

1 **Endogenous giant viruses contribute to intraspecies genomic variability in the model**
2 **green alga *Chlamydomonas reinhardtii***

3
4 Mohammad Moniruzzaman*¹ and Frank O. Aylward*^{1,2}
5

6 1 Department of Biological Sciences, 926 West Campus Drive, Virginia Tech, Blacksburg, VA
7 2 Center for Emerging, Zoonotic, and Arthropod-Borne Pathogens, Virginia Tech, Blacksburg,
8 VA

9 Email address for correspondence: monir@vt.edu, faylward@vt.edu
10

11 **Abstract:**
12

13 *Chlamydomonas reinhardtii* is an important eukaryotic alga that has been studied as a model
14 organism for decades. Despite extensive history as a model system, phylogenetic and genetic
15 characteristics of viruses infecting this alga have remained elusive. We analyzed high-
16 throughput genome sequence data of *C. reinhardtii* field isolates, and in six we discovered
17 sequences belonging to endogenous giant viruses that reach up to several hundred kilobases in
18 length. In addition, we have also discovered the entire genome of a closely related giant virus
19 that is endogenized within the genome of *Chlamydomonas incerta*, the closest sequenced
20 phylogenetic relatives of *C. reinhardtii*. Endogenous giant viruses add hundreds of new gene
21 families to the host strains, highlighting their contribution to the pangenome dynamics and inter-
22 strain genomic variability of *C. reinhardtii*. Our findings suggest that the endogenization of giant
23 viruses can have important implications for structuring the population dynamics and ecology of
24 protists in the environment.
25

26 **Introduction:**
27

28 *Chlamydomonas reinhardtii* is a widely studied unicellular green alga with a long history as a
29 model organism that dates back to the 1950s (Sasso *et al.*, 2018; Salomé & Merchant, 2019).
30 Despite this long history of research, no viruses that infect *C. reinhardtii* have yet been reported,
31 and as a result the diversity of viruses that infect this alga in nature remain unknown. In a recent
32 study, we identified widespread endogenization of “giant viruses” in numerous green algae,
33 which provides evidence of virus-host interactions that take place in nature (Moniruzzaman *et*
34 *al.*, 2020b). These Giant Endogenous Viral Elements (GEVEs) derive from giant viruses within
35 the phylum *Nucleocytoviricota*, which possess large and complex genomes that can reach up to
36 2.5 Mbp in length (Philippe *et al.*, 2013). Giant viruses often encode complex functional
37 repertoires in their genomes that include tRNA synthetases, rhodopsins, cytoskeletal
38 components, histones, and proteins involved in glycolysis, the TCA cycle, and other aspects of
39 central carbon metabolism (Aylward *et al.*). Recent studies have shown that giant viruses are
40 widespread in the environment and infect a wide range of eukaryotic hosts, including green
41 algae (Schulz *et al.*, 2020; Moniruzzaman *et al.*, 2020a; Endo *et al.*, 2020; Meng *et al.*, 2021).
42 The complex genomes of giant viruses coupled with their collectively broad host range and
43 ability to endogenize into the genomes of their hosts provides compelling evidence that they
44 may be important vectors of gene transfer in eukaryotes.

45 In our initial genomic survey of GEVEs we did not find evidence of endogenous giant viruses in
46 the type strain *C. reinhardtii* (CC-503 cw92). Several studies have recently reported draft
47 genomes of *C. reinhardtii* field isolates, however, and in this study we surveyed these strains for
48 evidence of GEVEs. We report that near-complete genomes of giant viruses are present in
49 several field isolates, and our results suggest that *C. reinhardtii* is a host to at least two distinct
50 lineages of giant viruses. These are the first insights into the diversity and genomic complexity
51 of viruses infecting *C. reinhardtii* in nature. We anticipate that this widely-studied green alga will
52 be a valuable model for future studies of virus-host interactions and the mechanistic aspects of
53 giant virus endogenization.

54

55

56 Results

57

58 We analyzed publicly available high-throughput genome sequencing data for 33 wild strains of
59 *C. reinhardtii*. This data was originally generated for population genomic studies of diverse *C.*
60 *reinhardtii* strains (Flowers *et al.*, 2015; Craig *et al.*, 2019; Hasan *et al.*, 2019). After *de novo*
61 assembly and annotation (see Methods for details), we identified GEVEs in six of the wild
62 strains (Figure 1A,B). In five of these (CC-2936, 2937, 2938, 3268, and GB-66), the GEVEs
63 range from 315-356 Kb in size and harbored all but one *Nucleocytoviricota* hallmark genes,
64 indicating that near-complete genomes of endogenous giant viruses have been retained in
65 these strains (Figure 1B, Dataset S1). In contrast, CC-3061 harbors a GEVE ~113 Kb in size
66 with 5 out of the 10 hallmark genes, indicating that part of the GEVE was lost over the course of
67 evolution (Supplementary Methods, Dataset S1). Moreover, to ensure that GEVEs were not
68 omitted due to assembly issues we also mapped reads from all genome sequencing projects
69 against the GEVEs, and we identified another highly fragmented GEVE in CC-3059 (see
70 Methods). Lastly, we also analyzed the assembled genome of *Chlamydomonas incerta*, a
71 species phylogenetically closest to *C. reinhardtii*, for which a long-read assembled genome has
72 been recently reported (Craig *et al.*, 2021). This analysis revealed a GEVE ~475 Kb long which
73 is integrated within a single 592 Kb contig of this alga (Figure 1B).

74

75 Using a newly established taxonomy of *Nucleocytoviricota* (Aylward *et al.* 2021), we determined
76 the phylogenetic position of the *C. reinhardtii* and *C. incerta* GEVEs and their relationships with
77 other chlorophyte GEVEs that were recently reported (Moniruzzaman *et al.*, 2020b) (Figure 1A).
78 Five of the strains harbored GEVEs that formed a cluster within the *Imitervirales* order,
79 consistent with their high pairwise average amino acid identity. The GEVE in *C. incerta* was the
80 closest phylogenetic relative of the *Imitevirales* GEVEs in *C. reinhardtii*, indicating that closely-
81 related giant viruses infect closely related *Chlamydomonas* species in nature. These GEVEs
82 formed a sister clade with the GEVEs present in six other volvocine algae and belonged to the
83 *Imitevirales* family 12 (Figure 1A). Although GEVE contigs could not be recovered from CC-
84 3059, read mapping revealed that this strain also harbors a fragmented *Imitervirales* GEVE (see
85 Methods). In contrast to the GEVEs that could be classified as *Imitervirales*, the GEVE in CC-
86 2938 strain belonged to the *Algavirales* (Figure 1A), indicating that *C. reinhardtii* is infected by
87 multiple phylogenetically distinct lineages of giant viruses in nature.

88

89 The coverage of the GEVE contigs was generally similar to those of the host *Chlamydomonas*
90 contigs (see Supplementary Information), consistent with their presence as endogenous
91 elements. The exception was the GEVE in CC-2938, in which two large contigs exhibited the
92 same coverage as those of the host (~8 reads per kilobase per million), while the remaining
93 GEVE contigs had coverage roughly twice that. This unusual pattern may be the product of
94 recent large-scale duplication which recently took place in part of this GEVE. Indeed, recent
95 work on other GEVEs in green algae found that large-scale duplications are common in GEVEs
96 (Moniruzzaman *et al.*, 2020b). This would explain why two large contigs with a summed length
97 of 109 kbp retain similar coverage compared to the host contigs, while the rest of the GEVE
98 contigs have roughly double that coverage.

99
100 The % GC-content of the *C. reinhardtii* GEVEs ranged from 58.27% (CC-2938) to 60.72% (CC-
101 3268), which is similar to the overall genomic GC content of *C. reinhardtii* (64%) (Merchant *et*
102 *al.*, 2007). Similarly, the GC content of the *C. incerta* GEVE was 64.8%, resembling the overall
103 GC content of the *C. incerta* genome (66%) (Craig *et al.*, 2021) (Figure 1B). The GEVEs also
104 contained several predicted spliceosomal introns, ranging from 25 (CC-3061) to 72 (*C. incerta*).
105 Spliceosomal introns are rare in free *Nucleocyotiviricota* but have been previously found in
106 GEVEs present in other members of the *Chlorophyta* (Moniruzzaman *et al.*, 2020b). It remains
107 unclear if the relatively high %GC content and spliceosomal introns are features of the viruses
108 themselves, or if the evolution of these features evolved after endogenization. In addition, the
109 GEVE in *C. incerta* was flanked by highly repetitive regions on both ends (Figure 2A). The
110 repetitive region at the 5'-end harbors several reverse transcriptases and transposases (Dataset
111 S1). These regions also have higher intron density compared to the GEVE region itself, and
112 lower number of Giant Virus Orthologous Group (GVOG) hits consistent with their eukaryotic
113 provenance (Figure 2A). This suggests that near-complete genomes of giant viruses can
114 integrate within highly repetitive regions of eukaryotic genomes, potentially with the facilitation of
115 transposable elements.

116
117 The GEVEs in *C. reinhardtii* encoded 99 (CC-3061) to 254 (CC-2937) genes, while the *C.*
118 *incerta* GEVE encoded 355 genes. Most of the genes were shared among the *Imitevirales C.*
119 *reinhartii* GEVEs, consistent with their high average amino acid identity (AAI) to each other
120 (>98.5% in all cases, Dataset S1). These GEVEs also shared a high number of orthogroups
121 with the *C. incerta* GEVE (Dataset S1). In contrast, only a few orthogroups were shared
122 between the *Imitevirales* and the *Algavirales* GEVEs consistent with the large phylogenetic
123 distance between these lineages. Between ~44-55% of the genes in the *C. reinhardtii* and *C.*
124 *incerta* GEVEs have matches to Giant Virus Orthologous Groups (GVOGs), confirming their
125 viral provenance (Figure 1B). In addition, different genes in these regions have best matches to
126 giant viruses, bacteria, and eukaryotes, which is a common feature of *Nucleocyotiviricota*
127 members given the diverse phylogenetic origin of the genes in these viruses (Filée *et al.*, 2008)
128 (Figure 2A). Based on the Cluster of Orthologous Group (COG) annotations, a high proportion
129 of the GEVE genes are involved in transcription, and DNA replication and repair; however,
130 genes encoding translation, nucleotide metabolism and transport, signal transduction, and
131 posttranslational modification were also present, consistent with the diverse functional potential
132 encoded by numerous *Nucleocyotiviricota* (Figure 1C).

133 A previous study has shown that several field strains of *C. reinhardtii* harbor many genes that
134 are absent in the reference genome (Flowers *et al.*, 2015), which were possibly acquired from
135 diverse sources. To quantify the amount of novel genetic material contributed by giant viruses to
136 *C. reinhardtii*, we estimated the number of unique gene families in the analyzed *C. reinhardtii*
137 field strains that are absent in the reference strain CC-503. On average ~1.78% of the genes in
138 the field strains were unique compared to the reference strain (Figure 2B). Moreover, the
139 GEVE-harboring field strains have significantly enriched in novel genes compared to those
140 without GEVEs (Two-sided Man-Whitney U-test p-value <0.05, Figure 2B). These results
141 suggest that endogenization of giant viruses is an important contributor to inter-strain genomic
142 variability in *C. reinhardtii*. Recent studies have highlighted the importance of horizontal gene
143 transfer (HGT) in structuring the pangenome of diverse eukaryotes (Fan *et al.*, 2020; Sibbald *et*
144 *al.*, 2020), and genes originating from endogenous *Nucleocyotoviricota* were found to shape the
145 genomes of many algal lineages, including members of the Chlorophyta (Moniruzzaman *et al.*,
146 2020b; Nelson *et al.*, 2021). Compared to the GEVE-free strains, GEVE-containing strains
147 harbored a significantly higher proportion of genes from two COG categories including
148 Transcription, and Replication and Repair (Two-sided Mann-Whitney U test p-value <0.05)
149 (Figure 2B). All together, these GEVEs contributed many genes with known functions, including
150 glycosyltransferases, proteins involved in DNA repair, oxidative stress, and heat shock
151 regulation (Dataset S1).

152
153 A recent comparative genomic analysis of *C. reinhardtii* analyzed the population structure of this
154 alga by comparing numerous field strains (Craig *et al.*, 2019). Interestingly, we found
155 *Imitervirales* GEVEs in both North America populations 1 and 2 (NA1 and NA2, respectively),
156 and in both cases the GEVE-harboring strains are members of populations that include strains
157 for which GEVEs could not be detected. Indeed, strains CC-2931, CC-2932, and CC-3268 were
158 all isolated from the same garden in North Carolina, yet a GEVE could only be detected in CC-
159 3268. This patchwork distribution of the *Imitervirales* GEVEs within *C. reinhardtii* populations
160 suggests that they are the product of independent endogenization events rather than a single
161 event in their shared evolutionary history. Moreover, the *Imitervirales* GEVEs we identified here
162 fall within the same clade as most of the GEVEs we previously identified in other green algae.
163 The prevalence of GEVEs within a particular lineage, together with their patchwork distribution
164 across *C. reinhardtii* strains in the same population, suggests that GEVEs are the product of an
165 active endogenization mechanism that takes place over short timescales rather than
166 “accidental” endogenization that may result from illegitimate recombination that occurs during
167 infection.

168 169 **Discussion**

170
171 While much work remains to elucidate the role of GEVEs in shaping the ecological and
172 evolutionary dynamics of *C. reinhardtii*, several possibilities remain open. Some genes
173 contributed by the GEVEs could be potentially co-opted by the host, leading to changes in
174 certain phenotypes compared to closely related strains without GEVEs. Strain-specific
175 endogenization can also potentially lead to intraspecific variations in chromosome structure,
176 partly mediated by the GEVE-encoded mobile elements (Filée, 2018). Finally, it is also possible

177 that some of these GEVE-loci can produce siRNAs that might participate in antiviral defense,
178 and similar phenomena has been suggested for the virus-like loci in the genome of moss
179 (*Physcomitrella patens*) (Lang *et al.*, 2018). Recent studies on the large-scale endogenization of
180 giant viruses into diverse green algal genomes suggest that interactions between giant viruses
181 and their algal hosts frequently shape eukaryotic genome evolution (Moniruzzaman *et al.*) and
182 leads to the introduction of large quantities of novel genetic material. Our results indicate that
183 these endogenization events can lead to genomic variability not only between algal species, but
184 also between strains within the same population. Results reported in this study advance our
185 understanding of how giant viruses shape the genome evolution of their hosts, while also
186 expanding the scope of *C. reinhardtii* as a model organism to study the evolutionary fate and
187 consequences of giant virus endogenization.

188

189 **Methods:**

190

191 All methods and relevant citations are available in the ‘Supplementary Information’ file.

192

193 **Data and Code availability:**

194

195 Dataset S1 contains information regarding the raw data source, GEVE functional annotations,
196 hallmark gene distribution in each GEVE and coverage information of the partial GEVE in CC-
197 3061.

198

199 All the GEVE fasta files, unique gene fasta in each of the strains and their annotations, and
200 concatenated alignment file used to build the phylogenetic tree in Figure 1 are available in
201 Zenodo: <https://zenodo.org/record/4958215>

202

203 Code and instructions for ViralRecall v2.0 and NCLDV marker search scripts are available at:
204 github.com/faylward.

205

206 **Acknowledgements**

207

208 We acknowledge the use of the Virginia Tech Advanced Research Computing Center for
209 bioinformatic analyses performed in this study. This work was supported by grants from the
210 Institute for Critical Technology and Applied Science and the NSF (IIBR-1918271) and a Simons
211 Early Career Award in Marine Microbial Ecology and Evolution to F.O.A.

212

213 **Conflict of interest statement**

214

215 The authors declare no conflict of interest relevant to the content of the manuscript.

216

217 **Figure legends:**

218

219 **Figure 1: General features and phylogeny of the GEVEs. A)** Maximum likelihood
220 phylogenetic tree of the GEVEs and representative members from diverse *Nucleocytoviricota*

221 families constructed from a concatenated alignment of seven *Nucleocytoviricota* hallmark genes
222 (see Methods). Individual families within each order are indicated with abbreviations (IM -
223 Imitievirales, AG - Algavirales) followed by family numbers, as specified in Aylward et al, 2019
224 (Aylward *et al.*). IDs of the GEVEs are indicated in bold-italic. **B)** Basic statistics of the GEVEs
225 present in various field strains of *C. reinhardtii* and the GEVE present in the *C. incerta* genome.
226 **C)** Functional potential of GEVEs as EggNOG categories. Categories of genes are normalized
227 across all the NOG categories except S (function unknown) and R (general function prediction).
228 Raw functional annotations are in Dataset S1. NOG categories: [J] Translation, [F] Nucleotide
229 metabolism, [T] Signal Transduction, [M] Cell wall/membrane biogenesis, [A] RNA processing
230 and modification, [O] Post-translational modification, protein turnover and chaperone, [G]
231 Carbohydrate metabolism, [Q] Secondary structure, [Y] Nuclear structure, [U] Intracellular
232 trafficking and secretion, [Z] Cytoskeleton, [E] Amino acid metabolism, [N] Cell motility, [B]
233 Chromatin structure and dynamics, [H] Coenzyme metabolism, [V] Defense mechanism, [C]
234 Energy production and conversion, [P] Inorganic ion transport and metabolism, [I] Lipid
235 metabolism, [D] Cell cycle control, [L] Replication and repair, [K] Transcription.
236 * *C. incerta* GEVE length includes flanking eukaryotic regions.

237
238 **Figure 2: GEVE genomic and functional characteristics. A)** Circular plots of two
239 representative GEVEs in *C. reinhardtii* and the GEVE present in *C. incerta*. For *C. reinhardtii*
240 one representative *Imitevirales* GEVE (CC-2937) and the Algavirales GEVE (CC-2938) are
241 shown. Circle plots show Giant Virus Orthologous Group (GVOG) hidden Markov model (HMM)
242 hits, spliceosomal introns and the best LAST hit matches (see Supplementary Methods).
243 Internal blue links delineate the duplicated regions. The eukaryotic regions flanking the *C.*
244 *incerta* GEVE are delineated with light blue stripes. **B)** Unique genes in the field strains of *C.*
245 *reinhardtii* compared to the reference strain CC-503. The heatmap represents % of unique
246 genes that can be classified in different EggNOG categories (except category [R] - general
247 function prediction and [S] - function unknown). Categories marked with ‘***’ are significantly
248 overrepresented in the GEVE-containing strains compared to those without GEVEs. The bar
249 plot on top of the heatmap represents % of unique genes in each strain. GEVE-containing
250 strains have significantly higher percentages of unique genes compared to the strains without
251 GEVEs.

252
253 **References:**

254 **Aylward FO, Moniruzzaman M, Ha AD, Koonin EV.** A Phylogenomic Framework for Charting
255 the Diversity and Evolution of Giant Viruses. *PLOS Biology*, 2021.

256 **Craig RJ, Böndel KB, Arakawa K, Nakada T, Ito T, Bell G, Colegrave N, Keightley PD,**
257 **Ness RW. 2019.** Patterns of population structure and complex haplotype sharing among field
258 isolates of the green alga *Chlamydomonas reinhardtii*. *Molecular ecology* **28**: 3977–3993.

259 **Craig RJ, Hasan AR, Ness RW, Keightley PD. 2021.** Comparative genomics of
260 *Chlamydomonas*. *The Plant cell*.

261 **Endo H, Blanc-Mathieu R, Li Y, Salazar G, Henry N, Labadie K, de Vargas C, Sullivan MB,**
262 **Bowler C, Wincker P, et al. 2020.** Biogeography of marine giant viruses reveals their interplay

- 263 with eukaryotes and ecological functions. *Nature ecology & evolution* **4**: 1639–1649.
- 264 **Fan X, Qiu H, Han W, Wang Y, Xu D, Zhang X, Bhattacharya D, Ye N. 2020.** Phytoplankton
265 pangenome reveals extensive prokaryotic horizontal gene transfer of diverse functions. *Science*
266 *advances* **6**: eaba0111.
- 267 **Filée J. 2018.** Giant viruses and their mobile genetic elements: the molecular symbiosis
268 hypothesis. *Current opinion in virology* **33**: 81–88.
- 269 **Filée J, Pouget N, Chandler M. 2008.** Phylogenetic evidence for extensive lateral acquisition of
270 cellular genes by Nucleocytoplasmic large DNA viruses. *BMC evolutionary biology* **8**: 320.
- 271 **Flowers JM, Hazzouri KM, Pham GM, Rosas U, Bahmani T, Khraiweh B, Nelson DR,**
272 **Jijakli K, Abdrabu R, Harris EH, et al. 2015.** Whole-Genome Resequencing Reveals
273 Extensive Natural Variation in the Model Green Alga *Chlamydomonas reinhardtii*. *The Plant cell*
274 **27**: 2353–2369.
- 275 **Hasan AR, Duggal JK, Ness RW. 2019.** Consequences of recombination for the evolution of
276 the mating type locus in *Chlamydomonas reinhardtii*. *The New phytologist* **224**: 1339–1348.
- 277 **Lang D, Ullrich KK, Murat F, Fuchs J, Jenkins J, Haas FB, Piednoel M, Gundlach H, Van**
278 **Bel M, Meyberg R, et al. 2018.** The *Physcomitrella patens* chromosome-scale assembly
279 reveals moss genome structure and evolution. *The Plant journal: for cell and molecular biology*
280 **93**: 515–533.
- 281 **Meng L, Endo H, Blanc-Mathieu R, Chaffron S, Hernández-Velázquez R, Kaneko H, Ogata**
282 **H. 2021.** Quantitative Assessment of Nucleocytoplasmic Large DNA Virus and Host Interactions
283 Predicted by Co-occurrence Analyses. *mSphere* **6**.
- 284 **Merchant SS, Prochnik SE, Vallon O, Harris EH, Karpowicz SJ, Witman GB, Terry A,**
285 **Salamov A, Fritz-Laylin LK, Maréchal-Drouard L, et al. 2007.** The *Chlamydomonas* genome
286 reveals the evolution of key animal and plant functions. *Science* **318**: 245–250.
- 287 **Moniruzzaman M, Martinez-Gutierrez CA, Weinheimer AR, Aylward FO. 2020a.** Dynamic
288 genome evolution and complex virocell metabolism of globally-distributed giant viruses. *Nature*
289 *communications* **11**: 1710.
- 290 **Moniruzzaman M, Weinheimer AR, Martinez-Gutierrez CA, Aylward FO. 2020b.** Widespread
291 endogenization of giant viruses shapes genomes of green algae. *Nature*.
- 292 **Nelson DR, Hazzouri KM, Lauersen KJ, Jaiswal A, Chaiboonchoe A, Mystikou A, Fu W,**
293 **Daakour S, Dohai B, Alzahmi A, et al. 2021.** Large-scale genome sequencing reveals the
294 driving forces of viruses in microalgal evolution. *Cell host & microbe* **29**: 250–266.e8.
- 295 **Philippe N, Legendre M, Doutré G, Couté Y, Poirot O, Lescot M, Arslan D, Seltzer V,**
296 **Bertaux L, Bruley C, et al. 2013.** Pandoraviruses: amoeba viruses with genomes up to 2.5 Mb
297 reaching that of parasitic eukaryotes. *Science* **341**: 281–286.
- 298 **Salomé PA, Merchant SS. 2019.** A Series of Fortunate Events: Introducing *Chlamydomonas*
299 as a Reference Organism. *The Plant cell* **31**: 1682–1707.
- 300 **Sasso S, Stibor H, Mittag M, Grossman AR. 2018.** From molecular manipulation of
301 domesticated to survival in nature. *eLife* **7**.

302 **Schulz F, Roux S, Paez-Espino D, Jungbluth S, Walsh DA, Denev VJ, McMahon KD,**
303 **Konstantinidis KT, Eloe-Fadrosh EA, Kyrpides NC, et al. 2020.** Giant virus diversity and host
304 interactions through global metagenomics. *Nature* **578**: 432–436.

305 **Sibbald SJ, Eme L, Archibald JM, Roger AJ. 2020.** Lateral Gene Transfer Mechanisms and
306 Pan-genomes in Eukaryotes. *Trends in parasitology* **36**: 927–941.

307

308

309

310

311

312

313

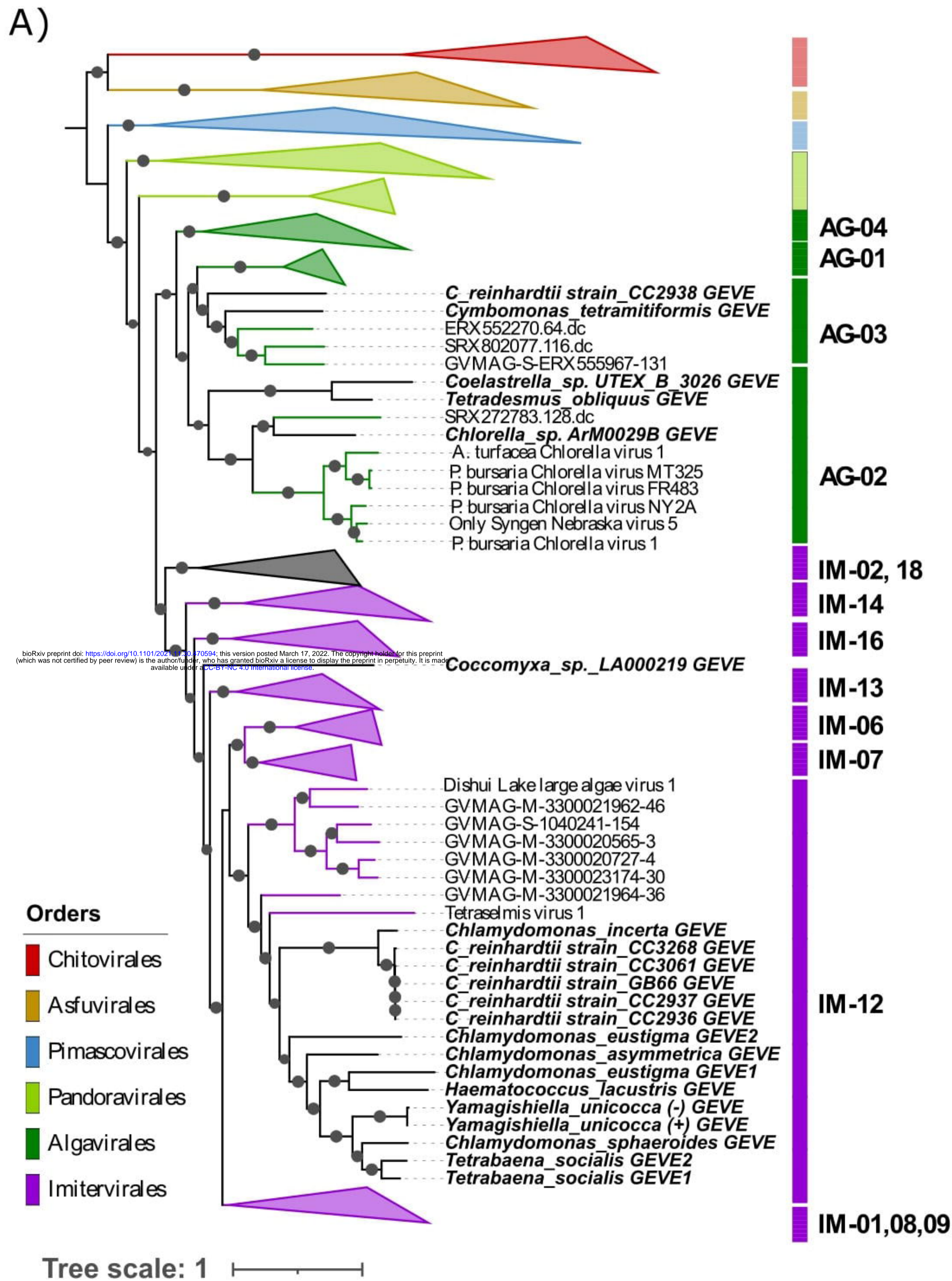
314

315

316

317

Figure 1



B)

Strain ID	GEVE Length (Kbp)	No. of Contigs	Protein count	Intron count	GC%	% of GVOG hits
CC-3061	112.7	8	99	25	60.51	56
CC-2938	315.2	11	214	47	58.27	44
GB-66	325.7	23	252	60	59.24	50
CC-3268	333.8	8	242	52	60.72	51
CC-2936	335.5	11	245	55	60.32	52
CC-2937	356.0	18	254	57	60.51	50
<i>C. incerta</i>	592.1*	1	355	72	64.78	48

