1 # Tracheophyte genomes keep track of the deep evolution of the
2 # *Caulimoviridae*
3
4 **Authors**
5 Seydina Diop[1], Andrew D.W. Geering[2], Françoise Alfama-Depauw[1], Mikaël Loaec[1], Pierre-Yves
6 Teycheney[3] and Florian Maumus[1*]
7
8 **Affiliations**
9 [1] URGI, INRA, Université Paris-Saclay, 78026 Versailles, France;
10 [2] Queensland Alliance for Agriculture and Food Innovation, The University of Queensland, GPO Box
11 267, Brisbane, Queensland 4001, Australia
12 [3] UMR AGAP, CIRAD, INRA, SupAgro, 97130 Capesterre Belle-Eau, France
13
14 **Corresponding author**
15 Florian Maumus
16 URGI-INRA
17 RD10 route de Saint Cyr
18 78026, Versailles
19 France
20 +33 1 30 83 31 74
21 florian.maumus@inra.fr
22
23
24

25  **Abstract**

26  Endogenous viral elements (EVEs) are viral sequences that are integrated in the nuclear genomes of
27  their hosts and are signatures of viral infections that may have occurred millions of years ago. The
28  study of EVEs, coined paleovirology, provides important insights into virus evolution. The
29  *Caulimoviridae* is the most common group of EVEs in plants, although their presence has often been
30  overlooked in plant genome studies. We have refined methods for the identification of caulimovirid
31  EVEs and interrogated the genomes of a broad diversity of plant taxa, from algae to advanced
32  flowering plants. Evidence is provided that almost every vascular plant (tracheophyte), including the
33  most primitive taxa (clubmosses, ferns and gymnosperms) contains caulimovirid EVEs, many of which
34  represent previously unrecognized evolutionary branches. In angiosperms, EVEs from at least one
35  and as many as five different caulimovirid genera were frequently detected, and florendoviruses
36  were the most widely distributed, followed by petuviruses. From the analysis of the distribution of
37  different caulimovirid genera within different plant species, we propose a working evolutionary
38  scenario in which this family of viruses emerged at latest during Devonian era (approx. 320 million
39  years ago) followed by vertical transmission and by several cross-division host swaps.

40

41

**Introduction**

Although the field of viral metagenomics is rapidly expanding the repertoire of viral genome sequences available for evolutionary studies [1], it only provides a picture of viral diversity over a very short geological time scale. However, viruses can leave molecular records in the genomes of their hosts in the form of endogenous viral elements (EVEs). EVEs are viral sequences that have been inserted in the nuclear genomes of their hosts by either active or passive integration mechanisms and in many cases have been retained for extended periods of time, sometimes millions of years. The study of EVEs does allow the evolution of viruses to be traced, much like a fossil record [2]. For example, the study of endogenous retroviruses has uncovered the extended diversity and host range of retroviruses, and has provided evidence that they have a marine origin, and that they developed in parallel with their vertebrate hosts more than 450 million years ago (MYA; [3 4 5]).

Plant EVEs were first discovered a little more than 20 years ago [6] but have only received a fraction of the research attention directed towards endogenous retroviruses in humans and other animals. Most characterized plant EVEs are derivatives of viruses in the family *Caulimoviridae* [7]. The *Caulimoviridae* is one of the five families of reverse-transcribing viruses or virus-like retrotransposons that occur in eukaryotes [8], and is the only family of viruses with a double-stranded DNA genome that infects plants (https://talk.ictvonline.org/). Unlike retroviruses, members of the *Caulimoviridae* do not integrate in the genome of their host to complete their replication cycle. However, caulimovirid DNA can occasionally integrate passively into their host genome. In fact, five of the eight officially recognized genera of the *Caulimoviridae* have EVE counterparts in at least one plant genome [7 9].

Recently, Geering et al. [10] showed that EVEs from an additional tentative genus of the *Caulimoviridae*, called 'Florendovirus', are widespread in the genomes of cultivated and wild angiosperms, and provided evidence for the oldest EVE integration event yet reported in plants, at 1.8 MYA [10]. Furthermore, sister taxa relationships between florendoviruses in South American and Australian plants suggested Gondwanan links and a minimum age of 34 MYA for this virus group based on estimates of when land bridges between these two continents were severed. About 65% of all angiosperm species that were examined contained endogenous florendoviruses and for five, these sequences contributed more than 0.5% of the total plant genome content. Furthermore, the discovery of endogenous florendoviruses in basal ANITA (*Amborella*, *Nymphaeales* and *Illiciales*, *Trimeniaceae-Austrobaileya*) grade angiosperm species also showed that beyond mesangiosperms, the host range of the *Caulimoviridae* extends or once extended to the most primitive known angiosperms. Furthermore, some reconstructed endogenous florendovirus genomes were bipartite in organization, a genome arrangement that is unique among viral retroelements [10]. Overall, the work of Geering et al. (2014) demonstrated that analyzing the genetic footprints left by viruses in plant genomes can contribute to a better understanding of the long-term evolution of the *Caulimoviridae*.

In this study, we hypothesized that other unrecognized groups of caulimovirid EVEs could exist, particularly in some of the more primitive plant taxa that have recently been subject to plant genome sequencing initiatives. Following an extensive search in the genomes of over 70 plant species, we discovered EVEs from several novel genera. We show that the *Caulimoviridae* host range extends throughout the Euphyllophyte and Lycopodiophyte clades, which constitute the Tracheophyta, and surpasses that of any other plant virus family. By analyzing the distribution of different genera of the *Caulimoviridae* within different plant species, we unveil a complex pattern of associations and propose a scenario in which the *Caulimoviridae* would have emerged approximately 320 million years ago.

3

87

## Results

### Augmenting the diversity of known endogenous caulimovirids

The reverse transcriptase (RT) domain is the most conserved domain in the genome of viral retroelements and is used for classification [11] [12]. The strong sequence conservation of this domain allows high quality alignments to be generated, even for distantly related taxa. We have thus used a collection of RT domains from known exogenous and endogenous caulimovirids to search for related sequences across the breadth of the Viridiplantae (four green algae, one moss, one lycopod, four gymnosperms, and 62 angiosperms; Supplementary Table 1) using tBLASTn.

Initially, over 8,400 protein-coding sequences were retrieved, all containing an RT domain with a best reciprocal hit against members of the *Caulimoviridae*, as opposed to the closely related *Metaviridae* (Ty3/Gypsy group LTR retrotransposons). To provide a preliminary classification, sequences with at least 55% amino acid identity to each other were clustered and then iteratively added to our reference set of RT domains to build a sequence similarity network. The successive networks were examined manually and representative sequences from each cluster were kept only when creating substantially divergent branches so as to cover an extended diversity of caulimovirid RT with a core sequence assortment. While this network-based approach cannot be taken as phylogenetic reconstruction, it provided a practical method to explore diversity.

In the final sequence similarity network (Figure 1), 17 groups with deep connections were identified, hereafter referred to as operational taxonomic units (OTUs). Remarkably, nine of these OTUs were distinct from recognized genera of the *Caulimoviridae*. Four of these novel OTUs were exclusively composed of sequences from gymnosperms, thereby representing a new and significant host range extension for the *Caulimoviridae*. These OTUs were named Gymnendovirus 1 to 4. Two other novel OTUs were composed of RTs from various angiosperms and were named Xendovirus and Yendovirus. The last three novel OTUs were lesser populated, comprising sequences from one or two plant species (*Petunia inflata* and *Petunia axillaris*; *Vitis vinifera*; *Glycine max*; named species-wise: Petunia-, Vitis-, and Glycine-endovirus). This initial search therefore enabled uncovering a significantly augmented diversity of caulimovirid RTs.

### Endogenous caulimovirid RT (ECRT) density across the Viridiplantae

To perform a more comprehensive search for ECRTs in our collection of plant genomes, we used the sequences from the final phylogenetic network (Figure 1) to search for ECRT nucleotide sequences that do not necessarily retain uninterrupted open reading frames. Using tBLASTn, we detected 14,895 genomic loci representing high-confidence ECRT candidates. Remarkably, ECRTs were found in nearly all seed plants, ranging from gymnosperms (ginkgo and conifers) to angiosperms. Quantitatively, over one-thousand ECRTs were detected in the genome assemblies of the gymnosperms *Picea glauca* (white spruce) and *Pinus taeda* (loblolly pine), as well as from the solanaceous plant species *Capsicum annuum* (bell pepper) (Figure 2A). In general, we observed a positive correlation between plant genome size and the number of ECRTs, although there were notable exceptions, such as the monocot *Zea mays* (maize), which has a relatively large genome at 2.1 Gb but no detectable ECRT. Five other seed plants from our sample also lacked ECRTs, including two other monocots (*Zostera marina* and *Oryza brachyantha*) and three dicots in the order Brassicales (*Arabidopsis thaliana*, *Schrenkiella parvula* and *Carica papaya*). When the number of

4

132    ECRTs was normalized against genome size, *Citrus sinensis* (sweet orange) and *Ricinus communis*
133    (castor bean) had the highest densities at 2.3 and 2 ECRTs per Mb, respectively (Figure 2B). The
134    primitive ANITA grade angiosperm *Amborella trichopoda* also had a relatively high density of ECRTs
135    (1 ECRT per Mb) compared to an average density of 0.2 ECRT per Mb across the 62 seed plant species
136    that were examined.

137

138            **Caulimovirid sequences also detected in ferns and a clubmoss**
139    From the plant genomes examined thus far, ECRTs were detected in gymnosperm genomes but not
140    in those from the spikemoss *Selaginella moellendorffii* and from the moss *Physcomitrella patens*,
141    which belong to the more basal land plant divisions Lycophyta and Bryophyta, respectively. Ferns
142    (class Polypodiopsida) represent a bifurcation between Lycophyta and seed plants in the evolution of
143    the Viridiplantae [13], but no high quality genome assemblies are publicly available for these plants.
144    However, six fern genomes have recently been sequenced at low coverage (approximately 0.4 to 2 x
145    genome size equivalent [14] and we therefore screened these datasets for the presence of ECRTs. A
146    total of twenty-one protein-coding ECRTs were detected in genomic contigs from five of the six fern
147    species examined (Supplementary Table 1). Sequence similarity network reconstruction using
148    representative fern ECRTs revealed that they form two novel OTUs that were named Fernendovirus 1
149    and 2 and numbered OTU #18 and 19, respectively (Supplementary Figure 2).

150

151    Additional basal lineages of the Viridiplantae are represented in the 1,000 plant transcriptomes
152    generated by the 1KP initiative [15] [16]. From this dataset, we found two transcript contigs (2.4 and 2.8
153    kilobases long, respectively) in the fern *Botrypus virginianus* (identifier BEGM-2004510) and *Lindsaea*
154    *linearis* (identifier NOKI-2097008), which contained ECRTs (Supplementary file 1). Remarkably, we
155    identified one more transcript contig (identified as ENQF-2084799, 2kb) that contained an ECRT in
156    the clubmoss *Lycopodium annotinum*, which belongs to the *Lycopoda*, the most basal radiation of
157    vascular plants (Tracheophyta). It is not possible to determine whether the mRNAs were transcribed
158    from exogenous viruses or from EVEs.

159

160            **Phylogenetic reconstruction**
161    Complete or near complete viral genomes were reconstructed from each novel OTU except
162    Fernendovirus 1 and 2 (Supplementary file 1). From the fern genomic data sets, we were able to
163    reconstruct fragments of Fernendovirus 1 and 2 genomes that contain sufficient genome coverage
164    for phylogenetic analysis. We also used the complete genomes of the type species of the eight
165    currently recognized genera in the family *Caulimoviridae* (*Badnavirus*, *Caulimovirus*, *Cavemovirus*,
166    *Petuvirus*, *Rosadnavirus*, *Solendovirus*, *Soymovirus* and *Tungrovirus*), those of two unassigned viruses,
167    Blueberry fruit drop-associated virus (BFDaV, [17]) and Rudbeckia flower distortion virus (RuFDV, [18]),
168    and caulimovirid EVEs from the tentative genera Orendovirus [19] and Florendovirus (Geering et al.,
169    2014). From this library of sequences, we aligned conserved protease, reverse transcriptase and
170    ribonuclease H1 domains to build a maximum likelihood phylogenetic tree (Figure 3). Importantly, all
171    newly identified EVEs grouped within the *Caulimoviridae* with strong bootstrap support.

172

173    In agreement with previous studies [10], the tree revealed two sister clades, hereafter referred to as
174    clade A and B. Clade A comprised sequences from representatives of Xendovirus and Yendovirus
175    OTUs and from members of the genera *Caulimovirus*, *Soymovirus*, *Rosadnavirus*, *Solendovirus*,
176    *Cavemovirus*, *Badnavirus*, *Tungrovirus* and Orendovirus, as well as RuFDV and BFDaV. The Xendovirus

177  OTU was found to be polyphyletic, hence a new taxon, Zendovirus, was raised to include the EVE
178  from *Fragaria vesca*, while the EVE from *Gossypium raimondii* (cotton) was retained in Xendovirus.
179  The Yendovirus OTU, comprising the EVE from *Capsicum annuum* (bell pepper), fell in the subclade
180  comprising bacilliform-shaped viruses in the genera *Badnavirus* and *Tungrovirus*. The reconstructed
181  genomes from novel OTUs found in single dicot species (Petunia-, Vitis-, and Glycine-endovirus) were
182  discarded from the phylogenetic reconstruction as they significantly weakened the robustness of the
183  tree. However, they unambiguously fell within clade A (data not shown).

184

185  Clade B comprised EVEs from the four gymnendovirus OTUs, the two fernendovirus OTUs, as well as
186  representatives of the genus *Petuvirus* and the tentative genus Florendovirus. Gymnendovirus 2 EVEs
187  were sister to all other clade B viruses, indicating that this group of viruses arose in gymnosperms.
188  Interestingly, the angiosperm-infecting caulimovirids in this clade were polyphyletic, indicating
189  independent origins and probable large host jumps of the most recent common ancestors.
190  Fernendovirus 1 and 2 were monophyletic, and the club moss EVE placed within Fernendovirus 1.
191  Again, the fernendoviruses appear to have arisen after a large host range swap of the most recent
192  common ancestor.

193

194          **ECRT distribution across seed plant genomes**
195  To address the distribution of caulimovirid EVEs in our collection of plant genomes, we determined
196  the most likely position within the reference *Caulimoviridae* phylogenetic tree proposed above
197  (Figure 3) for the 14,895 ECRTs that we collected from seed plant genomes using the pplacer
198  program [20]. For this, we extracted ECRT loci extending upstream and downstream so as to retrieve
199  potential sequences containing the contiguous fragment corresponding to the protease, RT and
200  ribonuclease H1 domains. Using more relaxed length criteria, we extracted a total of 134 ECRT loci
201  from the fern genomic data set that we also attempted to place on our reference tree.

202

203  Applying this strategy, we were able to assign unambiguous phylogenetic position on specific OTUs
204  to a total of 13,834 ECRTs (Figure 4), the remaining ECRT loci being placed on inner nodes of the
205  reference tree. Overall, we observed striking differences between *Caulimoviridae* genera for both the
206  number of ECRT loci and the number of plant species in which they were found. For instance,
207  Florendovirus ECRT loci were the most abundant, amounting to an overall total of 5,000 copies, and
208  they were also found in the highest number of host species (46 of the 62 seed plant species that
209  were screened). *Petuvirus* ECRT loci were also well represented, with an overall total of 1,900 copies
210  found in a total of 27/62 seed plant species, especially in dicots. Among the novel OTUs, ECRTs
211  classified as Yendovirus were found in the largest number of species, including monocots and dicots
212  (Figure 4).

213

214  Most importantly, the detailed distribution of *Caulimoviridae* ECRTs in plant genomes reveals striking
215  differences between ferns, gymnosperms and angiosperms (Figure 4). No single OTU spans more
216  than one plant division on Figure 4 (which describes plant genomic EVE contents). Fernendovirus 1
217  sequences were found in both fern genomes and a lycopod transcriptome, but were found in no
218  other genomes. Gymnosperm genomes include exclusively ECRT loci that are assigned to one of the
219  four Gymnendovirus OTUs, all of which were undetected outside of gymnosperms. Among
220  gymnosperms, the three conifer genomes analyzed contain a mixture ECRTs from the four
221  Gymnendovirus OTUs. By contrast, only ECRT loci classified as Gymnendovirus 2 were detected in

222  *Ginkgo biloba* (*Ginkgoales*). Within angiosperms, we also observed a dichotomy for the distribution
223  of ECRTs between monocots and dicots. On one hand, Yendovirus, *Badnavirus*, Orendovirus and
224  Florendovirus ECRTs are common in monocots, with Orendovirus ECRTs being the only monocot-
225  specific ones. On the other hand, *Petuvirus*, Florendovirus, Xendovirus, *Cavemovirus*/*Solendovirus*
226  and Yendovirus ECRTs are most widely distributed in dicots, Florendovirus and Yendovirus hence
227  being remarkably well represented in both dicots and monocots.
228
229
230  **Discussion**
231  Endogenous viral elements are considered relics of past infections, and an extrapolation of the
232  results from this study is that nearly every tracheophyte plant species in the world has at some point
233  in its evolutionary history been subject to infection by at least one, and sometimes five distinct viral
234  species/genera from the family *Caulimoviridae*. This finding attests to the tremendous adaptability of
235  the *Caulimoviridae*, surpassing any other group of plant viruses. Members of the *Caulimoviridae* have
236  likely also had a large influence on plant evolution, either as pathogens or as donors of novel genetic
237  material to the plant genome.
238
239  A defining moment in the evolution of the *Caulimoviridae* appears to be the development of
240  vasculature in plants. The presence of a 30K movement protein is an important feature of the
241  *Caulimoviridae* that distinguishes it from the LTR retrotransposon family *Metaviridae*, and this
242  protein is crucial for the formation of systemic infection by allowing intercellular trafficking of
243  macromolecules through increasing the size exclusion limit of plasmodesmata [21]. Although algae
244  contain plasmodesmata, which superficially resemble those of higher plants, they are not
245  homologous structures [22]. While the acquisition of a 30K movement protein would have provided a
246  selective advantage for ancestral caulimovirids to colonize the tracheophytes, it would not have
247  facilitated infection of more primitive plant forms.
248
249  When recapitulating the distribution of EVEs in plant genomes, the known host range of exogenous
250  viruses, and the phylogenetic relationships between caulimovirid OTUs and major groups of vascular
251  plants (Figure 5) to infer the evolutionary trajectories of plant-virus coevolution, we obtain a complex
252  pattern of host-virus associations. At the OTU level, the host distributions of petu- and xendovirus,
253  including dicots and the ANITA grade angiosperm (*Amborella trichopoda*) but not any of the monocot
254  species, is suggestive of horizontal transfer. In addition, although vertical transmission is overall well
255  supported by a co-evolutionary study of florendoviral EVEs and their host species [10], it could not be
256  confirmed for *A. trichopoda*. Together with the observation that *A. trichopoda* presents a high
257  density of ECRTs (Figure 2B), this may suggest that this species is permissive to infection by a range of
258  caulimovirid genera and/or that it represents a hotbed for the emergence of caulimovirid genera,
259  some of which swapped towards mesangiosperms.
260
261  At a deeper level within the caulimovirid tree, clade A caulimovirids were found exclusively in
262  mesangiosperm species. Clade B viruses, instead, were found to associate with plants from all the
263  major classes of Tracheophyta. Assuming the monophyly of clades A and B, the obtained plant-virus
264  associations could be explained by the emergence of these viruses in a common ancestor of the
265  gymnosperms and angiosperms followed by several major host swaps: in clade A, between monocots
266  and dicots, and in clade B between gymnosperms and angiosperms in the case of florendoviruses and

7

267    petuviruses (although the position of this latter genus in the tree is uncertain), and from
268    gymnosperms to ferns and clubmoss in the case of fernendoviruses. Following this scenario, the
269    *Caulimoviridae* could have emerged at the latest with the Spermatophyta, *i.e.* during the Devonian
270    era, about 320 MYA [23]. Such a scenario would however imply several host swaps that overlay plant
271    evolutionary history and incomplete sampling in clade A as the OTUs that associate with monocots
272    are not sister to but nested in those that associate with dicots. We therefore do not rule out the
273    possibility that the observed plant-virus associations reflect to some extent stochastic survival (or
274    sampling) of a diversity of ancestral members of the *Caulimoviridae* that would have arose at latest
275    during the emergence of clubmoss towards the end of the Silurian era 420 MYA.
276
277

278    **Methods**
279           **Discovery and clustering of novel Caulimoviridae OTUs**
280    We built a library containing an assortment of amino acid (aa) sequences from 54 RT domains
281    including four from *Retroviridae*, six from Ty3/Gypsy LTR retrotransposons, 41 from eight different
282    *Caulimoviridae* genera (Florendovirus, *Caulimovirus*, *Tungrovirus*, *Cavemovirus*, *Solendovirus*,
283    *Badnavirus*, *Soymovirus*, and *Petuvirus*), two from *Picea glauca*, and the one from the DIRS-1
284    element. We compared this library to a collection of 72 genome assemblies from the *Viridiplantae*
285    (listed in Supplementary Table 1) using tBLASTn with default parameters (except –e=1e-5). The hit
286    genomic loci were merged when overlapping and their coordinates were extended 120 bases
287    upstream and downstream. Extended hit loci were translated and the protein sequences of length
288    >=200aa were compared to the initial RT library using BLASTp with default parameters (except –
289    e=1e-5). Queries with best alignment score against Caulimoviridae over at least 170 residues were
290    selected for further analysis. For each plant species, the selected set of RT aa sequences have been
291    clustered following sequence similarity using the UCLUST program [24] with identity threshold set at
292    80%. The longest sequence from each resulting cluster was considered as the representative
293    sequence and it was appended to the initial RT library. To detect potential false positives, each set of
294    sequences (each consisting of the initial RT library and cluster representatives from one species) was
295    aligned using MUSCLE followed by filtering of lower fit sequences using two rounds of trimAl v1.2 [25]
296    to remove poorly aligned sequences (-resoverlap 0.75 -seqoverlap 50) separated by one round to
297    remove gaps from the alignment (-gt 0.5). The representative sequences from each plant species that
298    passed this selection were combined into a single file and appended to the initial RT library to be
299    clustered with UCLUST using an identity threshold of 55%. At this level of similarity, aa RT sequences
300    from every genus in the *Caulimoviridae* is separated except those from *Cavemovirus* and
301    *Solendovirus*.
302

303    Starting with the first cluster, one or more sequences presenting high quality alignment and
304    containing several conserved residues as determined contextually for each cluster were then
305    manually selected to be representative of the diversity observed within each cluster. The following
306    clusters were processed similarly while keeping the representative sequences selected from
307    previously processed clusters. Clusters containing ECRT sequences from only one plant species were
308    analyzed only when they contained at least three sequences. After processing each cluster
309    individually, a total of 56 ECRT sequences detected here and 20 RT from known genera have been
310    selected for their remarkable divergence. Together with four RT sequences from Ty3/Gypsy LTR
311    retrotransposons, these combined sequences (hereafter referred to as "diverse library") were

312  aligned with the GUIDANCE2 [26] program using MAFFT [27] to generate bootstrap supported MSA and to
313  remove columns (--colCutoff ) with confidence score below 0.95 (16/244 columns removed in the RT
314  sequence from Caulimovirus CaMV). The resulting MSA was then used to build the phylogenetic
315  network shown in Figure 1 and Supplementary Figure 1 with SplitsTree4 [28] applying the NeighborNet
316  method with uncorrected P distance model and 1,000 bootstrap tests. Manual analysis of this
317  network enabled the discrimination of 17 distinct groups sharing deep connections among
318  caulimovirid sequences.

320  In response to the discovery of several novel OTUs, we repeated ECRT mining in plant genomes using
321  the diverse library as query. This second search is also designed to be more sensitive as it takes into
322  account DNA sequences instead of uninterrupted ORFs. The workflow is identical to the one
323  employed for the initial search until obtaining the set of extended hit loci. These were directly
324  compared to the diverse library using BLASTx with default parameters (except –e=1e-5). Queries with
325  best alignment score against any Caulimoviridae with an alignment length above 80% of subject
326  length (set generically to 576 bp considering an average size of RT domains of 240 aa) were selected
327  for phylogenetic placement.

329  **Phylogenetic analysis**
330  Fragments of virus sequence were assembled using CodonCode aligner 6.0.2 using default settings or
331  using VECTOR NTI Advance 10.3.1 (Invitrogen) operated using default settings, except that the values
332  for maximum clearance for error rate and maximum gap length were increased to 500 and 200,
333  respectively.

335  Phylogenetic reconstruction was performed using the contiguous nucleotide sequences
336  corresponding to the protease, reverse transcriptase and ribonuclease H domains. Whole sequences
337  from Caulimoviridae genera representatives and Ty3 and Gypsy LTR retrotransposons were first
338  aligned with global method using MAFFT v7.3 [27]. The core genomes was extracted and re-aligned by
339  local method using MAFFT. The resulting alignment was tested for different evolutionary models
340  with pmodeltest v1.4 (from ETE 3 package [29]) which inferred the GTRGAMMA model. Phylogenetic
341  inference with maximum likelihood was then performed using RaxML v8.2 [30] under the predicted
342  model with 500 ML bootstrap replicates.

344  The resulting tree was then used as a reference to classify the ECRT loci mined from plant genomes.
345  We first added query sequences from each plant species separately to the reference alignment and
346  aligned each library using Mafft v7.3 (with options --addfragment, --keeplength and by reordering).
347  We then tested the most likely placement of each ECRT sequence on to the reference tree using
348  pplacer v1.1 alpha19 [20] with the option (--keep-at-most 1) with allows to keep one placement for
349  each query sequence. The python package Taxit was used to construct a reference package which we
350  used to run pplacer.

352  **Data availability**
353  The datasets and scripts generated during during the current study are available from the
354  corresponding author on request.

**References**

1    Roossinck, M. J., *Virus Res* (2016).

2    Aiewsakun, P. and Katzourakis, A., *Virology* **479-480**, 26 (2015).

3    Hayward, A., Grabherr, M., and Jern, P., *Proc Natl Acad Sci U S A* **110** (50), 20146 (2013).

4    Hayward, A., Cornwallis, C. K., and Jern, P., *Proc Natl Acad Sci U S A* **112** (2), 464 (2015).

5    Aiewsakun, P. and Katzourakis, A., *Nat Commun* **8**, 13954 (2017).

6    Bejarano, E. R., Khashoggi, A., Witty, M., and Lichtenstein, C., *Proc Natl Acad Sci U S A* **93** (2), 759 (1996).

7    Teycheney, P. Y. and Geering, A. D., in *Recent advances in plant virology* (Caister Academic Press, Norfolk, 2011), pp. 343.

8    Pringle, C. R., *Arch Virol* **143** (1), 203 (1998).

9    Mushegian, A. R. and Elena, S. F., *Virology* **476**, 304 (2015).

10   Geering, A. D. et al., *Nat Commun* **5**, 5269 (2014).

11   Xiong, Y. and Eickbush, T. H., *EMBO J* **9** (10), 3353 (1990).

12   Hansen, C. and Heslop-Harrison, J. S., *Advances in Botanical Research* **41**, 165 (2004).

13   Kenrick, P., *Philos Trans R Soc Lond B Biol Sci* **355** (1398), 847 (2000).

14   Wolf, P. G. et al., *Genome Biol Evol* **7** (9), 2533 (2015).

15   Matasci, N. et al., *Gigascience* **3**, 17 (2014).

16   Wickett, N. J. et al., *Proc Natl Acad Sci U S A* **111** (45), E4859 (2014).

17   Diaz-Lara, A. and Martin, R. R., *The American Phytopathological Society* **100** (11), 2211 (2016).

18   Lockhart, B., Mollov, D., Olszewski, N., and Goldsmith, N., *Virus Res* **In Press, Corrected Proof.** (2017).

19   Geering, A. D., Scharaschkin, T., and Teycheney, P. Y., *Arch Virol* **155** (1), 123 (2010).

20   Matsen, F. A., Kodner, R. B., and Armbrust, E. V., *BMC Bioinformatics* **11**, 538 (2010).

21   Link, K. and Sonnewald, U., in *Plant-Virus Interactions: Molecular Biology, Intra- and Intercellular Transport*, edited by T. Kleinow (Springer International Publishing, Cham, 2016), pp. 1.

22   Brunkard, J. O. and Zambryski, P. C., *Curr Opin Plant Biol* **35**, 76 (2017).

23   Hedges, S. B., Dudley, J., and Kumar, S., *Bioinformatics* **22** (23), 2971 (2006).

24   Edgar, R. C., *Bioinformatics* **26** (19), 2460 (2010).

25   Capella-Gutierrez, S., Silla-Martinez, J. M., and Gabaldon, T., *Bioinformatics* **25** (15), 1972 (2009).

26   Sela, I., Ashkenazy, H., Katoh, K., and Pupko, T., *Nucleic Acids Res* **43** (W1), W7 (2015).

27   Katoh, K. and Standley, D. M., *Mol Biol Evol* **30** (4), 772 (2013).

28   Huson, D. H. and Bryant, D., *Mol Biol Evol* **23** (2), 254 (2006).

29   Huerta-Cepas, J., Serra, F., and Bork, P., *Mol Biol Evol* **33** (6), 1635 (2016).

30   Stamatakis, A., *Bioinformatics* **30** (9), 1312 (2014).

396 **Acknowledgements**

399

400

401 **Author contributions**

402 SD, ADW, and FM performed research; FAD and ML provided research environment; SD, ADW, PYT,
403 and FM analyzed the data; ADW, PYT, and FM wrote the manuscript with contributions from SD; FM
404 designed the research.

405

406

407 **Competing interests**

408 The authors declare that they have no competing interests.

409

410

411 **Supplementary material**

412 Supplementary file 1: reconstructed ancestral sequences of members of the *Caulimoviridae* from
413 novel OTUs (fasta format) and onekp contigs.

414

415

416 **Figure legends**

417 Figure 1: Augmented diversity of the *Caulimoviridae*. Core of a phylogenetic network constructed
418 using an alignment of amino acid reverse transcriptase (RT) sequences from reference genera,
419 representative endogenous caulimovirid RTs (ECRTs) and Ty3/Gypsy LTR retrotransposons. The full
420 network is available in Supplementary Figure 1. This representation allows determining 17
421 Caulimoviridae OTUs. OTU names have dashed lime green outline when they include no known
422 reference genera (referred to as novel OTUs). Each fill color corresponds to a different OTU except
423 for OTUs comprising only a representative ECRT sequence that are colored with dark grey and named
424 after the only host plant genome they were detected in at this stage (Petunia-, Vitis-, and Glycine-
425 virus). * RT clustering at 55% identity groups *Cavemovirus* and *Solendovirus* into a single OTU (OTU
426 8). ** Sequences grouped in the Xendovirus OTU appeared to be paraphyletic after phylogenetic
427 reconstruction (see Figure 3).

428

429 Figure 2: Highly variable ECRT numbers and density across plants. (A) Number of ECRTs found in each
430 plant genome as function of Log10 genome size expressed in megabases (assembly gaps excluded).
431 Logarhitmic trendline indicates moderate correlation between the number of ECRT and genome size
432 ($R^2$=0.544). (B) Density of ECRTs per megabase in each plant genome as function of Log10 genome
433 size expressed in megabases (assembly gaps excluded). In (A) and (B), arrows indicate a sample of
434 outlier dots and the corresponding plant species name.

435

436 Figure 3: Phylogeny of the *Caulimoviridae*. Phylogenetic tree obtained by maximum likelihood search
437 from a multiple sequence alignment of the genomic regions containing protease, reverse
438 transcriptase and ribonuclease H1 domains from known (black) and novel (red) Caulimoviridae
439 genera. The sequences from Ty3/Gypsy LTR retrotransposons are used as outgroups. Bootstrap
440 support values below 50% are not shown. Sequences from members of the novel genera are

11

441 available in supplementary data. Closely related sequences were collapsed into branches. The
442 sequences contained in each branch are as follows. Orendovirus: Aegilops tauschii virus (AtV),
443 Brachypodium distachyon virus (BdV); *Tungrovirus*: *Rice tungro bacilliform virus* (RTBV), Rice tungro
444 bacilliform virus isolate west Bengal (RTBV); *Badnavirus*: *Commelina yellow mottle virus* (ComYMV),
445 *Banana streak OL virus* (BSOLV); Yendovirus: Capiscum annuum virus; Zendovirus: Fragaria vesca
446 virus; Blueberry: Blueberry fruit drop associated virus (BFDaV); *Caulimovirus*: *Cauliflower mosaic virus*
447 (CaMV), *Figwort mosaic virus* (FMV); Rudbeckia: Rudbeckia flower distortion virus (RuFDV);
448 *Soymovirus*: *Soybean chlorotic mottle virus* (SoyCaulimoviridae), *Peanut chlorotic streak virus* (PCSV);
449 *Solendovirus*: *Sweet potato vein clearing virus* (SPVCV), *Tobacco vein clearing virus* (TVCV);
450 *Cavemovirus*: *Cassava vein mosaic virus* (CsVMV), *Sweet potato collusive virus* (SPCV); *Petuvirus*:
451 *Petunia vein clearing virus* (PVCV); *Rosadnavirus*: *Rose yellow vein virus* (RYVV); Florendovirus:
452 Fragaria vesca virus (FvesV), Mimulus guttatus virus (MgutV); Gymnendovirus 1: Pinus taeda
453 Gymnendovirus 1, Picea glauca Gymnendovirus 1; Gymnendovirus 2: Pinus taeda Gymnendovirus 2,
454 Picea glauca Gymnendovirus 2, Ginkgo biloba Gymnendovirus 2; Gymnendovirus 3: Pinus taeda
455 Gymnendovirus 3; Gymnendovirus 4: Pinus taeda Gymnendovirus 4, Picea glauca Gymnendovirus 4;
456 Fernendovirus 1: Cystopteris protrusa Fernendovirus 1 contig 1, and the transcript scaffolds BEGM-
457 2004510 from *Botrypus virginianus*, NOKI-2097008 from *Lindsaea linearis*, and ENQF-2084799 from
458 *Lycopodium annotinum*; Fernendovirus 2: Dipteris conjugata Fernendovirus 2 Contigs 2, 4 and 1319.
459
460 Figure 4: Distribution of caulimovirid EVEs in Euphyllophyte. The left tree represents a cladogram of
461 Euphyllophyte species investigated in this study. The name of major branches and nodes is indicated.
462 The top tree represents the topology of the phylogenetic tree obtained in Figure 3. At the
463 intersection of these two trees, color code indicates the number of ECRT loci classified into each
464 Caulimoviridae genus for each plant species. Abbreviations of virus genera are as follows: Pe
465 (Petuvirus), Gy1 (Gymnendovirus 1), Gy2 (Gymnendovirus 2), Gy3 (Gymnendovirus 3), Gy4
466 (Gymnendovirus 4), Fe1 (Fernendovirus 1), Fe2 (Fernendovirus 2), Flo (Florendovirus), Soy
467 (*Soymovirus*), Rud (Rudbeckia flower distortion virus), Cau (*Caulimovirus*), Blu (Blueberry fruit drop-
468 associated virus), Zen (Zendovirus), Xen (Xendovirus), Yen (Yendovirus), CaS (*Cavemovirus*
469 +*Solendovirus*), Ros (Rosadnavirus), Bad (*Badnavirus*), Tun (*Tungrovirus*), Ore (Orendovirus).
470
471 Figure 5: Working scenario of *Caulimoviridae* deep evolution. The left tree is the same as in Figure 3
472 where the deepest *Caulimoviridae* node was annotated as LCA (last common ancestor). The top
473 cladogram indicates the evolutionary relationships between major classes of vascular plants. At the
474 intersection between both trees, color code indicates the presence of EVEs (green) and of known
475 associations with exogenous viruses (blue).
476
477 Supplementary Figure 1: Overview of the phylogenetic network used to build Figure 1.
478
479 Supplementary Figure 2: ECRT ORFs collected from ferns cluster as two novel OTUs. Representative
480 sequences identified in fern genomes were appended to the collection of sequences represented in
481 Figure 1. The resulting library has been re-aligned with MUSCLE and phylogenetic network was built
482 using SplitsTree. The branches containing fern sequences have been empirically grouped into two
483 novel OTUs (OTU 18 and OTU 19).
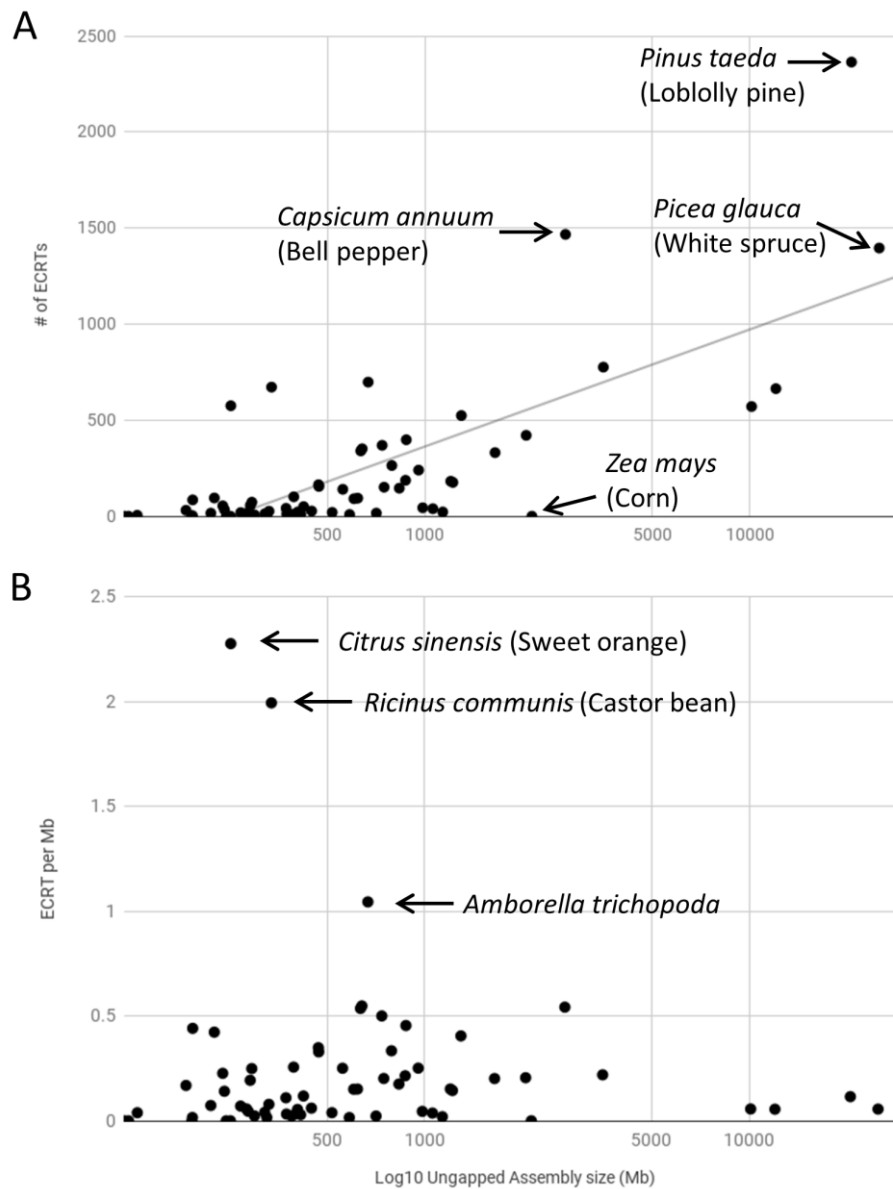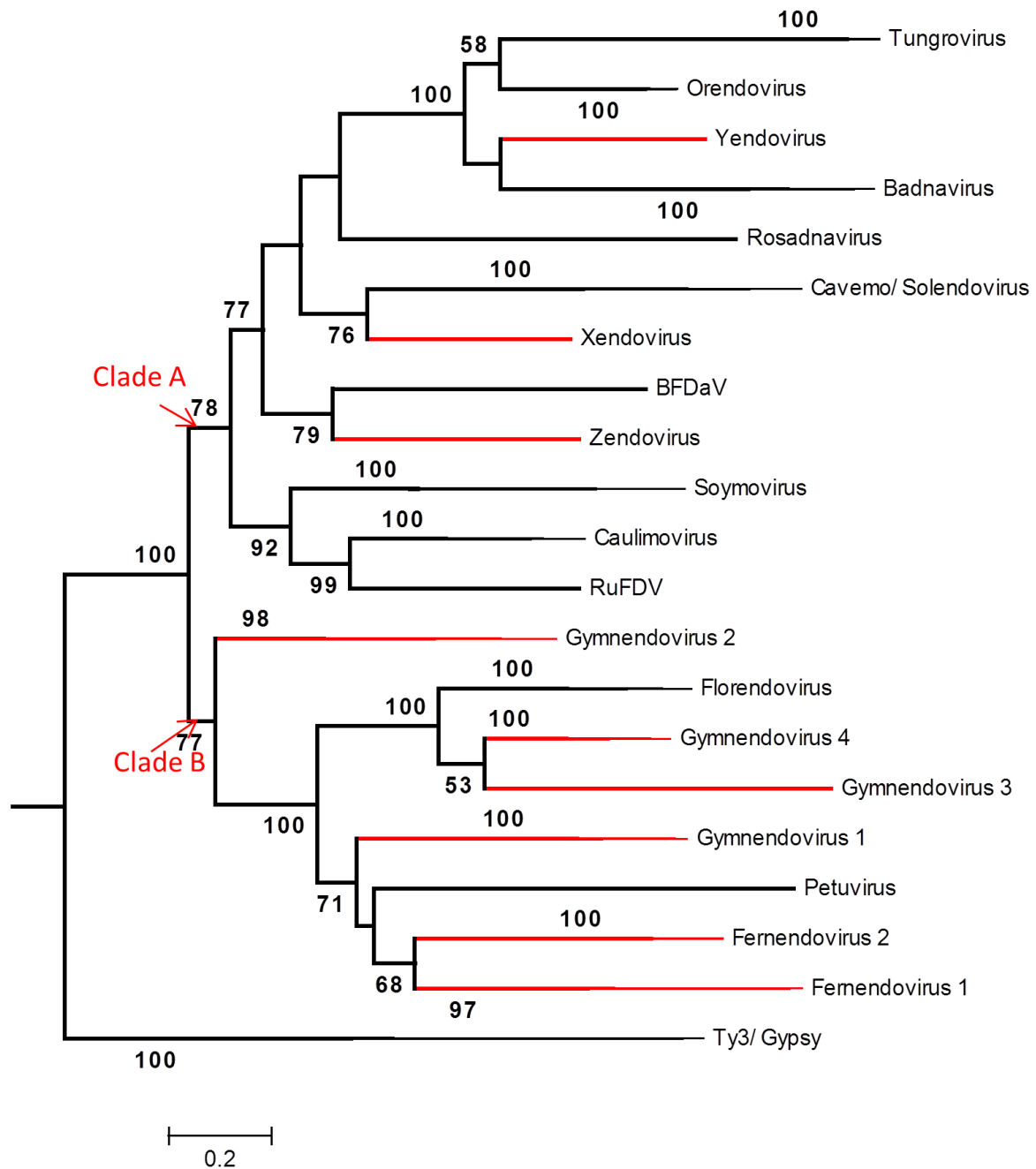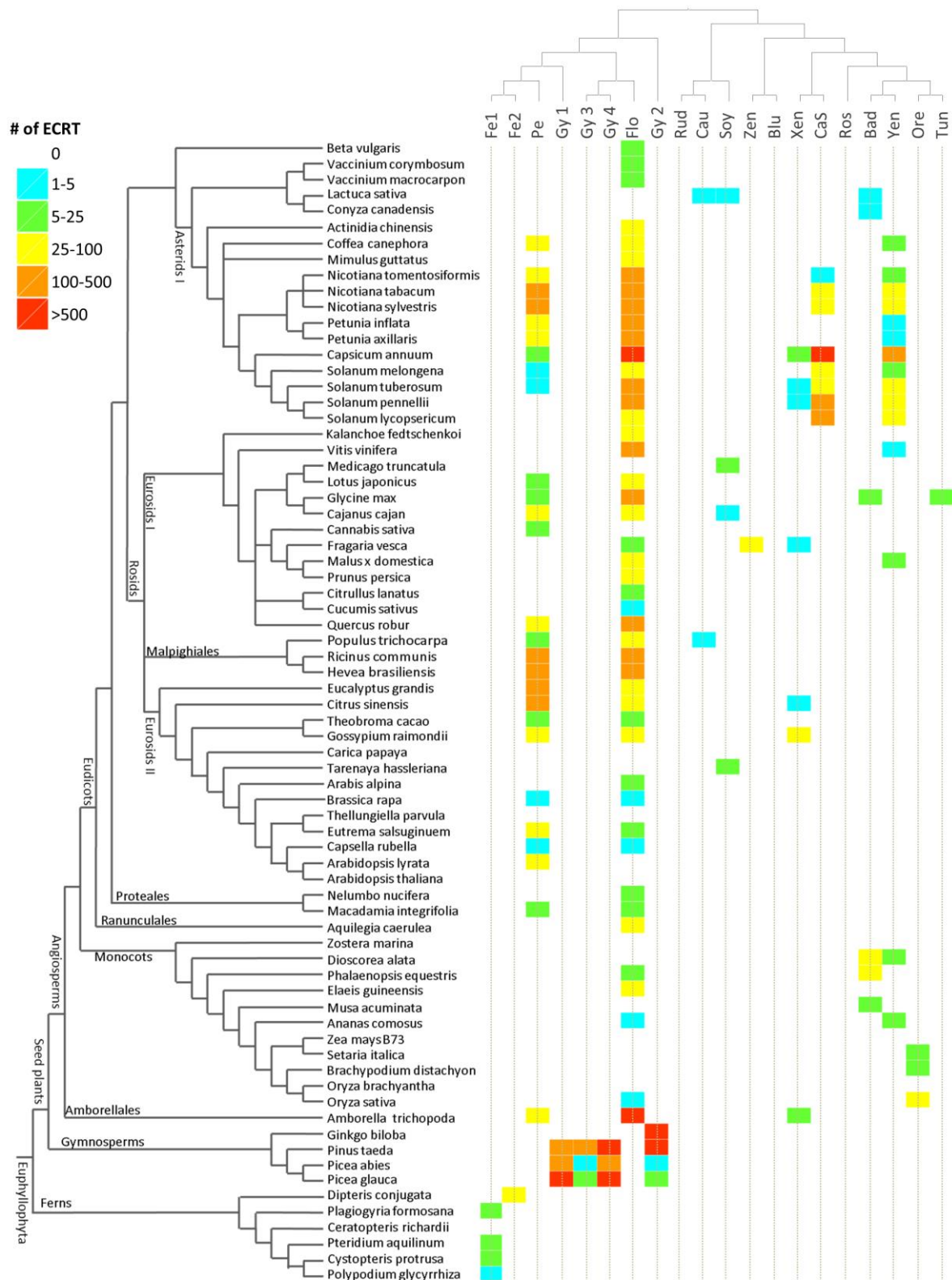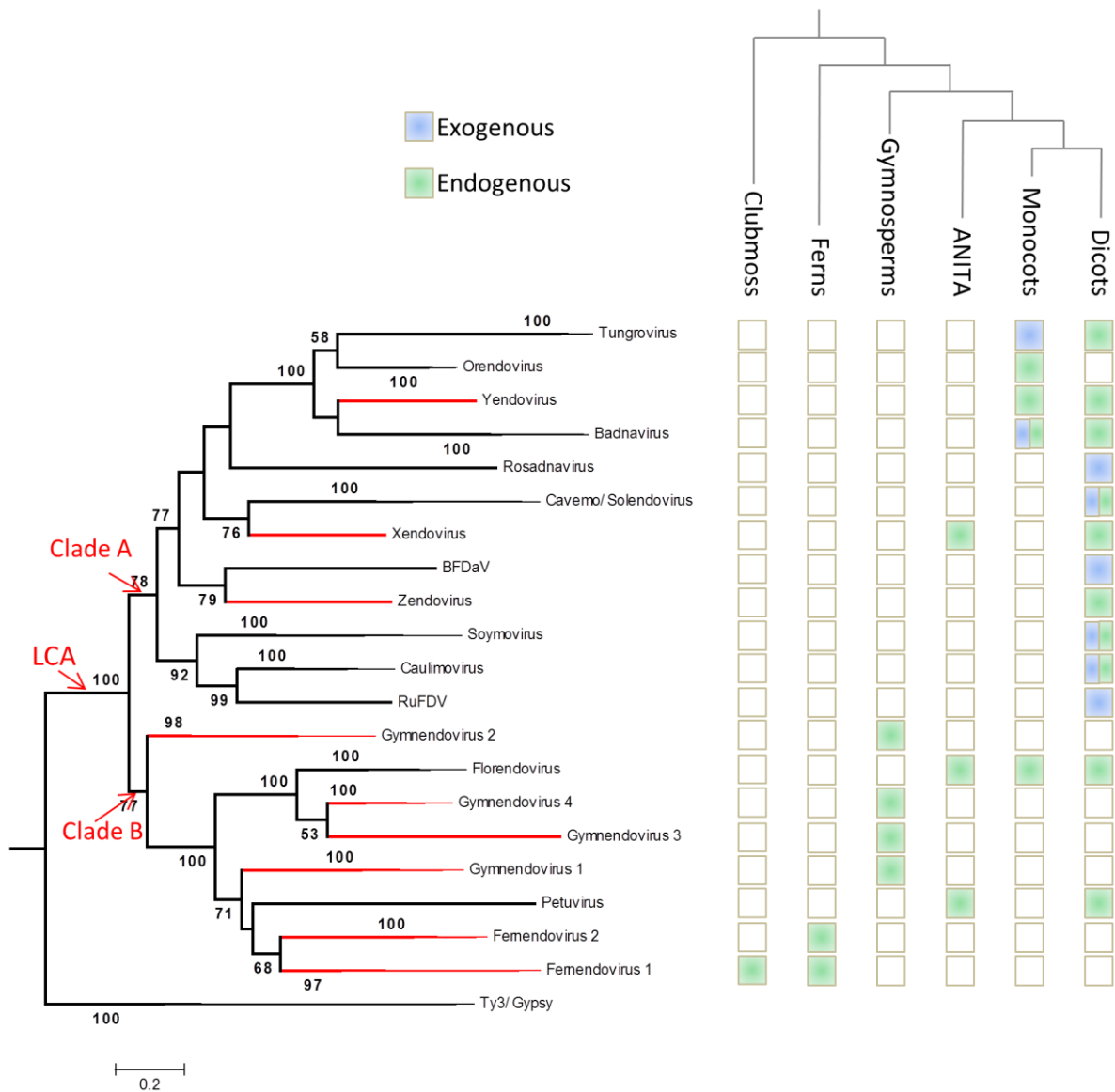484
485

12

**Figures**

Figure 1

Figure 2

Figure 3

Figure 4

Figure 5

**Supplementary figures**

Supplementary figure 1

Supplementary figure 2