

1 **Gene sharing networks to automate genome-based prokaryotic viral taxonomy**

2 Ho Bin Jang^{1*}, Benjamin Bolduc^{1*}, Olivier Zablocki¹, Jens H. Kuhn², Simon Roux³, Evelien M.
3 Adriaenssens⁴, J. Rodney Brister⁵, Andrew M Kropinski^{6,7}, Mart Krupovic⁸, Dann Turner⁹ & Matthew B.
4 Sullivan^{1,10#}

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

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

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

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

10 ⁵ National Center for Biotechnology Information, National Library of Medicine, National Institutes of
11 Health, Bethesda, MD, USA

12 ⁶ Department of Pathobiology, Ontario Veterinary College, University of Guelph, Guelph, ON, Canada
13 N1G 2W1

14 ⁷ Department of Food Science, University of Guelph, Guelph, ON, Canada, N1G 2W1

15 ⁸ Institut Pasteur, Unité Biologie Moléculaire du Gène chez les Extrêmophiles, Department of
16 Microbiology, Paris 75015, France

17 ⁹ Centre for Research in Biosciences, Department of Applied Sciences, Faculty of Health and Applied
18 Sciences, University of the West of England, Coldharbour Lane, Bristol, UK

19 ¹⁰ Department of Civil, Environmental and Geodetic Engineering, Ohio State University, Columbus, OH,
20 USA

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

22 * These authors contributed equally to this work.

23

24 **ABSTRACT**

25 Viruses of bacteria and archaea are likely to be critical to all natural, engineered and human ecosystems,
26 and yet their study is hampered by the lack of a universal or scalable taxonomic framework. Here, we
27 introduce vConTACT 2.0, a network-based application to establish prokaryotic virus taxonomy that
28 scales to thousands of uncultivated virus genomes, and integrates confidence scores for all taxonomic
29 predictions. Performance tests using vConTACT 2.0 demonstrate near-identical correspondence to the
30 current official viral taxonomy (>85% genus-rank assignments at 96% accuracy) through an integrated
31 distance-based hierarchical clustering approach. Beyond “known viruses”, we used vConTACT 2.0 to
32 automatically assign 1,364 previously unclassified reference viruses to tentative taxa, and scaled it to

33 modern metagenomic datasets for which the reference network was robust to adding 16,000 viral contigs.
34 Together these efforts provide a systematic reference network and an accurate, scalable taxonomic
35 analysis tool that is critically needed for the research community.

36 **Main text**

37 Bacteria and archaea modulate the nutrient and energy cycles that drive ocean and soil ecosystems¹⁻⁴,
38 and impact humans by producing metabolites that alter health, behavior, and susceptibility to disease⁵.
39 Viruses that infect these microbes modulate these ‘ecosystem roles’ via killing, metabolic reprogramming
40 and gene transfer^{6,7}, with substantial impacts predicted in the oceans⁸⁻¹⁰, soils^{11,12} and human
41 microbiome^{13,14}. However, ecosystem-scale understanding is bottlenecked by the lack of universal genes
42 or methods that could facilitate a formalized taxonomy and comparative surveys. In fact, viruses do not
43 share a single gene¹⁵, and, thus, no analog to microbial 16S rRNA-based phylogenies and OTUs are
44 possible¹⁶.

45 Another potential challenge is that some viruses are prone to high rates of gene exchange (i.e.,
46 ‘rampant mosaicism’¹⁷), which, if broadly true, would stymie genome-based prokaryotic virus
47 taxonomy¹⁸. Fortunately, explorations of viral sequence space are revealing structure^{19,20} and population
48 genetic support for a biological species definition²¹, and new hypotheses to explain variable evolution
49 among prokaryotic viruses²². Such findings, alongside rapidly expanding viral genome databases, led the
50 International Committee on Taxonomy of Viruses (ICTV) to present a consensus statement suggesting a
51 shift from the traditional (i.e., phenotypic and genotypic criteria to classify viruses within community-
52 curated taxonomic ranks) approach²³ towards a genome-centered, and perhaps one-day, largely
53 automated, viral taxonomy²⁴. This shift is particularly critical given the modern pace of viral discovery in
54 which, as of March 2018, hundreds of thousands of metagenome-derived viral reference genomes and
55 large genome fragments (369,518 at IMG/VR²⁵) now dwarf the 26,223 available from prokaryotic virus
56 sequences in the NCBI GenBank database²⁶. Thus, evaluation of approaches to establish a scalable,

57 genome-based viral taxonomy is needed as the implementation of a commonly agreed-upon approach
58 available to the community would be highly desirable.

59 Multiple genome-based strategies have been proposed to develop such a unified bacterial^{15,27-32},
60 archaeal³³ and eukaryotic³⁴ virus taxonomic framework. For bacterial viruses (“phages”), the first
61 approach targeted phage relationships only by using complete genome pairwise protein sequence
62 comparisons in a phylogenetic framework (the “phage proteomic tree”) and was broadly concordant with
63 ICTV-endorsed virus groupings of the time¹⁵. Such efforts were not widely adopted, presumably because
64 (i) need was low (few metagenomics studies existed), and (ii) the paradigm was that “rampant
65 mosaicism” would blur taxonomic boundaries and violate the assumptions of the underlying phylogenetic
66 algorithms used in the analyses¹⁷. Other efforts sought to establish percent of genes shared and percent
67 identity of-shared genes cut-offs to define genera and sub-family affiliations^{35,36}, but lacked taxonomic
68 resolution for several virus groups. This lack of resolution was due to the likelihood that the mode and
69 tempo of prokaryotic virus evolution could vary significantly across the viral sequence space²². Building
70 upon a prokaryotic classification algorithm, the Genome Blast Distance Phylogeny (GBDP)³⁷, a freely
71 accessible online tool (VICTOR) now provides phage genomes for classification via combined
72 phylogenetic and clustering methods from nucleotide and protein sequences³⁰. Although a key advance,
73 this method suffers from limited scalability (100-genomes limit) and taxonomic assignment challenges
74 for the many novel, environmental viruses that lack genes shared with reference genomes.

75 Alternatively, several groups reasoned that the highly variable evolutionary rates across phage
76 sequence space could be examined through gene sharing networks^{28,29,38} to determine whether a
77 meaningful structure, and therefore taxonomic signal, occurs in this space. These networks, based on
78 shared protein clusters (PCs) between viral genomes, were largely concordant with ICTV-endorsed taxa
79 independent of whether monopartite²⁸ (a single node type, i.e., viral genomes) or bipartite networks^{33,38}
80 (two node types, i.e., viral genomes and genes) were used. Given these successes, we previously revisited
81 the monopartite gene sharing network approach to establish an iVirus³⁹ app (vConTACT) to automate a

82 network-based classification pipeline for prokaryotic virus genomes. Performance tests indicated that the
83 network analytics used by vConTACT produced viral clusters (VCs) that are ~75% concordant with
84 accepted ICTV prokaryotic viral genera, even with seven times more genomes now available²⁹. The
85 capacity to incorporate these genomes and accuracy of the network-based analytics have resulted in viral
86 taxonomy applications across large-scale studies of ocean^{40,41}, freshwater⁴² and soil⁴³, and studies of
87 single-virus amplified genomes (vSAGs)^{44,45}. vConTACT 1.0 was an important step forward but could
88 not be used for automatic tentative taxonomic assignments because (i) it creates artefactual clusters of
89 both under-sampled genomes (i.e., low number of genomes in a VC) and highly-overlapped regions of
90 sequence space among some genomes²⁹, and (ii) lacks several key, community-desired features such as
91 confidence metrics for the resultant VCs, a metric for establishing hierarchical taxonomy, and scalability.

92 Here we introduce and evaluate vConTACT v2.0, which updates the network analytics and feature set
93 of the original program. We apply this program to (i) establish a centralized, ‘living’ taxonomic reference
94 network as a foundational community resource and (ii) demonstrate that the updated vConTACT is robust
95 and scalable to modern datasets.

96

97 **RESULTS AND DISCUSSION**

98

99 *vConTACT 2.0 key features and updates*

100 The underlying goal of vConTACT is to automatically assign viral genomes into relevant established
101 or tentative taxa, with performance assessed relative to ICTV-assigned, manually-curated taxa. Viral
102 reference genomes of a single ICTV genus that are correctly grouped by vConTACT into a single viral
103 cluster (VC) are deemed ‘concordant VCs’. The original vConTACT 1.0 performed well in this area, with
104 ~75% of VCs corresponding to ICTV genera²⁹. However, ~25% of VCs did not match ICTV genera
105 (termed ‘discordant VCs’). These mismatches broadly represented three scenarios: (i) VCs that
106 encompass ICTV genera represented by 1-2 genomes (termed ‘undersampled VCs’), (ii) VCs that

107 encompass ICTV genera represented by virus genomes that shared many genes and/or modules with other
108 VCs (termed ‘overlapping VCs’), and (iii) VCs that encompass ICTV genera represented by virus
109 genomes that shared many genes and/or gene modules across genomes within the VC, and within subsets
110 of the genomes in the VC (termed ‘structured VCs’). Further, vConTACT 1.0 lacked several key features
111 to enable broader adoption and utility as described above.

112 To address these issues and establish vConTACT v2.0, we (i) implemented a new clustering
113 algorithm, (ii) established confidence scores and measures of distance-based taxon separation that are
114 crucial for hierarchical taxonomy, and (iii) optimized expansion to a large-scale viral metagenomic
115 dataset. Briefly, the clustering algorithm was upgraded from Markov cluster (MCL) to ClusterONE⁴⁶
116 (CL1), resulting in single parameter optimization (i.e., the inflation factor, IF) to determine VC generation
117 being converted to three processes to better disentangle confounding signals across problematic regions of
118 the networks (Online Methods). All three processes consider edge weight, (i.e., degree of connection
119 between genomes), to (i) identify outlier genomes, (ii) detect and separate genomes that bridge
120 overlapping VCs, and (iii) break down structured VCs into concordant VCs through distance-based
121 hierarchical clustering. In addition, to help differentiate between meaningful taxonomic assignments and
122 those that might be artefacts, each VC now receives a topology-based confidence score (value range 0-1),
123 which aggregates information about network topological properties, and a taxonomic (genus) prediction
124 score (value range 0-1), which estimates the likelihood of VCs to be equivalent to a single ICTV genus
125 (Online Methods). In both scores, higher values indicate either more confident linkages (topology-based
126 confidence score) or higher taxonomic agreement (taxonomic prediction score). Therefore, vConTACT
127 2.0 assigns taxonomy by a two-step clustering approach, in which VCs are first defined using CL1, and
128 then VCs are further subdivided using hierarchical clustering to maximize the taxonomic prediction score.
129 In such cases where VCs were further sub-divided, these are referred to as sub-VCs (benchmarking
130 below).

131

132 *Performance comparison of vConTACT versions 1.0 and 2.0*

133 To assess clustering performance of vConTACT v1.0 and v2.0 (hereafter ‘v1.0’ and ‘v2.0’,
134 respectively), we quantified ICTV correspondence from 336 comparisons (Online Methods) against all
135 available ICTV-classified archaeal and bacterial virus genomes (n=2,304, accessed January 2018).
136 Notably, though some combination of family, order, genus and species designations were available for all
137 of these viruses, only 41% (n=940) had genus-level classifications (**Supplementary Table 1**). Our
138 performance comparisons focused on this subset of classified genomes. Composite performance, the sum
139 of six metrics (cluster-wise sensitivity, Sn ; positive prediction value, PPV ; geometric accuracy of Sn and
140 PPV , Acc ; cluster-wise separation, Sep_{cl} ; complex (ICTV taxon)-wise separation Sep_{co} ; and geometric
141 mean of Sep_{cl} and Sep_{co} , Sep) was used to assess overall performance of v1.0 and v2.0 (**Fig. 1a**). Each of
142 these metrics has values range from 0 to 1 with 1 indicating perfect clustering accuracy and/or coverage
143 (Online Methods). We found that v1.0 organized the 2,304 analysed viral genomes into 305 VCs at its
144 best inflation factor (IF=7), and 77.5% of these were concordant at the genus rank, whereas v2.0
145 identified 279 VCs, and 79.2% of these were concordant at the genus rank (**Supplementary Table 2**).
146 Moreover, we added to v2.0 a post-processing, Euclidean distance-based hierarchical clustering step to
147 split mismatched VCs. This step accurately and automatically classified 36 additional genera from
148 structured VCs (**Supplementary Table 1**), resulting in the highest composite score of 5.4 (maximum
149 achievable score of 6.0) at the genus rank, with a concordance of 85.0% and accuracy of 96.4%. (**Fig. 1a**
150 and **Supplementary Table 2**). Together, these findings suggest that both upgrading the clustering
151 algorithm and adding hierarchical clustering were critical to improve automatic VC designations.

152 Next, we assessed how v2.0 handled areas of the reference network that represented discordant VCs.
153 First, 55% of ICTV genera are undersampled (**Supplementary Table 1**), which in a gene-sharing
154 network manifests as weakly connected, small VCs prone to artefactual clustering. In v1.0, undersampled
155 VCs accounted for 64% (28/44) of all discordant VCs, and they could not be resolved by increasing IF
156 values (**Fig. 1b and d** and **Supplementary Table 1**). In contrast, v2.0 automatically and accurately

157 handled these same 28 undersampled VCs (comprising 60 genomes) by splitting the 37 problematic
158 genera into 22 outliers (i.e., genera with only one member) and correctly placing the remaining 38
159 genomes from 15 genera into 15 VCs (**Fig. 1c and d** and **Supplementary Table 1**). Thus, in instances in
160 which v1.0 performed poorly on undersampled VCs, v2.0 was able to resolve all undersampled VCs into
161 their appropriate ICTV genera.

162 Second, we evaluated the ability of v2.0 to handle overlapping VCs, which share more genes across
163 VCs than expected, presumably due to gene exchange that could erode structure in the network. In v1.0,
164 overlapping VCs could not be identified. In v2.0 we automated their detection via a ‘match coefficient’
165 between each VC that measured the connection within- and between- other VCs, and sensitivity analyses
166 established a maximum cluster overlap value of 0.8 as diagnostic (Online Methods). In this way, nine
167 overlapping VCs (ICTV-classified genera only) were detected. These clusters contained 30 viruses across
168 11 ICTV genera, which included viruses with known mosaic genomes⁴⁷ (e.g., lambdoid or mu-like phages
169 of the *P22virus*, *Lambdavirus*, *NI5virus*, and *Bcepnavirus* genera), temperate phages^{48,49} (i.e.,
170 *Mycobacterium* phages of the *Bignuzvirus*, *Phayoncevirus*, and *Fishburnevirus* genera and *Gordonia*
171 phages of the genus *Wizardvirus*), and three newly-established genera (i.e., *Cd119virus*, *P100virus* and
172 archaeal *Alphapleolipovirus*), all bearing low topology-based confidence scores (averages of 0.29 for
173 these VCs versus 0.50 for concordant VCs; P-value = 2.09e-08, Mann-Whitney U test) (**Supplementary**
174 **Fig. 1**). Interestingly, this set of viruses within overlapping VCs (74 in total, including non-classified
175 genomes from ICTV) contained 31 phages having a high gene content variation due to extensive gene
176 flow (HGCF, **Fig. 1e**), related to the recently proposed framework of phage evolutionary lifestyles²².
177 Further, these VCs contained highly recombinogenic temperate phages, more likely to exchange genes as
178 opposed to low gene content flux (LGCF) phages that follow a predominantly lytic life cycle
179 (**Supplementary Fig. 1b**). Thus, this observation may indicate a high linkage between overlapping
180 genomes and phages with high gene flow. Although unresolvable in v1.0, v2.0 could assign eight of the
181 11 ICTV genera (24 viruses) into eight ICTV-concordant VCs (**Supplementary Table 1**). The remaining

182 three ICTV genera, all comprised of *Mycobacterium* phages⁵⁰ (six genomes), could not be resolved. This
183 lack of resolution is presumably due to high gene flow resulting from a predominantly temperate lifestyle
184 that is associated with an exceptionally high fraction (avg = 69%) of genes shared across VCs
185 (**Supplementary Table 3**). Undoubtedly, these genomes are the most challenging to classify, and may
186 not be amenable to automated taxonomy. Whether such highly recombinogenic genomes are the
187 exception or the norm across environments is unknown.

188 Third, structured VCs contained genomes that our gene sharing networks placed into a single VC
189 (due to many shared genes and/or gene modules across all the member genomes), whereas ICTV
190 delineated multiple genera (due to subsets of the genomes also sharing additional genes). V1.0
191 qualitatively and selectively handled these structured VCs via decomposing hierarchical patterns of gene
192 sharing²⁷. In v2.0, we formalized an optimized, quantitative hierarchical decomposition distance measure
193 (9.0, Online Methods, **Fig. 2c**, and **Supplementary Fig. 2**) that maximized composite scores of two
194 geometric mean values of performance metrics (*Acc* and *Sep*; Online Methods) that divide discordant VCs
195 into concordant (to ICTV genera) sub-VCs, and used this distance as a generalized threshold. In the v2.0
196 network, 31 discordant VCs contained 101 phage and two archaeal virus genera, in which 23 (74%) were
197 structured VCs spanning 86 genera (**Fig. 2a,b** and **Supplementary Table 1**). This v2.0 approach resolved
198 30% (26 of 86) of these ICTV genera from 6 of the 23 structured VCs (**Fig. 2c**). Curiously, one such
199 structured VC was comprised of T4-like phages (of which nine out of ten T4-related genera were resolved;
200 **Supplementary Note 1**), in which hierarchical ‘T4 core’ and ‘cyano T4 core’ gene sets are well
201 documented⁵¹. In our networks, the T4-like phages represent a single VC, but with sub-VCs that are
202 consistent with ICTV-established genera (VC 1 in **Fig. 2c** and **Supplementary Table 1**). Extrapolating
203 from this network, we interpret structured VCs to represent areas of viral sequence space that are well-
204 sampled to the point that the core gene sets that define a virus (capsid, tail, replication machinery)
205 establish the VC in the network, whereas ecologically diverse viral genomes within the VC reveal
206 structure due to niche-defining genes that represent adaptation to diverse environments and/or hosts. We

207 posit that the 19 structured VCs that cannot be resolved towards ICTV concordance (**Fig. 2c** and
208 **Supplementary Table 1**), represent either regions of the network where niche-defining genomic
209 information is lacking or may require complementary phenotypic or evolutionary evidence to establish
210 ICTV genera, as done for the archaeal fuselloviruses (VC42) and bacterial microviruses (VCs 30 and 49).
211 Thus, whether these structured VCs result from lack of resolution in v2.0 or from genera needing ICTV
212 revision remains an outstanding question.

213 Finally, given such strong performance, we suggest that this gene sharing network already offers
214 significant new taxonomic insights. First, as described earlier, only 41% of the 2,304 reference virus
215 genomes are classified by ICTV at the genus rank. Thus, we propose that the remaining 1,364 currently
216 genus-unclassified reference viruses, which organized into 304 well-supported hierarchically decomposed
217 sub-VCs (**Supplementary Table 1**), represent genomes from *bona fide* novel virus genera. This finding,
218 if officialised, immediately doubles established viral taxonomy and invites a framework for manual
219 curation of these automatic assignments, which in itself will improve future vConTACT analytic
220 performance. As first evidence of the value of such an iterative process, we note that v2.0 clustering
221 suggested an alternative taxonomy among ten current ICTV genera: *Barnyardvirus*, *Bcep78virus*,
222 *Bpp1virus*, *Che8virus*, *Jerseyvirus*, *P68virus*, *Pbunavirus*, *Phietavirus*, *Phikmvvirus*, and *Yuavirus*
223 (**Supplementary Fig. 3** and **Note 2**), and manual inspection had already recommended some of these
224 ICTV genera be revised (e.g. *Phikmvvirus* viruses, ICTV proposal 2015.007a-Db). An automated
225 vConTACT-based approach would systematically identify such problematic taxa and drastically speed up
226 these critical revisions as new data become available.

227
228 *vConTACT v2.0 is scalable to modern virome datasets*

229 A major bottleneck regarding automated taxonomic assignments is the ability to robustly integrate
230 large sets of newly discovered virus genomes. To evaluate this concern, we added ~16K curated viral
231 genomes and large genome fragments from the Global Ocean Virome (GOV) dataset⁴⁰ to our reference

232 network. We added these genomes and genome fragments in successive 10% increments (i.e., 0%-
233 10%[,...], 0%-100%), to assess the impact of various data scales on the reference network stability of VC
234 assignments. Network changes were tracked by assessing (i) network performance metrics (*Sn*, *Acc* &
235 *PPV*, as above), (ii) ‘normalized mutual information’(NMI), as a measure of VC similarity (values range
236 from 0 to 1 with 0 indicating that none of the original member genomes within a VC remained in that
237 same VC and 1 indicating that all members in a VC remain in that same VC across time), and (iii)
238 ‘change centrality’ (CC), reflecting how much each node’s connections changed as more sequences were
239 added to the network (values range from 0-1 with 0 indicating no change and 1 indicating complete
240 change), classified over three ‘change intensity’ groups: low (0 - 0.283), medium (0.283 - 0.506) and high
241 (0.506 - 0.999) groups (Online Methods). Although CC indicates changes in connections between nodes,
242 these may still remain in a given VC, albeit re-shuffled. Together, NMI and CC assess the impact of
243 additional data on the network clusters and topology, respectively, while *Sn*, *Acc* and *PPV* assess
244 concordance with ICTV taxonomy.

245 All measures indicated that most network changes occurred with early additions of the novel GOV
246 data (up to 20-30% of the dataset), with the network largely stabilized after that (**Fig. 3**). For example,
247 *Acc* (mean value of *Sn* and *PPV*) is reduced by 12% when only using 20% of the GOV data, but stabilizes
248 at a ~7% decrease (**Supplementary Fig. 4**); similar responses were observed in NMI (**Fig. 3b**). This
249 initial drop appears driven by formation of novel, undersampled VCs, a disruptive effect similarly
250 observed with undersampled ICTV genera bearing low quality or confidence in VC membership. With
251 more data, undersampled VCs reach ‘saturation’, which increases confidence scores for these new VCs
252 and buffers from further disruption. This stabilization is likely due to strong intra-cluster forces (within
253 VCs) vastly out-weighing inter-cluster forces (between VCs).The lasting minimal decrease represents the
254 novelty of sequence space in GOV relative to RefSeq and the fact that these additions are commonly large
255 genome fragments rather than complete genomes. Sequential CC analysis showed minimal impact on the
256 RefSeq network structure and VC membership, as 85% of reference genomes had low-to-medium change,

257 whereas 0.05% of genomes experienced high change. The remaining 15% were classified as either
258 singleton, outlier, or overlaps. These data support a similar pattern as NMI fluctuations (**Fig. 3d**): as data
259 accumulated, fewer and fewer nodes or VCs were impacted due to new data influencing only pre-existing
260 areas in the network. Therefore, as a network grows in scale, adding new data mostly similar to pre-
261 existing data will have minimal impact on the underlying network structure (e.g., adding new marine data
262 to a marine network), as newly added data is already “represented,” whereas utterly novel data will
263 generate novel VCs and increase CC values. Indeed, most unaffected VCs (CC = 0) were non-marine or
264 soil in origin e.g., *Andromedavirus* viruses, *Saetivirus* viruses, two archaeal viruses (Methanobacterium
265 virus psiM2, Methanothermobacter virus psiM100), *Thermus* phages, or cyanobacterial mat viruses.

266

267 As contigs accumulate, the number of VCs also increases linearly ($R^2 = 0.998$, P-value = 1.2×10^{-12}).
268 We examined whether GOV data may partially resolve ICTV outlier and singleton genomes. More data
269 should create new connections to singletons, whereas outliers may get connected to new or existing VCs.
270 Out of 38 single-member VCs of singleton and outlier genomes (**Supplementary Fig. 5**), three
271 *Mycobacterium* phages clusters were improved, with two other *Mycobacterium* viruses genomes merged
272 into six-genera heterogeneous VCs. Together, this analysis suggests that v2.0’s underlying methodology
273 is sufficiently robust to handle large amounts of data. With 100% of GOV added (16,960 total contigs),
274 919 new VCs are created, representing potentially 919 new viral genera over existing RefSeq genomes.

275

276 *Community availability and future needs*

277 The utility of v2.0 depends upon its expert evaluation and community availability. To maximize this
278 evaluation, members of the ICTV Bacterial and Archaeal Viruses Subcommittee were invited as co-
279 authors to critique the work, and we made the resulting optimized tool available in two ways. First, the
280 source code is available through Bitbucket (<https://bitbucket.org/MAVERICLab/vcontact2> as a
281 downloadable python package. Second, v2.0 is available as an app through iVirus³⁹, the viral ecology

282 apps and data resource embedded in the CyVerse Cyberinfrastructure, with detailed usage protocols
283 available through Protocol Exchange (<https://www.nature.com/protocolexchange/>) and protocols.io
284 (<https://www.protocols.io/>). Finally, the curated reference network is available at each of these sites.

285 Although v2.0 performance metrics are strong and provide a critically needed, systematic reference
286 viral taxonomic network, limitations still remain. First, our reference network needs to be rebuilt each
287 time new data are added. Avoiding this reconstruction step will require the development of approximation
288 methods and/or a placement algorithm (akin to PPlacer for 16S phylogenies⁵²) to incorporate new data.
289 Second, although v2.0 handles reference prokaryotic virus genomes (including ssDNA or dsDNA phages)
290 and large GOV genome fragments, this framework has not been designed, tested or validated for
291 eukaryotic viruses, which pose unique computational challenges³⁴. Third, shorter prokaryotic virus
292 genomes and genome fragments (e.g., ≤ 3 PCs or ≤ 5 genes) are of low statistical power in the v2.0
293 framework, and will require new solutions to establish higher confidence VCs. Fourth, genomes
294 identified as singletons, outliers or overlapping are currently excluded from the gene-sharing network.
295 Although singletons and outliers can be resolved by the addition of new data, overlapping VCs can
296 remain challenging to resolve, particularly for the HGCF phages²² that are highly recombinogenic. Such
297 rampantly mosaic virus genomes are problematic for viral taxonomy. However, they are identifiable in
298 the networks and, at least to date, represent the minority of known viral sequence space. Most (~75%) are
299 LGCF viruses that remain amenable to automated genome-based viral taxonomy. Whether this situation
300 will remain so awaits further exploration of viral sequence space—particularly where temperate phages
301 may predominate (e.g., soils⁵³, human gut⁵⁴). For now, we propose vConTACT 2.0 as a tool that offers a
302 robust, systematic and automatic means to aid the classification of bacterial and archaeal viruses.

303

304 **METHODS**

305

306 **Data sets.** Full-length viral genomes were obtained from the National Center for Biotechnology
307 Information (NCBI) viral reference dataset^{26,55} ('ViralRefSeq', version 85, as of January, 2018),
308 downloaded from NCBI's viral genome page (<https://www.ncbi.nlm.nih.gov/genome/viruses/>) and
309 eukaryotic viruses were removed. The resulting file contained a total of 2,304 RefSeq viral genomes
310 including 2,213 bacterial viruses and 91 archaeal viruses (**Supplementary Table 1**). In parallel, the ICTV
311 taxonomy (ICTV Master Species List v1.3, as of February, 2018) was retrieved from the ICTV homepage
312 (<https://talk.ictvonline.org/files/master-species-lists/>). ICTV-classifications were available for a subset of
313 genomes at each taxonomic rank, and final dataset included; 884 viruses from two orders, 974 viruses
314 from 23 families, 363 viruses from 28 subfamilies, and 975 viruses from 264 genera. To maintain
315 hierarchical ranks of taxonomy, we manually incorporated 2016 and 2017 ICTV updates⁵⁶⁻⁵⁸ to NCBI
316 taxonomy when ICTV taxonomy was absent.

317
318 **Network construction.** A total of 231,165 protein sequences were extracted from the 2,304 viral
319 genomes (above). To group protein sequences into homologous protein clusters (PCs)²⁹, all proteins were
320 subjected to all-to-all BLASTP⁵⁹ searches (default parameters, cut-offs of $1E^{-5}$ on e-value and 50 on bit
321 score). A subsequent application of the MCL with inflation factor 2.0 grouped 204,540 protein sequences
322 into 25,510 PCs, with the remaining 26,625 proteins being to singletons (those that do not have close
323 relatives). The resulting output was parsed in the form of a matrix comprised of genomes and PCs (i.e.,
324 $2,304 \times 25,510$ matrix). We then determined the similarities between genomes by calculating the
325 probability of finding a common number of PCs between each pair of genomes, based on the following
326 hypergeometric equation as per Lima-Mendez et al²⁸:

327
328
$$P(X \geq c) = \sum_{i=c}^{\min(a,b)} \frac{C_a^i C_{n-a}^{b-i}}{C_n^b} \quad (1)$$

329

330 in which c is the number of PCs in common; a and b are the numbers of PCs and singletons in genomes A
331 and B, respectively; and n is the total number of PCs and singletons in the dataset. A score of similarity
332 between genomes was obtained by taking the negative logarithm (base 10) of the hypergeometric P-value
333 multiplied by the total number of pairwise genome comparisons (i.e., $2,304 \times 2,303$). Genome pairs
334 with a similarity score ≥ 1 were previously shown to be significantly similar through permutation test of
335 PCs and/or singletons between genomes²⁹. Afterwards, a gene (protein)-sharing network was constructed,
336 in which nodes are genomes and edges connect significantly similar genomes. This network was
337 visualized with Cytoscape software (version 3.6.0; <http://cytoscape.org/>), using an edge-weighted spring
338 embedded model, which places the genomes sharing more PCs closer to each other.

339
340 **Parameter optimization of vConTACT v1.0 and 2.0.** Due to different criteria for parameter
341 optimization between the clustering methods, different number and size of the clusters are often generated,
342 which can make objective performance comparisons difficult⁶⁰. Thus, to more comprehensively compare
343 performance, v1.0's MCL-based VCs were generated at inflation factors (IFs) of 2.0 to 7.0 by 1.0
344 increments, with an optimal IF of 1.4 showing the highest intra-cluster clustering coefficient (ICCC)²⁸
345 (**Supplementary Table 1** and **Supplementary Fig. 6**). CL1, which was incorporated into a new version
346 of vConTACT (v2.0), operates in multiple stages of complex detection⁴⁶. Unlike the MCL that uses a
347 single parameter²⁸, CL1 uses a set of parameters, which can act as the threshold for each stage of complex
348 detection. For example, as four main parameters of CL1, the minimum density, node penalty, the haircut,
349 and the overlap automatically quantifies (i) the cohesiveness of cluster, (ii) the boundaries of the clusters
350 (outliers), and (iii) the size of overlap between clusters, respectively⁴⁶. Of these parameters, the first two
351 are used to detect the coherent groups of VCs as follows:

352

$$353 \quad C = \frac{W_{in}(V)}{W_{in}(V) + W_{out}(V) + p|C|} \quad (2)$$

354

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

358 As another parameter of CL1, the haircut can find loosely connected regions of the network (outliers) by
359 measuring the ratio of connectivity of the node g within the cluster c to that of its neighbouring node h as:

360

$$361 \quad \Delta_{out} = k \sum_{j=1}^l W_{h,j} / \sum_{i=1}^k W_{g,i} \quad (3)$$

362

363 in which k is the number of edges of the node g , and W is the total weight of edges of the respective
364 nodes g and h . If the total weight of edges from a node (h) to the rest of the cluster (c) is less than x times
365 that we specified the average weight of nodes (g) within the given cluster, CL1 will remove the node (h)
366 from a given VC and place it into the outlier.

367 Additionally, CL1 can specify the maximum allowed overlap (ω) between two clusters, measured by the
368 match coefficient, as follow:

369

$$370 \quad \omega = i^2 / a * b \quad (4)$$

371

372 in which i is the size of overlap, which is divided by the product of the sizes of the two clusters under
373 consideration (a and b). Since CL1 identifies overlap between VCs, it can consequently find both
374 hierarchical and overlapping structures of viral groups. This capability is a significant improvement over
375 v1.0, given v1.0's MCL cannot handle modules with overlaps⁷. Specifically, CL1 (i) finds cluster(s)
376 having less than maximum value of specified overlap threshold (above) and (ii) merges these clusters
377 together with their interacting cluster(s) to make the results easier to interpret. Thus, in the resulting
378 output file, viral groups (or clusters) having the identical member viruses can be found in multiple
379 clusters, called 'overlapping clusters' (**Supplementary Table 1**). CL1 was run with varying conditions

380 for these four parameters (minimum density ranging from 0 to 1 by 0.1 increments; node penalty from 1
381 to 10 by 1.0; haircut from 0 to 1 by 0.05; overlap from 0 to 1 by 0.05) and default settings for other
382 parameters: 2 as minimum cluster size, weighted as edge weight, single-pass as merging, unused nodes as
383 seeding. We therefore obtained a total number of 53,361 clustering results, which we evaluated
384 individually to yield the highest performance on taxonomic data set (above), in terms of geometric mean
385 value of prediction accuracy (*Acc*) and clustering-wise separation (*Sep*, see next section), as previously
386 described⁶¹ We then used minimum density = 0.3, node penalty = 2.0, haircut = 0.65, and overlap = 0.8 to
387 derive the final set of clusters, resulting in a total of 279 VCs (**Supplementary Table 1**). As a post-
388 clustering step of v2.0, all VCs including discordant clusters (those comprising ≥ 2 taxa) were further
389 hierarchically separated into sub-clusters using the unweighted pair group method with arithmetic mean
390 (UPGMA) with pairwise Euclidean distances implemented in Scipy. To optimize the distance-based sub-
391 clustering of VCs, we assessed the distances of sub-clusters across all the VCs. These distances (ranging
392 from 1 to 20 in 0.5 increments) maximized the geometrical mean values of the prediction accuracy (*Acc*)
393 and clustering-wise separation (*Sep*) at the ICTV genus rank (see next section). This optimization resulted
394 in the distance of 9.0 yielding the highest composite score of *Acc* and *Sep* (**Supplementary Fig. 2**).
395 Notably, vConTACT v2.0 was designed to help users optimize (i) parameters for grouping of
396 genomes/contigs into VCs and (ii) distance for post-decomposition of VCs into sub-clusters. This tool
397 automatically evaluates the robustness of VCs and sub-clusters, respectively, based on the external
398 performance evaluation statistics (below).

399

400 **Performance comparison between vConTACT v1.0 and v2.0.** Since the external measures such as
401 precision, recall, and others often neglect overlapping clusters, which might not reflect the true
402 performance of CL1, we used 6 external quality metrics that were successfully used for performance
403 comparison between MCL and CL1⁶¹ (see below). Specifically, the performance of v1.0 (MCL) and v2.0
404 (CL1 alone and CL1 + hierarchical sub-clustering, respectively) were evaluated based on : (i) cluster-wise

405 sensitivity, Sn (ii) positive predictive value, PPV (iii) geometric accuracy of Sn and PPV , Acc (iv) cluster-
406 wise separation, Sep_{cl} (v) complex (ICTV taxon)-wise separation Sep_{co} , and (vi) geometric mean of Sep_{cl}
407 and Sep_{co} , Sep . As an internal parameter, we computed the intra- and inter-cluster proteome similarities
408 (fraction of shared genes between genome that are within the same VCs and different VCs, respectively).
409 For vConTACT v1.0, clustering result yielding the highest clustering accuracy value (inflation of 7.0)
410 was subsequently used for comparison to v2.0's clusters and sub-clusters.
411 To generate six external measures, we first built a contingency table T , in which row i corresponds to the
412 i^{th} annotated reference complex (i.e., ICTV-recognized order, family, subfamily, or genus), and column j
413 corresponds to the j^{th} predicted complex (i.e., sub-/clusters). The value of a cell T_{ij} denotes the number of
414 member viruses in common between the i^{th} reference complex and j^{th} predicted complex. Here, N_i is the
415 number of member viruses belonging to reference complex n . Sn and PPV are then defined as follows:

416

$$417 \quad Sn = \frac{\sum_{i=1}^n \max_j \{T_{ij}\}}{\sum_{i=1}^n N_i} \quad (5)$$

418

$$419 \quad PPV = \frac{\sum_{j=1}^n \max_i \{T_{ij}\}}{\sum_{j=1}^n T_{ij}} \quad (6)$$

420

421 Generally, higher Sn values indicate a better coverage of the member viruses in the real complexes,
422 whereas higher PPV values indicates that the predicted clusters are likely to be true positives. As a
423 summary metric, the Acc can be obtained by computing the geometrical mean of the Sn and PPV values:

424

$$425 \quad Acc = \sqrt{Sn \times PPV} \quad (7)$$

426

427

428 With the same contingency table used for S_n , PPV , and Acc , we calculated the averages of complex-wise
429 separation Sep_{co} , and cluster-wise separation Sep_{cl} , respectively, below:

430

$$431 \quad Sep_{co} = \frac{\sum_{i=1}^n Sep_{co_i}}{n} \quad (8)$$

432

$$433 \quad Sep_{cl} = \frac{\sum_{j=1}^m Sep_{cl_j}}{m} \quad (9)$$

434

435 High Sep_{co} and Sep_{cl} , (both have maximal values of 1.0) indicate how well a given complex is isolated
436 from the other complexes and a cluster from other clusters, respectively. To estimate these separation
437 results as a whole, the geometric mean (clustering-wise separation; Sep) of Sep_{co} and Sep_{cl} was computed:

438

$$439 \quad Sep = \sqrt{Sep_{co} \times Sep_{cl}} \quad (10)$$

440

441 High clustering-wise separation values indicate a bidirectional correspondence between a sub-/cluster and
442 each ICTV taxon: maximal value of 1.0 can be obtained when a sub-/cluster corresponds perfectly to each
443 taxon.

444 As an internal measure, the fraction of PCs²⁹ between two genomes (i.e., proteome similarity) was
445 computed by using the geometric index (G). The proteome similarity was estimated as:

446

$$447 \quad G_{AB} = \frac{|N(A) \cap N(B)|}{|N(A)| \times |N(B)|} \quad (11)$$

448

449 in which $N(A)$ and $N(B)$ indicate the number of PCs in the genomes of A and B, respectively. A total of
450 400,234 pairs of genomes with >1% proteome similarity are shown in **Supplementary Table 3**.

451

452 **Clustering-based confidence score.** To generate the confidence score per sub-cluster, we used four
453 confidence scoring methods, as previously described^{62,63}, with some modifications. Three of them exploit
454 the network topology properties by assessing (i) the significance of clustering coefficient, (ii) the weight
455 of cluster quality, and (iii) the probability of cluster quality. We then used combined these three values
456 into an aggregate topology-based confidence score.
457 Specifically, for the significance of the clustering coefficient, we quantified the fidelity (F) of the edge (p)
458 by calculating cumulative hypergeometric P- values using Equation 1 (above) between sub-clusters. The
459 fidelity values are lower (close to 0) for the genomes having the higher number of shared genes. We then
460 defined the confidence of sub-cluster cohesiveness as the product of the fidelity values of total edges (i.e.,
461 $p1$ and $p2$) within the sub-cluster c as below:

$$462$$
$$463 \text{ Confidence } (c) = F_{p1,c} \times F_{p2,c} \quad (12)$$
$$464$$

465 For the second scoring method, we computed the quality (Q) of sub-cluster (c) as:

$$466$$
$$467 Q_c = W_{in}/W_{in} + W_{out} \quad (13)$$
$$468$$

469 in which W_{in} and W_{out} are the total weight of edges that lie within sub-cluster c and across others,
470 respectively. For the third method, we evaluated the P-value of a one-sided Mann-Whitney U test for in-
471 weights and out-weights of sub-clusters. The rationale behind this test is that sub-clusters with a lower P-
472 value contains significantly higher in-weights than out-weights, thus indicative that a formed sub-cluster
473 is valid, and not a random fluctuation. All pairs of three values above were then incorporated into the
474 topology-based confidence score with the Spearman rank correlation coefficient by using in-house python
475 scripts and Scipy. Along with this confidence score, we quantified the likelihood that each sub-cluster

476 corresponds to an ICTV-sanctioned genus (or equivalent) by using distance threshold that are specified at
477 the ICTV genus rank, which we refer to as “taxon predictive score”. This score can be calculated as:

478

$$479 \quad prediction = \sum l_{i,j} / l_c \quad (14)$$

480

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

487

488 **Measuring effect of GOV on network structural changes.** GOV contigs (14,656) were added in 10%
489 increments (randomly selected at each iteration) to NCBI Viral RefSeq and processed using vConTACT
490 2.0 with one difference – Diamond⁶⁴ instead of BLASTp was used to construct the all-versus-all protein
491 comparison underlying the PC generation. Once generated, vConTACT 2.0 networks were post-processed
492 using a combination of the Scipy⁶⁵, Numpy, Pandas⁶⁶ and Scikit-learn⁶⁷ python 3.6 packages. Networks
493 were rendered using iGraph⁶⁸. To calculate NMI, each network’s genomes and their VC membership was
494 compared in pairwise fashion to all other networks using the “adjusted mutual info score” function of
495 Scikit-learn. Intra-cluster distances were calculated using the agglomerative clustering functions
496 “linkage” with distance calculated from shared PCs using the cluster average (also known as UPGMA),
497 and novel clusters identified using the “fcluster” function of Scipy’s hierarchical clustering. In parallel,
498 the method to calculate change centrality was calculated as described previously⁶⁹. CCs were calculated in
499 a successive way, in which each addition was compared to Viral RefSeq 85 independently of other
500 additions (0% versus 10%, 0% vs 20%, [...], 0% vs 100%).

501 **Code availability.** The vConTACT v2.0 package is freely distributed through Bit Bucket as a python
502 package (<https://bitbucket.org/MAVERICLab/vcontact2>).

503

504 REFERENCES

- 505 1. Falkowski, P. G., Fenchel, T. & Delong, E. F. The microbial engines that drive earth's
506 biogeochemical cycles. *Science* **320**, 1034–1039 (2008).
- 507 2. Sunagawa, S. *et al.* Structure and function of the global ocean microbiome. *Science* (80-.). **348**,
508 (2015).
- 509 3. Moran, M. A. The global ocean microbiome. *Science* **350**, (2015).
- 510 4. Zhao, M. *et al.* Microbial mediation of biogeochemical cycles revealed by simulation of global
511 changes with soil transplant and cropping. *ISME J.* **8**, 2045–2055 (2014).
- 512 5. Cho, I. & Blaser, M. J. The human microbiome: At the interface of health and disease. *Nature*
513 *Reviews Genetics* **13**, 260–270 (2012).
- 514 6. Fernández, L., Rodríguez, A. & García, P. Phage or foe: an insight into the impact of viral
515 predation on microbial communities. *ISME Journal* 1–9 (2018). doi:10.1038/s41396-018-0049-5
- 516 7. Hurwitz, B. L. & U'Ren, J. M. Viral metabolic reprogramming in marine ecosystems. *Current*
517 *Opinion in Microbiology* **31**, 161–168 (2016).
- 518 8. Suttle, C. a. Marine viruses-major players in the global ecosystem. *Nat. Rev. Microbiol.* **5**, 801–
519 812 (2007).
- 520 9. Brum, J. R. *et al.* Patterns and ecological drivers of ocean viral communities. *Science* (80-.). **348**,
521 (2015).
- 522 10. Danovaro, R. *et al.* Virus-mediated archaeal hecatomb in the deep seafloor. *Sci. Adv.* **2**, (2016).
- 523 11. Pratama, A. A. & van Elsas, J. D. The 'Neglected' Soil Virome - Potential Role and Impact.
524 *Trends in Microbiology* (2018). doi:10.1016/j.tim.2017.12.004
- 525 12. Gómez, P. & Buckling, A. Bacteria-phage antagonistic coevolution in soil. *Science* (80-.). **332**,
526 106–109 (2011).
- 527 13. Reyes, A., Semenkovich, N. P., Whiteson, K., Rohwer, F. & Gordon, J. I. Going viral: Next-
528 generation sequencing applied to phage populations in the human gut. *Nature Reviews*
529 *Microbiology* **10**, 607–617 (2012).
- 530 14. Abeles, S. R. & Pride, D. T. Molecular bases and role of viruses in the human microbiome.
531 *Journal of Molecular Biology* **426**, 3892–3906 (2014).
- 532 15. Rohwer, F. & Edwards, R. The phage proteomic tree: A genome-based taxonomy for phage. *J.*
533 *Bacteriol.* **184**, 4529–4535 (2002).
- 534 16. Yarza, P. *et al.* Uniting the classification of cultured and uncultured bacteria and archaea using
535 16S rRNA gene sequences. *Nat. Rev. Microbiol.* **12**, 635–645 (2014).
- 536 17. Lawrence, J. G., Hatfull, G. F. & Hendrix, R. W. Imbrolios of viral taxonomy: Genetic exchange
537 and failings of phenetic approaches. *J. Bacteriol.* **184**, 4891–4905 (2002).

- 538 18. Sullivan, M. B. Viromes, Not Gene Markers, for Studying Double-Stranded DNA Virus
539 Communities. *J. Virol.* **89**, 2459–2461 (2015).
- 540 19. Deng, L. *et al.* Viral tagging reveals discrete populations in Synechococcus viral genome sequence
541 space. *Nature* **513**, 242–245 (2014).
- 542 20. Gregory, A. C. *et al.* Genomic differentiation among wild cyanophages despite widespread
543 horizontal gene transfer. *BMC Genomics* **17**, (2016).
- 544 21. Bobay, L. & Ochman, H. Biological species in the viral world. **115**, (2018).
- 545 22. Mavrich, T. N. & Hatfull, G. F. Bacteriophage evolution differs by host, lifestyle and genome.
546 *Nat. Microbiol.* **2**, (2017).
- 547 23. Ackermann, H.-W. Phage Classification and Characterization BT - Bacteriophages: Methods and
548 Protocols, Volume 1: Isolation, Characterization, and Interactions. in (eds. Clokie, M. R. J. &
549 Kropinski, A. M.) 127–140 (Humana Press, 2009). doi:10.1007/978-1-60327-164-6_13
- 550 24. Simmonds, P. *et al.* Consensus statement: Virus taxonomy in the age of metagenomics. *Nat. Rev.*
551 *Microbiol.* **15**, 161–168 (2017).
- 552 25. Paez-Espino, D. *et al.* Uncovering Earth’s virome. *Nature* **536**, 425–430 (2016).
- 553 26. Brister, J. R., Ako-Adjei, D., Bao, Y. & Blinkova, O. NCBI viral Genomes resource. *Nucleic*
554 *Acids Res.* **43**, D571–D577 (2015).
- 555 27. Nishimura, Y. *et al.* ViPTree: The viral proteomic tree server. *Bioinformatics* **33**, 2379–2380
556 (2017).
- 557 28. Lima-Mendez, G., Van Helden, J., Toussaint, A. & Leplae, R. Reticulate representation of
558 evolutionary and functional relationships between phage genomes. *Mol. Biol. Evol.* **25**, 762–777
559 (2008).
- 560 29. Bolduc, B. *et al.* vConTACT: an iVirus tool to classify double-stranded DNA viruses that infect
561 *Archaea* and *Bacteria*. *PeerJ* **5**, e3243 (2017).
- 562 30. Meier-Kolthoff, J. P. & Göker, M. VICTOR: genome-based phylogeny and classification of
563 prokaryotic viruses. *Bioinformatics* (2017). doi:10.1093/bioinformatics/btx440
- 564 31. Yu, C. *et al.* Real Time Classification of Viruses in 12 Dimensions. *PLoS One* **8**, (2013).
- 565 32. Gao, Y. & Luo, L. Genome-based phylogeny of dsDNA viruses by a novel alignment-free method.
566 *Gene* **492**, 309–314 (2012).
- 567 33. Iranzo, J., Koonin, E. V., Prangishvili, D. & Krupovic, M. Bipartite Network Analysis of the
568 Archaeal Virosphere: Evolutionary Connections between Viruses and Capsidless Mobile
569 Elements. *J. Virol.* **90**, 11043–11055 (2016).
- 570 34. Aiewsakun, P. & Simmonds, P. The genomic underpinnings of eukaryotic virus taxonomy:
571 creating a sequence-based framework for family-level virus classification. *Microbiome* **6**, 38
572 (2018).
- 573 35. Lavigne, R. *et al.* Classification of myoviridae bacteriophages using protein sequence similarity.
574 *BMC Microbiol.* **9**, (2009).
- 575 36. Lavigne, R., Seto, D., Mahadevan, P., Ackermann, H. W. & Kropinski, A. M. Unifying classical
576 and molecular taxonomic classification: analysis of the Podoviridae using BLASTP-based tools.

- 577 *Res. Microbiol.* **159**, 406–414 (2008).
- 578 37. Henz, S. R., Huson, D. H., Auch, A. F., Nieselt-Struwe, K. & Schuster, S. C. Whole-genome
579 prokaryotic phylogeny. *Bioinformatics* **21**, 2329–2335 (2005).
- 580 38. Iranzo, J., Krupovic, M. & Koonin, E. V. The double-stranded DNA virosphere as a modular
581 hierarchical network of gene sharing. *MBio* **7**, (2016).
- 582 39. Bolduc, B., Youens-Clark, K., Roux, S., Hurwitz, B. L. & Sullivan, M. B. IVirus: Facilitating new
583 insights in viral ecology with software and community data sets imbedded in a cyberinfrastructure.
584 *ISME J.* **11**, 7–14 (2017).
- 585 40. Roux, S. *et al.* Ecogenomics and potential biogeochemical impacts of globally abundant ocean
586 viruses. *Nature* **537**, 689–693 (2016).
- 587 41. Vik, D. R. *et al.* Putative archaeal viruses from the mesopelagic ocean. *PeerJ* **5**, e3428 (2017).
- 588 42. Roux, S. *et al.* Ecogenomics of virophages and their giant virus hosts assessed through time series
589 metagenomics. *Nat. Commun.* **8**, (2017).
- 590 43. Emerson, J. B. *et al.* Host-linked soil viral ecology along a permafrost thaw gradient. *Nat.*
591 *Microbiol.* (2018). doi:10.1038/s41564-018-0190-y
- 592 44. Martinez-Hernandez, F. *et al.* Single-virus genomics reveals hidden cosmopolitan and abundant
593 viruses. *Nat. Commun.* **8**, (2017).
- 594 45. de la Cruz Peña, M. J. *et al.* Deciphering the Human Virome with Single-Virus Genomics and
595 Metagenomics. *Viruses* **10**, 113 (2018).
- 596 46. Nepusz, T., Yu, H. & Paccanaro, A. Detecting overlapping protein complexes in protein-protein
597 interaction networks. *Nat. Methods* **9**, 471–472 (2012).
- 598 47. Hulo, C., Masson, P., Le Mercier, P. & Toussaint, A. A structured annotation frame for the
599 transposable phages: A new proposed family ‘Saltoviridae’ within the Caudovirales. *Virology* **477**,
600 155–163 (2015).
- 601 48. Doyle, E. L. *et al.* Genome Sequences of Four Cluster P Mycobacteriophages. *Genome Announc.*
602 **6**, e01101-17 (2018).
- 603 49. Pope, W. H. *et al.* Bacteriophages of *Gordonia* spp. Display a spectrum of diversity and genetic
604 relationships. *MBio* **8**, (2017).
- 605 50. Pope, W. H. *et al.* Whole genome comparison of a large collection of mycobacteriophages reveals
606 a continuum of phage genetic diversity. *Elife* **4**, e06416 (2015).
- 607 51. Sullivan, M. B. *et al.* The genome and structural proteome of an ocean siphovirus: A new window
608 into the cyanobacterial ‘mobilome’. *Environ. Microbiol.* **11**, 2935–2951 (2009).
- 609 52. Matsen, F. A., Kodner, R. B. & Armbrust, E. V. pplacer: linear time maximum-likelihood and
610 Bayesian phylogenetic placement of sequences onto a fixed reference tree. *BMC Bioinformatics*
611 **11**, 538 (2010).
- 612 53. Howard-Varona, C., Hargreaves, K. R., Abedon, S. T. & Sullivan, M. B. Lysogeny in nature:
613 Mechanisms, impact and ecology of temperate phages. *ISME Journal* **11**, 1511–1520 (2017).
- 614 54. Mirzaei, M. K. & Maurice, C. F. Ménage à trois in the human gut: Interactions between host,
615 bacteria and phages. *Nature Reviews Microbiology* **15**, 397–408 (2017).

- 616 55. O’Leary, N. A. *et al.* Reference sequence (RefSeq) database at NCBI: Current status, taxonomic
617 expansion, and functional annotation. *Nucleic Acids Res.* **44**, D733–D745 (2016).
- 618 56. Krupovic, M. *et al.* Taxonomy of prokaryotic viruses: update from the ICTV bacterial and
619 archaeal viruses subcommittee. *Arch. Virol.* **161**, 1095–1099 (2016).
- 620 57. Adams, M. J. *et al.* Changes to taxonomy and the International Code of Virus Classification and
621 Nomenclature ratified by the International Committee on Taxonomy of Viruses (2017). *Arch.*
622 *Virol.* **162**, 2505–2538 (2017).
- 623 58. Adriaenssens, E. M. *et al.* Taxonomy of prokaryotic viruses: 2017 update from the ICTV Bacterial
624 and Archaeal Viruses Subcommittee. *Archives of Virology* 1–5 (2018). doi:10.1007/s00705-018-
625 3723-z
- 626 59. Altschul, S. F. *et al.* Gapped BLAST and PSI-BLAST: A new generation of protein database
627 search programs. *Nucleic Acids Research* **25**, 3389–3402 (1997).
- 628 60. Wiwie, C., Baumbach, J. & Röttger, R. Comparing the performance of biomedical clustering
629 methods. *Nat. Methods* **12**, 1033–1038 (2015).
- 630 61. Brohée, S. & van Helden, J. Evaluation of clustering algorithms for protein-protein interaction
631 networks. *BMC Bioinformatics* **7**, (2006).
- 632 62. Kamburov, A., Stelzl, U. & Herwig, R. IntScore: A web tool for confidence scoring of biological
633 interactions. *Nucleic Acids Res.* **40**, (2012).
- 634 63. Goldberg, D. S. & Roth, F. P. Assessing experimentally derived interactions in a small world.
635 *Proc. Natl. Acad. Sci.* **100**, 4372–4376 (2003).
- 636 64. Buchfink, B., Xie, C. & Huson, D. H. Fast and sensitive protein alignment using DIAMOND. *Nat.*
637 *Methods* **12**, 59–60 (2015).
- 638 65. Oliphant, T. E. SciPy: Open source scientific tools for Python. *Comput. Sci. Eng.* **9**, 10–20 (2007).
- 639 66. McKinney, W. Data Structures for Statistical Computing in Python. *Proc. 9th Python Sci. Conf.*
640 **1697900**, 51–56 (2010).
- 641 67. Pedregosa, F. *et al.* Scikit-learn: Machine Learning in Python. *J. Mach. Learn. Res.* **12**, 2825–2830
642 (2011).
- 643 68. Csárdi, G. & Nepusz, T. The igraph software package for complex network research. *InterJournal*
644 *Complex Syst.* **1695**, 1–9 (2006).
- 645 69. Federico, P., Pfeffer, J., Aigner, W., Miksch, S. & Zenk, L. Visual Analysis of Dynamic Networks
646 Using Change Centrality. in *2012 IEEE/ACM International Conference on Advances in Social*
647 *Networks Analysis and Mining* 179–183 (2012). doi:10.1109/ASONAM.2012.39

648

649 **ACKNOWLEDGEMENTS.** We thank Laura Bollinger, Gareth Trubl, and Igor Tolstoy for their
650 comments on improving the manuscript, as well as Wesley Zhi-Qiang You for helping push the network
651 analytics. High performance computational support was provided as an award from the Ohio
652 Supercomputer Center to MBS. Funding was provided in part by the Department of Energy’s Genome
653 Sciences Program Soil Microbiome Scientific Focus Area award (#SCW1632) to Lawrence Livermore
654 National Laboratory; an NSF Biological Oceanography award (OCE#1536989), and a Gordon and Betty

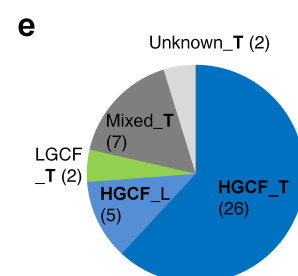
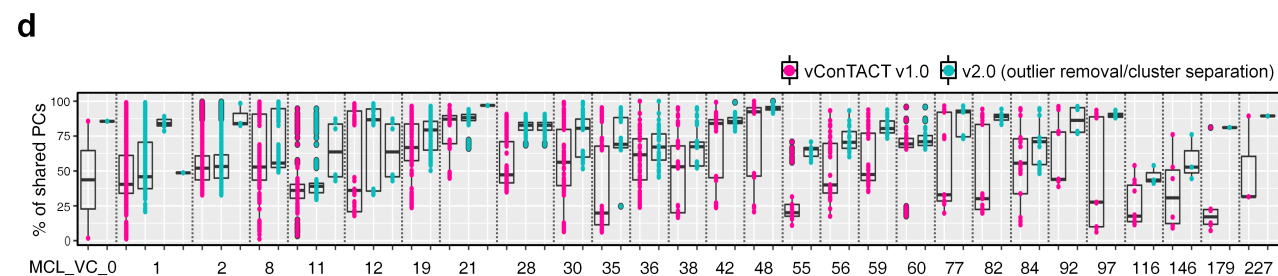
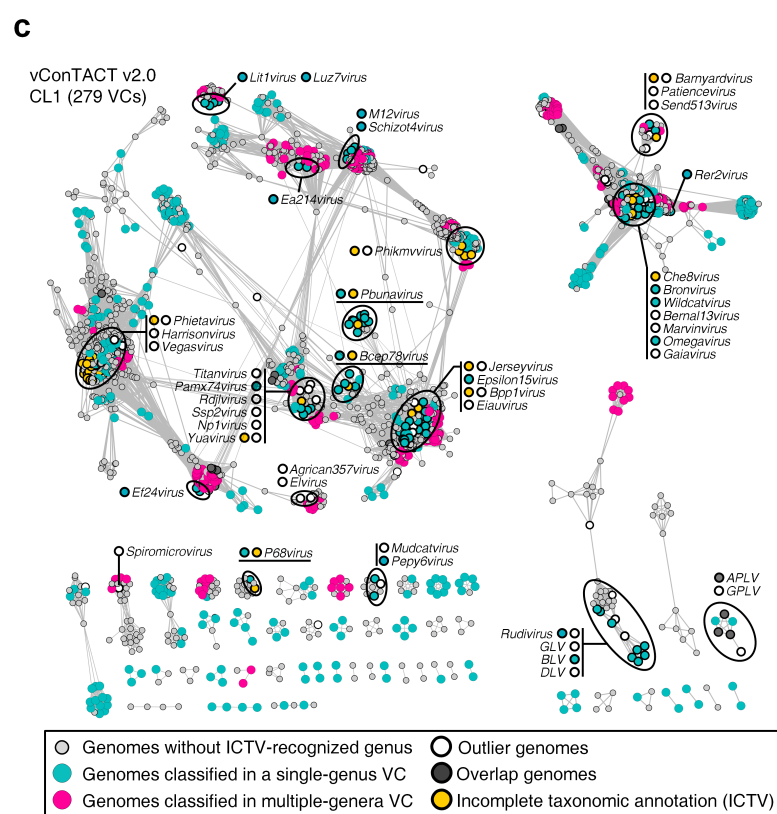
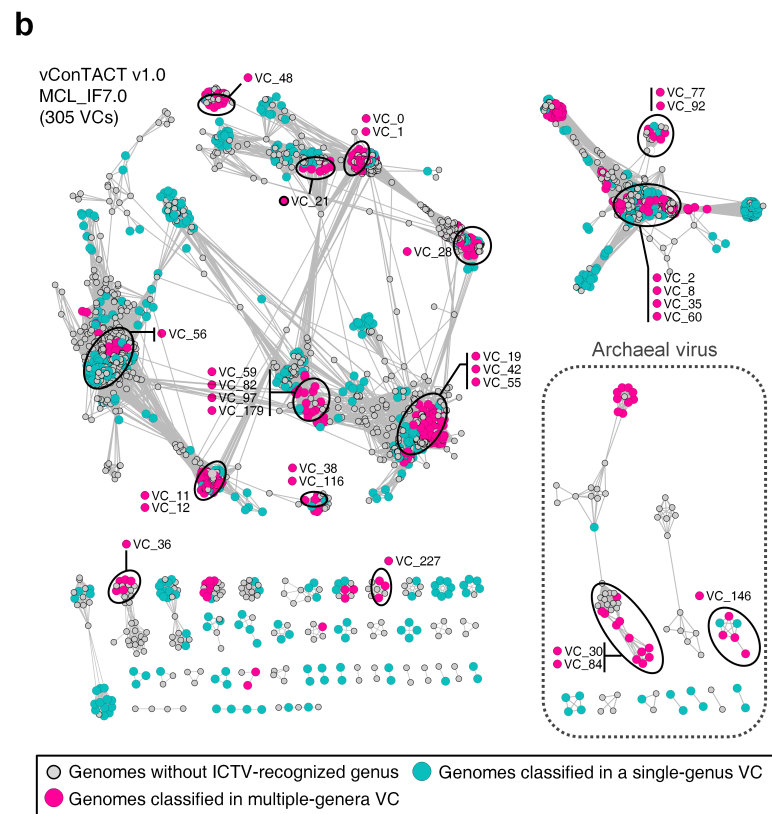
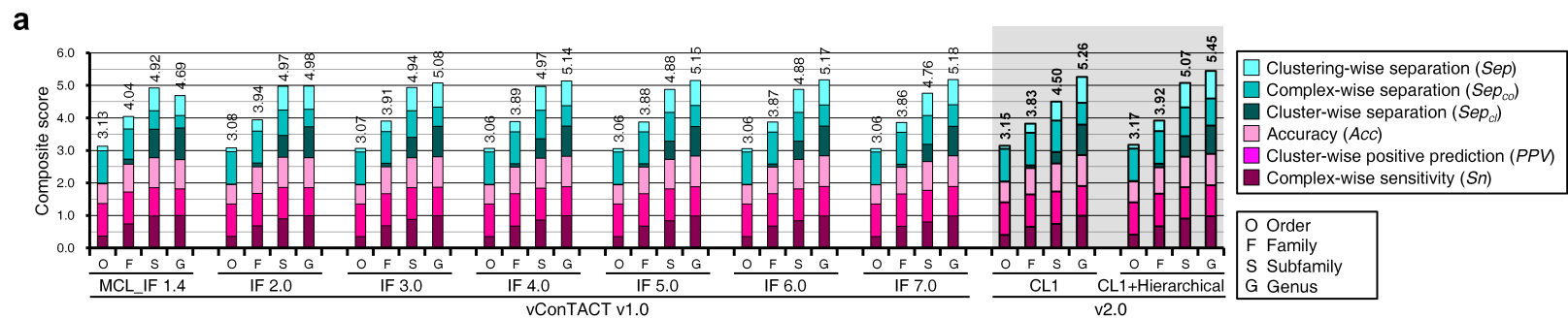
655 Moore Foundation Investigator Award (#3790) to MBS. Funding was provided to JRB by the Intramural
656 Research Program of the NIH, National Library of Medicine. The work conducted by the U.S.
657 Department of Energy Joint Genome Institute is supported by the Office of Science of the U.S.
658 Department of Energy under Contract DE-AC02-05CH11231 to SR. This work was funded in part
659 through Battelle Memorial Institute's prime contract with the US National Institute of Allergy and
660 Infectious Diseases (NIAID) under Contract No. HHSN272200700016I to JHK. The content of this
661 publication does not necessarily reflect the views or policies of the US Department of Health and Human
662 Services or of the institutions and companies affiliated with the authors.

663 **AUTHOR CONTRIBUTIONS.** HBJ, BB and MBS designed the study. OZ and MBS wrote the
664 manuscript with significant contributions from all co-authors. HBJ and BB performed the statistical and
665 network analyses.

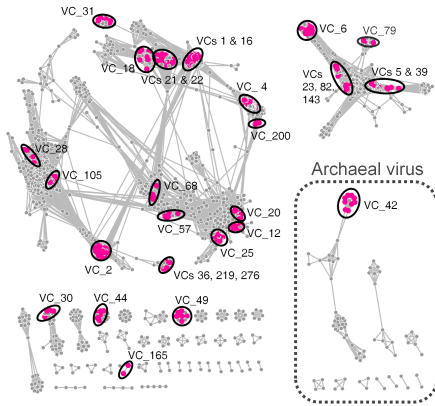
666 **COMPETING INTERESTS.** The authors declare no competing interests.

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

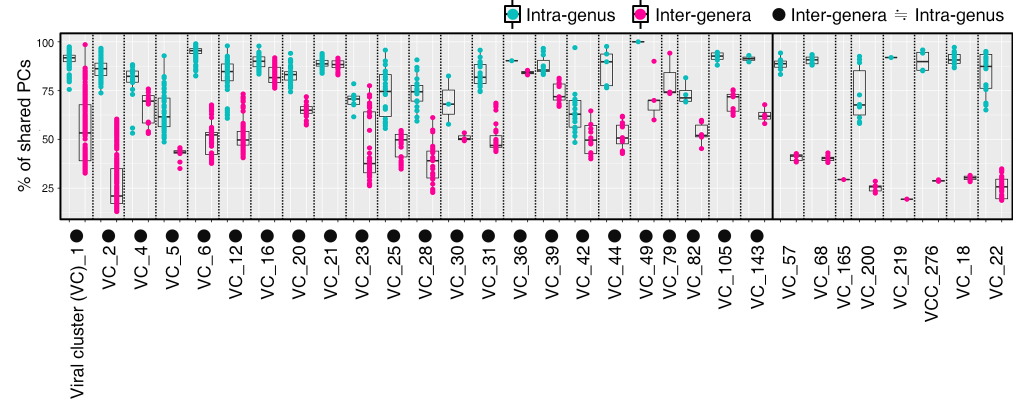
669



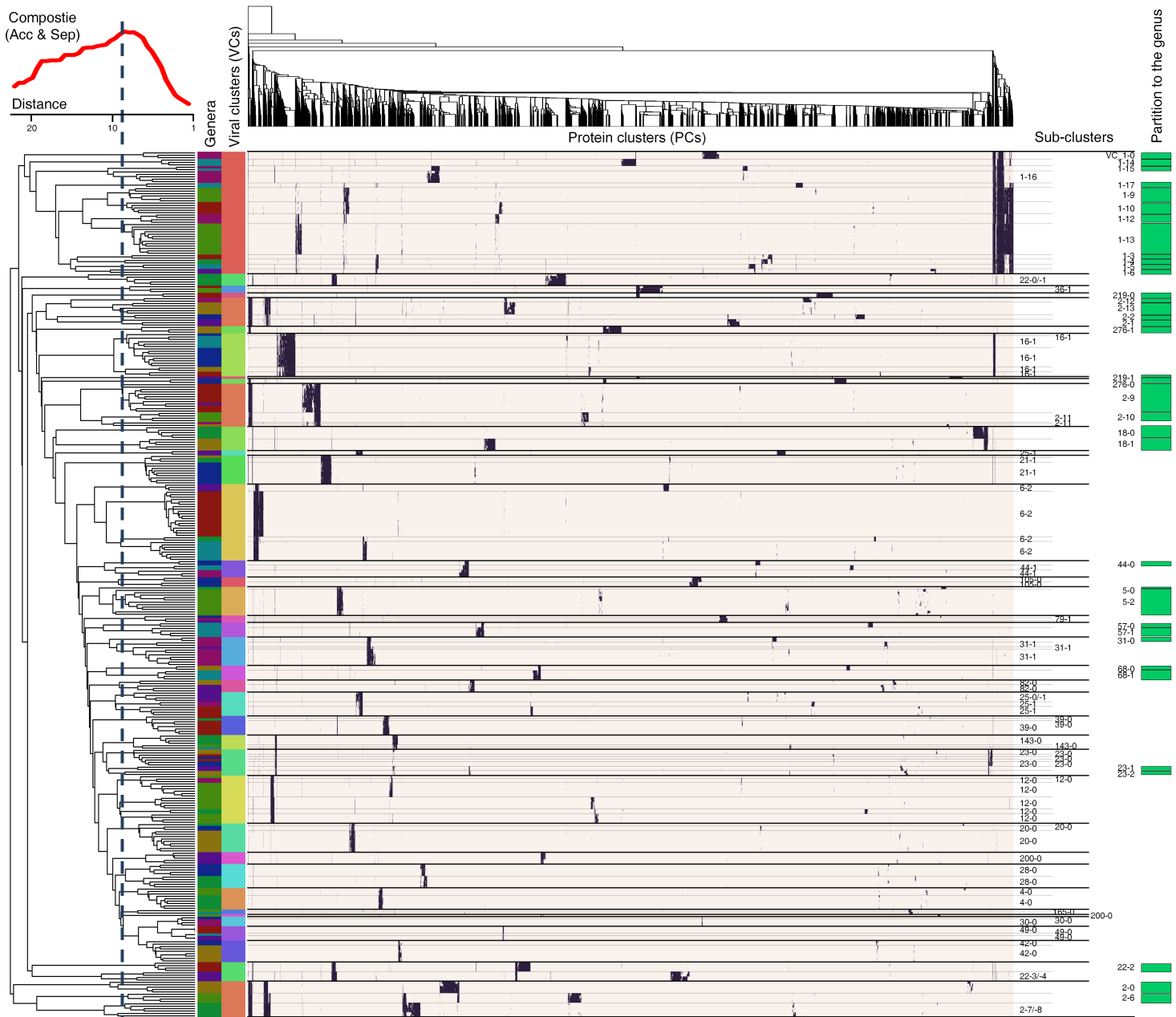
a



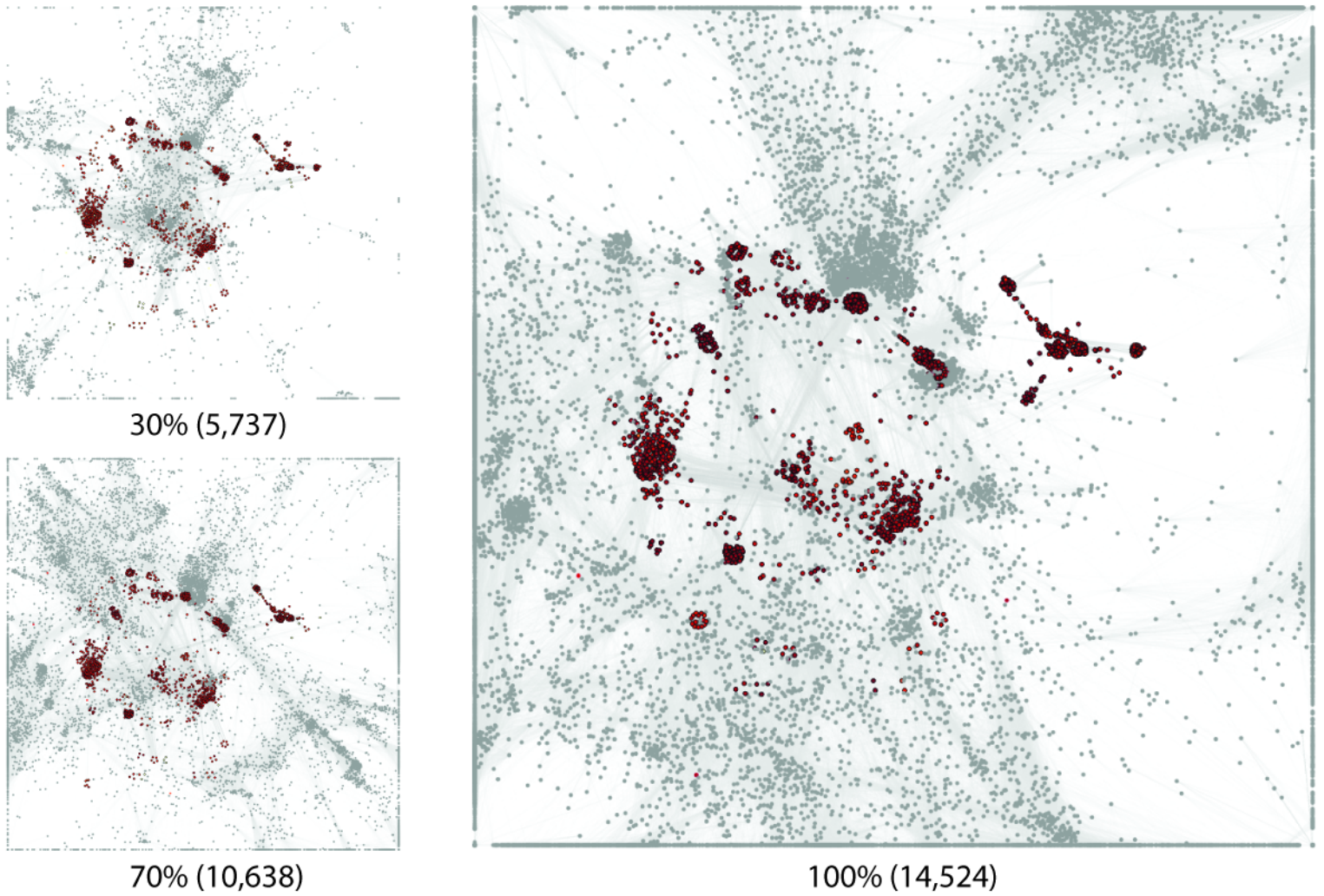
b



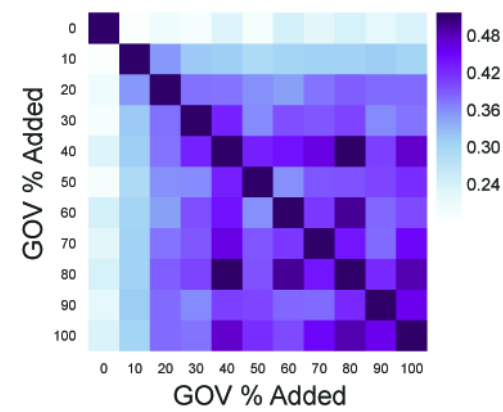
c



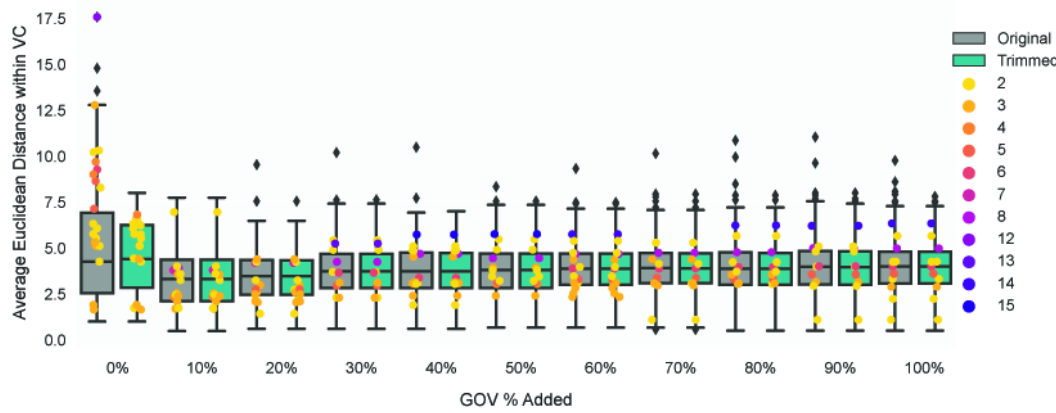
a



b



c



d

