# vConTACT: an iVirus tool to classify viruses that infect *Archaea* and *Bacteria*

Benjamin Bolduc [1] , Ho Bin Jang [1] , Guilhem Doulcier [2] , Zhi-Qiang You [3] , Simon Roux [1] , Matthew B Sullivan
Corresp. [1, 4]

[1] Department of Microbiology, Ohio State University, Columbus, Ohio, United States

[2] Department of Biology École Normale Supérieure, PSL Research University, Paris, France

[3] Department of Chemistry and Biochemistry, Ohio State University, Columbus, Ohio, United States

[4] Department of Civil, Environmental and Geodetic Engineering, Ohio State University, Columbus, Ohio, United States

Corresponding Author: Matthew B Sullivan
Email address: mbsulli@gmail.com

Taxonomic classification of archaeal and bacterial viruses is challenging, yet also fundamental for developing a predictive understanding of microbial ecosystems. Recent identification of hundreds of thousands of new viral genomes and genome fragments, whose hosts remain unknown, require a paradigm shift away from traditional classification approaches and towards the use of genomes for taxonomy. Here we revisited the use of genomes and their protein content as a means for developing a viral taxonomy for bacterial and archaeal viruses. A network-based analytic was optimized and benchmarked against authority-accepted taxonomic assignments and found to be largely concordant. Exceptions were manually examined and found to represent areas of viral genome 'sequence space' that are under-sampled or prone to excessive gene flow. While both cases are poorly resolved by genome-based taxonomic approaches, the former will improve as viral sequence space is better sampled and the latter are uncommon. Finally, given the largely robust taxonomic capabilities of this approach, we sought to enable researchers to easily and systematically classify new viruses. Thus, we established a tool, vConTACT, as an app at iVirus, where it operates as a fast, highly scalable, user-friendly app within the free and powerful CyVerse cyberinfrastructure.

1 **vConTACT: an iVirus tool to classify viruses that infect *Archaea* and *Bacteria***

2  Benjamin Bolduc[1&], Ho Bin Jang[1&], Guilhem Doulcier[3,4], Zhi-Qiang You[5], Simon Roux[1] &
3  Matthew B. Sullivan*[1,2]

4  [1]Department of Microbiology, The Ohio State University, Columbus, OH 43210

5  [2]Department of Civil, Environmental and Geodetic Engineering, The Ohio State University,
6  Columbus, OH 43210

7  [3]École normale supérieure, PSL Research University, IBENS, F-75005, Paris, France.

8  [4]ESPCI, PSL Research University, CBI, F-75005, Paris, France.

9  [5]Department of Chemistry and Biochemistry, The Ohio State University, Columbus, OH 43210

10  [&] These authors contributed equally to this work.

11

12  Corresponding Author:

13  Matthew B. Sullivan[1,2]

14  Riffe Building Rm 914, 496 W 12[th] Ave Columbus, OH 43210, USA

15  Email address: mbsulli@gmail.com

16

17

**Abstract.**

18

19      Taxonomic classification of archaeal and bacterial viruses is challenging, yet also

20   fundamental for developing a predictive understanding of microbial ecosystems. Recent

21   identification of hundreds of thousands of new viral genomes and genome fragments, whose

22   hosts remain unknown, require a paradigm shift away from traditional classification approaches

23   and towards the use of genomes for taxonomy. Here we revisited the use of genomes and their

24   protein content as a means for developing a viral taxonomy for bacterial and archaeal viruses. A

25   network-based analytic was optimized and benchmarked against authority-accepted taxonomic

26   assignments and found to be largely concordant. Exceptions were manually examined and found

27   to represent areas of viral genome 'sequence space' that are under-sampled or prone to excessive

28   gene flow. While both cases are poorly resolved by genome-based taxonomic approaches, the

29   former will improve as viral sequence space is better sampled and the latter are uncommon.

30   Finally, given the largely robust taxonomic capabilities of this approach, we sought to enable

31   researchers to easily and systematically classify new viruses. Thus, we established a tool,

32   vConTACT, as an app at iVirus, where it operates as a fast, highly scalable, user-friendly app

33   within the free and powerful CyVerse cyberinfrastructure.

34

**Introduction.**

35

36      Classification of viruses that infect Archaea and Bacteria remains challenging in

37   virology. Official viral taxonomy is handled by the International Committee for the Taxonomy of

38   Viruses (ICTV) and organizes viruses into order, family, subfamily, genus and species.

39   Historically, this organization derives from numerous viral features, such as morphology,

40   genome composition, segmentation, replication strategies and amino- and nucleic-acid

41   similarities – all of which is thought to roughly organize viruses according to their evolutionary

42   histories (Simmonds, 2015). As of 2015, the latest report issued, the ICTV has classified 7

43   orders, 111 families, 27 subfamilies, 609 genera and 3704 species

44   (http://ictvonline.org/virusTaxInfo.asp).

45          Problematically, however, current ICTV classification procedures cannot keep pace with

46   viral discovery and may need revision where viruses are not brought into culture. For example,

47   of the 4400 viral isolate genomes deposited into National Center for Biotechnology information

48   (NCBI) viral RefSeq, only 43% had been ICTV-classified by 2015. This is because the lengthy

49   'proposal' processes lags deposition of new viral genomes, in some cases for years (Fauquet &

50   Fargette, 2005). Concurrently, new computational approaches are providing access to viral

51   genomes and large genome fragments at unprecedented rates. One approach mines microbial

52   genomic datasets to provide virus sequences where the host is known – already adding 12,498

53   new prophages from publicly available bacterial and archaeal microbial genomes (Roux et al.,

54   2015a) and 89 (69 and 20, respectively) new virus sequences from single cell amplified genome

55   sequencing projects (Roux et al., 2014; Labonté et al., 2015). A second approach assembles viral

56   genomes and large genome fragments from metagenomics datasets. These studies have added

57   15,222 and 125,842 new virus sequences from oceanic viral metagenomes (Tara Oceans Virome,

58   or TOV (Brum et al., 2015b) and Global Oceans Virome, or GOV (Roux et al., 2016)), and from

59   microbial and viral metagenomes from a diversity of ecosystems (Paez-Espino et al., 2016),

60   respectively. A third approach leverages large-insert cloning and sequencing approaches to

61   identify new viral sequences – with 208 from a single seawater sample (Mizuno et al., 2013).

62   Such new virus genomes and large genome fragments will keep coming for the foreseeable

63    future and represent an incredible resource for viral ecology, but they also represent a daunting

64    challenge for taxonomy.

65        Currently such rapidly expanding genomic databases of the virosphere remain

66    unclassified and challenging to integrate into a systematic framework for three reasons. First,

67    viruses lack a universal marker gene, which prevents the taxonomic starting place that is so

68    valuable for microbes (Woese, Kandler & Wheelis, 1990). Second, though genomes and large

69    genome fragments are now much more readily available, researchers are reticent to use genomes

70    as a basis for taxonomy as a paradigm has emerged whereby viruses are rampantly mosaic and

71    therefore must exist as part of a genomic continuum such that any clustering in 'sequence space'

72    is an artifact of sampling. This is most well studied in the many genomes of mycobacteriophages

73    (Pope et al., 2015), but is contrasted by observations in cyanophages where efforts have been

74    made to more deeply sample variability in a single site with findings suggesting clear population

75    structure for naturally-occurring cyanophages (Deng et al., 2014) and that cyanophage

76    populations appear to fit a population genetics-based species definition (Marston & Amrich,

77    2009; Gregory et al., 2016). It is possible that gene flow in fast evolving RNA and ssDNA

78    viruses is rampant, but that slower evolving dsDNA viruses, particularly if obligately lytic rather

79    than temperate, could evolve into relatively stable populations that could be the basis of

80    taxonomy. Thus, it remains unclear whether viral genomes can serve as the sole basis for

81    taxonomy, or whether exploration of available data could help identify areas of viral genome

82    sequence space that are amenable to taxonomic 'rules' and others that are not.

83        Despite these challenges, numerous reference-independent, automated, genome-based

84    classification schemes have been proposed. An early effort recognized that more genes are

85    shared within related virus groups than between them (Lawrence, Hatfull & Hendrix, 2002),

86  which led to virologists grappling with natural diversity to use translated genomes as the basis of

87  whole genome phylogenomic tree classifications – e.g., the Phage Proteomic Tree (Edwards &

88  Rohwer 2002). Simulations showed this method to be very accurate for assigning fragmented

89  reads to the correct genomes (Edwards & Rohwer, 2005) but it suffers from the availability of

90  phage genomes. A second approach that has emerged for relatively well-studied virus groups, is

91  to use pairwise distances between aligned sequences to identify discontinuities that can indicate

92  classification thresholds. Two tools – Pairwise Sequence Comparison (PASC; Bao 2014) and

93  DEmARC (Lauber & Gorbalenya, 2012a) – are available to align the sequences either in context

94  of previous knowledge of the taxonomic affiliations (PASC) or naively using pairwise

95  distributions (DEmARC). These tools have worked well for several virus families, such as the

96  *Picornaviridae* (Lauber & Gorbalenya, 2012a) and *Filoviridae* (Lauber & Gorbalenya, 2012b)

97  for DEmARC, and *Bornaviridae* and arenavirus for PASC (Kuhn et al., 2014; Radoshitzky et al.,

98  2015). However, DEmARC and PASC suffer from several issues: (i) they are not generalizable

99  to the coming deluge of environmental viral genome sequences as they require *a priori* expert

100  knowledge to impose similarity thresholds at each level, (ii) ICTV subcommittees have

101  established varied sequence similarity thresholds across viral groups (Simmonds, 2015), which

102  would require a sliding threshold, and (iii) the methods can only classify sequences that are

103  similar to database references (Zanotto et al., 1996), which for the oceans at least represents <1%

104  of the viral genomes recovered (Brum et al., 2015a).

105      Complementarily, two network-based approaches have been utilized to organize virus

106  genome sequence space in a manner that enables classification without *a priori* knowledge. The

107  first, a gene sharing network (Lima-Mendez et al., 2008), predicts viral genes in all the genomes,

108  translates them into proteins, organizes these proteins into MCL-based protein families (protein

109    clusters, "PCs"), evaluates the number of shared protein clusters pairwise throughout the dataset

110    to establish a protein profile, and then represents this information as a weighted graph, with

111    nodes representing viral genomes and edges the similarity score of their shared protein content.

112    Given the 306 bacterial viruses (phages) known at the time, this method was precise as it

113    correctly placed 92% and 95% of these phages into their correct ICTV genus or family,

114    respectively (Lima-Mendez et al., 2008). The second, a bipartite network (Iranzo, Krupovic &

115    Koonin, 2016), incorporates both gene sharing as above and genomic similarity. In this work, all

116    dsDNA viruses along with mobile genetic elements were analyzed, which revealed a module-

117    based structure to the dsDNA virosphere. These two studies imply that even very distantly

118    related viruses can be organized into discrete populations by genomes alone and that there may

119    be hope for automated, genome-based viral taxonomy, at least for dsDNA viruses.

120          Here we optimize gene sharing networks and re-evaluate their efficacy for recapitulating

121    ICTV-based classifications using an expanded dataset of 2,010 bacterial and archaeal virus

122    genomes (available as of RefSeq v75), while also deeply exploring where network-based

123    methods have lower resolution and/or yield discontinuities with currently established

124    taxonomies. Further, we make these approaches accessible to researchers by developing a tool,

125    vConTACT (Viral CONTigs Automatic Clustering and Taxonomy), and deploy it as part of the

126    iVirus ecosystem of apps that leverages the CyVerse cyberinfrastructure (Bolduc et al., 2016).

127

128    **Materials and Methods.**

129          **Terminology**. Network topological parameters, their definitions and abbreviations are

130    available in Table 1.

131     **Reference datasets**. To test this methodology, we downloaded the entire NCBI viral

132     reference dataset ("ViralRefSeq", version 75, containing 5539 viruses) and removed eukaryotic

133     viruses by filtering against tables downloaded on NCBI's ViralRefSeq viral genome page

134     (http://www.ncbi.nlm.nih.gov/genomes/GenomesGroup.cgi?taxid=10239). The resulting file

135     ("Bacterial and Archaeal viruses"; BAV) contained 2010 total viruses; 1905 dsDNA, 88 ssDNA,

136     5 dsRNA and 12 ssRNA. All viruses contained taxonomic affiliation information, though not all

137     viruses had affiliations associated with each level of the taxonomy (i.e. not all viruses have a

138     "sub-family" designation). To improve taxonomic assignments, the ICTV taxonomy was also

139     retrieved (https://talk.ictvonline.org/files/master-species-lists/) and the ICTV affiliations were

140     used to supplement the NCBI version.

141     **Building protein cluster profiles**. To generate sequence profiles – with information

142     about the presence or absence of a sequence within one or more protein clusters (described

143     previously as protein *families* (Lima-Mendez et al., 2008)), proteins from each sequence were

144     first extracted from the ViralRefSeq proteins file. BLASTP (Altschul et al., 1997) was used to

145     compare all proteins (198,102) from the sequences in an all-verses-all pairwise comparison.

146     Protein clusters were subsequently identified using the Markov clustering algorithm (MCL) with

147     an inflation value of 2, resulting in 23,022 protein clusters ("PCs"). Finally, we generated protein

148     cluster profiles for each genome such that the presence of a gene within a protein cluster of a

149     viral genome was given a value of "1" and the absence "0". This resulted in a large 2,010 x

150     23,022 matrix.

151     **Generating the similarity network**. The similarity network is a graph where the nodes

152     (i.e. reference sequences) are linked by edges when the similarity between their pc-profiles is

153     considered significant. In other words, the network represents the overall similarity between

154    sequences based on the number of shared protein clusters. To calculate the similarity between the

155    profiles of two sequences (sequence A and sequence B), the hypergeometric formula was used to

156    estimate the probability that at least $c$ protein clusters would be in common:

$$P(X \geq c) = \sum_{i=c}^{\min(a,b)} \frac{C_a^i C_{n-a}^{b-i}}{C_n^b} \qquad (1)$$

157

158    Simply stated, the hypergeometric formula is the probability that genomes A and B would have $c$

159    protein clusters in common by chance. The probability can be converted to an expectation value

160    (E; for false positives) by multiplying the probability (P) by the total number of comparisons (T).

161    The expectation value can then be converted into a significance score:

$$S(A,B) = -\log(E) = -\log(P \times T) \qquad (2)$$

162

163    Genome pairs with significance scores greater than 1 (i.e. E-value $< 0.1$) are considered

164    sufficiently similar (see permutation test, below) and were joined by an edge in the similarity

165    network with a weight equal to their significance score. We refer to sequences within the

166    network as *nodes*, the relationships connecting them, *edges* and the strength of that relationship,

167    edge *weight*.

168        After generating the similarity network, groups of similar sequences (referred to as viral

169    clusters, "VCs") were clustered by applying MCL with an inflation of 2.

170        **Measuring the proportion of shared genes between genomes.** Given that genome

171    sizes between pairs can differ greatly, this can lead to large differences in the proportion of the

172    shared genes (Ågren et al., 2012). To counter this, we characterized the proportion of shared PCs

173    between two genomes using the geometric index (G) as a symmetric index:

$$G_{AB} = \frac{|N(A) \cap N(B)|}{|N(A)| \times |N(B)|} \qquad (3)$$

174

175    where *N(A)* and *N(B)* indicate the numbers of protein clusters (PCs) in the genomes of A and B,

176    respectively.

177        **Permutation test.** The stringency of the significant score was evaluated through

178    randomization of the original matrix where rows present viral genomes and columns PCs or

179    singletons that are not shared with any other protein sequences (Leplae et al., 2004). Briefly,

180    with an in-house R script, 1,000 matrices were generated by randomly rearranging PCs and/or

181    singletons within pairs of genomes having the significant score ≤ 1 (a negative control) and

182    calculated a new significant score each time. None of the genome pairs in this negative control

183    produced significant scores >1, indicating values above this significance threshold did not occur

184    by chance (Lima-Mendez et al., 2008).

185        **Affiliating sequence clusters with taxonomic groups**. To assign (in the case of

186    unknown sequences) or compare nodes (genomes) within clusters to their reference counterparts,

187    we first defined *membership* of a node *c* to a cluster *k* $B(c,k)$ according to two methods,

188    conservative and permissive. The conservative method 4) directly takes the result from the MCL

189    clustering to assign a node to a cluster:

$$B(c,k) = \begin{cases} 1 \text{ if Contig } c \in \text{Cluster } k, \\ 0 \qquad\qquad \text{otherwise} \end{cases} \tag{4}$$

191    while the permissive method takes the sum of all edge weights *w* linking the node to nodes of the

192    cluster, with the node becoming a member of its maximal membership cluster (5):

$$B'(c,k) = \frac{\sum_{i \in k} w_{c,i}}{\sum_{p \in \{\text{Clusters}\}} \sum_{j \in p} w_{g,j}} \tag{5}$$

194    The precision $P(k,t)$ of the taxonomic class *t* with respect to a cluster *k* as its membership of

195    nodes of the class *t* in the memberships of reference nodes in the cluster *k*.

196
$$P(k,t) = \frac{\sum_{\forall i \in \{\text{sequence of class } t\}} B(i,k)}{\sum_{\forall j \in \{\text{reference sequence}\}} B(j,k)} \qquad (6)$$

197  A cluster and all its node members are then affiliated with its maximal precision class. For the

198  conservative method, the cluster is affiliated with the taxonomic class associated with the

199  majority of its members (i.e. $\geq 50\%$). In cases where clusters do not contain at least half

200  reference sequences, the entire cluster will be unaffiliated.

201  **Measuring the connectivity of genomes to clusters.** The connection strength of a node

202  $g$ to cluster $c$ was calculated as the average edge weight linking it to nodes of cluster $c$:

203
$$W_{g,c} = \frac{1}{k}\sum_{i=1}^{k} wg,i \qquad (7)$$

204  where $k$ and $w$ are the number and total weight of edges of the node $g$ in the cluster $c$,

205  respectively. We refer to the average edge weight for node $g$ to the cluster it belongs to as its in-

206  VC average weight, and to other clusters within the network as out-VC average weight.

207  **Identifying sub-clusters**. To further subdivide heterogeneous clusters (those comprising

208  $\geq 2$ taxa), cluster-wise module profiles (i.e. a module profile only including viruses previously

209  identified as belonging to the same viral cluster) were hierarchically clustered using UPGMA

210  with pairwise distances calculated using Euclidean distance implemented in Scipy.

211  **Statistical calculations**. All calculations, statistics, network statistical analyses were

212  performed using in-house python scripts, with the Numpy, Scipy, Biopython and Pandas python-

213  packages. vConTACT is implemented in python with the same dependencies. The tool is

214  available at https://bitbucket.org/MAVERICLab/vcontact. Scripts used in the generational and

215  calculations of data are available at https://bitbucket.org/MAVERICLab/vcontact-SI.

216 **Network visualization and analysis**. The network was visualized with Cytoscape

217 (version 3.1.1; http://cytoscape.org/), using an edge-weighted spring embedded model, which

218 places the genomes or fragments sharing more PCs closer to each other. Topological properties

219 were estimated using a combination of python and the Network Analyzer 2.7 Cytoscape plug-in

220 (Assenov et al., 2008).

221

222 **Results and Discussion.**

223 ***vConTACT analytical workflow and terminology*:** The vConTACT analyses are based

224 on previously established gene sharing network methods (Lima-Mendez 2008). Briefly, PCs are

225 established across all genomes in the dataset; with vConTACT doing this by default using MCL

226 clustering from all-verses-all BLASTP comparisons (though user-specified clusters can also be

227 used). PC *profiles* of genomes or genome fragments (herein 'genome') are then calculated,

228 where the presence and absence of PCs (from the entire PC dataset) along a genome are

229 established and then compared pairwise between genomes (Fig. 1). The pairwise genome

230 comparisons are then mathematically adjusted (using the hypergeometric similarity formula) to

231 establish a probability that any genome pair would share $n$ PCs, given the total number of all

232 PCs. This probability is log-transformed (in similar fashion to BLAST E-values) into a

233 significance score and applied as a *weight* to an edge between the two paired genomes in a

234 similarity network. High significance scores represent a low probability that two genomes would

235 share $n$ PCs by chance, which can be interpreted as evidence of gene-sharing and presumably

236 evolutionary relatedness between the paired genomes. After evaluating all pairings in the dataset,

237 significance scores $\geq 1$ are retained, and a network of the remaining genome pairs is constructed.

238 MCL is subsequently applied to identify structure in the gene sharing network, but now the

239  clusters represent groups or related genomes and are termed *viral* clusters ("VCs"). MCL is also

240  applied against the network of PCs, whose members can be similar to members of other PCs.

241  This effectively organizes the PCs into a higher-order structure known as a protein module. The

242  relationship information identified from the genomes (organized into VCs) and PCs (organized

243  into protein modules) are used to create a *module profile*, which can then be mined for

244  taxonomic identification, functional profiling, etc.

245       ***Benchmarking network-based taxonomy:*** To benchmark the ability of network-based

246  taxonomy to capture 'known' viral relationships, we evaluated how vConTACT "re-classified"

247  viral sequences at various taxonomic levels using 2,010 bacterial and archaeal viral genomes

248  from VirRefSeq (v75). Of these reference genomes, ICTV-classifications were only available for

249  a subset; 654 viruses from 2 orders, 738 viruses from 19 families, 152 viruses from 11

250  subfamilies, and 562 viruses from 158 genera. The network was then decomposed into VCs

251  (described above) and a permutation test was used to establish significance score thresholds to

252  prevent random relationships from entering the network. This analysis used the initial network's

253  edge information to construct a matrix between genome pairs, and then permuted the edges 1,000

254  times. No edges were found to be significant during these tests, suggesting that relationships seen

255  within the network did not arise by chance and could be confidently used to establish taxonomic

256  groupings (see Materials and Methods, Table S1).

257       The resulting network, consisting of 1,964 viruses (nodes) and 65,393 relationships

258  (edges) between them (Fig. 2A), was then used as a basis for comparison to the ICTV-based

259  classifications. A total of 211 VCs were identified, spread among 46 components (unconnected

260  subnetworks), which more than doubles the 17 connected components identified previously

261  (Lima-Mendez et al., 2008). Of the 46 components, 38 included 1,891 phages representing 194

262    VCs (left, Fig. 2A), and 8 components included 73 archaeal viruses representing 17 VCs (right,

263    Fig. 2A). Most (87%) of the 1,891 phages belonged to the order *Caudovirales*, and comprised

264    the largest connected component (LCC) in the analysis (top left, Fig. 2A). At the VC level, the

265    network clustering performed well with average (across each taxonomic level) recall / precision

266    percentages of 100% / 100%, 90% / 86%, and 80% / 80% at the order, family and genus levels,

267    respectively (Fig. 2B). Of the 211 VCs resolved by the network, 76.4% contained a single ICTV-

268    accepted genus, suggesting concordance between the network VCs and accepted taxonomy,

269    whereas 15.1% and 8.5% of the VCs contained two and 3 or more genera, respectively (Fig. 2A

270    and C). Thus, roughly 4 out of 5 of the VCs correspond to ICTV genera.

271        Mechanistically, these discrepancies between network clustering and the ICTV

272    classification could derive from: (*i*) undersampling such that VCs with fewer members may not

273    represent the naturally-occurring diversity of that viral group, or (*ii*) gene flow between viral

274    genomes that ameliorates taxonomic boundaries by providing excessive access to genes outside

275    the VC's gene pool. While much of abundant viral genome sequence space is recently being

276    mapped (Paez-Espino et al., 2016; Roux et al., 2016), there remains contrasting paradigms about

277    the role of gene flow in structuring mycobacteriophage ("mycophages") vs cyanophage

278    populations (Gregory et al., 2016).

279        To discriminate between these possibilities, we identified these "ICTV-discordant" areas

280    of the network containing 2 or more ICTV genera (which we define as *heterogeneous VCs*),

281    focusing on three of the more well-populated (many member genomes) heterogeneous VCs, and

282    the archaeal virus heterogeneous VCs, which are among the least well-sampled taxa. Of the well-

283    sampled VCs, VCs containing the 2$^{nd}$, 3$^{rd}$, and 4$^{th}$ most members (i.e. genomes), included the

284    following: (i) VC1 contains the 8 genera belonging to the *Tevenvirinae* subfamily (*T4virus*,

285    *Cc31virus*, *Js98virus*, *Rb49virus*, *Rb69virus*, *S16virus*, *Sp18virus*, and *Schizot4virus*) and a

286    genus of the *Eucamyvirinae* (*Cp8virus*), as well as the *Tg1virus* and *Secunda5virus* that are not

287    assigned to a particular subfamily, (ii) VC2 contains three genera (*Biseptimavirus*, *Phietavirus*,

288    and *Triavirus*), and (iii) VC3 contains four genera (*Kayvirus*, *Silviavirus*, *Twortvirus*, and

289    *P100virus*) belonging to the *Spounavirinae* and the six *Bacillus* virus genera (*Agatevirus*,

290    *B4virus*, *Bc431virus*, *Bastillevirus*, *Nit1virus*, and *Wphvirus*). Finally, among the 73 archaeal

291    viruses, only the *Fuselloviridae* were accurately classified at the genus level, while most (63%)

292    archaeal viruses were incorrectly classified at the genus level.

293          ***Gene content analyses suggest ICTV classifications, not VC-based classifications***

294    ***should be revised***: A total of 23.6% of the VCs contained genomes from more than one ICTV-

295    recognized genus, which suggests 'lumping' by the network analyses (via MCL) or 'splitting'

296    during ICTV classification. To assess this, we computed the fraction of PCs that were shared

297    both within an ICTV genus and between the multiple ICTV genera found in each heterogeneous

298    VC and represented them as the percentage of intragenus similarity and intergenera similarity,

299    respectively. Of the 25 VCs, intragenus similarities of all but one (VC9) suggested they shared

300    more than 40% of their PCs (Fig. 3A, Table S2), which is consistent with the threshold

301    commonly used to define a new dsDNA viral genus (Lavigne et al., 2009). In contrast, the

302    intergenera similarities varied widely – some VCs (VCs 1-3, 9-11, 17, 20, 25, 33, 58, 91, 95)

303    shared 20-40% of their PCs (subfamily level), whereas others shared more than ~40% (VCs 12,

304    14, 24, 26, 37, 44, and 51) or less than ~20% (VCs 39, 55, 63, 74, and 77) of their PCs. Where

305    intergenera similarities are high (>40% of the PCs are shared), there may be a case to be made

306    for merging the currently recognized ICTV genera. Consistent with this, all 6 of these highly

307    (>40%) similar VCs (12, 14, 24, 26, 37 and 51) are suggested to be in need of revision, as these

308   include *G7Cvirus*, *N4virus*, *T1virus*, *HP34virus*, and *PhiKMVvirus* (Wittmann et al., 2015;

309   Eriksson et al., 2015; Niu et al., 2014; Krupovic et al., 2016). Additionally, we found that in

310   VC44, the phage CAjan, belonging to the *Seuratvirus*, shared 41.6-42.7% of its genes with three

311   phages (JenP1 and 2 and JenK1 of the *Nongaviru*s (Table S2)). Where intergenera similarities

312   are lower (<20%, or 20-40% of the PCs are shared), the appropriate taxonomic assignment may

313   require deeper sampling of viral genome sequence space and/or further network analytic

314   development.

315          To further assess these cases, we next examined four VCs (1-3, 14) that contained more

316   than 4 ICTV-recognized genera using hierarchical clustering of PC presence-absence data for

317   each genome (Fig. 3B). In parallel, we computed the actual connectivity of the genomes within

318   these heterogeneous VCs according to the average weight of edges that (i) are between genomes

319   of the same VC (in-VC avg. weight) and (ii) between the genomes of other VCs (out-VC avg.

320   weight) (Table S3; Materials and Methods). For example, within VC1, 8 genera of the

321   *Tevenvirinae* (*S16virus*, *Cc31virus*, *T4virus*, *Rb69virus*, *Sp18virus*, *Js98virus*, *Rb49virus* and

322   *Schizotvirus*) and their relatives (*Tg1virus* and *Secunda5virus*) share, on average, 61% and 38%

323   of their total PCs, respectively, and 39% between all 10 genera (Table S2). Outside VC1, they

324   share ~11.2% of genes with other viral groups (Table S2). We found that the 10 genera within

325   VC1 are more tightly interconnected than those of the 210 VCs overall, with average in-cluster

326   values of 223.7 and 131.9 and average out-cluster values of 13.1 and 9.0, respectively (Table

327   S3). These observations indicate that higher cross-similarities of 10 genera can be attributed to a

328   large fraction of their shared genes, whereas only a small fraction of gene shared by other groups

329   can hold them together.

330    Upon closer inspection, some of this 'lumping' appeared to be due to poorly sampled

331    regions of sequence space. For example, VC1 also contained the *Cp8virus* of the subfamily

332    *Eucampyvirinae*, which is odd to be placed alongside the *Tevevirinae*, given that other ICTV-

333    recognized genus (*Cp220virus*) of the *Eucampyvrinae* is grouped into the separate cluster (VC

334    87). Since both genera (*Cp8virus* and *Cp220virus*) are distantly related to the *Tevenvirinae*

335    (Javed et al., 2014), containing only ~11% the shared genes to (an average weight of 18.5) and

336    ~6% (11.8) , respectively (Tables S2 and S3), these groupings might be driven by the lack of

337    diversity among the *Cp220virus* where only 2 reference genomes (i.e., *Campylobacter* phages

338    CPX and NCTC12673) available in our ViralRefSeq dataset. To test this, we computationally

339    doubled the number of the genomes for this group while holding the number and weight of their

340    connections constant, finding that the *Cp220virus* genomes clearly separated from VC1 and

341    instead were correctly placed alongside VC 87 (Table S4). Consistently, among the

342    heterogeneous VCs 39, 55, 63, 74, and 77 showing < ~20% intergenera similarities (Figs. 3A and

343    S1), increasing the genome numbers of poorly-sampled ICTV genera led to clustering of

344    members of those genera into their correct VCs (Table S4). Together these findings suggest that

345    additional sampling in poorly sampled areas of viral sequence space will be required to most

346    accurately establish genome-based taxonomy – issues that parallel those presented by long

347    branch attraction for phylogenies (Bergsten, 2005).

348    Similar structure emerged from hierarchical clustering of PC presence / absence data

349    from the 3 other well-represented heterogeneous VCs. In VC2, the three known subgroups of the

350    *Phietavirus* (Gutiérrez et al., 2014) were resolved, sharing 44.9% of their PCs, and separate from

351    two other subgroups – the *Biseptimavirus* and *Triavirus*, which shared 22.3% of their PCs (Fig.

352    3B, Table S2). In VC3, containing the *Spounavirinae* (Krupovic et al., 2016), each sub-cluster

353 has a corresponding ICTV genus with largely overlapping sets of genes while also showing a

354 clearly distinct set(s) of genes. Of these, the six *Bacillus* virus genera (*Wphvirus*, *Bastillevirus*,

355 *B4virus*, *Bc431virus*, *Agatevirus*, and *Nit1virus*) appear to be closely related to the

356 *Spounavirinae*, with ~20% of total PCs in common (Fig. 3B, Table S2). Finally, VC14 produced

357 a clear division of the *Tunavirinae* (Krupovic et al., 2016), in which the *Escherichia* virus Jk06 is

358 placed in a separate branch due to its less shared common genes (~56%) to the other

359 *Rogue1virus* members (~82%); their highly-overlapped genes between genera above the genus

360 boundary (40%) are associated with "taxonomic lumping" as described above (Niu et al., 2014;

361 Krupovic et al., 2016).

362 We next evaluated three phage groups which were poorly represented in the S277

363 network (Lima-Mendez et al 2008) and also represent some of the most abundant, widespread,

364 and/or extensively studied phage groups (Grose & Casjens, 2014; Pope et al., 2015; Roux et al.,

365 2015b)) – the mycobacteriophages, *Teveniviriae*, *Autographivirinae* and the archaeal viruses.

366 ***Mycobacterium phages.*** The largest viral group covering 16.1% of the total population of

367 the LCC (mostly *Caudovirales*, top left Fig. 1A) includes phages infecting *Mycobacteria*. The

368 318 mycophage genomes were assigned to 14 VCs (Fig. 4A), 13 of which were composed of

369 reference genomes belonging to a single ICTV-recognized genus for each VC. The 14th

370 mycophage VC, VC25, contained three ICTV-recognized genera – the *Bignuzvirus*,

371 *Charlievirus*, and *Che9cvirus*. Although the module-based approach discerned the structure in

372 this VC, which would group them into the known genera (Fig. S1), this "lumping" into a single

373 VC reflects (*i*) their undersampling (i.e., each genus has 1 to at most 3 viruses) and/or (*ii*) highly-

374 overlapped genes between genera. Indeed, of the 3 phages belonging to the *Che9cvirus*, phages

375 Babsiella and Che9c shared 45% of their genes, but also shared 35% and 36% of their genes with

376    the *Bignuzvirus* and 28% and 32% with the *Charlievirus*, respectively (Table S2), which results

377    in higher connectivity between three genera than other viral groups to which they linked (Table

378    S3). These findings contrast those in the rest of the network, and suggest that some phage groups

379    (e.g., mycophages) may more frequently exchange genes than others.

380       To quantify this, we next examined features of the network. For example, many VC59

381    mycophages were broadly linked to nine VCs that contain other mycophages and phages from

382    diverse hosts (Fig. 4A). To characterize this further, we analyzed the topological properties using

383    the betweenness centrality (BC), as it can identify the node residing in the shortest path between

384    other nodes (Halary et al., 2009). Specifically, in the shared-gene network, high-betweenness

385    nodes (phages) can act as bridges between phages that would remain disconnected, due to their

386    mosaic content of genes (Lima-Mendez et al., 2008). Indeed, these eight VC 59 phages had 42-

387    fold higher average BC than those of other mycophages and their relatives (0.04 vs. 9.45E-04)

388    (Fig. S2), strongly indicating they may be prone to increased gene flow and thus exceptionally

389    'mosaic' genomes (Halary et al., 2009; Pope et al., 2015a).

390       ***The Tevenvirinae*** As the second-largest group, containing 94 viruses in the

391    heterogeneous VC1, which were further connected to 74 distant relatives and taxonomically

392    unclassified myo-/siphovirus(es), appeared to be restricted to a densely interconnected region

393    (Fig. 4). A subsequent hierarchical clustering within VC1 grouped these 168 viral genomes into

394    5 subgroups (Fig. S3). Interestingly, three phages infecting cyanobacteria (P-SSM2, P-SSM4,

395    and S-PM2) and T4-like phages that were initially found in a single cluster (Lima-Mendez et al.,

396    2008) are separated into two clusters of the Exo T-evens (VC_8) and T-evens/Pseudo/Schizo T-

397    evens (VC_1), respectively (Filee, 2006) (upper in Fig. 4B; Fig. S3). This network grouping can

398    identify the Exo T-evens including cyano- and pelagiphages, which the literature suggests to be

399    only distantly related to other T4 superfamily viruses (Comeau & Krisch, 2008; Roux et al.,

400    2015b).

401        ***The Autographivirinae*** We further identified 8 VCs associated with the

402    *Autographivirinae*. Of four genera defined by the NCBI and/or ICTV, the *T7virus*, *SP6virus*,

403    *Kp34virus* were found in VCs 4, 28, and 37, respectively, whereas the *PhiKMVvirus* were spread

404    across VCs 13 and 37 (Fig. 4B; also Fig. S4). Notably, a previous phylogenetic study based on

405    three conserved proteins (i.e., RNA polymerase, head-tail connector and the DNA maturase B)

406    showed considerable diversity of the *phiKMVvirus* (Eriksson et al., 2015). We also observed

407    distinct patterns of PC sharing between the PhiKMV-related genome(s) and other viruses in each

408    cluster (Fig. S4), suggesting that the *PhiKMVvirus* should likely be divided into two new

409    subgroups.

410        In addition, as the recently emerged groups, nine *Acinetobacter* phages (Huang et al.,

411    2013), as well as phage vB_CsaP_GAP227 (Abbasifar et al., 2013) and its close relatives were

412    found in VCs 54 and 93, respectively (Fig. S4); all of them encode T7-specific RNA polymerase

413    (Lavigne et al., 2009), which suggest that they fall within the *Autographivirnae* subfamily.

414        Finally, many viruses are now thought to co-opt host genes to improve viral fitness;

415    these stolen 'auxiliary metabolic genes' are known from cyanophage genomes (photosynthesis

416    genes; (Sullivan et al., 2006; Millard et al., 2009; Labrie et al., 2013), but also from ocean viral

417    metagenomes where viruses are now shown to contain genes involved in central carbon

418    metabolism (Hurwitz, Hallam & Sullivan, 2013) and nitrogen and sulfur cycling (Roux et al.,

419    2016) in ways that likely drive niche differentiation (Hurwitz, Brum & Sullivan, 2014). Thus it is

420    striking that VC22 in our network, which contains 19 cyanopodoviruses, had many linkages to

421    taxonomically disparate *Tevenvirinae*, which turned out the be driven by photosynthesis genes

422  shared across these viral taxa (Fig. 4B). Such "host" genes in viruses can bring taxonomically

423  disparate viral groups closer together, and the network can thus help identify such niche defining

424  viral genes for viruses infecting well studied hosts.

425  **The Archaeal Viruses**. Of the 72 archaeal viruses, 66 were associated with 18 VCs,

426  while 6 viruses (Haloviruses HHTV-1 and VNH-1, Hyperthermophilic Archaeal Virus 1 & 2,

427  Pyrococcous abyssi virus 1, and His 1 virus) were not included in the network, due to lack of

428  statistically significant similarity to any other virus. Of the 25 heterogeneous VCs, archaeal

429  viruses comprise 3 of them (VCs 51, 74 and 77), likely owing to their gene products showing

430  little similarity to published viruses outside of other archaeal viruses (Prangishvili, Garrett &

431  Koonin, 2006). All 3 VCs show considerable sharing of PCs within each VC (61.3 %, 50.2 %

432  and 67.6 %, respectively). VCs 74 and 77, each consisting of 2 genera

433  (*Gammalipothrixvirus*/*Rudivirus* and *Betalipothrixvirus*/*Deltalipothrixvirus*) unify the entire

434  *Ligamenvirales* order (2 families). Though the genera are distinguished mainly by their virion

435  morphology (Prangishvili & Krupovič, 2012), it can be argued that some lipothrixviruses share

436  as much similarity within the *Lipothrixviridae* family as to the rudiviruses, exemplified by the 10

437  genes shared between AFV-1 (a lipothrixvirus) and SIRV1 (a rudivirus) (Prangishvili &

438  Krupovič, 2012) and that they likely derive from a common ancestor (Goulet et al., 2009). In

439  addition to the number of PCs shared between AFV-1 and the rudivirus in VC74 (Fig. S1), the

440  more "distal" position between AFV-2 (*Deltalipothrixvirus*) and the other VC77 members

441  (*Betalipothrixvirus*) (Fig. S1), the order-level separation is easily seen in the overall network

442  (Fig. 2). VC55 (*Alphafusellovirus*/*Betafusellovirus*) consists of all known *Fuselloviridae*

443  members. Like VCs 74 and 77, their genera are separated mainly through virion morphology,

444  with Alphafusellovirus lemon-shaped and Betafusellovirus pleomorphic, and also through their

445   attachment structures (Redder et al., 2009). The large number of "core" genes (13) shared among

446   all family members argues for frequent recombination events, with even distant fuselloviruses

447   potentially capable of recombination during integration. Furthermore, some fuselloviruses

448   exhibit regions >70% pairwise identity on the nucleotide level, including ASV-1

449   (*Betafusellovirus*) and SSV-K1 (*Alphafusellovirus*) (Redder et al., 2009). Despite shared non-

450   core regions between the *fuselloviridae*, the high similarity between the two genera is also

451   revealed in the network through unification into a single VC. The most recently identified

452   member of the *Fuselloviridae*, *Sulfolobales* Mexican fusellovirus 1 (SMF1) has no official ICTV

453   classification between family, though clustering within the VC shows clear association with the

454   *Betafusellovirus*. It is remarkable that nearly all known archaeal viruses not only fall within the

455   network, but that most of their VCs follow a genus-level affiliation.

456          **vConTACT, *an iVirus tool for network-based viral taxonomy***: Given the strong and

457   robust performance of these network classification methods (Lima-Mendez et al., 2008) to

458   largely capture known viral taxonomy from genomes alone, we sought to democratize the

459   analytical capability. To this end, we developed a tool named "vContact" (overview of its logic

460   in Fig. 1) and integrated it into iVirus, a virus ecology-focused set of tools also known as "apps"

461   and databases (Bolduc et al., 2016). Such implementation at iVirus enables any user to run the

462   application simply by inputting viral sequences with all compute, storage and data repository

463   happening via the CyVerse cyberinfrastructure (formerly the iPlant Collaborative (Goff et al.,

464   2011).

465

466   **Conclusions**

467     Network-based approaches have been widely used to explore mathematical, statistical,

468     biological, and structural properties of a set of entities (nodes) and the connections between them

469     (edges) in a variety of biological and social systems (Dagan, 2011; Barberán et al., 2012). Such

470     approaches are invaluable for developing a quantitative framework to evaluate if and where

471     taxonomically meaningful classifications can be made in viral sequence space (Simmonds et al.,

472     2016). By expanding upon prior large-scale analyses (Lima-Mendez et al., 2008; Koonin,

473     Krupovic & Yutin, 2015; Roux et al., 2015a) we sought here to quantitatively evaluate when and

474     where such network-based classifications will perform poorly. These findings suggest that under

475     sampled viral sequence space and some phage groups with exceptionally high gene flow (e.g.,

476     mosaic genomes of the mycophages) are currently challenging and represent about 1 in 4

477     publicly-available, dsDNA viral genomes. While these problematic viral genomes await

478     improved representation of viral sequence space and/or improvements in network analytics to

479     best resolve their taxonomy, the remaining ¾ of viral genomes appear ready for gene sharing

480     network-based viral taxonomy. To this end, we present vConTACT as a publicly-available tool

481     for researchers to effectively enable large-scale, automated virus classification. Given the scale

482     of thousands of new virus genomes and genome fragments discovered through increasingly used

483     metagenomics approaches (Roux 2016 and Paez et al 2016), such step-wise progress towards an

484     automated taxonomic classifier will be foundational to most rapidly integrate viruses into models

485     that seek to make predictions in ecosystems ranging from the oceans and soils to bioreactors and

486     humans.

487

488     **Acknowledgements.**

493

**References**

495   Abbasifar R., Kropinski AM., Sabour PM., Ackermann H-W., Alanis Villa A., Abbasifar A.,

496     Griffiths MW. 2013. The Genome of Cronobacter sakazakii Bacteriophage

497     vB_CsaP_GAP227 Suggests a New Genus within the Autographivirinae. *Genome*

498     *Announcements* 1:e00122-12-e00122-12. DOI: 10.1128/genomeA.00122-12.

499   Ågren J., Sundström A., Håfström T., Segerman B. 2012. Gegenees: Fragmented alignment of

500     multiple genomes for determining phylogenomic distances and genetic signatures unique

501     for specified target groups. *PLoS ONE* 7. DOI: 10.1371/journal.pone.0039107.

502   Altschul SF., Madden TL., Schäffer AA., Zhang J., Zhang Z., Miller W., Lipman DJ. 1997.

503     Gapped BLAST and PSI-BLAST: a new generation of protein database search programs.

504     *Nucleic Acids Research* 25:3389–3402.

505   Assenov Y., Ramirez F., Schelhorn S-E., Lengauer T., Albrecht M. 2008. Computing topological

506     parameters of biological networks. *Bioinformatics* 24:282–284. DOI:

507     10.1093/bioinformatics/btm554.

508   Barberán A., Bates ST., Casamayor EO., Fierer N. 2012. Using network analysis to explore co-

509     occurrence patterns in soil microbial communities. *The ISME Journal* 6:343–351.

510   Bergsten J. 2005. A review of long-branch attraction. *Cladistics* 21:163–193. DOI:

511     10.1111/j.1096-0031.2005.00059.x.

512   Bolduc B., Youens-Clark K., Roux S., Hurwitz BL., Sullivan MB. 2016. iVirus: facilitating new

513     insights in viral ecology with software and community data sets imbedded in a

514     cyberinfrastructure. *The ISME Journal*:1–8. DOI: 10.1038/ismej.2016.89.

515    Brum JR., Ignacio-Espinoza JC., Roux S., Doulcier G., Acinas SG., Alberti A., Chaffron S.,
516        Cruaud C., de Vargas C., Gasol JM., Gorsky G., Gregory AC., Guidi L., Hingamp P.,
517        Iudicone D., Not F., Ogata H., Pesant S., Poulos BT., Schwenck SM., Speich S., Dimier C.,
518        Kandels-Lewis S., Picheral M., Searson S., Tara Oceans Coordinators., Bork P., Bowler C.,
519        Sunagawa S., Wincker P., Karsenti E., Sullivan MB. 2015a. Ocean plankton. Patterns and
520        ecological drivers of ocean viral communities. *Science (New York, N.Y.)* 348:1261498. DOI:
521        10.1126/science.1261498.

522    Brum JR., Ignacio-Espinoza JC., Roux S., Doulcier G., Acinas SG., Alberti A., Chaffron S.,
523        Cruaud C., de Vargas C., Gasol JM., Gorsky G., Gregory AC., Guidi L., Hingamp P.,
524        Iudicone D., Not F., Ogata H., Pesant S., Poulos BT., Schwenck SM., Speich S., Dimier C.,
525        Kandels-Lewis S., Picheral M., Searson S., Bork P., Bowler C., Sunagawa S., Wincker P.,
526        Karsenti E., Sullivan MB. 2015b. Patterns and ecological drivers of ocean viral
527        communities. *Science* 348:1261498–1261498. DOI: 10.1126/science.1261498.

528    Comeau AM., Krisch HM. 2008. The Capsid of the T4 Phage Superfamily: The Evolution,
529        Diversity, and Structure of Some of the Most Prevalent Proteins in the Biosphere.
530        *Molecular Biology and Evolution* 25:1321–1332. DOI: 10.1093/molbev/msn080.

531    Dagan T. 2011. Phylogenomic networks. *Trends in Microbiology* 19:483–491. DOI:
532        10.1016/j.tim.2011.07.001.

533    Deng L., Ignacio-Espinoza JC., Gregory AC., Poulos BT., Weitz JS., Hugenholtz P., Sullivan
534        MB. 2014. Viral tagging reveals discrete populations in Synechococcus viral genome
535        sequence space. *Nature* advance on. DOI: 10.1038/nature13459.

536    Edwards RA., Rohwer F. 2005. Viral metagenomics. *Nature Reviews Microbiology* 3:504–510.

537    Eriksson H., Maciejewska B., Latka A., Majkowska-Skrobek G., Hellstrand M., Melefors Ö.,
538        Wang JT., Kropinski AM., Drulis-Kawa Z., Nilsson AS. 2015. A suggested new
539        bacteriophage genus, "Kp34likevirus", within the Autographivirinae subfamily of
540        podoviridae. *Viruses* 7:1804–1822. DOI: 10.3390/v7041804.

541    Fauquet CM., Fargette D. 2005. International Committee on Taxonomy of Viruses and the 3,142
542        unassigned species. *Virology journal* 2:64. DOI: 10.1186/1743-422X-2-64.

543     Filee J. 2006. A Selective Barrier to Horizontal Gene Transfer in the T4-Type Bacteriophages
544         That Has Preserved a Core Genome with the Viral Replication and Structural Genes.
545         *Molecular Biology and Evolution* 23:1688–1696. DOI: 10.1093/molbev/msl036.

546     Goff S a., Vaughn M., McKay S., Lyons E., Stapleton AE., Gessler D., Matasci N., Wang L.,
547         Hanlon M., Lenards A., Muir A., Merchant N., Lowry S., Mock S., Helmke M., Kubach A.,
548         Narro M., Hopkins N., Micklos D., Hilgert U., Gonzales M., Jordan C., Skidmore E.,
549         Dooley R., Cazes J., McLay R., Lu Z., Pasternak S., Koesterke L., Piel WH., Grene R.,
550         Noutsos C., Gendler K., Feng X., Tang C., Lent M., Kim S-J., Kvilekval K., Manjunath
551         BS., Tannen V., Stamatakis A., Sanderson M., Welch SM., Cranston K a., Soltis P., Soltis
552         D., O'Meara B., Ane C., Brutnell T., Kleibenstein DJ., White JW., Leebens-Mack J.,
553         Donoghue MJ., Spalding EP., Vision TJ., Myers CR., Lowenthal D., Enquist BJ., Boyle B.,
554         Akoglu A., Andrews G., Ram S., Ware D., Stein L., Stanzione D. 2011. The iPlant
555         Collaborative: Cyberinfrastructure for Plant Biology. *Frontiers in Plant Science* 2:1–16.
556         DOI: 10.3389/fpls.2011.00034.

557     Goulet A., Blangy S., Redder P., Prangishvili D., Felisberto-Rodrigues C., Forterre P.,
558         Campanacci V., Cambillau C. 2009. Acidianus filamentous virus 1 coat proteins display a
559         helical fold spanning the filamentous archaeal viruses lineage. *Proceedings of the National*
560         *Academy of Sciences* 106:21155–21160.

561     Gregory AC., Solonenko SA., Ignacio-Espinoza JC., LaButti K., Copeland A., Sudek S.,
562         Maitland A., Chittick L., dos Santos F., Weitz JS., Worden AZ., Woyke T., Sullivan MB.
563         2016. Genomic differentiation among wild cyanophages despite widespread horizontal gene
564         transfer. *BMC Genomics* 17:930. DOI: 10.1186/s12864-016-3286-x.

565     Grose JH., Casjens SR. 2014. Understanding the enormous diversity of bacteriophages: The
566         tailed phages that infect the bacterial family Enterobacteriaceae. *Virology* 468:421–443.
567         DOI: 10.1016/j.virol.2014.08.024.

568     Halary S., Leigh JW., Cheaib B., Lopez P., Bapteste E. 2009. Network analyses structure genetic
569         diversity in independent genetic worlds. *Proceedings of the National Academy of Sciences*
570         107:127–132. DOI: 10.1073/pnas.0908978107.

571     Huang G., Le S., Peng Y., Zhao Y., Yin S., Zhang L., Yao X., Tan Y., Li M., Hu F. 2013.

572     Characterization and genome sequencing of phage Abp1, a new phiKMV-like virus

573     infecting multidrug-resistant acinetobacter baumannii. *Current Microbiology* 66:535–543.

574     DOI: 10.1007/s00284-013-0308-7.

575   Hurwitz BL., Brum JR., Sullivan MB. 2014. Depth-stratified functional and taxonomic niche

576     specialization in the "core" and "flexible" Pacific Ocean Virome. *The ISME journal*:1–13.

577     DOI: 10.1038/ismej.2014.143.

578   Hurwitz BL., Hallam SJ., Sullivan MB. 2013. Metabolic reprogramming by viruses in the sunlit

579     and dark ocean. *Genome Biology* 14:R123. DOI: 10.1186/gb-2013-14-11-r123.

580   Iranzo J., Krupovic M., Koonin E V. 2016. The Double-Stranded DNA Virosphere as a Modular

581     Hierarchical Network of Gene Sharing. *mBio* 7:e00978-16. DOI: 10.1128/mBio.00978-16.

582   Koonin E V., Krupovic M., Yutin N. 2015. Evolution of double-stranded DNA viruses of

583     eukaryotes: from bacteriophages to transposons to giant viruses. *Annals of the New York*

584     *Academy of Sciences*:n/a-n/a. DOI: 10.1111/nyas.12728.

585   Krupovic M., Dutilh BE., Adriaenssens EM., Wittmann J., Vogensen FK., Sullivan MB.,

586     Rumnieks J., Prangishvili D., Lavigne R., Kropinski AM., Klumpp J., Gillis A., Enault F.,

587     Edwards RA., Duffy S., Clokie MRC., Barylski J., Ackermann H-W., Kuhn JH. 2016.

588     Taxonomy of prokaryotic viruses: update from the ICTV bacterial and archaeal viruses

589     subcommittee. *Archives of Virology* 161:1095–1099. DOI: 10.1007/s00705-015-2728-0.

590   Kuhn JH., Dürrwald R., Bào Y., Briese T., Carbone K., Clawson AN., deRisi JL., Garten W.,

591     Jahrling PB., Kolodziejek J., Rubbenstroth D., Schwemmle M., Stenglein M., Tomonaga

592     K., Weissenböck H., Nowotny N. 2014. Taxonomic reorganization of the family

593     Bornaviridae. *Archives of Virology* 160:621–632. DOI: 10.1007/s00705-014-2276-z.

594   Labonté JM., Swan BK., Poulos B., Luo H., Koren S., Hallam SJ., Sullivan MB., Woyke T., Eric

595     Wommack K., Stepanauskas R. 2015. Single-cell genomics-based analysis of virus–host

596     interactions in marine surface bacterioplankton. *The ISME Journal* 9:2386–2399. DOI:

597     10.1038/ismej.2015.48.

598   Labrie SJ., Frois-Moniz K., Osburne MS., Kelly L., Roggensack SE., Sullivan MB., Gearin G.,

599     Zeng Q., Fitzgerald M., Henn MR., Chisholm SW. 2013. Genomes of marine

600 cyanopodoviruses reveal multiple origins of diversity. *Environmental Microbiology*
601 15:1356–1376. DOI: 10.1111/1462-2920.12053.

602 Lauber C., Gorbalenya AE. 2012a. Partitioning the Genetic Diversity of a Virus Family:
603 Approach and Evaluation through a Case Study of Picornaviruses. *Journal of virology*
604 86:3890–3904.

605 Lauber C., Gorbalenya AE. 2012b. Genetics-based classification of filoviruses calls for
606 expanded sampling of genomic sequences. *Viruses* 4:1425–1437. DOI: 10.3390/v4091425.

607 Lavigne R., Darius P., Summer EJ., Seto D., Mahadevan P., Nilsson AS., Ackermann HW.,
608 Kropinski AM. 2009. Classification of Myoviridae bacteriophages using protein sequence
609 similarity. *BMC Microbiology* 9:224. DOI: 10.1186/1471-2180-9-224.

610 Lawrence JG., Hatfull GF., Hendrix RW. 2002. Imbroglios of Viral Taxonomy: Genetic
611 Exchange and Failings of Phenetic Approaches. *Journal of Bacteriology* 184:4891–4905.
612 DOI: 10.1128/JB.184.17.4891-4905.2002.

613 Leplae R., Hebrant A., Wodak SJ., Toussaint A. 2004. ACLAME: a CLAssification of Mobile
614 genetic Elements. *Nucleic acids research* 32:D45-9. DOI: 10.1093/nar/gkh084.

615 Lima-Mendez G., Van Helden J., Toussaint A., Leplae R. 2008. Reticulate representation of
616 evolutionary and functional relationships between phage genomes. *Molecular Biology and*
617 *Evolution* 25:762–777. DOI: 10.1093/molbev/msn023.

618 Marston MF., Amrich CG. 2009. Recombination and microdiversity in coastal marine
619 cyanophages. *Environmental Microbiology* 11:2893–2903. DOI: 10.1111/j.1462-
620 2920.2009.02037.x.

621 Millard AD., Zwirglmaier K., Downey MJ., Mann NH., Scanlan DJ. 2009. Comparative
622 genomics of marine cyanomyoviruses reveals the widespread occurrence of Synechococcus
623 host genes localized to a hyperplastic region: Implications for mechanisms of cyanophage
624 evolution. *Environmental Microbiology* 11:2370–2387. DOI: 10.1111/j.1462-
625 2920.2009.01966.x.

626 Mizuno CM., Rodriguez-Valera F., Kimes NE., Ghai R. 2013. Expanding the Marine Virosphere
627 Using Metagenomics. *PLoS Genetics* 9. DOI: 10.1371/journal.pgen.1003987.

628   Niu YD., McAllister TA., Nash JHE., Kropinski AM., Stanford K. 2014. Four Escherichia coli
629        O157:H7 Phages: A New Bacteriophage Genus and Taxonomic Classification of T1-Like
630        Phages. *PLoS ONE* 9:e100426. DOI: 10.1371/journal.pone.0100426.

631   Paez-Espino D., Eloe-Fadrosh EA., Pavlopoulos GA., Thomas AD., Huntemann M., Mikhailova
632        N., Rubin E., Ivanova NN., Kyrpides NC. 2016. Uncovering Earth's virome. *Nature*
633        536:425–430. DOI: 10.1038/nature19094.

634   Pope WH., Bowman C a., Russell D a., Jacobs-Sera D., Asai DJ., Cresawn SG., Jacobs WR.,
635        Hendrix RW., Lawrence JG., Hatfull GF. 2015. Whole genome comparison of a large
636        collection of mycobacteriophages reveals a continuum of phage genetic diversity. *eLife*
637        4:e06416. DOI: 10.7554/eLife.06416.

638   Prangishvili D., Garrett RA., Koonin E V. 2006. Evolutionary genomics of archaeal viruses:
639        Unique viral genomes in the third domain of life. *Virus Research* 117:52–67.

640   Prangishvili D., Krupovič M. 2012. A new proposed taxon for double-stranded DNA viruses, the
641        order "Ligamenvirales." *Archives of Virology* 157:791–795.

642   Radoshitzky SR., Bào Y., Buchmeier MJ., Charrel RN., Clawson AN., Clegg CS., DeRisi JL.,
643        Emonet S., Gonzalez JP., Kuhn JH., Lukashevich IS., Peters CJ., Romanowski V., Salvato
644        MS., Stenglein MD., de la Torre JC arlos. 2015. Past, present, and future of arenavirus
645        taxonomy. *Archives of virology* 160:1851–1874. DOI: 10.1007/s00705-015-2418-y.

646   Redder P., Peng X., Brugger K., Shah SA., Roesch F., Greve B., She Q., Schleper C., Forterre P.,
647        Garrett RA., Prangishvili D. 2009. Four newly isolated fuselloviruses from extreme
648        geothermal environments reveal unusual morphologies and a possible interviral
649        recombination mechanism. *Environmental Microbiology* 11:2849–2862.

650   Roux S., Hawley AK., Torres Beltran M., Scofield M., Schwientek P., Stepanauskas R., Woyke
651        T., Hallam SJ., Sullivan MB. 2014. Ecology and evolution of viruses infecting uncultivated
652        SUP05 bacteria as revealed by single-cell- and meta- genomics. *eLife*:e03125. DOI:
653        10.7554/eLife.03125.

654   Roux S., Hallam SJ., Woyke T., Sullivan MB. 2015a. Viral dark matter and virus-host
655        interactions resolved from publicly available microbial genomes. *eLife* 4:e08490. DOI:

656          10.7554/eLife.08490.

657    Roux S., Enault F., Ravet V., Pereira O., Sullivan MB. 2015b. Genomic characteristics and

658          environmental distributions of the uncultivated Far-T4 phages. *Frontiers in microbiology*

659          6:199. DOI: 10.3389/fmicb.2015.00199.

660    Roux S., Brum JR., Dutilh BE., Sunagawa S., Duhaime MB., Loy A., Poulos BT., Solonenko N.,

661          Lara E., Poulain J., Pesant S., Kandels-Lewis S., Dimier C., Picheral M., Searson S., Cruaud

662          C., Alberti A., Duarte CM., Gasol JM., Vaqué D., Bork P., Acinas SG., Wincker P.,

663          Sullivan MB. 2016. Ecogenomics and potential biogeochemical impacts of globally

664          abundant ocean viruses. *Nature* 537:689–693. DOI: 10.1038/nature19366.

665    Simmonds P. 2015. Methods for virus classification and the challenge of incorporating

666          metagenomic sequence data. *Journal of General Virology* 96:1193–1206. DOI:

667          10.1099/vir.0.000016.

668    Sullivan MB., Lindell D., Lee J a., Thompson LR., Bielawski JP., Chisholm SW. 2006.

669          Prevalence and evolution of core photosystem II genes in marine cyanobacterial viruses and

670          their hosts. *PLoS Biology* 4:1344–1357. DOI: 10.1371/journal.pbio.0040234.

671    Woese CR., Kandler O., Wheelis ML. 1990. Towards a natural system of organisms: proposal

672          for the domains Archaea, Bacteria, and Eucarya. *Proceedings of the National Academy of*

673          *Sciences* 87:4576–4579.

674    Zanotto PM de., Gibbs MJ., Gould EA., Holmes EC. 1996. A reevaluation of the higher

675          taxonomy of viruses based on RNA polymerases. *Journal of virology* 70:6083–6096.
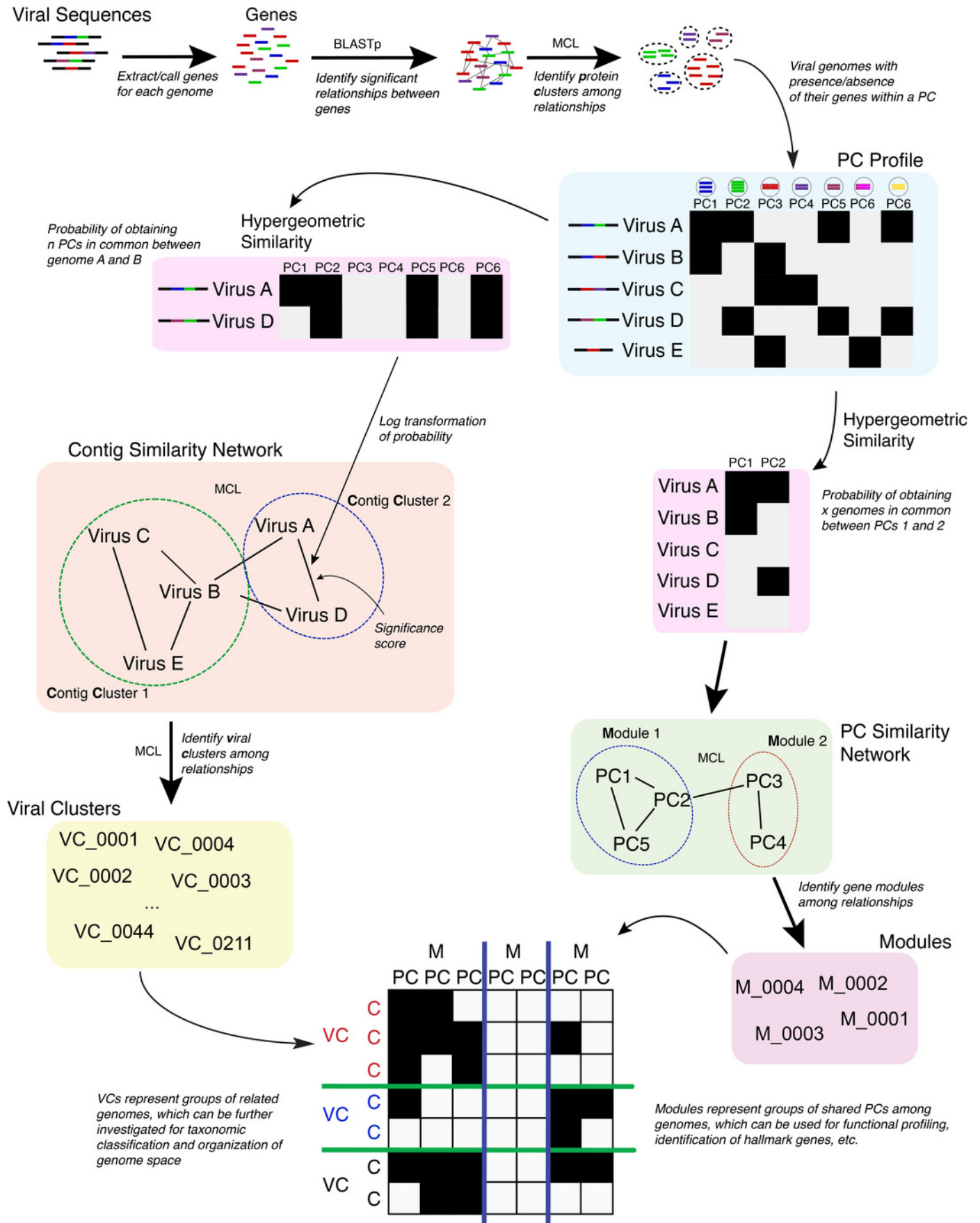
676

**Table 1**(on next page)

Terminology used.

1

| Terminology | Definition |
|---|---|
| **Nodes** | Also known as *vertices*, these are points within a network. In this work, they are viral genomes. |
| **Edges** | Also known as *arcs*, these lines connect nodes in the network. In this work, edges have a property called *weight*, which represents the strength (as measured by significance score) between two genomes. |
| **Betweenness centrality (BC)** | Measure of how influential a node is within a network, measured by the number of shortest paths that pass through the node from all other nodes. |
| **Connected component** | A subgraph in which any two nodes are connected to each other directly (to each other) or indirectly (through other nodes). |
| **Largest connected component (LCC)** | The connected component with the greatest number of nodes. |
| **Viral cluster (VC)** | A group of viral sequences sharing a significant number of genes. |
| **Protein cluster (PC)** | A group of highly similar and related proteins, defined in this work using MCL on BLAST E-values between proteins. |
| **Module Profile** | A table-like representation of the presence/absence data between groups of protein clusters (modules) and groups of genomes (viral clusters). |
| **Precision (P)** | Also known as the *positive predictive value*, is a measure of how many true positives are identified. |
| **Recall (R)** | Also known as *sensitivity*, is a measure of how many of the total positives are identified. |

2

# Figure 1

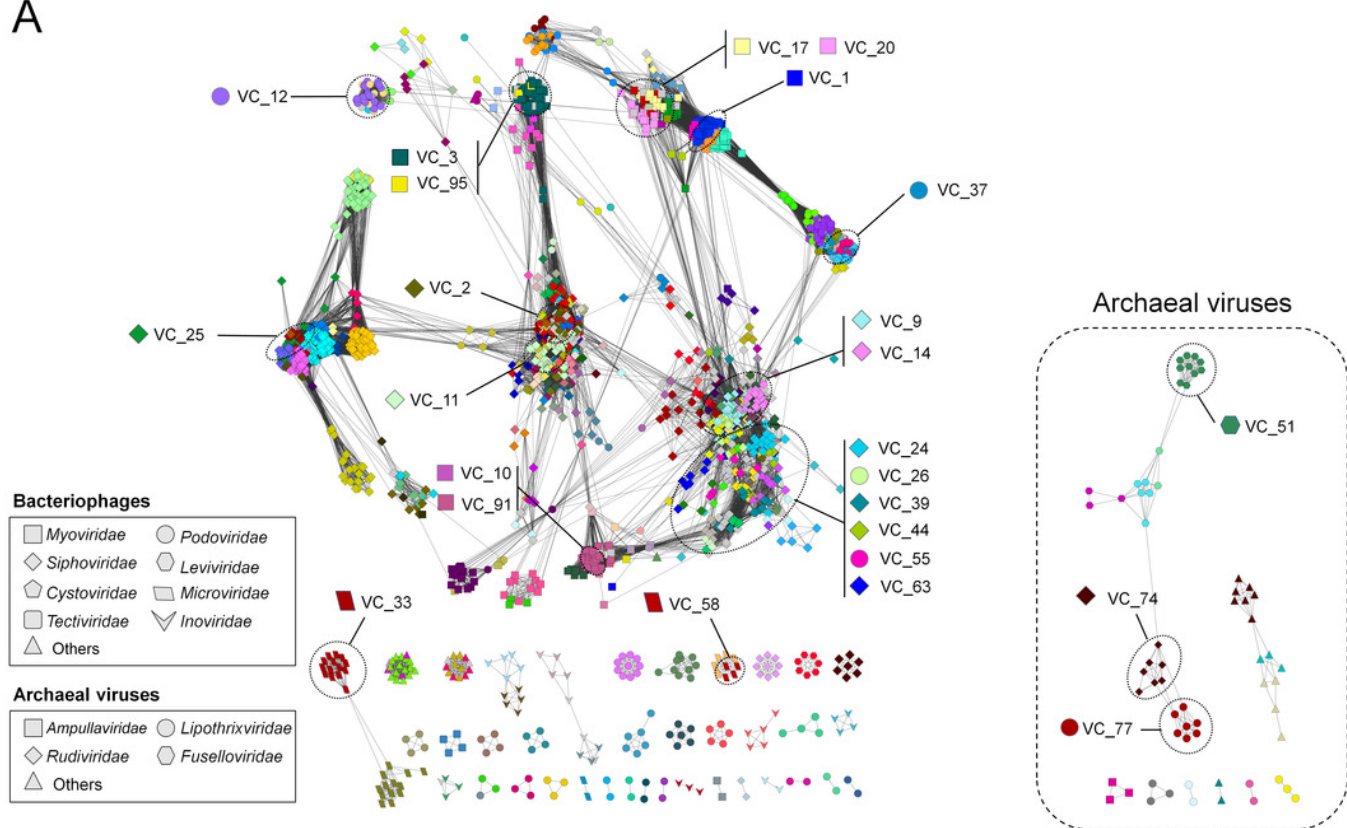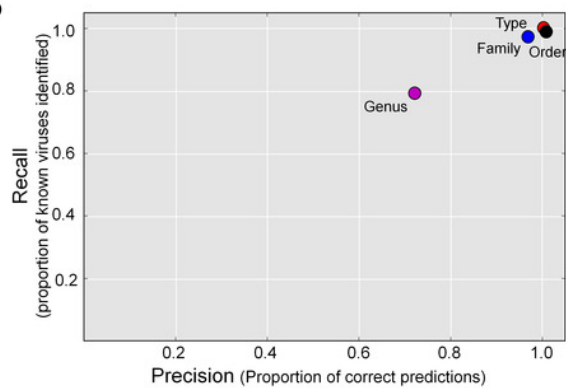Overview of the vContact processing pipeline.

# Figure 2

Protein-sharing network.

Protein-sharing network for 1,964 archaeal and bacterial virus genomes benchmarked against ICTV-accepted viral taxonomy. (A) Each node represents a viral genome from RefSeq, with its shape representing the viral family (as indicated in the legend) and each distinct color the node's viral cluster (VC). Edges between nodes indicate a statistically significant relationship between the protein profiles of their viral genomes, with edge colors (darker = more significant) corresponding to their weighted similarity scores (threshold of ≥1). VCs within the network are discriminated using the MCL algorithm (Materials and Methods) and denoted as separate colors. The position of 26 heterogeneous VCs that contain 2 or more genera is indicated. (B) Precision and recall of network-based assignments as compared to ICTV assignments for each taxonomic level (genus, family, order, and type). (C) Percentage (Y-axis) of VCs that contain the number (X-axis) of each ICTV taxonomic level (genus, family, and order).
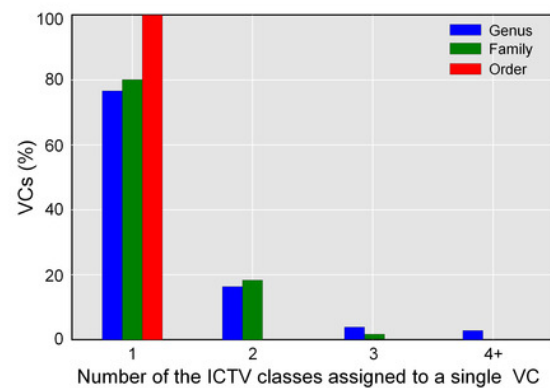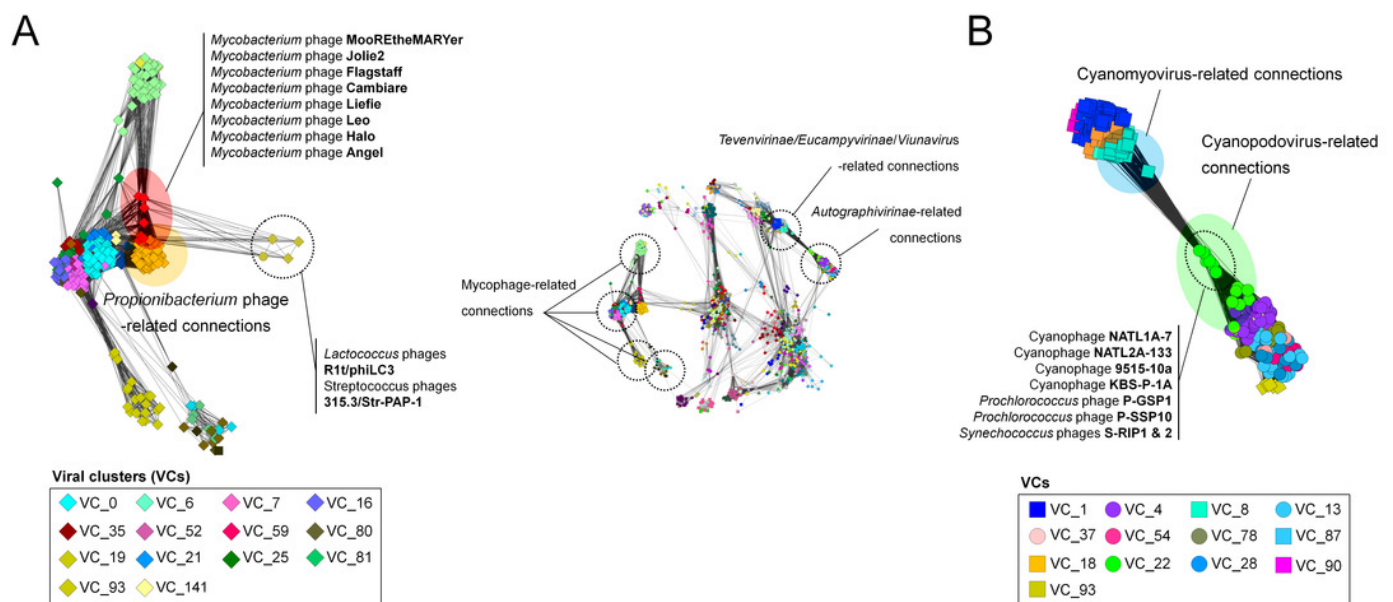
# Figure 3

Detailed view of three major viral groups with relatives.

A detailed view of network regions containing three major viral groups and their relatives. Viruses (nodes) are grouped by the MCL clustering. The different shapes and colors of the nodes represent different viral families (Figure 2) and viral clusters (VCs, legends), respectively. The location of viral groups is indicated for illustrative purpose.
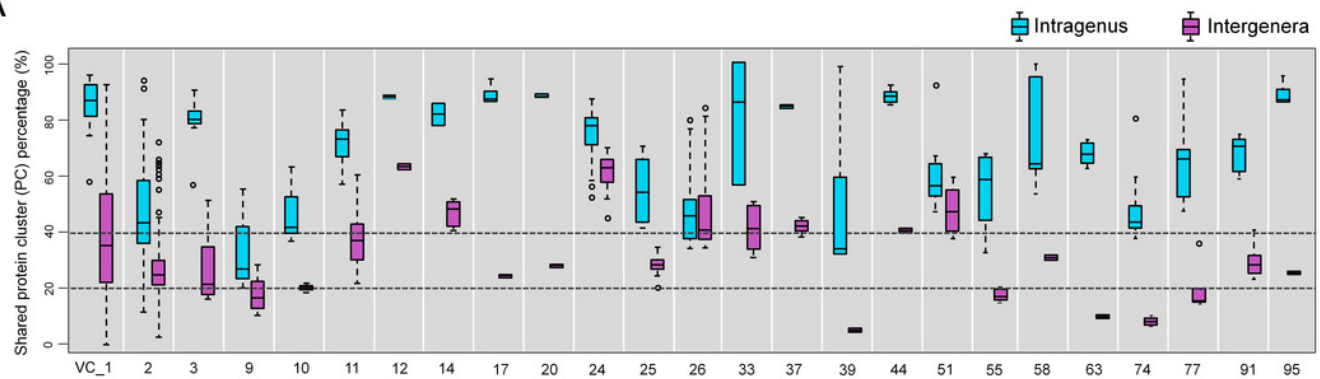
# Figure 4

Heterogeneous VCs

Evaluation of VCs which contained taxon representatives from more than one ICTV genus. (A) Box plots show the percent inter- and intra-genus proteome similarities in the heterogeneous VCs. Dotted lines indicate the cut-off values of 20% and 40% proteome similarities to define the subfamily and genus, respectively, which have been ratified by the ICTV Bacterial and Archaeal Viruses Subcommittee. (B) Module profiles showing the presence and absence of PCs across genomes. Presence (dark box) denotes a gene that is present within a protein cluster. Genes from related genomes often cluster into the same PC, with alignments of highly related genomes showing large groups of PCs. Genomes are further partitioned using hierarchical clustering (see materials and methods).