

1 **Genetic diversity and characterization of circular replication(Rep)-encoding**
2 **single-stranded (CRESS) DNA viruses**

3 Perumal Arumugam Desingu^{1*}, K. Nagarajan²

4 ¹Department of Microbiology and Cell Biology, Indian Institute of Science, Bengaluru, India

5 ²Department of Veterinary Pathology, Madras Veterinary College, Vepery, Chennai, 600007, Tamil Nadu
6 Veterinary and Animal Sciences University (TANUVAS)

7

8

9

10

11

12

13 **Keywords:** CRESS DNA virus; Rep gene; Cap gene; genetic diversity; evolution;
14 classification

15

16

17

18

19 **Running title:** Genetic diversity and evolution of CRESS DNA viruses

20

21

22

23

24

25

26 * Corresponding authors,

27 Perumal Arumugam Desingu, M.V.Sc., Ph.D.,

28 DST-INSPIRE Faculty, Department of Microbiology and Cell Biology, Division of

29 Biological Sciences, Indian Institute of Science, Bangalore – 560012, India. Email:

30 perumald@iisc.ac.in; padesingu@gmail.com Phone: +91 80 2293 2068, Fax: +91 80 2360

31 2697

32 Abstract

33 The CRESS-DNA viruses are the ubiquitous virus detected in almost all eukaryotic life trees
34 and play an essential role in the maintaining ecosystem of the globe. Still, their genetic
35 diversity is not fully understood. Here we bring to light the genetic diversity of Replication
36 (Rep) and Capsid (Cap) proteins of CRESS-DNA viruses. We divided the Rep protein of the
37 CRESS-DNA virus into ten clusters using CLANS and phylogenetic analyzes. Also, most of
38 the Rep protein in Rep cluster 1 (R1) and R2 (*Circoviridae*, *Smacoviridae*, *Nanoviridae*, and
39 CRESSV1-5) contain the Viral_Rep superfamily and P-loop_NTPase superfamily domains,
40 while the Rep protein of viruses in other clusters has no such characterized functional
41 domain. The *Circoviridae*, *Nanoviridae*, and CRESSV1-3 viruses contain two domains, such
42 as Viral_Rep and P-loop_NTPase; the CRESSV4 and CRESSV5 viruses have only the
43 Viral_Rep domain, and most of the sequences in the pCRESS-related group have only P-
44 loop_NTPase, and *Smacoviridae* do not have these two domains. Further, we divided the Cap
45 protein of the CRESS-DNA virus into 20 clusters using CLANS and phylogenetic analyzes.
46 The Rep and Cap proteins of *Circoviridae* and *Smacoviridae* are grouped into a specific
47 cluster. Cap protein of CRESS-DNA viruses grouped with one cluster and Rep protein with
48 another cluster. Further, our study reveals that selection pressure plays a significant role in
49 the evolution of CRESS-DNA viruses' Rep and Cap genes rather than mutational pressure.
50 We hope this study will help determine the genetic diversity of CRESS-DNA viruses as more
51 sequences are discovered in the future.

52 Importance

53 The genetic diversity of CRESS-DNA viruses is not fully understood. CRESS-DNA viruses
54 are classified as CRESSV1 to CRESSV6 using only Rep protein. This study revealed that the
55 Rep protein of the CRESS-DNA viruses is classified as CRESSV1 to CRESSV6 groups and
56 the new *Smacoviridae*-related, CRESSV2-related pCRESS-related, *Circoviridae*-related, and
57 1 to 4 outgroups, according to the Viral_Rep and P-loop_NTPase domain organization,
58 CLANS, and phylogenetic analysis. Furthermore, for the first time in this study, the Cap
59 protein of CRESS-DNA viruses was classified into 20 distinct clusters by CLANS and
60 phylogenetic analysis. Through this classification, the genetic diversity of CRESS-DNA
61 viruses clarifies the possibility of recombinations in Cap and Rep proteins. Finally, it has
62 been shown that selection pressure plays a significant role in the evolution and genetic
63 diversity of Cap and Rep proteins. This study explains the genetic diversity of CRESS-DNA
64 viruses and hopes that it will help classify future detected viruses.

65 Introduction

66 Circular replication(Rep)-encoding single-stranded (CRESS)-DNA viruses are ubiquitous
67 viruses that are reported to spread worldwide and infect almost all eukaryotic tree of life¹⁻³.
68 CRESS-DNA viruses have also been found in environmental samples such as sewage,
69 seawater, lakes, and springs⁴⁻¹¹. Recently, ssDNA viruses have been classified into 13
70 families¹; ten families (*Anelloviridae*, *Bacilladnaviridae*, *Bidnaviridae*, *Circoviridae*,
71 *Geminiviridae*, *Genomoviridae*, *Nanoviridae*, *Parvoviridae*, *Redondoviridae*, and

72 *Smacoviridae*) are reported from the eukaryotes¹². These viruses are commonly found with
73 replication initiation protein (Rep) and structural capsid protein (Cap)^{1,12}. Of the ten ssDNA
74 virus families found in eukaryotes, the *Bidnaviridae* and *Parvoviridae* families have the
75 linear Genome topology, and the *Anelloviridae* family have a different Rep protein, with the
76 remaining seven families containing circular ssDNA with Rep protein containing the
77 preserved HUH endonuclease motif and superfamily 3 helicase (S3H) domain¹².

78 Recently, these characterized seven families of ssDNA viruses infect eukaryotes
79 (*Bacilladnaviridae*, *Circoviridae*, *Geminiviridae*, *Genomoviridae*, *Nanoviridae*,
80 *Redondoviridae*, and *Smacoviridae*), and uncharacterized CRESS-DNA viruses have been
81 classified into separate groups using this characteristic and conserved two-domain Rep
82 protein¹². Thus, unclassified CRESS-DNA viruses are classified as CRESSV1 through
83 CRESSV6¹². So far, the Rep protein of CRESS-DNA viruses has been characterized to
84 contain the HUH motif and S3H domain^{1,12}. It is also not widely known what other domains
85 are present in the rep protein of CRESS-DNA viruses that accumulate day by day through
86 metagenomic sequencing in different environmental samples and how they help classify
87 CRESS-DNA viruses. Furthermore, the classification of CRESS DNA viruses by capsid
88 proteins is challenging due to the lack of conserved portions of the capsid proteins of the
89 CRESS DNA viruses as found in the Rep protein¹². In particular, the capsid proteins of
90 CRESS DNA viruses are reported to be derived from a number of RNA viruses¹³⁻¹⁶. It is also
91 largely unknown which of the Cap proteins of the CRESS-DNA viruses that accumulate day
92 by day through metagenomic sequencing in different environmental samples are related to the
93 RNA viruses and the diversity in the Cap proteins of the CRESS-DNA viruses. A recent
94 study found that capsid proteins in Cruciviruses (CRESS DNA virus) are highly conserved
95 and possibly acquired from RNA viruses, but the Rep protein is more diversified than Cap
96 protein¹⁷. From these, it is speculated that Cruciviruses may have obtained Rep protein from
97 different CRESS-DNA viruses by recombination¹⁷. Therefore, it appears that the genetic
98 variation and recombination of CRESS-DNA viruses can be detected by dividing the capsid
99 proteins of almost identical CRESS-DNA viruses into groups. However, it should be noted
100 that there is no mechanism for classifying the capsid proteins of CRESS-DNA viruses so far.

101 The present study systematically classified the CRESS-DNA viruses Rep and Cap proteins
102 and reported the presence of different group-specific various domain organizations in the Rep
103 protein. Further, it explains the recombination-mediated evolution of the CRESS-DNA virus
104 and reveals that selection pressure plays a significant role in the evolution of CRESS-DNA
105 viruses' Rep and Cap genes rather than mutational pressure

106 **Results**

107 **CLANS based classification of CRESS-DNA Rep protein**

108 As a first step towards understanding the genetic diversity of the CRESS DNA viruses, we
109 analyzed the inter-relationship between the core viral proteins such as Rep and Cap proteins
110 of various isolates of CRESS-DNA viruses. We first chose the Rep protein for our analysis
111 since it shows a high degree of conservation among the CRESS-DNA viruses^{2,18-20}. To

112 explore the sequence diversity of the Rep protein of CRESS-DNA viruses, we collected 1160
113 (sequences details are provided in **Supplementary Data 1**) amino acid sequences of CRESS-
114 DNA viruses from the NCBI Database and grouped them based on pairwise sequence
115 similarity using the CLANS (CLuster ANALysis of Sequences) tool^{21,22}. The analysis grouped
116 the CRESS-DNA virus Rep protein sequences into ten different clusters (R1 to R10) [Rep
117 Cluster 1 (R1)] (a minimum of 10 viral sequences to a maximum of 487 sequences per group)
118 (**Figure 1A**, the individual sequence details in the different clusters are listed in
119 **Supplementary Data 2**). The majority of the clusters, except clusters R7 and R8, showed
120 inter-connections at a p -value threshold of $1e^{-2}$ (**Figure 1A**). Further, we also observed three
121 different superclusters (Super-cluster 1 - clusters R1, R2, R4, and R5; Super-cluster 2 -
122 clusters R3, R6, and R9; Super-cluster 3 - clusters R7 and R8) at a p -value threshold of $1e^{-5}$
123 (**Supplementary Figure 1A**). For a better understanding of the genetic diversity of the
124 CRESS-DNA virus, we classified the CRESS-DNA viruses into three broad groups as
125 follows (i) culturable CRESS-DNA viruses (*Circoviridae*, *Geminiviridae*, *Smacoviridae*,
126 *Cruciviridae*, etc.)²³ which are infective, (ii) replication-competent circular DNA (rccDNA)
127 which include the bovine meat and milk factors (BMMF) and Sphinx infective DNA
128 molecule²⁴ and (iii) uncharacterized and uncultivated CRESS-DNA¹ which were detected as
129 a DNA molecule in the viral metagenomic analysis. Interestingly, most of the sequences in
130 clusters R1 and R2 grouped with highly characterized and culturable viral families of
131 *Circoviridae*, *Smacoviridae*, and *Cruciviridae*. Further, cluster R8 sequences exclusively
132 belonged to BMMF of rccDNA, while all other clusters included uncultured CRESS-DNA
133 viruses. The remaining clusters (R3, R4, R5, R6, R7, R9, and R10) were classified as
134 uncharacterized and uncultivated CRESS-DNA.

135 **Domain organization in the CRESS-DNA virus Rep protein**

136 We were interested to find out if these different Rep protein clusters have any significant
137 differences in the organization of functional domains. The Rep genes of CRESS-DNA
138 viruses have been reported to contain two main functional domains/motif, HUH endonuclease
139 motif and superfamily 3 helicase domains^{1,25}. In this context, we analyzed the domain
140 organization of the Rep protein of viruses from different clusters using the Conserved
141 Domain search tool (<https://www.ncbi.nlm.nih.gov/Structure/cdd/wrpsb.cgi?>)²⁶⁻²⁹.
142 Interestingly, we found that only the Rep protein of viruses in clusters R1 and R2 displayed
143 functional domains such as Viral_Rep superfamily (Cdd:pfam02407) and P-loop_NTPase
144 superfamily (Cdd:pfam00910). On the other hand, cluster R8 (BMMF) sequences contained
145 Rep_1 superfamily (Cdd:pfam01446), a homologous domain to the Rep1 domain of bacteria
146 involved in plasmid replication. Moreover, we did not find any known putative functional
147 domains in our conserved domain analysis of other clusters consisting of uncultured viruses
148 (Rep cluster R3, R4, R5, R6, R7, R9, and R10). However, it should be noted that R4 and R5
149 in clusters of Rep protein that do not express these putative functional domains have
150 evolutionary links with R1 and R2 clusters that express functional domains (**Figure 1A**).
151 Similarly, cluster R7 has evolutionary links with cluster R8 that holds the Rep_1 domain
152 (**Figure 1A**).

153 We then analyzed the diversity of the Rep protein domains in depth belonging to clusters R1
154 and R2 to further classify the viruses in these clusters, which are highly related to culturable
155 viruses (sequences details are provided in **Supplementary Data 3**). To explore the different
156 domain organizations present in the viruses of clusters R1 and R2, we re-clustered them into
157 10 sub-clusters (cluster **a** to **j**) at a p -value threshold of $1e^{-38}$ (**Supplementary Figure 1B**). Of
158 these ten sub-clusters, we noted that sub-clusters such as **a**, **b**, **g**, and **h** formed a single group
159 (group 1), and **c**, **d**, **i**, and **j** sub-clusters formed a separate group (group 2) (**Supplementary**
160 **Figure 1B**). Furthermore, it can be seen that there are some evolutionary links between these
161 two groups (group 1 and group 2), but the sub-clusters **e** and **f** together as a separate group
162 (group 3) (**Supplementary Figure 1B**).

163 The viruses in the sub-cluster **a** majorly contain two main domains: Viral_Rep and P-
164 loop_NTPase domains. Some sequences had one of the following additional domains in
165 between Viral_Rep and P-loop_NTPase domain such as the AAA ATPase domain, Penta-EF
166 hand, DNA-binding ATP-dependent protease La, Type III secretion system protein PrgH-
167 EprH (PrgH), and Parvovirus non-structural protein NS1 (**Supplementary Data 4**).
168 Similarly, cluster **b** viruses also contained the Viral_Rep and P-loop_NTPase domains. In
169 addition, few sequences had a third functional domain between Viral_Rep and P-
170 loop_NTPase domain such as AAA+-type ATPase, SpoVK/Ycf46/Vps4 family, or Type VI
171 protein secretion system component VasK. Moreover, the cluster **b** viruses also contained a
172 combination of domains such as (i) incomplete Viral_Rep + P-loop_NTPase, and (ii)
173 Viral_Rep + incomplete P-loop_NTPase (**Supplementary Data 4**). Also, most sequences in
174 sub-cluster **g** contain Viral_Rep + P-loop_NTPase domains, and some sequences are
175 incomplete with these domains or possess only one of the two domains (**Supplementary**
176 **Data 4**). Significantly, most sequences in the sub-cluster **h** contain only the P-loop_NTPase
177 domains (**Supplementary Data 4**). More interestingly, it was revealed that most of the
178 sequences in sub-cluster **c** and **d** have only Viral_Rep domains (**Supplementary Data 4**).
179 Also, sub-cluster **i**, which is grouped with sub-cluster **c** and **d**, contains the Viral_Rep+P-
180 loop_NTPase domains, and sub-cluster **j** contains the Viral_Rep domain+incomplete P-
181 loop_NTPase domains (**Supplementary Data 4**). Finally, it is essential to note that the
182 sequences in sub-clusters **e** and **f** have no known putative functional domains
183 (**Supplementary Data 4**). Collectively, our analyses reveal a vast diversity of domains in the
184 viral Rep protein of CRESS-DNA viruses ranging from lack of any known functional
185 domains to the combination of multiple functional domains.

186 **Phylogenetic tree based classification of CRESS-DNA virus Rep protein and group-** 187 **specific domain organization**

188 Recently CRESS-DNA viruses have been classified into different groups CRESSV1 to
189 CRESSV6 using Rep protein^{1,12}. Therefore, we are interested in finding out which CRESSV
190 groups the clusters of Rep protein with varying organizations of domain identified in this
191 current study belong to. To find out, we performed a phylogenetic analysis of the sequences
192 of the Rep protein used in this present study with the sequences used to classify the CRESS-
193 DNA viruses in the previous study¹ (**Supplementary Data 5**). In this phylogenetic analysis,
194 we observed that CRESSV6, *P.pulchra*, pCRESS9, *Genomoviridae*, and *Geminiviridae* were

195 grouped together, and CRESSV4, CRESSV5, and *Nanoviridae* have formed another group
196 (**Figure 1B**), as in the previous study^{1,12}. As in the previous study¹, in plasmid CRESS
197 sequences (pCRESS), CRESS1, pCRESS2, and pCRESS3 formed a separate group, and
198 CRESS4, pCRESS5, CRESS6, pCRESS7, and pCRESS8 formed another group (**Figure 1B**).
199 Furthermore, CRESSV1 and CRESSV3 revealed a close association with *Circoviridae*
200 (**Figure 1B**). In addition, the group that showed a relationship with Smacoviridae was called
201 Smacoviridae-related; the group that showed contact with the CRESSV2 sequences was also
202 called CRESSV2-related; the group that showed a relationship with the pCRESS sequences
203 was called pCRESS-related; the group that showed contact with *Circoviridae* was called
204 *Circoviridae*-related; also the groups formed an outgroup were named as outgroup 1 to 4
205 (**Figure 1B**).

206 We first explored the sub-cluster **a** to **j** created by the clusters R1 and R2 with domain
207 organizations. Notably, we observed the sub-cluster **a** and **b** sequences that revealed the
208 domain organization Viral_Rep+P-loop_NTPase grouped into the CRESSV1, CRESSV2,
209 CRESSV3, *Circoviridae*, and *Circoviridae*-related groups (**Figure 1B; Supplementary Data**
210 **4**). Interestingly, sub-clusters **c** and **d**, which contain only Viral_Rep domains, are grouped
211 with CRESSV4 and CRESSV5, respectively (**Figure 1B; Supplementary Data 4**). We
212 observed that sub-clusters **e** and **f** grouped with *Smacoviridae* without any known putative
213 functional domains (**Figure 1B; Supplementary Data 4**). Significantly, sub-cluster **g**, which
214 display mostly Viral_Rep +P-loop_NTPase domains and some sequences with these domains
215 incomplete or with only one of the two domains, formed the CRESSV2-related group
216 (**Figure 1B; Supplementary Data 4**). Similarly, it should be noted that the sub-cluster **h**,
217 which contains most of the sequences only P-loop_NTPase domains, formed the pCRESS-
218 related group (**Figure 1B; Supplementary Data 4**). Also, sub-cluster **i** often have
219 Viral_Rep+P-loop_NTPase domains and sub-clusters **j** with Viral_Rep domain+incomplete
220 P-loop_NTPase domains grouped with *Nanoviridae* (**Figure 1B; Supplementary Data 4**).

221 Next, we explored clusters R3, R4, R5, R6, and R9 without any known putative functional
222 domains. Of these clusters, R4 and R5 combined with clusters R1 and R2 to form Super-
223 cluster 1 (**Supplementary Figure 1A**). Note that cluster R4 forms the Smacoviridae-related
224 group, and cluster R5 forms the outgroup 1 (**Figure 1B; Supplementary Data 4**). We
225 observed that the R3, R6, and R9 clusters formed Super-cluster 2 created outgroup 4,
226 outgroup 3, and outgroup 2, respectively (**Figure 1B; Supplementary Data 4**). These results
227 show that CRESS-DNA virus Rep proteins group into the phylogenetic tree, as is the case
228 with CLANS clustering and domain organizations.

229 **Classification of CRESS-DNA virus Cap protein using CLANS**

230 While the Rep protein of CRESS-DNA viruses is evolutionarily conserved, the Cap protein is
231 highly diverse^{2,18-20}. Therefore, previous studies analyzed the evolution of capsid proteins
232 primarily by structural fold comparisons rather than sequence comparisons^{23,30-32}. However,
233 we took advantage of the recent explosion in the metagenomic data from CRESS-DNA
234 viruses. We employed a sequence comparison method to classify and identify the genetic
235 diversity of CRESS-DNA virus Cap protein. We collected 1823 amino acid sequences of

236 CRESS-DNA viruses from the NCBI Database and grouped them based on pairwise
237 similarity (CLANS analysis) (sequences details are provided in **Supplementary Data 6**). The
238 analysis classified the CRESS-DNA virus Cap gene sequences into 20 different clusters
239 (minimum of ten sequences per group was considered to classify them as an individual
240 cluster) (the individual sequence details in the different clusters are listed in **Supplementary**
241 **Data 7**). Most of the clusters show interconnections with other clusters, except the clusters
242 C3 (Cap cluster 3), C4, C19, and C20, which were isolated from other clusters (orphan
243 clusters) in a pairwise similarity network (**Supplementary Figure 2**) (p -value threshold of
244 $1e^{-02}$). Cluster C1 of CRESS-DNA virus sequences clustered with *Circoviridae* viruses, while
245 cluster C2 showed a relationship with *Geminiviridae* viruses, cluster C3 sequences clustered
246 with *Smacoviridae* viruses, and cluster C6 sequences clustered with *Cruciviridae* virus
247 sequences (**Supplementary Data 7**). Among the 20 clusters identified for the Cap protein of
248 the CRESS-DNA viruses (**Figure 2A**), the clusters C2, C7, C8, C11, C12, C15, and C18
249 form a supercluster (**Figure 2A**) in the sequence similarity network analysis at a p -value
250 threshold of $>1e^{-04}$. Similarly, the supercluster consists of clusters C1, C14, and C17 in
251 CLANS analysis (**Figure 2A; Supplementary Data 7**). In addition, clusters C3, C4, C5, C6,
252 C9, C10, C13, C16, C19, and C20 were isolated from other clusters (orphan clusters) in a
253 pairwise similarity network (**Figure 2A; Supplementary Data 7**) (p -value threshold of $1e^{-$
254 04). Collectively, CRESS-DNA virus Cap proteins also split into separate groups in CLANS
255 analysis and are thought to support the classification of Cap proteins.

256 **Only *Cruciviridae* Cap proteins related to RNA viruses**

257 In previous studies, it has been reported that the cap protein of the CRESS-DNA virus is
258 related to the RNA virus¹³⁻¹⁶, so we were interested to find out which of these 20 clusters is
259 related to the RNA virus. To do this, we retrieved the RNA virus sequences associated with
260 the Cap protein of the CRESS-DNA virus from the NCBI Database and performed CLANS
261 analysis (sequences details are provided in **Supplementary Data 8**). This analysis noted that
262 RNA viruses revealed association only with *Cruciviridae* virus sequences belonging to
263 cluster C6 at a p -value threshold of $1e^{-02}$ (**Figure 2B; Supplementary Data 9**)

264 **Phylogenetic tree based classification of CRESS-DNA virus Cap protein**

265 We examined whether CLANS analysis-based clustering of CRESS-DNA virus cap protein
266 sequences also grouped into the phylogenetic tree. Because cap proteins do not have common
267 domains as seen in CRESS-DNA virus Rep proteins, and the sequence alignments are low
268 from most genetic variants, we performed separate phylogenetic analysis for (i) supercluster
269 C1, C14, C17; (ii) supercluster C2, C7, C8, C11, C12, C15, and C18; and (iii) orphan clusters
270 such as C3, C4, C5, C6, C9, C10, C13, C16, C19, and C20. To do this, we first performed
271 phylogenetic analysis using sequences from the C1, C14, and C17 clusters that formed the
272 Cap protein supercluster. These clusters C1, C14, and C17 are well aligned (**Supplementary**
273 **Data 10**) and split into separate groups for the phylogenetic tree (**Figure 3A**). Similarly, C2,
274 C7, C8, C11, C12, C15, and C18 clusters are well aligned (**Supplementary Data 11**) and
275 split into separate groups for the phylogenetic tree (**Figure 3B**). In particular, C8, C11, and
276 C12 formed an outgroup, and this outgroup group C8 was somewhat detached, and C11 and

277 C12 grouped slightly closer together into the phylogenetic tree (**Figure 3B**) as seen in the
278 CLANS analysis (**Figure 2A**). Similarly, the C7 and C15 clusters grouped in the
279 phylogenetic tree (**Figure 3B**), as seen in the CLANS analysis (**Figure 2A**), and the C18 and
280 some C2 sequences grouped together with this (C7 and C15) group (**Figure 3B**). Also,
281 although the cluster C2 sequences are majorly grouped together, it is noteworthy that some
282 sequences are grouped together with a group formed by C8, C11, and C12 and a group
283 created by C7, C15, and C18 (**Figure 3B**). We then performed phylogenetic analysis
284 separately for the orphan clusters C3, C4, C5, C6, C9, C10, C13, C16, C19, and C20 clusters.
285 Thus, the sequences in these clusters are well-aligned C3 (**Supplementary Data 12**), C4
286 (**Supplementary Data 13**), C5 (**Supplementary Data 14**), C6 (**Supplementary Data 15**),
287 C9 (**Supplementary Data 16**), C10 (**Supplementary Data 17**), C13, C16, C19 and C20
288 (**Supplementary Data 18**), to form the phylogenetic tree C3 (**Supplementary Figure 3A**),
289 C4 (**Supplementary Figure 3B**), C5 (**Supplementary Figure 4A**), C6 (**Supplementary**
290 **Figure 4B**), C9 (**Supplementary Figure 5A**), C10 (**Supplementary Figure 4B**), C13, C16,
291 C19, and C20 (**Supplementary Figure 5C**).

292 **Recombination mediated evolution of CRESS-DNA viruses**

293 Recently, it has been reported that the cap protein of Cruciviruses is very similar, but the Rep
294 protein may be derived from different sources with greater diversity¹⁷; we examined whether
295 the sequences in the cluster of these 20 Cap proteins received the Rep protein from the same
296 group or from different groups. To do this, we took the representative sequences in each Cap-
297 cluster and identified the phylogenetic tree groups that contain its Rep protein
298 (**Supplementary Data 19**). In this analysis, it appears that the sequences in the same cap-
299 cluster have different groups of rep proteins (**Supplementary Data 19**). From these, it can be
300 inferred that the CRESS-DNA virus has the potential to acquire genetic diversity through
301 recombination in the Cap and Rep genes.

302 **Role of host codon usage selection pressure on Rep gene evolution**

303 Since we observe homology at amino acid levels between the Rep gene of CRESS-DNA
304 viruses but not any significant identity at the nucleotide sequence level, we suspected this
305 might be due to this virus's host codon usage bias-based selection pressure. To explore this,
306 we first analyzed the base composition of 1115 nucleotide sequence of CRESS-DNA viruses'
307 Rep genes (the details of nucleotide sequences used in the analysis are presented in
308 **Supplementary data 20**) as AT to GC ratio can affect codon usage in microbes^{33,34}. Our
309 study revealed that the Rep gene of CRESS-DNA viruses contains A>T>G>C with
310 AT%>GC% (average GC content is 43.6±SD7.06) (**Figure 4A; Supplementary Data 21**).
311 We next analyzed the codon usage bias using the effective number of codon usage (ENc)
312 analysis. ENc values <35 indicate high codon bias, and values >50 show general random
313 codon usage^{35,36}. The Rep gene of CRESS-DNA viruses has ENc values ranging from 31 to
314 61, while most of the ENc values fall between 40 and 60 (average ENc 51.004±SD5.73)
315 (**Figure 4B; Supplementary data 21**), indicating weak to strong codon usage bias. In
316 addition, we calculated the relative synonymous codon usage (RSCU) value which is the ratio
317 between the observed to the expected value of synonymous codons for a given amino acid. A

318 RSCU value of one indicates that there is no bias for that codon. In contrast, RSCU values
319 >1.0 have positive codon usage bias (defined as abundant codons), and RSCU values <1.0
320 have negative codon usage bias (defined as less-abundant codons)^{36,37}. Our analysis revealed
321 that the RSCU values of 28 codons were >1 and 31 codons were <1 in the Rep gene of all the
322 CRESS-DNA viruses (**Figure 4C; Supplementary Data 21**), clearly indicating a codon
323 usage bias (both positive and negative).

324 Next, we performed ENc-GC3s plot analysis where the ENc values are plotted against the
325 GC3s values (GC content at the third position in the codon) to determine the significant
326 factors such as selection or mutation pressure affecting the codon usage bias³⁸. In this
327 analysis, genes whose codon bias is affected by mutations will lie on or around the expected
328 curve. In contrast, genes whose codon bias is affected by selection and other factors will lie
329 beneath the expected curve^{36,38}. Interestingly, we observed that most of the points fall below
330 the expected curve in the ENc-GC3s plot analysis (**Figure 4D; Supplementary Data 21**),
331 indicating the strong presence of selection pressure rather than mutation pressure. Similarly,
332 neutrality plot analysis where GC12 values (average of the GC content percentage at the first
333 and second position in the codon) are plotted against GC3 values to evaluate the degree of
334 influence of mutation pressure and natural selection on the codon usage patterns, displayed a
335 slope of 0.2899 ($Y=0.2899*X+30.61$, $r=0.662$; $p<0.0001$) (**Figure 4E; Supplementary**
336 **Data 21**), indicating that the mutation pressure and natural selection were 28.9% and 71.1%,
337 respectively. Moreover, we performed Parity rule 2 bias analysis, where the AT bias
338 $[A3/(A3+T3)]$ is plotted against GC-bias $[G3/(G3+C3)]$ to determine whether mutation
339 pressure and natural selection affect the codon usage bias³⁸. If $A=T$ and $G=C$, it indicates
340 no mutation pressure and natural selection, while any discrepancies indicate mutation
341 pressure and natural selection. Our analysis of CRESS-DNA Rep gene sequences shows
342 unequal A to T and G to C numbers, indicating the presence of mutation and selection
343 pressure (**Figure 4F, Supplementary Data 21**). Taken together, these results suggest that
344 CRESS-DNA has wide host-range adaptation, maintaining better codon usage pattern with
345 bacteria, and further selection pressure has played a significant role in the evolution of the
346 CRESS-DNA viruses Rep gene rather than mutational pressure.

347 **Role of host codon usage selection pressure on Cap gene evolution**

348 Similar to the Rep gene, our NCBI nucleotide BLAST analysis of the Cap gene also showed
349 limited homology between the Cap gene of CRESS-DNA viruses. Since we observed a strong
350 codon-bias-based evolution in the Rep gene of CRESS-DNA viruses (**Figure 4A-F**), we
351 tested whether the Cap gene of the CRESS-DNA viruses also shows codon-bias-based
352 evolution to explore whether the evolution of Cap gene was influenced by mutation pressure
353 or selection pressure, we retrieved 1134 nucleotide sequences of Cap genes of CRESS-DNA
354 viruses (**Supplementary Data 22**) from NCBI public database. Our analysis of the
355 nucleotide base composition of the Cap gene revealed that the Cap gene contains $AT\%>GC\%$
356 (average GC content is $44.72\pm SD 5.96$) (**Figure 5A; Supplementary Data 23**). Further, the
357 Cap protein of the CRESS-DNA virus has ENc value ranging from 33 to 61, while most of
358 the sequence ENc values fall between 40 to 60 (average ENc $51.54\pm SD 5.36$) (**Figure 5B;**
359 **Supplementary Data 23**). Similarly, the RSCU values of 27 codons were >1 , and 31 codons

360 were <1 in all the Cap genes of CRESS-DNA viruses (**Figure 5C; Supplementary Data 23**).
361 Also, nine codons showed RSCU values <0.7, and 6 codons showed RSCU values>1.5,
362 indicating the presence of under-represented and over-represented codon bias in the Cap
363 gene, respectively (**Supplementary Data 23**). Moreover, we performed ENc-GC3s plot
364 analysis and found that most points fall below the expected curve in the ENC-GC3s plot
365 (**Figure 5D; Supplementary Data 23**). In line with this, the neutrality plot displayed a slope
366 of 0.1404 ($Y=0.1404*X+40.64$; $r= 0.512$; $p<0.0001$) (**Figure 5E; Supplementary Data 23**),
367 indicating 14% of mutation pressure and 86% of selection pressure in this gene and Parity
368 rule 2 bias analysis showed discrepancies in the A to T and G to C numbers in the third
369 position of the codon (**Figure 5F; Supplementary Data 23**). Taken together, these results
370 indicate that the selection pressure played a more significant role in the Cap gene than the
371 Rep genes of CRESS-DNA viruses.

372 Discussion

373 The genetic diversity of CRESS-DNA viruses so far is known only to be the tip of the
374 iceberg. Many novel CRESS-DNA viruses have recently been detected by metagenomic
375 sequencing^{8,39-41}. The rapid development of metagenomic sequencing suggests that in the
376 future, most CRESS-DNA viruses will be detected from different sources and that these
377 CRESS-DNA viruses will be divided into different virus families. Therefore, it is hoped that
378 identifying and classifying genetic diversity in CRESS-DNA viruses will help determine their
379 importance in transmission and pathogenesis and design antivirals and vaccines for
380 appropriate control and prevention. However, the classification of CRESS-DNA viruses has
381 been determined using only the Rep protein^{1,12}. This is because Rep protein contains
382 conserved HUH motif, and S3H domains, while Cap protein is unclassified because it has
383 high genetic diversity without being conserved^{1,12}. However, of the cruciviruses that classify
384 Cap protein well, the report that Cap proteins are nearly identical and that the highly diverse
385 Rep protein may be derived from different CRESS-DNA virus sources is critical here¹⁷.
386 Therefore, it can be expected that the genetic diversity and genetic recombination events of
387 CRESS-DNA viruses can be determined by detecting and classifying the diversity in both
388 Rep and Cap proteins.

389 It is noteworthy that recently, unclassified CRESS-DNA viruses using the Rep protein of
390 CRESS-DNA viruses were grouped into six groups called CRESSV1 to CRESSV6^{1,12}. The
391 present study reveals that there are not only CRESSV1 to CRESSV6 groups but also groups
392 with Smacoviridae-related, CRESSV2-related, pCRESS-related, Circoviridae-related, and 1
393 to 4 outgroups are there. So far, it has been reported that the Rep protein of the CRESS-DNA
394 virus contains the HUH motif and S3H domains^{1,12}. In this study, we report the presence of
395 domains such as Viral_Rep superfamily (Cdd: pfam02407) and P-loop_NTPase superfamily
396 (Cdd: pfam00910) in the Rep protein of most CRESS-DNA viruses. Furthermore, this present
397 study revealed the presence of these two domains in the CRESSV1, CRESSV2, CRESSV3,
398 *Circoviridae*, and *Circoviridae*-related groups and the *Nanoviridae* group. However, CLANS
399 and phylogenetic analyses clarify the viral_Rep and P-loop_NTPase domains in the
400 CRESSV1, CRESSV2, CRESSV3, *Circoviridae*, and *Circoviridae*-related groups are very
401 close and distinct from the *Nanoviridae* group. It is noteworthy that CRESSV1, CRESSV2,

402 CRESSV3, Circoviridae, and Circoviridae-related groups together formed the Rep sub-
403 cluster **a** and **b** and the sequences in the *Nanoviridae* group Rep sub-cluster **i** and **j**
404 (**Supplementary Figure 1B**). Our phylogenetic tree (**Figure 1B**) and previous study ¹²
405 reflect this diversity. In particular, some sequences in the rep sub-cluster **a** and **b** appear to
406 have an additional domain (Penta-EF hand, DNA-binding ATP-dependent protease La, Type
407 III secretion system protein PrgH-EprH (PrgH), etc.) between the Viral_Rep and P-
408 loop_NTPase domains. From the acquisition of such additional functional domains, it is clear
409 that these viruses are stepping into the next stage of evolution, and when more sequences are
410 found later, it is possible to speculate that they are likely to be classified as separate virus
411 families. However, since these sequences are detected by metagenomic sequencing from
412 uncultured viruses and maybe sequence alignment error, it may be imperative to isolate the
413 viruses and identify the significance of these additional functional domains.

414 Furthermore, in CLANS analysis, the sub-cluster **c** and **d** were grouped with the sub-cluster **i**
415 and **j** reacting with *Nanoviridae* (**Supplementary Figure 2B**), of which the sub-cluster **c** was
416 CRESSV4, and the sub-cluster **d** were CRESSV5 viruses (**Supplementary Data 4**); and
417 reflect in our phylogenetic tree (**Figure 1B**) and previous study ¹². In particular, the
418 CRESSV4 and CRESSV5 viruses have only the Viral_Rep domain, and the sequences in the
419 sub-cluster **j** related to *Nanoviridae* are the Viral_Rep+incomplete P-loop_NTPase, and the
420 sequences in the sub-cluster **i** are the Viral_Rep+P-loop_NTPase domains. Of these, it can
421 speculate that the CRESSV4 and CRESSV5 viruses, which have only the Viral_Rep domain
422 only, may have appeared first, followed by the sub-cluster **j** with the viral_Rep+incomplete
423 P-loop_NTPase, and finally the sub-cluster **i** virus with the Viral_Rep+P-loop_NTPase
424 domains. Similarly, sub-cluster **g** (CRESSV2-related group) that are often Viral_Rep+P-
425 loop_NTPase domains and some sequences where these domains are incomplete or show
426 only one of the two domains may have led to the emergence of CRESSV2 viruses with
427 Viral_Rep+P-loop_NTPase domains. Furthermore, it is essential to note that sub-cluster **h**,
428 which usually contains only the P-loop_NTPase domain, formed the pCRESS-related group.
429 Interestingly, no functional domains were found in Smacoviridae's Rep protein in the
430 Conserved Domain search tool, which revealed links between *Circoviridae* and *Nanoviridae*
431 in CLANS and biogenetic analyzes. Similarly, no functional domains were found in the
432 sequences in group R5 (CLANS) or Smacoviridae-related group (phylogenetic tree). It is
433 noteworthy that the sequences of Rep protein that formed the outgroups in this phylogenetic
434 analysis are the clusters of CLANS analysis, R3, R5, R6, and R9, forming separate groups.
435 R3, R5, R6, and R9 clusters formed Super-cluster 2 created outgroup 4, outgroup 1, outgroup
436 3, and outgroup 2, respectively. Remarkably, no functional domains are found in the
437 sequences in the R3, R5, R6, and R9 clusters that make up the outgroups. However, the
438 sequences that make up the outgroups are detected from uncultured viruses by metagenomic
439 sequencing, and it can be expected that the functional significance will be revealed by
440 isolating these viruses and characterizing the Rep protein.

441 Cap protein was high in genetic diversity, making it challenging to align and phylogenetically
442 classify correctly. Therefore, in this study, we subdivided the closest sequences into clusters
443 using CLANS analysis and then did phylogenetic classification by aligning them well using

444 the corresponding clusters. First, in this study, the Rep proteins were divided into clusters in
445 the CLANS analysis and then phylogenetic classification using the related clusters, which is
446 consistent with phylogenetic classification in the previous studies^{1,12}. Accordingly, we divide
447 the cap protein into 20 clusters using CLANS analysis, and (i) supercluster C1, C14, C17; (ii)
448 supercluster C2, C7, C8, C11, C12, C15, and C18; and (iii) orphan clusters such as C3, C4,
449 C5, C6, C9, C10, C13, C16, C19, and C20 became well aligned and led to phylogenetic
450 classification. Furthermore, only the Cruciviridae virus sequences in Cap-cluster C6 revealed
451 evolutionary relationships with RNA viruses, but future studies need to determine the
452 evolutionary origins of the sequences in other Cap-clusters. Remarkably, this study revealed
453 that viruses in the same Cap-cluster derive their Rep protein from groups of different Rep
454 proteins, which can be speculated to be generated by genetic recombination. These can be
455 believed to underscore the importance of classifying Cap protein. Finally, this study makes it
456 clear that selection pressure plays a more significant role than mutational pressure in the
457 genetic diversity and evolution of CRESS-DNA virus Cap and Rep protein. Therefore, it can
458 be expected that there will be more opportunities to detect CRESS-DNA viruses with greater
459 genetic diversity and/or recombination in the future. We hope this study will help determine
460 the genetic diversity/recombination of CRESS-DNA viruses as more sequences are
461 discovered in the future.

462 In conclusion, to the best of our knowledge, this is the first report on the CRESS-DNA virus
463 Rep protein classification using a different domain organization pattern; and CLANS and
464 phylogenetic analysis based on the classification of Cap protein. Furthermore, this study also
465 clarifies the genetic diversity in CRESS-DNA viruses formed by recombination and selection
466 pressures in Cap and Rep proteins. It is widely expected that CRESS-DNA viruses, which
467 have tremendous genetic diversity in the future, will be able to be detected from different
468 sources in different parts of the world through rapidly growing metagenomic sequences. We
469 hope this study will help you determine and accurately classify using CLANS, phylogenetic
470 groups, the domain organization pattern, genetic diversity, and recombination of those
471 CRESS-DNA viruses.

472 **Materials and Methods**

473 **I. Databases search, collection, and curation**

474 Complete genome sequences of CRESS-DNA viruses were retrieved from the NCBI
475 nucleotide database (<https://www.ncbi.nlm.nih.gov/nucleotide/>). Rep and Cap genes'
476 characterized protein-coding sequence (CDS) region and their corresponding amino acid
477 sequences were retrieved from the database in the available complete genome sequence of
478 CRESS-DNA viruses. The uncharacterized CDS of CRESS-DNA viruses were classified as a
479 Cap/Rep protein using NCBI protein BLAST
480 (<https://blast.ncbi.nlm.nih.gov/Blast.cgi?PAGE=Protein>) analysis, and the sequences were
481 retrieved. Further, complete genome sequences which contain Cap/Rep of every CRESS-
482 DNA were individually used to perform separate NCBI protein BLAST analysis
483 (<https://blast.ncbi.nlm.nih.gov/Blast.cgi?PAGE=Protein>). Their BLAST aligned sequences of

484 other ssDNA viruses (example: Circoviridae, Smacoviridae, Cruciviridae, etc.) were
485 retrieved.

486 **II. CLANS (CLuster ANalysis of Sequences) analysis**

487 The CLANS analysis was performed in the online Toolkit software
488 (<https://toolkit.tuebingen.mpg.de/tools/clans>). The protein sequences retrieved from the NCBI
489 database were subjected to the pairwise sequence similarity calculation using the online
490 CLANS analysis in the Toolkit²¹ with a scoring matrix of BLOSUM45 and BLAST HSP's
491 (High Scoring Pair) up to an E-value of $1e^{-2}$. Next, the CLANS files obtained from the
492 Toolkit were visualized in a Java application (clans.jar)²². A minimum of 1,00,000 rounds
493 was used to show the sequences connection and clusters in the clans.jar application. The
494 clusters were classified based on the Network method using offset values and global average
495 with maximum rounds of 10000 in clans.jar analysis.

496 **III. Analysis of functional domain organization in the protein**

497 We determined the domain organizations in the Rep protein of the CRESS-DNA virus using
498 the Conserved Domain search tool (<https://www.ncbi.nlm.nih.gov/Structure/cdd/wrpsb.cgi>).
499 For this, we used the CDD v3.19-58235 PSSMs database, Expect Value threshold line 0.01,
500 Composition-based statistics adjustment applied, and Performed by the maximum number of
501 hits to 500 in the Conserved Domain Search tool²⁶⁻²⁹.

502 **IV. Phylogenetic analyses**

503 The phylogenetic analysis was performed in PhyML 3.3_1 using the amino acid sequences of
504 the Rep/Cap protein of the CRESS-DNA virus clustered into clusters in the CLANS analysis,
505 which was retrieved from the NCBI public database. The phylogenetic analysis is in PhyML
506 3.3_1, Evolutionary model LG, Equilibrium frequencies Empirical ML- Model, discrete
507 gamma model [number of categories (n=4)], tree topology search with SPR (Subtree Pruning
508 and Regraphing), tree topology, branch length, and model parameters are optimizing
509 parameters, and SH-like statistics are used to test the branch support⁴²⁻⁴⁴. Further, the
510 phylogenetic trees were visualized through the interactive tree of life (iTOL) v5⁴⁵.

511 **V. Codon usage bias analysis**

512 **a) Nucleotide sequence composition analysis**

513 The nucleotide composition of CDSs, specifically the A%, T%, G%, and C% composition of
514 the Rep/Cap genes of CRESS-DNA viruses, were analyzed using Automated Codon Usage
515 Analysis (ACUA) Software⁴⁶.

516 **b) Relative Synonymous Codon Usage (RSCU) Analysis**

517 RSCU value is the ratio between the observed to the expected value of synonymous codons
518 for a given amino acid. When the RSCU value is one, it indicates that there is no bias for that
519 codon^{36,37}. This study determined the RSCU values using the ACUA Software⁴⁶. The

520 nucleotide sequences of Rep/Cap genes of CRESS-DNA viruses obtained from the NCBI
521 nucleotide public database were used for this analysis.

522 **c) Effective Number of Codons (ENc)**

523 The effective number of codon usage from 61 codons for the 20 amino acids is one method
524 that determines the codon usage bias and may range from 20 to 61. ENc values <35 indicate
525 high codon bias, and values >50 show general random codon usage^{35,36}. In this study, the
526 ENc values were determined on the online server (<http://ppuigbo.me/programs/CAIcal/>)⁴⁷,
527 and the input nucleotide sequences used in the CAI calculation were used in this analysis.

528 **VI. Determining the selection and mutation pressure**

529 **a) ENc-GC3s plot**

530 In this analysis, the ENc values are plotted against the third position of GC3s of codon values
531 to determine the significant factors such as selection or mutation pressure affecting the codon
532 usage bias³⁸. The expected curve was determined by estimating the expected ENc values for
533 each GC3s as recommended in previous publications^{36,38}. The ENc and GC3s for every gene
534 were obtained from an online CAI analysis server (<http://ppuigbo.me/programs/CAIcal/>)⁴⁷.
535 The genes would lie on or around the expected curve when mutation pressure only affects
536 codon bias. In contrast, they would fall considerably below the expected curve if codon bias
537 is influenced by selection and other factors^{36,38}.

538 **b) Neutrality plot analysis**

539 In a neutrality plot, GC12 values of the codon are plotted against GC3 values to evaluate the
540 degree of influence of mutation pressure and natural selection on the codon usage patterns.
541 The GC12 and GC3 values for the nucleotide sequences of Rep/Cap genes of CRESS-DNA
542 viruses were obtained from an online CAI analysis server
543 (<http://ppuigbo.me/programs/CAIcal/>)⁴⁷.

544 **c) Parity Rule 2 (PR2)-bias plot**

545 The PR2-bias, the AT bias [$A3/(A3+ T3)$] is plotted against GC-bias [$G3/(G3 + C3)$] to
546 mutation pressure and natural selection affecting the codon usage bias³⁸. The A3, T3, G3, and
547 C3 values of nucleotide sequences of Rep/Cap genes of CRESS-DNA viruses were obtained
548 using the ACUA Software⁴⁶.

549 **Acknowledgments**

550 PAD is a DST-INSPIRE faculty is supported by research funding from the Department of
551 Science and Technology (DST/INSPIRE/04/2016/001067), Government of India, and Core
552 grant from the Science and Engineering Research Board (SERB) (CRG/2018/002192),
553 Department of Science and Technology (DST), Government of India.

554

555

556 **Data Availability Statement**

557 We have retrieved the nucleotide sequences from publically available NCBI databases.
558 Further, all the nucleotide sequences accession numbers and names are indicated in the
559 respective figures and supplementary data.

560 **Conflict of interest**

561 There is no potential conflict of interest.

562

563 **FIGURE LEGENDS**

564 **Figure 1: CLANS analysis-based classification of CRESS-DNA virus Rep protein.** (A)
565 Representative CRESS-DNA virus Rep protein sequences were clustered using CLANS
566 Toolkit by their pairwise sequence similarity network. A total of 1160 amino acid sequences
567 of Rep protein (**Supplementary Data 1**) of CRESS-DNA viruses were used in this analysis
568 Classification of clusters was carried out by a Network-based method using offset values and
569 global average with maximum rounds 10000 in CLANS Toolkit analysis. The P -value $\leq 1e^{-02}$
570 was used to show the lines connecting the sequences. (B) Phylogenetic relationship of Rep
571 protein of CRESS-DNA viruses. The maximum-likelihood method inferred the evolutionary
572 history using the Subtree-Pruning-Regrafting algorithm in PhyML 3.3_1. A total of 1509
573 amino acid sequences of Rep protein (**Supplementary Data 5**) of CRESS-DNA viruses were
574 used in this analysis.

575 **Figure 2: Sequence similarities (CLANS) analysis-based CRESS-DNA virus capsid**
576 **protein clustering.** (A) A total of 1823 amino acid sequences of Cap protein of CRESS-
577 DNA viruses (**Supplementary Data 6**) were used and classified by their pairwise sequence
578 similarity network using CLANS. The clusters were classified using the Network-based
579 method using offset values and global average with a maximum of 10000 in CLANS Toolkit
580 analysis. The P -value $\leq 1e^{-05}$ was used to show the lines connecting the sequences. (B)
581 Pairwise sequence similarity based on CRESS-DNA virus capsid protein and +RNA viruses
582 relationship. Representative CRESS-DNA virus capsid protein sequences and their
583 relationship with RNA viruses using CLANS Toolkit. A total of 1967 amino acid sequences
584 of Cap protein of CRESS-DNA viruses and +RNA viruses were used in this analysis
585 (**Supplementary Data 8**). The clusters were classified using the Network-based method
586 using offset values and global average with a maximum of 10000 in CLANS Toolkit
587 analysis. The P -value $\leq 1e^{-02}$ was used to show the lines connecting the sequences.

588 **Figure 3: Phylogenetic relationship of CRESS-DNA virus Cap protein superclusters.**
589 (A) Phylogenetic tree depicting the genetic relationship between the CRESS-DNA virus Cap
590 protein supercluster formed by the clusters C1, C14, and C17. The details of sequences in
591 each cluster (**Supplementary Data 7**) and alignment are provided in **Supplementary Data**
592 **10.** (B) The phylogenetic tree represents the genetic relationship between the CRESS-DNA
593 virus Cap protein supercluster created by the clusters C2, C7, C8, C11, C12, C15, and C18.
594 The details of sequences in each cluster (**Supplementary Data 7**) and alignment are provided
595 in **Supplementary Data 11.** The maximum-likelihood method inferred the evolutionary

596 history using the Subtree-Pruning-Regrafting algorithm and bootstrap values in PhyML
597 3.3_1.

598 **Figure 4: Host codon usage selection pressure on Rep gene of CRESS-DNA virus**
599 **evolution.** (A) Representing the A, T, G, and C fraction; (B) Represent the ENc values; (C)
600 represent the codon usage fraction and RSCU values; (D) Represents ENc plotted against
601 GC3s; (E) Neutrality plot analysis of the GC12 and that of the GC3; and (F) Parity Rule 2
602 (PR2)-bias plot (Total of 1115 nucleotide sequences of Rep gene of CRESS-DNA viruses
603 were used in this analysis).

604 **Figure 5: Host codon usage selection pressure on Cap gene of CRESS-DNA virus**
605 **evolution.** (A) A, T, G, and C fraction; (B) ENc values; (C) Codon usage fraction and
606 RSCU values; (D) ENc plotted against GC3s; (E) Neutrality plot analysis of the GC12 and
607 that of the GC3; and (F) Parity Rule 2 (PR2)-bias plot (Total of 1134 nucleotide sequences
608 of Cap gene of CRESS-DNA viruses used in this analysis).

609 **Supplementary Figure 1: Pairwise amino acid sequence similarity network-based**
610 **CRESS-DNA virus classification of Rep protein of CRESS-DNA viruses.** (A) A total of
611 1160 amino acid sequences of Rep protein of CRESS-DNA viruses (**Supplementary Data 1**)
612 were clustered by CLANS. The clusters were classified using the Network-based method
613 using offset values and global average with a maximum of 10000 in CLANS Toolkit
614 analysis. The P -value $\leq 1e^{-05}$ was used to show the lines connecting the sequences. (B) Sub-
615 clustering of the cluster 1 and 2 protein sequences of CRESS-DNA viruses Rep proteins into
616 10 different clusters (cluster a to j) at a P -value threshold of $1e^{-38}$.

617 **Supplementary Figure 2:** A total of 1823 amino acid sequences of Cap protein of CRESS-
618 DNA viruses were (**Supplementary Data 6**) used and classified by their pairwise sequence
619 similarity network using CLANS. The clusters were classified using the Network-based
620 method using offset values and global average with a maximum of 10000 in CLANS Toolkit
621 analysis. The P -value $\leq 1e^{-02}$ was used to show the lines connecting the sequences.

622 **Supplementary Figure 3: Phylogenetic relationship of CRESS-DNA virus Cap protein**
623 **cluster C3 and C4.** (A) Phylogenetic tree depicting the genetic relationship between the
624 CRESS-DNA virus Cap protein cluster C3. The details of sequences in the cluster
625 (**Supplementary Data 7**) and alignment are provided in **Supplementary Data 12.** (B)
626 Phylogenetic tree depicting the genetic relationship between the CRESS-DNA virus Cap
627 protein cluster C4. The details of sequences in the cluster (**Supplementary Data 7**) and
628 alignment are provided in **Supplementary Data 13.** The maximum-likelihood method
629 inferred the evolutionary history using the Subtree-Pruning-Regrafting algorithm and
630 bootstrap values in PhyML 3.3_1.

631 **Supplementary Figure 4: Phylogenetic relationship of CRESS-DNA virus Cap protein**
632 **cluster C6 and C7.** (A) Phylogenetic tree depicting the genetic relationship between the
633 CRESS-DNA virus Cap protein cluster C6. The details of sequences in the cluster
634 (**Supplementary Data 7**) and alignment are provided in **Supplementary Data 14.** (B)
635 Phylogenetic tree depicting the genetic relationship between the CRESS-DNA virus Cap

636 protein cluster C7. The details of sequences in the cluster (**Supplementary Data 7**) and
637 alignment are provided in **Supplementary Data 15**. The maximum-likelihood method
638 inferred the evolutionary history using the Subtree-Pruning-Regrafting algorithm and
639 bootstrap values in PhyML 3.3_1.

640 **Supplementary Figure 5: Phylogenetic relationship of CRESS-DNA virus Cap protein**
641 **cluster C10, C11 C13, C16, C19, and C20.** (A) Phylogenetic tree depicting the genetic
642 relationship between the CRESS-DNA virus Cap protein cluster C10. The details of
643 sequences in the cluster (**Supplementary Data 7**) and alignment are provided in
644 **Supplementary Data 16.** (B) Phylogenetic tree depicting the genetic relationship between
645 the CRESS-DNA virus Cap protein cluster C11. The details of sequences in the cluster
646 (**Supplementary Data 7**) and alignment are provided in **Supplementary Data 17.** (B)
647 Phylogenetic tree depicting the genetic relationship between the CRESS-DNA virus Cap
648 protein cluster C13, C16, C19, and C20. The details of sequences in each cluster
649 (**Supplementary Data 7**) and alignment are provided in **Supplementary Data 18.** The
650 maximum-likelihood method inferred the evolutionary history using the Subtree-Pruning-
651 Regrafting algorithm and bootstrap values in PhyML 3.3_1.

652 References

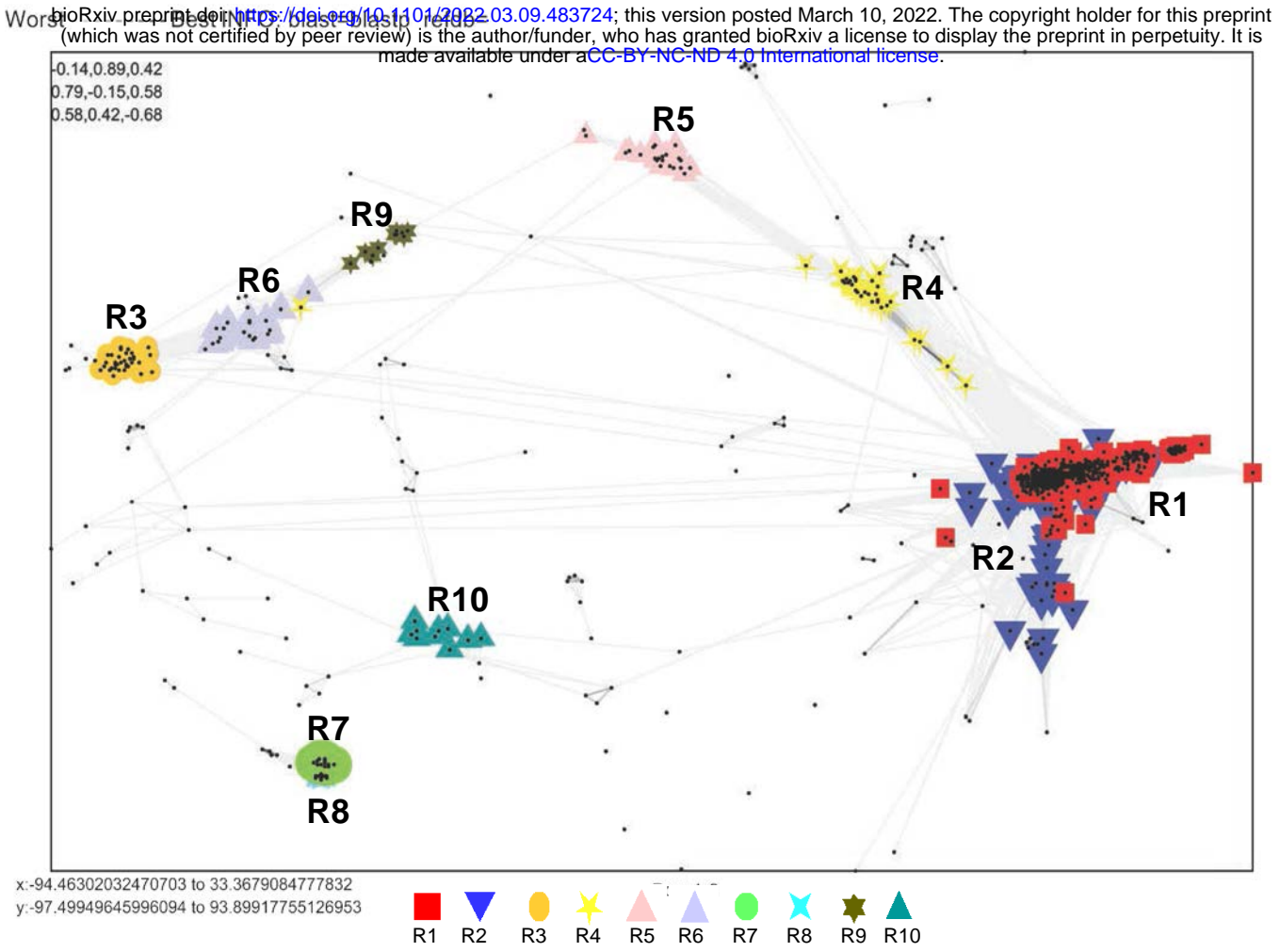
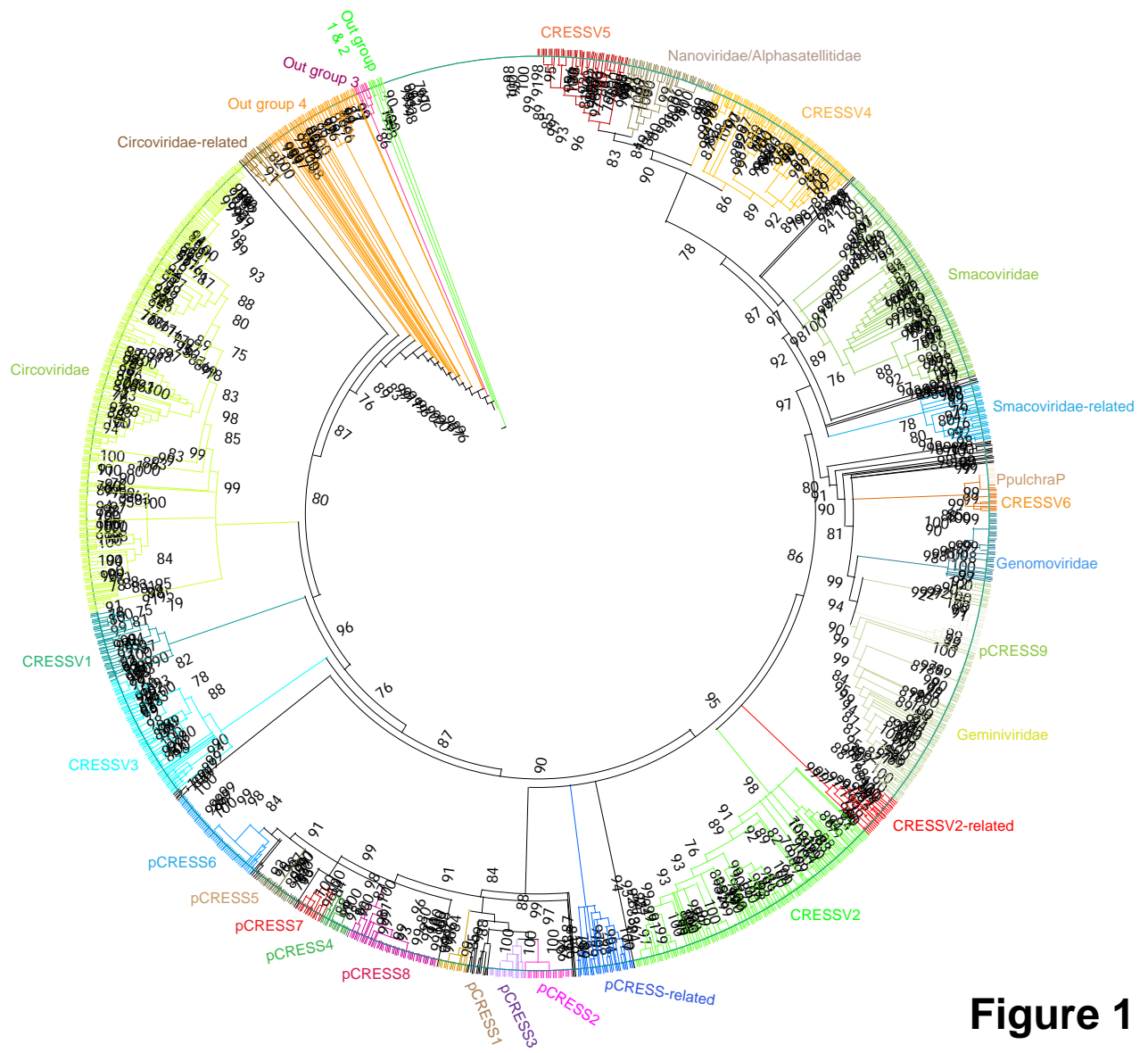
- 653 1 Kazlauskas, D., Varsani, A., Koonin, E. V. & Krupovic, M. Multiple origins of prokaryotic and
654 eukaryotic single-stranded DNA viruses from bacterial and archaeal plasmids. *Nature*
655 *communications* **10**, 3425, doi:10.1038/s41467-019-11433-0 (2019).
- 656 2 Zhao, L., Rosario, K., Breitbart, M. & Duffy, S. Eukaryotic Circular Rep-Encoding Single-
657 Stranded DNA (CRESS DNA) Viruses: Ubiquitous Viruses With Small Genomes and a Diverse
658 Host Range. *Advances in virus research* **103**, 71-133, doi:10.1016/bs.aivir.2018.10.001
659 (2019).
- 660 3 Krupovic, M. Networks of evolutionary interactions underlying the polyphyletic origin of
661 ssDNA viruses. *Current opinion in virology* **3**, 578-586, doi:10.1016/j.coviro.2013.06.010
662 (2013).
- 663 4 Chow, C. E. & Suttle, C. A. Biogeography of Viruses in the Sea. *Annual review of virology* **2**,
664 41-66, doi:10.1146/annurev-virology-031413-085540 (2015).
- 665 5 Labonte, J. M. & Suttle, C. A. Previously unknown and highly divergent ssDNA viruses
666 populate the oceans. *The ISME journal* **7**, 2169-2177, doi:10.1038/ismej.2013.110 (2013).
- 667 6 Ng, T. F. *et al.* High variety of known and new RNA and DNA viruses of diverse origins in
668 untreated sewage. *Journal of virology* **86**, 12161-12175, doi:10.1128/JVI.00869-12 (2012).
- 669 7 Dayaram, A. *et al.* Diverse circular replication-associated protein encoding viruses circulating
670 in invertebrates within a lake ecosystem. *Infection, genetics and evolution : journal of*
671 *molecular epidemiology and evolutionary genetics in infectious diseases* **39**, 304-316,
672 doi:10.1016/j.meegid.2016.02.011 (2016).
- 673 8 Rosario, K., Schenck, R. O., Harbeitner, R. C., Lawler, S. N. & Breitbart, M. Novel circular
674 single-stranded DNA viruses identified in marine invertebrates reveal high sequence
675 diversity and consistent predicted intrinsic disorder patterns within putative structural
676 proteins. *Frontiers in microbiology* **6**, 696, doi:10.3389/fmicb.2015.00696 (2015).
- 677 9 Dayaram, A. *et al.* Diverse small circular DNA viruses circulating amongst estuarine molluscs.
678 *Infection, genetics and evolution : journal of molecular epidemiology and evolutionary*
679 *genetics in infectious diseases* **31**, 284-295, doi:10.1016/j.meegid.2015.02.010 (2015).

- 680 10 Bistolas, K. S. I., Rudstam, L. G. & Hewson, I. Gene expression of benthic amphipods (genus:
681 Diporeia) in relation to a circular ssDNA virus across two Laurentian Great Lakes. *PeerJ* **5**,
682 e3810, doi:10.7717/peerj.3810 (2017).
- 683 11 Blinkova, O. *et al.* Frequent detection of highly diverse variants of cardiovirus, cosavirus,
684 bocavirus, and circovirus in sewage samples collected in the United States. *Journal of clinical*
685 *microbiology* **47**, 3507-3513, doi:10.1128/JCM.01062-09 (2009).
- 686 12 Krupovic, M. *et al.* Cressdnaviricota: a Virus Phylum Unifying Seven Families of Rep-Encoding
687 Viruses with Single-Stranded, Circular DNA Genomes. *Journal of virology* **94**,
688 doi:10.1128/JVI.00582-20 (2020).
- 689 13 Kazlauskas, D. *et al.* Evolutionary history of ssDNA bacilladnaviruses features horizontal
690 acquisition of the capsid gene from ssRNA nodaviruses. *Virology* **504**, 114-121,
691 doi:10.1016/j.virol.2017.02.001 (2017).
- 692 14 Diemer, G. S. & Stedman, K. M. A novel virus genome discovered in an extreme environment
693 suggests recombination between unrelated groups of RNA and DNA viruses. *Biology direct* **7**,
694 13, doi:10.1186/1745-6150-7-13 (2012).
- 695 15 Roux, S. *et al.* Chimeric viruses blur the borders between the major groups of eukaryotic
696 single-stranded DNA viruses. *Nature communications* **4**, 2700, doi:10.1038/ncomms3700
697 (2013).
- 698 16 Krupovic, M., Ravantti, J. J. & Bamford, D. H. Geminiviruses: a tale of a plasmid becoming a
699 virus. *BMC evolutionary biology* **9**, 112, doi:10.1186/1471-2148-9-112 (2009).
- 700 17 de la Higuera, I. *et al.* Unveiling Crucivirus Diversity by Mining Metagenomic Data. *mBio* **11**,
701 doi:10.1128/mBio.01410-20 (2020).
- 702 18 Yoon, H. S. *et al.* Single-cell genomics reveals organismal interactions in uncultivated marine
703 protists. *Science* **332**, 714-717, doi:10.1126/science.1203163 (2011).
- 704 19 Kazlauskas, D., Varsani, A. & Krupovic, M. Pervasive Chimerism in the Replication-Associated
705 Proteins of Uncultured Single-Stranded DNA Viruses. *Viruses* **10**, doi:10.3390/v10040187
706 (2018).
- 707 20 Simmonds, P. *et al.* Consensus statement: Virus taxonomy in the age of metagenomics.
708 *Nature reviews. Microbiology* **15**, 161-168, doi:10.1038/nrmicro.2016.177 (2017).
- 709 21 Zimmermann, L. *et al.* A Completely Reimplemented MPI Bioinformatics Toolkit with a New
710 HHpred Server at its Core. *J Mol Biol* **430**, 2237-2243, doi:10.1016/j.jmb.2017.12.007 (2018).
- 711 22 Frickey, T. & Lupas, A. CLANS: a Java application for visualizing protein families based on
712 pairwise similarity. *Bioinformatics* **20**, 3702-3704, doi:10.1093/bioinformatics/bth444
713 (2004).
- 714 23 Krupovic, M. & Koonin, E. V. Multiple origins of viral capsid proteins from cellular ancestors.
715 *Proc Natl Acad Sci U S A* **114**, E2401-E2410, doi:10.1073/pnas.1621061114 (2017).
- 716 24 Whitley, C. *et al.* Novel replication-competent circular DNA molecules from healthy cattle
717 serum and milk and multiple sclerosis-affected human brain tissue. *Genome announcements*
718 **2**, doi:10.1128/genomeA.00849-14 (2014).
- 719 25 Gorbalenya, A. E., Koonin, E. V. & Wolf, Y. I. A new superfamily of putative NTP-binding
720 domains encoded by genomes of small DNA and RNA viruses. *FEBS letters* **262**, 145-148,
721 doi:10.1016/0014-5793(90)80175-i (1990).
- 722 26 Lu, S. *et al.* CDD/SPARCLE: the conserved domain database in 2020. *Nucleic acids research*
723 **48**, D265-D268, doi:10.1093/nar/gkz991 (2020).
- 724 27 Marchler-Bauer, A. *et al.* CDD/SPARCLE: functional classification of proteins via subfamily
725 domain architectures. *Nucleic acids research* **45**, D200-D203, doi:10.1093/nar/gkw1129
726 (2017).
- 727 28 Marchler-Bauer, A. *et al.* CDD: NCBI's conserved domain database. *Nucleic acids research* **43**,
728 D222-226, doi:10.1093/nar/gku1221 (2015).
- 729 29 Marchler-Bauer, A. *et al.* CDD: a Conserved Domain Database for the functional annotation
730 of proteins. *Nucleic acids research* **39**, D225-229, doi:10.1093/nar/gkq1189 (2011).

- 731 30 Abrescia, N. G., Bamford, D. H., Grimes, J. M. & Stuart, D. I. Structure unifies the viral
732 universe. *Annu Rev Biochem* **81**, 795-822, doi:10.1146/annurev-biochem-060910-095130
733 (2012).
- 734 31 Greene, L. H. *et al.* The CATH domain structure database: new protocols and classification
735 levels give a more comprehensive resource for exploring evolution. *Nucleic acids research*
736 **35**, D291-297, doi:10.1093/nar/gkl959 (2007).
- 737 32 Krupovic, M. & Bamford, D. H. Double-stranded DNA viruses: 20 families and only five
738 different architectural principles for virion assembly. *Current opinion in virology* **1**, 118-124,
739 doi:10.1016/j.coviro.2011.06.001 (2011).
- 740 33 Du, M. Z. *et al.* The GC Content as a Main Factor Shaping the Amino Acid Usage During
741 Bacterial Evolution Process. *Front Microbiol* **9**, 2948, doi:10.3389/fmicb.2018.02948 (2018).
- 742 34 Bohlin, J., Brynildsrud, O., Vesth, T., Skjerve, E. & Ussery, D. W. Amino acid usage is
743 asymmetrically biased in AT- and GC-rich microbial genomes. *PLoS One* **8**, e69878,
744 doi:10.1371/journal.pone.0069878 (2013).
- 745 35 Zhao, Y. *et al.* Analysis of codon usage bias of envelope glycoprotein genes in nuclear
746 polyhedrosis virus (NPV) and its relation to evolution. *BMC Genomics* **17**, 677,
747 doi:10.1186/s12864-016-3021-7 (2016).
- 748 36 Wang, L. *et al.* Genome-wide analysis of codon usage bias in four sequenced cotton species.
749 *PLoS One* **13**, e0194372, doi:10.1371/journal.pone.0194372 (2018).
- 750 37 Gun, L., Yumiao, R., Haixian, P. & Liang, Z. Comprehensive Analysis and Comparison on the
751 Codon Usage Pattern of Whole Mycobacterium tuberculosis Coding Genome from Different
752 Area. *Biomed Res Int* **2018**, 3574976, doi:10.1155/2018/3574976 (2018).
- 753 38 Tian, H. F. *et al.* Genetic and codon usage bias analyses of major capsid protein gene in
754 Ranavirus. *Infect Genet Evol* **84**, 104379, doi:10.1016/j.meegid.2020.104379 (2020).
- 755 39 Guo, Z., He, Q., Tang, C., Zhang, B. & Yue, H. Identification and genomic characterization of a
756 novel CRESS DNA virus from a calf with severe hemorrhagic enteritis in China. *Virus research*
757 **255**, 141-146, doi:10.1016/j.virusres.2018.07.015 (2018).
- 758 40 Moens, M. A. J., Perez-Tris, J., Cortey, M. & Benitez, L. Identification of two novel CRESS DNA
759 viruses associated with an Avipoxvirus lesion of a blue-and-gray Tanager (*Thraupis*
760 *episcopus*). *Infection, genetics and evolution : journal of molecular epidemiology and*
761 *evolutionary genetics in infectious diseases* **60**, 89-96, doi:10.1016/j.meegid.2018.02.015
762 (2018).
- 763 41 Liu, Q. *et al.* Viral metagenomics revealed diverse CRESS-DNA virus genomes in faeces of
764 forest musk deer. *Virology journal* **17**, 61, doi:10.1186/s12985-020-01332-y (2020).
- 765 42 Lemoine, F. *et al.* Renewing Felsenstein's phylogenetic bootstrap in the era of big data.
766 *Nature* **556**, 452-456, doi:10.1038/s41586-018-0043-0 (2018).
- 767 43 Guindon, S. *et al.* New algorithms and methods to estimate maximum-likelihood
768 phylogenies: assessing the performance of PhyML 3.0. *Systematic biology* **59**, 307-321,
769 doi:10.1093/sysbio/syq010 (2010).
- 770 44 Lemoine, F. *et al.* NGPhylogeny.fr: new generation phylogenetic services for non-specialists.
771 *Nucleic acids research* **47**, W260-W265, doi:10.1093/nar/gkz303 (2019).
- 772 45 Letunic, I. & Bork, P. Interactive Tree Of Life (iTOL) v5: an online tool for phylogenetic tree
773 display and annotation. *Nucleic acids research* **49**, W293-W296, doi:10.1093/nar/gkab301
774 (2021).
- 775 46 Vetrivel, U., Arunkumar, V. & Dorairaj, S. ACUA: a software tool for automated codon usage
776 analysis. *Bioinformatics* **2**, 62-63, doi:10.6026/97320630002062 (2007).
- 777 47 Puigbo, P., Bravo, I. G. & Garcia-Vallve, S. CALcal: a combined set of tools to assess codon
778 usage adaptation. *Biol Direct* **3**, 38, doi:10.1186/1745-6150-3-38 (2008).

779

780

A**B****Figure 1**

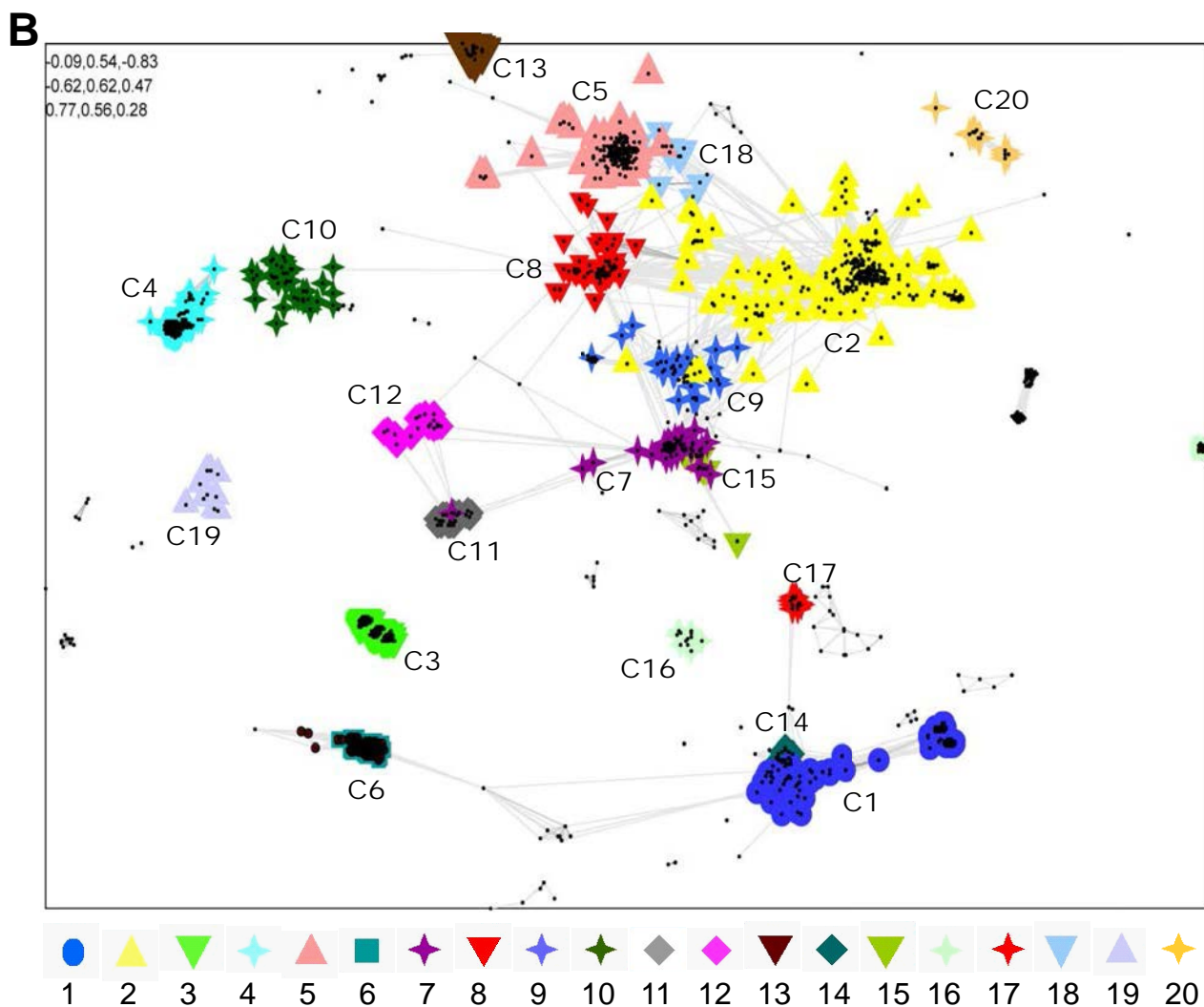
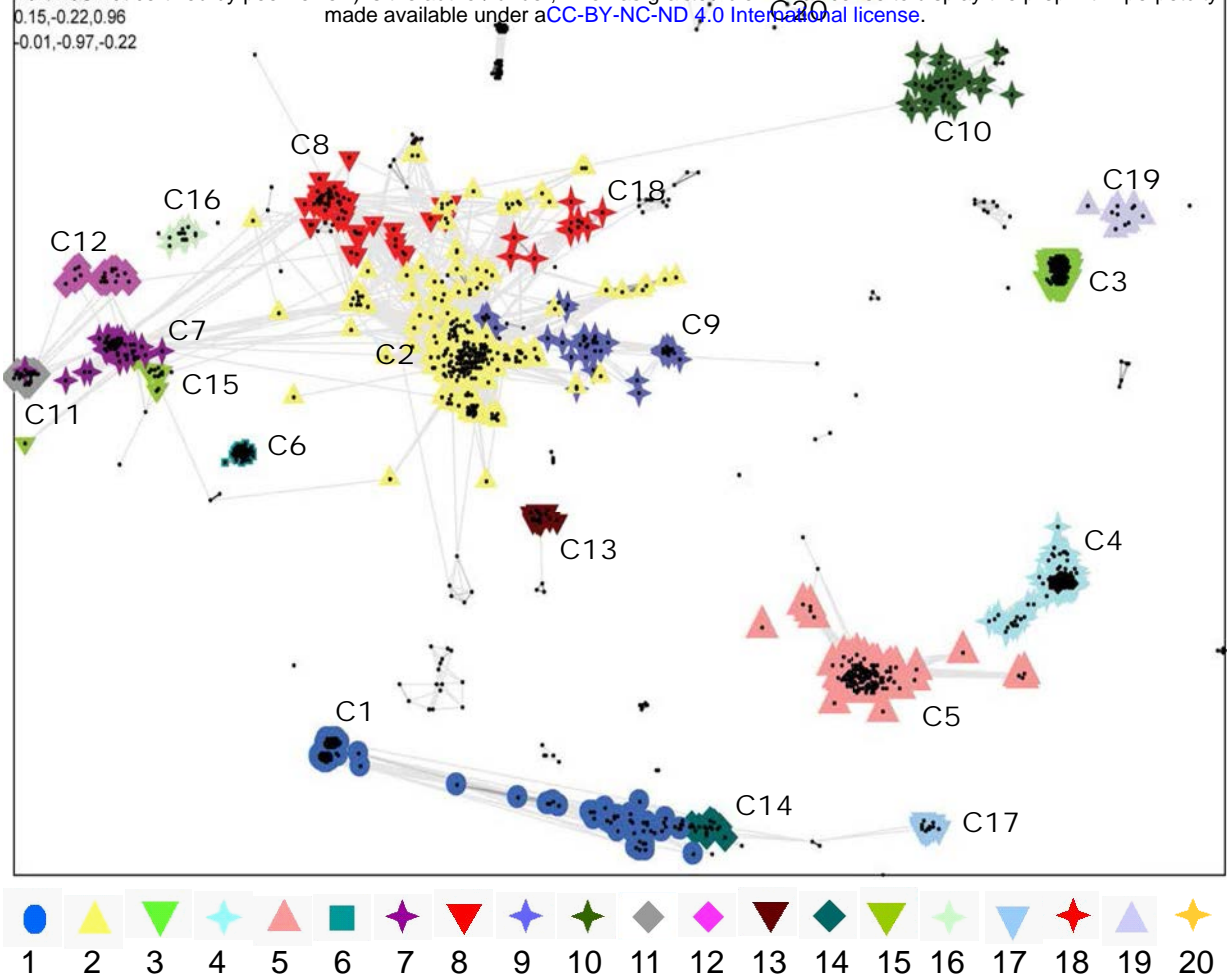
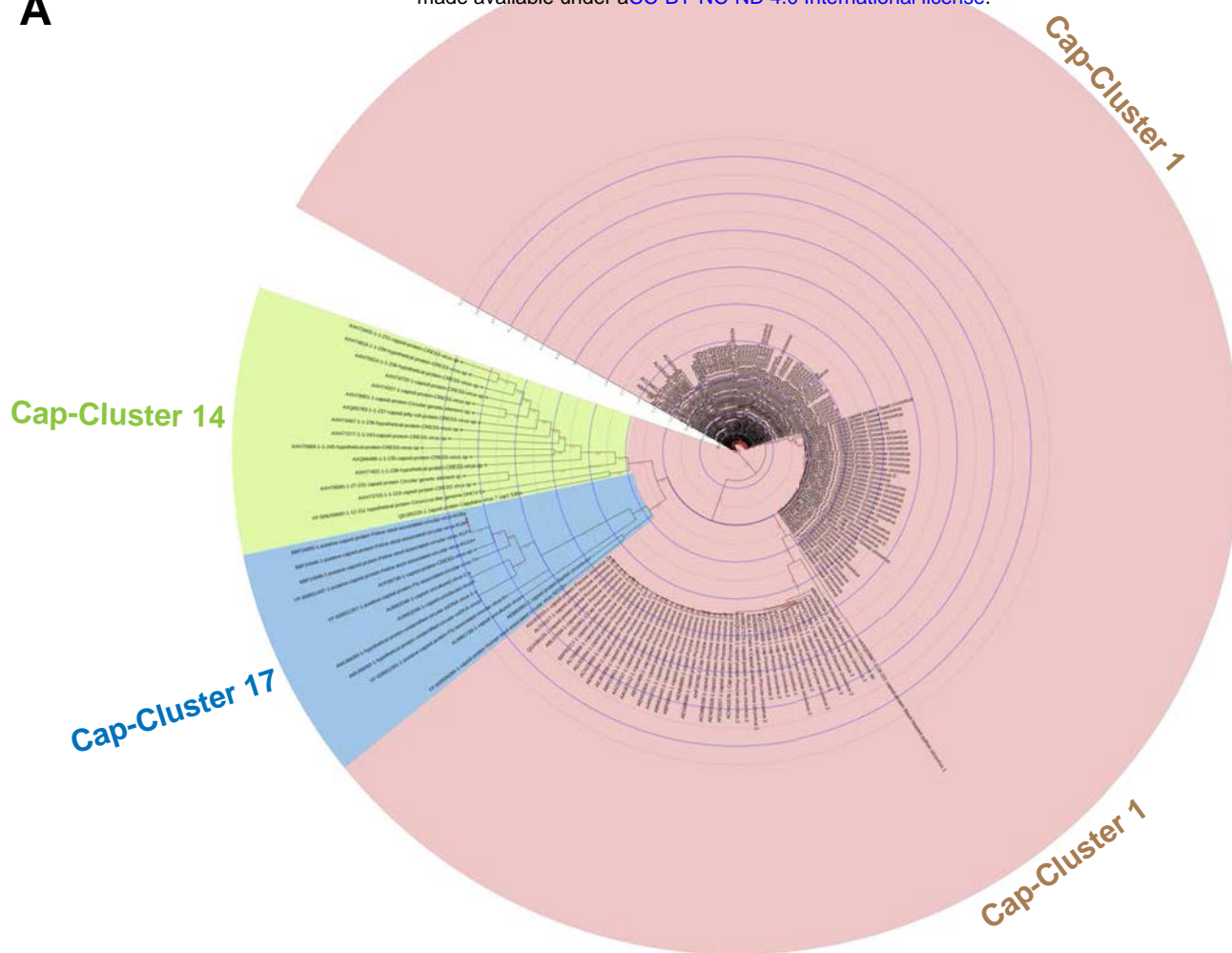


Figure 2

A



B

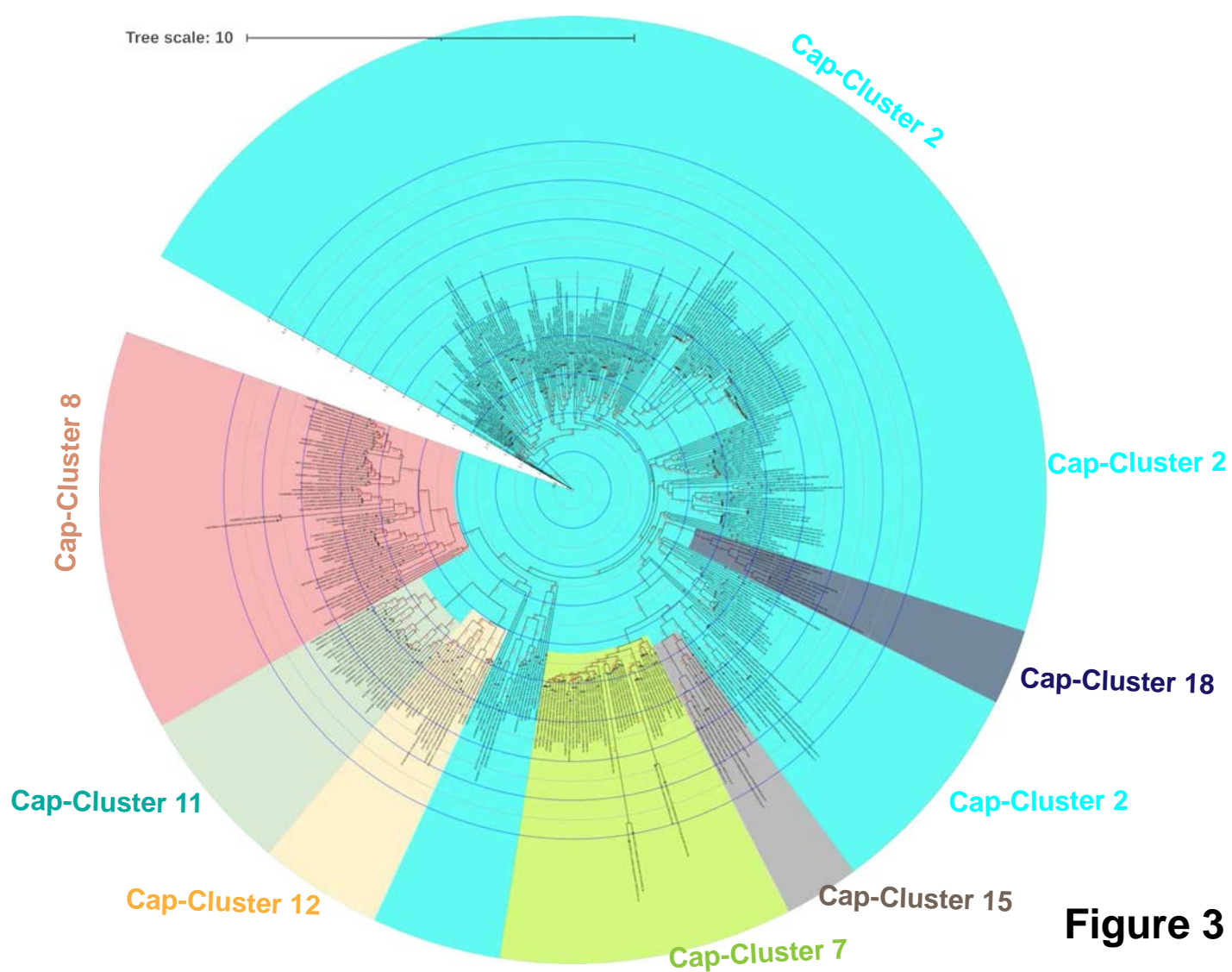


Figure 3

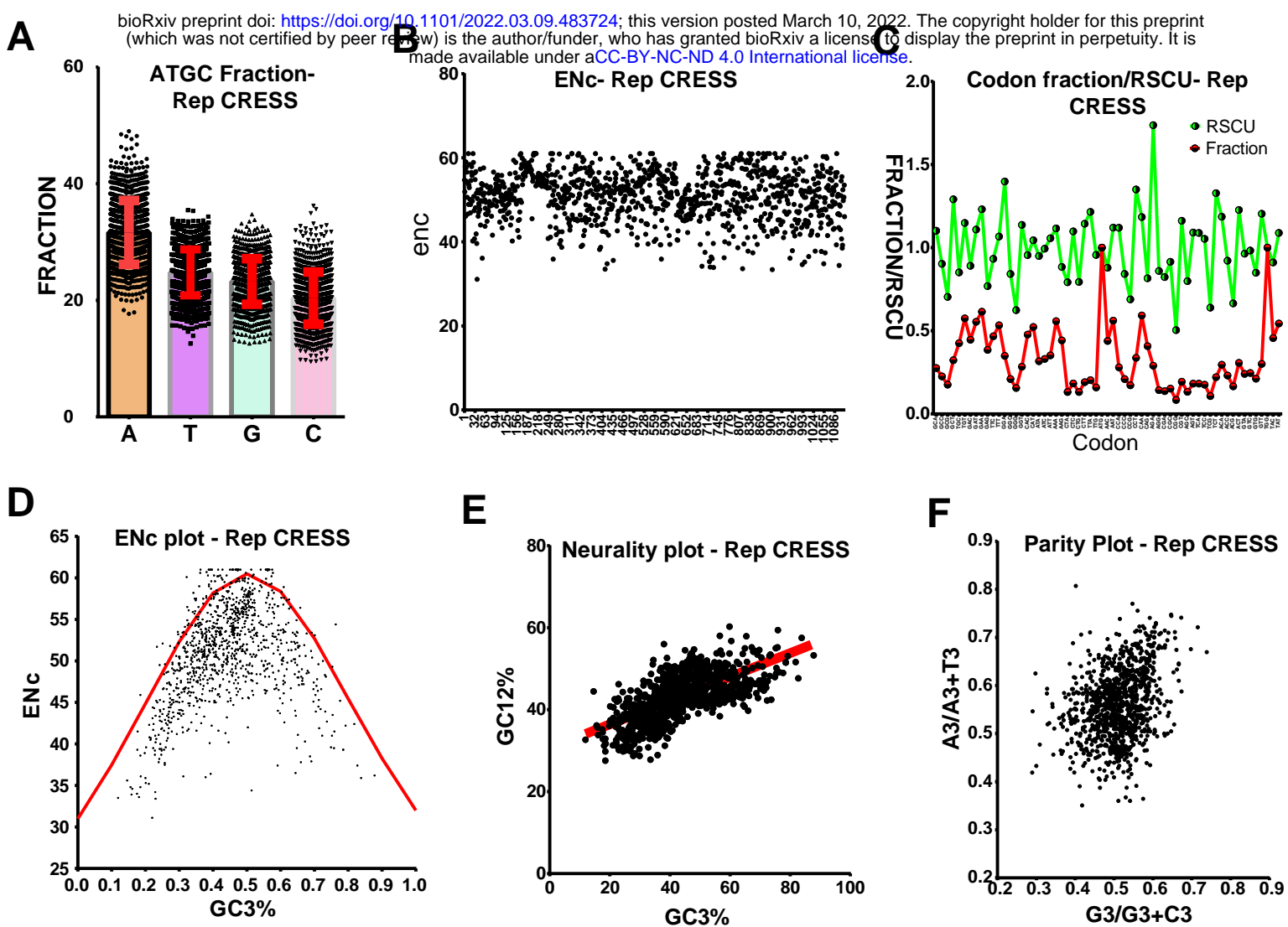


Figure 4

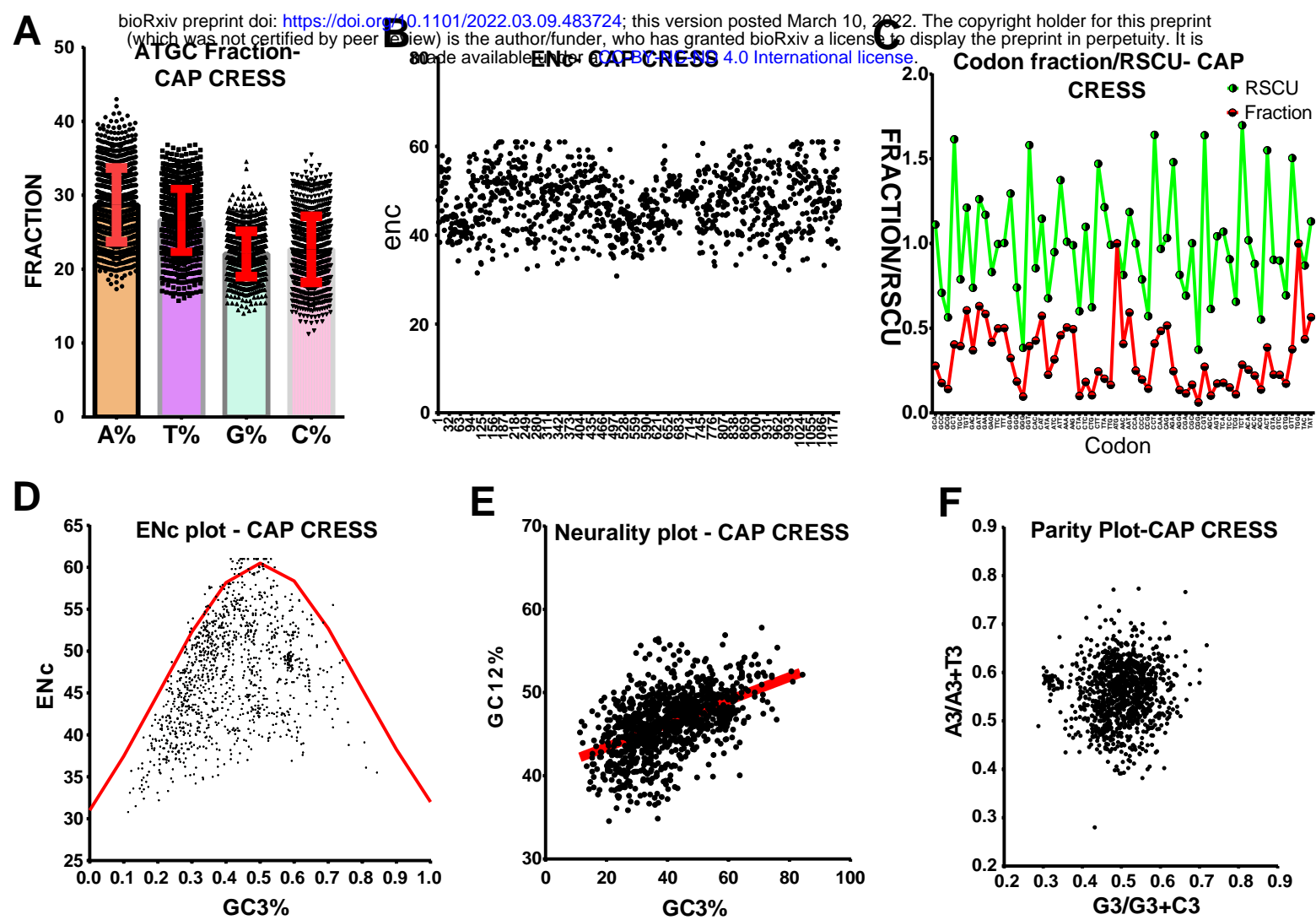
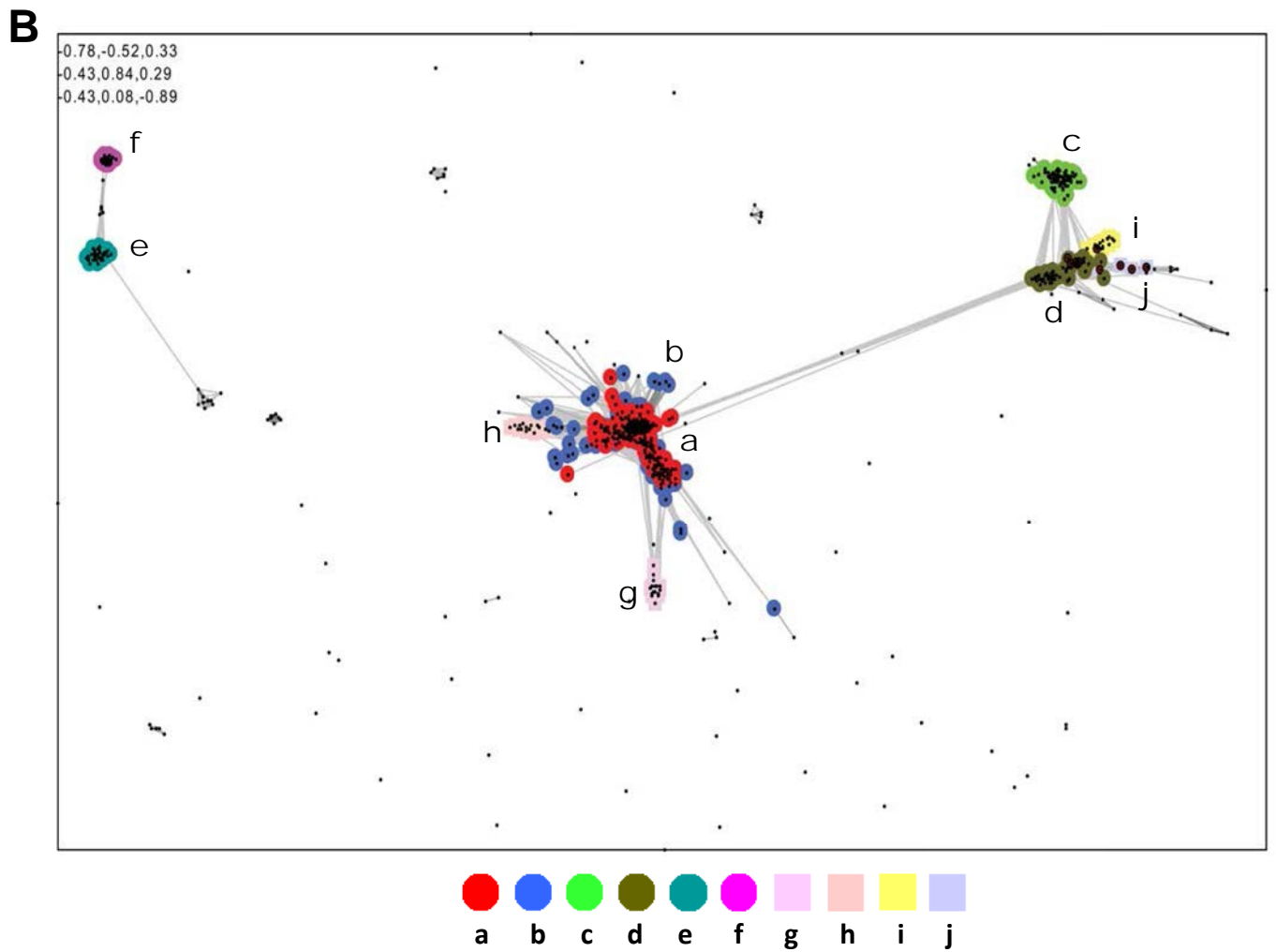
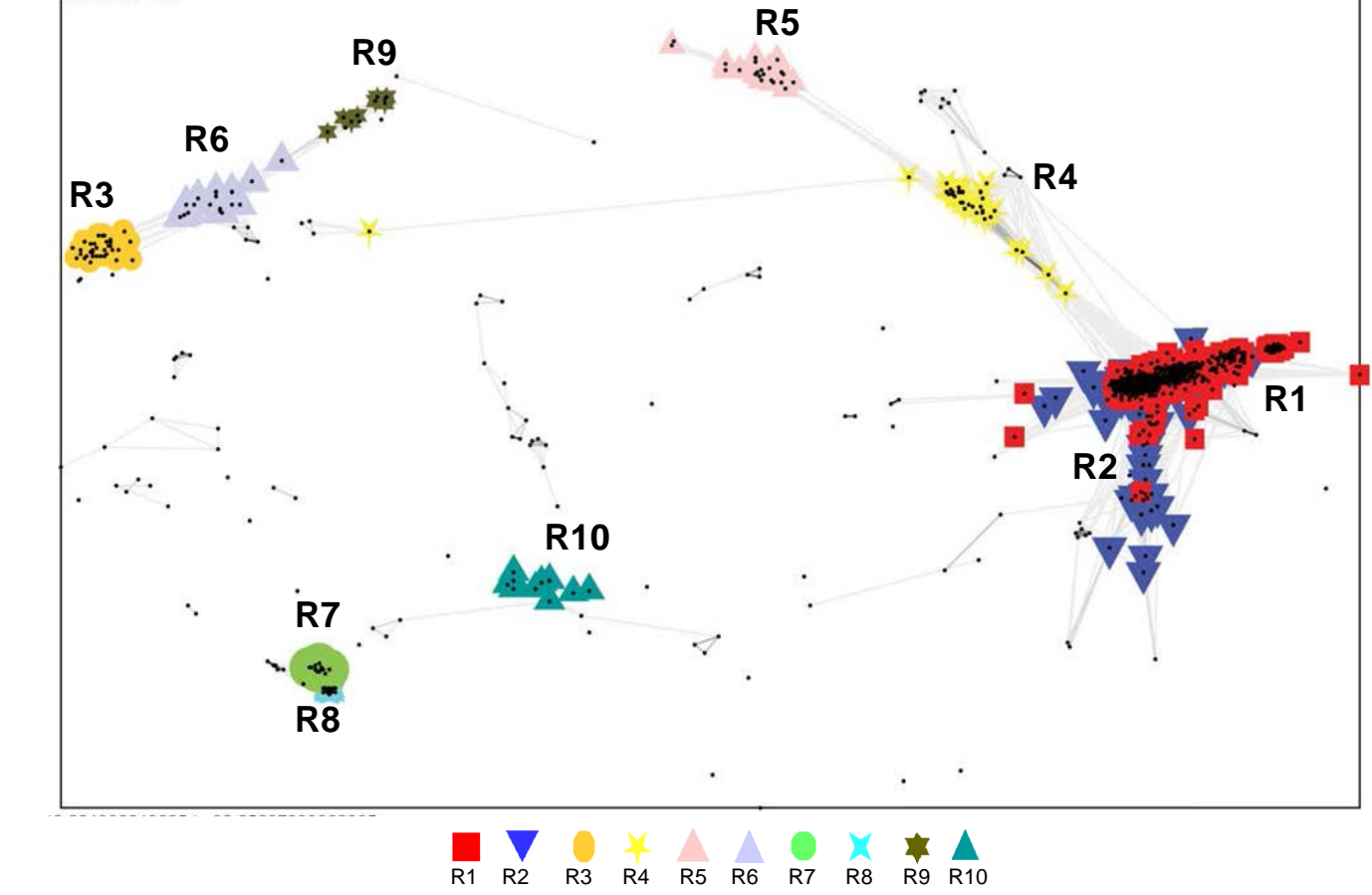
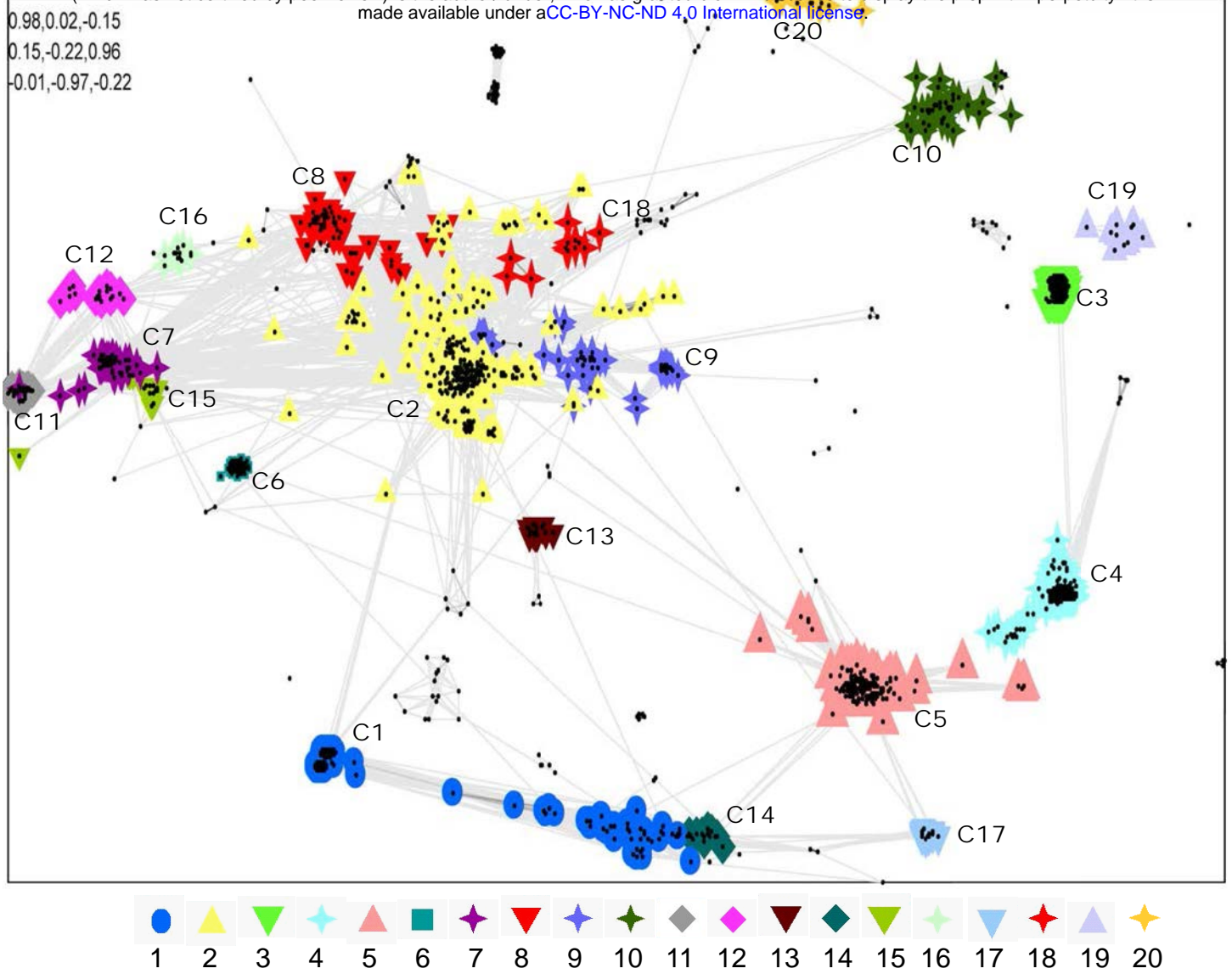


Figure 5

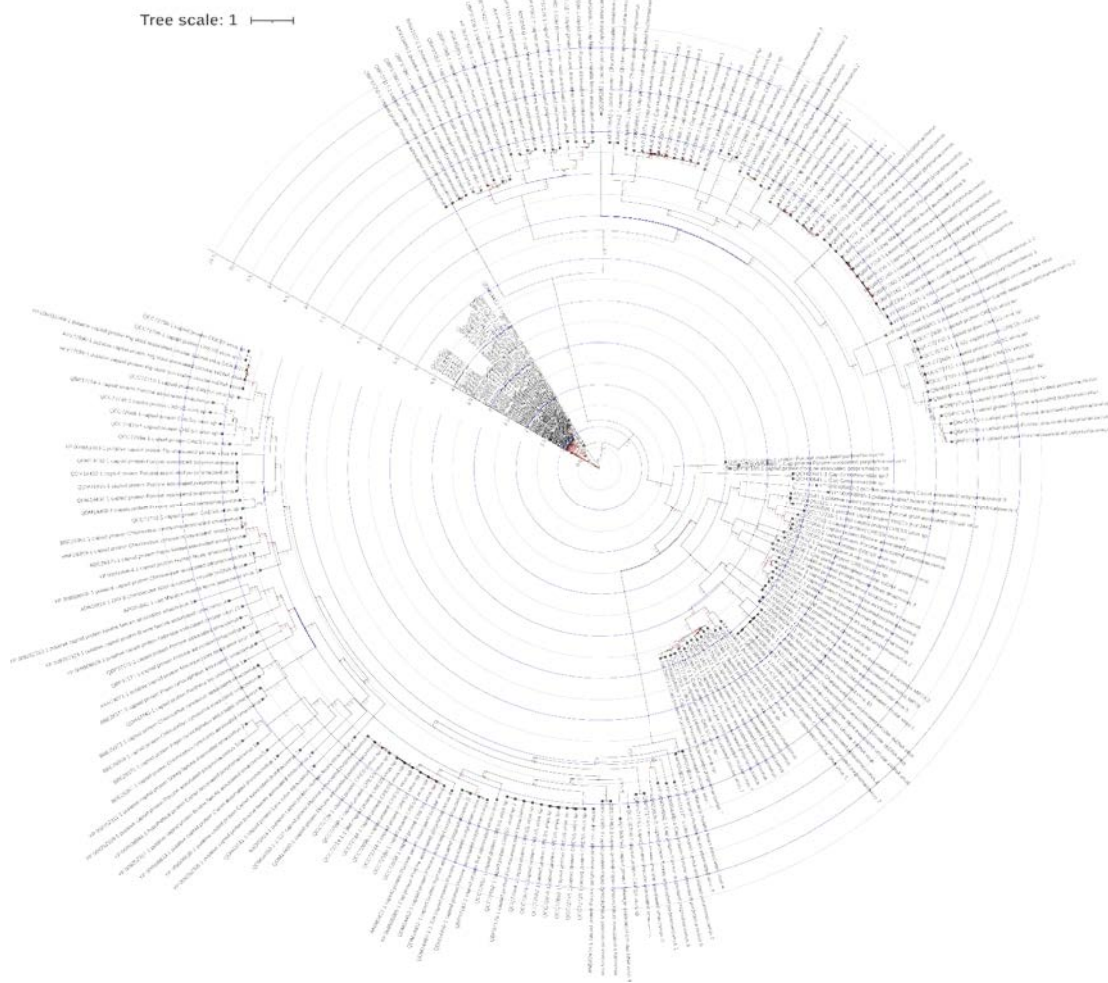


Supplementary Figure 1

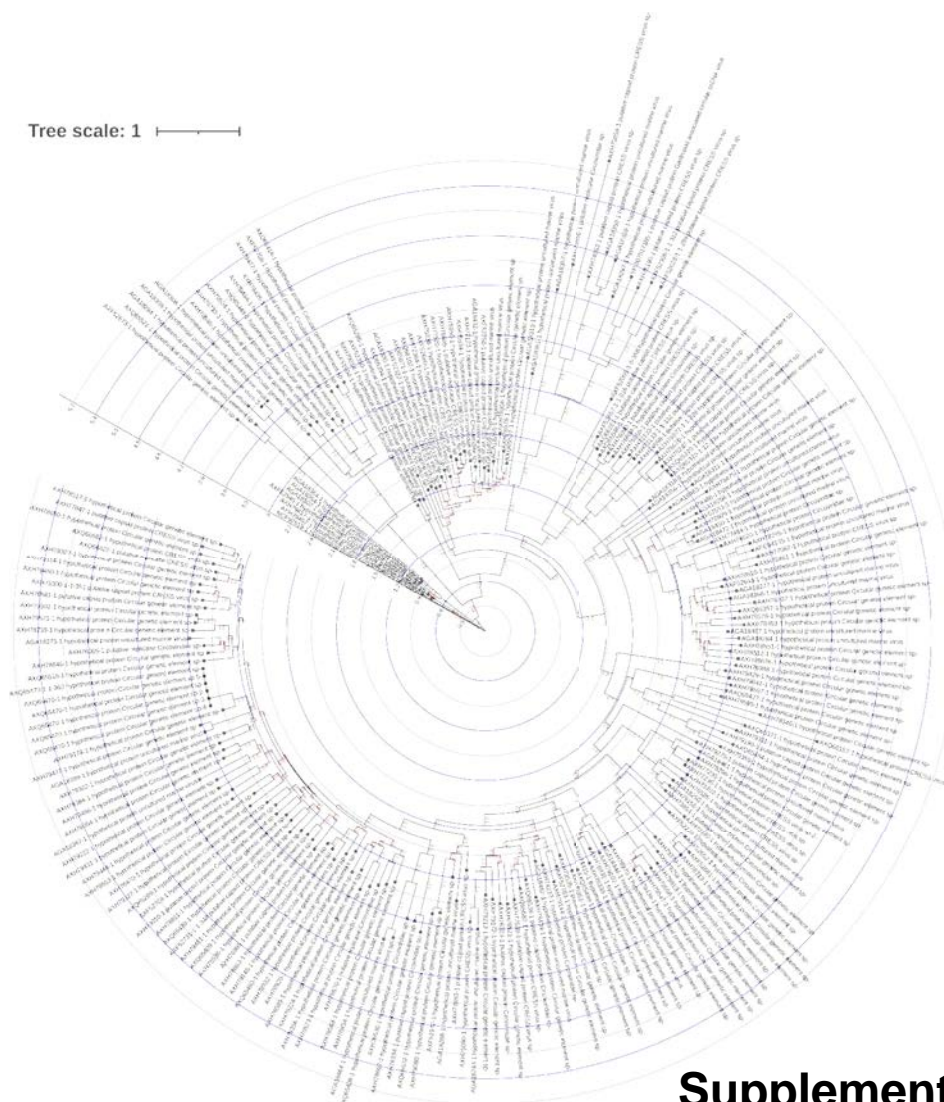


Supplementary Figure 2

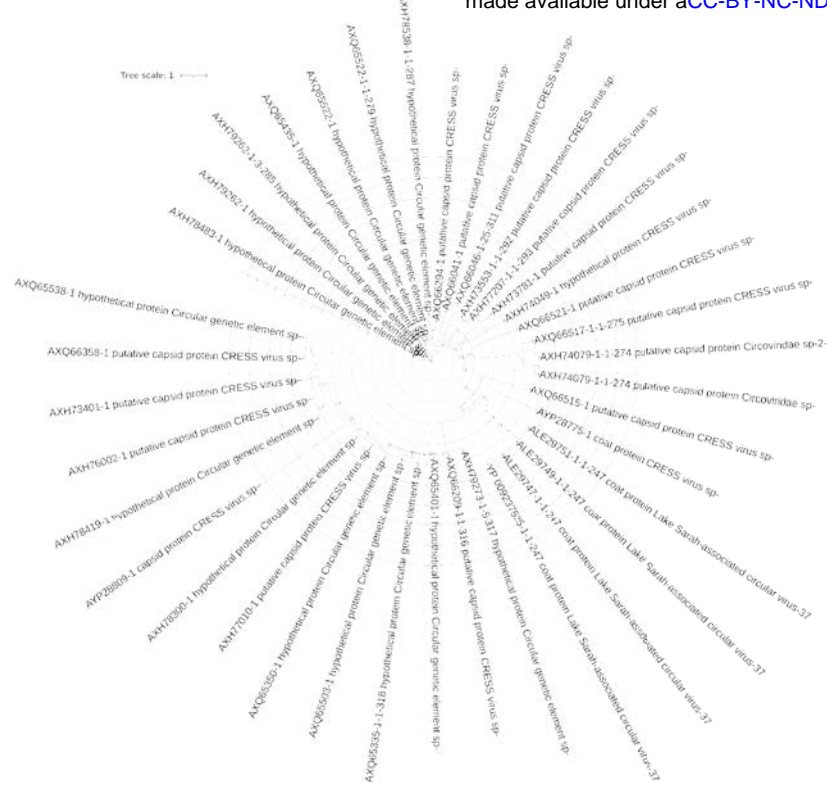
A



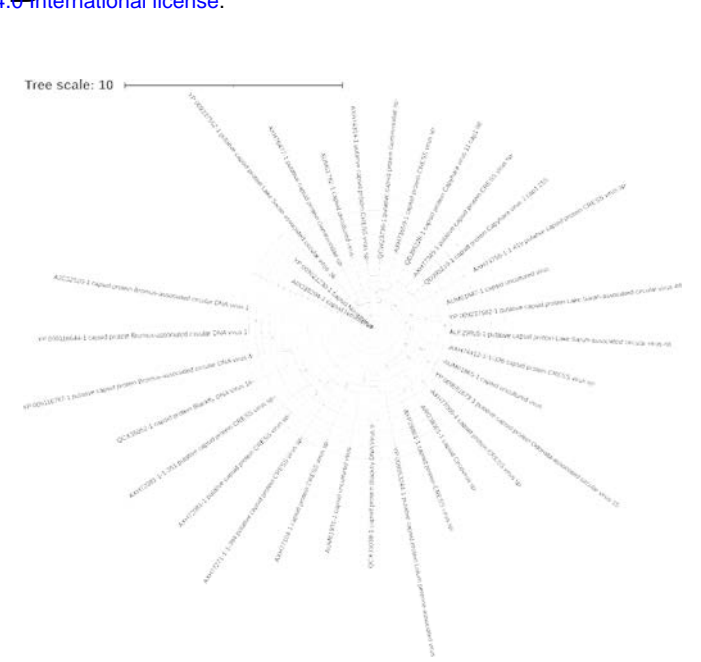
B



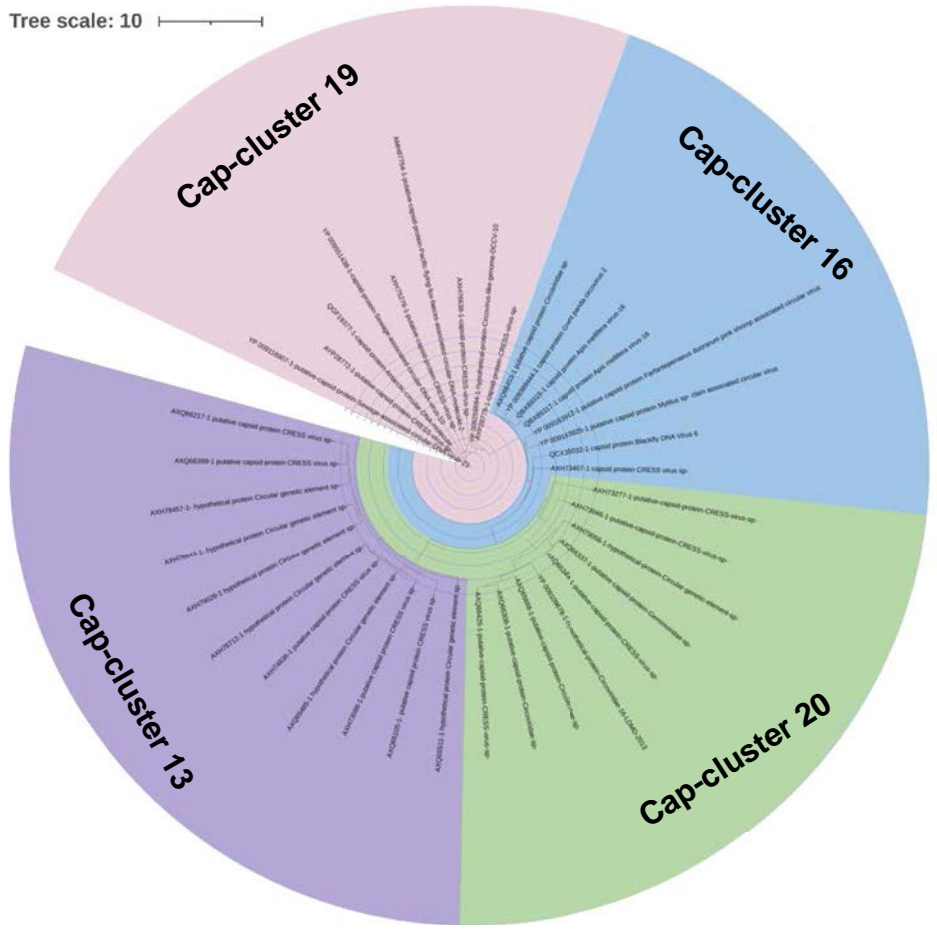
A



B



C



Supplementary Figure 5