

1 **Title:** iVirus: facilitating new insights in viral ecology with software and community
2 datasets imbedded in a cyberinfrastructure

3

4 **Authors:** Benjamin Bolduc¹, Ken Youens-Clark², Simon Roux¹, Bonnie L. Hurwitz^{*2},
5 Matthew B. Sullivan^{*1,3}

6

7 **Affiliations:** ¹Department of Microbiology, The Ohio State University, Columbus, OH
8 43210

9 ²Department of Agricultural and Biosystems Engineering, University of Arizona, Tucson,
10 AZ 85721

11 ³Department of Civil, Environmental and Geodetic Engineering, The Ohio State
12 University, Columbus, OH 43210

13 * Co-corresponding authors

14

15 Corresponding Authors: Matthew B. Sullivan (mbsulli@gmail.com; Riffe Building Rm
16 914, 496 W 12th Ave Columbus, OH 43210, USA; 1-614-247-1616) and Bonnie L.
17 Hurwitz (bhurwitz@email.arizona.edu; 1177 E. 4th Street, Shantz Building, Room 624,
18 Tucson, AZ, USA 85721-0038; 1-520-626-9819)

19

20 **Abstract**

21 Microbes impact nutrient and energy transformations throughout the world's
22 ecosystems, yet they do so under viral constraints. In complex communities, viral
23 metagenome (virome) sequencing is transforming our ability to quantify viral diversity
24 and impacts. While some bottlenecks, e.g., few reference genomes and non-quantitative
25 viromics, have been overcome, the void of centralized datasets and specialized tools now
26 prevents viromics from being broadly applied to answer fundamental ecological
27 questions. Here we present iVirus, a community resource that leverages the CyVerse
28 cyberinfrastructure to provide access to viromic tools and datasets. The iVirus Data
29 Commons contains both raw and processed data from 1866 samples and 73 projects
30 derived from global ocean expeditions, as well as existing and legacy public repositories.
31 Through the CyVerse Discovery Environment, users can interrogate these datasets using
32 existing analytical tools (software applications known as “Apps”) for assembly, ORF
33 prediction, and annotation, as well as several new Apps specifically developed for
34 analyzing viromes. Because Apps are web-based and powered by CyVerse super-
35 computing resources, they enable scalable analyses for a broad user base. Finally, a use-
36 case scenario documents how to apply these advances towards new data. This growing
37 iVirus resource should help researchers utilize viromics as yet another tool to elucidate
38 viral roles in nature.

39

40 **Viral metagenomics – an ecological tool with increasing impact**

41 Since the first viral metagenomic study was conducted in marine systems over a
42 decade ago (Breitbart *et al.*, 2002), the field has now expanded to include ecological

43 studies of viral communities throughout the oceans globally, as well as diverse lakes and
44 eukaryote-associated samples, including humans (Djikeng *et al.*, 2009; Hannigan *et al.*,
45 2015; Hurwitz and Sullivan, 2013; Roux *et al.*, 2012; Stern *et al.*, 2012). Highlights of
46 some of the ecological advances enabled by these studies include revealing (i) that virus-
47 encoded ‘auxiliary metabolic genes’ (AMGs) extend far beyond the photosynthesis genes
48 known from cyanobacterial cultures (Hurwitz *et al.*, 2013; Hurwitz, Westveld, *et al.*,
49 2014; Sharon *et al.*, 2011), (ii) long-term co-evolutionary features between viruses and
50 their microbial hosts in both the human gut (Minot *et al.*, 2013) and the oceans (Hurwitz,
51 Westveld, *et al.*, 2014), and (iii) the ecological drivers of viral community structure
52 throughout the Pacific Ocean (Hurwitz, Westveld, *et al.*, 2014) and global surface oceans
53 (Brum *et al.*, 2015).

54 Many technological advances have enabled these discoveries – including
55 optimized sampling strategies specific for viruses (reviewed in Duhaime and Sullivan,
56 2012; Solonenko *et al.*, 2013) and improvements in low-input library preparation
57 methods and decreased sequencing costs (Reuter *et al.*, 2015; Caporaso *et al.*, 2012) –
58 and this has led to a data deluge whereby analytical limitations now represent the major
59 bottleneck for virome-enabled viral ecology as follows. First, the number and size of
60 newly generated metagenomes necessitates large-scale data storage and compute needs
61 that require the development of community-available infrastructures specialized to viral
62 sequence data. Second, the lack of tools for analyzing these large-scale datasets requires
63 development by programmers, who speak a different ‘language’ from researchers
64 generating much of the data. This often results in published code in public repositories
65 that can be difficult to install or use without computational training. Finally, finding and

66 analyzing viral datasets for comparative metagenomics is laborious and time-consuming
67 as raw and processed viral metagenomic datasets are deposited across a diverse array of
68 data repositories, such as Genbank (Benson *et al.*, 1998), EMBL (Kanz *et al.*, 2005),
69 NCBI's genomes project (Wheeler *et al.*, 2003), Metavir (Roux *et al.*, 2014), and
70 VIROME (Wommack *et al.*, 2012).

71 Together, these technological limitations impede researchers from applying new
72 tools to their data or leave them dependent on outsourcing data analysis to those
73 unfamiliar with the ecosystems being studied. To enable scalability and accessibility of
74 viral ecology research, we developed iVirus, a collection of software and datasets
75 leveraging the CyVerse cyberinfrastructure (formerly iPlant Collaborative) to provides
76 users with free access to computing, data management and storage, and analysis toolkits
77 (Goff *et al.*, 2011). Briefly, iVirus seeks to collect viral sequence datasets in its Data
78 Commons, adapt pre-existing metagenomic tools as software applications (referred
79 henceforth as *Apps*), and develop new analytical capabilities *within* the CyVerse
80 cyberinfrastructure. Together, these advances help consolidate cutting-edge tools and
81 curated datasets to empower researchers seeking to incorporate viral ecology into their
82 own work.

83 Here we summarize iVirus's current capabilities, and invite community feedback,
84 through the protocols.io interface (described later), to allow us to improve iVirus so that
85 it becomes an indispensable tool for ecologists seeking to include viruses in their studies.

86

87 **What is CyVerse and how does it help my research?**

88 CyVerse (Goff *et al.*, 2011) is an NSF-funded platform that seeks to bring
89 together biologists and computer scientists to solve ‘big data’ problems in biology.
90 Within the CyVerse cyberinfrastructure, users conduct research by navigating the
91 Discovery Environment (DE) to identify datasets from the Data Commons (next section)
92 and conduct analyses using Apps. Apps are like a computer software program, except
93 that they (i) have been pre-installed, (ii) leverage large-scale CyVerse compute resources,
94 and (iii) can be integrated within a larger data context and workflow. This can improve
95 biological research in multiple ways. First, it helps minimize installation issues and local
96 systems administration needs that often impede biological research. Second, Apps can be
97 linked together to create analytical *workflows* where the output from one App is used as
98 the input for the next, in a linear manner. Because users can select which App they use at
99 each stage in the workflow, the user can copy and update new workflows as new analysis
100 tools emerge. Third, Apps ensure reproducibility and validity of research studies because
101 they can be encoded into versioned Apps, along with raw or processed example datasets
102 directly from the author. CyVerse can also assign digital object identifier (DOI) to Apps
103 or datasets to allow longer term digital preservation and citable referencing in research
104 articles. Finally, because all Apps and their output are tied to the user’s home directory in
105 the CyVerse Discovery Environment, all data are collected in one place. This avoids the
106 common problem of data being scattered among multiple systems (HPC, personal/lab
107 computer, cloud computing) due to system-level requirements for implementing myriad
108 bioinformatics software used for modern metagenomic analyses.

109 CyVerse Apps can be built by any developer using any source programming
110 language to create community-specific tools. Moreover, the developer can encode

111 hardware or software requirements within the App to lessen the burden on the user in
112 installing and implementing that App. Once an App has been developed it can be shared
113 with the community through the CyVerse cyberinfrastructure by a request to the CyVerse
114 team via a simple public submission form. Once an App is public it can be vetted by the
115 research community via a 5-star rating system and feedback forms. Further, community
116 developers can refine an App by duplicating and modifying it, and then re-publish the
117 new App for recognition and further vetting and use by the community.

118 For developers, the process of creating an App follows one of two routes. First,
119 Apps can be developed using CyVerse's API (application program interface) called
120 AGAVE (<http://agaveapi.co/>) that provides the developer with simple commands to
121 access input data and write logs and results to the CyVerse Data Store. The developer can
122 use the API to specify computational requirements for running code on HPC resources at
123 the Texas Advanced Computer Center (TACC) that are integrated in the CyVerse
124 cyberinfrastructure. This process allows the developer to match the code to HPC compute
125 resources and circumvent difficulties users might experience in installing and re-using the
126 code on different systems. Alternatively, Apps can be deployed using Docker images
127 (www.docker.com), where the code is packaged with additional software dependences
128 and can be run on CyVerse Docker-dedicated servers. Docker's compute autonomy and
129 portability alone have made Docker images a mainstay for releasing open source code to
130 the user community (Merkel, 2014), in addition to traditional code repositories such as
131 Github (<https://github.com>). CyVerse extends Docker by allowing developers to deploy
132 Docker images on CyVerse compute resources and attach these images to easy-to-use,
133 web-based Apps in the CyVerse Discovery Environment. This provides developers with a

134 means to publish code in an accessible format to a growing research community and to
135 gain feedback on its utility rather than be drowned in inquiries about installation minutia.

136

137 **Centralized viral metagenomic data resources in the iVirus Data Commons**

138 The CyVerse cyberinfrastructure provides a common ecosystem for data, big or
139 small, by providing a mechanism for communities to share data through a Community
140 Data Commons. The iVirus Data Commons leverages these CyVerse resources to make
141 datasets accessible from the Pacific Ocean Virome (POV; Hurwitz and Sullivan, 2013),
142 *Tara* Oceans Virome (TOV; Brum *et al.*, 2015), Southern Ocean Virome (Brum *et al.*,
143 2016), Virsorter Curated Dataset (Roux, Enault, *et al.*, 2015), and legacy viral datasets
144 from the retired Cyberinfrastructure for Advanced Microbial Ecology Research and
145 Analysis (CAMERA) project (Seshadri *et al.*, 2007). Beyond these, viral data were also
146 mined and hand-curated from Genbank's Sequence Read Archive (SRA; Benson *et al.*,
147 1998) and MG-RAST (Wilke *et al.*, 2015). In total, the iVirus Data Commons now
148 contains data from 73 projects, 1 866 samples, and 5.5 billion reads, including contigs
149 assembled from 75 viromes and 121 viral genomes.

150

151 **iVirus Apps are Geared to Viral Metagenomics and Community Ecology**

152 iVirus Apps are developed using protocols defined by CyVerse – either through
153 the Agave API or Docker – with focus on those needed for viral metagenomics. One
154 operational goal of iVirus is to collect and deploy the most commonly used tools for viral
155 metagenome pipelines – from raw read processing to assembly and analysis (see
156 overview in Fig. 1, and Use Case Scenario presented in next section). This includes tools

157 for read quality control, assemblers adapted to different input read types and various tools
158 for analyzing assembled viral sequence data.

159 Because iVirus exists *within* the CyVerse environment, Apps developed through
160 other community-based CyVerse efforts are also available to iVirus. These include Apps
161 relevant for viral metagenomics such as read pre-filtering for assembly, gene calling,
162 taxonomic identification and sequence alignment tools. This list also includes microbial
163 metagenomic analysis Apps from iMicrobe that can be used to assemble contigs from
164 metagenomes, predict and functionally annotated genes in these contigs, and align them
165 to each other and reference genomes (where available) for comparative genomic
166 analyses. In general, Apps developed for HPC are located in a separate category within
167 CyVerse, under “High Performance Computing” while Docker and non-HPC-enabled
168 Apps are organized into folders appropriate for their “theme” (i.e. iMicrobe, iVirus,
169 Functional Analysis, etc). Because Apps are constantly being updated and developed, a
170 full list of iVirus Apps are maintained at <http://ivirus.us/available-tools/>, along with
171 computational protocols, App descriptions, relevant articles on tools, and news in
172 protocols.io at <http://protocols.io/groups/ivirus>

173 iVirus is a sub-project of iMicrobe and uses open source software developed by
174 that project (<http://www.imicrobe.us> and <http://protocols.io/groups/imicrobe>) to query,
175 search, and download data in the iVirus Data Commons from a project specific data
176 website (<http://data.ivirus.us/>). This web-based resource allows users to perform
177 advanced searches on top of the iVirus Data Commons to discover viral datasets based on
178 related metadata. To better enable search capabilities, iVirus metadata are mapped to the
179 iMicrobe ontology that interconnects existing standards and terminology from the

180 Minimal Information about any (x) Sequence Ontology (MIxS) and community specific
181 ontologies (BCO-DMO, ENVO, CheBI, BCO, OBOE). As such, the location of project
182 specific datasets can be easily discovered and re-used within CyVerse.

183 Beyond these more generally usable Apps, we have developed several iVirus
184 Apps specifically for viral metagenomic study, with more being added as needs arise and
185 development opportunities become available. In many cases, Apps developed for iVirus
186 and iMicrobe also handle file manipulations, such as compression, separating reads,
187 converting formats, so as to eliminate whenever possible the “minor” details that are
188 often time-consuming and rate-limiting for users.

189 A brief summary of selected iVirus Apps follows, along with reference to their
190 source tools and how they have already enabled viral ecology where available. A broader
191 analytical pipeline for processing viral metagenomes is overviewed in Figure 1 and
192 described in a User Case Scenario below and at

193 <https://www.protocols.io/view/Processing-a-Viral-Metagenome-Using-iVirus-ev3be8n>

194

195 **PCPipe**: This App compares open reading frames (ORFs) from a user-defined dataset to
196 existing viral protein clusters (PCs) as a means to organize proteins derived from viral
197 metagenomics into functional units that can be used as (i) a universal functional diversity
198 metric for viruses, (ii) a scaffold for iterative functional annotations, and (iii) input to
199 ecological comparisons through software such as QIIME, <http://qiime.org> (Caporaso *et*
200 *al.*, 2010). This is necessary as viral metagenomes are often dominated by novel
201 sequences, where only 10-20% of reads map to known proteins in reference databases. In
202 contrast, up to 50-70% of reads will typically map to PCs (Hurwitz and Sullivan, 2013).

203 The PCPipe App accepts user-generated ORFs from viral metagenomic assemblies as
204 input, matches them to ORFs in a user-supplied PC database, and then self-clusters the
205 remaining unclustered ORFs to capture the PCs unique to that dataset. Reference
206 sequences from new PCs are annotated using a collection of the non-redundant proteins
207 and associated annotation from the SIMAP database (Rattei *et al.*, 2009). PCs were
208 originally developed for analyzing unknown proteins from the Global Ocean Survey that
209 doubled the known protein universe at the time (GOS; Yooseph *et al.*, 2007). This
210 approach has proved similarly valuable for organizing viral protein sequence space
211 (Brum *et al.*, 2015; Hurwitz and Sullivan, 2013). Such an organizational tool has served
212 as a means to estimate the size of the global virome at a few million proteins (Cesar
213 Ignacio-Espinoza *et al.*, 2013), as well as to make ecological inferences about viral
214 communities with regards to their diversity (Hurwitz and Sullivan, 2013; Brum *et al.*,
215 2015; Roux *et al.*, 2012), niche differentiating genes (Hurwitz, Brum, *et al.*, 2014) and
216 ecological drivers (Brum *et al.*, 2015).

217

218 **VirSorter**: This App identifies viral sequences in microbial genomes and metagenomic
219 datasets (Roux, Hallam, *et al.*, 2015). This is necessary as viral genomes are
220 underrepresented in databases – e.g., 92% of 1,659 genome-sequenced phages derive
221 from only 4 of 54 known bacterial phyla (Roux, Enault, *et al.*, 2015). VirSorter can
222 identify diverse viral sequences from microbial datasets, both integrated in the host
223 chromosome and extrachromosomal. Briefly, VirSorter compares a dataset of nucleotide
224 sequences against a user-defined, pre-computed viral database that includes viral
225 sequences from RefSeq and (if desired) contigs assembled from viral metagenomes. The

226 comparison also takes into account viral hallmark genes, as well as statistical enrichment
227 of viral genes, depletion in hits to the PFAM database, and strand bias. VirSorter output
228 includes a summary file with “confidence” categories for each identified sequence, as
229 well as predicted proteins, PFAM domain hits, suspected circular sequences and metrics
230 files. This tool is powerful and highly scalable – its first application was to nearly 15 000
231 publically available archaeal and bacterial genomes, where VirSorter identified 12 498
232 new host-associated viruses and their genomes, which augmented publicly available viral
233 genome reference datasets approximately 10-fold (Roux, Hallam, *et al.*, 2015). Further,
234 VirSorter scales to handle contigs derived from metagenomic datasets (Roux, Enault, *et*
235 *al.*, 2015). VirSorter has since been used to identify viruses out of boiling hot springs in
236 Yellowstone National Park (Munson-McGee *et al.*, 2015), the *Tara* Oceans Viromes
237 (Lima-Mendez *et al.*, 2015), and hypersaline environments in the Atacama Desert (Crits-
238 Christoph *et al.*, 2016).

239

240 **vContact**: This App assigns contigs to taxonomic groups using the presence or absence
241 of shared PCs along the length of the contig. This is critical as viruses lack a universal
242 gene marker (Edwards and Rohwer, 2005) and less than 0.1% of viruses in natural
243 environments are represented in public databases (Brum *et al.*, 2015), which necessitates
244 new approaches to taxonomically classify surveyed viral genomes. Inspired by
245 algorithms to detect prophage in microbial genomes (Lima-Mendez *et al.*, 2008),
246 vContact clusters contigs by their PC profiles (note: see preferred method for PC
247 generation as vContact-PCs, but the user could generate PCs however they prefer).
248 Reference sequences and their taxonomic lineages can be seeded within the analysis to

249 improve clustering and taxonomic predictions. The vContact-generated network can be
250 mined for its contig clusters (VCs: viral clusters), which roughly correspond to sub-
251 family level viral taxonomy. This approach has been used to organize nearly all (99.3%)
252 of the 12,498 new viral sequences identified from publicly available microbial genomes
253 into 614 ‘viral clusters’, which represent approximately genus-level groupings (Roux,
254 Hallam, *et al.*, 2015). vContact can also incorporate annotations associated with contig
255 and PCs, allowing users to examine the relationship of any annotated contig/PC in
256 context of its vContact cluster.

257

258 **vContact-PCs**: This App serves as a companion tool for vContact to generate PCs using
259 a Markov clustering algorithm (MCL; Enright *et al.*, 2002). Users provide a BLASTP file
260 of an all-against-all protein comparison, and vContact-PCs parses the BLAST file and
261 applies MCL against its similarity scores. vContact-PCs then exports files formatted for
262 use with vContact.

263

264 **Fizkin**: This App performs Bayesian network analyses based on the amount of shared
265 sequence content in viromes and relevant environmental metadata. Specifically, the App
266 randomly subsets 300K reads (or the lowest common denominator for reads in viromes)
267 from up to 15 viromes and performs a pairwise all-vs-all kmer-based sequence
268 comparison between all virome pairs as previously described (Hurwitz, Westveld, *et al.*,
269 2014). The kmer search in Fizkin is implemented in Jellyfish (Marcais and Kingsford,
270 2011) and is used to rapidly generate a matrix of shared sequence counts between each
271 virome pair. This matrix is then used as input into a Bayesian network analysis (Hoff,

272 2005). The user can also input environmental metadata (continuous or discrete
273 measurements, or latitude and longitude) that are used as part of the analysis to determine
274 which environmental factors are significant in defining the structure of the network.
275 Output includes: (i) a table indicating which environmental factors are significant and (ii)
276 a social network graph to visually represent the distance between viromes where
277 statistical samples are taken from the marginal posterior distributions, and samples names
278 are placed at the posterior mean. This type of social network analysis has been used to
279 evaluate viral community structure and ecological drivers (Hurwitz, Westveld, *et al.*,
280 2014), as well as to quantify lysogeny through comparative analysis of experimentally
281 induced and non-induced viromes (Brum *et al.*, 2016).

282

283 **BatchBowtie**: This App runs bowtie2 on any number of reads files within a directory the
284 user selects. Read recruitment against reference sequences (i.e viral genomes) can be
285 employed as a means to visualize spatial or temporal distributions of genomes via the
286 reads relative abundance across samples (Brum *et al.*, 2015). The App additionally offers
287 the ability to convert between interleaved and non-interleaved fastq files, compressed
288 files, and can generate SAM and BAM-formatted outputs (typically used by
289 Read2RefMapper).

290

291 **Read2RefMapper**: This App consumes BAM alignment files (i.e. from BatchBowtie),
292 generating coverage tables and relative abundance plots - useful for identifying the
293 abundance of reads against a set of reference sequences. Users can select a variety of
294 filtering options based on the percent of read *and* reference covered (i.e. 75% of a

295 reference sequence must be covered to be considered “present”), alignment identities, as
296 well as numerous coverage calculations. If users provide a file with the size of each
297 metagenome, Read2RefMapper will also normalize the coverage between samples.

298

299 **Finding and using iVirus: A use case scenario ‘live’ at protocols.io**

300 Taken together – the Apps mentioned above, in addition to the Apps already
301 available in CyVerse – can be used to process a viral metagenome from “raw” sequence
302 to minimally characterized viral assemblies. The guide(s) for this *and other* examples are
303 available at protocols.io: [https://www.protocols.io/view/Processing-a-Viral-Metagenome-](https://www.protocols.io/view/Processing-a-Viral-Metagenome-Using-iVirus-ev3be8n)
304 [Using-iVirus-ev3be8n](https://www.protocols.io/view/Processing-a-Viral-Metagenome-Using-iVirus-ev3be8n). These protocols are organized as collections, and available within
305 the iVirus and iMicrobe groups at protocols.io, at <https://www.protocols.io/groups/ivirus>
306 and <https://www.protocols.io/groups/imicrobe>. These groups serve as a centralized
307 location that offers additional documentation, feedback, as well as citations using these
308 tools and protocols. We have utilized protocols.io here so as to keep evolving processing
309 steps up-to-date, as well as include images and annotations for each step. Further,
310 protocols.io is ideal for obtaining community feedback as it provides users the
311 opportunity to ask questions and/or interact with the protocol’s author through a simple to
312 user interface.

313 The use case scenario starts with test data, which are reads from publicly available
314 Ocean Sampling Day 2014 samples ([https://github.com/MicroB3-IS/osd-](https://github.com/MicroB3-IS/osd-analysis/wiki/Guide-to-OSD-2014-data)
315 [analysis/wiki/Guide-to-OSD-2014-data](https://github.com/MicroB3-IS/osd-analysis/wiki/Guide-to-OSD-2014-data)), a subset of which are already available on
316 CyVerse’s data store. Next a user must first register for a *free* account in CyVerse
317 (<http://user.iplantcollaborative.org/>) and then proceed through the process using available

318 iVirus Apps summarized in Fig. 1. While only a few Apps are highlighted for each step,
319 there are wide selections of choices for Apps available to the user. Example data is
320 provided for each step in the iVirus Data Commons found under the
321 `/iplant/home/shared/iVirus/ExampleData/`
322 (<https://de.iplantcollaborative.org/de/?type=data&folder=/iplant/home/shared/iVirus/ExampleData>) folder within the CyVerse Discovery Environment. Output from each stage in
323 this scenario is used as input for the next, though users can test any step individually, as
324 each stage is organized separately and contains its own folders with inputs and outputs to
325 help users identify which files are associated with each step.

327

328 **Step 1: Upload Read Data.** Before processing can begin, reads to be analyzed
329 must be uploaded to CyVerse's Data Store in the user's account. Small files can be
330 uploaded from the user's computer through the Discovery Environment's (DE) upload
331 menu, with larger files transferrable through popular SFTP software, such as Cyberduck
332 (<https://cyberduck.io>) and iRODS (<http://irods.org>). Data uploaded to a user's CyVerse
333 home directory are private and accessible only to the user until they grant access to
334 collaborators (via private invitation) or publish their data to the larger CyVerse
335 community by transferring files to a shared Community folder. Users can also analyze or
336 leverage growing publicly available datasets at iVirus in the Data Commons as
337 previously described. Use case read files are already uploaded to the iVirus Data
338 Commons.

339 **Step 2: Quality Control (QC) of Read Data using Trimmomatic.**

340 Protocol available at protocols.io: <https://www.protocols.io/view/Quality-Control->
341 [of-Reads-Using-Trimomatic-Cyverse-ewbbfan](https://www.protocols.io/view/Quality-Control-of-Reads-Using-Trimomatic-Cyverse-ewbbfan)

342 Once raw read data has been uploaded, reads need to be quality filtered to ensure
343 high quality reads for assembly. FastQC
344 (<http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>) is one such App already
345 available through CyVerse, and provides a visualization of read lengths, quality scores,
346 duplicate reads, N and GC content of raw reads that can be used to determine the
347 appropriate parameters for quality control. Once quality-filtering parameters have been
348 determined for a given sequencing run via FastQC, reads files can be trimmed and quality
349 controlled by using the Trimmomatic App.

350 **Step 3: Assembly of QC Read Data using SPAdes.**

351 Protocol available at protocols.io: <https://www.protocols.io/view/Assembling->
352 [Viral-Metagenomic-Data-with-SPAdes-Cyve-evzbe76](https://www.protocols.io/view/Assembling-Viral-Metagenomic-Data-with-SPAdes-Cyve-evzbe76)

353 Following QC, reads are then assembled using one of the assemblers available in
354 CyVerse. Most frequently, assembler selection is based on read type (Sanger, 454,
355 Illumina, PacBio, etc) and to a lesser extent, its performance for a particular sample type.
356 IDBA-UD, SOAPDenovo, Trinity and SPAdes are available for viral metagenomic
357 assembly. Some assemblers have high memory variants for larger data sets, and should
358 only be used when the standard versions fail to assemble.

359 **Step 4: Identification of Viral Sequences from Assembled Data using VirSorter.**

360 Protocol available at protocols.io: <https://www.protocols.io/view/Identifying->
361 [Viral-Sequences-Using-VirSorter-Cyvers-ev2be8e](https://www.protocols.io/view/Identifying-Viral-Sequences-Using-VirSorter-Cyvers-ev2be8e)

362 To identify viral sequences from the contigs file, VirSorter is used. Generally,
363 contigs larger than 3-kb can be successfully used as input – and can be single cell
364 genomes, microbial or viral metagenomes, fragmented or complete genomes.

365 **Step 5: Characterizing Viral Sequence Data through Protein Clustering with**
366 ***PCPipe and vContact.***

367 Protocol available at protocols.io: [https://www.protocols.io/view/Preparing-Data-](https://www.protocols.io/view/Preparing-Data-for-vContact-from-Proteins-Cyverse-ev7be9n)
368 [for-vContact-from-Proteins-Cyverse-ev7be9n](https://www.protocols.io/view/Preparing-Data-for-vContact-from-Proteins-Cyverse-ev7be9n)

369 Protocol available at protocols.io: [https://www.protocols.io/view/Applying-](https://www.protocols.io/view/Applying-vContact-to-Viral-Sequences-and-Visualizi-ev8be9w)
370 [vContact-to-Viral-Sequences-and-Visualizi-ev8be9w](https://www.protocols.io/view/Applying-vContact-to-Viral-Sequences-and-Visualizi-ev8be9w)

371 Large-scale characterization of viral genomic data remains one of the most
372 daunting challenges in viral ecology. A relatively recent method of analyzing complex
373 viral data is through *organizing* viral sequence space using PCs (protein clusters),
374 reducing the problems associated with data complexity as a byproduct. Regardless of the
375 type of analysis, iVirus has access to a number of tools to characterize viral data.

376

377 **Concluding Remarks**

378 While viruses are increasingly recognized for their roles in microbial-dominated
379 ecosystems, they remain understudied, particularly due to challenges stemming from the
380 lack of centralized viral metagenomic resources. iVirus offers a community-focused
381 resource, built on the CyVerse cyberinfrastructure and designed to directly address the
382 challenges of viral ecology in the era of next-generation sequencing, high performance
383 computing and big data analytics. This is done through 1) leveraging CyVerse's Data
384 Store to provide large data storage capacity and a centralized location for collecting data,

385 2) developing Apps, or software applications designed to take advantage of HPC
386 resources that require limited bioinformatics training on part of the researcher, 3)
387 collecting viral datasets in the iVirus Data Commons to provide a centralized location for
388 discovering datasets via environmental metadata and collaborating within the field, and
389 4) positioning these resources to maximize community exposure and feedback through
390 extensive and ‘live’ documentation at protocols.io.

391

392 **Acknowledgements**

393 The authors wish to thank Joanne B. Emerson, Dean R. Vik, Gareth G. Trubl,
394 Consuelo Gazitua, Pilar Manrique, and Jacob H. Munson-McGee for their helpful
395 feedback in designing community-minded tools; Matthew Vaughn from the Texas
396 Advanced Computer Center (TACC) and Nirav Merchant from the BIO5 Institute whose
397 assistance in App development, troubleshooting and making Apps publicly available
398 were invaluable. We thank Lenny Teytelman, Lori Kindler, and Alexei Stoliartchouk for
399 their agile development of protocols.io to accommodate requests for the iVirus group.
400 The development of iVirus and this publication was partially funded in part by Gordon
401 and Betty Moore Foundation grants GBMF4491, GBMF4733, GBMF3305, GBMF3790,
402 by the US Department of Energy Office of Biological and Environmental Research under
403 the Genomic Science program (Award DE- SC0010580) and by a TRIF award to the
404 University of Arizona Ecosystem Genomics Initiative.

405

406 **Conflict of Interest**

407 The authors declare no conflict of interest.

408

409 **References**

- 410 Benson DA, Boguski MS, Lipman DJ, Ostell J, Francis BF. (1998). GenBank. *Nucleic*
411 *Acids Res* **26**: 1–7.
- 412 Breitbart M, Salamon P, Andresen B, Mahaffy JM, Segall AM, Mead D, *et al.* (2002).
413 Genomic analysis of uncultured marine viral communities. *Proc Natl Acad Sci* **99**:
414 14250–14255.
- 415 Brum JR, Hurwitz BL, Schofield O, Ducklow HW, Sullivan MB. (2016). Seasonal time
416 bombs: dominant temperate viruses affect Southern Ocean microbial dynamics.
417 *ISME J* **10**: 437–449.
- 418 Brum JR, Ignacio-Espinoza JC, Roux S, Doucier G, Acinas SG, Alberti A, *et al.* (2015).
419 Patterns and ecological drivers of ocean viral communities. *Science (80-)* **348**:
420 1261498–1261498.
- 421 Caporaso JG, Kuczynski J, Stombaugh J, Bittinger K, Bushman FD, Costello EK, *et al.*
422 (2010). QIIME allows analysis of high-throughput community sequencing data.
423 *Nat Methods* **7**: 335–336.
- 424 Caporaso JG, Lauber CL, Walters W a, Berg-Lyons D, Huntley J, Fierer N, *et al.* (2012).
425 Ultra-high-throughput microbial community analysis on the Illumina HiSeq and
426 MiSeq platforms. *ISME J* **6**: 1621–4.
- 427 Cesar Ignacio-Espinoza J, Solonenko SA, Sullivan MB. (2013). The global virome: Not
428 as big as we thought? *Curr Opin Virol* **3**: 566–571.
- 429 Crits-Christoph A, Gelsinger DR, Ma B, Wierchos J, Ravel J, Davila A, *et al.* (2016).
430 Functional interactions of archaea, bacteria and viruses in a hypersaline endolithic
431 community. *Environ Microbiol* n/a–n/a.
- 432 Djikeng A, Kuzmickas R, Anderson NG, Spiro DJ. (2009). Metagenomic analysis of
433 RNA viruses in a fresh water lake. *PLoS One* **4**: e7264.
- 434 Duhaime MB, Sullivan MB. (2012). Ocean viruses: Rigorously evaluating the
435 metagenomic sample-to-sequence pipeline. *Virology* **434**: 181–186.
- 436 Edwards RA, Rohwer F. (2005). Viral metagenomics. *Nat Rev Microbiol* **3**: 504–510.
- 437 Enright a J, Van Dongen S, Ouzounis C a. (2002). An efficient algorithm for large-scale
438 detection of protein families. *Nucleic Acids Res* **30**: 1575–84.
- 439 Goff S a., Vaughn M, McKay S, Lyons E, Stapleton AE, Gessler D, *et al.* (2011). The
440 iPlant Collaborative: Cyberinfrastructure for Plant Biology. *Front Plant Sci* **2**: 1–
441 16.
- 442 Hannigan GD, Meisel JS, Tyldsley AS, Zheng Q, Hodgkinson BP, SanMiguel AJ, *et al.*
443 (2015). The Human Skin Double-Stranded DNA Virome: Topographical and
444 Temporal Diversity, Genetic Enrichment, and Dynamic Associations with the
445 Host Microbiome. *MBio* **6**: e01578–15.
- 446 Hoff PD. (2005). Bilinear Mixed-Effects Models for Dyadic Data. *J Am Stat Assoc* **100**:
447 286–295.
- 448 Hurwitz BL, Brum JR, Sullivan MB. (2014). Depth-stratified functional and taxonomic
449 niche specialization in the “core” and “flexible” Pacific Ocean Virome. *ISME J*
450 1–13.
- 451 Hurwitz BL, Hallam SJ, Sullivan MB. (2013). Metabolic reprogramming by viruses in
452 the sunlit and dark ocean. *Genome Biol* **14**: R123.

- 453 Hurwitz BL, Sullivan MB. (2013). The Pacific Ocean Virome (POV): A Marine Viral
454 Metagenomic Dataset and Associated Protein Clusters for Quantitative Viral
455 Ecology Thompson, F (ed). *PLoS One* **8**: e57355.
- 456 Hurwitz BL, Westveld a. H, Brum JR, Sullivan MB. (2014). Modeling ecological drivers
457 in marine viral communities using comparative metagenomics and network
458 analyses. *Proc Natl Acad Sci* **111**: 10714–10719.
- 459 Kanz C, Aldebert P, Althorpe N, Baker W, Baldwin A, Bates K, *et al.* (2005). The EMBL
460 nucleotide sequence database. *Nucleic Acids Res* **33**: 29–33.
- 461 Lima-Mendez G, Faust K, Henry N, Decelle J, Colin S, Carcillo F, *et al.* (2015).
462 Determinants of community structure in the global plankton interactome. *Science*
463 (80-) **348**: 1262073_1–1262073_9.
- 464 Lima-Mendez G, Van Helden J, Toussaint A, Leplae R. (2008). Reticulate representation
465 of evolutionary and functional relationships between phage genomes. *Mol Biol*
466 *Evol* **25**: 762–777.
- 467 Marcais G, Kingsford C. (2011). A fast, lock-free approach for efficient parallel counting
468 of occurrences of k-mers. *Bioinformatics* **27**: 764–770.
- 469 Minot S, Bryson A, Chehoud C, Wu GD, Lewis JD, Bushman FD. (2013). Rapid
470 evolution of the human gut virome. *Proc Natl Acad Sci* **110**: 12450–12455.
- 471 Munson-McGee JH, Field EK, Bateson M, Rooney C, Stepanauskas R, Young MJ.
472 (2015). Nanoarchaeota, their Sulfolobales Host, and Nanoarchaeota virus
473 distribution across Yellowstone National Park hot springs. *Appl Environ*
474 *Microbiol* **81**: 7860–7868.
- 475 Rattei T, Tischler P, Götz S, Jehl MA, Hoser J, Arnold R, *et al.* (2009). SIMAP-A
476 comprehensive database of pre-calculated protein sequence similarities, domains,
477 annotations and clusters. *Nucleic Acids Res* **38**: 223–226.
- 478 Reuter JA, Spacek D V., Snyder MP. (2015). High-Throughput Sequencing
479 Technologies. *Mol Cell* **58**: 586–597.
- 480 Roux S, Enault F, Hurwitz BL, Sullivan MB. (2015). VirSorter: mining viral signal from
481 microbial genomic data. *PeerJ* **3**: e985.
- 482 Roux S, Enault F, Robin A, Ravet V, Personnic S, Theil S, *et al.* (2012). Assessing the
483 Diversity and Specificity of Two Freshwater Viral Communities through
484 Metagenomics Martin, DP (ed). *PLoS One* **7**: e33641.
- 485 Roux S, Hallam SJ, Woyke T, Sullivan MB. (2015). Viral dark matter and virus-host
486 interactions resolved from publicly available microbial genomes. *Elife* **4**: e08490.
- 487 Roux S, Tournayre J, Mahul A, Debroas D, Enault F. (2014). Metavir 2: new tools for
488 viral metagenome comparison and assembled virome analysis. *BMC*
489 *Bioinformatics* **15**: 76.
- 490 Seshadri R, Kravitz S a., Smarr L, Gilna P, Frazier M. (2007). CAMERA: A community
491 resource for metagenomics. *PLoS Biol* **5**: 0394–0397.
- 492 Sharon I, Battchikova N, Aro E-M, Giglione C, Meinel T, Glaser F, *et al.* (2011).
493 Comparative metagenomics of microbial traits within oceanic viral communities.
494 *ISME J* **5**: 1178–1190.
- 495 Solonenko SA, Ignacio-Espinoza JC, Alberti A, Cruaud C, Hallam S, Konstantinidis K,
496 *et al.* (2013). Sequencing platform and library preparation choices impact viral
497 metagenomes. *BMC Genomics* **14**: 320.
- 498 Stern A, Mick E, Tirosh I, Sagy O, Sorek R. (2012). CRISPR targeting reveals a reservoir

499 of common phages associated with the human gut microbiome. *Genome Res* **22**:
500 1985–1994.

501 Wheeler DL, Church DM, Federhen S, Lash AE, Madden TL, Pontius JU, *et al.* (2003).
502 Database resources of the national center for biotechnology. *Nucleic Acids Res*
503 **31**: 28–33.

504 Wilke A, Bischof J, Gerlach W, Glass E, Harrison T, Keegan KP, *et al.* (2015). The MG-
505 RAST metagenomics database and portal in 2015. *Nucleic Acids Res* **44**:
506 gkv1322.

507 Wommack KE, Bhavsar J, Polson SW, Chen J, Dumas M, Srinivasiah S, *et al.* (2012).
508 VIROME: a standard operating procedure for analysis of viral metagenome
509 sequences. *Stand Genomic Sci* **6**: 427–439.

510 Yooseph S, Sutton G, Rusch DB, Halpern AL, Williamson SJ, Remington K, *et al.*
511 (2007). The Sorcerer II global ocean sampling expedition: Expanding the universe
512 of protein families. *PLoS Biol* **5**: 0432–0466.

513

514 **Figure Legends**

515

516 Figure 1. An organizational overview of how a user might leverage iVirus, iMicrobe, and

517 CyVerse Apps to analyze a viral metagenomic dataset. Arrows indicate direction of use

518 case scenario. Processing stages are represented by text in blocks, with bolded text

519 indicating the Apps used in the use case scenario.

