# 1 Transcriptome Mining Reveals a Spectrum of RNA

# 2 Viruses in Primitive Plants

3

4 Jonathon C.O. Mifsud[a,b#], Rachael V. Gallagher[b.c], Edward C. Holmes[a], Jemma L. Geoghegan[d,e]

5

6 [a]Sydney Institute for Infectious Diseases, School of Life and Environmental Sciences and

7 School of Medical Sciences, The University of Sydney, Sydney, NSW 2006, Australia.

8 [b]Department of Biological Sciences, Macquarie University, North Ryde, NSW 2109 Australia

9 [c]Hawkesbury Institute for the Environment, Western Sydney University, Locked Bag 1797,

10 Penrith, NSW 2751, Australia

11 [d]Department of Microbiology and Immunology, University of Otago, Dunedin 9016, New

12 Zealand.

13 [e]Institute of Environmental Science and Research, Wellington 5018, New Zealand.

14

15 Running Head: Mining Reveals RNA Viruses in Primitive Plants

16 #Address correspondence to: Jonathon C.O. Mifsud, jmif9945@uni.sydney.edu.au

17 Word Counts: Abstract 250, Importance 125, Text 9074.

18 Keywords: plant virus, virus discovery, algae virus, *Benyviridae*, *Bunyavirales*, *Secoviridae*,

19 evolution

**Abstract**

20

21    Current knowledge of plant viruses stems largely from those affecting economically important

22    plants. Yet, plant species in cultivation represent a small and bias subset of the plant kingdom.

23    Here, we describe virus diversity and abundance from a survey of 1079 transcriptomes from

24    species across the breadth of the plant kingdom (Archaeplastida) by analysing open-source

25    data from the One Thousand Plant Transcriptomes Initiative (1KP). We identified 104 potentially

26    novel viruses, of which 40% comprised single-stranded positive-sense RNA viruses across eight

27    orders, including members of the *Hepelivirales*, *Tymovirales*, *Cryppavirales*, *Martellivirales* and

28    *Picornavirales*. One-third of the newly described viruses comprised double-stranded RNA

29    viruses from the orders *Durnavirales* and *Ghabrivirales*. The remaining were negative-sense

30    RNA viruses from the *Rhabdoviridae*, *Aspiviridae, Yueviridae, Phenuiviridae* and the newly

31    proposed *Viridisbunyaviridae.* Our analysis considerably expands the known host range of 13

32    virus families to include lower plants (e.g., *Benyviridae* and *Secoviridae*) and four virus families

33    to include algae hosts (e.g., *Tymoviridae* and *Chrysoviridae)*. The discovery of the first 30 kDa

34    movement protein in a non-vascular plant, suggests that the acquisition of plant virus movement

35    proteins occurred prior to the emergence of the plant vascular system. More broadly, however,

36    a co-phylogeny analysis revealed that the evolutionary history of these families is largely driven

37    by cross-species transmission events. Together, these data highlight that numerous RNA virus

38    families are associated with older evolutionary plant lineages than previously thought and that

39    the scarcity of RNA viruses found in lower plants to date likely reflects a lack of investigation

40    rather than their absence.

**Importance**

Our knowledge of plant viruses is mainly limited to those infecting economically important host species. In particular, we know little about those viruses infecting primitive plant lineages such as the ferns, lycophytes, bryophytes and charophytes. To expand this understanding, we conducted a broad-scale viral survey of species across the breadth of the plant kingdom. We find that primitive plants harbour a wide diversity of RNA viruses including some that are sufficiently divergent to comprise a new virus family. The primitive plant virome we reveal offers key insights into the evolutionary history of core plant virus gene modules and genome segments. More broadly, this work emphasises that the scarcity of viruses found in these species to date likely reflects the absence of research in this area.

## 1. Introduction

Viruses are responsible for almost 50% of all emerging plant disease (1). Historically, virus identification and characterisation have focused on pathogenic viruses that infect species of economic importance with 69% of the current phytovirosphere — the total assemblage of viruses across the plant kingdom — discovered in cultivated plant species even though they represent less than 0.17% of all known plant diversity (2, 3). Importantly, the advent of metagenomic sequencing technology enables the comprehensive screening of plant tissues for novel and known viruses (4). Despite this, virus diversity in the vast majority of plants remains unquantified (5).

Our ability to infer the origins and diversification of the phytovirosphere from genomic data requires adequate sampling of the viruses across the plant kingdom. Several key plant groups are severely underrepresented or absent in previous studies of the phytovirosphere, including green algae (excluding the Chlorophytes), lower plants, gymnosperms and several angiosperm orders (5, 6). Improving knowledge across these groups will undoubtedly help uncover the evolutionary history of plant virus lineages. For instance, an analysis of the evolutionary history of viruses from algal ancestors might reveal deep associations that shaped the trajectory of plant evolution, including how the key evolutionary transitions of plants – such as terrestrialisation – have shaped the contemporary land plant virome (5). Similarly, through broad sampling across the plant kingdom, we can gain a stronger understanding of the acquisition of viruses through cross-species transmission from plant-associated organisms such as invertebrates, fungi, or protists (5).

The majority (68%) of the currently documented genera of plant viruses have positive-sense single-stranded RNA (+ssRNA) genomes and the majority of virus diversity is known only from angiosperms (7) (Figure 1). Currently, 16 viruses belonging to 12 virus families have been found

4

75     in gymnosperms (8-12). Outside of several viruses found in ferns, we know little of the diversity

76     of viruses in the lycophytes, bryophytes and charophytes that together encompass ~27,000

77     species (13-16) (Figure 1). A partial analysis of published transcriptome data detected

78     homologs of the canonical RNA virus RNA-dependent RNA polymerase (RdRp) in algae,

79     several lower plants and gymnosperms (17). However, it is yet to be determined whether

80     viruses that infect freshwater algae – that include the *Zygnematophyceae* ancestors of land

81     plants – resemble those infecting angiosperms or that of the green algae (chlorophytes) which

82     are dominated by double-stranded DNA (dsDNA) viruses particularly from the *Phycodnaviridae*

83     (18). To date, two +ssRNA viruses related to the benyvirids have been identified in freshwater

84     algae (19, 20). Unlike the Chlorophyta, the Charophyta characteristically contain

85     plasmodesmata and homologs of the key components of the land plant innate immune system,

86     both of which have been speculated to explain the absence of double-strand (ds) DNA viruses

87     in land plants (5, 21, 22). An understanding of the viruses infecting the Charophyta and other

88     lower plants is required to effectively test these ideas.

89     Transcriptome mining has become an inexpensive and efficient method of virus discovery that

90     leverages previous investment (23-29). To this end, we mined the transcriptome data generated

91     by the One Thousand Plant Transcriptomes Initiative (1KP) using sequence homology searches

92     of known plant viruses. The 1KP project provides a major untapped source of polyA-selected

93     transcriptome data for virus discovery drawn from species across the breadth of the plants in a

94     broad sense including green plants (Viridiplantae), glaucophytes (Glaucophyta), red algae

95     (Rhodophyta) (30, 31). Our broad aim was to revise our understanding of the phytovirosphere

96     using data across the plant kingdom and undertake phylogenetic analyses of plant viruses to

97     provide insights into their origins and diversification.

98     **2. Methods**

5

## 2.1 Transcriptome data generation

99

100 The 1KP generated RNA sequencing libraries from 1,143 species across the breadth of the

101 plant kingdom (30). In addition, 30 Chromista and red alga species were also included. Due to

102 the diversity of species examined, samples were obtained from multiple sources including field

103 collections, greenhouses, culture collections and laboratory specimens (32). For the majority of

104 species, young leaves or shoots were collected, although occasionally a mix of vegetative and

105 reproductive tissues was used. To avoid RNA degradation, RNA extraction was performed

106 immediately after tissue collection or tissue was frozen in liquid nitrogen and stored in a -80°C

107 until extraction (32). Several extraction protocols were used including CTAB and TRIzol (see

108 (32) for complete details). All sequencing was conducted at BGI-Shenzhen, China, using a

109 combination of in-house protocols or TruSeq chemistry (32). All libraries were prepared from

110 polyA RNA. Paired-end sequencing was initially completed using Illumina GAII machines (11%

111 of libraries) with a ~72bp read length but later the HiSeq platform was used (89% of libraries)

112 with a 90 bp read length (32).

## 2.2 Surveying for viruses in the 1KP

113

114 Raw transcriptomes (n = 1079, belonging to 960 plants species) from the 1KP major release

115 were downloaded from the NCBI Short Read Archive (SRA) database (BioProject accession

116 PRJEB21674) and converted to FASTQ format using the SRA Toolkit program fastq-dump in

117 combination with the parallel-fastq-dump wrapper (https://github.com/rvalieris/parallel-fastq-

118 dump) (33). One hundred transcriptomes within the BioProject were not publicly available

119 (released 22/08/2019) at the commencement of this study and thus not analysed.

120 Transcriptomes from the 1KP pilot study (BioProject accession PRJEB4921) and secondary

121 project (BioProject accession PRJEB8056) were similarly not analysed. To reduce the

122 downstream computing resources needed, raw sequences were mapped to their respective

123    host genome scaffold using bowtie2 (34). Genome scaffolds were assembled as part of a

124    previous study (30). Where genome scaffolds were not available (n = 2) all reads were

125    assembled *de novo*. Trinity RNA-seq (v2.1.1) was used to quality trim and assemble *de novo*

126    the unaligned reads captured from mapping (35). The assembled contigs were then assigned to

127    known virus families and annotated through similarity searches against the NCBI nucleotide

128    database (nt), the non-redundant protein database (nr) and a custom viral RdRp database using

129    BLASTN and Diamond (BLASTX) (36, 37). To filter out weak BLAST sequence matches an e-

130    value cut-off of $1 \times 10^{-10}$ was employed. To identify potential false positives, putative viral

131    contigs were manually compared across the three BLAST searches (nt, nr and RdRp) to ensure

132    matches to virus-associated sequences were consistent.

133    **2.3 Virus filtering and abundance calculations**

134    For all analyses, we focused on virus families known to infect plants or algae. As our analyses

135    rely on sequence-based similarity searches for virus detection it is necessarily biased towards

136    viruses that exhibit to existing virus families. Together, the Virus-Host database (38) and the

137    International Committee on Taxonomy of Viruses (39) were used to develop a list of plant virus

138    families and genera to filter out virus-like contigs associated with vertebrate, invertebrate or

139    fungi hosts based upon their top BLASTx and BLASTn matches. Packages within the Tidyverse

140    collection (v1.3.0) in RStudio were used to complete these tasks (40-42). Where the host was

141    ambiguous (e.g., belonged to a family or genera known to infect both plant and fungal species)

142    the contig was inspected manually.

143    The relative abundance of each transcript within the host transcriptome was calculated using

144    RNA-Seq by Expectation-Maximization (v1.2.28) (43). To account for variation in the number of

145    unaligned reads between libraries after mapping, contig abundance was standardised by the

7

146    total number of unaligned paired reads. Contigs under 200 nucleotides in length were excluded

147    from further analysis.

148    **2.4 Genome extension and annotation**

149    Where a novel virus-like contig was discovered, we re-assembled the complete library – without

150    removing host reads – in an attempt to recover a complete virus genome. For all re-assembled

151    libraries, we recalculated abundance measurements to account for both host and non-host

152    reads. The recalculated abundance measurements are shown in Supplementary Table 4. We

153    further re-assembled all libraries belonging to non-flowering plants (n = 402). Reads were

154    mapped onto virus-like contigs using Bbmap and heterogeneous coverage and potential

155    misassemblies were manually resolved using Geneious (v11.0.9) (44, 45).

156    To determine whether a virus was novel, we followed the criteria as specified by The

157    International Committee on Taxonomy of Viruses (39) (http://www.ictvonline.org/). Novel viruses

158    were named using a combination of the host common name - if documented – and the

159    associated virus taxonomic group (e.g., *Interrupted club-moss deltapartitivirus*). In cases where

160    host assignment proved difficult the suffix "associated" was added to the host name to signify

161    this (e.g., *Calypogeia fissa associated deltaflexivirus*). Where the taxonomic position of a virus

162    was ambiguous the suffix "-like" was used (e.g., *Goldenrod fern qin-like virus*). Virus acronyms

163    were created using a combination of the first and/or second letters of the host common name - if

164    documented – and virus taxonomic group (e.g., *Leucodon julaceus beny-like virus* (LjBV)).

165    Where multiple related viruses were found in the same host, we assigned each a number (e.g.,

166    *Odontoschisma prostratum bunyavirus 3* (OdprBV3)).

167    The percentage identity among virus sequences was calculated via multiple sequence

168    alignments using Clustal Omega (v1.2.3) (46). The RdRp protein coding domain was used for

8

169    all sequence alignments. Percentage identity matrices were converted to heat map plots using a

170    custom R script provided by (28).

171    To characterise functional domains, predicted protein sequences along with their closest viral

172    relatives were subjected to a domain-based search using the Conserved Domain Database

173    (v3.18) (https://www.ncbi.nlm.nih.gov/Structure/cdd/cdd.shtml) and cross-referenced with the

174    PFAM (v34.0) and Uniclust30 (v2018_08) databases available within the MMseqs2 webserver

175    (47). To recover additional annotations, we used HHpred within the MPI Bioinformatics Toolkit

176    webserver to query the PDB_mmCIF70 (v.12_Oct), SCOPe70 (v2.08), UniProt-SwissProt-

177    viral70 (v3_Nov_2021) and TIGRFAMs (v15.0) databases (48). Virus genome diagrams were

178    produced using the program littlegenomes (49). Where available NCBI/GenBank CDS

179    information was used to annotate reference virus sequences (50).

180    **2.5 Detection of endogenous virus elements**

181    All genome scaffolds produced by the 1KP were used as a database in which we queried using

182    the protein translations of the viruses discovered in this study. Endogenous viral elements (i.e.,

183    EVEs) were detected using the tblastn algorithm (51). The search threshold was limited to 100

184    amino acids in length with an e-value cut off of $1 \times 10^{-20}$. Where multiple hits across several plant

185    scaffolds were observed we manually examined the sequence. Suspected endogenous virus

186    sequences were queried against a subset whole-genome shotgun contig database which

187    included green plants (taxid: 33090) and red algae (taxid: 2763). In addition, the virus-like

188    sequences discovered in this study were checked for host gene contamination using the

189    contamination function implemented in CheckV (v0.8.1) (52). All potential endogenous

190    sequences were removed from further analyses.

191    **2.6 Assessing library contamination by eukaryotes, bacteria, and protozoa**

192   For libraries in which a novel virus was discovered we investigated whether reads belonging to

193   other eukaryotes were also present in the sequencing libraries. To achieve this, we obtained

194   taxonomic identification for raw reads in each library – without the removal of host reads – by

195   aligning them to the NCBI nt database using the KMA aligner and the CCMetagen program (53,

196   54). Sequence abundance was calculated by counting the number of nucleotides matching the

197   reference sequence with an additional correction for template length (the default parameter in

198   KMA). Krona charts generated by CCMetagen were edited were further edited in Adobe

199   Illustrator (https://www.adobe.com) (55). Library contamination was also assessed by the 1KP

200   and used to inform our host-virus assignments (31).

201   **2.7 Phylogenetic analysis of plant viruses**

202   Phylogenetic trees of the plant-associated viruses discovered here were inferred using a

203   maximum likelihood approach. We combined our translated virus contigs with known virus

204   protein sequences from each respective virus family taken from NCBI/GenBank (50).

205   Sequences were then aligned with the program Clustal Omega (v1.2.3) with default parameters

206   (46). Sites of ambiguity were removed using trimAl (v1.2) (56). To estimate phylogenetic trees,

207   selection of the best-fit model of amino acid substitution was determined using the Akaike

208   information criterion, corrected AIC, and the Bayesian information criterion with the ModelFinder

209   function (-m MFP) in IQ-TREE (57, 58). All phylogenetic trees were created using IQ-TREE with

210   1000 bootstrap replicates. Phylogenetic trees were annotated with FigTree (v1.4.4) (59) and

211   further edited in Adobe Illustrator (https://www.adobe.com).

212   To visualise the occurrence of cross-species transmission and virus-host co-divergence across

213   plant virus families, we reconciled the co-phylogenetic relationship between viruses and their

214   hosts. For each select plant virus family, a vascular plant host cladogram was constructed using

215   trees from (60) and (61), using the R package V.PhyloMaker (v0.1.0) (62). As lower plants and

216   non-plant species are not present in the V.PhyloMaker megatree, these hosts were added to the

217   cladogram using the software phyloT, a phylogenetic tree generator based on NCBI taxonomy

218   (http://phylot.biobyte.de/) as well as topologies available in the appropriate literature. The host

219   information was obtained from the NCBI Virus database (accessed 14/12/2021) and available

220   literature (63) A tanglegram that graphically represents the correspondence between host and

221   virus trees was created using the R packages phytools (v0.7-80) and APE (v5.5) (64, 65). Virus

222   sequences from each family were obtained through a broad survey of all virus genomic data

223   available on GenBank. The virus phylogenies used in the co-phylogenies were constructed as

224   detailed above. To quantify the relative frequencies of cross-species transmission versus virus-

225   host co-divergence we reconciled the co-phylogenetic relationship between viruses and their

226   hosts using the Jane co-phylogenetic software package (66). Jane employs a maximum

227   parsimony approach to determine the best 'map' of the virus phylogeny onto the host

228   phylogeny. The cost of duplication, host-jump and extinction event types were set to one, while

229   host-virus co-divergence was set to zero as it was considered the likely null event. Following the

230   parsimony principle, the reconciliation proceeds by minimising the total event cost. The number

231   of generations and the population size was both set to 100. Jane was chosen over its successor

232   eMPRess as it allows for a virus to be associated with multiple host species and handle

233   polytomies (67). For a multi-host virus, we represented each association as a polytomy on the

234   virus phylogeny.

235   **2.8 Assigning plant host clades**

236   Each plant host was assigned to each clade in a previous study based upon their phylogenetic

237   positioning and lineage information (30). To improve clarity when colouring the phylogenies

238   (although not the tanglegrams) we reduced the number of clades from 25 to ten (core eudicots,

11

239    basal eudicots, monocots, basalmost angiosperms, gymnosperms, fern and fern allies, non-

240    vascular, green algae, red algae and lastly Chromista) by combining those that were closely

241    related or potentially overlapping to increase the number of species in each group (SI Table 1).

242    **2.9 Data availability**

243    The raw One Thousand Plant Transcriptomes Initiative sequence reads are available at

244    BioProject PRJEB21674. All viral genomes and corresponding sequences assembled in this

245    study have been deposited in NCBI GenBank and assigned accession numbers xxxx-yyyy.
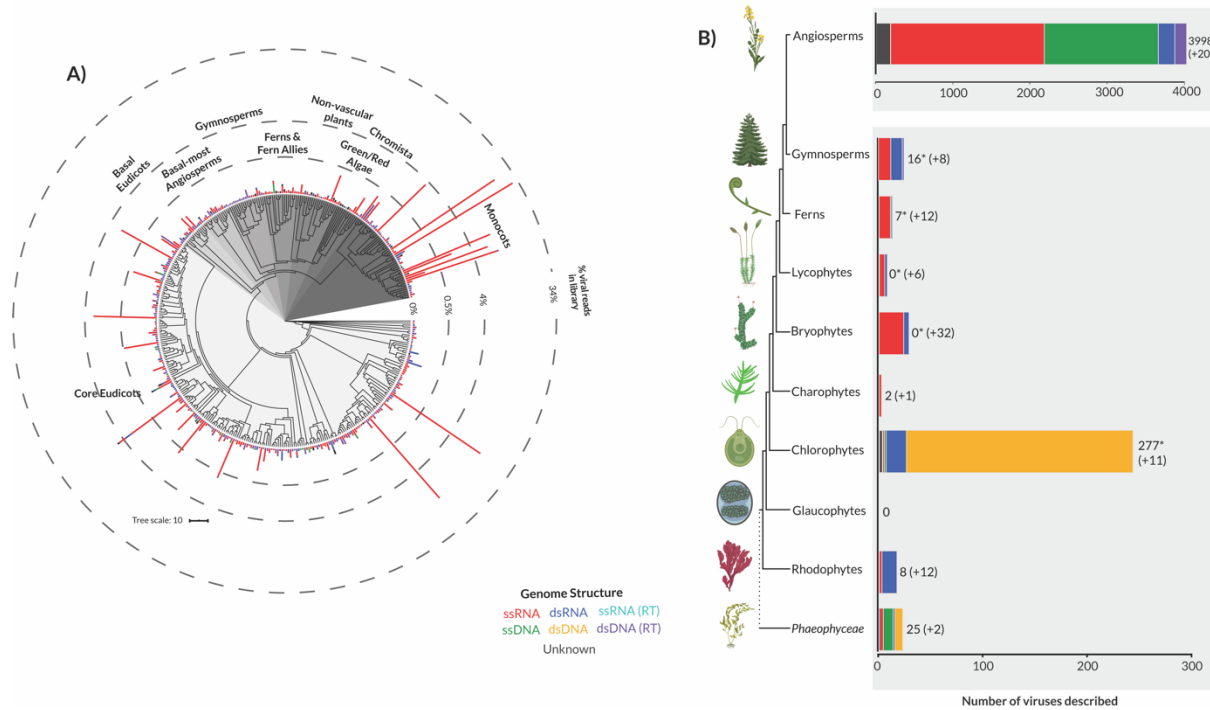
246    **3. Results**

247    We characterised the viruses found in the transcriptomes of 960 plant species within the 1KP

248    major release. The transcriptomes represented a broad taxonomic sampling across the

249    Archaeplastida (green plants, glaucophytes and red algae). Sequencing libraries had a median

250    of 25,187,714 paired reads (range 10,156,464–46,650,336). A median of 82% of reads (range

251    1%-96%) in these libraries mapped to host genome scaffolds and were subsequently removed.

252    *De novo* assembly of the sequencing reads resulted in a median of 36,015 contigs (range

253    1,396–146,217) per library, with a total of 41,256,176 contigs generated (SI Table 2).

254    **3.1 Diversity and abundance of plant viruses**

255    In total, virus-like transcripts were found for 603 plant species; 69% of these were plant-

256    associated while numerous identified sequences shared high similarity to non-plant associated

257    viruses including those known to infect fungi, invertebrate and vertebrate hosts. Among the non-

258    plant-associated virus transcripts, 34% were unclassified (10% of total virus-like transcripts)

259    such that they were most closely related to a virus sequence with little to no taxonomic

260    information (i.e., a virus sequence classified as only belonging to the *Riboviria*). If an RdRp-like

261    region was detected in an unclassified virus-like transcript we further assessed whether it could

262     be plant-associated (see Phylogenetic analysis of identified viruses). The remaining non-plant-

263     associated virus transcripts were largely classified within the *Orthomyxoviridae* (vertebrate

264     associated) (25%), *Rhabdoviridae* (invertebrate associated) (17%), *Partitiviridae* (fungus

265     associated) (10%), *Mimiviridae* (amoeboid associated) (10%) and *Adenoviridae* (vertebrate

266     associated) (7%) and excluded from the remainder of this study. These sequences are

267     discussed in more detail in the section on "Presence of contaminants in sequencing libraries"

268     below. Although some of these viruses could represent plant infection it remains challenging to

269     discern and we, therefore, made the conservative decision to remove them from the analysis.

270     We detected transcripts closely associated with viruses containing single and double-stranded

271     DNA and RNA genomes. The majority of virus-like sequences belonged to families with

272     +ssRNA genomes (61%) or reverse-transcribing dsDNA viruses (22%) (Figure 1). The +ssRNA

273     virus transcripts were predominately classified within the *Betaflexiviridae* (30%), *Potyviridae*

274     (19%), *Secoviridae* (16%) and *Alphaflexiviridae* (10%) (SI Table 3). Negative-sense single-

275     stranded RNA (-ssRNA) virus transcripts were classified within the *Aspiviridae* (0.04%),

276     *Rhabdoviridae* (6%) and *Tospoviridae* (3%) (*Phenuiviridae* and *Yueviridae* transcripts were later

277     detected in the unclassified virus-like transcripts) (SI Table 3). dsDNA virus transcripts with

278     sequence similarities to the *Phycodnaviridae* were detected across the algae samples. These

279     phycodna-like virus transcripts frequently encoded the chitinase and DNA ligase genes which

280     are homologous to those in distantly related host organisms including fungi and bacteria. Due to

281     the difficulties discerning whether these transcripts represent *Phycodnaviridae* sequences or

282     contamination, we excluded all phycodnavirus-related sequences. All remaining dsDNA viruses

283     were exclusively reverse-transcribing viruses from the *Caulimoviridae*. We failed to detect any

284     sequences that shared homology with several plant virus families including *Reoviridae*,

285     *Nanoviridae* and *Fimoviridae* (although see the Discussion for caveats).

286

**Figure 1.** (A) Phylogram of virus composition across the One Thousand Plant Transcriptomes Initiative (1KP) samples. Plant-associated virus abundance was summarised for each plant species and normalised using a Box-Cox transformation. The height of each bar represents the percentage of virus reads detected in each plant species (after the removal of host reads). Plant clades are labelled and differentiated by shades of grey. The 1KP ASTRAL tree was used as the basis for this tree (30). Clade and abundance annotations were added using the Interactive Tree of Life (iTOL) web-based tool (109). (B) The phytovirosphere across the Plantae and *Phaeophyceae*. A schematic tree of the evolution of major plant groups. Each bar represents the number of total viruses formally or likely associated with each host group and is coloured by virus genome composition. The total number of viruses for each plant group plus those found in this study is also shown at the end of each bar. The Virus-Host (38) and NCBI virus databases (110) combined with literature searches were used to obtain virus counts. Lineage branches are not drawn to scale. To our know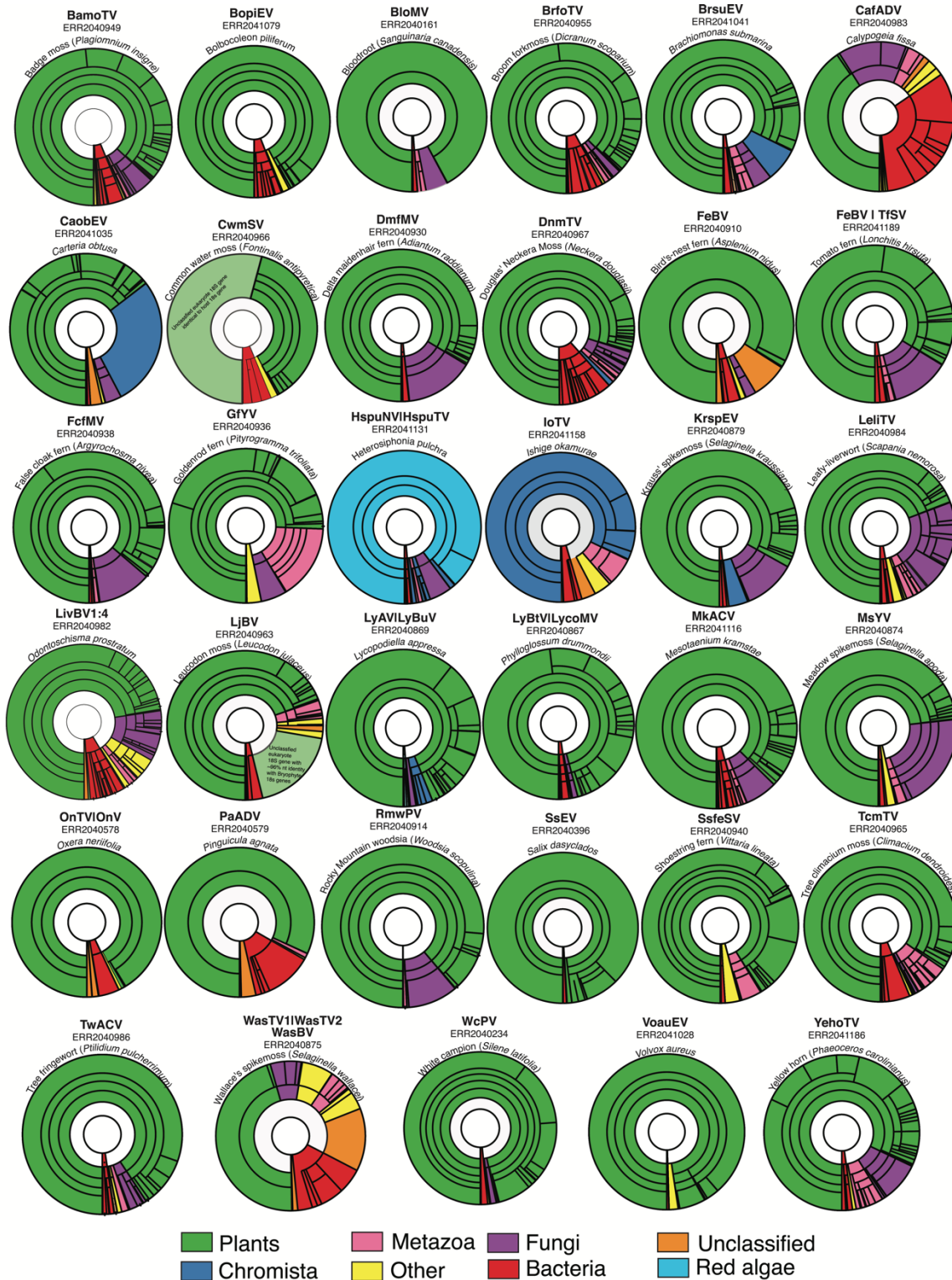ledge, no viruses have been found in the Glaucophytes. Plant and algae images were obtained from BioRender.com or drawn in Adobe Illustrator

14

301    (https://www.adobe.com). *Transcriptome scaffolds from libraries belonging to these host

302    groups shared homology to virus RdRps and were partially analysed but not assembled or

303    deposited to GenBank (17).

304    There was a large range of total viral abundance in each library ($5\times10^{-6}$% – 31% reads after

305    host-associated reads were removed). Viruses with +ssRNA genomes accounted for the vast

306    majority (99.8%) of virus abundance detected (Figure 1, SI Table 3). As expected, virus

307    discovery was concentrated in the flowering plants (angiosperms), which have the highest

308    number of previously classified viruses. For instance, plant virus-like sequences were frequently

309    discovered in the core eudicots and monocots (i.e., 73% of libraries in which plant virus

310    transcripts were found). The detection rate of plant viruses was highest in the most basal

311    angiosperms (57%) and monocots (50%). No significant difference in virus abundance was

312    observed between sequencing platforms (Genome Analyzer II and Illumina HiSeq 2000;

313    p=0.327).

**3.2 Presence of contaminants in sequencing libraries**

315    The bacterial, fungal and insect species that live in or on plant tissues are commonly sampled

316    within plant sequencing libraries (31), although contamination from other plants is also a

317    possibility during sample preparation or sequencing. To quantify the extent of library

318    contamination we used the KMA and CCMetagen tools (Figure 2). Among the libraries analysed

319    (n = 95), bacteria were consistently detected representing a median of 1.5% of total abundance

320    (range 0.01%-33%). A median of 2% of library abundance was associated with fungi sequences

321    (range 0%-53%). Arthropods and chordates were also commonly detected across libraries

322    (found in 87 and 89 libraries, respectively) but at lower abundance (median 0.15%, range 0%-

323    11.4%). The presence of chordate associated reads is likely attributed to various routes of

324    sample contamination (e.g., faeces) or during sample processing and sequencing.

15

**Figure 2. Taxonomic assignments of reads in select One Thousand Plant Transcriptomes**

**Initiative (1KP) libraries.** Each Krona graph illustrates the relative abundance of taxa in a

328    metatranscriptome at varying taxonomic levels. For clarity, a maximum depth of five taxonomic

329    levels was chosen for each graph. The library Sequence Read Archive accession number, host

330    species, and the corresponding virus of interest are annotated above each graph. Segments are

331    highlighted based upon the species taxonomic grouping (plants = green, Chromista = blue,

332    unclassified = orange, bacteria = red, metazoa = pink, fungi = purple, red algae = light blue,

333    other = yellow). Here "plants" encompasses the Viridiplantae. Reads without any match in the nt

334    database are not shown.

335    The detection of four vertebrate associated viruses across several libraries provided further

336    evidence of library contamination. Sequences belonging to these viruses - *Influenza A virus* (16

337    libraries), *Human mastadenovirus C* (30 libraries), *Human immunodeficiency virus* (15 libraries)

338    and *Parainfluenza virus 5* (3 libraries) – were present at low abundance and showed little

339    genetic variation between libraries. Notably, chordate-associated reads were only present in

340    66% of libraries in which these viruses were found. The failure to consistently detect potential

341    hosts for these viruses suggests contamination during sequencing. The four vertebrate

342    associated viruses were largely absent in libraries in which novel plant-associated viruses were

343    discovered, except for the *Larix speciosa*, *Brachiomonas submarina*, *Climacium dendroides*,

344    *Silene latifolia and Oxera neriifolia* transcriptomes.

345    In addition, the 1KP compared all assembled sequences to a reference set of nuclear 18S

346    ribosomal RNA sequences from the SILVA small subunit rRNA database using BLASTn (31,

347    68). Where a sample had several alignments to any other plant sequences outside of the

348    expected source family the sample was described as having "worrisome contamination" (31).

349    This applied to eleven plant libraries in which novel viruses were identified. Below, we discuss

350    library contaminates from viewpoint of virus-host associations.

351    **3.3 Phylogenetic analysis of identified viruses**

352    To infer phylogenetic relationships between identified viruses, order and family-level

353    phylogenetic trees were estimated using the highly conserved viral region that comprises the

354    RdRp. In total, we assembled 104 RdRp contigs that likely represent novel virus species, of

355    which 41 were considered as unclassified or non-plant associated due to their similarities to

356    virus groups known to infect non-plant hosts (SI Table 4). Further analysis of these contigs

357    revealed that they are likely plant-associated.

358    **3.3.1 Positive-sense single-stranded RNA ((+)ssRNA) viruses**
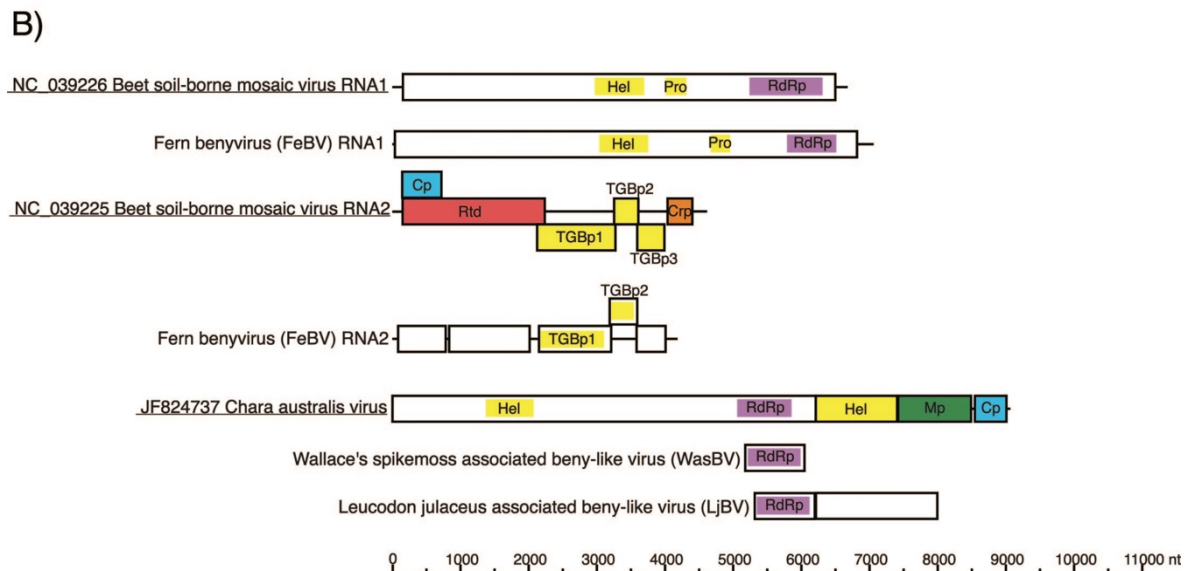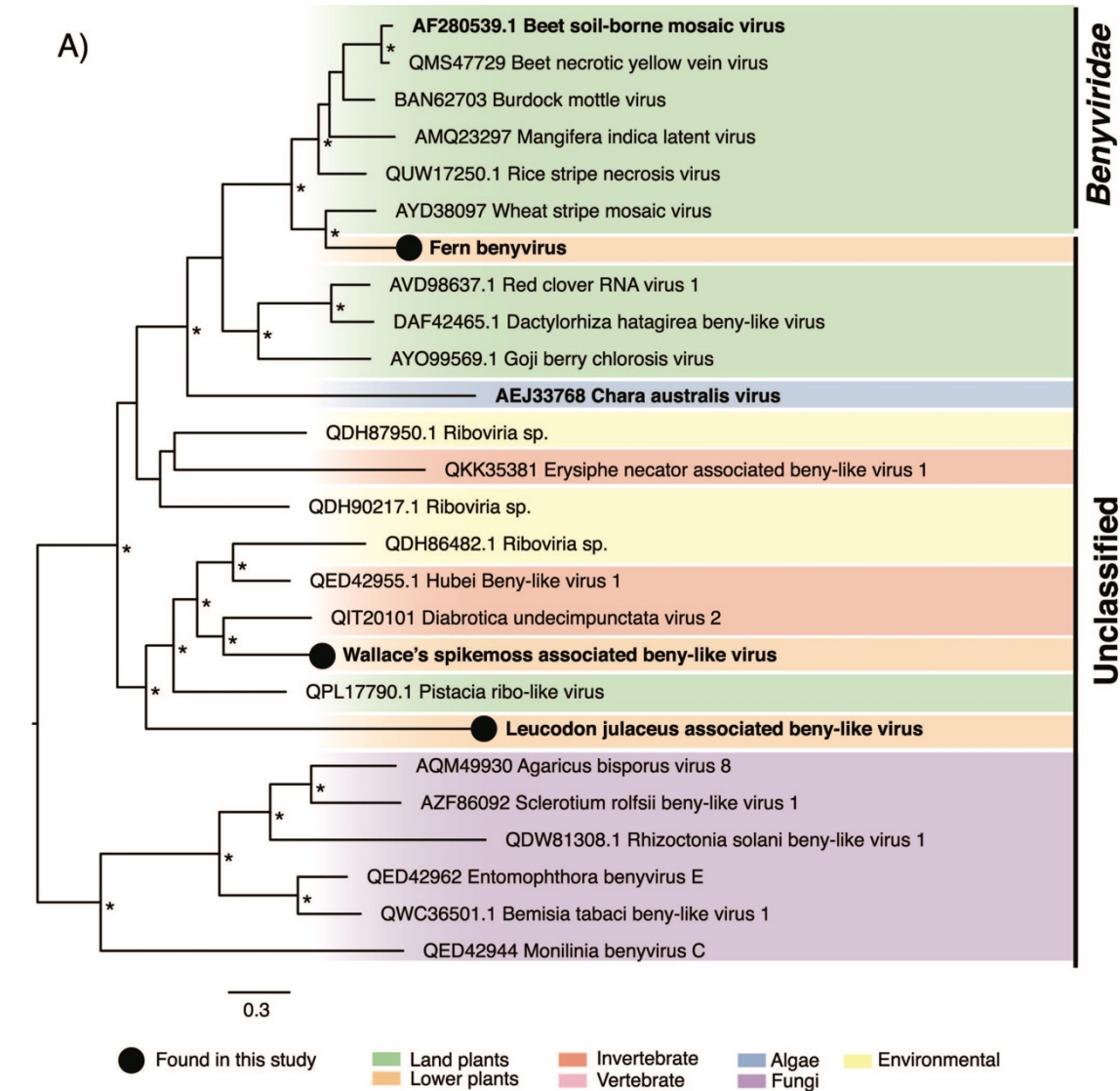
359    ***Hepelivirales***

360    ***Benyviridae.*** We identified three beny-like sequences that to our knowledge represent the first

361    benyvirid found in lower plants. The first sequence, tentatively named *Fern benyvirus* (FeBV),

362    was found in both the bird's-nest fern (*Asplenium nidus*) and tomato fern (*Lonchitis hirsuta*).

363    Together with *Wheat stripe mosaic virus*, FeBV represents a well-supported clade separate

364    from the remaining plant benyviruses (Figure 3).

365    The triple gene block (TGB) is a hallmark gene module of the *Benyviridae* among several other

366    virus families in the class *Alsuviricetes* (69)). In both fern libraries, proteins resembling the TGB

367    were assembled (Figure 3). The TGB proteins shared ~34% amino acid identity with the TGB

368    protein of other benyvirids. To our knowledge, this is the first TGB protein found outside of

369    flowering plants. Phylogenetic analysis placed the TGB1 protein of FeBV basal to the

370    *Benyviridae* (SI figure 1).

371    Two additional beny-like viruses, named here *Leucodon julaceus associated beny-like virus*

372    (LjBV) and *Wallace's spikemoss associated beny-like virus* (WasBV) were assembled. LjBV and

373    WasBV cluster with unclassified algae, invertebrates, fungi and soil-derived viruses forming a

374    group basal to all plant benyvirids and potentially constitute a novel virus group (Figure 3). LjBV

18

375    contains a second open reading frame (ORF) with no detectable homology to known sequences

376    (Figure 3).

377    Due to the phylogenetic placement of LjBV and WasBV close to viruses infecting distant hosts

378    (e.g., invertebrates and fungi), we investigated the potential of contamination from other

379    eukaryotes as the source of these viruses. Of note, the Wallace's spikemoss metatranscriptome

380    contained reads that matched various fungi orders (7% of all reads) as well as those matching

381    the plant-parasitic oomycete *Albugo laibachii* (7%) which makes inferring virus-host

382    relationships challenging (Figure 2). Reads belonging to various fungi species accounted for

383    10% of the bird's-nest fern transcriptome and 12% of the tomato fern transcriptome (Figure 2).

384    Despite the presence of fungi-associated reads, the phylogenetic position of FeBV suggests

385    that FeBV is likely plant-associated (Figure 3). No concerning contaminants were detected in

386    the *Leucodon julaceus* transcriptome.

A) Phylogenetic tree of Benyviridae and unclassified beny-like viruses

*Benyviridae:*
- AF280539.1 Beet soil-borne mosaic virus
- QMS47729 Beet necrotic yellow vein virus
- BAN62703 Burdock mottle virus
- AMQ23297 Mangifera indica latent virus
- QUW17250.1 Rice stripe necrosis virus
- AYD38097 Wheat stripe mosaic virus
- Fern benyvirus
- AVD98637.1 Red clover RNA virus 1
- DAF42465.1 Dactylorhiza hatagirea beny-like virus
- AYO99569.1 Goji berry chlorosis virus

*Unclassified:*
- AEJ33768 Chara australis virus
- QDH87950.1 Riboviria sp.
- QKK35381 Erysiphe necator associated beny-like virus 1
- QDH90217.1 Riboviria sp.
- QDH86482.1 Riboviria sp.
- QED42955.1 Hubei Beny-like virus 1
- QIT20101 Diabrotica undecimpunctata virus 2
- Wallace's spikemoss associated beny-like virus
- QPL17790.1 Pistacia ribo-like virus
- Leucodon julaceus associated beny-like virus
- AQM49930 Agaricus bisporus virus 8
- AZF86092 Sclerotium rolfsii beny-like virus 1
- QDW81308.1 Rhizoctonia solani beny-like virus 1
- QED42962 Entomophthora benyvirus E
- QWC36501.1 Bemisia tabaci beny-like virus 1
- QED42944 Monilinia benyvirus C

Scale: 0.3

Legend:
● Found in this study
- Land plants
- Lower plants
- Invertebrate
- Vertebrate
- Algae
- Fungi
- Environmental

B) Genome organization

- NC_039226 Beet soil-borne mosaic virus RNA1 — Hel, Pro, RdRp
- Fern benyvirus (FeBV) RNA1 — Hel, Pro, RdRp
- NC_039225 Beet soil-borne mosaic virus RNA2 — Cp, Rtd, TGBp1, TGBp2, TGBp3, Crp
- Fern benyvirus (FeBV) RNA2 — TGBp1, TGBp2
- JF824737 Chara australis virus — Hel, RdRp, Hel, Mp, Cp
- Wallace's spikemoss associated beny-like virus (WasBV) — RdRp
- Leucodon julaceus associated beny-like virus (LjBV) — RdRp

Scale: 0, 1000, 2000, 3000, 4000, 5000, 6000, 7000, 8000, 9000, 10000, 11000 nt

387

20

388     **Figure 3.** (A) Phylogenetic relationships of the beny-like viruses identified in this study. ML

389     phylogenetic tree based on the RNA-1 replicase protein shows the topological position of virus-

390     like sequences discovered in this study (black circles) in the context of their closest relatives.

391     Branches are highlighted to represent host clade (land plants = green, lower plants = orange,

392     invertebrate = red, vertebrate = pink, algae = blue, fungi = purple, yellow = environmental,

393     Chromista = light blue, red algae = dark green). Here "Land plants" encompasses both

394     angiosperms and gymnosperms while "Lower plants" includes the bryophytes, lycophytes, and

395     ferns. All branches are scaled to the number of amino acid substitutions per site and trees were

396     mid-point rooted for clarity only. An asterisk indicates node support of >70% bootstrap support.

397     Tip labels are bolded when the genome structure is shown on the right. (B) Genomic

398     organization of the beny-like virus sequences identified in this study and representative species

399     used in the phylogeny. Beet soil-borne mosaic virus RNA three and four are not pictured here.

400     The data underlying this figure and definitions of acronyms used are presented in SI Table 5.

401     ***Tymovirales***

402     **_Betaflexiviridae._** We identified 18 virus sequences that fell within the order *Tymovirales*. Four

403     virus transcripts were associated with the *Betaflexiviridae*. The first, named *Sea beet*

404     *betaflexivirus* (SbBV) clusters with *Agapanthus virus A*, an unclassified betaflexivirus (Figure 4).

405     The remaining sequences denoted *Iranian poppy betaflexivirus* (IpBV), *Linum macraei*

406     *betaflexivirus* (LimBV) and *Lycopod associated betaflexivirus* (LyBtV) resemble capilloviruses.

407     Notably, LyBtV may extend the known host range of the *Betaflexiviridae* from angiosperms to

408     lower plants. All sequences phylogenetically cluster with known capilloviruses and potentially

409     represent novel virus species (Figure 4). The *Phylloglossum drummondii* library in which LyBtV

410     was assembled had contamination from lycopod and dicot species (Figure 2). As the majority of

21

411    plant-associated reads were assigned to lycophytes (50%), LyBtV has been tentatively assigned

412    to this group.

**Tymovirales**

0  1000  2000  3000  4000  5000  6000  7000  8000  10000  12000 nt

RdRp   Cp
Mp   Hel
Mtr   Pro

Tymovirus
NP_047920 Erysimum latent virus
AVW89220 Passion fruit yellow mosaic virus
YP_004464924 Asclepias asymptomatic virus
QHB15484 Insect-associated tymovirus 1
ACX94288 Mertensia leaf curl virus
NP_044328 Kennedya yellow mosaic virus
YP_001285472 Okra mosaic virus
AFC95826 Watercress white vein virus
AAB92649 Turnip yellow mosaic virus
DAF42470 Glehnia littoralis marafivirus
Bloodroot marafivirus
Oxera neriifolia tymo-like virus
Marafivirus
Maculavirus
NP_542612 Grapevine fleck virus
QSC42383 Mutum virus
AYW01753 Ek Balam virus
AFT60444 Culex originated Tymoviridae-like virus
AGE84283 Culex originated Tymoviridae-like virus
YP_009272813 Bat tymo-like virus
QDW81035 Cattle tick tymovirus-like virus 1
AYV61005 Guarapuava tymovirus-like 2
QJX60166 Antioquia tymovirus-like 1
AYV61001 Guarapuava tymovirus-like 1
AYV61007 Guarapuava tymovirus-like 3
AKQ48574 Bee Macula-like virus
YP_009160324 Bee Macula-like virus
YP_009159826 Varroa Tymo-like virus
YP_009553654 Grapevine associated tymo-like virus
QED94454 Prunus yellow spot-associated virus
QQG34678 Liriodendron tulipifera tymovirus 1
Aulacomnium heterostichum tymo-like virus
QQG34657 Sinomenium acutum tymovirus 1
Tree climacium moss associated tymo-like virus
Leafy-liverwort associated tymo-like virus
Broom forkmoss associated tymo-like virus
AMN92730 Fusarium graminearum mycotymovirus 1
QQG34653 Triticum polonicum mycotymovirus 1
QDF82047 Sclerotinia sclerotiorum mycotymovirus 1
QOE77941 Sclerotinia sclerotiorum mycotymovirus 2
AWY10995 Sclerotinia sclerotiorum tymo-like RNA virus 4
AWY10994 Sclerotinia sclerotiorum tymo-like RNA virus 3
QJQ28892 Botrytis cinerea mycotymovirus 1
Douglas Neckera Moss associated tymo-like virus
Wallace spikemoss associated tymo-like virus 1
QDH90244 Riboviria sp.
QDH87807 Riboviria sp.
Ishige okamurae tymovirus
QOX06053 Lentinula edodes tymo-like virus 1
AQM49940 Agaricus bisporus virus 12
QDH91497 Riboviria sp.
AMM45289 Rhizoctonia solani positive-strand RNA virus 1
Badge moss associated tymo-like virus
Wallace spikemoss associated tymo-like virus 2
Yellow horn associated tymo-like virus
Alphaflexiviridae

*Tymoviridae*

Q6PLS1 Pear black necrotic leaf spot virus
AYN44466 Citrus tatter leaf virus
AAP80757 Apple stem grooving virus
AOR84129 Apple stem grooving virus
QQG34607 Breadfruit capillovirus 1
YP_009268859 Yacon virus A
QQG34669 Iris germanica capillovirus 1
QQG34585 Avellana capillovirus 1
Lycopod associated betaflexivirus
QQG34587 Rhodiola betaflexivirus 1
AWK77906 Hobart betaflexivirus 1
Linum macraei betaflexivirus
QQG34583 Silene betaflexivirus 1
YP_009551984 Mume virus A
QDK55539 Cherry virus A
QGW49048 Loquat virus A
QQG34645 Suaeda fruticosa betaflexivirus 1
Iranian poppy betaflexivirus
AGU09694 Apple stem pitting virus
ADT91615 Apricot latent virus
QED43569 Garlic latent virus
AAY18409 Blueberry scorch virus
BAU25805 Cherry necrotic rusty mottle virus
AHJ80314 Cherry twisted leaf associated virus
AFA43536 Citrus leaf blotch virus
AWK77900 Darwin betaflexivirus 1
QVY19268 Agapanthus virus A
Sea beet betaflexivirus
YP_009373228 Agave tequilana leaf virus
ABG37965 Grapevine virus A
AFV39891 Potato virus T
YP_009103996 Carrot Ch virus 2
YP_002308565 Peach mosaic virus
CAA68080 Apple chlorotic leaf spot virus
QUE49151 Sclerotinia sclerotiorum gammaflexivirus 1
AWV68780 Botrytis virus F
AVD68668 Entoleuca gammaflexivirus 2
QDO72745 Pistacia-associated flexivirus 1
AVD68667 Entoleuca gammaflexivirus 1
YP_009508363 Sclerotinia sclerotiorum deltaflexivirus 1
QKN22695 Erysiphe necator associated deltaflexivirus 4
ALM62223 Soybean leaf-associated mycoflexivirus 1
ANS13830 Fusarium graminearum deltaflexivirus 1
YP_009552771 Sclerotinia sclerotiorum deltaflexivirus 2
Calypogeia fissa associated deltaflexivirus
YP_009268715 Rhizoctonia solani flexivirus 1
Pinguicula agnata associated deltaflexivirus

*Betaflexiviridae*
*Gammaflexiviridae*
*Deltaflexiviridae*

0.4   ● Found in this study

Land plants   Chromista
Lower plants   Fungi
Green algae   Vertebrate
Invertebrate   Environmental

413

23

414 **Figure 4.** Left: Phylogenetic relationships of the viruses within the order *Tymovirales*. ML

415 phylogenetic tree based on the replication protein shows the topological position of virus-like

416 sequences discovered in this study (black circles) in the context of their closest relatives. See

417 Figure 3 for the colour scheme. All branches are scaled to the number of amino acid

418 substitutions per site and trees were mid-point rooted for clarity only. An asterisk indicates node

419 support of >70% bootstrap support. Tip labels are bolded when the genome structure is shown

420 on the right. Right: Genomic organization of the virus sequences identified in this study and

421 representative species used in the phylogeny.

422 ***Tymoviridae.*** We identified 12 virus-like sequences that clustered within the *Tymoviridae* and

423 related viruses. *Ishige okamurae associated tymo-like virus* (IoTV) was detected in the brown

424 alga *Ishige okamurae* and likely represents the first virus in the order *Tymovirales* from brown

425 algae. IoTV, along with ten sequences assembled from hornworts, liverworts and bryophytes

426 grouped with tymo-like viruses from fungus and environmental samples (Figure 4). It is

427 uncertain whether the true hosts of the novel tymo-like viruses discovered here are plants.

428 Fungi contaminates were detected across these libraries but varied in abundance (range 1%-

429 21%, mean = 6%). Despite their clustering with mycotymoviruses, *Broom forkmoss associated*

430 *tymo-like virus* (BrfoTV) and *Tree climacium moss associated tymo-like virus* (TcmTV) were

431 assembled from libraries with ~1% fungal reads, highlighting the inherent difficulties in host-

432 virus assignment. Importantly, <1% of reads in *Ishige okamurae* transcriptome belonged to

433 species of fungi (Figure 2).

434 We assembled two tymo-like virus sequences denoted *Oxera neriifolia tymo-like virus* (OnTV)

435 and *Bloodroot marafivirus* (BloMV). BloMV and OnTV grouped with the unclassified *Glehnia*

436 *littoralis marafivirus* (Figure 4). Marafiviruses and tymoviruses are commonly distinguished from

437 each other based upon a highly conserved 16 nucleotide (nt) sequence known as the "tymobox"

24

438    [GAGUCUGAAUUGCUUC] in tymoviruses and the "marafibox" [CA(G/A)GGUGAAUUGCUUC]

439    in marafiviruses (70, 71). While these two novel viruses cluster together phylogenetically, they

440    differ in terms of genome structure and motifs. A "marafibox" like sequence appears to be

441    present in BloMV (CAACGCGAAUUGCUUU) (5606-5621 nt) albeit differing by several

442    residues. This finding, combined with the BloMV genome likely consisting of a single large ORF,

443    supports the assignment of BloMV as a *Marafivirus*. OnTV, like members of the *Tymovirus*

444    genera, contains both a second ORF – likely encoding a coat protein (CP) – and a tymobox

445    (1493-1508 nt) (Figure 4). Phylogenetic analysis of the coat protein sequence places OnTV and

446    BloMV in a clade with macula- and marafi-like viruses (SI Figure 1).

447    **Deltaflexiviridae.** We assembled two sequences that share similarities to members of the

448    mycotymovirus family, *Deltaflexiviridae*. The first sequence was detected in the liverwort

449    *Calypogeia fissa*, tentatively named *Calypogeia fissa associated deltaflexivirus* (CafADV) and

450    appeared distantly related to delta- and gammaflexiviruses. A second related partial sequence,

451    named here *Pinguicula agnata virus* (PaV), shared 32% amino acid identity with mycoflexivirus,

452    *Botrytis virus F*. In a phylogenetic analysis with members of the *Tymovirales,* CafADV and

453    PaAGV are placed with the deltaflexivirids (Figure 4).

454    It is unclear whether the source of these virus sequences is from plants or contamination from

455    other eukaryotes. The *C. fissa* library contained numerous contaminates including algae, fungi

456    and bacteria representing 1%, 15% and 33% of total reads respectively, which make discerning

457    the host association for CafAV challenging (Figure 2). Interestingly, no fungi-associated reads

458    were found in the *P. agnata* library suggesting a potential plant origin (Figure 2).

459    **Picornavirales**

460    **_Secoviridae._** We identified four sequences that shared similarities to members of the

461    _Secoviridae_ denoted _Common water moss secovirus_ (CwmSV), _Salix dasyclados secovirus_

462    (SadSV), _Tomato fern secovirus_ (TfSV) and _Shostring fern secovirus_ (SfSV). CwmSV, TfSV and

463    SfSV cluster within the nepoviruses and likely represented the first seco-like virus detected in

464    the bryophytes and ferns (Figure 5).

**Figure 5.** (A) Left: Phylogenetic relationships of the viruses identified within the virus family

*Secoviridae*. ML phylogenetic trees based on the Pro-pol region show the topological position of

virus-like sequences discovered in this study (black circles) in the context of their closest

relatives. Right: Genomic organization of the seco-like sequences identified in this study and

470    representative species used in the phylogeny. (B) Multiple amino acid sequence alignment of

471    the 30K movement protein "LPL" motifs which are highly conserved throughout the nepoviruses.

472    (C) Phylogenetic relationships of the Nepovirus 30K movement proteins. D) Phylogenetic

473    relationships of the Nepovirus coat proteins. For all trees, branches are scaled to the number of

474    amino acid substitutions per site and trees were mid-point rooted for clarity only. An asterisk

475    indicates node support of >70% bootstrap support. Tip labels are bolded when the genome

476    structure is shown on the right. See Figure 3 for the colour scheme. Viruses discovered in this

477    study are signified using a black circle on the tree tip.

478    A putative RNA2 ORF was assembled for the three nepovirus-like sequences each containing a

479    complete CP (Figure 5). The CPs fall within the nepovirus subgroup C (Figure 5D). While a

480    movement protein (MP) domain was not formally detected, we predict that the region upstream

481    of the CP contains a putative movement-like protein. For CwmSV, this region (amino acid

482    position 312-883) displayed sequence homology to the MP of *Blackcurrant reversion virus* (E-

483    value: 5.42e-86, amino acid identity: 46%). Both TfSV and SfSV displayed similar levels of

484    homology in this region. We detected the LPL motif which is commonly found in nepovirus MPs

485    in all three viruses (Figure 5B). Phylogenetic analysis of the putative MPs placed these viruses

486    with *Blackcurrant reversion virus* in the genera Nepovirus (Figure 5C).

487    We found little evidence that these viruses were detected due to contamination by land plants or

488    other eukaryotes. The *F. antipyretica* transcriptome was composed of reads closely related to a

489    feather moss belonging to the order Hypnales to which *F. antipyretica* is also found.

490    Furthermore, a large proportion of reads were assigned to an uncultured eukaryote 18S rRNA

491    gene (54%) (HG421124.1) that was identical to the *F. antipyretica* 18S rRNA (AF023714.1)

492    among other bryophyte 18S rRNA genes in a blastn search (e-value = 2e-102, nucleotide

493    identity = 100%) (Figure 2). Fungi represented 12% of reads in the *L. hirsute* transcriptome.

28

494    Despite this, it is unlikely that TfSV is fungi-associated as no fungal contamination was detected

495    in the *Vittaria lineata* transcriptome in which the closely related SfSV sequence (amino acid

496    identity: 78%) was assembled (Figure 2).

497    **Lenarviricota**

498    **Mitoviridae.** We identified six virus sequences that cluster within the *Mitoviridae* - denoted

499    *Chinese swamp cypress mitovirus* (CscMV)*, Asian bayberry mitovirus* (AsbaMV)*, False cloak*

500    *ferns mitovirus* (FcfMV)*, Delta maidenhair fern mitovirus* (DmfMV) and *Lycopod associated*

501    *mitovirus* (LycoMV). The fern (FcfMV and DmfMV) and lycophyte (LycoMV) associated

502    sequences cluster with the fern *Azolla filiculoides mitovirus 1* and form a sister group to the

503    plant mitoviruses and non-retroviral endogenous RNA viral elements (NERVEs) (Figure 6) (16).

504    The gymnosperm associated sequences form a sister all the plant-associated mitoviruses and

505    NERVEs. KpcMV extends the known host range of plant mitoviruses from ferns to lycophytes.

506    Another mito-virus sequence was detected in the green alga *Bolbocoleon piliferum*, denoted

507    *Bolbocoleon piliferum mito-like virus* (BopiMV). BopiMV falls basal to the mitoviruses, distinct

508    from various unclassified mito-like viruses including the green algae associated *mito-like*

509    *picolinusvirus* (QOW97241) (Figure 6). All novel sequences show strong conservation of the

510    motifs characteristic of mitovirus RdRps (SI Figure 2) (72).

**Figure 6.** Left: Phylogenetic relationships of the viruses within the virus families *Narnaviridae* and *Mitoviridae*. ML phylogenetic trees based on the replication protein show the topological position of virus-like sequences discovered in this study (black circles) in the context of their

515    closest relatives. See Figure 3 for the colour scheme. Blue stars signify mitovirus sequences

516    identified in (16). Red stars signify non-retroviral endogenous RNA viral elements (NERVEs). All

517    branches are scaled to the number of amino acid substitutions per site and trees were mid-point

518    rooted for clarity only. An asterisk indicates node support of >70% bootstrap support. Tip labels

519    are bolded when the genome structure is shown on the right. Right: Genomic organization of the

520    virus sequences identified in this study and representative species used in the phylogeny.

521    There is little evidence to suggest that these sequences are derived from a non-plant organism.

522    While the FcfMV and DmfMV libraries were contaminated with fungi, (12% and 15% of reads

523    respectively) fungi-associated reads were absent in the libraries of all other mitoviruses. As the

524    codon UGA encodes tryptophan (Trp) in fungal mitochondria this codon assignment is also

525    present in fungal mitoviruses (73-75). In contrast, the UGA codon in plant mitochondria is a stop

526    codon and hence absent from plant mitovirus sequences except as a stop codon (16). The

527    absence of internal UGA codons in these sequences is further evidence that these sequences

528    are plant-derived (16, 76). Although additional analyses are required, we found no evidence

529    through searches of the 1KP genome scaffolds and the WGS shotgun database that these

530    sequences are mitochondrial or nuclear NERVEs. Furthermore, CscMV, AsbaMV and LycoMV

531    contain complete RdRps and their UTRs share similarities in length and identity with plant

532    mitoviruses.

533    ***Narnaviridae*** A partial narna-like virus sequence was identified in the red alga *Heterosiphonia*

534    *pulchra* denoted Heterosiphonia pulchra narna-like virus (HspuNV). HspuNV clusters with

535    unclassified trypanosomatid associated viruses. While ~5% of reads in this library were

536    associated with fungi the phylogenetic position of this virus suggests that it is not derived from

537    fungi (Figure 2, Figure 6).

538    ***Tolivirales***

31

539    ***Tombusviridae.*** An alphacarmo-like virus tentatively named *Ihi tombusvirus* (IhiTV) was

540    identified in an Ihi (*Portulaca molokiniensis*) sample. IhiTV is phylogenetically positioned within

541    the alphacarmoviruses (SI Figure 3).

542    ***Patatavirales***

543    ***Potyviridae.*** We identified three virus-like sequences that clustered with plant viruses in the

544    family *Potyviridae – Traubia modesta potyvirus* (TramPV), *Common milkweed potyvirus*

545    (ComPV) and *Salt wort potyvirus* (SawPV). TramPV and ComPV shared 87% amino acid

546    identity and may therefore represent a single virus species. The potyvirus-like sequences

547    discovered all group with known potyviruses in a phylogenetic analysis of the Nib gene (SI

548    Figure 3).

549    ***Martellivirales***

550    ***Endornaviridae.*** Six alphaendorna-like virus sequences were detected in the four green algae

551    species and one lycophyte. The green algae and lycophyte associated alphaendorna-like

552    viruses termed *Bolbocoleon piliferum endorna-like virus* (BopiEV), *Volvox aureus endorna-like*

553    *virus* (VoauEV), *Carteria obtusa associated endorna-like virus* (CaobEV), *Brachiomonas*

554    *submarina associated endorna-like virus* (BrsuEV), *Staurastrum sebaldi endornavirus* (SsEV)

555    and *Krauss' spikemoss associated endorna-like virus* (KrspEV) fall across the

556    alphaendornavirus phylogeny and predominately cluster with algae and fungi associated viruses

557    (Figure 7). There was little evidence of algae (non-host) or fungi contamination in the *S. sebaldi*,

558    *B. piliferum* and *V. aureus* transcriptomes with <1% of all reads associated with these groups

559    (Figure 2). Non-green algae contaminants were present in the *C. obtus* (28%)*, B. submarina*

560    (7%) and *S. kraussiana* (4%) transcriptomes where fungi also appeared as a notable

561    contaminate representing 11% of all reads (Figure 2). To our knowledge, these sequences

32

562    represent the first endornavirus associated with charophytes, chlorophytes and lycophytes

563    although further work is needed to confirm the virus-host associations.

564    **Unclassified.** We identified a virus-like sequence in an *Oxera neriifolia* library, termed *Oxera*

565    *neriifolia associated virus*. The sequence, 10,214 nt in length contained four ORFs. The first

566    ORF (7,536 nt) comprised of a viral methyltransferase, helicase, and RNA polymerase while the

567    third ORF (513 nt) most closely resembled a CP. ORF one and ORF three shared the greatest

568    sequence similarity with *Culex pipiens associated Tunisia virus* (32% amino acid identity). The

569    second and fourth ORF share no homology to known viruses. The genome organization of OnV

570    is distinct from the other related plant virus families (Figure 7). OnV forms a distinct and well-

571    supported outgroup to the *Closteriviridae, Bromoviridae* and *Mayoviridae* families. As such, OnV

572    may potentially constitute a new virus family (Figure 7). We found little evidence that OnV was

573    detected due to contamination by other eukaryotes (Figure 2).

575 **Figure 7.** Left: Phylogenetic relationships of the (A) endorna-like and (B) unclassified (+)ssRNA

576 virus identified in this study. ML phylogenetic trees based on the replication protein show the

577 topological position of virus-like sequences discovered in this study in the context of those

578 obtained previously. Right: Genomic organization of the (A) endorna-like and (B) unclassified

579 ssRNA virus sequence identified in this study and representative species used in the phylogeny.

580 For all trees, branches are scaled to the number of amino acid substitutions per site and trees

581 were mid-point rooted for clarity only. An asterisk indicates node support of >70% bootstrap

582 support. Tip labels are bolded when the genome structure is shown on the right. See Figure 3

583 for the colour scheme. Viruses discovered in this study are signified using a black circle on the

584 tree tip.

585 **3.3.2 Negative-sense single-stranded RNA ((-)ssRNA) viruses**

586 ***Bunyavirales***

587 ***Phenuiviridae.*** A phenui-like virus sequence termed *Brown algae phenui-like virus* (BralPV)

588 was recovered from a *Sargassum thunbergii* transcriptome. The partial L segment clusters with

589 the unclassified plant and fungi viruses (Figure 8). No additional phenui-like virus segments

590 were recovered. There were no concerning contaminants were detected in the *S. thunbergii*

591 transcriptome (Figure 2).

592 ***Viridisbunyaviridae.*** We identified 16 bunya-like virus sequences from eight liverwort, moss

593 and lycophyte libraries. Three libraries contained multiple distinct putative complete and partial

594 viruses. The overall pairwise nucleotide identity was <70% between each sequence (Figure 8).

595 As such we consider each a different bunya-like viruses. These sequences group together to

596 form a novel clade of unclassified bunya-like viruses distantly related to oomycete, fungi, and

597 invertebrate viruses (Figure 8). Bunyaviruses typically comprise three segments (L, M, and S),

598    although only the L segment was recovered for these sequences. These sequences represent

599    the first plant-associated viruses that cluster near the unofficially named *Deltamycobunyaviridae*

600    (77) (Figure 8). As the complete coding sequences of the viruses discovered share <30% amino

601    acid identity to the nearest relatives in the *Deltamycobunyaviridae,* they may constitute a new

602    virus family. We tentatively name this virus family the *Viridisbunyaviridae, (Viridis* meaning

603    green, while bunya is derived from the virus *order Bunyavirales* in which this clade falls within).

604    There was no evidence suggesting that these sequences originated from non-plant

605    contaminants. Host assignment was unclear for *Lycopod associated bunyavirus* and *Liverwort*

606    *associated bunyavirus 1:4* as reads belonging to several lycophyte and liverwort species,

607    respectively, were found in the source transcriptomes (Figure 2).

609     **Figure 8.** Phylogenetic relationships of the viruses (A) Left: A phylogeny depicting a novel clade

610     of viruses related to the *Deltamycobunyaviridae* in the context of the *Bunyavirales*. Right:

611     Genomic organization of the virus sequences identified in this study and representative species

612     used in the phylogeny (B) Percent identity matrix of the novel bunya-like viruses. Identity scores

613     are calculated from an alignment of the RdRp protein coding sequence. For clarity, the 100%

614     identity along the diagonal has been removed. Where sequence identity is >= 60% the value is

615     shown. (C) Left: A phylogeny depicting the phenui-like virus identified in this study in the context

616     of the *Phenuiviridae*. Right: Genomic organization of the virus sequences identified in this study

617     and representative species used in the phylogeny. For all trees, branches are scaled to the

618     number of amino acid substitutions per site and trees were mid-point rooted for clarity only. An

619     asterisk indicates node support of >70% bootstrap support. Tip labels are bolded when the

620     genome structure is shown on the right. See Figure 3 for the colour scheme. Viruses discovered

621     in this study are signified using a black circle on the tree tip.

622     ***Mononegavirales***

623     ***Rhabdoviridae.*** We identified seven sequences that clustered with plant viruses in the family

624     *Rhabdoviridae* denoted *Canadian violet rhabdovirus 1* (CvRV1), Canadian violet rhabdovirus 2

625     (CvRV2), *Common ivy rhabdovirus* (CoiRV) and *Indian pipe rhabdovirus* (InpRV), *Tree fern*

626     *varicosa-like virus* (TfVV), *Monoclea gottschei varicosa-like virus* (MgVV) and *Bug moss*

627     *associated rhabdo-like virus* (BmRV). Notably, TfVV and MgVV expand the host range of the

628     rhabdoviruses from angiosperms and gymnosperm to ferns and liverworts. RNA2 segments

629     were recovered for both viruses, TfVV RNA2 contained five genes while MgVV contained four

630     (Figure 9C). Two partial segments sharing similarities to the nucleocapsid (N) of *Black grass*

631     *varicosavirus-like virus* (YP_009130620.1) were found in the Indian pipe library and share 50%

632    amino acid identity. All sequences likely represent novel species within known plant infecting

633    genera (Figure 9C).

634    BmRV is a partial sequence (693 nt) most closely related to the unclassified *Hubei rhabdo-like*

635    *virus 2* (44% amino acid identity). Further evidence is needed to confirm BmRV as the first moss

636    rhabdovirus, but the relatively low proportion of contaminates in this library (3% algae and 3%

637    fungi) suggests that this virus is plant-associated (Figure 2). While 53% of reads in the MgVV

638    library were fungi associated the phylogenetic position of MgVV suggests it is derived from

639    plants (Figure 2, Figure 9C).

**Figure 9.** Left: Phylogenetic relationships of the viruses within the families, (A) *Aspiviridae*, (B) *Yue-* and *Qinviridae* and (C) *Rhabdoviridae.* Right: Genomic organization of the virus sequences identified in this study and representative species used in the phylogeny. For all trees, branches are scaled to the number of amino acid substitutions per site and trees were

645    mid-point rooted for clarity only. An asterisk indicates node support of >70% bootstrap support.

646    Tip labels are bolded when the genome structure is shown on the right. See Figure 3 for the

647    colour scheme. Viruses discovered in this study are signified using a black circle on the tree tip.

648    For trees (A) and (B), the RdRp motif C trimer of each sequence is shown in brackets at the end

649    of the tip label.

650    ***Serpentovirales***

651    **Aspiviridae.** We identified an aspi-like sequence termed *Nees' Pellia aspi-like virus* (NpAV). A

652    complete RNA1 segment (6989 nt) was assembled, although no other segments were

653    recovered (Figure 9A). NpAV most closely resembles *Rhizoctonia solani negative-stranded*

654    *virus 3* (amino acid identity: 22%) and falls basal to all the unclassified aspi-like viruses

655    including those found in fungi, invertebrates, and oomycetes. NpAV is the first aspi-like virus

656    identified in plants outside of the angiosperms and may constitute a novel virus group (Figure

657    9A). Notably, unlike the other aspiviruses that possess a SDD sequence in motif C of the RdRp

658    – a known signature for segmented negative-stranded RNA viruses – NpAV has a GDD

659    sequence (Figure 9A).

660    ***Goujianvirales***

661    **Yueviridae.** An yue-like virus sequence termed *Meadow spikemoss associated yue-like virus*

662    (MsYV) was found in the lycophyte *Selaginella apoda* and most closely resembles algae

663    associated *Bremia lactucae associated yuevirus-like virus 1* (amino acid identity: 26%).

664    Phylogenetic analysis supports the assignment of MsYV as the first plant yuevirus (Figure 9B)

665    A second partial yue-like virus sequence was detected in a *Pityrogramma trifoliata* library and

666    termed *Goldenrod fern associated yue-like virus* (GfYV). GfYV falls with a group of oomycete

667    associated viruses. Consistent with the qin-like viruses, GfYV has an IDD (Ile-Asp-Asp)

41

668   sequence motif instead of the common GDD (Gly-Asp-Asp) in the catalytic core of its RdRp,

669   while MsYV contains SDD (Ser-Asp-Asp) in the same manner as many yue-like viruses (Figure

670   9B). The libraries from which GfYV and MsYV were assembled are contaminated with fungal

671   reads (5% and 21%, respectively) as such host assignment is made with caution (Figure 2).

672   Reads belonging to oomycetes were not found in either library.

673   **3.3.3 Double-stranded RNA (dsRNA) viruses**

674   ***Durnavirales***

675   ***Amalgaviridae.*** We detected five sequences that cluster with amalga-like viruses. *Lycopod*

676   *associated amalgavirus* (LycoAV) is a partial RdRp containing a sequence that falls basal to the

677   *Amalgaviridae* and represents the first amalga-like virus in the lycophytes (Figure 10). Three

678   amalga-like sequences were discovered in green and red algae transcriptomes and cluster with

679   *Diatom colony associated dsRNA virus 2* (Figure 10). As noted in the case of *Bryopsis*

680   *mitochondria-associated dsRNA virus* and several green algae associated viruses (78) when

681   translated into amino acids using the protozoan mitochondrial code, two overlapping ORFs are

682   present: the first, encoding a hypothetical protein, while the second, a replicase through a –1

683   ribosomal frameshift (79). For two of the amalga-like sequences identified in this study –

684   *Nucleotaenium eifelense virus* (NueiV) and *Rhodella violacea virus* (RhviV) – a similar structure

685   was observed but we were unable to identify any ribosomal frameshift motifs in either sequence

686   (Figure 10). Further work is needed to confirm if these sequences should be translated through

687   the mitochondrial genetic code.

688   A contig containing what appears to be a complete coding sequence (3259 nt) and RdRp motifs

689   was assembled in the *Woodsia scopulina* transcriptome and tentatively named *Rocky Mountain*

690   *woodsia associated virus* (RmwPV). The predicted RdRp region (918-1921 nt) of RmvPV

691    shares similarity to both partiti-like viruses (e.g., *Ustilaginoidea virens nonsegmented virus 2*,

692    26% aa identity) and the unclassified *Phytophthora infestans RNA virus 1* (42% aa identity)

693    which has been shown to likely constitutes a novel virus family (80). The resemblance RmwPV

694    shares with two seemingly distantly related virus groups suggest its position within the

695    *Durnavirales* should be treated with caution (Figure 10).

696    The transcriptome in which RmwPV was discovered is contaminated with fungal reads (10%

697    (Figure 2). If RmwPV was derived from fungal contaminates this could potentially explain the

698    phylogenetic placement of RmwPV (Figure 10). The *Lycopodiella appressa* transcriptome in

699    which LycoAV was discovered is contaminated by reads belonging to species across various

700    land plant groups. Reads belonging to land plants comprised 35% of plant-associated reads

701    while lycopod associated reads comprised 65% (Figure 2).

703 **Figure 10.** Left: Phylogenetic relationships of the viruses within the order *Durnavirales*. ML

704 phylogenetic trees based on the replication protein show the topological position of the virus-like

705 sequences discovered in this study (black circles) in the context of their closest relatives. See

706 Figure 3 for the colour scheme. Green stars are used to signify sequences that have been

707 translated using the protozoan mitochondrial genetic code. Red stars are used to signify

708 sequences for which a dsRNA3 coat protein-like segment has been described. All branches are

709 scaled to the number of amino acid substitutions per site and trees were mid-point rooted for

710 clarity only. An asterisk indicates node support of >70% bootstrap support. Tip labels are bolded

711 when the genome structure is shown on the right. Right: Genomic organization of the virus

712 sequences identified in this study and representative species used in the phylogeny.

713 ***Partitiviridae.*** We detected 14 sequences that share a resemblance with members of the

714 *Partitiviridae*. For each of these sequences, complete dsRNA1 and dsRNA2 segments were

715 recovered. Ten sequences were found in non-flowering plants and cluster within the

716 deltapartitiviruses. A clade within the deltpartitiviruses is known to encode a third segment

717 comprising of a divergent dsRNA2 full-length capsid protein with unknown function (Figure 10).

718 We identified dsRNA3 segments in related conifer associated sequences but not in those found

719 in moss and lycophyte libraries (Figure 10). Phylogenies estimated on the coding sequences of

720 dsRNA2 and dsRNA3 reveal essentially the same grouping which is largely consistent with the

721 host phylogeny (SI Figure 4). We extend the known host range of the deltapartitiviruses to

722 include liverworts, mosses, and lycophytes. The remaining sequences were found in eudicots

723 and cluster with known plant partitiviruses (Figure 10). The white campion was judged to have

724 contamination from ginseng and chickweed (31). However, the relatively low proportion of the

725 library these contaminates compose (<1%) suggests that it is unlikely these species are the

45

726     host of WcPV (Figure 2). There is no evidence that the other partiti-like sequences discovered

727     are derived from contaminates.

728     **Ghabrivirales**

729     **Chrysoviridae.** We identified two partial sequences that share resemblance with members of

730     the alphachrysoviruses denoted *Mesotaenium kramstae alphachyrso-like virus* (MkACV) and

731     *Tree fringewort alphachyrso-like virus* (TwACV) (Figure 11). A complete RNA2 segment was

732     recovered for MkACV which shared similarity with the p98 of various chrysoviruses (Figure 11).

733     The MkACV RNA2 segment did not contain the "PGDGXCXXHX" motif commonly found in this

734     protein (81). To our knowledge, these sequences represent the first chyrsoviruses in liverworts

735     and algae. While reads belonging to fungi were found in the libraries MkACV and TwACV were

736     assembled from, the phylogenetic positioning of the viruses suggest that they are plant-derived

737     (Figure 2, Figure 11).

739    **Figure 11.** Left: Phylogenetic relationships of the viruses within the order *Ghabrivirales*. (A) A

740    phylogeny of the *Chrysoviridae,* (B) an order level phylogeny. ML phylogenetic trees based on

741    the replication protein show the topological position of the virus-like sequences discovered in

742    this study (black circles) in the context of their closest relatives. See Figure 3 for the colour

743    scheme. Green stars are used to signify sequences that have been translated using the

744    protozoan mitochondrial genetic code. All branches are scaled to the number of amino acid

745    substitutions per site and trees were mid-point rooted for clarity only. An asterisk indicates node

746    support of >70% bootstrap support. Virus taxonomic names are labelled to the right. Right:

747    Genomic organization of the virus sequences identified in this study and representative species

748    used in the phylogeny.

749    ***Totiviridae.*** Thirteen sequences sharing similarities to toti-like viruses were discovered in eight

750    red and green algae transcriptomes. All sequences share less than 50% amino acid identity

751    across their coding sequence, as such we consider each a putative toti-like viruses. Among

752    these sequences four cluster with *Delisea pulchra totivirus IndA* (AMB17469.1) to form a red

753    alga associated clade basal to the totiviruses (Figure 11). *Gracilaria vermiculophylla toti-like*

754    *virus* (GrveTV) along with *Red algae totivirus 1* (BBZ90082) form a sister group to the protozoan

755    infecting leishmaniaviruses (Figure 11). The remaining sequences are phylogenetically

756    positioned across the tree of toti-like viruses, commonly occupying basal positions (Figure 11).

757    *Prasiola crispa* is contaminated by reads from the fungi, *Candida albicans*. *Prasiola crispa toti-*

758    *like virus* (PrcrTV), clusters with unclassified protist, fungi, invertebrate and algae viruses

759    including *Elkhorn sea moss toti-like virus* (EsmTV) (Figure 2, Figure 11). The Kappaphycus

760    alvarezii transcriptome in which EsmTV was found showed no evidence of contamination

761    suggesting that PrcrTV may also be derived from algae (Figure 2). The *Mazzaella japonica*

762    transcriptome in which *Red algae toti-like virus 2:3* (RedTV2/3) were discovered was

48

763 predominantly composed of reads associated with the red algae genera Chondrus. As >99% of

764 reads in this library belong to red algae species RedTV2 and RedTV3 have been assigned to

765 this group. The *Porphyridium purpureum* transcriptome is highly contaminated by reads

766 belonging to flowering plants and an unidentified cloning vector (M10197.1) (Figure 2). The

767 phylogenetic positioning of the viruses discovered from this transcriptome (*Porphyridium*

768 *purpureum toti-like virus* 1 & 2) point towards being derived from red algae rather than flowering

769 plants (Figure 11).

**3.4 Long-term virus-host evolutionary relationships**

771 To examine the frequency of cross-species transmission and co-divergence among plant

772 viruses, we estimated tanglegrams that depict pairs of rooted phylogenetic trees displaying the

773 evolutionary relationship between a virus family and their hosts. This revealed cross-species

774 transmission as the predominate evolutionary event predicted among all the RNA virus groups

775 analysed (median 65%, range 46%-79%) (Figure 12). Cross-species transmission was most

776 frequent in the *Betaflexiviridae* (79%) and the subfamily *Betarhabdovirinae* (79%). Virus-host

777 co-divergence (median 23%, range 14%-29%) and to a lesser extent duplication (i.e.,

778 speciation) (median 4.6%, range 1.4%-24%) and extinction events (median 2.9%, range 0%-

779 11%) were detected across plant virus families (Figure 12). Co-divergence was most frequently

780 predicted in the *Benyviridae* and *Tymoviridae* representing 29% and 26% of events respectively.

781 Importantly, however, the results of our co-phylogenetic analysis are undoubtedly influenced by

782 the sample of plant viruses and will likely change as the number of plant viruses identified

783 increases.

785    **Figure 12.** (A) Tanglegram of rooted phylogenetic trees for select virus groups and their hosts.

786    Lines and branches are coloured to represent host clade. The cophylo function implemented in

787    phytools (v0.7-80) was used to maximise the congruence between the host (left) and virus

788    (right) phylogenies. Supplementary Figure 5 provides the names of the hosts and viruses along

789    with additional tanglegrams for the *Secoviridae* and *Rhabdoviridae*. (B). Reconciliation analysis

790    of select virus groups. Barplots illustrate the range of the proportion of possible events and are

791    coloured by event type.

792

## 4. Discussion

794    Our ability to reconstruct the evolutionary history of plant viruses and understand the drivers of

795    their emergence has been constrained by inadequate sampling across the enormous, extant

796    diversity of plant species. Here, we provide a large-scale virus discovery project based on

797    mining transcriptomes from across the entire breadth of the plant kingdom. In doing so we have

798    identified 104 potentially novel virus species. We considerably expand upon the known host

799    range of 13 virus families to now include lower plants and expand a further four virus families to

800    include host associations with algae. We also find the first evidence of a movement protein with

801    a predicted molecular weight of ~30 kDa (herein referred to as a "30K MP") in a virus of non-

802    vascular plants. Collectively, this new knowledge advances our understanding of RNA virus

803    diversity across the Archaeplastida.

### 4.1 RNA viruses are widespread across lower plant lineages

805    To date, viral surveys in basal plant lineages (namely ferns, bryophytes and algae) have

806    revealed only the minimal occurrence of (+)RNA viruses (5, 17, 20, 78, 82, 83), supporting the

807    idea that RNA viromes in angiosperms evolved as they diversified during the Cretaceous (84).

51

808    However, our results potentially challenge this paradigm as we detected the first evidence of

809    sets of (+)ssRNA viruses in lower plants and algae, implying that these groups are associated

810    with older lineages of plants. Several of these viruses are deep branching and sit basal to

811    angiosperm infecting viruses (e.g., LycoAV and LyBuV) in phylogenetic trees. Other viruses

812    discovered here occupy ambiguous positions between established plant virus families (e.g.,

813    OnV) or cluster in large numbers to form novel plant-associated clades (e.g., the

814    *Viridisbunyaviridae* in the *Bunyavirales*). Benyviruses are typically transmitted by the root-

815    infecting plasmodiophorids *Polymyxa betae* and *Polymyxa graminis* (85, 86). The Phytomyxids

816    (plasmodiophorids and phagomyxids) are parasites of plants, diatoms, oomycetes and brown

817    algae and have been shown to demonstrate cross-kingdom host shifts (e.g., between

818    angiosperms and oomycetes) (87). As such the plasmodiophorids may be a vehicle for cross-

819    species transmission between aquatic protists and land plants (5). FeBV, a beny-like virus

820    identified in this study, formed a clade along with *Wheat stripe mosaic virus* distinct from

821    members of the genus *Benyvirus*. Deciphering the evolutionary history and mode of

822    transmission for the lower plant beny-like viruses will require further studies with particular

823    emphasis on these taxa. Interestingly, no plasmodiophorid-associated reads were detected in

824    any of the libraries from which we assembled a beny-like virus. LjBV and WasBV appear

825    distantly related to the benyviruses. These viruses group with a suite of unclassified viruses

826    assembled from a soil metatranscriptome study suggesting that, like the benyviruses, this larger

827    group of unclassified viruses may involve soil-borne parasites like the plasmodiophorids (88).

828    Our detection of tymovirid-like sequences in the lycophytes, bryophytes and brown algae

829    dramatically expands the known host range of the *Tymovirales*. Several of these viruses were

830    similar to unclassified Riboviria species assembled from a recent survey of common wild oat soil

831    rhizosphere and detritosphere (88) (Figure 4). The metatranscriptome of the sequenced soil

832    samples from the common wild oat study was largely composed of Viridiplantae, fungi,

833    Amoebozoa, protists, nematodes, and other eukaryotes. As such, using phylogenetic clusters to

834    infer host associations of our viruses remains challenging. Indeed, these viruses may result

835    from contamination from other eukaryotes (e.g., fungi or invertebrates) although we found no

836    consistent evidence among these viruses (Figure 2). Assuming these viruses are plant-

837    associated, their phylogenetic pattern suggests that they may have resulted from cross-kingdom

838    transmission events that frequent the Alsuviricetes.

839    The partial deltaflexi-like virus we detected in *P. agnata* (PaADV) is particularly noteworthy. The

840    deltaflexiviruses are only known to infect fungi, although no fungi associated reads were found

841    in the *P. agnata* metatranscriptome (Figure 2). The mycovirus families *Delta*- and

842    *Gammaflexiviridae* are thought to have been derived from the plant alpha- and betaflexivirids

843    through cross-species transmission (5, 89)). As such PaAGV could potentially represent an

844    intermediate between the plant and fungi flexiviruses or perhaps a more recent fungus to plant

845    transmission. As only a fragment of the polymerase gene was assembled for this virus future

846    work should confirm the presence of PaAGV and its phylogenetic position.

847    **4.2 The extension of the *Mitoviridae* to a lycophyte host**

848    Through the analysis of mitoviruses-like, non-retroviral endogenous RNA viral elements

849    (NERVEs), it was argued that the origin of plant mitovirus NERVEs was a single horizontal

850    transfer from a fungal mitovirus before the origin of vascular plants in the early Silurian, ~400

851    MYA (90). Evidence of contemporary mitoviruses in flowering plants and a fern have challenged

852    this view, suggesting that a lineage of plant rather than fungal mitoviruses are the immediate

853    ancestors of plant mitovirus NERVEs (16). Indeed, plant-to-fungus transmission would eliminate

854    code conflicts between fungi and plant mitochondrial genetic codes (76). Herein, we

855    demonstrate the existence of a lower plant-associated sister clade to the angiosperm

856    mitoviruses and NERVEs. This clade includes a clubmoss associated mitovirus, the most

857 primitive plant mitovirus sequence to date. This finding aligns with the estimation of the origin of

858 plant mitovirus NERVEs occurring as early as the evolution of the clubmoss (90). The recent

859 finding of mitoviruses in green algae – including BopiMV in this study – highlight the broad host

860 range of mitoviruses (78, 83). The phylogenetic position of these viruses and the absence of

861 NERVEs from these groups suggest that they are not the ancestors of land plant mitoviruses

862 and NERVEs.

863 **4.4 Establishment of a new virus family in the *Bunyavirales*: *Viridisbunyaviridae***

864 We identified 16 bunya-like viruses assembled from six non-vascular plant libraries including

865 liverwort, moss, and lycophyte species. These viruses form a novel clade within the

866 *Bunyavirales* and represent the first viruses in this order to be associated with lower plants. This

867 clade likely represents a novel virus family which we have tentatively named the

868 *Viridisbunyaviridae.* Several libraries contained up to five distinct viruses (each sharing <70%

869 nucleotide identity). Virus co-infections are frequently observed in plants and have been

870 reported in the closest relatives of these viruses, the *Deltamycobunyaviridae* (91, 92). As with

871 previous studies we were only able to recover the bunyavirus L segment (92, 93). Further

872 studies are needed to recover the missing small and medium-sized segments and to confirm the

873 presence of mixed infections in plants.

874 **4.5 Discovery of the first 30 kDa movement protein in non-vascular plants**

875 Through the discovery of lower plant-associated viruses, we have gained insights into how the

876 genome structure and composition of contemporary flowering plants viruses have evolved. The

877 detection of secovirid-like sequences in bryophytes and ferns represents the first occurrence of

878 plant secoviruses outside of angiosperms and the first evidence of a 30K MP homolog in non-

879 vascular plants. These proteins aid the cell-to-cell movement of viruses in plants. For example,

54

880     the MP of *Cucumber mosaic virus* increases the size exclusion limit of plasmodesmata allowing

881     virus particles to pass through cell walls (94). To date, homologs of 30K MP have only been

882     detected in plant viruses infecting angiosperms to the lycophytes (17, 95). Further work is

883     needed to confirm the presence and function of 30K MPs in viruses infecting the bryophytes and

884     other lower plants.

885     **4.6 Detection of Deltaparitivirus dsRNA3 segments in gymnosperms but not in non-**

886     **vascular plants**

887     Our discovery of six tri-segmented deltapartitivirus species provides insights into the evolution of

888     the deltapartitivirus dsRNA3 segment. dsRNA3 segments have been found in several alpha-

889     and deltapartitiviruses infecting flowering plants (96-99). These segments typically encode

890     seemingly full-length capsid protein or in the case of alphapartitivirus *Rosellinia necatrix*

891     *partitivirus 2*, a truncated version of the RdRp which may serve as an interfering RNA (100).

892     There is some debate as to the source of dsRNA3 segments, particularly whether they are

893     satellite viruses that co-opt the RdRp of the co-infecting helper viruses or that the additional

894     segment is a result of coinfection of two different plant partitiviruses and the second RdRp-

895     encoding segment is lost after the initial infection (101). For the first time, we find dsRNA3

896     segments in conifer associated viruses but not in those found in lower plants including

897     bryophytes and lycophytes. The absence of dsRNA3 in non-vascular plants means that it is

898     possible that this segment evolved after the divergence of vascular and non-vascular plants in

899     the Silurian period (102). It is possible that dsRNA3 segments exist for the non-vascular plant

900     infecting deltapartitiviruses but was not detected due to the large degree of divergence between

901     this segment and reference sequences (including those found in this study). However, the

902     dsRNA1 and dsRNA2 segments of the putative lower plant deltapartitiviruses shared >50% aa

903     identity with the tri-segmented deltapartitiviruses - well above the detection limit for tools such

904    as Diamond BlastX (37). dsRNA3 segments typically appear no more divergent than dsRNA2

905    segments therefore it is unlikely that we would be able to detect both the dsRNA1 and dsRNA2

906    segments without detecting dsRNA3. Further work is needed to confirm the presence of

907    deltapartitivirus dsRNA3 segments.

908    **4.7 Discovery of an unsegmented varicosavirus-like viruses in ferns and liverworts**

909    Finally, the recently discovered gymnosperm varicose-like *Pinus flexilis virus 1* in the family

910    *Rhabdoviridae* contains an unsegmented genome organisation that differs from the typical bi-

911    segmented structure of the varicosaviruses (25, 103). We find the bi-segmented structure in

912    varicosavirus-like viruses for the first time in ferns and liverworts (TfVV and MgVV) which

913    predate the gymnosperms.

914    **4.8 Caveats**

915    Importantly, the data generated under the 1KP were not explicitly created for virus discovery,

916    such that there are important caveats associated with the methods and metatranscriptomic data

917    mined for virus contigs. For instance, as axenic cultures are not a viable option in most

918    instances, the 1KP samples are commonly contaminated by nucleic acids belonging to

919    bacterial, fungal, and insect species. We addressed this by using a combination of host/virus

920    abundance measurements and phylogenetic analyses to improve the accuracy of virus-host

921    assignments. For most of the viruses described, phylogenetic placement within plant infecting

922    virus families strongly supports their association with plants. However, several of the viruses

923    found in algae and lower plants were associated with lineages known to infect invertebrates and

924    fungi or unclassified viruses recovered from environmental samples. The association between

925    the viruses of lower plants and algae with that of fungi and invertebrate viruses may reflect the

926    absence of algal and lower plant viruses in reference sequence databases. Experimental

927    confirmation is needed to formally assign the viruses discovered in this study to their hosts.

928    The average sequencing depth of the 1KP libraries was 1.99 gigabases of sequence per

929    sample (range 1.3-3.0), lower than many other virus discovery studies (6, 104, 105).

930    Sequencing depth has been shown to correlate with the ability to detect viruses present at low

931    abundance (106, 107). Further, a large proportion of the virus transcripts detected were from

932    viruses whose full-length genomic or subgenomic mRNAs were polyadenylated at the 3′ end (SI

933    Table 4, Figure 1). Although this was anticipated (i.e. the libraries generated by the 1KP

934    initiative were prepared from polyA+ RNA), it limited the detection of non-polyadenylated viruses

935    (e.g., dsRNA, dsDNA) and may have contributed to the lack of phycodnavirus sequences

936    detected in algae (107).

937    To reduce the computational burden of assembly, we attempted to remove host-associated

938    reads before contig assembly by mapping them to the host scaffolds provided by the 1KP

939    initiative. While this step reduces the occurrence of false-positive virus detection it also risks

940    removing virus reads, particularly reverse-transcribing plant viruses (108). While we frequently

941    detected transcripts associated with the reverse-transcribing family *Caulimoviridae,* no members

942    of the *Metaviridae* or *Pseudoviridae* were detected.

943    **Acknowledgements**

947

948

949 **Funding**

**References**

1. Anderson JT. 2016. Plant fitness in a rapidly changing world. New Phytologist 210:81-87.

2. Wren JD, Roossinck MJ, Nelson RS, Scheets K, Palmer MW, Melcher U. 2006. Plant virus biodiversity and ecology. PLoS Biology 4:80.

3. Mifsud JCO. 2020. Explorations of the plant virosphere.

4. Roossinck MJ, Martin DP, Roumagnac P. 2015. Plant virus metagenomics: advances in virus discovery. Phytopathology 105:716-727.

5. Dolja VV, Krupovic M, Koonin EV. 2020. Deep roots and splendid boughs of the global plant virome. Annual review of phytopathology 58:23-53.

6. Shates TM, Sun P, Malmstrom CM, Dominguez C, Mauck KE. 2019. Addressing Research Needs in the Field of Plant Virus Ecology by Defining Knowledge Gaps and Developing Wild Dicot Study Systems. Frontiers in Microbiology 9.

7. Dolja VV, Koonin EV. 2018. Metagenomics reshapes the concepts of RNA virus evolution by revealing extensive horizontal virus transfer. Virus Research 244:36-52.

8. Veliceasa D, Enünlü N, Kós PB, Köster S, Beuther E, Morgun B, Deshmukh SD, Lukács N. 2006. Searching for a new putative cryptic virus in Pinus sylvestris L. Virus Genes 32:177-186.

9. Sidharthan VK, Kalaivanan NS, Baranwal VK. 2021. Discovery of putative novel viruses in the transcriptomes of endangered plant species native to India and China. Gene 786:145626.

10. Han S, Karasev A, Ieki H, Iwanami T. 2002. Nucleotide sequence and taxonomy of Cycas necrotic stunt virus. Archives of virology 147:2207-2214.

978    11.   Yang S, Shan T, Wang Y, Yang J, Chen X, Xiao Y, You Z, He Y, Zhao M, Lu J. 2020.

979          Virome of riverside phytocommunity ecosystem of an ancient canal.

980    12.   Nibert ML, Pyle JD, Firth AE. 2016. A +1 ribosomal frameshifting motif prevalent among

981          plant amalgaviruses. Virology 498:201-208.

982    13.   Dawes C. 2016. Chapter 4 - Macroalgae Systematics, p 107-148. *In* Fleurence J, Levine

983          I (ed), Seaweed in Health and Disease Prevention doi:https://doi.org/10.1016/B978-0-

984          12-802772-1.00004-X. Academic Press, San Diego.

985    14.   Christenhusz MJ, Byng JW. 2016. The number of known plants species in the world and

986          its annual increase. Phytotaxa 261:201-217.

987    15.   Valverde RA, Sabanadzovic S. 2009. A novel plant virus with unique properties infecting

988          Japanese holly fern. Journal of General Virology 90:2542-2549.

989    16.   Nibert ML, Vong M, Fugate KK, Debat HJ. 2018. Evidence for contemporary plant

990          mitoviruses. Virology 518:14-24.

991    17.   Mushegian A, Shipunov A, Elena SF. 2016. Changes in the composition of the RNA

992          virome mark evolutionary transitions in green plants. BMC biology 14:68.

993    18.   Short SM, Staniewski MA, Chaban YV, Long AM, Wang DL. 2020. Diversity of Viruses

994          Infecting Eukaryotic Algae. Current Issues in Molecular Biology 39:29-61.

995    19.   Gibbs AJ, Torronen M, Mackenzie AM, Wood JT, Armstrong JS, Kondo H, Tamada T,

996          Keese PL. 2011. The enigmatic genome of Chara australis virus. Journal of General

997          Virology 92:2679-2690.

998    20.   Vlok M, Gibbs AJ, Suttle CA. 2019. Metagenomes of a freshwater charavirus from British

999          Columbia provide a window into ancient lineages of viruses. Viruses 11.

1000   21.   Han G-Z. 2019. Origin and evolution of the plant immune system. New Phytologist

1001          222:70-83.

1002    22.    Brunkard JO, Zambryski PC. 2017. Plasmodesmata enable multicellularity: new insights

1003            into their evolution, biogenesis, and functions in development and immunity. Current

1004            Opinion in Plant Biology 35:76-83.

1005    23.    Greninger AL. 2018. A decade of RNA virus metagenomics is (not) enough. Virus

1006            Research 244:218-229.

1007    24.    Miller AK, Mifsud JCO, Costa VA, Grimwood RM, Kitson J, Baker C, Brosnahan CL,

1008            Pande A, Holmes EC, Gemmell NJ, Geoghegan JL. 2021. Slippery when wet: cross-

1009            species transmission of divergent coronaviruses in bony and jawless fish and the

1010            evolutionary history of the Coronaviridae. Virus Evolution doi:10.1093/ve/veab050.

1011    25.    Bejerman N, Dietzgen RG, Debat H. 2021. Illuminating the Plant Rhabdovirus

1012            Landscape through Metatranscriptomics Data. Viruses 13:1304.

1013    26.    Parry R, Wille M, Turnbull OMH, Geoghegan JL, Holmes EC. 2020. Divergent Influenza-

1014            Like Viruses of Amphibians and Fish Support an Ancient Evolutionary Association.

1015            Viruses 12:1042.

1016    27.    Grimwood RM, Holmes EC, Geoghegan JL. 2021. A Novel Rubi-Like Virus in the Pacific

1017            Electric Ray (Tetronarce californica) Reveals the Complex Evolutionary History of the

1018            Matonaviridae. Viruses 13.

1019    28.    Gilbert KB, Holcomb EE, Allscheid RL, Carrington JC. 2019. Hiding in plain sight: New

1020            virus genomes discovered via a systematic analysis of fungal public transcriptomes.

1021            PLoS One 14:e0219207.

1022    29.    Lauber C, Seitz S, Mattei S, Suh A, Beck J, Herstein J, Börold J, Salzburger W, Kaderali

1023            L, Briggs JAG, Bartenschlager R. 2017. Deciphering the Origin and Evolution of

1024            Hepatitis B Viruses by Means of a Family of Non-enveloped Fish Viruses. Cell host &

1025            microbe 22:387-399.e6.

30.  Leebens-Mack JH, Barker MS, Carpenter EJ, Deyholos MK, Gitzendanner MA, Graham SW, Grosse I, Li Z, Melkonian M, Mirarab S, Porsch M, Quint M, Rensing SA, Soltis DE, Soltis PS, Stevenson DW, Ullrich KK, Wickett NJ, DeGironimo L, Edger PP, Jordon-Thaden IE, Joya S, Liu T, Melkonian B, Miles NW, Pokorny L, Quigley C, Thomas P, Villarreal JC, Augustin MM, Barrett MD, Baucom RS, Beerling DJ, Benstein RM, Biffin E, Brockington SF, Burge DO, Burris JN, Burris KP, Burtet-Sarramegna V, Caicedo AL, Cannon SB, Çebi Z, Chang Y, Chater C, Cheeseman JM, Chen T, Clarke ND, Clayton H, Covshoff S, et al. 2019. One thousand plant transcriptomes and the phylogenomics of green plants. Nature 574:679-685.

31.  Carpenter EJ, Matasci N, Ayyampalayam S, Wu S, Sun J, Yu J, Jimenez Vieira FR, Bowler C, Dorrell RG, Gitzendanner MA, Li L, Du W, K. Ullrich K, Wickett NJ, Barkmann TJ, Barker MS, Leebens-Mack JH, Wong GK-S. 2019. Access to RNA-sequencing data from 1,173 plant species: The 1000 Plant transcriptomes initiative (1KP). GigaScience 8.

32.  Johnson MT, Carpenter EJ, Tian Z, Bruskiewich R, Burris JN, Carrigan CT, Chase MW, Clarke ND, Covshoff S, dePamphilis CW. 2012. Evaluating methods for isolating total RNA and predicting the success of sequencing phylogenetically diverse plant transcriptomes. PLoS One 7:e50226.

33.  Leinonen R, Sugawara H, Shumway M, Collaboration INSD. 2010. The sequence read archive. Nucleic acids research 39:19-21.

34.  Langmead B, Salzberg SL. 2012. Fast gapped-read alignment with Bowtie 2. Nature Methods 9:357-359.

35.  Haas BJ, Papanicolaou A, Yassour M, Grabherr M, Blood PD, Bowden J, Couger MB, Eccles D, Li B, Lieber M, MacManes MD, Ott M, Orvis J, Pochet N, Strozzi F, Weeks N, Westerman R, William T, Dewey CN, Henschel R, Leduc RD, Friedman N, Regev A.

1050        2013. De novo transcript sequence reconstruction from RNA-seq using the Trinity

1051        platform for reference generation and analysis. Nature Protocols 8:1494-1512.

1052   36.   Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. 1990. Basic local alignment

1053        search tool. Journal of Molecular Biology 215:403-410.

1054   37.   Buchfink B, Xie C, Huson DH. 2015. Fast and sensitive protein alignment using

1055        DIAMOND. Nature Methods 12:59-60.

1056   38.   Mihara T, Nishimura Y, Shimizu Y, Nishiyama H, Yoshikawa G, Uehara H, Hingamp P,

1057        Goto S, Ogata H. 2016. Linking Virus Genomes with Host Taxonomy. Viruses 8:66.

1058   39.   Gilmer D, Ratti C, Consortium IR. 2017. ICTV Virus taxonomy profile: Benyviridae. The

1059        Journal of general virology 98:1571.

1060   40.   RStudio T. 2020. RStudio: integrated development for R.

1061   41.   Team RC. 2013. R: A language and environment for statistical computing. Vienna,

1062        Austria.

1063   42.   Wickham H, Averick M, Bryan J, Chang W, McGowan LDA, François R, Grolemund G,

1064        Hayes A, Henry L, Hester J. 2019. Welcome to the Tidyverse. Journal of Open Source

1065        Software 4:1686.

1066   43.   Li B, Dewey CN. 2011. RSEM: accurate transcript quantification from RNA-Seq data

1067        with or without a reference genome. BMC Bioinformatics 12:323.

1068   44.   Bushnell B. 2014. BBMap: a fast, accurate, splice-aware aligner.  Lawrence Berkeley

1069        National Lab.(LBNL), Berkeley, CA (United States),

1070   45.   Kearse M, Moir R, Wilson A, Stones-Havas S, Cheung M, Sturrock S, Buxton S, Cooper

1071        A, Markowitz S, Duran C. 2012. Geneious Basic: an integrated and extendable desktop

1072        software platform for the organization and analysis of sequence data. Bioinformatics

1073        28:1647-1649.

1074    46.    Sievers F, Wilm A, Dineen D, Gibson TJ, Karplus K, Li W, Lopez R, McWilliam H,

1075           Remmert M, Söding J. 2011. Fast, scalable generation of high-quality protein multiple

1076           sequence alignments using Clustal Omega. Molecular systems biology 7:539.

1077    47.    Mirdita M, Steinegger M, Söding J. 2019. MMseqs2 desktop and local web server app

1078           for fast, interactive sequence searches. Bioinformatics 35:2856-2858.

1079    48.    Zimmermann L, Stephens A, Nam S-Z, Rau D, Kübler J, Lozajic M, Gabler F, Söding J,

1080           Lupas AN, Alva V. 2018. A Completely Reimplemented MPI Bioinformatics Toolkit with a

1081           New HHpred Server at its Core. Journal of Molecular Biology 430:2237-2243.

1082    49.    Lay CL. 2021. biolumber/littlegenomes: First release.  doi:10.5281/ZENODO.5081375.

1083    50.    Sayers EW, Cavanaugh M, Clark K, Pruitt KD, Schoch CL, Sherry ST, Karsch-Mizrachi I.

1084           2021. GenBank. Nucleic acids research 49:D92-D96.

1085    51.    Gertz EM, Yu Y-K, Agarwala R, Schäffer AA, Altschul SF. 2006. Composition-based

1086           statistics and translated nucleotide searches: improving the TBLASTN module of

1087           BLAST. BMC biology 4:1-14.

1088    52.    Nayfach S, Camargo AP, Schulz F, Eloe-Fadrosh E, Roux S, Kyrpides NC. 2021.

1089           CheckV assesses the quality and completeness of metagenome-assembled viral

1090           genomes. Nature biotechnology 39:578-585.

1091    53.    Marcelino VR, Clausen PTLC, Buchmann JP, Wille M, Iredell JR, Meyer W, Lund O,

1092           Sorrell TC, Holmes EC. 2020. CCMetagen: comprehensive and accurate identification of

1093           eukaryotes and prokaryotes in metagenomic data. Genome Biology 21:103.

1094    54.    Clausen PT, Aarestrup FM, Lund O. 2018. Rapid and precise alignment of raw reads

1095           against redundant databases with KMA. BMC bioinformatics 19:1-8.

1096    55.    Ondov BD, Bergman NH, Phillippy AM. 2011. Interactive metagenomic visualization in a

1097           Web browser. BMC Bioinformatics 12:385.

1098 56. Capella-Gutiérrez S, Silla-Martínez JM, Gabaldón T. 2009. trimAl: a tool for automated

1099 alignment trimming in large-scale phylogenetic analyses. Bioinformatics 25:1972-1973.

1100 57. Nguyen L-T, Schmidt HA, Von Haeseler A, Minh BQ. 2015. IQ-TREE: a fast and

1101 effective stochastic algorithm for estimating maximum-likelihood phylogenies. Molecular

1102 biology and evolution 32:268-274.

1103 58. Kalyaanamoorthy S, Minh BQ, Wong TK, Von Haeseler A, Jermiin LS. 2017.

1104 ModelFinder: fast model selection for accurate phylogenetic estimates. Nature methods

1105 14:587-589.

1106 59. Rambaut A, Drummond A. 2012. FigTree: Tree figure drawing tool, version 1.4.0.

1107 Institute of Evolutionary Biology, University of Edinburgh.

1108 60. Zanne AE, Tank DC, Cornwell WK, Eastman JM, Smith SA, FitzJohn RG, McGlinn DJ,

1109 O'Meara BC, Moles AT, Reich PB. 2014. Three keys to the radiation of angiosperms into

1110 freezing environments. Nature 506:89-92.

1111 61. Smith SA, Brown JW. 2018. Constructing a broadly inclusive seed plant phylogeny. Am

1112 J Bot 105:302-314.

1113 62. Jin Y, Qian H. 2019. V.PhyloMaker: an R package that can generate very large

1114 phylogenies for vascular plants. Ecography 42:1353-1359.

1115 63. Hatcher EL, Zhdanov SA, Bao Y, Blinkova O, Nawrocki EP, Ostapchuck Y, Schäffer AA,

1116 Brister JR. 2017. Virus Variation Resource–improved response to emergent viral

1117 outbreaks. Nucleic acids research 45:D482-D490.

1118 64. Paradis E, Schliep K. 2019. ape 5.0: an environment for modern phylogenetics and

1119 evolutionary analyses in R. Bioinformatics 35:526-528.

1120 65. Revell LJ. 2012. phytools: an R package for phylogenetic comparative biology (and other

1121 things). Methods in Ecology and Evolution 3:217-223.

1122    66.    Conow C, Fielder D, Ovadia Y, Libeskind-Hadas R. 2010. Jane: a new tool for the

1123             cophylogeny reconstruction problem. Algorithms for Molecular Biology 5:1-10.

1124    67.    Santichaivekin S, Yang Q, Liu J, Mawhorter R, Jiang J, Wesley T, Wu Y-C, Libeskind-

1125             Hadas R. 2020. eMPRess: a systematic cophylogeny reconciliation tool. Bioinformatics

1126             doi:10.1093/bioinformatics/btaa978.

1127    68.    Quast C, Pruesse E, Yilmaz P, Gerken J, Schweer T, Yarza P, Peplies J, Glöckner FO.

1128             2012. The SILVA ribosomal RNA gene database project: improved data processing and

1129             web-based tools. Nucleic acids research 41:D590-D596.

1130    69.    Morozov SY, Solovyev AG. 2003. Triple gene block: modular design of a multifunctional

1131             machine for plant virus movement. Journal of General Virology 84:1351-1366.

1132    70.    Hammond R, Ramirez P. 2001. Molecular characterization of the genome of Maize

1133             rayado fino virus, the type member of the genus Marafivirus. Virology 282:338-347.

1134    71.    Ding S, Howe J, Keese P, Mackenzie A, Meek D, Osorlo-Keese M, Skotnicki M, Srifah

1135             P, Torronen M, Gibbs A. 1990. The tymobox, a sequence shared by most tymoviruses:

1136             its use in molecular studies of tymoviruses. Nucleic acids research 18:1181-1187.

1137    72.    Xie J, Ghabrial SA. 2012. Molecular characterizations of two mitoviruses co-infecting a

1138             hyovirulent isolate of the plant pathogenic fungus Sclerotinia sclerotiorum. Virology

1139             428:77-85.

1140    73.    Heinze C. 2012. A novel mycovirus from Clitocybe odora. Archives of virology 157:1831-

1141             1834.

1142    74.    Nibert ML. 2017. Mitovirus UGA (Trp) codon usage parallels that of host mitochondria.

1143             Virology 507:96-100.

1144    75.    Zhang T, Li W, Chen H, Yu H. 2015. Full genome sequence of a putative novel mitovirus

1145             isolated from Rhizoctonia cerealis. Archives of virology 160:1815-1818.

1146 76. Shackelton LA, Holmes EC. 2008. The role of alternative genetic codes in viral evolution

1147 and emergence. Journal of Theoretical Biology 254:128-134.

1148 77. Nerva L, Turina M, Zanzotto A, Gardiman M, Gaiotti F, Gambino G, Chitarra W. 2019.

1149 Isolation, molecular characterization and virome analysis of culturable wood fungal

1150 endophytes in esca symptomatic and asymptomatic grapevine plants. Environmental

1151 microbiology 21:2886-2904.

1152 78. Charon J, Marcelino VR, Wetherbee R, Verbruggen H, Holmes EC. 2020.

1153 Metatranscriptomic identification of diverse and divergent RNA viruses in green and

1154 Chlorarachniophyte algae cultures. Viruses 12.

1155 79. Koga R, Horiuchi H, Fukuhara T. 2003. Double-stranded RNA replicons associated with

1156 chloroplasts of a green alga, Bryopsis cinicola. Plant molecular biology 51:991-999.

1157 80. Cai G, Myers K, Hillman BI, Fry WE. 2009. A novel virus of the late blight pathogen,

1158 Phytophthora infestans, with two RNA segments and a supergroup 1 RNA-dependent

1159 RNA polymerase. Virology 392:52-61.

1160 81. Covelli L, Coutts RH, Di Serio F, Citir A, Açıkgöz S, Hernandez C, Ragozzino A, Flores

1161 R. 2004. Cherry chlorotic rusty spot and Amasya cherry diseases are associated with a

1162 complex pattern of mycoviral-like double-stranded RNAs. I. Characterization of a new

1163 species in the genus Chrysovirus. Journal of General Virology 85:3389-3397.

1164 82. Rousvoal S, Bouyer B, López-Cristoffanini C, Boyen C, Collén J. 2016. Mutant swarms

1165 of a totivirus-like entities are present in the red macroalga Chondrus crispus and have

1166 been partially transferred to the nuclear genome. Journal of phycology 52:493-504.

1167 83. Charon J, Murray S, Holmes EC. 2021. Revealing RNA virus diversity and evolution in

1168 unicellular algae transcriptomes. Virus Evolution 7.

1169 84. Kenrick P, Crane PR. 1997. The origin and early evolution of plants on land. Nature

1170 389:33-39.

1171  85.  Valente JB, Pereira FS, Stempkowski LA, Farias M, Kuhnem P, Lau D, Fajardo TVM,

1172       Nhani Junior A, Casa RT, Bogo A, da Silva FN. 2019. A novel putative member of the

1173       family Benyviridae is associated with soilborne wheat mosaic disease in Brazil. Plant

1174       Pathology 68:588-600.

1175  86.  Tamada T, Schmitt C, Saito M, Guilley H, Richards K, Jonard G. 1996. High resolution

1176       analysis of the readthrough domain of beet necrotic yellow vein virus readthrough

1177       protein: a KTER motif is important for efficient transmission of the virus by Polymyxa

1178       betae. Journal of General Virology 77:1359-1367.

1179  87.  Neuhauser S, Kirchmair M, Bulman S, Bass D. 2014. Cross-kingdom host shifts of

1180       phytomyxid parasites. BMC Evolutionary Biology 14:33.

1181  88.  Starr EP, Nuccio EE, Pett-Ridge J, Banfield JF, Firestone MK. 2019. Metatranscriptomic

1182       reconstruction reveals RNA viruses with the potential to shape carbon cycling in soil.

1183       Proceedings of the National Academy of Sciences 116:25900.

1184  89.  Ghabrial SA, Caston JR, Jiang D, Nibert ML, Suzuki N. 2015. 50-plus years of fungal

1185       viruses. Virology 479-480:356-68.

1186  90.  Bruenn JA, Warner BE, Yerramsetty P. 2015. Widespread mitovirus sequences in plant

1187       genomes. Peerj 3.

1188  91.  Moreno Goncalves AB, Lopez-Moya JJ. 2019. When viruses play team sports: mixed

1189       infections in plants. Phytopathology 110:29-48.

1190  92.  Botella L, Jung T. 2021. Multiple Viral Infections Detected in Phytophthora condilina by

1191       Total and Small RNA Sequencing. Viruses 13:620.

1192  93.  Botella L, Janoušek J, Maia C, Jung MH, Raco M, Jung T. 2020. Marine Oomycetes of

1193       the Genus Halophytophthora Harbor Viruses Related to Bunyaviruses. Frontiers in

1194       Microbiology 11.

94.  Su S, Liu Z, Chen C, Zhang Y, Wang X, Zhu L, Miao L, Wang X-C, Yuan M. 2010. Cucumber Mosaic Virus Movement Protein Severs Actin Filaments to Increase the Plasmodesmal Size Exclusion Limit in Tobacco  The Plant Cell 22:1373-1387.

95.  Mushegian AR, Elena SF. 2015. Evolution of plant virus movement proteins from the 30K superfamily and of their homologs integrated in plant genomes. Virology 476:304-315.

96.  Kumar S, Subbarao BL, Kumari R, Hallan V. 2017. Molecular characterization of a novel cryptic virus infecting pigeonpea plants. PloS one 12:e0181829.

97.  Sabanadzovic S, Ghanem-Sabanadzovic NA. 2008. Molecular characterization and detection of a tripartite cryptic virus from rose. Journal of Plant Pathology:287-293.

98.  Chen L, Chen J, Liu L, Yu X, Yu S, Fu T, Liu W. 2006. Complete nucleotide sequences and genome characterization of double-stranded RNA 1 and RNA 2 in the Raphanus sativus-root cv. Yipinghong. Archives of virology 151:849-859.

99.  Wu LP, Du YM, Xiao H, Peng L, Li R. 2020. Complete genomic sequence of tea-oil camellia deltapartitivirus 1, a novel virus from Camellia oleifera. Archives of Virology 165:227-231.

100.  Chiba S, Lin YH, Kondo H, Kanematsu S, Suzuki N. 2013. Effects of Defective Interfering RNA on Symptom Induction by, and Replication of, a Novel Partitivirus from a Phytopathogenic Fungus, Rosellinia necatrix. Journal of Virology 87:2330-2341.

101.  Nibert ML, Ghabrial SA, Maiss E, Lesker T, Vainio EJ, Jiang D, Suzuki N. 2014. Taxonomic reorganization of family Partitiviridae and other recent progress in partitivirus research. Virus Research 188:128-141.

102.  Harrison CJ, Morris JL. 2018. The origin and early evolution of vascular plant shoots and leaves. Philosophical Transactions of the Royal Society B: Biological Sciences 373:20160496.

1220    103.    Walker PJ, Blasdell KR, Calisher CH, Dietzgen RG, Kondo H, Kurath G, Longdon B,

1221            Stone DM, Tesh RB, Tordo N. 2018. ICTV virus taxonomy profile: Rhabdoviridae.

1222            Journal of General Virology 99:447-448.

1223    104.    Shi M, Lin X-D, Tian J-H, Chen L-J, Chen X, Li C-X, Qin X-C, Li J, Cao J-P, Eden J-S,

1224            Buchmann J, Wang W, Xu J, Holmes EC, Zhang Y-Z. 2016. Redefining the invertebrate

1225            RNA virosphere. Nature 540:539.

1226    105.    Hao X, Zhang W, Zhao F, Liu Y, Qian W, Wang Y, Wang L, Zeng J, Yang Y, Wang X.

1227            2018. Discovery of plant viruses from tea plant (Camellia sinensis (L.) O. Kuntze) by

1228            metagenomic sequencing. Frontiers in Microbiology 9:2175.

1229    106.    Maclot F, Candresse T, Filloux D, Malmstrom CM, Roumagnac P, van der Vlugt R,

1230            Massart S. 2020. Illuminating an ecological blackbox: using high throughput Sequencing

1231            to characterize the plant virome across scales. Frontiers in Microbiology 11:2575.

1232    107.    Visser M, Bester R, Burger JT, Maree HJ. 2016. Next-generation sequencing for virus

1233            detection: covering all the bases. Virology Journal 13:85.

1234    108.    Llorens C, Muñoz-Pomer A, Bernad L, Botella H, Moya A. 2009. Network dynamics of

1235            eukaryotic LTR retroelements beyond phylogenetic trees. Biology Direct 4:41.

1236    109.    Letunic I, Bork P. 2019. Interactive Tree Of Life (iTOL) v4: recent updates and new

1237            developments. Nucleic acids research 47:W256-W259.

1238    110.    Brister JR, Ako-Adjei D, Bao Y, Blinkova O. 2015. NCBI viral genomes resource. Nucleic

1239            Acids Res 43:D571-7.

1240

1241 **Supplementary Information**

1242 **Supplementary Figure 1.** (A) Phylogram of the triple gene block (TGB) protein 1. ML

1243 phylogenetic trees show the topological position of the newly discovered TGB sequence in the

1244 tomato fern (black circle) in the context of the closest relatives. (B) Phylogram of the

1245 Tymoviridae virus coat proteins (CP). ML phylogenetic trees show the topological position of the

1246 newly discovered CP sequences in (black circle) in the context of the closest relatives.

1247 Branches are highlighted to represent virus taxonomy (Maculavirus = green, Marafivirus =

1248 orange, Tymovirus = red and unclassified = grey). For each colour, a lighter shade signifies that

1249 this virus is related to but has not formally been assigned to this genus. (C) Phylogram of the

1250 *Oxera neriifolia associated virus* coat protein (CP). ML phylogenetic trees show the topological

1251 position of the newly discovered CP sequence (black circle) in the context of the closest

1252 relatives. For all trees, branches are scaled to the number of amino acid substitutions per site

1253 and trees were mid-point rooted for clarity only. Numbers at the nodes indicate bootstrap

1254 support over 70% (1000 replicates).

1255 **Supplementary Figure 2.** Multiple sequence alignment conserved amino acid motifs in RNA-

1256 dependent RNA polymerase (RdRp) regions of the mitoviruses discovered in this study along

1257 with reference mitoviruses. The bar above each residue is green if 100% of residues in that

1258 column are identical, green-brown if they are 30%-99%, and red if under 30%. The numbers

1259 under each section correspond to regions containing motifs identified in (72).

1260 **Supplementary Figure 3.** (A) Phylogenetic relationships of the viruses identified within the

1261 virus families *Potyviridae and Tombusviridae*. ML phylogenetic trees based upon alignments of

1262 the amino acid sequences of the RdRp protein show the topological position of discovered

1263 virus-like sequences (black circles) from this study in the context of their closest relatives. See

1264 Figure 3 for the colour scheme. All branches are scaled to the number of amino acid

1265    substitutions per site and trees were mid-point rooted for clarity only. An asterisk indicates node

1266    support of >70% bootstrap support.

1267    **Supplementary Figure 4.** Phylogram of the deltapartitii-like virus (A) coat protein/RNA2 (CP)

1268    and (B) RNA3/coat protein 2. ML phylogenetic trees show the topological position of the newly

1269    discovered CP sequences in (black circle) in the context of the closest relatives. All branches

1270    are scaled to the number of amino acid substitutions per site and trees were mid-point rooted for

1271    clarity only. Numbers at the nodes indicate bootstrap support over 70% (1000 replicates).

1272    **Supplementary Figure 5.** Tanglegram of rooted phylogenetic trees for select virus families and

1273    their hosts. Lines and branches are coloured to represent host clade. The cophylo function

1274    implemented in phytools (v0.7-80) was used to maximise the congruence between the host (left)

1275    and virus (right) phylogenies.

1276    **Supplementary Table 1**. Clade assignment for all One Thousand Plant Transcriptomes

1277    Initiative (1KP) species for which a virus was detected.

1278    **Supplementary Table 2.** Summary information for each One Thousand Plant Transcriptomes

1279    Initiative (1KP) libraries analysed.

1280    **Supplementary Table 3.** Proportion of transcripts and abundance assigned to each plant virus

1281    family.

1282    **Supplementary Table 4.** Summary table of the viruses discovered in this study

1283    **Supplementary Table 5.** Genome annotation information underlying the annotation graphs

## Supplementary References

1284

1285    1.    Xie J, Ghabrial SA. 2012. Molecular characterizations of two mitoviruses co-infecting a

1286          hyovirulent isolate of the plant pathogenic fungus Sclerotinia sclerotiorum. Virology

1287          428:77-85.