For submission to: *Applied and Environmental Microbiology*

**TITLE: Diversity of active viral infections within the *Sphagnum* microbiome**

Authors:

**Joshua M.A. Stough[1*], Max Kolton[2], Joel E. Kostka[2], David J. Weston[3,4], Dale A. Pelletier[3], and Steven W. Wilhelm[1#]**

Addresses:

[1] Department of Microbiology, University of Tennessee, Knoxville, Tennessee, United States of America 37996

[2] School of Biology and School of Earth and Atmospheric Sciences, Georgia Institute of Technology, Atlanta, GA, USA

[3] Biosciences Division, Oak Ridge National Laboratory, Oak Ridge, Tennessee, USA

[4] Climate Change Science Institute, Oak Ridge National Laboratory, Oak Ridge, Tennessee, USA

#Address correspondence to Steven W. Wilhelm: wilhelm@utk.edu

*present address: Department of Microbiology & Immunology, University of Michigan, Ann Arbor 48109

Keywords:  Viruses, RNA-seq, Sphagnum, Peat bogs, microbial ecology

1

1  **Abstract**

2  *Sphagnum*-dominated peatlands play an important role in global carbon storage and represent

3  significant sources of economic and ecological value. While recent efforts to describe microbial

4  diversity and metabolic potential of the *Sphagnum* microbiome have demonstrated the

5  importance of its microbial community, little is known about the viral constituents. We used

6  metatranscriptomics to describe the diversity and activity of viruses infecting microbes within

7  the *Sphagnum* peat bog. The vegetative portions of 6 *Sphagnum* plants were obtained from a

8  peatland in northern Minnesota and total RNA extracted and sequenced. Metatranscriptomes

9  were assembled and contigs screened for the presence of conserved virus marker genes. Using

10  bacteriophage capsid protein, gp23, as a marker for phage diversity, we identified 33 contigs

11  representing undocumented phage s that were active in the community at the time of sampling.

12  Similarly, RNA-dependent RNA polymerase and the Nucleo-Cytoplasmic Large DNA Virus

13  (NCLDV) major capsid protein were used as markers for ssRNA viruses and NCLDV,

14  respectively. In total 114 contigs were identified as originating from undescribed ssRNA viruses,

15  22 of which represent near-complete genomes. An additional 64 contigs were identified as being

16  from NCLDVs. Finally, 7 contigs were identified as putative virophage or polinto-like viruses.

17  We developed co-occurrence networks with these markers in relation to the expression of

18  potential-host housekeeping gene *rpb1* to predict virus-host relationships, identifying 13 groups.

19  Together, our approach offers new tools for the identification of virus diversity and interactions

20  in understudied clades, and suggest viruses may play a considerable role in the ecology of the

21  *Sphagnum* microbiome.

22

2

23    **Significance**

24    *Sphagnum*-dominated peatlands play an important role in maintaining atmospheric carbon

25    dioxide levels by modifying conditions in the surrounding soil to favor its own growth over other

26    plant species. This slows rates of decomposition and facilitates the accumulation of fixed carbon

27    in the form of partially decomposed biomass. The unique environment produced by *Sphagnum*

28    enriches for the growth of a diverse microbial consortia that benefit from and support the moss's

29    growth, while also maintaining the hostile soil conditions. While a growing body of research has

30    begun to characterize the microbial groups that colonize *Sphagnum*, little is currently known

31    about the ecological factors that constrain community structure and define ecosystem function.

32    Top-down population control by viruses is almost completely undescribed. This study provides

33    insight into the significant viral influence on the *Sphagnum* microbiome, and identifying new

34    potential model systems to study virus-host interactions in the peatland ecosystem.

35

36

**Introduction**

Peatlands represent one of the most significant biological carbon sinks on the planet, storing an estimated 25% of terrestrial carbon in the form of partially decomposed organic matter (1-3). This accumulation of carbon is achieved through much slower rates of respiration and decomposition than observed in soil, due in large part to the low pH, nutrient-poor, and anaerobic environments created by the dominant moss population (4, 5), of which the genus *Sphagnum* is most prevalent (6, 7). As these environmental conditions appear to favor the growth of *Sphagnum* over vascular plants, primary production is dominated by the moss, which further retards decomposition due to production of antimicrobial compounds such as sphagnic acid (8-10) and sphagnan (11, 12). Despite this, *Sphagnum* and other peat mosses cultivate a diverse, symbiotic microbiome that appears to abate nutritional gaps for the moss and also contribute to the unique biogeochemical characteristics of the peatland ecosystem (13-15). In addition to their value as reservoirs of microbial diversity, the partially decomposed organic matter, known as *Sphagnum* peat, serves as an important economic resource for use in horticulture. As many peat bogs have begun to experience stress due to anthropogenic disturbances (16-18) and possibly climate change (19), the *Sphagnum* microbiome is of interest in peatland conservation and the ecosystem's services to the surrounding environment.

While there is a growing body of research characterizing the microbial groups that colonize *Sphagnum* (15), little is currently known about the ecological factors that define community structure and ecosystem function. Studies suggest that subtle differences in pH and available nutrients, manipulated by different *Sphagnum* species and strains, create distinct microbial consortia (14, 20, 21). Other observations suggest a more homogenous community (22), highlighting a need for further study. Culture-dependent experiments isolating endophytic

4

60    bacteria indicate *Sphagnum* cultivates symbionts with abilities that include antifungal activity

61    (20, 23) and nitrogen fixation (14), and that these microbiomes may be passed vertically to the

62    moss progeny (21). Yet while examinations of how environmental conditions and host-microbe

63    symbiotic interactions shape the structure and function of microbial communities, the influence

64    of virus populations on the *Sphagnum* microbiome remains unexplored.

65        Viruses are the most abundant biological entities on Earth, and central to global

66    ecosystems as they can drive the host evolution through predator-prey interactions and horizontal

67    gene transfer (24). Moreover, viruses can lyse single-celled primary producers and heterotrophs,

68    releasing nutrient elements from the biomass of prokaryotes and eukaryotic protists (25, 26).

69    Viruses may also act as a top-down control on the composition and evenness of microbial

70    communities, targeting hosts that reach higher cell densities, a phenomenon referred to as the

71    "*kill-the-winner*" model (27).

72        As lab studies of viruses require hosts that can be grown in culture, many

73    environmentally relevant viruses are poorly understood and their representation in reference

74    databases is often skewed. Previous efforts to describe environmental viromes have focused

75    primarily on the sequencing of shotgun or PCR-targeted metagenomes. While these methods

76    have proven powerful, rapidly expanding available reference material for bacteriophage (28, 29),

77    it leaves the considerable diversity of RNA viruses largely untapped (30). Moreover, the

78    common approach of selecting for viruses based on size-exclusion with filters removes many of

79    the Nucleo-Cytoplasmic Large DNA Viruses (NCLDVs, or commonly "giant viruses") that are

80    also environmentally relevant and phylogenetically informative (31, 32). Metagenomic

81    sequencing also limits observations to virus particles: from these data inferences on viral activity

82    require tenuous assumptions. The advent of high-throughput RNA sequencing offers viral

5

83   ecologists the opportunity to study active infections in the environment, as DNA viruses only

84   produce transcripts inside a host. Moreover, this approach also captures fragments of RNA virus

85   genomes. When sequencing is of sufficient depth and multiple samples are collected with spatial

86   and temporal variability, these data present an opportunity to develop hypothetical relationships

87   between virus and host markers (33) for subsequent in lab testing.

88       In this study, we analyzed metatranscriptomes from the microbial community inhabiting

89   the vegetative portion of *Sphagnum fallax* and *S. magellanicum* plants in Northern Minnesota,

90   with the goal of describing active viral infections within the *Sphagnum* microbiome. Using

91   marker genes conserved within several viral taxa, we identified an active and diverse

92   bacteriophage population, largely undescribed in previous studies. We also identified ongoing

93   infections by a diverse consortium of "giant" viruses and potentially corresponding

94   virophage/polinton-like viruses (hereafter referred to as virophage), including several giant

95   viruses closely related to the recently discovered Klosneuviruses (34). Finally, a number of novel

96   positive-sense single-stranded RNA viruses, some of which assembled into near complete

97   genomes, were observed. With this information in hand we developed statistical network

98   analyses, correlating co-expression of viral marker genes with housekeeping transcripts from

99   potential hosts. The resulting observations propose several virus-host pairings that, moving

100  forward, can be tested in a laboratory setting. Together, these results demonstrate new potential

101  model systems to study virus-host interactions in the peat bog ecosystem, and provide insight

102  into the significant viral influence on the *Sphagnum* microbiome.

103

104  **Results**

6

105    *Identification of resident phage populations*

106    To identify active virus populations in the *Sphagnum* phyllosphere, we obtained *S. fallax*

107    *and S. magellanicum* plant matter samples (3 from each species) from peatland terrariums as a

108    part of the Spruce and Peatlands Under Changing Environments (SPRUCE) project for

109    metatranscriptomic sequencing. Across all six *Sphagnum* phyllosphere samples, 33 contigs were

110    identified as transcripts encoding major capsid protein (*gp23*) originating from bacteriophage,

111    while only 6 contigs were identified using three other marker genes. Concurrent with this, more

112    reads mapped to *gp23* contigs than to the other marker genes combined, the most abundant of

113    which were three ribonucleotide reductase contigs.

114    Of the 33 contigs, 18 were assigned to the *Eucampyvirinae* subfamily with

115    *Campylobacter* viruses CP220 and PC18, while the rest were spread amongst the other Myovirus

116    taxa, predominantly the *Tevenvirinae* (Fig 1).  SS4 contig 77559 was the most abundant, with

117    consistently high expression across all samples, whereas other contigs dominated just one or two

118    samples. Of the 6 contigs identified using the other 3 viral marker genes, one was identified as a

119    potential *gp20* homologue, originating within *Myoviridae* with *Clostridium* virus phiCD119 as

120    the closest relative (SFig 1).  Two contigs were identified as *recA* contigs, likely originating in

121    myovirus and siphovirus relatives (SFig 2), while the remaining three contigs were identified as

122    ribonucleotide reductase transcripts (SFig 3).

123

124    *Single-stranded RNA virus diversity and abundance*

125    Within our samples, 114 contigs originated from RNA viruses, the majority of which

126    belonged to the currently unassigned *Barnaviridae* and Astrovirus-like families (Fig 2).

7

127    Additionally, a large number of *Picornaviruses* were observed, most of which were closely

128    related to the unclassified marine *Aurantiochytrium* single-stranded RNA virus, and *Secoviridae*

129    plant viruses. Lastly, several contigs were closely related to the *Nidovirales* clade, which

130    generally infect animal species.

131      Among these, 22 contigs were found to be near complete ssRNA virus genomes (based

132    on gene content and size), encoding multiple viral genes in addition to RDRP. Gene regions were

133    identified and annotated using the NCBI conserved domain and PFam HMM search tools, and

134    the full-length RDRP sequence was used to construct a maximum likelihood phylogenetic tree

135    (Fig 3). Of the partial ssRNA genomes that were assembled, 2 were missing the conserved Rhv

136    structural genes, while one was missing a RNA virus Helicase. The majority of these contigs fall

137    under the *Picornavirales* order, which also included the most complete viral genomes. As was

138    observed with the shorter RDRP contigs above, most of the *Picornavirales* contigs were most

139    closely related to the unclassified marine species, or members of the family *Secoviridae* clade,

140    whose membership includes the Parsnip yellow fleck virus. A number of partial *Picornavirus*

141    genomes were also identified as members of the family *Dicistroviridae*. Outside the

142    *Picornavirales*, most contigs clustered closely with the unassigned Astrovirus-like *Phytophthora*

143    *infestans* RNA virus. To determine the relative abundance of different RNA virus genomes in the

144    peat bog samples, we mapped reads back to contigs and calculated transcripts per million (TPM)

145    values to account for contig length and library size. The most abundant contig across all samples

146    was SS4 contig 3964, which was most closely related to the Rotifer birnavirus. All other contigs

147    appear to be abundant prominently in one or two samples, and absent or in low abundance in the

148    others, with no patterns of abundance apparent.

149    *Giant viruses and virophage in Sphagnum microbiome*

8

150        Of the 10 gene markers tested to identify Nucleo-Cytoplasmic Large DNA Viruses

151    (NCLDVs), only the giant virus major capsid protein (MCP) was detected in the

152    metatranscriptome. 64 contigs were observed with homology to MCP, representing every known

153    group of NCLDVs (Fig 4). Out of the 64 MCP contigs, 46 were placed within the *Mimiviridae*

154    taxa. Most contigs (25) closely aligned with the recently discovered Klosneuviruses, with the

155    Indivirus and Catovirus representing the most diversity in these samples. The next most abundant

156    group were the "extended *Mimiviridae*" (7 contigs), species with known similarity to

157    Mimiviruses but that infect eukaryotic algae. Six contigs phylogenetically were similar to the

158    *Asfarviridae,* here represented by the African swine fever Virus. Potential relatives of the giant

159    virus outliers, Pandoravirus and Pithovirus, were not observed (due to methodological

160    limitations), and the *Iridoviriae* were poorly represented (1 contig). Using the virophage MCP

161    and packaging ATPase as markers, we identified 7 contigs as transcripts originating in putative

162    virophage or polinton-like viruses, all of which were phylogenetically placed amongst isolates

163    identified from freshwater ecosystems (Fig 5).

164        As was observed with the other major viral taxa described, the majority of contigs were

165    most abundantly expressed in one or two samples and present at very low levels in the rest. The

166    most abundant NCLDV-MCP contig in the samples was SS2 contig 73240, most closely related

167    to *Megavirus chilensis*, which was the most highly expressed giant virus contig across all

168    samples. Four other contigs (SS6 contig 110585, SS4 contigs 55722 and 141177, and SS5 contig

169    119519) were highly expressed across all six samples.

170    *Prediction of virus-host pairs*

171        By comparing and correlating expression of virus marker genes to *rpb1* expression from

172    cellular organisms, we endeavored to predict potential virus-host groups in the *Sphagnum*

9

173    phyllosphere. Fig 6 shows statistically robust networks containing at least one virus and one host,

174    where co-occurrence and correlation were observed in more than one sample. A total of 13 virus-

175    host groups were detected, spread across the major viral taxa detected in this dataset. We note

176    that no networks containing the virophage/polinton-like viruses emerged. Four relationships

177    were predicted from bacteriophage *gp23* abundance, the simplest of which was a *Tevenvirinae*

178    phage-*Metazoa-Rhizaria* group with moderate correlations (Fig 6a). The other 3 relationships are

179    more complicated, containing multiple potential hosts and, for the largest predicted group,

180    multiple virus transcripts. The majority of potential hosts in these groups were identified as

181    eukaryotic, with only one putative bacterium and two archaea. Correlation coefficients for the

182    phage-prokaryote clusters were lower than was observed in the other major viral taxa, with low

183    to moderate correlations between viruses and bacteria.

184    We observed 4 predicted RNA virus-host clusters, all of which contained multiple hosts

185    grouped with a single virus (Fig 6b). Most of the predicted hosts appear closely related to

186    eukaryotic single-celled protists, within the *Excavata* and *Rhizaria* supergroups. Correlation

187    coefficients observed in these relationships are generally higher than observed in the phage-host

188    clusters. The 5 predicted NCLDV-host clusters (Fig 6c) were the most highly correlated and

189    complex. Predicted hosts were highly varied, ranging from diatoms to animals, though all virus

190    members were placed either within *Mimiviridae* or the extended Mimivirus group. MCP contigs

191    originating in close relatives of the recently discovered Klosneuviruses are present in both the 7-

192    and 10-member clusters, in addition to a pair of contigs most closely related to *Aureococcus*

193    *anophagefferens* Virus (AaV). An additional 15 statistically significant clusters across all three

194    viral taxa were observed where the virus and host were present in only one sample (not shown).

195

10

196    **Discussion**

197         Understanding the virus burden on microbial communities in ecologically-rich

198    ecosystems is an important step forward in resolving their function and predicting how they

199    might respond to various drivers of ecosystem scale change.  In the present study we used

200    metatranscriptomes to describe the diversity and activity of the resident virus populations in a

201    peat moss (*Sphagnum*) microbiome. We identified previously undescribed virus activity from

202    multiple taxa, most of which are poorly represented in either the literature or reference sequence

203    databases. We used read mapping to quantify the relative abundance of active viral infections.

204    Lastly, we compared expression of viral transcripts to that of potential hosts, using a correlation

205    co-occurrence networks approach (33) to predict putative hosts for the observed virus

206    populations. Together, our results suggest that the *Sphagnum* phyllosphere represents a

207    significant and largely untapped source of virus diversity and activity. Viruses were highly active

208    across all samples, with some individual viruses exhibiting abundant activity in single samples,

209    while others were more pervasive. Given that our observations were based on RNA sequencing

210    data, they do not represent a full accounting of the virus particles present in the community.

211    However, metatranscriptomic data, allows us to distinguish virus populations active at the time

212    of sampling. In addition, as viruses only transcribe their genes during infection, virus and host

213    transcripts are expected to co-occur, and it is possible that the abundance of transcripts (at least

214    for DNA viruses) could be used to predict natural hosts of viruses observed in the ecosystem

215    which can be tested in a laboratory or field setting. Ultimately, this study identifies from within a

216    complex community a number of candidate virus-host model systems for future study.

217    *Viral diversity and activity in Sphagnum plants*

11

218    As viruses lack a universal genetic marker like the bacterial 16S rRNA gene, we opted to

219    screen metatranscriptome assemblies for genes previously demonstrated to be largely or wholly

220    conserved amongst individual viral taxa. Within the expanded and diverse genetic potential of

221    giant viruses, only a handful of genes are currently conserved amongst all members (32, 35) and

222    these, in addition to several markers conserved amongst a large portion of giant viruses were

223    used to identify activity in the *Sphagnum* phyllosphere. Out of the 10 genes used to screen the

224    metatranscriptomes, we only MCP transcripts. This is not surprising given the number of capsid

225    proteins needed for viral assembly: indeed this transcriptional pattern was previously observed in

226    both cultures (36) and marine systems by Moniruzzaman *et al.* (2017). It should be noted that the

227    RNA-seq dataset used in those studies was poly-A selected, enriching for eukaryotic transcripts,

228    and thus coverage of eukaryotic virus gene expression would be much higher than in the

229    *Sphagnum* metatranscriptome. That we observed MCP expression in abundance suggests a

230    significant number of infections occurred at the time of sampling. While the magnitude of giant

231    virus diversity in *Sphagnum* dominated ecosystems is, to our knowledge, completely unexplored,

232    the richness observed here is considerably larger than expected compared to better documented

233    systems. 64 distinct MCP genotypes were identified in the *Sphagnum* phyllosphere

234    metatranscriptomes, which is high when compared to one recent survey that identified 30 novel

235    MCP transcripts from multiple environmental datasets (37), and another which observed 107

236    NCLDV sequences in 16 publicly available environmental metagenomes of comparable

237    sequencing depth isolated from different ecosystems (38). Most of the MCP contigs identified

238    were placed in clusters around a small number of virus relatives, highlighting the under-sampled

239    diversity of giant viruses in the literature, poor representation in reference databases, and the

240    considerable diversity present in *Sphagnum* peat bogs. The significant giant virus diversity

12

241  observed here implies a corresponding eukaryotic richness that is also under-described (39).

242  Additionally, a series of virophage transcripts were detected, indicating a significant response to

243  infections by giant viruses in the system. Many of these are phylogenetically grouped with the

244  polintoviruses, transposable elements that produce virion particles that can exploit the replication

245  machinery of actively infecting giant viruses to reproduce, often at the expense of the giant (40,

246  41).  These observations suggest that while an active picoeukaryotic population may persist,

247  mortality mechanisms beyond grazer-driven losses are at play and likely important to carbon

248  flow in the system.

249       The use of RNA-seq presents a unique opportunity to capture the genomic material of

250  RNA viruses that is lost in metagenomic sequencing. As such, RNA virus representation in

251  sequencing databases and the literature is largely constrained to culture-based studies. All known

252  RNA viruses require a functional RNA-dependent RNA polymerase (RdRP) to copy their

253  genome inside the host cell, a function exclusive to viruses, making it a highly specific marker

254  for RNA virus discovery (42, 43). Recent attempts to use metatranscriptomes to describe

255  environmental RNA viruses have proven successful, not only identifying marker gene fragments

256  in datasets, but assembling complete and near-complete genomes (33, 43). The diversity and

257  composition of RNA virus populations in *Sphagnum* peatlands is largely unknown: it is currently

258  limited to the small group of RNA-DNA hybrid chimeric Cruciviruses (44). Here, as was

259  observed with the giant viruses, most RNA virus contigs were placed in clusters with a single

260  represented species, suggesting a significant degree of uncharacterized diversity. This is not

261  entirely surprising, as RNA viruses are expected to make up as much as half of the virus particles

262  in the Earth's oceans, and yet they are almost as poorly understood and represented in

263  sequencing databases as giant viruses (30).  Similarly, we assembled and identified 22 near-

13

264    complete RNA virus genomes, where completeness was determined primarily by size and the

265    presence of the 6 core genes. As there are currently only 265 sequenced genomes within the

266    *Picornavirales*, most of which grouped within the *Picornaviridae*, this represents a sizeable

267    addition to the known diversity of ssRNA viruses. This is especially true for the unassigned and

268    unclassified taxa, and establishes a strong foundation for future efforts to describe RNA virus

269    populations in *Sphagnum*.

270         Description of bacteriophage populations in *Sphagnum* peatlands is currently limited to

271    the ssDNA viruses of the *Microviridae* (45) and *Caudovirales* (46) observed in metagenomics

272    data, though it appears that phage are the most abundant biological entities in the *Sphagnum*

273    phyllosphere (46). Given this, and the dominance of bacteria in the *Sphagnum* microbiome as

274    previously described (15), the relatively low abundance of active bacteriophage in our samples

275    was a surprise. Marker genes to identify bacteriophage were chosen based on their conservation

276    across phage taxa and their success in other environmental datasets. Gp20 (phage portal protein)

277    and Gp23 (major capsid protein) have been shown previously to be highly conserved and

278    effective for phylogenetic assignment of members of the *Myoviridae* (47-49). RecA  is conserved

279    across all three bacteriophage taxa and could illuminate lysogeny, and ribonucleotide reductase

280    (RNR) has been used as an effective marker for screening novel viruses from marine sequencing

281    datasets (50). As such, we identified 39 bacteriophage contigs using these markers, 33 of which

282    were from Gp23. This may represent a similar phenomenon as MCP in the giant viruses above,

283    where transcripts encoding structural proteins are much more abundant than other genes and

284    sequencing lacked the depth to detect them. For the purpose of discovering novel phage species,

285    DNA sequencing through metagenomics may prove more successful.

286    *Virus-host predictions*

14

287     Future study of viral dynamics in peatlands will require the establishment of model

288     virus/host pairs for *in vitro* experimentation and *in situ* tracking. While culture-based techniques

289     can yield model systems, it is not always clear whether the isolated organisms are

290     environmentally relevant. In order to address this, we attempted to use statistical methods to

291     propose virus/host pairs as potential future model systems based on their cooccurrence in

292     samples and the correlation of their abundance. As viruses produce transcripts only when

293     actively infecting a host, positive correlation and co-occurrence between virus and host

294     transcripts is expected, and might be used to predict host-virus relationships, provided an

295     appropriate transcriptional proxy for growth and activity is available (33). In this study, we used

296     the eukaryotic RNA-polymerase gene *rpb1* as a marker for abundance and activity in potential

297     hosts, as it is conserved amongst all eukaryotic organisms, is phylogenetically informative, and

298     has been previously described as one of the more consistently expressed eukaryotic genes in

299     marine systems, scaling well with the activity of the organism (51), though the stability of its

300     expression has not been evaluated in terrestrial ecosystems. We used NCLDV MCP abundance

301     as a proxy for giant virus production, Gp23 for phage production, as transcription is necessary

302     for the assembly of new virus particles and transcript abundance in some appears to be closely

303     linked to viral replication. We also used RdRP as a proxy for RNA virus production,

304     acknowledging the caveat that we cannot distinguish between abundance of free virus particles

305     and active infections (33).

306     Correlation and co-occurrence matrices, clustered into groups by similarity and tested

307     with the SIMPROF permutation test, yielded 13 predicted groups of viruses and hosts. For

308     ssRNA and giant viruses, several of the networks produced in the analysis included multiple

309     bacterial and archaeal sequences picked up in the RNA polymerase screen. As we have no reason

15

310    to believe bacterial species are infected by NCLDVs or Picornaviruses, it is likely these

311    predictions represent a confounding relationship between prokaryotes and potential eukaryotic

312    hosts, observed in network analyses for all three viral taxa described here, where a beneficial

313    interaction results in an indirect correlation with viral infection. Indeed, previous use of this

314    method in marine systems showed a similar phenomenon, where an algal Mimivirus and a

315    known host were grouped with a fungal species and another virus (33). Even after the

316    consideration of bacterial species within the predicted groups, some remain complicated with

317    multiple viruses and potential eukaryotic hosts, which may be explained by a broader host range

318    amongst giant viruses enabled by the expansion of genetic material and increased independence

319    from host machinery. Similar relationships were observed amongst RNA viruses, though these

320    are more tenuous, as we are unable to distinguish whether sequencing reads originated transcripts

321    or genomic material.

322        All together, we have identified a considerable amount of viral diversity from several

323    major viral taxa active within a poorly understood microbial ecosystem. As they were identified

324    from transcript sequencing data, the viruses described here likely only represent a fraction of the

325    whole virus community, which may be elucidated through further culture-independent work. We

326    have also used transcript abundance within a statistical framework to predict several host-virus

327    relationships which can be sought out and tested in culture. These results establish an important

328    and much needed foundation for future research into the microbial ecology in *Sphagnum* peat

329    bogs.

330

331    **Materials and Methods**

16

332    *Sample collection and Survey of Environmental Conditions*

333    Triplicate individual plants of *Sphagnum magellanicum* and *Sphagnum fallax* were

334    collected on August 2015 from the SPRUCE experiment site at the S1 bog on the Marcell

335    Experimental Forest (U.S. Forest Service, http://mnspruce.ornl.gov/). The S1 Bog is an acidic

336    and nutrient-deficient ombrotrophic *Sphagnum*-dominated peatland bog (surface pH≤4.0) located

337    approximately 40 km north of Grand Rapids, Minnesota, USA (47°30.476′ N; 93°27.162′ W; 418

338    m above mean sea level) (52-54). To characterize the *Sphagnum* virome, *Sphagnum* samples

339    were collected as previously described (54). Only green living plants were sampled: samples

340    focused on the capitulum plus about 2-3 cm of green living stem. B *Sphagnum* stems

341    (phyllosphere) were cleaned from unrelated plant debris, and frozen immediately on dry ice.

342    Frozen samples were overnight shipped to the Georgia Institute of Technology for RNA

343    extraction.

344

345    *RNA Extraction and Sequencing*

346    One gram of *Sphagnum* phyllosphere tissue was ground with a mortar and pestle under liquid

347    nitrogen. The fine powder was transferred to 10 extraction tubes and total RNA isolated using the

348    PowerPlant RNA Isolation Kit with DNase according to the manufacturer's protocol (MoBio

349    Laboratories, Carlsbad, CA, USA). DNA-depleted RNA was quantified using the Qubit RNA HS

350    Assay Kit (Invitrogen, Carlsbad, CA, USA) and quality was assessed on the Agilent 2100

351    BioAnalyzer using the Agilent RNA 6000 Pico Kit (Agilent Technologies). Additionally, the

352    absence of DNA contamination was confirmed by running a polymerase chain reaction using

353    universal bacterial 16S rRNA primers 515F and 806R. Finally, RNA samples without detectible

17

354    DNA contamination and exhibiting an RNA integrity number (RIN) > 6 were pooled. Extracted

355    total environmental RNA samples were was sent on dry ice to the Joint Genome Institute (JGI)

356    facilities for meta-transcriptomes libraries construction and sequencing. All protocols employed

357    were standard JGI protocols Ribosomal RNA subtraction from total environmental RNA was

358    completed using the Ribo-Zero rRNA Removal Kit (Illumina, San Diego, CA). rRNA depleted

359    environmental RNA were used to construct paired end metatranscriptomes libraries using TruSeq

360    kit and sequenced on the Illumina HiSeq2000 platform at the JGI facilities using a single-end

361    250bp flow cell.

362    *RNA-seq Data Processing*

363        Raw sequences (see Supplementary Table 2) were downloaded from the Department of

364    Energy Joint Genome Institute server and processed using the CLC Genomics Workbench v.

365    10.0.1 (QIAGEN, Hilden, Germany). Reads below a 0.03 quality score cutoff were removed

366    from subsequent analyses, and the remaining reads were trimmed of any ambiguous and low

367    quality 5' bases. Samples were subjected to a subsequent *in silico* rRNA reduction using the

368    SortmeRNA 2.0 software package (55). Filtered reads were *de novo* assembled with cutoffs of

369    300 base minimum contig length and average coverage of 2, leaving a total of 705,526 contigs

370    across all samples.

371    *Screening Assemblies for Marker Genes*

372        Marker genes to identify bacteriophage were chosen based on their conservation across

373    phage taxa and their success in other environmental datasets. Gp20 (phage portal protein) and

374    Gp23 (major capsid protein) have been shown previously to be highly conserved and effective

375    for phylogenetic assignment of members of the *Myoviridae* (47-49). RecA  is conserved across

18

376    all three bacteriophage taxa and could illuminate lysogeny, and ribonucleotide reductase (RNR)

377    has been used as an effective marker for screening novel viruses from marine sequencing

378    datasets (50). To identify contigs specific to the NucleoCytoplasmic Large DNA Virus

379    (NCLDV) clade, contig libraries were screened for the presence of 10 genes previously identified

380    as core NCLDV genes as previously described (33). Briefly, contig libraries were queried against

381    Nucleo-Cytoplasmic Virus Orthologous Groups (NCVOG) protein databases for each of the

382    following 10 marker genes in a Blastx search with a minimum e-value cutoff of $10^{-3}$: A32 virion

383    packaging ATPase (NCVOG0249), VLFT-like transcription factor (NCVOG0262), Superfamily

384    II Helicase II (NCVOG0024), mRNA capping enzyme (NCVOG1117), D5 helicase-primase

385    (NCVOG0023), ribonucleotide reductase small subunit (NCVOG0276), RNA polymerase large

386    subunit (NCVOG0271), RNA polymerase small subunit (NCVOG0274), B-family DNA

387    polymerase (NCVOG0038), and major capsid protein (NCVOG0022). Resulting hits were then

388    queried against the NCBI refseq protein database (56) and only contigs with top hits to virus

389    genes were maintained for subsequent analyses. A similar method was used to identify virophage

390    transcripts, where the virophage major capsid protein and packaging ATPase genes were used as

391    markers.

392         Contigs derived from ssRNA viruses were identified by screening the contig library for

393    RNA-dependent RNA Polymerase (RDRP), a distinctive and wholly conserved RNA virus gene

394    and a strong phylogenetic marker (57). A BLAST database of RDRP sequences was downloaded

395    from the pfam database (58) under code pf00680. Contigs were aligned using Blastx with a

396    minimum evalue of $10^{-4}$. Hits were queried against the NCBI refseq protein database and only

397    hits to viral RDRP genes were retained for downstream analyses.

19

398    To identify RNA virus genome fragments, contig libraries were screened as described

399    above using the following core set of genes observed in RNA viruses: CRPV capsid (Pfam

400    08762), VP4 (Pfam 11492), RdRP (Pfam 00680), Peptidase C3 (Pfam 00548), Peptidase C3G

401    (Pfam 12381), Rhv (Pfam 00073), and RNA Helicase (Pfam 00910). BLAST databases for core

402    RNA virus genes were constructed from reference sequences downloaded from pfam. Query

403    sequences were then cross-referenced to identify contigs with hits to multiple RNA virus core

404    genes. Only contigs > 1000 bases with at least one viral RDRP region were retained for further

405    analysis. ORFs were predicted on these putative partial genomes using the CLC Genomics

406    Workbench. Features on the partial genomes were predicted using the Pfam HMM domain and

407    the NCBI Conserved Domain Database searches (59, 60). Genome architecture was visualized

408    using the Illustrator for Biological Sequences (IBS) software package (61).

409    *Phylogenetic Analysis*

410    Reference sequences for viral marker genes and Rpb1 were downloaded from the

411    InterPro and RefSeq databases (STable 1) (62). Reference sequences were aligned using the

412    MUSCLE alignment algorithm (63) in the MEGA v7.0.26 software package (64). Maximum

413    likelihood phylogenetic trees were constructed in PhyML (65) with the LG substitution model

414    and the aLRT SH-like likelihood method. Putative viral and Rpb1 contigs assembled from the

415    metatranscriptomes were translated into proteins according to the reading frame of the top

416    BLAST hit. Translated proteins were placed on the reference trees in a maximum likelihood

417    framework in pplacer (66). Trees with abundance data were visualized using the iToL web

418    interface (67).

419    *Statistical Analysis*

20

420    Quality filtered and trimmed reads were stringently mapped to the selected contigs (0.97

421    identity fraction, 0.7 length fraction) in CLC Genomics Workbench 10.0.1. Expression values

422    were calculated as a modification of the transcript per million (TPM) metric. Read counts were

423    normalized by contig length in kb to determine the reads per kilobase (RPK) values for every

424    contig within each library. These RPK values were then summed and divided by 1 million, to

425    determine the sequencing depth scaling factor for each library. TPM for a contig was calculated

426    by dividing its RPK value by the scaling factor for the library.

427    Expression values for contigs were imported into the PRIMER7 (68) statistical software

428    package and $\log_2$ transformed. Expression values from each contig were correlated (Pearson's

429    rho) to one another and statistically grouped by co-occurrence using group average hierarchical

430    clustering. The SIMPROF test (69) was used to determine the statistical significance level of

431    resulting clusters (alpha = 0.05, 1000 permutations). Statistically significant clusters with at least

432    one viral contig, one *rpb*1 contig and less than 10 total members were visualized and annotated

433    in Cytoscape 3.5.1 (70).

434    *Accession Numbers*

435    Full RNA-seq libraries have been made publicly available on the JGI website under

436    accession number Gp0146911.

437

21

438    *Acknowledgements*

22

446 **References**

447

448 1. **Post WM, Emanuel WR, Zinke PJ, Stangenberger AG.** 1982. Soil carbon pools and
449 world life zones. Nature **298:**156-159.
450 2. **Gorham E.** 1991. Northern peatlands: role in the carbon cycle and probable responses to
451 climatic warming. Ecological applications **1:**182-195.
452 3. **Bridgham SD, Patrick Megonigal J, Keller JK, Bliss NB, Trettin C.** 2006. The carbon
453 balance of North American wetlands. Wetlands **26:**889-916.
454 4. **van Breemen N.** 1995. How Sphagnum bogs down other plants. Trends in Ecology &
455 Evolution **10:**270-275.
456 5. **Lamers LPM, Bobbink R, Roelofs JGM.** 2000. Natural nitrogen filter fails in polluted
457 raised bogs. Global Change Biology **6:**583-586.
458 6. **Turetsky MR.** 2003. The role of bryophytes in carbon and nitrogen cycling. Bryologist
459 **106:**395-409.
460 7. **Turetsky MR, Bond-Lamberty B, Euskirchen E, Talbot J, Frolking S, McGuire AD,**
461 **Tuittila ES.** 2012. The resilience and functional role of moss in boreal and arctic
462 ecosystems. New Phytologist **196:**49-67.
463 8. **Verhoeven JTA, Liefveld WM.** 1997. The ecological significance of organochemical
464 compounds in Sphagnum. Acta Botanica Neerlandica **46:**117-130.
465 9. **Mellegard H, Stalheim T, Hormazabal V, Granum PE, Hardy SP.** 2009.
466 Antibacterial activity of sphagnum acid and other phenolic compounds found in
467 Sphagnum papillosum against food-borne bacteria. Letters in Applied Microbiology
468 **49:**85-90.
469 10. **Freeman C, Ostle N, Kang H.** 2001. An enzymic 'latch' on a global carbon store - A
470 shortage of oxygen locks up carbon in peatlands by restraining a single enzyme. Nature
471 **409:**149-149.
472 11. **Stalheim T, Ballance S, Christensen BE, Granum PE.** 2009. Sphagnan - a pectin-like
473 polymer isolated from Sphagnum moss can inhibit the growth of some typical food
474 spoilage and food poisoning bacteria by lowering the pH. Journal of Applied
475 Microbiology **106:**967-976.
476 12. **Hajek T, Ballance S, Limpens J, Zijlstra M, Verhoeven JTA.** 2011. Cell-wall
477 polysaccharides play an important role in decay resistance of Sphagnum and actively
478 depressed decomposition in vitro. Biogeochemistry **103:**45-57.
479 13. **Lin X, Tfaily MM, Green SJ, Steinweg JM, Chanton P, Imvittaya A, Chanton JP,**
480 **Cooper W, Schadt C, Kostka JE.** 2014. Microbial Metabolic Potential for Carbon
481 Degradation and Nutrient (Nitrogen and Phosphorus) Acquisition in an Ombrotrophic
482 Peatland. Applied and Environmental Microbiology **80:**3531-3540.
483 14. **Leppanen S, Rissanen A, Tiirola M.** 2015. Nitrogen fixation in Sphagnum mosses is
484 affected by moss species and water table level. Plant and Soil **389:**185-196.
485 15. **Kostka JE, Weston DJ, Glass JB, Lilleskov EA, Shaw AJ, Turetsky MR.** 2016. The
486 Sphagnum microbiome: new insights from an ancient plant lineage. New Phytologist
487 **211:**57-64.

23

488    16.    **Dudova L, Hajkova P, Buchtova H, Opravilova V.** 2013. Formation, succession and
489         landscape history of Central-European summit raised bogs: A multiproxy study from the
490         Hruby Jesenik Mountains. Holocene **23:**230-242.

491    17.    **Ireland AW, Clifford MJ, Booth RK.** 2014. Widespread dust deposition on North
492         American peatlands coincident with European land-clearance. Vegetation History and
493         Archaeobotany **23:**693-700.

494    18.    **Swindles GT, Turner TE, Roe HM, Hall VA, Rea HA.** 2015. Testing the cause of the
495         Sphagnum austinii (Sull. ex Aust.) decline: Multiproxy evidence from a raised bog in
496         Northern Ireland. Review of Palaeobotany and Palynology **213:**17-26.

497    19.    **Galka M, Tobolski K, Gorska A, Lamentowicz M.** 2017. Resilience of plant and
498         testate amoeba communities after climatic and anthropogenic disturbances in a Baltic bog
499         in Northern Poland: Implications for ecological restoration. Holocene **27:**130-141.

500    20.    **Opelt K, Chobot V, Hadacek F, Schonmann S, Eberl L, Berg G.** 2007. Investigations
501         of the structure and function of bacterial communities associated with Sphagnum mosses.
502         Environmental Microbiology **9:**2795-2809.

503    21.    **Bragina A, Cardinale M, Berg C, Berg G.** 2013. Vertical transmission explains the
504         specific Burkholderia pattern in Sphagnum mosses at multi-geographic scale. Frontiers in
505         Microbiology **4:**10.

506    22.    **Bragina A, Maier S, Berg C, Muller H, Chobot V, Hadacek F, Berg G.** 2012. Similar
507         diversity of Alphaproteobacteria and nitrogenase gene amplicons on two related
508         Sphagnum mosses. Frontiers in Microbiology **3:**10.

509    23.    **Opelt K, Berg G.** 2004. Diversity and antagonistic potential of bacteria associated with
510         bryophytes from nutrient-poor habitats of the Baltic Sea coast. Applied and
511         Environmental Microbiology **70:**6569-6579.

512    24.    **Brussaard CPD, Wilhelm SW, Thingstad F, Weinbauer MG, Bratbak G, Heldal M,**
513         **Kimmance SA, Middelboe M, Nagasaki K, Paul JH, Schroeder DC, Suttle CA,**
514         **Vaque D, Wommack KE.** 2008. Global-scale processes with a nanoscale drive: the role
515         of marine viruses. ISME J **2:**575-578.

516    25.    **Jover LF, Effler TC, Buchan A, Wilhelm SW, Weitz JS.** 2014. The elemental
517         composition of virus particles: implications for marine biogeochemical cycles. Nat Rev
518         Micro **12:**519-528.

519    26.    **Wilhelm SW, Suttle CA.** 1999. Viruses and Nutrient Cycles in the SeaViruses play
520         critical roles in the structure and function of aquatic food webs. BioScience **49:**781-788.

521    27.    **Thingstad TF, Lignell R.** 1997. Theoretical models for the control of bacterial growth
522         rate, abundance, diversity and carbon demand. Aquatic Microbial Ecology **13:**19-27.

523    28.    **Roux S, Brum JR, Dutilh BE, Sunagawa S, Duhaime MB, Loy A, Poulos BT,**
524         **Solonenko N, Lara E, Poulain J, Pesant S, Kandels-Lewis S, Dimier C, Picheral M,**
525         **Searson S, Cruaud C, Alberti A, Duarte CM, Gasol JM, Vaque D, Bork P, Acinas**
526         **SG, Wincker P, Sullivan MB, Tara Oceans C.** 2016. Ecogenomics and potential
527         biogeochemical impacts of globally abundant ocean viruses. Nature **537:**689-+.

528    29.    **Simmonds P, Adams MJ, Benko M, Breitbart M, Brister JR, Carstens EB, Davison**
529         **AJ, Delwart E, Gorbalenya AE, Harrach B, Hull R, King AMQ, Koonin EV,**
530         **Krupovic M, Kuhn JH, Lefkowitz EJ, Nibert ML, Orton R, Roossinck MJ,**
531         **Sabanadzovic S, Sullivan MB, Suttle CA, Tesh RB, van der Vlugt RA, Varsani A,**
532         **Zerbini M.** 2017. Virus taxonomy in the age of metagenomics. Nature Reviews
533         Microbiology **15:**161-168.

24

534 30.   **Steward GF, Culley AI, Mueller JA, Wood-Charlson EM, Belcaid M, Poisson G.**
535       2013. Are we missing half of the viruses in the ocean? The ISME Journal **7:**672-679.
536 31.   **Wilhelm SW, Bird JT, Bonifer KS, Calfee BC, Chen T, Coy SR, Gainer PJ, Gann**
537       **ER, Heatherly HT, Lee J, Liang XL, Liu J, Armes AC, Moniruzzaman M, Rice JH,**
538       **Stough JMA, Tams RN, Williams EP, LeCleir GR.** 2017. A Student's Guide to Giant
539       Viruses Infecting Small Eukaryotes: From Acanthamoeba to Zooxanthellae. Viruses-
540       Basel **9:**18.
541 32.   **Yutin N, Wolf YI, Raoult D, Koonin EV.** 2009. Eukaryotic large nucleo-cytoplasmic
542       DNA viruses: Clusters of orthologous genes and reconstruction of viral genome
543       evolution. Virology Journal **6:**13.
544 33.   **Moniruzzaman M, Wurch LL, Alexander H, Dyhrman ST, Gobler CJ, Wilhelm**
545       **SW.** 2017. Virus-host relationships of marine single-celled eukaryotes resolved from
546       metatranscriptomics. Nature Communications **8:**10.
547 34.   **Schulz F, Yutin N, Ivanova NN, Ortega DR, Lee TK, Vierheilig J, Daims H, Horn**
548       **M, Wagner M, Jensen GJ, Kyrpides NC, Koonin EV, Woyke T.** 2017. Giant viruses
549       with an expanded complement of translation system components. Science **356:**82-+.
550 35.   **Moniruzzaman M, LeCleir GR, Brown CM, Gobler CJ, Bidle KD, Wilson WH,**
551       **Wilhelm SW.** 2014. Genome of brown tide virus (AaV), the little giant of the
552       Megaviridae, elucidates NCLDV genome expansion and host–virus coevolution.
553       Virology **466-467:**60-70.
554 36.   **Moniruzzaman M, Gann ER, Wilhelm SW.** 2018. Infection by a Giant Virus (AaV)
555       Induces Widespread Physiological Reprogramming in Aureococcus anophagefferens
556       CCMP1984 – A Harmful Bloom Algae. Frontiers in Microbiology **9**.
557 37.   **Wilhelm SW, Coy SR, Gann ER, Moniruzzaman M, Stough JMA.** 2016. Standing on
558       the shoulders of giant viruses: five lessons learned about large viruses infecting small
559       eukaryotes and the opportunities they create. Plos Pathogens **12:**5.
560 38.   **Kerepesi C, Grolmusz V.** 2017. The "Giant Virus Finder" discovers an abundance of
561       giant viruses in the Antarctic dry valleys. Archives of Virology **162:**1671-1676.
562 39.   **Rusin LY.** 2016. Metagenomics and biodiversity of sphagnum bogs. Molecular Biology
563       **50:**645-648.
564 40.   **Krupovic M, Koonin EV.** 2014. Evolution of eukaryotic single-stranded DNA viruses of
565       the Bidnaviridae family from genes of four other groups of widely different viruses.
566       Scientific Reports **4:**5347.
567 41.   **Krupovic M, Koonin EV.** 2015. Polintons: a hotbed of eukaryotic virus, transposon and
568       plasmid evolution. Nature Reviews Microbiology **13:**105.
569 42.   **Tomaru Y, Nagasaki K.** 2007. Flow cytometric detection and enumeration of DNA and
570       RNA viruses infecting marine eukaryotic microalgae. Journal of Oceanography **63:**215-
571       221.
572 43.   **Miranda JA, Culley AI, Schvarcz CR, Steward GF.** 2016. RNA viruses as major
573       contributors to Antarctic virioplankton. Environmental Microbiology **18:**3714-3727.
574 44.   **Quaiser A, Krupovic M, Dufresne A, Francez A-J, Roux S.** 2016. Diversity and
575       comparative genomics of chimeric viruses in Sphagnum-dominated peatlands. Virus
576       Evolution **2:**vew025-vew025.
577 45.   **Quaiser A, Dufresne A, Ballaud F, Roux S, Zivanovic Y, Colombet J, Sime-Ngando**
578       **T, Francez A-J.** 2015. Diversity and comparative genomics of Microviridae in
579       Sphagnum- dominated peatlands. Frontiers in Microbiology **6**.

25

580   46.   **Ballaud F, Dufresne A, Francez A-J, Colombet J, Sime-Ngando T, Quaiser A.** 2015.
581         Dynamics of Viral Abundance and Diversity in a Sphagnum-Dominated Peatland:
582         Temporal Fluctuations Prevail Over Habitat. Frontiers in Microbiology **6:**1494.

583   47.   **Dorigo U, Jacquet S, Humbert JF.** 2004. Cyanophage diversity, inferred from g20 gene
584         analyses, in the largest natural lake in France, Lake Bourget. Applied and Environmental
585         Microbiology **70:**1017-1022.

586   48.   **Roux S, Enault F, Robin A, Ravet V, Personnic S, Theil S, Colombet J, Sime-**
587         **Ngando T, Debroas D.** 2012. Assessing the Diversity and Specificity of Two Freshwater
588         Viral Communities through Metagenomics. Plos One **7:**12.

589   49.   **Comeau AM, Krisch HM.** 2008. The capsid of the T4 phage superfamily: The
590         evolution, diversity, and structure of some of the most prevalent proteins in the biosphere.
591         Molecular Biology and Evolution **25:**1321-1332.

592   50.   **Sakowski EG, Munsell EV, Hyatt M, Kress W, Williamson SJ, Nasko DJ, Polson**
593         **SW, Wommack KE.** 2014. Ribonucleotide reductases reveal novel viral diversity and
594         predict biological and ecological features of unknown marine viruses. Proceedings of the
595         National Academy of Sciences of the United States of America **111:**15786-15791.

596   51.   **Alexander H, Jenkins BD, Rynearson TA, Dyhrman ST.** 2015. Metatranscriptome
597         analyses indicate resource partitioning between diatoms in the field. Proceedings of the
598         National Academy of Sciences **112:**E2182-E2190.

599   52.   **Wilson RM, Hopple AM, Tfaily MM, Sebestyen SD, Schadt CW, Pfeifer-Meister L,**
600         **Medvedeff C, McFarlane KJ, Kostka JE, Kolton M, Kolka RK, Kluber LA, Keller**
601         **JK, Guilderson TP, Griffiths NA, Chanton JP, Bridgham SD, Hanson PJ.** 2016.
602         Stability of peatland carbon to rising temperatures. Nature Communications **7:**10.

603   53.   **Hanson PJ, Riggs JS, Nettles WR, Phillips JR, Krassovski MB, Hook LA, Gu L,**
604         **Richardson AD, Aubrecht DM, Ricciuto DM.** 2017. Attaining whole-ecosystem
605         warming using air and deep-soil heating methods with an elevated $CO_2$ atmosphere.
606         Biogeosciences **14:**861.

607   54.   **Warren MJ, Lin XJ, Gaby JC, Kretz CB, Kolton M, Morton PL, Pett-Ridge J,**
608         **Weston DJ, Schadt CW, Kostka JE, Glass JB.** 2017. Molybdenum-Based Diazotrophy
609         in a Sphagnum Peatland in Northern Minnesota. Applied and Environmental
610         Microbiology **83:**14.

611   55.   **Kopylova E, Noe L, Touzet H.** 2012. SortMeRNA: fast and accurate filtering of
612         ribosomal RNAs in metatranscriptomic data. Bioinformatics **28:**3211-3217.

613   56.   **O'Leary NA, Wright MW, Brister JR, Ciufo S, Haddad D, McVeigh R, Rajput B,**
614         **Robbertse B, Smith-White B, Ako-Adjei D, Astashyn A, Badretdin A, Bao Y,**
615         **Blinkova O, Brover V, Chetvernin V, Choi J, Cox E, Ermolaeva O, Farrell CM,**
616         **Goldfarb T, Gupta T, Haft D, Hatcher E, Hlavina W, Joardar VS, Kodali VK, Li**
617         **W, Maglott D, Masterson P, McGarvey KM, Murphy MR, O'Neill K, Pujar S,**
618         **Rangwala SH, Rausch D, Riddick LD, Schoch C, Shkeda A, Storz SS, Sun H,**
619         **Thibaud-Nissen F, Tolstoy I, Tully RE, Vatsan AR, Wallin C, Webb D, Wu W,**
620         **Landrum MJ, Kimchi A, et al.** 2016. Reference sequence (RefSeq) database at NCBI:
621         current status, taxonomic expansion, and functional annotation. Nucleic Acids Res
622         **44:**D733-745.

623   57.   **Koonin EV.** 1991. The phylogeny of RNA-dependent RNA polymerases of positive-
624         strand RNA viruses. Journal of General Virology **72:**2197-2206.

26

625 58. **Finn RD, Coggill P, Eberhardt RY, Eddy SR, Mistry J, Mitchell AL, Potter SC,**
626 **Punta M, Qureshi M, Sangrador-Vegas A, Salazar GA, Tate J, Bateman A.** 2016.
627 The Pfam protein families database: towards a more sustainable future. Nucleic Acids
628 Research **44:**D279-D285.
629 59. **Finn RD, Clements J, Arndt W, Miller BL, Wheeler TJ, Schreiber F, Bateman A,**
630 **Eddy SR.** 2015. HMMER web server: 2015 update. Nucleic Acids Research **43:**W30-
631 W38.
632 60. **Marchler-Bauer A, Derbyshire MK, Gonzales NR, Lu SN, Chitsaz F, Geer LY, Geer**
633 **RC, He J, Gwadz M, Hurwitz DI, Lanczycki CJ, Lu F, Marchler GH, Song JS,**
634 **Thanki N, Wang ZX, Yamashita RA, Zhang DC, Zheng CJ, Bryant SH.** 2015. CDD:
635 NCBI's conserved domain database. Nucleic Acids Research **43:**D222-D226.
636 61. **Liu WZ, Xie YB, Ma JY, Luo XT, Nie P, Zuo ZX, Lahrmann U, Zhao Q, Zheng YY,**
637 **Zhao Y, Xue Y, Ren J.** 2015. IBS: an illustrator for the presentation and visualization of
638 biological sequences. Bioinformatics **31:**3359-3361.
639 62. **Finn RD, Attwood TK, Babbitt PC, Bateman A, Bork P, Bridge AJ, Chang HY,**
640 **Dosztanyi Z, El-Gebali S, Fraser M, Gough J, Haft D, Holliday GL, Huang HZ,**
641 **Huang XS, Letunic I, Lopez R, Lu SN, Marchler-Bauer A, Mi HY, Mistry J, Natale**
642 **DA, Necci M, Nuka G, Orengo CA, Park Y, Pesseat S, Piovesan D, Potter SC,**
643 **Rawlings ND, Redaschi N, Richardson L, Rivoire C, Sangrador-Vegas A, Sigrist C,**
644 **Sillitoe I, Smithers B, Squizzato S, Sutton G, Thanki N, Thomas PD, Tosatto SCE,**
645 **Wu CH, Xenarios I, Yeh LS, Young SY, Mitchell AL.** 2017. InterPro in 2017-beyond
646 protein family and domain annotations. Nucleic Acids Research **45:**D190-D199.
647 63. **Edgar RC.** 2004. MUSCLE: a multiple sequence alignment method with reduced time
648 and space complexity. Bmc Bioinformatics **5:**1-19.
649 64. **Kumar S, Stecher G, Tamura K.** 2016. MEGA7: Molecular evolutionary genetics
650 analysis version 7.0 for bigger datasets. Molecular Biology and Evolution **33:**1870-1874.
651 65. **Guindon S, Dufayard JF, Lefort V, Anisimova M, Hordijk W, Gascuel O.** 2010.
652 New algorithms and methods to estimate maximum-likelihood phylogenies: assessing the
653 performance of PhyML 3.0. Systematic Biology **59:**307-321.
654 66. **Matsen FA, Kodner RB, Armbrust EV.** 2010. pplacer: linear time maximum-
655 likelihood and Bayesian phylogenetic placement of sequences onto a fixed reference tree.
656 Bmc Bioinformatics **11:**16.
657 67. **Letunic I, Bork P.** 2016. Interactive tree of life (iTOL) v3: an online tool for the display
658 and annotation of phylogenetic and other trees. Nucleic Acids Research **44:**W242-W245.
659 68. **Clark KR, Gorley RN.** 2015. PRIMER v7: User manual/tutorial. PRIMER-E, Plymouth.
660 69. **Clarke KR, Somerfield PJ, Gorley RN.** 2008. Testing of null hypotheses in exploratory
661 community analyses: similarity profiles and biota-environment linkage. Journal of
662 Experimental Marine Biology and Ecology **366:**56-69.
663 70. **Shannon P, Markiel A, Ozier O, Baliga NS, Wang JT, Ramage D, Amin N,**
664 **Schwikowski B, Ideker T.** 2003. Cytoscape: A software environment for integrated
665 models of biomolecular interaction networks. Genome Research **13:**2498-2504.

666

667

668    **Figure Legends**

669    Figure 1: Phylogenetic placement of identified phage major capsid protein contigs (red) on a

670    Myovirus *gp23* maximum likelihood reference tree (references in black). Node support (aLRT-

671    SH statistic) > 50% are shown. Contigs are shown with their abundance (log$_2$ transformed TPM)

672    in a heatmap surrounding the tree. Sample order on the heatmap is provided in the inset.

673    Figure 2: Phylogenetic placement of identified ssRNA virus RNA-dependent RNA polymerase

674    contigs on maximum likelihood reference tree. Branch width represents the number of contigs

675    placed on the reference branch. Node support (aLRT-SH statistic) >50% are shown.
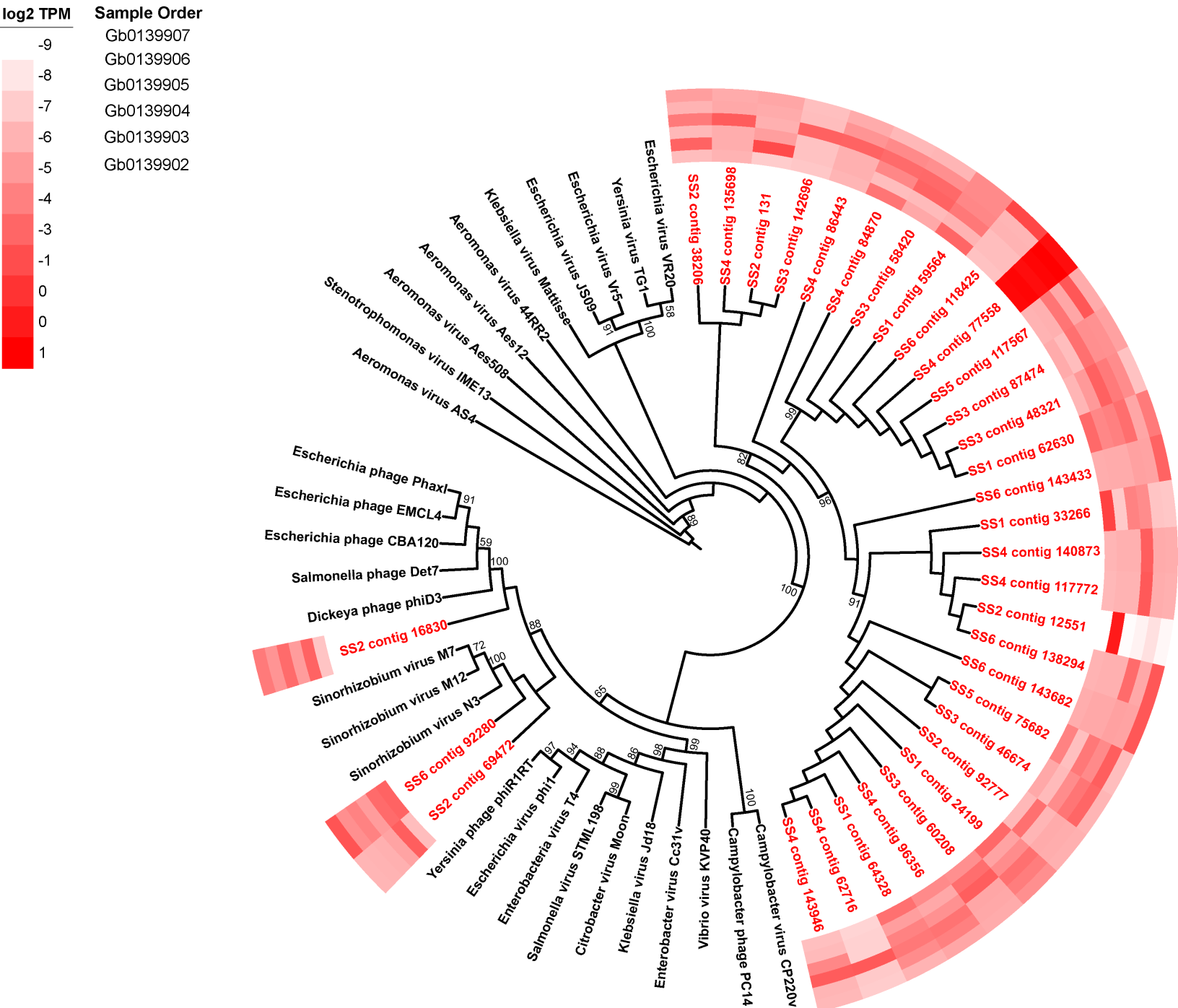
676    Figure 3: Phylogeny, genome architecture, and abundance of partial ssRNA virus genomes. Tree

677    represents phylogenetic placement of RDRP gene regions from partial ssRNA virus genome

678    contigs (red) on a maximum likelihood reference tree (references in black). Node support (aLRT-

679    SH statistic) >50% are shown. Center panel represents genome architecture determined by

680    conserved domain search and ORF prediction. Length of contigs and gene regions is measured in

681    kb. Heat map in right panel shows abundance of reads mapped to partial genome contigs in log$_2$

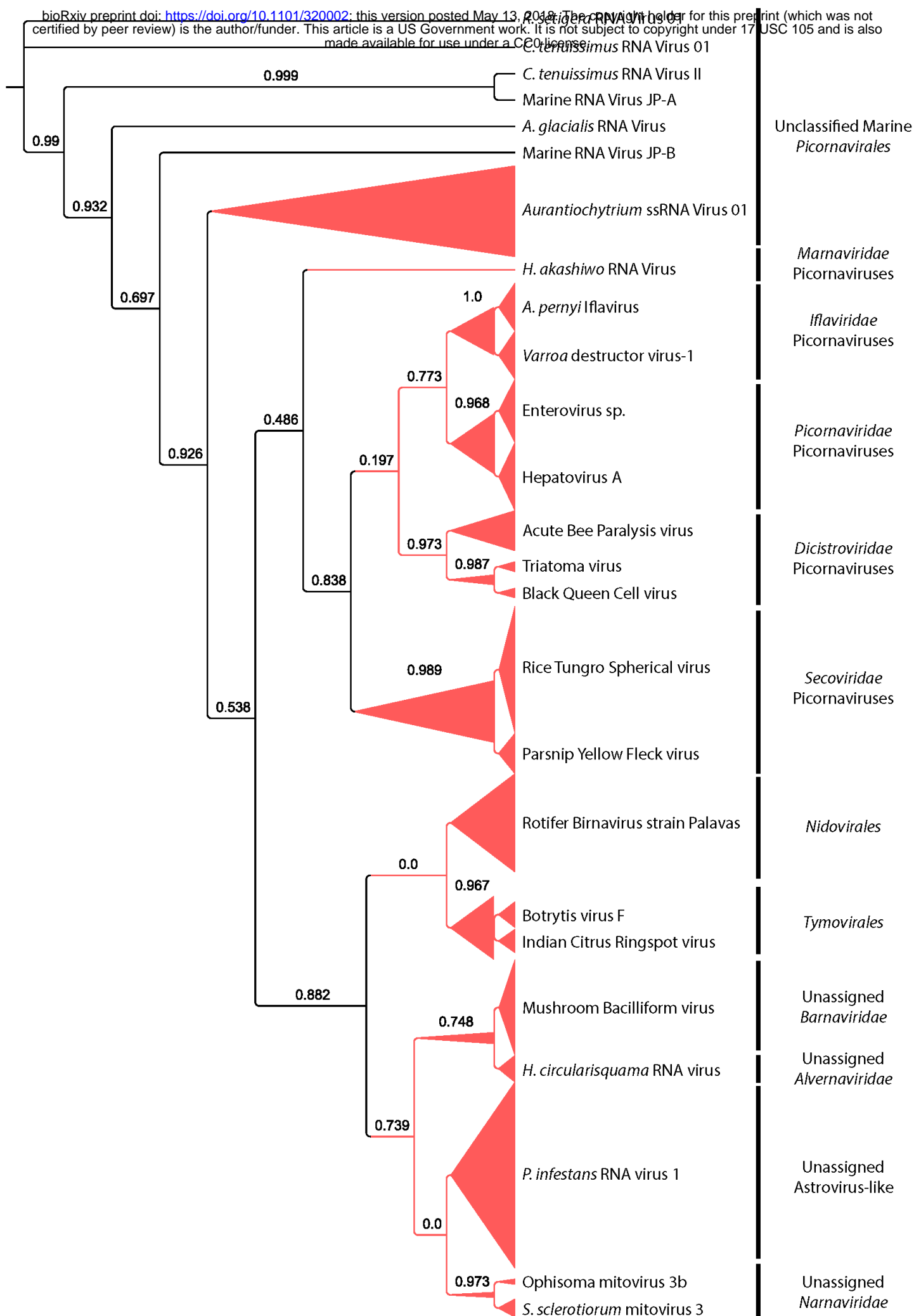682    TPM from each of the 6 metatranscriptome libraries.

683    Figure 4. Phylogenetic placement of identified NCLDV major capsid protein contigs (red) on a

684    maximum likelihood reference tree (references in black). Node support (aLRT-SH statistic)

685    >50% are shown. Contigs are shown with their abundance (log$_2$ transformed TPM) in a heatmap

686    surrounding the tree.

687    Figure 5: Phylogenetic placement of identified virophage A.) major capsid protein and B.)

688    ATPase contigs (red) on a maximum likelihood reference tree (references in black). Node

689    support (aLRT-SH statistic) >50% are shown.

28

690    Figure 6: Correlation co-occurrence network analysis of conserved viral gene and host RNA

691    polymerase expression for A.) bacteriophage (Gp23), B.) ssRNA viruses (RDRP), and C.)

692    NCLDVs (NCLDV MCP). Nodes in red represent virus contigs and blue nodes represent

693    potential hosts. Nodes are connected by edges colored according to the Pearson correlation

694    coefficient values between to contigs. Only relationships with contigs expressed in more than
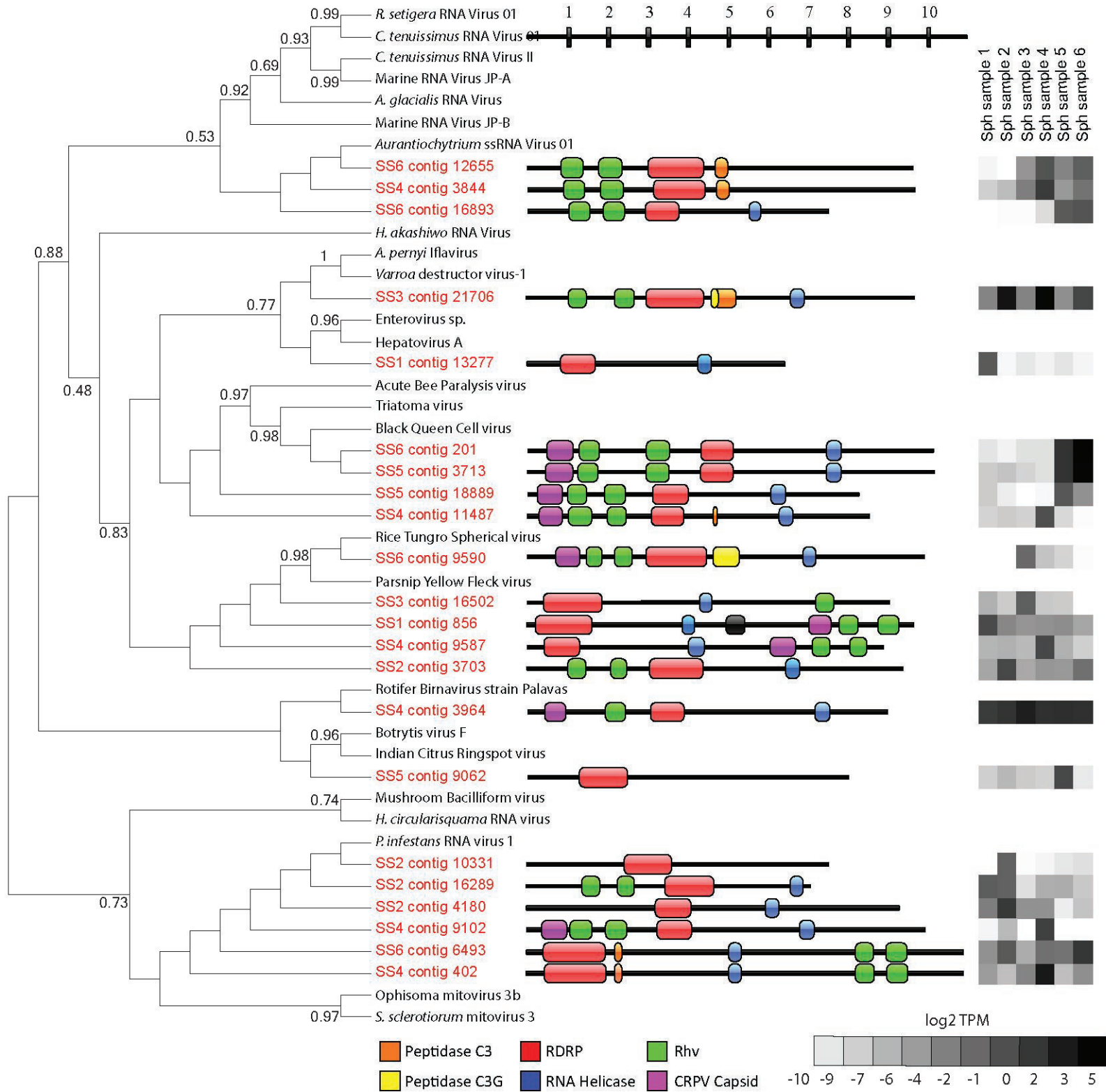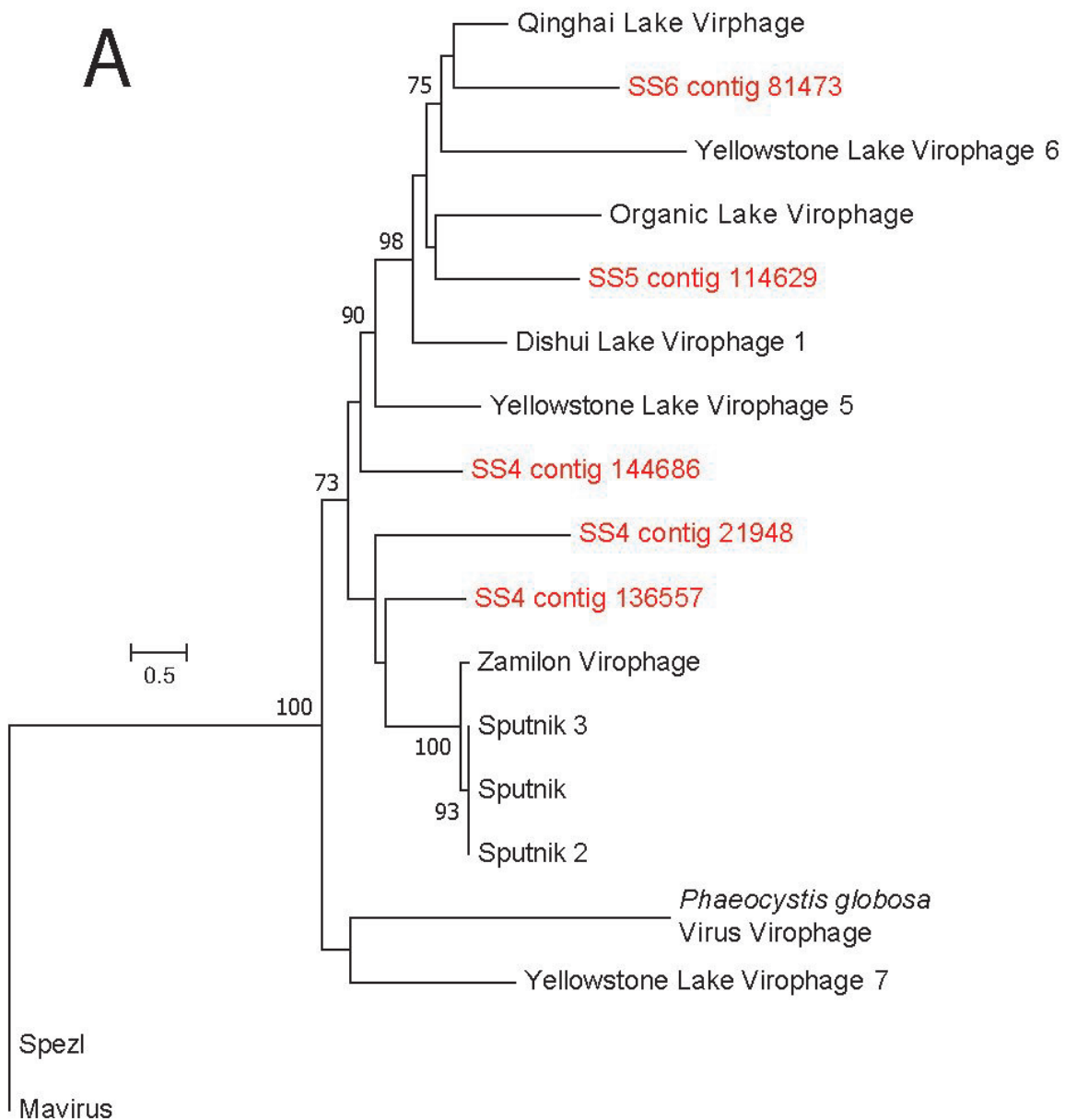
695    one sample are shown.

696

Figure 3: Phylogeny, genome architecture, and abundance of partial ssRNA virus genomes. Tree represents phylogenetic placement of RDRP gene regions from partial ssRNA virus genome contigs on a maximum likelihood reference tree. Node support (aLRT-SH statistic) >50% are shown. Center panel represents genome architecture determined by conserved domain search and ORF prediction. Length of contigs and gene regions is measured in kb. Heat map in right panel shows abundance of reads mapped to partial genome contigs in log2 TPM from each of the 6 metatranscriptome libraries.

A

- Qinghai Lake Virphage
- 75 — SS6 contig 81473
- Yellowstone Lake Virophage 6
- 98 — Organic Lake Virophage
- SS5 contig 114629
- 90 — Dishui Lake Virophage 1
- Yellowstone Lake Virophage 5
- 73 — SS4 contig 144686
- SS4 contig 21948
- SS4 contig 136557
- Zamilon Virophage
- 100 — Sputnik 3
- Sputnik
- 93 — Sputnik 2
- *Phaeocystis globosa* Virus Virophage
- Yellowstone Lake Virophage 7
- 100
- Spezl
- Mavirus

0.5

B

- Organic Lake Virophage
- Yellowstone Lake Virophage 6
- SS5 contig 73836
- 87 — Yellowstone Virophage 7
- SS6 contig 28119
- Qinghai Lake Virophage
- 92 — 87 — Dishui Lake Virophage1
- *Phaeocystis globosa* Virus Virophage
- Zamilon Virophage
- 92 — Sputnik
- 89 — Sputnik 2
- Sputnik 3
- 100
- Yellowstone Lake Virophage 5
- Mavirus
- Spezl

0.5