

Taxonomic assignment of uncultivated prokaryotic virus genomes is enabled by gene-sharing networks

Ho Bin Jang^{1*}, Benjamin Bolduc^{1*}, Olivier Zablocki¹, Jens H. Kuhn², Simon Roux³, Evelien M. Adriaenssens^{4,5}, J. Rodney Brister⁶, Andrew M Kropinski^{7,8}, Mart Krupovic⁹, Rob Lavigne¹⁰, Dann Turner¹¹, & Matthew B. Sullivan^{1,12#}

¹ Department of Microbiology, Ohio State University, Columbus, OH, USA

² Integrated Research Facility at Fort Detrick, National Institute of Allergy and Infectious Diseases, National Institutes of Health, Fort Detrick, Frederick, MD, USA

³ U.S. Department of Energy Joint Genome Institute, Walnut Creek, CA, USA

⁴ Institute of Integrative Biology, University of Liverpool, Liverpool, UK

⁵ Quadram Institute Bioscience, Norwich Research Park, Norwich, UK

⁶ National Center for Biotechnology Information, National Library of Medicine, National Institutes of Health, Bethesda, MD, USA

⁷ Department of Pathobiology, Ontario Veterinary College, University of Guelph, Guelph, ON, Canada N1G 2W1

⁸ Department of Food Science, University of Guelph, Guelph, ON, Canada, N1G 2W1

⁹ Institut Pasteur, Unité Biologie Moléculaire du Gène chez les Extrêmophiles, Department of Microbiology, Paris 75015, France

¹⁰ Laboratory of Gene Technology, Department of Biosystems, Faculty of BioScience Engineering, KU Leuven, Leuven, Belgium

¹¹ Centre for Research in Biosciences, Department of Applied Sciences, Faculty of Health and Applied Sciences, University of the West of England, Coldharbour Lane, Bristol, UK

¹² Department of Civil, Environmental and Geodetic Engineering, Ohio State University, Columbus, OH, USA

correspondence to: Matthew Sullivan, sullivan.948@osu.edu

* These authors contributed equally to this work.

ABSTRACT

Microbiomes from every environment contain a myriad of uncultivated archaeal and bacterial viruses, but studying these viruses is hampered by the lack of a universal, scalable, taxonomic framework. We present vConTACT v2.0, a network-based application utilizing whole genome gene-sharing profiles for virus taxonomy that integrates distance-based hierarchical clustering and confidence scores for all taxonomic predictions. We report near-identical (96%) replication of existing genus-level viral taxonomy assignments from the International Committee on Taxonomy of Viruses (ICTV) for NCBI virus Refseq. Application of vConTACT v2.0 to 1,364 previously unclassified viruses deposited in virus

39 **RefSeq as reference genomes produced automatic, high confidence genus assignments for**
40 **820/1364. We applied vconTACTv2.0 to analyse 15,280 Global Ocean Virome genome**
41 **fragments and were able to provide taxonomic assignments for 31% of these data, which**
42 **shows that our algorithm is scalable to very large metagenomic datasets. Our taxonomy**
43 **tool can be automated and applied to metagenomes from any environment for virus**
44 **classification.**

45

46

47 **Editors summary**

48 **Classification of archaeal and bacterial viruses can be automated with an algorithm that**
49 **identifies relationships based on shared gene content.**

50

51

52 Bacteria and archaea have roles in nutrient and energy cycles in ocean and soil ecosystems¹⁻
53 ⁴, as well as playing a vital part in human and health⁵. Viruses that infect bacteria and archaea
54 modulate these ‘ecosystem roles’ by killing, metabolic reprogramming or gene transfer^{6,7}, with
55 substantial effects of viral predation predicted in ocean⁸⁻¹⁰, soil^{11,12} and human microbiomes^{13,14}.
56 However, ecosystem-scale understanding of virus dynamics is hampered by the lack of universal
57 viral genes, or methods that enable a formalized taxonomy or comparative surveys. For example,
58 viruses do not have a single, universal marker gene¹⁵ so microbial-style 16S rRNA-based
59 phylogenies and operational taxonomic units (OTUs) are impossible¹⁶.

60 Virus sequencing has revealed structure^{17,18} and population genetic support for a species
61 definition¹⁹, and hypotheses have been put forward to explain variable evolution among
62 prokaryotic viruses²⁰. Together with rapidly expanding viral genome databases, these advances
63 have led the International Committee on Taxonomy of Viruses (ICTV) to present a consensus

64 statement suggesting a shift from the ‘traditional’ classification criteria²¹ e.g. virion morphology,
65 single/multiple gene phylogenies, towards a genome-centered, and perhaps one-day, largely
66 automated, viral taxonomy²².

67 Given the pace of viral discovery a virus taxonomy is urgently needed. Hundreds of
68 thousands of metagenome-derived viral genomes and large genome fragments (more than
69 700,000 at IMG/VR²³) dwarf the 34,091 prokaryotic virus genomes present in the NCBI
70 GenBank database²⁴. Together with the recently proposed ‘minimum information about
71 uncultivated virus genomes’ (MIUViGs) community guidelines²⁵, evaluation of approaches to
72 establish a scalable, genome-based viral taxonomy is needed to enable a universal classification
73 framework.

74 Multiple genome-based strategies have been proposed to develop a taxonomic framework for
75 viruses of bacteria^{15,26–31}, archaea³² or eukaryotes³³. For bacterial viruses (“phages”), one early
76 approach used complete genome pairwise protein sequence comparisons in a phylogenetic
77 framework (the “phage proteomic tree”) and was broadly concordant with ICTV-endorsed virus
78 groupings at the time¹⁵. However, this approach was not widely adopted as it was thought that
79 “rampant mosaicism” might blur taxonomic boundaries and violate the assumptions of the
80 underlying phylogenetic algorithms used in the analyses³⁴. Other approaches estimated the
81 fraction of genes shared, and percent identity of shared gene cut-offs, to define genera and sub-
82 family affiliations^{35,36}, but this approach failed to define taxonomic classification for several
83 known virus groups due to the likelihood that the mode and tempo of prokaryotic virus evolution
84 is highly variable²⁰. Building on a prokaryotic classification algorithm, the Genome Blast
85 Distance Phylogeny (GBDP)³⁷, which comes with a freely accessible online tool (VICTOR),
86 classifies phage genomes by combining phylogenetic and clustering methods²⁹. This method has

87 insufficient scalability (100 genomes limit) and limited taxonomic assignment for viruses that
88 lack reference genomes.

89 Gene sharing networks, based on shared protein clusters (PCs) between viral genomes, have
90 been shown to be largely concordant with ICTV-endorsed taxa, independent of whether
91 monopartite^{27,28,38} (a single node type, i.e., viral genomes) or bipartite networks^{32,38} (two node
92 types, i.e., viral genomes and genes) were used. We used a monopartite gene sharing network to
93 build an iVirus³⁹ app (vConTACT v1.0, hereafter v1.0) to automate network-based classification
94 of prokaryotic viruses. v1.0 produced viral clusters ('VCs') that were ~75% concordant with
95 ICTV prokaryotic viral genera²⁸. Network-based analytics have been applied to viral taxonomy
96 in large-scale studies of ocean^{40,41}, freshwater⁴² and soil⁴³ and studies of single-virus amplified
97 genomes (vSAGs)^{44,45}. In all of these environments the viruses could only be classified upon
98 application of a gene sharing network method. v1.0 cannot, however, make tentative taxonomic
99 assignments. This is because v1.0 creates artifactual VCs of both undersampled genomes and
100 highly overlapped regions of viral sequence space²⁸, and lacks per-VC confidence metrics,
101 necessary for establishing hierarchical taxonomy.

102 Here we present vConTACT v2.0 (hereafter v2.0), which has a new clustering algorithm,
103 confidence scoring of clusters and network analytics that together enable automation, improved
104 taxonomy assignments and scalability to much larger datasets. We apply v2.0 to establish a
105 centralized, 'living' taxonomic reference network as a community resource and show that v2.0 is
106 robust and scalable to large metagenomic datasets.

107

108 RESULTS

109 **Description of vConTACTv2.0**

110 The aim of vConTACT is to automatically assign viral genomes into established or new taxa,
111 with performance assessed relative to ICTV-assigned, manually-curated taxa (**Fig. 1**). However,
112 in the current ICTV taxonomy for prokaryotic viruses, taxonomic classifications above the genus
113 level are only sporadically available for sub-family and order ranks. For example, of the 2,304
114 prokaryotic virus genomes available in RefSeq, 84.2% are unclassified at the subfamily level,
115 and 61.6% are unclassified at the order level with virtually all of the remaining 38% lumped into
116 a single “*Caudovirales*” order. Moreover, among the *Caudovirales*, the three phenotypically
117 recognized and dominant bacterial virus family level designations – *Podoviridae*, *Myoviridae*
118 and *Siphoviridae* – are being called into question by genome-based taxonomy methods^{46–48} and
119 are thus in flux. Therefore, we focused specifically on assigning viruses at the genus level, as it
120 constitutes the principal taxon of molecular classification in the ICTV taxonomy.

121 In a network-based genome taxonomy framework (**Fig. 1a**), related genomes emerge as a
122 group of nodes strongly connected through multiple edges, here termed a Viral Cluster, or ‘VC’.
123 In a taxonomic context and based on the clustering of viral reference genomes, we have
124 previously demonstrated that the network parameters can be tuned such that the VCs best
125 represent genus-level grouping of viral genomes²⁸. In v1.0, ~75% of VCs corresponded to
126 established ICTV genera²⁸ (‘concordant VCs’), but ~25% ‘discordant VCs’ were present.
127 Discordant VCs can occur by production of outlier cluster genomes with no close relatives from
128 ‘undersampled VCs’, or by incorrect overlapping of multiple ICTV genera that share many genes
129 or by misassignment of multiple ICTV genera into a structured VC (**Fig. 1b**).

130 To address these problems we developed a new clustering algorithm, established confidence
131 scores and distance-based taxon separation for hierarchical taxonomy, and optimized and
132 evaluated scalability and robustness using a large-scale viral metagenomic dataset. Briefly, after

133 the MCL-clustered protein clusters are generated, we optimized the protein-cluster-based gene-
134 sharing information to establish an automated two-step process whereby VCs are defined using
135 ClusterONE⁴⁹ (CL1), rather than MCL which is used in v1.0, and then subdivided using
136 hierarchical clustering to disentangle problematic regions of the networks (**Fig. 1b**, Online
137 Methods). This approach considers edge weight (degree of connection between genomes) to
138 identify outlier genomes that are weakly connected with members of their VC compared to
139 neighbour genomes, detect and separate genomes that ‘bridge’ overlapping VCs, and break down
140 structured VCs into concordant VCs through distance-based hierarchical clustering (**Fig. 1b**).

141 Additionally, v2.0 incorporates confidence scores for each VC to help differentiate between
142 meaningful taxonomic assignments and those that might be artefacts. Briefly, each VC receives
143 two types of confidence scores: a topology-based score (value range 0-1), which aggregates
144 information about network topological properties, and a taxonomy-based score (value range 0-1),
145 which estimates the likelihood of predicted VCs to be equivalent to a single ICTV genus (Online
146 Methods). Higher values indicate either more confident linkages within the VC or better
147 taxonomic agreement for the topology and taxonomy-based scores, respectively, and the
148 taxonomy-based score is used to automatically optimize the hierarchical clustering of structured
149 VCs into ICTV-concordant ‘subclusters’.

150 Finally, although we present v2.0 as a monopartite (one type of node) network tool, it
151 produces the necessary output to be visualized as a bipartite network (**Supplementary Fig. 1**). In
152 bipartite visualizations, two types of nodes are used to display genomes and their connecting,
153 shared protein clusters (PCs). Information about which PCs link a given set a viruses together is
154 also provided (**Supplementary Table 1**; Online Methods), as it can enable identification of core
155 virus group genes that might be useful for downstream analyses.

156

157

158 **Comparison of vConTACT 1.0 and 2.0**

159 To assess clustering performance of v1.0 and v2.0, we quantified concordance with the set of
160 940 prokaryotic virus genomes that have ICTV genus-level classification (accessed January
161 2018, Online Methods and **Supplementary Table 2**). Clustering performance was evaluated by
162 a composite performance score of Accuracy (*Acc*) and Separation (*Sep*). Both *Acc* and *Sep* are
163 aggregate measures themselves (Online Methods), and report clustering precision, and how
164 resulting clusters (or VCs) correspond to a single ICTV genus, respectively (**Fig. 2a**). Each
165 metric has a value between 0 and 1, with 1 indicating perfect clustering accuracy and/or coverage.

166 v2.0's CL1, combined with hierarchical clustering, resulted in an overall performance
167 improvement of 28.8% (**Fig. 2a**). To assess which changes in v2.0 contributed to improved
168 performance we further optimized v1.0's MCL-based VC clustering and found that, at an IF of 7,
169 we could achieve nearly equivalent performance (**Fig. 2a, Supplementary Table 3**) and more
170 VCs predicted by the optimized MCL-based configuration as it organized the 940 viral genomes
171 into 180 VCs, whereas v2.0's CL1 identified 157 VCs. However, higher values in *Sep* for CL1
172 indicate better performance for assigning single genera into single VCs, even though MCL at its
173 optimal IF value (i.e., 7) generated more VCs (**Supplementary Table 1**). Thus, although more
174 VCs were assigned to ICTV genera by the optimized MCL configuration, they were largely
175 discordant VCs of either lumped or split ICTV genera, or both; whereas this behavior was ~50%
176 reduced using CL1 (see **Supplementary Fig. 2a and b**). Among these 22 lumped or split VCs
177 from the optimized MCL configuration, the virus genomes shared very few proteins (average =
178 17% range: 1-30%; **Supplementary Fig. 1b**) similarities, which modern cut-offs would suggest

179 should have been separated as separate genera, here outliers in the network. To better resolve
180 these issues, we added a post-processing, Euclidean distance-based hierarchical clustering step to
181 split mismatched VCs in v2.0. This step accurately classified 36 additional genera from the
182 problematic structured VCs (**Supplementary Table 2**), which increased v2.0's *Sep* value by 7%.
183 Together, these findings suggested that both upgrading the clustering algorithm and adding
184 hierarchical clustering were critical to improve automatic VC assignments.

185

186 **vConTACT v2.0 can analyse genomic relationships**

187 Next, we tested whether v2.0 could resolve discordant VCs (**Fig. 1b**). First, 55% of ICTV genera
188 are undersampled (**Supplementary Table 2**), which in a gene-sharing network manifests as
189 weakly connected, small VCs prone to artifactual clustering (**Fig. 1b**, top row) due to outlier
190 genomes only weakly connected to any given VC. In v1.0, undersampled VCs accounted for
191 64% (28/44) of all discordant VCs, and could not be resolved by increasing IF values (**Fig. 2b**
192 **and d** and **Supplementary Table 2**). v2.0 correctly places 38 genomes from 15 genera into 15
193 now concordant VCs using the same input data (**Fig. 2c and d** and **Supplementary Table 2**).
194 Second, we evaluated the ability of v2.0 to resolve overlapping VCs (**Fig. 1b**). We detected
195 overlapping VCs using a 'match coefficient' that measures the connection within- and between-
196 other VCs (Online Methods). This approach identified nine overlapping VCs (ICTV-classified
197 genera only) containing 30 viruses in 11 ICTV genera. These included viruses with known
198 mosaic genomes⁴⁷ (lambdoid or mu-like phages of the *P22virus*, *Lambdavirus*, *NI5virus*, and
199 *Bcepmyovirus* genera), recombinogenic temperate phages^{50,51} (*Mycobacterium* phages of the
200 *Bignuzvirus*, *Phayoncevirus*, and *Fishburnevirus* genera and *Gordonia* phages of the genus
201 *Wizardvirus*), and three newly-established genera (*Cd119virus*, *P100virus* and archaeal

202 *Alphapleolipovirus*), all bearing low topology-based confidence scores (averages of 0.32 for
203 these VCs versus 0.52 for concordant VCs; P-value = 6.12e-09, Mann-Whitney U test)
204 (**Supplementary Fig. 3a**). Overlapping VCs are linked to high horizontal gene flow, since most
205 viruses in these VCs were classified as having high gene content variation (HGCF, **Fig. 2e**,
206 **Supplementary Fig. 3b**) as assigned by a recently proposed framework of phage evolutionary
207 lifestyles²⁰. Though unresolvable in v1.0, v2.0 could assign eight of the 11 ICTV genera (24
208 viruses) into 8 ICTV-concordant VCs (**Supplementary Table 2**). The remaining 3 ICTV genera,
209 all comprised of *Mycobacterium* phages⁵² (6 genomes), could not be resolved (**Supplementary**
210 **Table 2**), and may not be amenable to automated taxonomy.

211 Third, structured VCs (**Fig. 1b**, bottom row) contained genomes that both gene sharing
212 networks placed into a single VC due to many shared genes and/or gene modules across all the
213 member genomes, but distributed into several ICTV genera due to subsets of the genomes also
214 sharing additional genes (**Supplementary Note 1**). For v1.0 we previously reported that these
215 structured VCs could be decomposed through hierarchical clustering²⁷, but in v2.0, we
216 formalized an optimized, quantitative hierarchical decomposition distance measure for this
217 process (Online Methods and **Supplementary Fig. 4**). In the v2.0 network, 23 of the 31
218 discordant VCs (74%) were structured VCs, spanning 86 genera (**Fig. 3a,b** and **Supplementary**
219 **Table 2**). Automated v2.0 resolved 30% (26 of 86) of these ICTV genera from 6 of the 23
220 structured VCs (**Fig. 3c**).

221 Of the 2,304 reference virus genomes classified by ICTV at the genus rank, 1,364 are
222 currently unassigned to a genus. This set of 1,364 reference viruses was organized into 404 well-
223 supported VCs with v2.0 (**Supplementary Table 2**). 544/1364 were placed in 104 VCs with
224 genomes from known ICTV taxa, whereas 820/1364 formed 200 separate VCs. We propose that

225 these 820 genomes can be as 200 *bona fide* novel virus genera and have submitted these to the
226 ICTV for consideration. If ratified, application of vContact2.0 will double the number of
227 prokaryotic viral genera (which is currently 264).

228 v2.0 clustering changed the taxonomy of ten established ICTV genera: *Barnyardvirus*,
229 *Bcep78virus*, *Bpp1virus*, *Che8virus*, *Jerseyvirus*, *P68virus*, *Pbunavirus*, *Phietavirus*,
230 *Phikmvvirus*, and *Yuavirus* (**Supplementary Fig. 5** and **Supplementary Note 2**), and manual
231 inspection by ICTV members involved in this study has recommended revision of *Phikmvvirus*
232 viruses (ICTV proposal 2015.007a-Db). Hierarchical decomposition of structured VCs into
233 subclusters indicated that the gene content-based distance correctly recapitulated the ICTV
234 taxonomy, but the cut-offs used to define subclusters are different from those currently used to
235 delineate established genera (**Fig. 3c** and **Supplementary Fig. 4**). Universal cut-offs are known
236 to be of limited use. Manual curation by experts has resulted in different cut-offs across viral
237 sequence space⁵³. A standardized taxonomy has been proposed for bacteria and archaea⁵⁴ and for
238 viruses standardization would be invaluable for automating virus taxonomy. v2.0 VCs and
239 subclusters will provide a reference baseline for the ICTV to translate network-derived cut-offs
240 into systematic taxonomic demarcation criteria.

241 Some taxon assignments are not amenable to being resolved by gene-sharing networks. For
242 example, when genera are defined on phenotypic or evolutionary evidence, e.g., archaeal
243 fuselloviruses⁵⁵ (VC42) or bacterial microviruses⁵⁶ (VCs 30 and 49), a gene-sharing network
244 approach will not be suitable (see **Fig. 3c** and **Supplementary Table 2**). An automated
245 vConTACT-based approach can however identify problematic taxa and speed up revisions to the
246 taxonomy.

247

248 **vConTACT v2.0 can scale to large virome datasets**

249 To evaluate scalability of our algorithm, we added 15,280 curated viral genomes and large
250 genome fragments (≥ 10 kb) from the Global Ocean Virome (GOV) dataset⁴⁰ to our reference
251 network in 10% increments (i.e., 0%, 10%, ..., 100% of the total dataset). The final network
252 comprised 16,960 sequences (**Fig. 4a**). We evaluated whether the incremental addition of GOV
253 data to the network led to changes in node connections, as estimated by the ‘change centrality’
254 metrics (CC, values range from 0-1 with 0 indicating no change and 1 indicating complete
255 change; **Fig. 4b**). We also evaluated concordance between v2.0 clustering and ICTV genera
256 using the Sn, Acc & PPV performance metrics (**Fig. 4c**). A large fraction of added data initially
257 experiences a moderate change (CC = 0.4), but the entire dataset eventually stabilized, as CC
258 values for most of the data ranged from 0 to 0.1. A similar trend was observed for accuracy (Acc,
259 **Fig. 4c**). This indicated that v2.0 can scale to thousands of input sequences, and that our
260 reference network clustering is robust to large-scale data additions.

261 We assessed whether GOV data can resolve ICTV outlier and singleton genomes as a proxy
262 for assessing taxonomic ramifications of adding data. We reasoned that more data might connect
263 outliers to new or existing VCs. Of 38 single-member VCs (**Supplementary Fig. 6**) three
264 *Mycobacterium* phage VCs were improved, while two *Mycobacterium* virus genomes were
265 merged into larger heterogeneous VCs composed of six ICTV genera, which did not constitute
266 an improvement. We observed that 919 new VCs were created with the full GOV dataset (15,
267 280 total contigs). We propose that these new VCs represent 919 viral genera that are not
268 represented in the existing 264 ICTV genera. According to a recent consensus statement, any
269 taxonomic reference network must be constrained to complete genomes²², and large genome

270 fragments commonly derived from metagenome-based studies must be utilized in a relevant
271 manner to address questions specific to that study, so these results remain preliminary.

272

273 **Discussion**

274 vConTACT v2.0 offers a scalable, robust, systematic and automated means to classify bacterial
275 and archaeal virus sequences. V2.0 is a highly scalable tool. Overall, there is a strong linear ($R^2 =$
276 0.99, see **Supplementary Fig. 7**) correlation between number of sequences and runtimes. For
277 example, running the full virus dataset RefSeq with Diamond would take ~10 minutes on a
278 regular laptop, while a GOV-sized dataset would run for several hours.

279 There are limitation of v2.0. First, the complete reference network needs to be rebuilt each
280 time new data are added. Avoiding this reconstruction step will require the development of
281 approximation methods and/or a placement algorithm (akin to PPlacer for 16S phylogenies⁵⁷) to
282 incorporate new data. Second, CL1-based VC generation may require manual parameter
283 optimization if datasets with overlapping genomes are included. We have added an auto-
284 optimization option for determining the optimal distance for hierarchical decomposition of
285 structured VCs in v2.0. v2.0 can run with prokaryotic viruses but has not been designed, tested or
286 validated for eukaryotic viruses. These viruses will require new algorithms for classification as
287 they have more diverse genomic configurations (segmentation, overlapping genes and
288 ambisense transcriptional gene configurations) that pose unique computational challenges^{33,58}.
289 Short, complete prokaryotic virus genomes and small fragments of larger genomes (e.g., ≤ 3 PCs
290 or ≤ 5 genes) have low statistical power in gene-sharing networks, and will require new solutions
291 to establish higher confidence VCs, and remain taxonomically inaccessible using v2.0. Finally,

292 genomes identified as singletons, outliers or overlapping are currently excluded from the gene-
293 sharing network, which leaves a large fraction of viral sequence space unclassified.

294 Assuming broad acceptance of vConTACT v2.0, and parallel efforts with eukaryotic
295 viruses³³, we may finally have the foundation to realize the consensus statement goals^{22,25} of
296 establishing a genome-based viral taxonomy to better capture the broader viral sequence
297 landscape emerging from environmental surveys.

298

299 **ACKNOWLEDGEMENTS.**

300 We thank Laura Bollinger, G. Trubl, and I. Tolstoy for their comments on improving the
301 manuscript, as well as Z-Q. You for helping push the network analytics. High performance
302 computational support was provided as an award from the Ohio Supercomputer Center to MBS.
303 Funding was provided in part by the Department of Energy's Genome Sciences Program Soil
304 Microbiome Scientific Focus Area award (#SCW1632) to Lawrence Livermore National
305 Laboratory; NSF Biological Oceanography awards (OCE#1536989 and OCE#1756314), and a
306 Gordon and Betty Moore Foundation Investigator Award (#3790) to MBS. Funding was
307 provided to JRB by the Intramural Research Program of the NIH, National Library of Medicine.
308 The work conducted by the U.S. Department of Energy Joint Genome Institute is supported by
309 the Office of Science of the U.S. Department of Energy under Contract DE-AC02-05CH11231
310 to SR. This work was funded in part through Battelle Memorial Institute's prime contract with
311 the US National Institute of Allergy and Infectious Diseases (NIAID) under Contract No.
312 HHSN272200700016I to JHK. The content of this publication does not necessarily reflect the
313 views or policies of the US Department of Health and Human Services or of the institutions and
314 companies affiliated with the authors.

315 **AUTHOR CONTRIBUTIONS.**

316 HBJ, BB and MBS designed the study. OZ and MBS wrote the manuscript with significant
317 contributions from all co-authors. HBJ and BB performed the statistical and network analyses.

318

319 **COMPETING INTERESTS.**

320 The authors declare no competing interests.

321

322 **MATERIALS & CORRESPONDENCE.** Correspondence and material requests should be
323 addressed to Matthew B. Sullivan at sullivan.948@osu.edu.

324

325

326 **REFERENCES**

- 327 1. Falkowski, P. G., Fenchel, T. & Delong, E. F. The microbial engines that drive earth's
328 biogeochemical cycles. *Science* **320**, 1034–1039 (2008).
329 2. Sunagawa, S. *et al.* Structure and function of the global ocean microbiome. *Science* (80-
330). **348**, (2015).

- 331 3. Moran, M. A. The global ocean microbiome. *Science* **350**, (2015).
- 332 4. Zhao, M. *et al.* Microbial mediation of biogeochemical cycles revealed by simulation of
333 global changes with soil transplant and cropping. *ISME J.* **8**, 2045–2055 (2014).
- 334 5. Cho, I. & Blaser, M. J. The human microbiome: At the interface of health and disease.
335 *Nature Reviews Genetics* **13**, 260–270 (2012).
- 336 6. Fernández, L., Rodríguez, A. & García, P. Phage or foe: an insight into the impact of viral
337 predation on microbial communities. *ISME Journal* 1–9 (2018). doi:10.1038/s41396-018-
338 0049-5
- 339 7. Hurwitz, B. L. & U'Ren, J. M. Viral metabolic reprogramming in marine ecosystems.
340 *Current Opinion in Microbiology* **31**, 161–168 (2016).
- 341 8. Suttle, C. a. Marine viruses-major players in the global ecosystem. *Nat. Rev. Microbiol.* **5**,
342 801–812 (2007).
- 343 9. Brum, J. R. *et al.* Patterns and ecological drivers of ocean viral communities. *Science* (80-
344). **348**, (2015).
- 345 10. Danovaro, R. *et al.* Virus-mediated archaeal hecatomb in the deep seafloor. *Sci. Adv.* **2**,
346 (2016).
- 347 11. Pratama, A. A. & van Elsas, J. D. The ‘Neglected’ Soil Virome - Potential Role and
348 Impact. *Trends in Microbiology* (2018). doi:10.1016/j.tim.2017.12.004
- 349 12. Gómez, P. & Buckling, A. Bacteria-phage antagonistic coevolution in soil. *Science* (80-).
350 **332**, 106–109 (2011).
- 351 13. Reyes, A., Semenkovich, N. P., Whiteson, K., Rohwer, F. & Gordon, J. I. Going viral:
352 Next-generation sequencing applied to phage populations in the human gut. *Nature*
353 *Reviews Microbiology* **10**, 607–617 (2012).
- 354 14. Abeles, S. R. & Pride, D. T. Molecular bases and role of viruses in the human
355 microbiome. *Journal of Molecular Biology* **426**, 3892–3906 (2014).
- 356 15. Rohwer, F. & Edwards, R. The phage proteomic tree: A genome-based taxonomy for
357 phage. *J. Bacteriol.* **184**, 4529–4535 (2002).
- 358 16. Yarza, P. *et al.* Uniting the classification of cultured and uncultured bacteria and archaea
359 using 16S rRNA gene sequences. *Nat. Rev. Microbiol.* **12**, 635–645 (2014).
- 360 17. Deng, L. *et al.* Viral tagging reveals discrete populations in *Synechococcus* viral genome
361 sequence space. *Nature* **513**, 242–245 (2014).
- 362 18. Gregory, A. C. *et al.* Genomic differentiation among wild cyanophages despite
363 widespread horizontal gene transfer. *BMC Genomics* **17**, (2016).
- 364 19. Bobay, L. & Ochman, H. Biological species in the viral world. **115**, (2018).
- 365 20. Mavrich, T. N. & Hatfull, G. F. Bacteriophage evolution differs by host, lifestyle and
366 genome. *Nat. Microbiol.* **2**, (2017).
- 367 21. Ackermann, H.-W. Phage Classification and Characterization BT - Bacteriophages:
368 Methods and Protocols, Volume 1: Isolation, Characterization, and Interactions. in (eds.
369 Clokie, M. R. J. & Kropinski, A. M.) 127–140 (Humana Press, 2009). doi:10.1007/978-1-
370 60327-164-6_13
- 371 22. Simmonds, P. *et al.* Consensus statement: Virus taxonomy in the age of metagenomics.
372 *Nat. Rev. Microbiol.* **15**, 161–168 (2017).
- 373 23. Paez-Espino, D. *et al.* IMG/VR v.2.0: an integrated data management and analysis system
374 for cultivated and environmental viral genomes. *Nucleic Acids Res.* gky1127-gky1127
375 (2018).
- 376 24. Brister, J. R., Ako-Adjei, D., Bao, Y. & Blinkova, O. NCBI viral Genomes resource.

- 377 *Nucleic Acids Res.* **43**, D571–D577 (2015).
- 378 25. Roux, S. *et al.* Minimum Information about an Uncultivated Virus Genome (MIUViG): a
379 community consensus on standards and best practices for describing genome sequences
380 from uncultivated viruses. *Nat. Biotechnol.* (2018).
- 381 26. Nishimura, Y. *et al.* ViPTree: The viral proteomic tree server. *Bioinformatics* **33**, 2379–
382 2380 (2017).
- 383 27. Lima-Mendez, G., Van Helden, J., Toussaint, A. & Leplae, R. Reticulate representation of
384 evolutionary and functional relationships between phage genomes. *Mol. Biol. Evol.* **25**,
385 762–777 (2008).
- 386 28. Bolduc, B. *et al.* vConTACT: an iVirus tool to classify double-stranded DNA viruses that
387 infect *Archaea* and *Bacteria*. *PeerJ* **5**, e3243 (2017).
- 388 29. Meier-Kolthoff, J. P. & Göker, M. VICTOR: genome-based phylogeny and classification
389 of prokaryotic viruses. *Bioinformatics* (2017). doi:10.1093/bioinformatics/btx440
- 390 30. Yu, C. *et al.* Real Time Classification of Viruses in 12 Dimensions. *PLoS One* **8**, (2013).
- 391 31. Gao, Y. & Luo, L. Genome-based phylogeny of dsDNA viruses by a novel alignment-free
392 method. *Gene* **492**, 309–314 (2012).
- 393 32. Iranzo, J., Koonin, E. V., Prangishvili, D. & Krupovic, M. Bipartite Network Analysis of
394 the Archaeal Virosphere: Evolutionary Connections between Viruses and Capsidless
395 Mobile Elements. *J. Virol.* **90**, 11043–11055 (2016).
- 396 33. Aiewsakun, P. & Simmonds, P. The genomic underpinnings of eukaryotic virus
397 taxonomy: creating a sequence-based framework for family-level virus classification.
398 *Microbiome* **6**, 38 (2018).
- 399 34. Lawrence, J. G., Hatfull, G. F. & Hendrix, R. W. Imbroglios of viral taxonomy: Genetic
400 exchange and failings of phenetic approaches. *J. Bacteriol.* **184**, 4891–4905 (2002).
- 401 35. Lavigne, R. *et al.* Classification of myoviridae bacteriophages using protein sequence
402 similarity. *BMC Microbiol.* **9**, (2009).
- 403 36. Lavigne, R., Seto, D., Mahadevan, P., Ackermann, H. W. & Kropinski, A. M. Unifying
404 classical and molecular taxonomic classification: analysis of the Podoviridae using
405 BLASTP-based tools. *Res. Microbiol.* **159**, 406–414 (2008).
- 406 37. Henz, S. R., Huson, D. H., Auch, A. F., Nieselt-Struwe, K. & Schuster, S. C. Whole-
407 genome prokaryotic phylogeny. *Bioinformatics* **21**, 2329–2335 (2005).
- 408 38. Iranzo, J., Krupovic, M. & Koonin, E. V. The double-stranded DNA virosphere as a
409 modular hierarchical network of gene sharing. *MBio* **7**, (2016).
- 410 39. Bolduc, B., Youens-Clark, K., Roux, S., Hurwitz, B. L. & Sullivan, M. B. iVirus:
411 Facilitating new insights in viral ecology with software and community data sets
412 imbedded in a cyberinfrastructure. *ISME J.* **11**, 7–14 (2017).
- 413 40. Roux, S. *et al.* Ecogenomics and potential biogeochemical impacts of globally abundant
414 ocean viruses. *Nature* **537**, 689–693 (2016).
- 415 41. Vik, D. R. *et al.* Putative archaeal viruses from the mesopelagic ocean. *PeerJ* **5**, e3428
416 (2017).
- 417 42. Roux, S. *et al.* Ecogenomics of virophages and their giant virus hosts assessed through
418 time series metagenomics. *Nat. Commun.* **8**, (2017).
- 419 43. Emerson, J. B. *et al.* Host-linked soil viral ecology along a permafrost thaw gradient. *Nat.*
420 *Microbiol.* (2018). doi:10.1038/s41564-018-0190-y
- 421 44. Martinez-Hernandez, F. *et al.* Single-virus genomics reveals hidden cosmopolitan and
422 abundant viruses. *Nat. Commun.* **8**, (2017).

- 423 45. de la Cruz Peña, M. J. *et al.* Deciphering the Human Virome with Single-Virus Genomics
424 and Metagenomics. *Viruses* **10**, 113 (2018).
- 425 46. Aiewsakun, P., Adriaenssens, E. M., Lavigne, R., Kropinski, A. M. & Simmonds, P.
426 Evaluation of the genomic diversity of viruses infecting bacteria, archaea and eukaryotes
427 using a common bioinformatic platform: Steps towards a unified taxonomy. *J. Gen. Virol.*
428 **99**, 1331–1343 (2018).
- 429 47. Hulo, C., Masson, P., Le Mercier, P. & Toussaint, A. A structured annotation frame for
430 the transposable phages: A new proposed family ‘Saltoviridae’ within the Caudovirales.
431 *Virology* **477**, 155–163 (2015).
- 432 48. Adriaenssens, E. M. *et al.* Taxonomy of prokaryotic viruses: 2017 update from the ICTV
433 Bacterial and Archaeal Viruses Subcommittee. *Archives of Virology* 1–5 (2018).
434 doi:10.1007/s00705-018-3723-z
- 435 49. Nepusz, T., Yu, H. & Paccanaro, A. Detecting overlapping protein complexes in protein-
436 protein interaction networks. *Nat. Methods* **9**, 471–472 (2012).
- 437 50. Doyle, E. L. *et al.* Genome Sequences of Four Cluster P Mycobacteriophages. *Genome*
438 *Announc.* **6**, e01101-17 (2018).
- 439 51. Pope, W. H. *et al.* Bacteriophages of *Gordonia* spp. Display a spectrum of diversity and
440 genetic relationships. *MBio* **8**, (2017).
- 441 52. Pope, W. H. *et al.* Whole genome comparison of a large collection of mycobacteriophages
442 reveals a continuum of phage genetic diversity. *Elife* **4**, e06416 (2015).
- 443 53. Nelson, D. Phage taxonomy: We agree to disagree. *Journal of Bacteriology* **186**, 7029–
444 7031 (2004).
- 445 54. Parks, D. H. *et al.* A standardized bacterial taxonomy based on genome phylogeny
446 substantially revises the tree of life. *Nat. Biotechnol.* **36**, 996 (2018).
- 447 55. Krupovic, M., Quemin, E. R. J., Bamford, D. H., Forterre, P. & Prangishvili, D.
448 Unification of the Globally Distributed Spindle-Shaped Viruses of the Archaea. *J. Virol.*
449 **88**, 2354–2358 (2014).
- 450 56. Rokyta, D. R., Burch, C. L., Caudle, S. B. & Wichman, H. A. Horizontal gene transfer and
451 the evolution of microvirid coliphage genomes. *J. Bacteriol.* **188**, 1134–1142 (2006).
- 452 57. Matsen, F. A., Kodner, R. B. & Armbrust, E. V. pplacer: linear time maximum-likelihood
453 and Bayesian phylogenetic placement of sequences onto a fixed reference tree. *BMC*
454 *Bioinformatics* **11**, 538 (2010).
- 455 58. Marz, M. *et al.* Challenges in RNA virus bioinformatics. *Bioinformatics* **30**, 1793–1799
456 (2014).

457
458
459
460

461 Figure Legends

462 **Figure 1. Virus genome classification visualized as networks.** (a) Left side panel: matrix of
463 shared protein clusters (PCs, grey blocks) between a set of virus genomes can be visualized as a
464 network of interconnected nodes, as shown on the right-side of the panel. Each node in this
465 sample 6-node network represents a virus genome that may be connected to other nodes through
466 edges. The edge value represents the strength of connectivity between nodes. If a set of nodes
467 have considerably higher edge weights than the rest of the network they are linked to, these are
468 grouped together to form a viral cluster, or ‘VC’.

(b) Each row depicts a node clustering scenario

469 in which vConTACT v2.0 has improved upon. On the left side, each scenario is first depicted as
470 a genome-PC matrix highlighting how shared protein clusters between certain genomes may
471 induce erroneous virus groupings due to outlier genomes, overlapping viral groups or VCs
472 containing multiple viral groups. On the right side of the matrices, the topology of each
473 clustering scenario is depicted as small networks of nodes (color-coded according to the ICTV
474 genera colors next to the matrices), and shows how vConTACT version 1 and 2 handled
475 clustering of problematic genomes and/or VCs. (c) Heatmap key corresponding to the various
476 values related to edge weight in (a) and (b), which serve to connect the nodes in the networks
477 and shows how closely related each connected node is to other nodes based on the number of
478 common PCs between genomes.

479
480 **Figure 2. Performance of vConTACT 1.0 and 2.0 on prokaryotic virus genomes.** (a) The
481 same colors denote individual performance metrics for the ICTV genera (G) including 940 viral
482 genomes, which are achieved by the Markov clustering (MCL) algorithm at each inflation factor
483 (v1.0) as well as ClusterONE (CL1) and CL1 followed by distance-based hierarchical clustering
484 (CL1 + H) (v2.0), respectively. For more objective comparisons, MCL at an inflation factor of
485 7.0 followed by hierarchical clustering (IF 7.0 + H) with the same distance (i.e., 9.0) used for
486 v2.0 was included. The total height and number of each bar indicate the composite score for
487 overall performance comparison. For details, see Online Methods. (b, c) Gene-sharing networks
488 were built using 2,304 archaeal and bacterial virus genomes retrieved from Viral RefSeq v85.
489 Viral clusters (VCs) were obtained by vConTACT v1.0 (b) and v2.0 (c) that used MCL with
490 inflation factor (IF) of 7.0 and CL1, respectively (see Online Methods). For both networks,
491 genomes (nodes) are color-coded according to their taxonomic assignments. For example,
492 genomes (only members of the ICTV-recognized genus) that are classified in VCs containing a
493 single ICTV genus are colored in cyan, while genomes found in VCs containing more than two
494 genera are colored in pink. Genomes without ICTV genus affiliation are in grey. Nodes with
495 bold borders indicate those that were correctly identified either as outlier, overlap genomes or
496 separate VCs through v1.0 (b), compared to v2.0 (c). Genomes whose taxonomic assignments
497 and/or annotation are incomplete are colored in yellow, identified through v2.0 (c). For details,
498 see **Supplementary Figs 3 and 5** and **Supplementary Table 2**. (d) Box plots of the percentage
499 of shared protein clusters (PCs) between member viruses within 28 v1.0-generated undersampled
500 VCs having ≥ 2 genera before (pink), and after (cyan) removal of outlier and/or separation into
501 individual clusters by v2.0. All box plots (n=61) were defined in terms of the minima, center,
502 maxima, percentiles and sample size (**Supplementary Table 5**). (e) Pie charts depicting the
503 number of overlapping genomes that belong to the high (HGCF) or low (LGCF) gene content
504 flux evolutionary modes or mixed and lytic or temperate phages. Data on the lifestyle and
505 evolutionary modes of 74 viruses were collected from Mavrigh and Hatfull²². For details, see
506 **Supplementary Fig. 3**.

507
508 **Figure 3. Application of the hierarchical decomposition to discordant VCs.** (a) Distribution
509 of all 31 discordant VCs across the archaeal and bacterial virus gene sharing network, where
510 genomes (nodes) of the given VCs are highlighted in pink and others in grey. (b) Box plots show
511 the fraction (%) of protein clusters (PCs) that were shared within an ICTV genus (i.e., intra-
512 genus proteome similarity) and between multiple genera (i.e., inter-genera similarity) found in
513 each discordant VC including structured clusters whose member genera have similar inter-genera
514 and intra-genus similarities (black dot). All box plots (n=60) were defined in terms of the

515 minima, center, maxima, percentiles and sample size (**Supplementary Table 6**). (c) Left, A full
516 link dendrogram is represented. Note that the Euclidean distance of nine yielded the highest
517 composite score of accuracy (*Acc*) and clustering-wise separation (*Sep*) for sub-clusters from all
518 v2.0-generated VCs, which was used to split the discordant clusters (**Online Methods** and
519 **Supplementary Fig. 4**). Right, module profiles showing the presence and absence of 7,662 total
520 protein clusters (PCs) across 362 genomes. Each row represents a phage and each column
521 represents a PC, with a unique color (left of the module) representing the genome's VC and
522 ICTV genus, respectively. Sub-clusters, which are generated by distance-based hierarchical
523 grouping, are represented across all discordant VCs on the right side of the heat map. From the
524 12 discordant VCs, 37 sub-clusters (corresponding to a single ICTV genus), are highlighted as
525 green boxes. For details, see **Supplementary Table 2**.

526
527 **Figure 4. Adding the Global Ocean Virome to NCBI Viral RefSeq.** (a) Selected network
528 images from the largest connected component of GOV additions. Red nodes are virus RefSeq
529 genomes, and grey nodes are GOV. Despite adding 15,280 new genomes, the network maintains
530 its overall structure. (b) Change centralities on a per-genome (grey) and per-VC (aqua) basis
531 through successive, 10% increments of GOV data. A value of zero in change centrality (Y-axis)
532 represent no change in any of the nodes connected to the origin node (or that the node was
533 removed), while a value of one represents origin node creation. High change centrality scores
534 imply that nodes are being created adjacent to the origin node, with the further a node's creation
535 is from the origin node, the less of an impact it has on the origin node's centrality. Dotted lines in
536 each violin represent quartiles, whereas the width of each violin plot is scaled to be equal
537 between GOV % (X-axis), such that distributions can be compared between datasets. Numbers in
538 parentheses indicate the number of genomes corresponding the GOV % above and numbers in
539 bracket indicate the number of corresponding VCs. Pairwise heatmap comparison at all GOV
540 incremental additions using normalized mutual information (NMI) values. NMI measures VC
541 similarity to other VCs by comparing genome content changes across incremental additions of
542 data. Darker blue hues correspond to more similar information content (i.e. genomes maintaining
543 the same VC membership). (c) GOV network performance through successive data
544 accumulations. As GOV sequences are added (X-axis), individual performance score (ranging
545 from 0 to 1, Y-axis; calculated from the clustering-wise positive predictive value (PPV),
546 clustering-wise sensitivity and accuracy) across genus- and family-level predictions (represented
547 by circular and square data points, respectively) generally trend towards stabilization. Boxplots
548 depicting the average Euclidean distance within VCs across GOV data increments. Grey boxes
549 are samples prior to hierarchical trimming, while blue boxes are post-trimming. Points represent
550 discordant VCs, with darker hues representing increasing discordance (i.e., more genera per VC).

551
552

553

554 ONLINE METHODS

555

556 **Data sets.** Full-length viral genomes were obtained from the National Center for Biotechnology
557 Information (NCBI) viral reference dataset^{24,59} ('ViralRefSeq', version 85, as of January, 2018),
558 downloaded from NCBI's viral genome page (<https://www.ncbi.nlm.nih.gov/genome/viruses/>)
559 and eukaryotic viruses were removed. The resulting file contained a total of 2,304 RefSeq viral
560 genomes including 2,213 bacterial viruses and 91 archaeal viruses (**Supplementary Table 2**). In
561 parallel, the ICTV taxonomy (ICTV Master Species List v1.3, as of February, 2018) was
562 retrieved from the ICTV homepage (<https://talk.ictvonline.org/files/master-species-lists/>). ICTV-
563 classifications were available for a subset of genomes at each taxonomic rank, and the final
564 dataset included: 884 viruses from two orders, 974 viruses from 23 families, 363 viruses from 28
565 subfamilies, and 940 viruses from 264 genera. To maintain hierarchical ranks of taxonomy, we
566 manually incorporated 2016 and 2017 ICTV updates^{48,60,61} to NCBI taxonomy when ICTV
567 taxonomy was absent.

568

569 **Generation of viral protein clusters.** Both version 1 and 2 of vConTACT share an identical
570 protein clustering initial step, in which viral proteins are grouped in protein clusters (PCs)
571 through MCL, followed by the formation of viral clusters (VCs) using either MCL (version 1) or
572 ClusterOne (version 2). First, a total of 231,166 protein sequences were extracted from the 2,304
573 viral genomes (above). Second, to group protein sequences into homologous protein clusters
574 (PCs)²⁸, all proteins were subjected to all-versus-all BLASTP⁶² searches (default parameters, cut-
575 offs of $1E^{-5}$ on e-value and 50 on bit score). Third, PCs were generated by applying MCL
576 (inflation factor of 2.0), and resulted into all the proteins being organized into 25,513 PCs, with a
577 fraction of proteins (26,625 or 11.5%) as singletons (i.e. isolated protein with no relatives).

578

579 **Calculating genome similarity between viruses.** The resulting output was parsed in the form of
580 a matrix comprised of genomes, PCs and singleton proteins (i.e., 2,304 × 52,138 matrix)
581 (**Supplementary Table 1**). We then determined the similarities between genomes by calculating
582 a one-tailed *P* value of observing at least *c* PCs in common between each pair of genomes, based
583 on the following hypergeometric equation as per Lima-Mendez et al²⁷:

584

$$585 \quad P(X \geq c) = \sum_{i=c}^{\min(a,b)} \frac{C_a^i C_{n-a}^{b-i}}{C_n^b}$$

586 (1)

587

588 in which *c* is the number of PCs in common; *a* and *b* are the numbers of PCs and singletons in
589 genomes A and B, respectively; and *n* is the total number of PCs and singletons in the dataset.
590 The hypergeometric formula calculates the probability of sharing a number of common PCs
591 between two genomes at or above the number (*c*) under the null hypothesis that the observed
592 result is likely to occur by chance. A score of similarity between genomes was obtained by
593 taking the negative logarithm (base 10) of the hypergeometric P-value multiplied by the total
594 number of pairwise genome comparisons (i.e., (2,304 × 2,303)/2). Genome pairs with a
595 similarity score ≥1 were previously shown to be significantly similar through permutation test
596 where PCs and singleton proteins with genome pairs having a similarity score below the given
597 threshold (negative control) were randomly rearranged. None of the genome pairs in this
598 negative control produced similarity score >1, indicating values above this threshold did not
599 occur by chance²⁸.

600

601 **Network visualization.** The gene (protein)-sharing network was constructed, in which nodes are
602 genomes and edges connect significantly similar genomes. This network was visualized with
603 Cytoscape software (version 3.6.0; <http://cytoscape.org/>), using an edge-weighted spring
604 embedded model, which places the genomes sharing more PCs closer to each other.

605

606 **Parameter optimization for viral cluster formation of vConTACT v1.0 and 2.0.** Due to
607 different criteria for parameter optimization between the clustering methods, different number
608 and size of the clusters are often generated, which can make objective performance comparisons
609 difficult⁶³. Thus, to more comprehensively compare performance, v1.0's MCL-based VCs were
610 generated at inflation factors (IFs) of 2.0 to 7.0 by 1.0 increments, with an optimal IF of 1.4
611 showing the highest intra-cluster clustering coefficient (ICCC)²⁷ (**Supplementary Table 2** and
612 **Supplementary Fig. 8**). Unlike MCL, which uses a single parameter²⁷ (i.e., the inflation factor),
613 VC formation with CL1 (used in vConTACT v2.0), involves multiple parameters that can detect
614 complex network relationships⁴⁹. The three main parameters of CL1, minimum density/node
615 penalty, haircut, and overlap, automatically quantify (i) the cohesiveness of a cluster, (ii) the
616 boundaries of the clusters (i.e. outlier genomes), and (iii) the size of overlap between clusters ,
617 respectively⁴⁹. Of these parameters, the first one is used to detect the coherent groups of VCs as
618 follows:

619

$$620 \quad C = \frac{W_{in}(V)}{W_{in}(V) + W_{out}(V) + p|C|}$$

621 (2)

622

623 in which $W_{in}(V)$ and $W_{out}(V)$ are the total weight of edges that lie within cluster V and that
624 connect the cluster V and the rest of the network, respectively, $|C|$ is the size of the cluster, p is a
625 penalty that counts the possibility of uncharted connections for each node.

626 The second parameter, the haircut, can find loosely connected regions of the network (outliers)
627 by measuring the ratio of connectivity of the node g within the cluster c to that of its
628 neighbouring node h as:

629

$$630 \quad \Delta_{out} = k \sum_{j=1}^l W_{h,j} / \sum_{i=1}^k W_{g,i}$$

631 (3)

632

633 in which k is the number of edges of the node g , and W is the total weight of edges of the
634 respective nodes g and h . If the total weight of edges from a node (h) to the rest of the cluster (c)
635 is less than x times that we specified the average weight of nodes (g) within the given cluster,
636 CL1 will remove the node (h) from a given VC and consider it an outlier.

637 The third CL1 parameter, the overlap size, determines the maximum allowed overlap (ω)
638 between two clusters, measured by the match coefficient, as follows:

639

$$640 \quad \omega = i^2 / a * b$$

641 (4)

642

643 in which i is the size of overlap, which is divided by the product of the sizes of the two clusters
644 under consideration (a and b). Since CL1 identifies overlap between VCs, it can find both
645 hierarchical and overlapping structures within viral groups. This ability is a significant

646 improvement over v1.0, as v1.0's MCL cannot handle modules with overlaps⁷. Specifically, for
647 each pair of clusters, CL1 calculates the overlap score between them (above) and merges these
648 clusters if the overlap is larger than a given threshold. Thus, in the resulting output file, viral
649 groups (or clusters) having the identical member viruses can be found in multiple clusters, called
650 'overlapping viral clusters' (**Supplementary Table 2 and Fig. 1b, middle row**).

651 To determine the best parameter combination to use for CL1, we tested a wide range of
652 values for the three aforementioned parameters: minimum density ranging from 0 to 1 by 0.1
653 increments; node penalty from 1 to 10 by 1.0; haircut from 0 to 1 by 0.05; overlap from 0 to 1 by
654 0.05) and default settings for the other parameters: 2 as minimum cluster size, weighted as edge
655 weight, single-pass as merging, unused nodes as seeding. This resulted in 53,361 clustering
656 results, which we evaluated individually to determine the highest performance on our genome
657 data set (above) To identify the best parameter combination, we used the geometric mean value
658 of prediction accuracy (*Acc*) and clustering-wise separation (*Sep*, see next section), as previously
659 described⁶⁴. The final, optimized CL1 parameters were a minimum density of 0.3, a node penalty
660 of 2.0, a haircut of 0.65, and an overlap of 0.8, which resulted in 280 VCs (**Supplementary
661 Table 2**).

662 Next, to further decompose 'discordant VCs', we added as a post-clustering step in v2.0,
663 which allows additional hierarchical separation of such VCs into sub-clusters using the
664 unweighted pair group method with arithmetic mean (UPGMA) with pairwise Euclidean
665 distances (implemented in Scipy). To determine the optimal distance for sub-clustering of VCs,
666 we assessed the distances of sub-clusters across all the VCs in the network. We tested the effect
667 of these distances (ranging from 1 to 20 in 0.5 increments) and picked as optimal distance the
668 one which maximized the composite score by multiplying the prediction accuracy (*Acc*) and

669 clustering-wise separation (*Sep*) at the ICTV genus rank (see next section). A distance of 9.0
670 yielded the highest composite score of *Acc* and *Sep* (**Supplementary Fig. 4**). Notably,
671 vConTACT v2.0 was designed to help users optimize these parameters for grouping of
672 genomes/contigs into VCs and distance for post-decomposition of VCs into sub-clusters. This
673 tool automatically evaluates the robustness of each VCs and sub-clusters, based on the external
674 performance evaluation statistics (below).

675

676 **Performance comparison between vConTACT v1.0 and v2.0.** Six external quality metrics
677 were used to compare clustering performance between MCL and CL1⁶⁴ (**Fig. 2a**). Specifically,
678 the performance of v1.0 (MCL) and v2.0 (CL1 alone and CL1 + hierarchical sub-clustering)
679 were evaluated based on : (i) cluster-wise sensitivity, *Sn* (ii) positive predictive value, *PPV* (iii)
680 geometric mean of *Sn* and *PPV*, *Acc* (iv) cluster-wise separation, *Sep_{cl}* (v) complex (ICTV
681 taxon)-wise separation *Sep_{co}*, and (vi) geometric mean of *Sep_{cl}* and *Sep_{co}*, *Sep*. As an internal
682 parameter, we computed the intra- and inter-cluster proteome similarities (fraction of shared
683 genes between genome that are within the same VCs and different VCs, respectively). For
684 vConTACT v1.0, we only included clustering results which had been determined to yield the
685 highest clustering accuracy value (i.e., inflation factor of 7.0), and this configuration was used
686 for comparison to v2.0's clustering. Therefore, testing each parameter combination (6
687 performance metrics, for one taxon rank, for 10 clustering results, all cross-compared; i.e., 6 x 1
688 x 45) resulted in 270 comparisons.

689 To generate six external measures, we first built a contingency table *T*, in which row *i*
690 corresponds to the *i*th annotated reference complex (i.e., ICTV-recognized order, family,
691 subfamily, or genus), and column *j* corresponds to the *j*th predicted complex (i.e., sub-/clusters).

692 The value of a cell T_{ij} denotes the number of member viruses in common between the i^{th}
 693 reference complex and j^{th} predicted complex.

694 **Sensitivity**: The sensitivity can be defined as the fraction of member viruses of complex i which
 695 are found in sub-/cluster j .

$$696 \quad \quad \quad Sn_{i,j} = T_{i,j}/N_i$$

697 (5)

698 In the formula above, N_i is the number of member viruses of complex i . We then calculated the
 699 coverage of complex i by its best-matching cluster Sn_{co_i} , as the maximal fraction of member
 700 viruses of complex i assigned to the same sub-/cluster by the formula below:

$$701 \quad \quad \quad Sn_{co_i} = \max_{j=1}^m Sn_{i,j}$$

702 (6)

703 The clustering-wise sensitivity was computed as the weighted average of Sn_{co_i} over all
 704 complexes. Higher Sn values indicate a better coverage of the member viruses in the real
 705 complexes as:

$$706 \quad \quad \quad Sn = \frac{\sum_{i=1}^n Ni Sn_{co_i}}{\sum_{i=1}^n Ni}$$

707 (7)

708 **Positive predictive value**: The positive predictive value (*PPV*) indicates the proportion of
 709 member viruses of the sub-/cluster j which belong to complex i , relative to the total number of
 710 member viruses of the sub-/cluster assigned to all complexes by:

$$711 \quad \quad \quad PPV_{i,j} = T_{i,j}/\sum_{i=1}^n T_{i,j} = T_{i,j}/T_{.j}$$

712 (8)

713 where T_j is the marginal sum of a column j . We calculated the maximal fraction of member
 714 viruses of sub-/cluster j found in the same annotated complex PPV_{cl_j} , as the prediction
 715 reliability of sub-/cluster j to belong to its best-matching complex as:

$$716 \quad PPV_{cl_j} = \max_{i=1}^n PPV_{i,j}$$

717 (9)

718 The clustering-wise PPV was then computed as the weighted average of PPV_{cl_j} over all
 719 sub/clusters by:

$$720 \quad PPV = \frac{\sum_{j=1}^m T_j PPV_{cl_j}}{\sum_{j=1}^m T_j}$$

721 (10)

722 Higher PPV values indicate that the predicted sub-/clusters are likely to be true positives.

723 **Accuracy:** As a summary metric, the Acc can be obtained by computing the geometrical mean of
 724 the Sn and PPV values as:

$$725 \quad Acc = \sqrt{Sn \times PPV}$$

726 (11)

727
 728
 729 **Complex- and Cluster-wise separations:** With the same contingency table used for Sn , PPV , and
 730 Acc , we calculated the relative frequencies with respect to the marginal sums for each row
 731 ($F_{row_{i,j}}$) and each column ($F_{col_{i,j}}$), respectively:

$$732 \quad F_{row_{i,j}} = T_{i,j} / \sum_{j=1}^m T_{i,j}$$

733 (12)

734
$$F_{col_{i,j}} = T_{i,j} / \sum_{i=1}^n T_{i,j}$$

735 (13)

736 Then the separation is computed as the product of column-wise and row-wise frequencies as:

737
$$Sep_{i,j} = F_{col_{i,j}} \times F_{row_{i,j}}$$

738 (14)

739 The separation values range from 0 to 1, with 1 indicating a perfect correspondence between
 740 complex j and sub-/cluster i (i.e., the cluster contains all the members of the complex and only
 741 them). Additionally, the separation penalizes the case when member viruses of a given complex
 742 are split into multiple sub-/clusters. The complex-wise Sep_{co} and cluster-wise Sep_{cl} values are
 743 calculated as the average of Sep_{co_i} over all complexes, and of Sep_{cl_j} over all sub-/cluster,
 744 respectively:

745
 746
$$Sep_{co} = \frac{\sum_{i=1}^n Sep_{co_i}}{n}$$

747 (15)

748

749
$$Sep_{cl} = \frac{\sum_{j=1}^m Sep_{cl_j}}{m}$$

750 (16)

751 To estimate these separation results as a whole, the geometric mean (clustering-wise separation;
 752 Sep) of Sep_{co} and Sep_{cl} was computed:

753

754
$$Sep = \sqrt{Sep_{co} \times Sep_{cl}}$$

755 (17)

756 High clustering-wise separation values indicate a bidirectional correspondence between a sub-
757 /cluster and each ICTV taxon: a score of 1.0 indicates that a cluster corresponds perfectly to each
758 taxon. For overall comparison, we used a composite score⁴⁹, calculated by multiplying *Acc* by
759 *Sep*.

760 As an internal measure, the fraction of PCs²⁸ between two genomes (i.e., proteome similarity)
761 was computed by using the geometric index (G). The proteome similarity was estimated as:

$$763 \quad G_{AB} = \frac{|N(A) \cap N(B)|}{|N(A)| \times |N(B)|}$$

764 (18)

765

766 in which $N(A)$ and $N(B)$ indicate the number of PCs in the genomes of A and B, respectively. A
767 total of 400,234 pairs of genomes with >1% proteome similarity are shown in **Supplementary**
768 **Table 4**.

769

770 **Clustering-based confidence score.** To generate confidence scores for each viral cluster
771 prediction, we used three previously described confidence scoring methods^{65,66}, with some
772 modifications. Two of them exploit the network topology properties by assessing the weight of
773 cluster quality and the probability of cluster quality. We then combined these two values as an
774 aggregate topology-based confidence score per VC. For the first scoring method, we computed
775 the quality (Q) of sub-cluster (c) as:

776

$$777 \quad Q_c = W_{in} / (W_{in} + W_{out})$$

778 (19)

779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800

in which W_{in} and W_{out} are the total weight of edges that lie within sub-cluster c and across others, respectively. For the second method, we evaluated the P-value of a one-sided Mann-Whitney U test for in-weights and out-weights of sub-clusters. The rationale behind this test is that sub-clusters with a lower P-value contains significantly higher in-weights than out-weights, thus indicative that a formed sub-cluster is valid, and not a random fluctuation. These two independent values, weight of cluster quality and the probability of cluster quality are then multiplied to derive a topology-based confidence score for each cluster. Along with this confidence score, we quantified the likelihood that each sub-cluster corresponds to an ICTV-approved genus (or equivalent) by using distance threshold that are specified at the ICTV genus rank, which we refer to as “taxon predictive score”. This score can be calculated as:

$$prediction = \sum l_{i,j} / l_c \tag{20}$$

Specifically, for a sub-cluster (c) having the genus-level assignment, vConTACT v2.0 automatically measures the maximum distance between taxonomically-known member viruses and calculate the scores by dividing the sum of links having less than the given maximum distance threshold between nodes (i and j) by the total number of links (l_c) between all nodes. For a sub-cluster that does not have the genus-level assignment, v2.0 uses Euclidean distance of 9.0 that can maximize the prediction accuracy and clustering-wise separation (see above) as distance threshold.

801 **Measuring effect of GOV on network structural changes.** GOV contigs (14,656 sequences)
802 were added in 10% increments (randomly selected at each iteration) to NCBI Viral RefSeq and
803 processed using vConTACT v2.0 with one difference – Diamond⁶⁷ instead of BLASTp was used
804 to construct the all-versus-all protein comparison underlying the PC generation. For running this
805 large number of sequences, high-memory computer nodes from the Ohio State supercomputer
806 Center⁶⁸ were used. Once generated, vConTACT v2.0 networks were post-processed using a
807 combination of the Scipy⁶⁹, Numpy, Pandas⁷⁰ and Scikit-learn⁷¹ python 3.6 packages. Networks
808 were rendered using iGraph⁷². The method to calculate change centrality was calculated as
809 described previously⁷³. CCs were calculated in a successive way, in which each addition was
810 compared to Viral RefSeq 85 independently of other additions (0% versus 10%, 0% vs 20%,
811 [...], 0% vs 100%).

812

813 **Data and code availability statement**

814 The set of reference genomes used to evaluate vConTACT were retrieved from
815 <https://www.ncbi.nlm.nih.gov/genome/viruses/>. The Global Ocean Virome (GOV) contigs were
816 retrieved from the publically available CyVerse data commons repository, accessible at
817 <http://datacommons.cyverse.org/browse/iplant/home/shared/iVirus/GOV>. The utility of v2.0
818 depends upon its expert evaluation and community availability. The tool is available through
819 Bitbucket (<https://bitbucket.org/MAVERICLab/vcontact2>) as a downloadable python package,
820 and usable as an app through iVirus³⁹, the viral ecology apps and data resource embedded in the
821 CyVerse Cyberinfrastructure, with detailed usage protocols available through Protocol Exchange
822 (<https://www.nature.com/protocolexchange/>) and protocols.io (<https://www.protocols.io/>).
823 Finally, the curated reference network is available at each of these sites, and will be updated

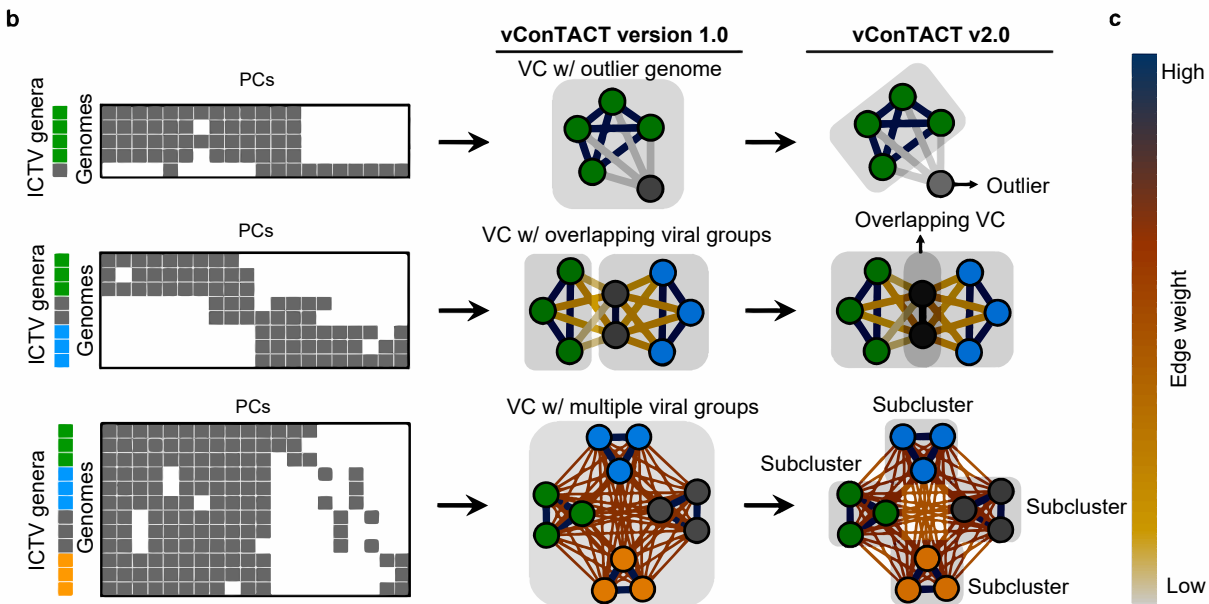
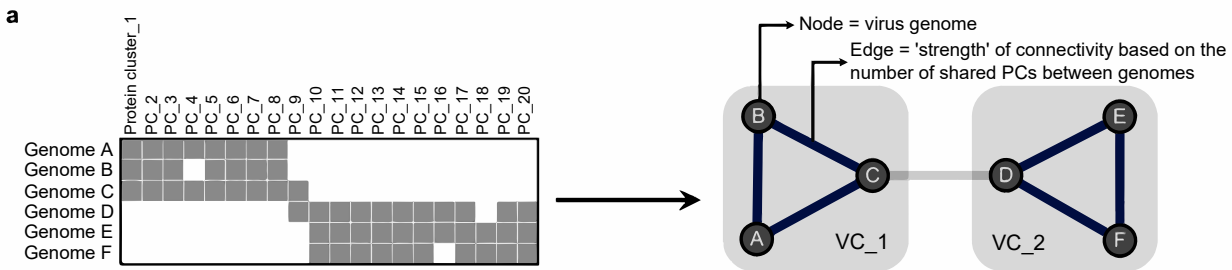
824 approximately bi-yearly with as complete genomes become available and resources exist to
825 support this effort.

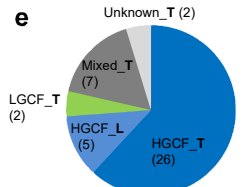
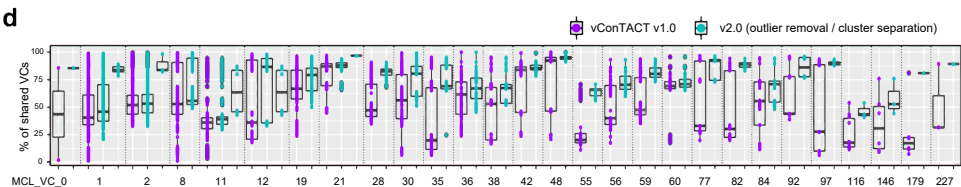
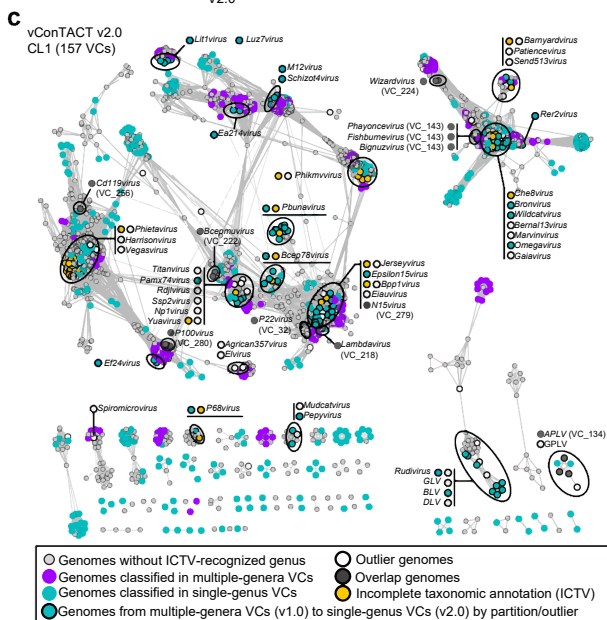
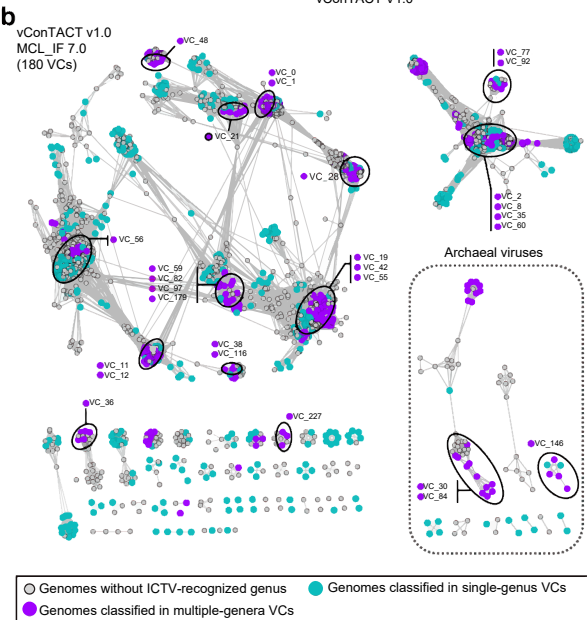
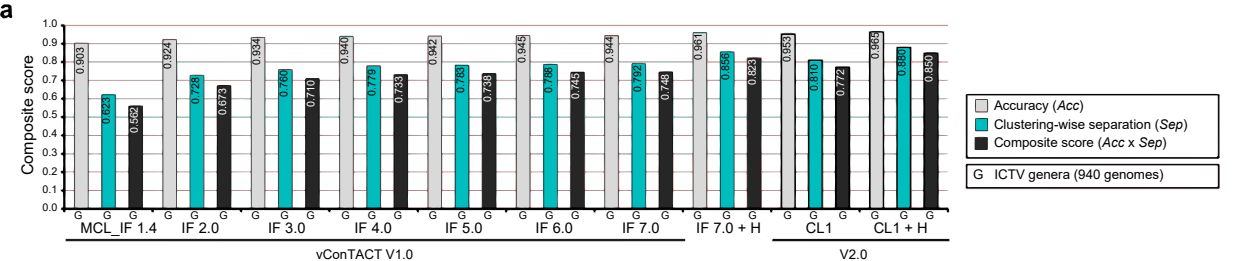
826 Methods-only References

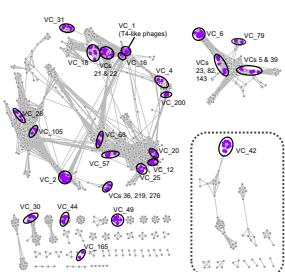
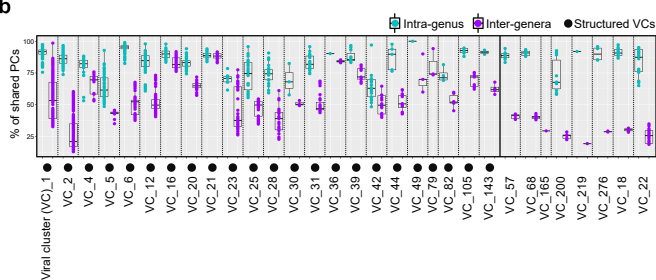
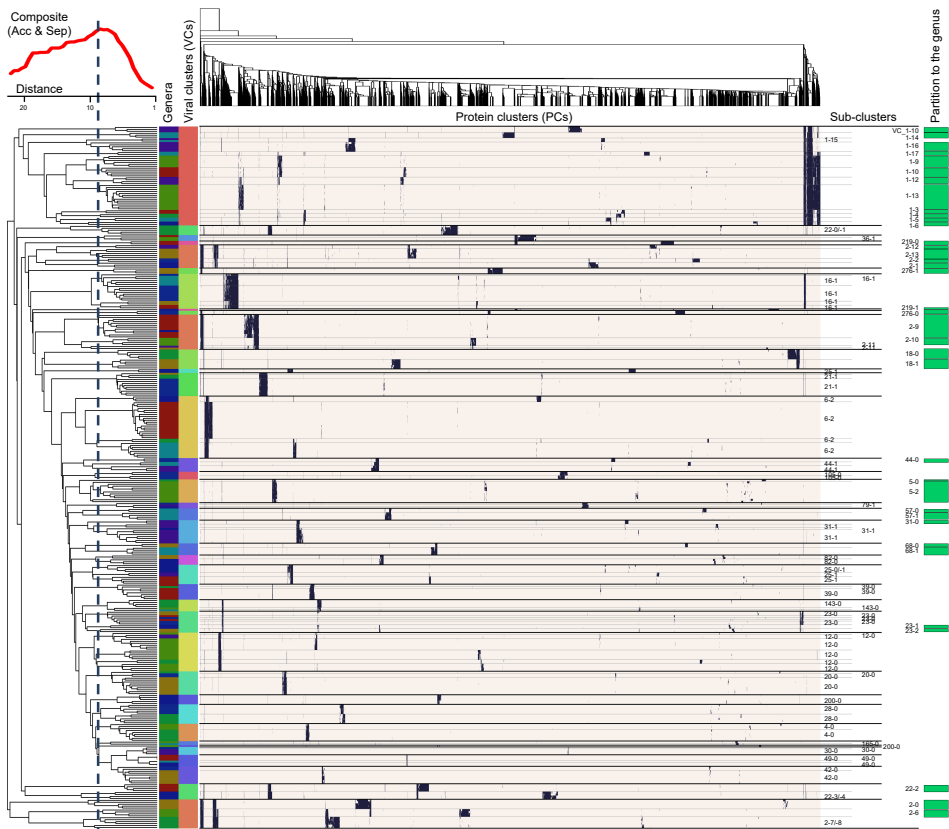
- 827 59. O’Leary, N. A. *et al.* Reference sequence (RefSeq) database at NCBI: Current status,
828 taxonomic expansion, and functional annotation. *Nucleic Acids Res.* **44**, D733–D745
829 (2016).
- 830 60. Krupovic, M. *et al.* Taxonomy of prokaryotic viruses: update from the ICTV bacterial and
831 archaeal viruses subcommittee. *Arch. Virol.* **161**, 1095–1099 (2016).
- 832 61. Adams, M. J. *et al.* Changes to taxonomy and the International Code of Virus
833 Classification and Nomenclature ratified by the International Committee on Taxonomy of
834 Viruses (2017). *Arch. Virol.* **162**, 2505–2538 (2017).
- 835 62. Altschul, S. F. *et al.* Gapped BLAST and PSI-BLAST: A new generation of protein
836 database search programs. *Nucleic Acids Research* **25**, 3389–3402 (1997).
- 837 63. Wiwie, C., Baumbach, J. & Röttger, R. Comparing the performance of biomedical
838 clustering methods. *Nat. Methods* **12**, 1033–1038 (2015).
- 839 64. Brohée, S. & van Helden, J. Evaluation of clustering algorithms for protein-protein
840 interaction networks. *BMC Bioinformatics* **7**, (2006).
- 841 65. Kamburov, A., Stelzl, U. & Herwig, R. IntScore: A web tool for confidence scoring of
842 biological interactions. *Nucleic Acids Res.* **40**, (2012).
- 843 66. Goldberg, D. S. & Roth, F. P. Assessing experimentally derived interactions in a small
844 world. *Proc. Natl. Acad. Sci.* **100**, 4372–4376 (2003).
- 845 67. Buchfink, B., Xie, C. & Huson, D. H. Fast and sensitive protein alignment using
846 DIAMOND. *Nat. Methods* **12**, 59–60 (2015).
- 847 68. Ohio Supercomputer Center . (1987).
- 848 69. Oliphant, T. E. SciPy: Open source scientific tools for Python. *Comput. Sci. Eng.* **9**, 10–20
849 (2007).
- 850 70. McKinney, W. Data Structures for Statistical Computing in Python. *Proc. 9th Python Sci.*
851 *Conf.* **1697900**, 51–56 (2010).
- 852 71. Pedregosa, F. *et al.* Scikit-learn: Machine Learning in Python. *J. Mach. Learn. Res.* **12**,
853 2825–2830 (2011).
- 854 72. Csárdi, G. & Nepusz, T. The igraph software package for complex network research.
855 *InterJournal Complex Syst.* **1695**, 1–9 (2006).
- 856 73. Federico, P., Pfeffer, J., Aigner, W., Miksch, S. & Zenk, L. Visual Analysis of Dynamic
857 Networks Using Change Centrality. in *2012 IEEE/ACM International Conference on*
858 *Advances in Social Networks Analysis and Mining* 179–183 (2012).
859 doi:10.1109/ASONAM.2012.39
860

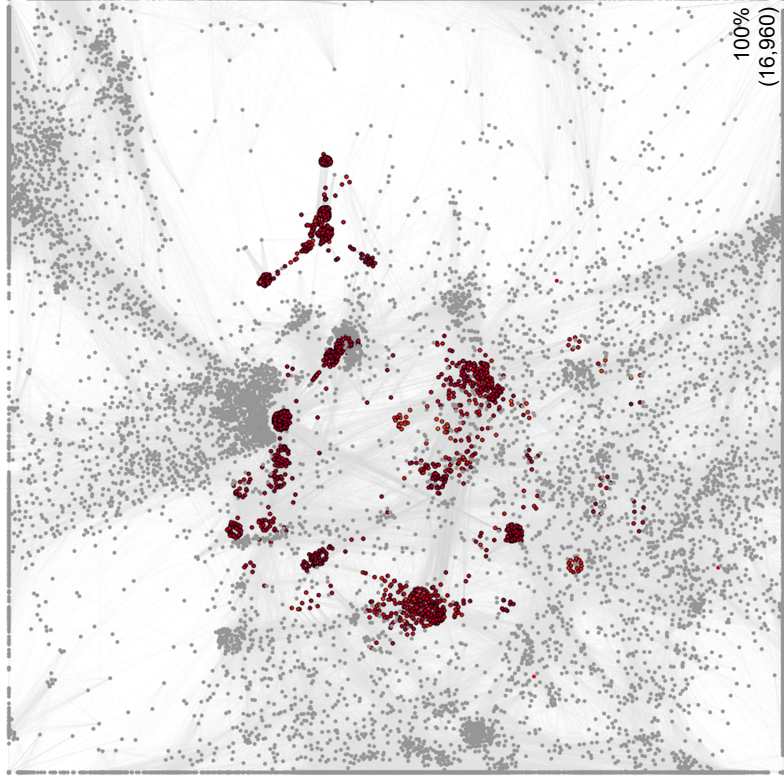
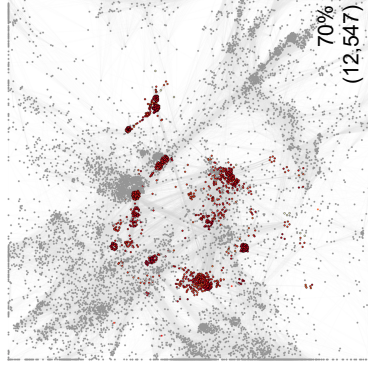
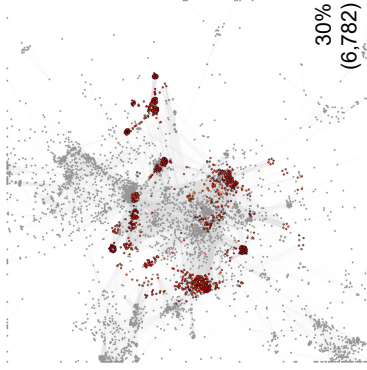
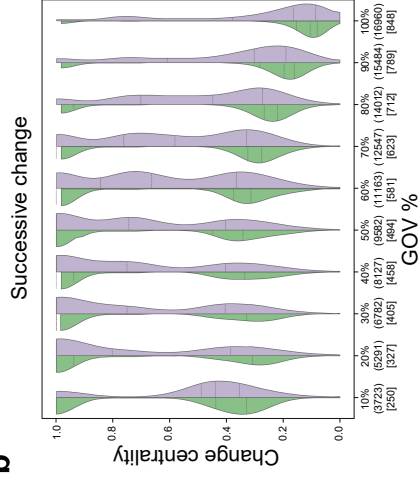
861
862
863
864
865

866
867
868





a**b****c**

a**b****c**