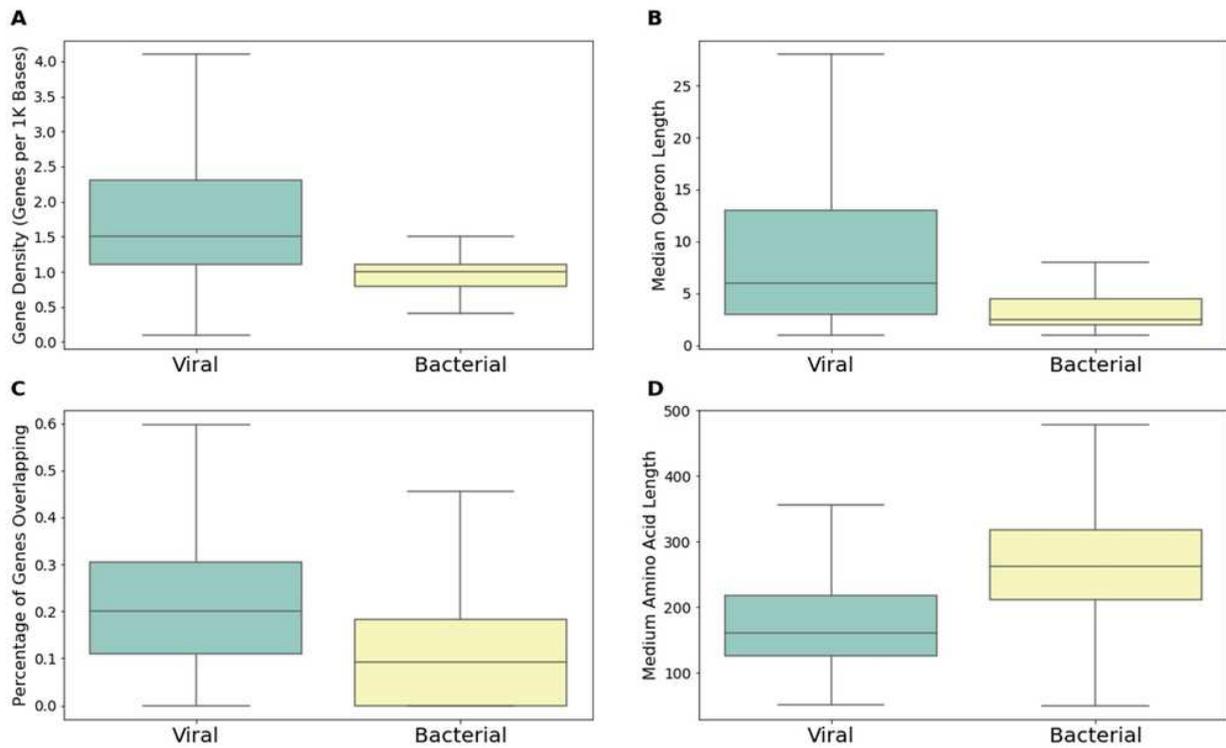


Figure 1

Protein compositional feature performance The performance of the 4 protein compositional features in the 10KB training dataset. A) Gene density represented by number of genes per 1KB. B) Median operon length is a representative measure of strand switching frequency. An operon is defined as a set of closely linked genes on the same strand. C) Percentage of overlapping peptides measured as a percentage of all predicted genes. Viruses that have a lysogeny phase are known to overlap genes for different life cycles. D) Median amino acid length as viral peptides are commonly shorter than bacterial peptides.

Figure 2

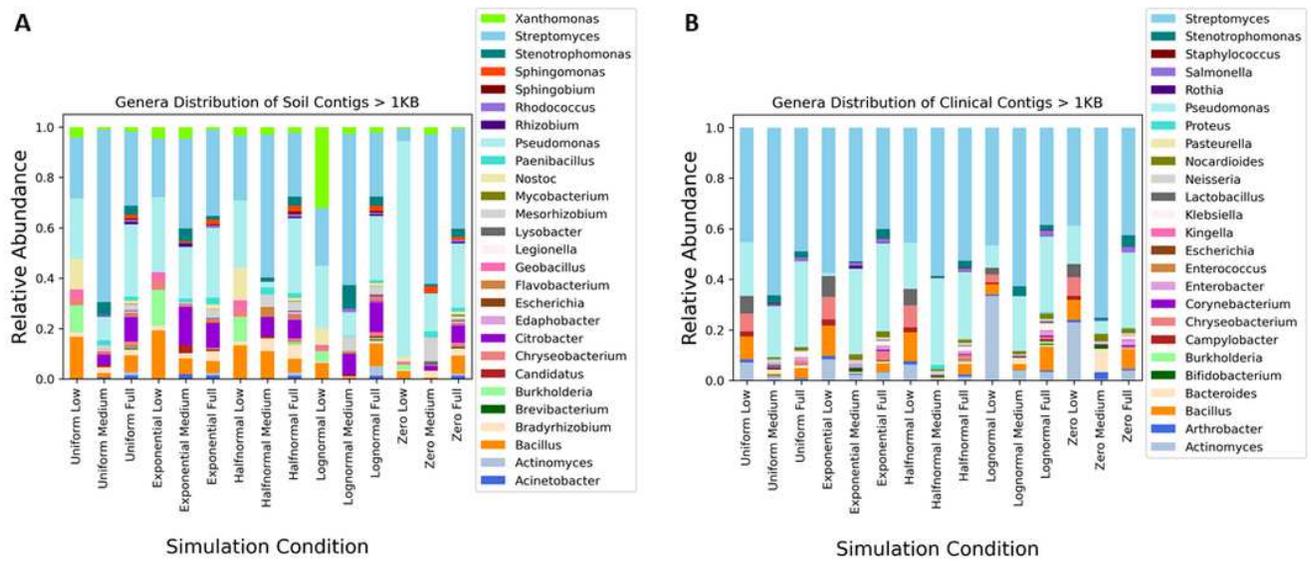


Figure 2

Relative abundance of genera in simulations. These figures highlight the relative abundance of the contigs greater than 1KB. Bacterial contigs represent 98.46% of contigs, while phages and prophages combine for the remaining 1.54%. A) The contig distribution within 15 soil simulations. B) The contig distribution within 15 clinical simulations.

Figure 3

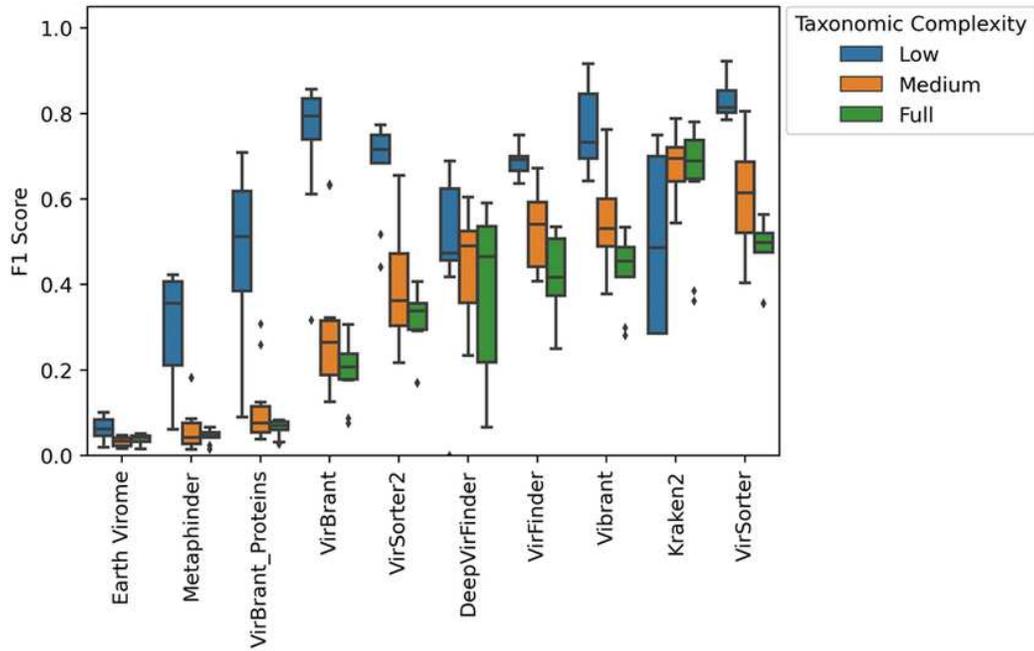


Figure 3

F1 scores of tools by taxonomic conditions Dodge boxplots by taxonomic complexity arranged by average F1 performance with the best performing tools on the right side of the x axis.

Figure 4

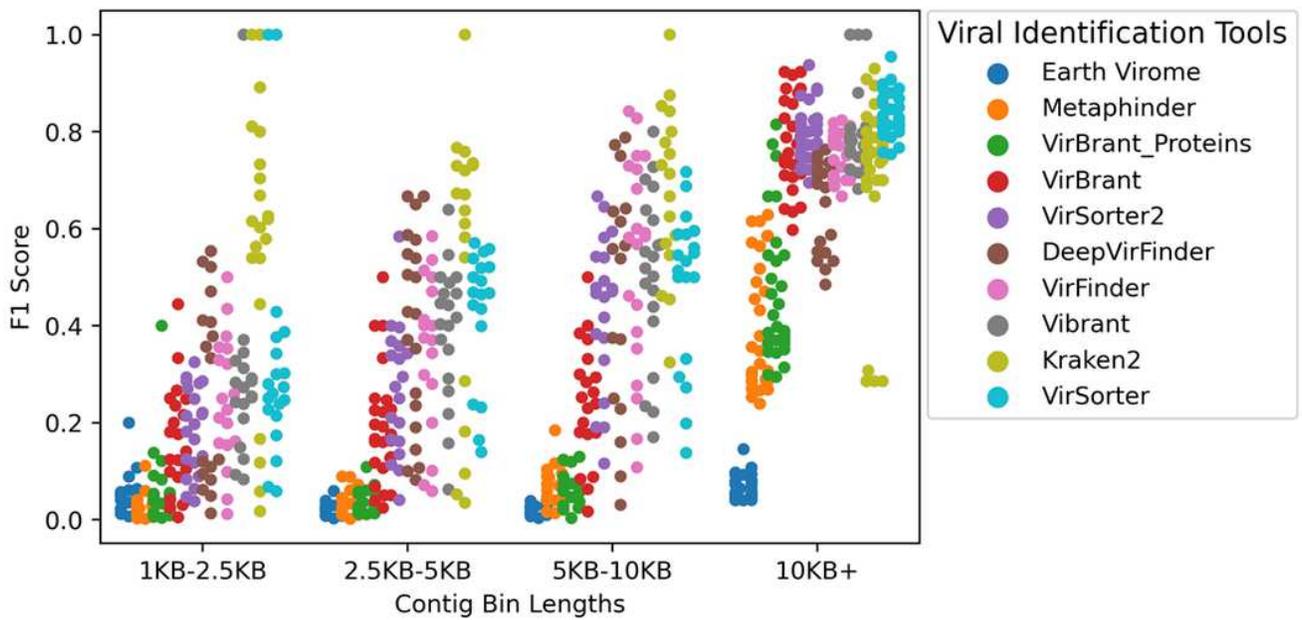


Figure 4

F1 scores of tools across contig length bins in all simulations. The average F1 performance of all tools increases as the bin representing contig lengths increases. All thirty simulations are included as part of this figure, however in some simulations, predicted viral contigs of a specific length are absent. This may cause some tools to have more data points than others.

Figure 5

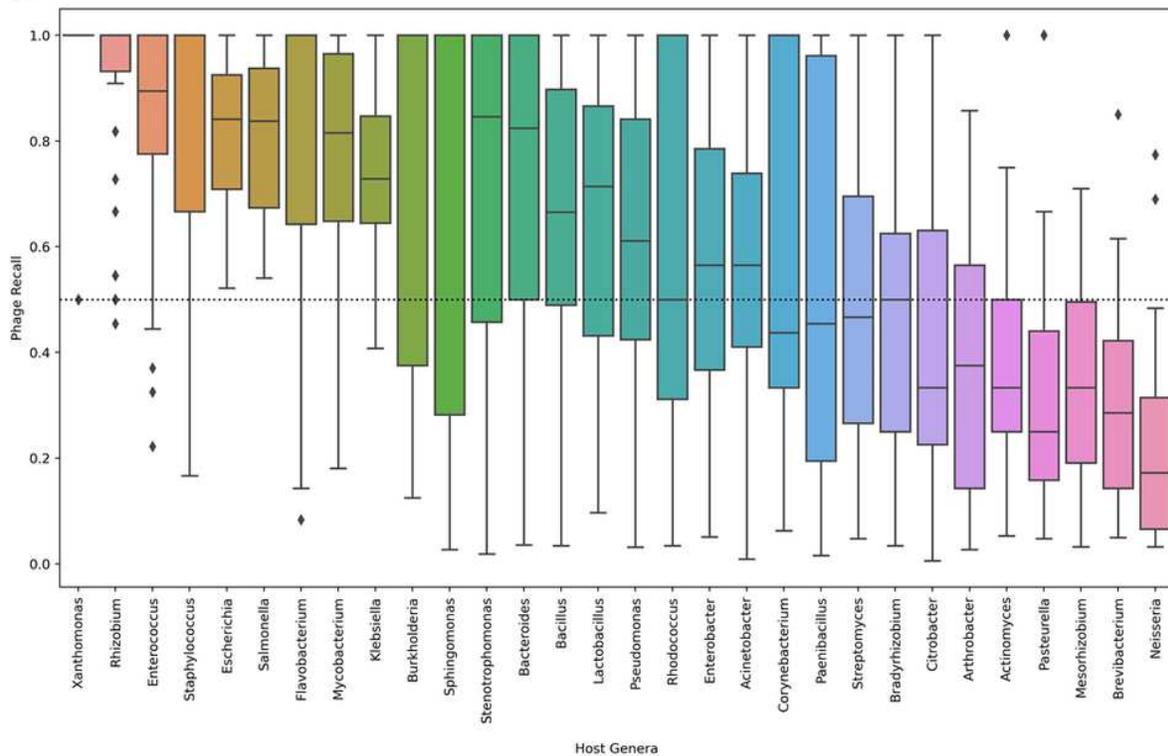


Figure 5

Viral recall by host genera in medium and full complexity simulations. The 30 host genera of phage are listed in order of mean recall along the x-axis. The dotted line in the figure demarcates 0.5 recall.

Supplementary Files

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

- [GlickmanBMCSupplemental.pptx](#)
- [SupplementalGlickmanBMC.pdf](#)