# Dynamic Genome Evolution and Blueprint of Complex Virocell Metabolism in Globally-Distributed Giant Viruses

Mohammad Moniruzzaman, Carolina A. Martinez-Gutierrez, Alaina R. Weinheimer, Frank O. Aylward*

Department of Biological Sciences, Virginia Tech, Blacksburg VA

*Email for correspondence: faylward@vt.edu

Key words: Giant viruses, NCLDV, Megavirales, virocell, viral diversity, marine viruses

**Abstract**

The discovery of giant viruses with large genomes has transformed our understanding of the limits of viral complexity in the biosphere, and subsequent research in model virus-host systems has advanced our knowledge of intricate mechanisms used by these viruses to take over host cells during infection. The extent of the metabolic diversity encoded by these viruses in the environment is less well-understood, however, and their potential impact on global biogeochemical cycles remains unclear. To address this, we generated 501 metagenome-assembled genomes (MAGs) of NCLDVs from diverse environments around the globe and analyzed their encoded functional diversity and potential for reprogramming host physiology. We found that 476 (95%) of the MAGs belonged to the *Mimiviridae* and *Phycodnaviridae* families, and of these we recovered 96% from aquatic environments, highlighting the diversity of these viral families in global freshwater and marine systems. MAGs encoded diverse genes predicted to be involved in nutrient uptake and processing, light harvesting, central nitrogen metabolism, and the manipulation of cell death, underscoring the complex interplay between these viruses and their hosts. Surprisingly, numerous genomes encoded genes involved in glycolysis, gluconeogenesis, and the TCA cycle, including one genome with a 70%-complete glycolytic pathway, suggesting that many of these viruses can even reprogram fundamental aspects of their host's central carbon metabolism. Phylogenetic trees of NCLDV metabolic genes together with their cellular homologs revealed distinct clustering of viral sequences into divergent clades, indicating these metabolic genes are virus-specific and were acquired in the distant past. Our findings reveal that diverse NCLDV genomes encode complex, cell-like metabolic capabilities with evolutionary histories that are largely independent of cellular life, strongly implicating them as distinct drivers of biogeochemical cycles in their own right.

**Abbreviations:** MAGs: Metagenome assembled genomes, NCLDV: Nucleo-cytoplasmic large DNA viruses

34 **Introduction**

35 Nucleocytoplasmic large DNA viruses (NCLDV) are a diverse group of eukaryotic viruses that
36 include several families of "giants" known for both their large virion size, reaching up to 1.5
37 μm, and genomes reaching ~2.5 million base-pairs in length (1). The discovery of the first giant
38 virus, *Acanthaomeba polyphaga mimivirus*, led to a paradigm shift in the field of virology by
39 showing that, contrary to the traditional view of viruses as "filterable infectious agents", viruses
40 could be larger than even some cellular lineages both in terms of physical size and genomic
41 contents (2, 3). Several subsequent studies have continued to expand our knowledge of
42 NCLDV diversity through the discovery of *Pithoviridae*, *Marseilleviridae*, *Pandoraviruses*, and
43 several other members of the *Mimiviridae* family, all encoding large genomes with diverse
44 genomic repertoires (4), and evolutionary genomic analysis has revealed common ancestry
45 of all these groups together with algal viruses of the family *Phycodnaviridae* and vertebrate
46 viruses of the families *Iridoviridae*, *Poxviridae*, and *Asfarviridae* (5).

47 Despite the genomic novelty of many NCLDVs, their diversity in the environment, mode of
48 genome evolution, and potential role in shaping ecological processes remains poorly
49 understood. Many pioneering discoveries of NCLDV over the last decade have leveraged
50 Amoebozoa as a model host for isolation (1), but it is likely that a variety of other unicellular
51 eukaryotes in the environment are infected by these viruses. Some *Phycodnaviridae* members
52 that infect algae have been studied for decades (6), and several algae-infecting Mimiviruses
53 have recently been isolated from diverse aquatic systems (7, 8). Moreover, recent cultivation-
54 independent analyses have provided tantalizing evidence suggesting that some NCLDV
55 groups are broadly distributed in nature and are potentially playing critical roles in the
56 ecological and evolutionary dynamics of unicellular eukaryotes, particularly in aquatic
57 environments (9, 10).

58 Understanding the phylogenetic and genomic diversity of NCLDV in the environment is
59 especially critical given recent findings on the biogeochemical significance of virus-mediated
60 metabolic reprogramming of host cells into "virocells" (11–13). The large and dynamic
61 genomes of NCLDVs have been shown to encode a variety of metabolic genes, including
62 those involved in nitrogen metabolism (14), fermentation (8, 14), and sphingolipid biosynthesis
63 (11), which likely contribute to shifts in host physiology during infection. Although these genes
64 are thought to be acquired by NCLDV from diverse sources through Lateral Gene Transfer
65 (LGT), the origin of many of these genes and the extent to which they are characteristic of
66 NCLDV genomes more broadly remains obscure. Given the emerging scientific consensus on
67 their ability to rewire the physiology of globally-abundant protists and impact marine
68 biogeochemistry (15), it is imperative to obtain a comprehensive view of the genome diversity,
69 evolutionary dynamics, and potential metabolic activities of these 'giants' of the virosphere.

## Results and Discussion

To address critical questions regarding the genomic diversity, evolutionary relationships, and virocell metabolism of NCLDVs in the environment, we developed a workflow to generate metagenome-assembled genomes (MAGs) of NCLDVs from publicly-available metagenomic data (see Methods, Figure S1). We surveyed 1,545 metagenomes and generated 501 novel NCLDV MAGs that ranged in size from 100-1,400 Kbp. Our workflow included steps to remove potential contamination from cellular organisms and bacteriophage and minimize possible strain heterogeneity in each MAG (see Methods). To ensure our NCLDV MAGs represented nearly-complete genomes, we only retained MAGs that contained at least 4 of 5 key NCLDV marker genes that are known to be highly conserved in these viruses (5) and had a total length > 100 Kbp (see Methods for details and rationale). Most of the MAGs were generated from marine and freshwater environments (444 and 36, respectively), but we also found 21 in metagenomes from bioreactors, wastewater treatment plants, oil fields, and soil samples (labeled "other" in Figure 1, S2; details in Dataset S1).

We constructed a multi-locus phylogenetic tree of the NCLDV MAGs together with 121 reference genomes using five highly conserved genes that have been used previously for phylogenetic analysis of these viruses (5) (Figure 1). The majority of our MAGs placed within the *Mimiviridae* and *Phycodnaviridae* families (350 and 126, respectively), but we also identified new genomes in the *Iridoviridae* (16), *Asfarviridae* (7), *Marseillviridae* (1), and *Pithoviridae* (1). Our phylogeny revealed that the *Phycodnaviridae* are polyphyletic and consist of at least two distinct monophyletic groups, one of which is sister to the *Mimiviridae* (Late *Phycodnaviridae*, 108 MAGs), and one which is basal branching to the *Mimiviridae*-Late *Phycodnaviridae* clade (Early *Phycodnaviridae*, 18 MAGs). In addition to the phylogeny, we evaluated the pairwise Average Amino Acid Identity (AAI) between NCLDV genomes to assess genomic divergence. AAI values provided results that were largely consistent with our phylogenetic analysis, with intra-family AAI values ranging from 26-100% (Figure 1b), highlighting the vast sequence divergence between even NCLDV genomes within the same family.

Given the large diversity within each of the NCLDV families, we sought to identify major clades within these groups that could be used for finer-grained classification. Using the rooted NCLDV phylogeny we calculated optimal clades within each family using the Dunn index (see Methods; (16)), resulting in 54 total clades, including 18 from the *Mimiviridae*, 13 for the Early *Phycodnaviridae*, and 6 for the Late *Phycodnaviridae* (Figure 1, Dataset S1). No cultured representatives were present in 31 of the clades (57%), including 2 from the *Asfarviridae*, 9 from the Early *Phycodnaviridae*, 1 in the *Iridoviridae*, 3 in the Late *Phycodnaviridae*, 1 in the *Marseilleviridae*, 14 in the *Mimiviridae*, and 1 in the *Pithoviridae*. Compared to references

106 available in GenBank, this increases the number of available genomes in the *Mimiviridae* over
107 eightfold (from 47 to 397) and the Phycodnaviridae by over threefold (from 46 to 121),
108 highlighting the vast diversity of environmental NCLDV that have not been sampled using
109 culture-based methods.

110 Analysis of the genome size distribution across the NCLDV phylogeny provided results that
111 are consistent with the current knowledge of these viruses. For example, Late
112 *Phycodnaviridae* clade 1 contained sequenced representatives of the Prasinoviruses,
113 including known *Ostreococcus* and *Micromonas* viruses, which encode the smallest
114 *Phycodnaviridae* genomes known (17). Consistent with this finding, the MAGs belonging to
115 this clade were also smaller and ranged in size from 100-225 Kbp, suggesting that small
116 genome size is broadly characteristic of this group. By comparison, genomes in the Early
117 *Phycodnaviridae* were larger and formed more divergent groups with long branches,
118 suggesting a large amount of untapped diversity in this clade (Figure 1a). The Pandoraviruses
119 and *Mollivirus sibericum*, notable for their particularly large genomes, formed a distinct clade
120 in the early *Phycodnaviridae*. The Coccolithoviruses and Phaeoviruses (18) were also placed
121 in the Early *Phycodnaviridae*, and we identified 7 and 2 new members of these groups,
122 respectively. Compared to the Late *Phycodnaviridae*, genome sizes of our MAGs were also
123 notably higher in the *Mimiviridae*, which are known to encode among the largest viral
124 genomes. In Mimivirus Clade 16, which includes *Acanthaeomeba polyphaga mimivirus*, we
125 identified 19 new MAGs, 13 of which have genomes > 500 Kbp. Taken together, these results
126 are consistent with the larger genomes that have been observed in the *Mimiviridae* compared
127 to the Late *Phycodnaviridae* (5). Although our NCLDV MAGs contain most marker genes we
128 would expect to find in these genomes, it is likely that many are not complete, and these
129 genome size estimates are therefore best interpreted as underestimates.

130 To assess the diversity of protein families across the NCLDV families, we calculated
131 orthologous groups (OGs) between our MAGs and 126 reference genomes, resulting in
132 81,411 OGs (Dataset S2). Of these, only 21,927 (27%) shared homology to known protein
133 families, highlighting the high level of novel genes in NCLDV that has been observed in other
134 studies. Moreover, 55,692 (68%) of the OGs were present in only one NCLDV genome
135 (singleton OGs), and overall the degree distribution of protein family membership revealed
136 only a small number of widely-shared protein families (Figure 2a,b), consistent with what has
137 been shown for dsDNA viruses in general (19). To visualize patterns of gene sharing across
138 the NCLDV we constructed a bipartite network in which both genomes and OGs can be
139 represented (Figure 2c). Analysis of this network revealed primarily family-level clustering,
140 with the *Mimiviridae* and early and late *Phycodnaviridae* clustering near each other, and the
141 *Pithoviridae, Marseiviridae*, and *Poxviridae* clustering separately. Interestingly, although

142  Pandoraviruses are members of the Early *Phycodnaviridae* clade, they clustered
143  independently in a small sub-network, indicating that the particularly large genomes and novel
144  genomic repertoires in this group are distinct from all other NCLDVs. These patterns suggest
145  that genomic content in the NCLDV is shaped in part by evolutionary history, but that large-
146  scale gains or losses of genomic content can occur over short evolutionary timescales, as has
147  occurred in the *Pandoraviruses*. Overall this is consistent with the "accordion-like" genome
148  evolution that has been postulated for NCLDV (20), whereby lineages evolve through a
149  balanced process of gene gain and loss over long evolutionary timescales. In many respects
150  this mode of genome evolution is not unlike that of *Bacteria* and *Archaea*, where genomic
151  repertoires are shaped by a mixture of vertical inheritance and HGT (21).

152  To further elucidate the evolutionary history of the large number of genes in NCLDVs, we
153  investigated clade-specific patterns in gene sharing. We found distinct clustering of NCLDV
154  OGs based on their presence in NCLDV clades, indicating that the majority of the OGs are
155  unevenly distributed across clades (Figure 3a). This was confirmed by an enrichment analysis,
156  where we identified sets of enriched OGs in each of the major NCLDV clades (Mann-Whitney
157  U test, corrected p-value < 0.01). The most common functional categories among the clade-
158  specific OGs are predicted to be involved in DNA replication, translation, and transcription.
159  Translational machinery was particularly enriched in Mimivirus clade 16, which contains many
160  cultivated representatives known to have the highest proportion of translation-associated
161  genes of any virus (22, 23). The clade-specific genomic repertoires of NCLDV suggest that
162  this is an appropriate phylogenetic scale for examining functional diversity across the NCLDV,
163  and we anticipate these clades will be useful groupings that can be used in future studies
164  examining spatiotemporal trends in viral diversity in the environment.

165  Relatively recent studies on model NCLDV-host systems have pointed out the presence of
166  genes involved in rewiring key aspects of cell physiology during infection, such as apoptosis,
167  nutrient processing and acquisition, and oxidative stress regulation (14, 24–26). We found a
168  number of genes involved in such processes to be broadly encoded across NCLDVs,
169  particularly in the *Mimiviridae* and *Phycodnaviridae* families (Figure 3b). Superoxide
170  dismutase (SOD) and Glutathione peroxidase (GPx), key players in regulating cellular
171  oxidative stress, are prevalent in phylogenetically divergent NCLDVs. Giant virus replication
172  possibly occurs under high oxidative stress inside the host cells (26) and thus, the presence
173  of enzymes with antioxidant activity might be crucial in preventing damage to the viral
174  machineries. SOD was biochemically characterized in *Megavirus chilensis*, and was
175  suggested to reduce the oxidative stress induced early in the infection (25). In addition, GPx
176  was found to be upregulated during infection by algal giant viruses (25, 26). Genes putatively
177  involved in the regulation of cellular apoptosis are also widespread in giant viruses, including

178    C14-family caspase-like proteins and several classes of apoptosis inhibitors, such as Bax1

179    (27). C14-family metacaspases were reported in a giant virus obtained through single virus

180    genomics approach, while viral activation and recruitment of cellular metacaspase was found

181    during *Emiliania huxleyi virus* (*EhV*) replication (24, 25). In *Chlorella* viruses, a K+ channel

182    (KcV) protein mediates host cell membrane depolarization, facilitating genome delivery within

183    the host (28). We identified KcV in genomes from all the major clades of late *Phycodnaviridae*

184    and *Mimiviridae*, suggesting that host membrane depolarization is a widely-adopted aspect of

185    NCLDV infection strategy. Lastly, in almost all the major *Mimiviridae* and *Phycodnaviridae*

186    clades we detected genes involved in DNA repair and processing, such as photolyases,

187    mismatch repair (*mutS)*, histones, and histone acetyl transferases, of which the latter two have

188    previously been reported in a number of giant virus families, with a possible role of viral

189    histones in packaging of DNA within the capsid (7, 29–31). All together, these results

190    demonstrate that many important aspects of viral reproduction and infection found in cultivated

191    NCLDV are widespread in nature and a common feature of virocell metabolism during giant

192    virus infection.

193    Viruses are thought to restructure host metabolism during infection to align with virion

194    production rather than cell growth, leading to altered nutrient demands inside the cell (12, 32).

195    We found that genes involved in nutrient acquisition and light-driven energy generation are

196    widespread in several NCLDV clades, including rhodopsins, chlorophyll a/b binding proteins,

197    ferritin, central nitrogen metabolism, and diverse nutrient transporters (Figure 3b). Recent

198    studies on the structure and mechanism of rhodopsin present in two giant viruses have

199    revealed that these are light-driven proton pumps, with potential to reshape energy transfer

200    within the infected host (33, 34). Similarly, widely-distributed chlorophyll a/b binding proteins

201    in giant viruses might increase photosynthetic light-harvesting capacity of infected cells, since

202    protists and plants are known to suppress their photosynthetic machineries, including the

203    chlorophyll binding antenna proteins (25) in response to virus infection (25, 35, 36).

204    Additionally, the presence of the key eukaryotic iron storage protein ferritin (37) and

205    transporters predicted to target ammonium (AmT), phosphorus (Phosphate permease and

206    Phosphate:Na+ symporters), sulfur (TauE/SafE family), and iron (Fe2+/Mn2+ transporters)

207    highlights the shifting nutrient demands of virocells compared to their uninfected counterparts.

208    Most of the MAGs were found in aquatic environments where nutrient availability may be

209    limiting for cellular growth, and alteration of nutrient acquisition strategies during infection may

210    be a key mechanism for increasing viral production. For example, although iron is crucial for

211    photosynthesis and myriad other cellular processes (38), it is often present in low

212    concentrations in marine environments (39, 40), and the production of viral ferritin may aid in

213    regulating the availability of this key micronutrient during virion production. Moreover, nitrogen

214  and phosphorus are limiting for microbial growth in many marine ecosystems, and given the
215  N:C and N:P ratios of viral biomass are relatively higher than that of cellular material (41), it is
216  likely crucial for viruses to boost acquisition of these nutrients with their own transporters.
217  Indeed, a recent study has revealed that an NCLDV-encoded ammonium transporter  (AmT)
218  can influence the nutrient flux in host cells by altering the dynamics of ammonium uptake (14).

219  Strikingly, many NCLDV genomes encode genes involved in central carbon metabolism,
220  including most of the enzymes for glycolysis, gluconeogenesis, the TCA cycle, and the
221  glyoxylate shunt (Figures 3b, 4a, S5). Central carbon metabolism is generally regarded as a
222  fundamental feature of cellular life, and so it is remarkable to consider that giant viruses
223  cumulatively encode nearly every step of these pathways. These genes were particularly
224  enriched in Mimivirus clades 1, 9, and 16, but a few of them were also present in several
225  *Phycodnaviridae* members (Figure 3b,4a). The glycolytic enzymes glyceraldehyde-3-
226  phosphate dehydrogenase (G3P), phosphoglycerate mutase (PGM), and phosphoglycerate
227  kinase (PGK) as well as the TCA cycle enzymes aconitase and succinate dehydrogenase
228  (SDH) were particularly prevalent. Additionally, we identified a fused gene in 16 MAGs that
229  encode the glycolytic enzymes G3P and PGK, which carry out adjacent steps in glycolysis
230  (Figure 4a, b), representing a unique domain architecture that has not been reported in cellular
231  lineages before. Interestingly, in many MAGs, TCA cycle genes were co-localized on viral
232  contigs, suggesting possible co-regulation of these genes during infection (Figure 4c).
233  Remarkably, one NCLDV MAG (ERX552257.96) encoded enzymes for 7 out of 10 steps of
234  glycolysis (Figure 4d), highlighting the high degree of metabolic independence that some giant
235  viruses can achieve from their hosts. The fact that viruses encode these diverse central
236  metabolic pathways underscores their ability to fundamentally reprogram virocell metabolism
237  through manipulation of intracellular carbon fluxes.

238  Phylogenies of a number of viral metabolic genes identified here together with their cellular
239  homologues revealed that NCLDV sequences tended to group together in deep-branching
240  clades, except for a few cases were multiple acquisitions from cellular sources was evident
241  (Figure 5, Figure S6). For example, aconitase, succinate dehydrogenase subunits B and C,
242  PhoH, glyceraldehyde-3-phosphate dehydrogenase, and superoxide dismutase all showed
243  distinct deep-branching viral clusters and were present in members of multiple NCLDV
244  families, suggesting they diverged from their cellular homologs in the distant past (Figure 5,
245  S6). This pattern was also observed for rhodopsin, similar to previous reports that NCLDV
246  rhodopsins represent a virus-specific clade (33), although our study suggests that at least
247  some NCLDVs independently acquired a bacterial rhodopsin. Phosphoglycerate kinase,
248  chlorophyll a-b binding proteins, and ammonium transporter (AmT) also appear to have been
249  acquired multiple times, but nonetheless show several deep-branching viral clades. These

250 results demonstrate that while NCLDVs have acquired numerous central metabolic genes
251 from cellular hosts, many of these metabolic genes have subsequently diversified into virus-
252 specific lineages. Indeed, detailed functional characterization of viral rhodopsin and Cu-Zn
253 superoxide dismutases has revealed that they have different structural and mechanistic
254 properties compared to the cellular homologs (33, 42), indicating that many metabolic genes
255 in giant viruses evolved to have specific functions in the context of host-virus interactions. Our
256 finding of the fused G3P-PGK glycolytic enzyme in many *Mimivirus* MAGs further reinforces
257 this view and demonstrates that NCLDV are unique drivers of evolutionary innovation in
258 metabolic genes. These results run contrary to a canonical view of viral evolution in which
259 viruses are seen as "pickpockets" that sporadically acquire genes from their cellular hosts
260 rather than encoding their own virus-specific metabolic machinery (43). Although these
261 metabolic enzymes were likely acquired from cellular lineages at some point, their distinct
262 evolutionary trajectory differentiates them from their cellular counterparts and demonstrates
263 that NCLDV are themselves a driver of evolutionary innovation in core metabolic pathways.

264 Viruses have historically been viewed as "accessories" to cellular life, and as such their
265 influence on biogeochemical cycles has largely been viewed through the lens of their impact
266 on host mortality, rather than any direct metabolic activities of their own. The large number of
267 cellular metabolic genes encoded in NCLDV genomes that we reveal in this study brings to
268 light an alternative view in which virus-specific enzymes have a direct role in shaping virocell
269 physiology. Scaled across viral infections in global aquatic environments, this raises the
270 possibility that viral enzymes can substantially alter global biogeochemical fluxes in their own
271 right. Moreover, the distinct evolutionary lineages of viral metabolic genes implicate NCLDV
272 as unique drivers of metabolic innovation, in stark contrast to the traditional view in which they
273 are merely occasional "pickpockets" of cellular genes rather than *de facto* evolutionary
274 innovators. Taken together, these findings argue that just as microbes are considered the
275 "engines that shape global biogeochemical cycles" (44), viruses must be considered alongside
276 their cellular counterparts as agents of metabolic fluxes with their own encoded physiology.

277

278 **Figure Legends**

279 **Figure 1.** A) Phylogeny of the 501 NCLDV MAGs presented in this study together with 121
280 reference genomes. The phylogeny was constructed from a concatenated alignment of 5
281 highly conserved marker genes that are present throughout the NCLDV families using the
282 VT+F+I+G4 model in IQ-TREE. The tree is rooted at *Poxviridae/Asfarviridae* branch,
283 consistent with previous studies (5). The inner strip is colored according to the phylogeny of
284 the MAGs, while the outer strip is colored according to the habitat in which they were found.

285  The bar chart represents genome size, which ranges from 100 - 2,474 Kbp, and the dotted
286  line denotes the 500 Kbp mark. Clades with >5 genomes are indicated with two letter
287  abbreviations and clade numbers. Abbreviations: MM: *Mimiviridae*, EP: Early
288  *Phycodnaviridae*, LP: Late *Phycodnaviridae*, IR: *Iridoviridae*, MR: *Marseilleviridae*, PT:
289  *Pithoviridae*. For the list of all the clades, see Dataset 1.  B) Average amino acid identity (AAI)
290  heatmap of the MAGs and reference genomes, with rows and columns clustered according to
291  the phylogeny.

292  **Figure 2.** A) The distribution of the orthologous groups (OGs) in the NCLDV MAGs and
293  reference genomes. The barplot on the left shows the proportion of OGs in each frequency
294  category that could be assigned an annotation, while the barplot on the right shows the total
295  number of OGs in each frequency category (log scale). B) The degree distribution of the OG
296  occurrence in the genomes analyzed. The best fit to a power law distribution is also shown.
297  C) A bipartite network of the OGs, with large nodes corresponding to genomes and small
298  nodes corresponding to OGs. The size of the genome nodes is proportional to their genome
299  size, and they are colored according to their family-level classification.

300  **Figure 3.** A) The barplot shows the number of enriched OGs in each of the major NCLDV
301  clades analyzed in this study. Only a subset of total functional categories are shown here; a
302  full table can be found in Dataset S2. B) A heatmap showing the occurrence of OGs with >5
303  total members across the major NCLDV clades, with shading corresponding to the percent of
304  MAGs in that clade that encode a given OG. C) A bubble plot of select metabolic genes
305  detected in the NCLDV clades, with bubble size proportional to the percent of genomes in a
306  clade that encode that protein. Abbreviations: G3P: glycerol-3-phosphate; LCM: Large
307  conductance mechanosensitive; SCM; small conductance mechanosensitive.

308  **Figure 4.** A) Presence of central carbon metabolism enzymes in the NCLDVs. The number of
309  genomes harboring a particular enzyme is provided beside its abbreviated name. Enzymes
310  that were not detected in any of the studied NCLDVs are in grey. B) Representative CDS from
311  genome ERX552243.92 illustrating the domain organization (PFAM and Interpro) of the fused-
312  domain gene (G3P + PGK) involved in glycolysis, that was detected in 16 of the NCLDV
313  MAGs. C) Example of co-localization of genes involved in TCA cycle on genomic contigs from
314  five representative NCLDV MAGs. Location of a number of other genes commonly present in
315  NCLDVs are also shown. D) Presence/absence of genes involved in central-carbon
316  metabolism in NCLDV genomes assembled in this study. Only the genomes harboring 3 or
317  more enzymes are shown. G3P+PGK indicates the fused-domain gene illustrated in panel B.
318  Blue arrow indicates the genome that harbors 7 out of 10 enzymes involved in glycolysis.
319  Abbreviations: **HK**: hexokinase, **PGI**: Phosphoglucoisomerase, **PFK**: Phosphofructokinase,
320  **ALD**: aldolase, **TPI:** Triose-phosphate isomerase, **G3P**: Glyceraldehyde 3-phosphate

321    dehydrogenase, **PGK**: Phosphoglycerate kinase, **PGM**: Phosphoglycerate mutase, **ENO**:

322    Enolase, **PYK**: Pyruvate kinase, **PEPCK**: PEP carboxykinase, **FBP**: Fructose 1,6-

323    bisphosphatase, **G6P**: Glucose 6-phosphatase, **PDH**: Pyruvate dehydrogenase, **PC**: Pyruvate

324    carboxylase, **CS**: Citrate synthase, **ACON**: Aconitase, **ICL**: Isocitrate lyase, **ICD**: Isocitrate

325    dehydrogenase, **αKDH**: α-ketoglutarate dehydrogenase, **SCS**: Succinyl-CoA synthetase, **SD**:

326    Succinate dehydrogenase (subunits A, B and C), **FH**: Fumarate hydratase, **MS**: Malate

327    synthase, **MDH**: Malate dehydrogenase.

328    **Figure 5.** Phylogenetic reconstruction of a number of representative NCLDV genes likely

329    involved in carbon and nutrient metabolism and light harvesting. NCLDV-specific clusters are

330    encircled with dashed ovals in each of the trees, while number of genes from different NCLDV-

331    clades contributing to these monophyletic groups are also provided (MM: *Mimiviridae*, EP:

332    Early *Phycodnaviridae*, LP: Late *Phycodnaviridae*) Colors of the clade names correspond to

333    those in Figure 1. Although node support values are not provided for better visual clarity, all

334    the NCLDV-specific nodes are supported by >90% ultrafast bootstrap values (see Methods

335    and Data availability statement for details).

336    * - Bacteriophage sequences are only present in the PhoH tree.

337    ** - Unclassified sequences (environmental) are only present in the Rhodopsin tree.

338

339

340    **Methods**

341    **Data Availability**

342    Nucleotide sequences of NCLDV MAGs, predicted proteins, alignments used for phylogenies,

343    raw orthologous groups membership files, and other major data products are available on the

344    Aylward Lab Figshare account: https://figshare.com/authors/Frank_Aylward/5008391 in the

345    project titled "NCLDV".

346    **Assembling NCLDV Genomes from Metagenomes.** Although phylogenetic binning of

347    metagenomic contigs belonging to Archaea and Bacteria is now commonplace (45), this

348    approach is rarely used for viruses, and it was therefore necessary for us to develop a novel

349    workflow to recover high-confidence NCLDV genomes from metagenomic data. Moreover,

350    methods for assessing the completeness and potential contamination of prokaryotic bins, such

351    as employed by the popular tool CheckM, rely on knowledge of shared single-copy protein

352    families in different lineages, but this information is not applicable to viral bins given their

353    fundamentally distinct genomic repertoires. We therefore also developed a workflow for

354    quality-checking NCLDV bins. The overall process can be divided into three main stages: 1)

355  initial binning of contigs, 2) identification of bins corresponding to NCLDV, and 3) quality-
356  checking bins to ensure contamination is not present. An overview of this workflow can be
357  found in Figure S1.

358   1) **Initial binning of contigs.** We obtained assembled contigs (> 10 Kbp) and coverage
359       files for 1545 metagenomes that had been previously assembled in a large-scale study
360       that examined bacterial and archaeal diversity (46). We chose to use the program
361       Metabat2 (47) for binning because this program bins contigs based on sequence
362       coverage and tetranucleotide frequencies, which are metrics that would be expected
363       to be consistent in viral genomes and therefore useful for binning. Moreover, although
364       some binning tools rely on marker gene sets that are specific to cellular lineages,
365       MetaBat2 does not use marker genes for binning and is therefore more appropriate for
366       binning viral sequences. We used the parameters -s 100000, -m 10000, --minS 75, --
367       maxEdges 75, which are more stringent than the default parameters and would be
368       expected to yield more conservative, high-confidence binning results.

369   2) **Screening Bins.** After binning contigs it is necessary to screen the bins to identify
370       ones correspond to putative NCLDV genomes. It has previously been shown that 5
371       highly conserved and NCLDV-specific protein families are present in almost all known
372       NCLDV genomes, and we therefore used these for screening. The protein families
373       correspond to the Major Capsid Protein (MCP), Superfamily II helicase (SFII), Virus-
374       like transcription factor (VLTF3), DNA Polymerase B (PolB), and packaging ATPase
375       (A32). We predicted protein sequences from all bins using Prodigal (48) (default
376       parameters), and matched proteins to the 5 NCLDV marker genes using HMMER3
377       (49) with custom Hidden Markov Models (HMMS, see section "protein families used
378       for screening" below). We only considered bins that had >= 4 of the markers for further
379       analysis in order to exclude NCLDV genomes that were mostly incomplete. Moreover,
380       we only considered bins > 100 Kbp on the grounds that the smallest NCLDV genome
381       is 103 Kbp (50), and that a higher cutoff may bias recovery against NCLDV clades with
382       smaller genome sizes. After implementing these screens we recovered 517 candidate
383       bins.

384   3) **Quality-checking Bins.** To ensure that the bins were indeed viral and did not include
385       contamination for cellular sources, we screened all contigs using ViralRecall
386       (https://github.com/faylward/viralrecall). This tool compares encoded proteins to virus-
387       specific (from the VOG database) and cellular-specific HMMs to assess their
388       provenance, and is therefore useful in determining contigs that may represent
389       contamination from a cellular organism. ViralRecall uses custom subsets of the Viral
390       Orthologous Groups (VOG: vogdb.org) and Pfam databases (51) for the virus-specific

391    and cellular-specific HMMs, respectively, and generates a score for each contig
392    (negative scores indicating more hits to cellular HMMs, positive scores indicating more
393    hits to viral HMMs). In addition to ViralRecall we also used LAST (parameter -m 500;
394    (51, 52)) to compare all of the encoded proteins in each contig to RefSeq 92 (53), and
395    recorded the top 5 hits for each protein. To remove contigs that derived from cellular
396    organisms, we removed contigs that had a ViralRecall score < 0 (indicting a net cellular
397    signal), contained < 3 encoded proteins with hits total to HMMs in the VOG database,
398    and had no LAST hits to known NCLDV proteins. Additionally, to exclude possible
399    bacteriophage sequences, we removed contigs for which the encoded proteins had at
400    least one LAST hit to a bacteriophage and zero hits to a known NCLDV genome. A
401    summary of all of the contigs removed in this way can be found in Dataset S1.

402    To uncover strain heterogeneity, we identified cases where marker genes were found
403    in multiple copies. We used SFII, VLTF3, A32, and PolB, for this, excluding MCP
404    because multiple copies of this gene is commonplace in NCLDV genomes. We
405    identified 16 bins where more than one marker gene was found to be present in
406    multiple copies, and excluded these bins from further analysis. Of the remaining 501
407    bins, 62 contained one marker gene that was not single-copy, but these were retained
408    because some complete NCLDV genomes contain multiple copies, and overly strict
409    thresholds would exclude potentially novel NCLDV lineages with genomic repertoires
410    distinct from what has been observed. Overall 501 bins passed all screening
411    procedures and are subsequently referred to as NCLDV Metagenome Assembled
412    Genomes (MAGs). Overall by using strict binning parameters and excluding possible
413    cellular or bacteriophage contigs, we recovered high quality MAGs that consisted of
414    relatively few contigs (422 bins with < 20 contigs, including 2 bins with only a single
415    contig; mean N50 contig size of 37.4 Kbp across MAGs).

416

417    **Protein families used for screening.** To screen preliminary bins and identify metagenome-
418    assembled NCLDV genomes, we used a custom set of HMMs created for 5 NCLDV-specific
419    protein families: The Major Capsid Protein (MCP), Superfamily II helicase (SFII), Virus-like
420    transcription factor (VLTF3), DNA Polymerase B (PolB), and packaging ATPase (A32). These
421    5 protein families have previously been used for phylogenetic analysis of NCLDV and are
422    typically not found in cellular organisms (5). To generate these models, we manually annotated
423    proteins from 126 complete NCLDV genomes available in NCBI that span the 7 major families.
424    We then generated model-specific HMMER3 score cutoffs based on the scores recovered
425    from matching known protein family members to these HMMs (Figure S4). These scores were
426    used in determining the presence/absence of these protein families in the metaBAT2 bins.

**Phylogenetic Reconstruction of NCLDV MAGs.** To assess the phylogeny of the NCLDV MAGs we generated a concatenated tree of all 501 MAGs together with 121 reference NCLDV genomes using the marker genes PolB, VLTF3, MCP, A32, and SFII. These proteins have previously been shown to be useful for phylogenetic analysis of NCLDV. In some cases, NCLDV are known to have introns or split genes, and we generated a Python script to identify these cases, check to ensure the proteins hit to the same HMM and had no sequence overlap, and subsequently concatenate the proteins. Alignments were created using ClustalOmega, and trimAl was used for trimming (parameter -gt 0.1). We ran IQ-TREE (54) with the "-m TEST" ModelFinder option (55), which identified VT+F+I+G4 as the optimal model. We then ran IQ-TREE on the alignment with 1000 ultrafast bootstraps to assess confidence (56).

**Clade Delineation**. Given the large phylogenetic diversity of NCLDV examined in this study, we sought to identify clades of closely related viruses within each of the major families. To this end we used the Dunn index to identify optimal clade-level delineations in our multi-locus phylogenetic tree of the NCLDV. We first generated a rooted ultrametric phylogenetic tree in R using the "ape" package and generated clades at different tree cut heights. For each cut height we calculated the Dunn index (16) using the "cluster.stats" package in R. We found that a height of 2.45 had the lowest Dunn index and therefore provided the best clustering (Figure S3). All clusters were then manually inspected and edited to ensure that they represented monophyletic groups, and these final clusters were used as clades.

**Generation of Orthologous Groups and Annotation.** We used ProteinOrtho v6.06 (57) to calculate the orthologous groups shared between the 501 NCLDV MAGs and 127 reference genomes. Protein files were generated using Prodigal with default parameters. Because of the accelerated evolutionary rate of viruses, we used the relaxed parameters "-e=1e-3 --identity=15 -p=blastp+ --selfblast --cov=40" for proteinortho. This resulted in 81,412 orthologous groups (OGs). For each OG, we randomly selected a representative for annotation. Representatives were compared to the EggNOG 4.5, TIGRFam, Pfam, VOG, and COG databases (e-value 1e-5), and best his were recorded.

**Bipartite Network Analysis.** Bipartite networks of NCLDV genomes and their protein families were created in igraph. For visualization purposes, only protein families present in > 5 genomes were analyzed. In the bipartite graph two node types were present: Genome nodes and Protein Family nodes. Each protein family node was connected to a genome node if it was encoded in that genome. The spring-directed layout was generated using the layout.fruchterman.reingold() command with 10,000 iterations.

**Average Amino Acid Identity Calculation.** AAI between the NCLDV genomes was calculated using a custom Python script available on GitHub

462   (https://github.com/faylward/pangenomics/blob/master/lastp_aai.py). The script employs
463   pairwise LASTP searches (52) (parameter -m 500) and calculates the AAI and alignment
464   fraction (AF) between all genome pairs. For visualization purposes, genome pairs in which
465   one member had an AF < 10 were considered to have an AAI of 0.

466   **Clade-specific OG enrichment.** To evaluate if specific NCLDV clades were enriched in
467   particular OGs, we performed an enrichment analysis on all clades with >5 members. For each
468   OG, a Mann-Whitney U test was performed on overall OG membership in that clade compared
469   to all other NCLDV (including reference genomes). Only OGs with a known annotation that
470   were present in >6 genomes were used. P-values were corrected using the Benjamini-
471   Hochberg procedure in R (58), and values < 0.01 were considered significant.

472   **Phylogenetic reconstruction of NCLDV metabolic genes:** We collected the reference
473   sequences for most of gene trees from the EggNOG database (59) with the following
474   exceptions: for phosphate permease and chlorophyll a/b binding proteins we used sequences
475   from the Pfam database (51), while rhodopsin sequences were collected from the MicRhoDE,
476   a dedicated server for rhodopsins from different domains of life and environments (60). In case
477   of the superoxide dismutase (SOD) and phoH genes, we curated additional sequences from
478   viruses other than NCLDVs from the NCBI Refseq database.  For each gene, a diagnostic
479   tree was built using FastTree (61) implemented in the ETE3 package. The diagnostic trees
480   were used to select a smaller set of reference sequences and remove short or redundant
481   sequences. For construction of the final trees, we aligned the sequences using ClustalOmega
482   and trimmed using trimAl (parameter -gt 0.1). IQ-TREE (54) was used to build maximum
483   likelihood phylogenetic trees with the model 'LG + I +G4' and 1000 ultrafast bootstrap
484   replicates (56).

485

### References

487   1.   N. Brandes, M. Linial, Giant Viruses – Big Surprises
488        https:/doi.org/10.20944/preprints201904.0172.v1.

489   2.   D. Raoult, P. Forterre, Redefining viruses: lessons from Mimivirus. *Nat. Rev. Microbiol.*
490        **6**, 315–319 (2008).

491   3.   B. La Scola, *et al.*, A giant virus in amoebae. *Science* **299**, 2033 (2003).

492   4.   S. Aherfi, P. Colson, B. La Scola, D. Raoult, Giant Viruses of Amoebas: An Update.
493        *Front. Microbiol.* **7**, 349 (2016).

494   5.   E. V. Koonin, N. Yutin, Multiple evolutionary origins of giant viruses. *F1000Res.* **7**
495        (2018).

496   6.   A. Jeanniard, *et al.*, Towards defining the chloroviruses: a genomic journey through a

497        genus of large DNA viruses. *BMC Genomics* **14**, 158 (2013).

498    7.    M. Moniruzzaman, *et al.*, Genome of brown tide virus (AaV), the little giant of the
499        Megaviridae, elucidates NCLDV genome expansion and host-virus coevolution. *Virology*
500        **466-467**, 60–70 (2014).

501    8.    C. R. Schvarcz, G. F. Steward, A giant virus infecting green algae encodes key
502        fermentation genes. *Virology* **518**, 423–433 (2018).

503    9.    Y. Li, *et al.*, Degenerate PCR Primers to Reveal the Diversity of Giant Viruses in Coastal
504        Waters. *Viruses* **10** (2018).

505   10.   M. Moniruzzaman, *et al.*, Diversity and dynamics of algal Megaviridae members during
506        a harmful brown tide caused by the pelagophyte, Aureococcus anophagefferens. *FEMS*
507        *Microbiol. Ecol.* **92**, fiw058 (2016).

508   11.   A. Vardi, *et al.*, Host-virus dynamics and subcellular controls of cell fate in a natural
509        coccolithophore population. *Proceedings of the National Academy of Sciences* **109**,
510        19327–19332 (2012).

511   12.   P. Forterre, The virocell concept and environmental microbiology. *ISME J.* **7**, 233–236
512        (2013).

513   13.   B. M. Schieler, *et al.*, Nitric oxide production and antioxidant function during viral
514        infection of the coccolithophore Emiliania huxleyi. *ISME J.* **13**, 1019–1031 (2019).

515   14.   A. Monier, *et al.*, Host-derived viral transporter protein for nitrogen uptake in infected
516        marine phytoplankton. *Proc. Natl. Acad. Sci. U. S. A.* **114**, E7489–E7498 (2017).

517   15.   A. E. Zimmerman, *et al.*, Metabolic and biogeochemical consequences of viral infection
518        in aquatic ecosystems. *Nature Reviews Microbiology* (2019)
519        https:/doi.org/10.1038/s41579-019-0270-x.

520   16.   J. C. Dunn†, Well-Separated Clusters and Optimal Fuzzy Partitions. *Journal of*
521        *Cybernetics* **4**, 95–104 (1974).

522   17.   K. D. Weynberg, M. J. Allen, W. H. Wilson, Marine Prasinoviruses and Their Tiny
523        Plankton Hosts: A Review. *Viruses* **9** (2017).

524   18.   W. H. Wilson, J. L. Van Etten, M. J. Allen, The Phycodnaviridae: the story of how tiny
525        giants rule the world. *Curr. Top. Microbiol. Immunol.* **328**, 1–42 (2009).

526   19.   J. Iranzo, M. Krupovic, E. V. Koonin, The Double-Stranded DNA Virosphere as a
527        Modular Hierarchical Network of Gene Sharing. *MBio* **7** (2016).

528   20.   J. Filée, Route of NCLDV evolution: the genomic accordion. *Curr. Opin. Virol.* **3**, 595–
529        599 (2013).

530   21.   V. Kunin, C. A. Ouzounis, The balance of driving forces during genome evolution in
531        prokaryotes. *Genome Res.* **13**, 1589–1594 (2003).

532   22.   F. Schulz, *et al.*, Giant viruses with an expanded complement of translation system
533        components. *Science* **356**, 82–85 (2017).

534   23.   J. S. Abrahão, R. Araújo, P. Colson, B. La Scola, The analysis of translation-related
535        gene set boosts debates around origin and evolution of mimiviruses. *PLoS Genet.* **13**,
536        e1006532 (2017).

537    24. K. D. Bidle, L. Haramaty, J. Barcelos E Ramos, P. Falkowski, Viral activation and
538         recruitment of metacaspases in the unicellular coccolithophore, Emiliania huxleyi. *Proc.*
539         *Natl. Acad. Sci. U. S. A.* **104**, 6049–6054 (2007).

540    25. M. Moniruzzaman, E. R. Gann, S. W. Wilhelm, Infection by a Giant Virus Induces
541         Widespread Physiological Reprogramming in Aureococcus Anophagefferens – A
542         Harmful Bloom Algae https:/doi.org/10.1101/256149.

543    26. U. Sheyn, S. Rosenwasser, S. Ben-Dor, Z. Porat, A. Vardi, Modulation of host ROS
544         metabolism is essential for viral infection of a bloom-forming coccolithophore in the
545         ocean. *ISME J.* **10**, 1742–1754 (2016).

546    27. K. S. Robinson, A. Clements, A. C. Williams, C. N. Berger, G. Frankel, Bax Inhibitor 1 in
547         apoptosis and disease. *Oncogene* **30**, 2391–2400 (2011).

548    28. M. Neupärtl, *et al.*, Chlorella viruses evoke a rapid release of K+ from host cells during
549         the early phase of infection. *Virology* **372**, 340–348 (2008).

550    29. M. G. Fischer, M. J. Allen, W. H. Wilson, C. A. Suttle, Giant virus with a remarkable
551         complement of genes infects marine zooplankton. *Proceedings of the National Academy*
552         *of Sciences* **107**, 19508–19513 (2010).

553    30. V. Thomas, *et al.*, Lausannevirus, a giant amoebal virus encoding histone doublets.
554         *Environ. Microbiol.* **13**, 1454–1466 (2011).

555    31. K. Okamoto, *et al.*, Cryo-EM structure of a Marseilleviridae virus particle reveals a large
556         internal microassembly. *Virology* **516**, 239–245 (2018).

557    32. S. Rosenwasser, C. Ziv, S. G. van Creveld, A. Vardi, Virocell Metabolism: Metabolic
558         Innovations During Host–Virus Interactions in the Ocean. *Trends in Microbiology* **24**,
559         821–832 (2016).

560    33. D. M. Needham, *et al.*, A distinct lineage of giant viruses brings a rhodopsin
561         photosystem to unicellular marine predators. *Proc. Natl. Acad. Sci. U. S. A.* **116**, 20574–
562         20583 (2019).

563    34. D. Bratanov, *et al.*, Unique structure and function of viral rhodopsins. *Nat. Commun.* **10**,
564         4939 (2019).

565    35. P. Juneau, J. E. Lawrence, C. A. Suttle, P. J. Harrison, Effects of viral infection on
566         photosynthetic processes in the bloom-forming alga Heterosigma akashiwo. *Aquatic*
567         *Microbial Ecology* **31**, 9–17 (2003).

568    36. J. Rahoutei, I. Garcia-Luque, M. Baron, Inhibition of photosynthesis by viral infection:
569         Effect on PSII structure and function. *Physiologia Plantarum* **110**, 286–292 (2000).

570    37. H. Botebol, *et al.*, Central role for ferritin in the day/night regulation of iron homeostasis
571         in marine phytoplankton. *Proc. Natl. Acad. Sci. U. S. A.* **112**, 14652–14657 (2015).

572    38. J. Morrissey, C. Bowler, Iron utilization in marine cyanobacteria and eukaryotic algae.
573         *Front. Microbiol.* **3**, 43 (2012).

574    39. S. L. Hogle, *et al.*, Pervasive iron limitation at subsurface chlorophyll maxima of the
575         California Current. *Proc. Natl. Acad. Sci. U. S. A.* **115**, 13300–13305 (2018).

576    40. M. J. Behrenfeld, Z. S. Kolber, Widespread iron limitation of phytoplankton in the south
577         pacific ocean. *Science* **283**, 840–843 (1999).

578  41.  L. F. Jover, T. C. Effler, A. Buchan, S. W. Wilhelm, J. S. Weitz, The elemental
579       composition of virus particles: implications for marine biogeochemical cycles. *Nat. Rev.*
580       *Microbiol.* **12**, 519–528 (2014).

581  42.  A. Lartigue, *et al.*, The megavirus chilensis Cu,Zn-superoxide dismutase: the first viral
582       structure of a typical cellular copper chaperone-independent hyperstable dimeric
583       enzyme. *J. Virol.* **89**, 824–832 (2015).

584  43.  D. Moreira, P. López-García, Ten reasons to exclude viruses from the tree of life. *Nat.*
585       *Rev. Microbiol.* **7**, 306–311 (2009).

586  44.  P. G. Falkowski, T. Fenchel, E. F. Delong, The microbial engines that drive Earth's
587       biogeochemical cycles. *Science* **320**, 1034–1039 (2008).

588  45.  N. Sangwan, F. Xia, J. A. Gilbert, Recovering complete and draft population genomes
589       from metagenome datasets. *Microbiome* **4**, 8 (2016).

590  46.  D. H. Parks, *et al.*, Recovery of nearly 8,000 metagenome-assembled genomes
591       substantially expands the tree of life. *Nat Microbiol* **2**, 1533–1542 (2017).

592  47.  D. D. Kang, *et al.*, MetaBAT 2: an adaptive binning algorithm for robust and efficient
593       genome reconstruction from metagenome assemblies. *PeerJ* **7**, e7359 (2019).

594  48.  D. Hyatt, *et al.*, Prodigal: prokaryotic gene recognition and translation initiation site
595       identification. *BMC Bioinformatics* **11**, 119 (2010).

596  49.  S. R. Eddy, Accelerated Profile HMM Searches. *PLoS Comput. Biol.* **7**, e1002195
597       (2011).

598  50.  C. A. Tidona, G. Darai, The complete DNA sequence of lymphocystis disease virus.
599       *Virology* **230**, 207–216 (1997).

600  51.  S. El-Gebali, *et al.*, The Pfam protein families database in 2019. *Nucleic Acids Res.* **47**,
601       D427–D432 (2019).

602  52.  S. M. Kiełbasa, R. Wan, K. Sato, P. Horton, M. C. Frith, Adaptive seeds tame genomic
603       sequence comparison. *Genome Res.* **21**, 487–493 (2011).

604  53.  N. A. O'Leary, *et al.*, Reference sequence (RefSeq) database at NCBI: current status,
605       taxonomic expansion, and functional annotation. *Nucleic Acids Res.* **44**, D733–45
606       (2016).

607  54.  L.-T. Nguyen, H. A. Schmidt, A. von Haeseler, B. Q. Minh, IQ-TREE: a fast and effective
608       stochastic algorithm for estimating maximum-likelihood phylogenies. *Mol. Biol. Evol.* **32**,
609       268–274 (2015).

610  55.  S. Kalyaanamoorthy, B. Q. Minh, T. K. F. Wong, A. von Haeseler, L. S. Jermiin,
611       ModelFinder: fast model selection for accurate phylogenetic estimates. *Nat. Methods*
612       **14**, 587–589 (2017).

613  56.  D. T. Hoang, O. Chernomor, A. von Haeseler, B. Q. Minh, L. S. Vinh, UFBoot2:
614       Improving the Ultrafast Bootstrap Approximation. *Mol. Biol. Evol.* **35**, 518–522 (2018).

615  57.  M. Lechner, *et al.*, Proteinortho: detection of (co-)orthologs in large-scale analysis. *BMC*
616       *Bioinformatics* **12**, 124 (2011).

617  58.  Y. Benjamini, Y. Hochberg, Controlling the False Discovery Rate: A Practical and

618  Powerful Approach to Multiple Testing. *Journal of the Royal Statistical Society: Series B*
619  *(Methodological)* **57**, 289–300 (1995).

620  59.  J. Huerta-Cepas, *et al.*, eggNOG 4.5: a hierarchical orthology framework with improved
621  functional annotations for eukaryotic, prokaryotic and viral sequences. *Nucleic Acids*
622  *Res.* **44**, D286–93 (2016).

623  60.  D. Boeuf, S. Audic, L. Brillet-Guéguen, C. Caron, C. Jeanthon, MicRhoDE: a curated
624  database for the analysis of microbial rhodopsin diversity and evolution. *Database* **2015**
625  (2015).

626  61.  M. N. Price, P. S. Dehal, A. P. Arkin, FastTree 2--approximately maximum-likelihood
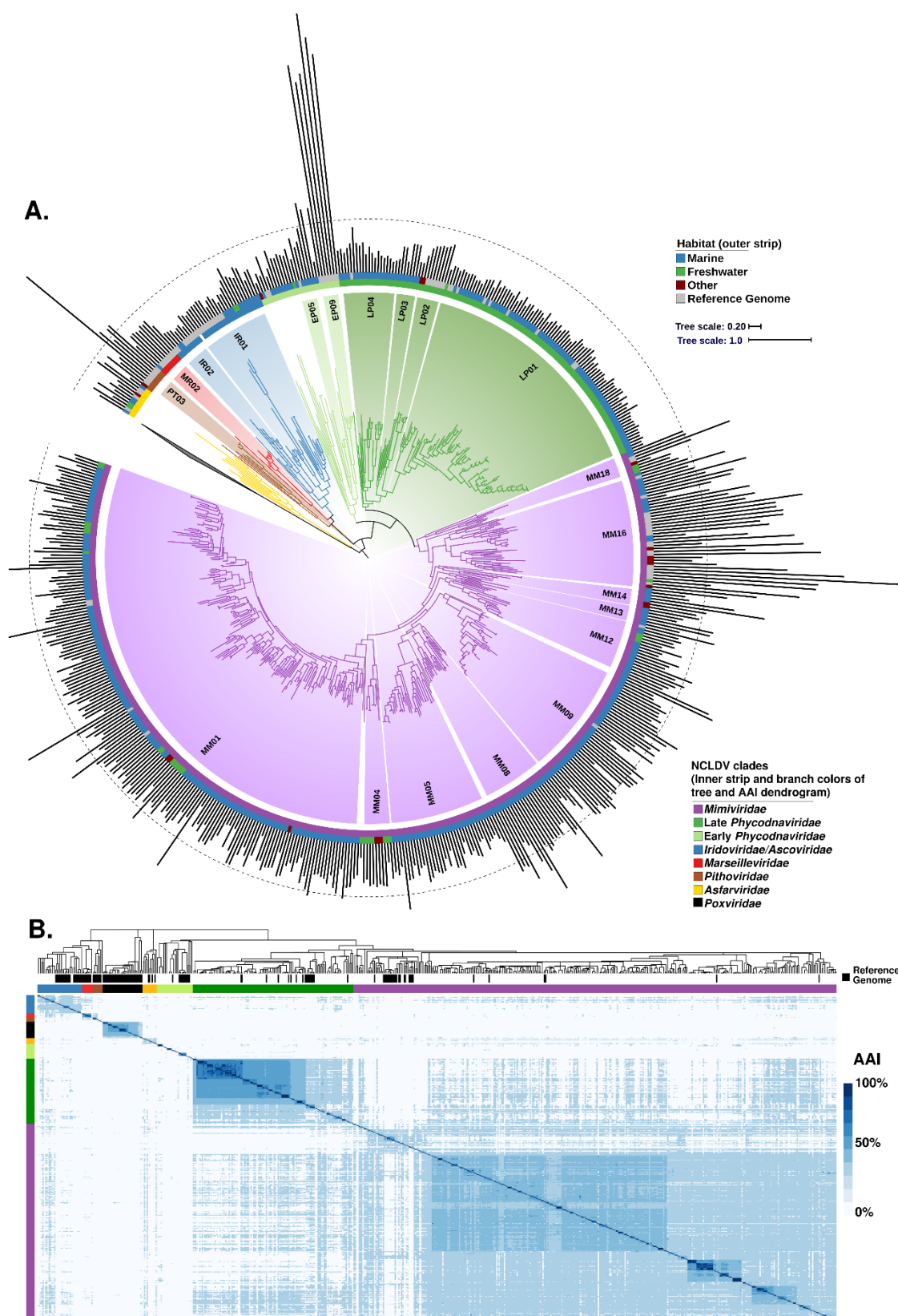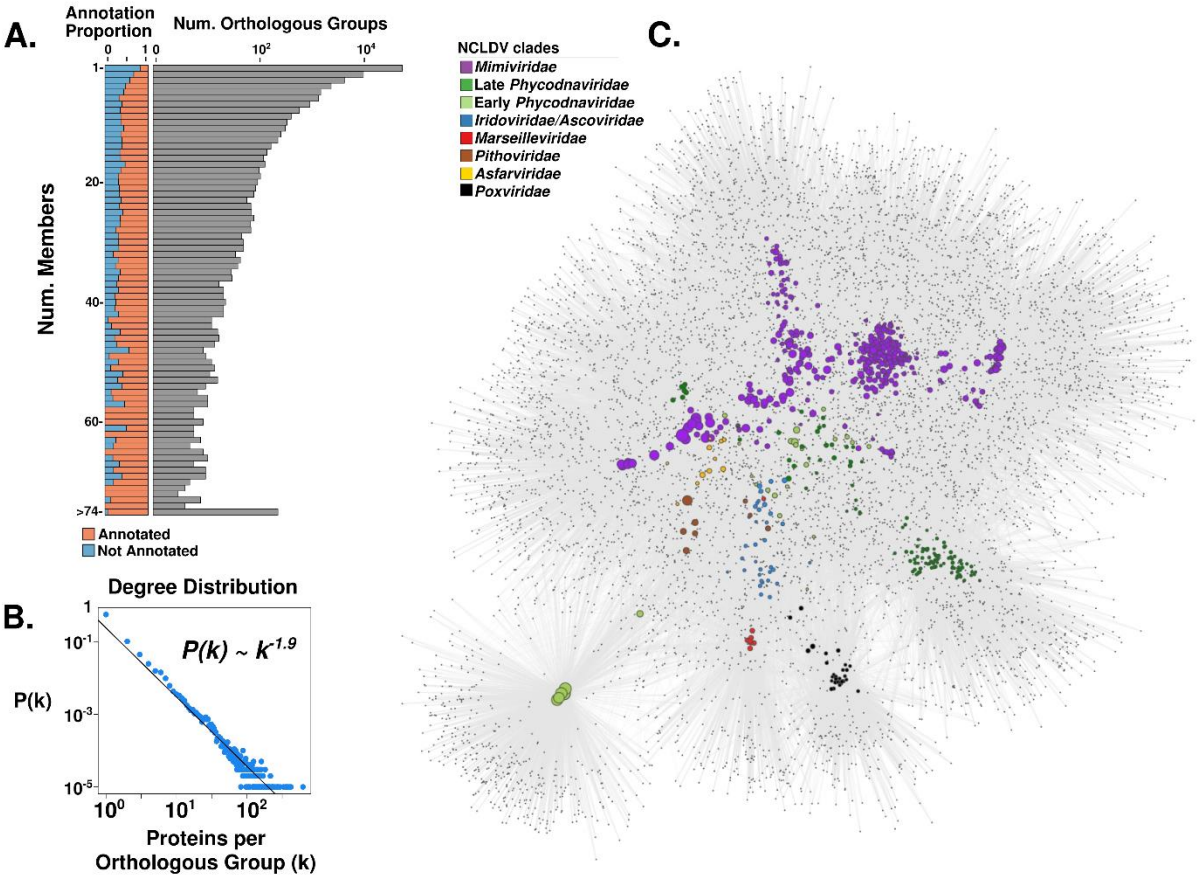627  trees for large alignments. *PLoS One* **5**, e9490 (2010).

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646    **Figure 1.**



647

648

649

650 **Figure 2.**



651

652

653

654

655

656

657

658

659

660

661

662

663

664    **Figure 3.**



665

666

667

668

669

670

671

672

673

674

675    **Figure 4.**



676

677

678

679

680

681

682

683

684

685

686 **Figure 5.**



687