



Multiple waves of viral invasions in Symbiodiniaceae algal genomes

Journal:	<i>Virus Evolution</i>
Manuscript ID	VEVOLI-2022-084
Manuscript Type:	Research Article
Date Submitted by the Author:	09-Apr-2022
Complete List of Authors:	de Almeida Benites, Luiz Felipe; Rutgers University New Brunswick, Department of Biochemistry and Microbiology Stephens, Timothy G.; Rutgers University New Brunswick, Department of Biochemistry and Microbiology Bhattacharya, Debashish; Rutgers University New Brunswick, Department of Biochemistry and Microbiology
Keywords:	Symbiodiniaceae evolution, virus horizontal gene transfer, dinoflagellate virus, Symbiodinium virus, algal virus

SCHOLARONE™
Manuscripts

Multiple waves of viral invasions in Symbiodiniaceae algal genomes

Author list

L. Felipe Benites¹

Timothy G. Stephens¹

Debashish Bhattacharya¹

Affiliations

¹ Department of Biochemistry and Microbiology, Rutgers University, New Brunswick, NJ
08901, USA

Corresponding author

L. Felipe Benites - Department of Biochemistry and Microbiology, Rutgers University, New
Brunswick, NJ 08901, USA (l.felipebenites@gmail.com)

Multiple waves of viral invasions in Symbiodiniaceae algal genomes

Author list

L. Felipe Benites¹

Timothy G. Stephens¹

Debashish Bhattacharya¹

Affiliations

¹ Department of Biochemistry and Microbiology, Rutgers University, New Brunswick, NJ
08901, USA

Abstract

Dinoflagellates from the family Symbiodiniaceae are phototrophic marine protists that engage in symbiosis with diverse hosts. Their large and distinct genomes show pervasive gene duplication and large-scale retroposition events. However, little is known about the role and scale of horizontal gene transfer (HGT) in the genomic evolution of this algal family. In other dinoflagellates, higher levels of HGTs have been observed, linked to major genomic transitions, such as the appearance of a viral acquired nucleoprotein that originated *via* HGT from a large DNA algal virus. Previous work showed Symbiodiniaceae from different hosts being actively infected by several viral groups, such as giant DNA viruses and ssRNA viruses, that may play an important role in coral health. This includes a hypothetical latent viral infection, whereby viruses could persist in the cytoplasm or integrate into the host genome as a provirus. This hypothesis received some experimental support however, the cellular localization of putative latent viruses

1
2
3 and their taxonomic affiliation are still unknown. In addition, despite the finding of viral
4 sequences in some genomes of Symbiodiniaceae, viral origin, taxonomic breadth, and metabolic
5 potential have not been explored. To address these questions, we searched for evidence of
6 protein sequences of putative viral origin in 13 Symbiodiniaceae genomes. We found 59
7 candidate viral-derived HGTs that give rise to 12 phylogenies across 10 genomes. We also
8 describe the taxonomic affiliation of these virus-related sequences, their structure, and genomic
9 context. These results lead us to propose a model to explain the origin and fate of
10 Symbiodiniaceae viral acquisitions.
11
12
13
14
15
16
17
18
19
20
21
22
23

24 **Keywords**

25
26 Symbiodiniaceae evolution; virus horizontal gene transfer; dinoflagellate virus; *Symbiodinium*
27 virus; algal virus
28
29
30
31
32

33 **1. Introduction**

34
35 Dinoflagellates from the family Symbiodiniaceae are phototrophic marine protists that
36 engage in symbiosis with diverse hosts such as foraminifera, ciliates, radiolarians, mollusks, and
37 cnidarians such as anemones and most notably corals (LaJeunesse et al. 2018). Their large and
38 distinct genomes show pervasive gene duplication (Aranda et al. 2016; Nand et al. 2021),
39 expansion of lineage specific gene families with unknown functions (González-Pech et al. 2021),
40 and large-scale retroposition events (Song et al. 2017). However, little is known about the role of
41 horizontal gene transfer (HGT) in Symbiodiniaceae evolution and the scale of gene acquisition
42 events (González-Pech et al. 2021). In *Fugacium kawagutii* and *Breviolum minutum* (formerly
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 *Symbiodinium* clades F and B), ~0.7% of their genes (246 and 244 sequences respectively), have
4
5 putatively originated from prokaryote-derived HGT events (Fan et al. 2020).
6

7
8 Higher levels of HGTs have been observed in other dinoflagellates and appear to be
9
10 linked to major genomic transitions (Wisecaver et al. 2013; Janouškovec et al. 2017).
11
12 Remarkably, all sequenced dinoflagellate genomes contain a family of viral acquired
13
14 nucleoproteins, known as the dinoflagellate/viral nucleoproteins (DVNPs), which have a putative
15
16 role in chromatin regulation (Irwin et al. 2018). In *B. minutum*, 19 divergent DVNPs have been
17
18 identified and it has been suggested that they are involved in regulating chromosome structure
19
20 (Shoguchi et al. 2013). Furthermore, it has been proposed that the DVNP protein family
21
22 originated in the ancestor of dinoflagellates by HGT from a large DNA algal virus, possibly from
23
24 the family *Phycodnaviridae* (Gornik et al. 2012). Acquisition of the DVNP gene family
25
26 coincides with massive genomic expansion in this group and may have conferred immunity to
27
28 infection by the same group of exogenous viruses (Gornik et al. 2019).
29
30
31
32

33 Previous work has found cases of Symbiodiniaceae isolated from corals being actively
34
35 infected by several viral groups, ranging from giant DNA viruses to ssRNA viruses (A.M.S.
36
37 Correa et al. 2013; Weynberg et al. 2014; Wood-Charlson et al. 2015; A.M.S. Correa et al. 2016;
38
39 Levin et al. 2017; Messyasz et al. 2020). Most importantly, the presence and abundance of viral
40
41 sequences were found to be associated with bleached and diseased corals, suggesting that viruses
42
43 may play a role in coral health and bleaching events by infecting (subsequently killing or
44
45 disrupting the function of) Symbiodiniaceae cells in the host (Thurber and A.M.S. Correa 2011).
46
47 Recently, Grupstra et al. (2022) demonstrated the rapid formation of diverse viral major capsid
48
49 protein (MCP) “aminotypes” (unique amino acid sequences) from a Symbiodiniaceae-infecting
50
51 dinoflagellate RNA virus in coral colonies exposed to high temperatures (30.3 °C) that are
52
53
54
55
56
57
58
59
60

1
2
3 known to cause coral bleaching. This suggests that the switch from a persistent to productive
4
5 viral infection was triggered by exposure to high temperatures.
6

7
8 Early observations linking symbiotic dinoflagellates from anemones with viruses,
9
10 stresses, and coral diseases (Wilson et al. 2001) led to the hypothesis that the algal cells harbored
11
12 latent proviruses. In this type of infection, a virus persists in the cytoplasm as an episome (an
13
14 extrachromosomal molecule) or becomes integrated into the host genome, replicating as a
15
16 (latent) provirus, synchronized with the host cell (Hyman and Abedon 2012). After exposure to
17
18 stress, the latent provirus virus enters a lytic cycle, lysing the host cell and propagating into other
19
20 susceptible hosts. This hypothesis was partially corroborated by experiments using multiple
21
22 Symbiodiniaceae isolates from various cnidarian hosts, showing the production of viral-like
23
24 particles (VLPs) in thermally and ultraviolet (UV) light stressed algal cells (Wilson et al. 2005;
25
26 Davy et al. 2006; Lohr et al. 2007; Lawrence et al. 2014; Lawrence et al. 2017; Weynberg et al.
27
28 2017; Benites et al. 2018). However, the cellular localization of putative latent viruses (i.e.,
29
30 cytoplasmic or genomic) and their precise taxonomic affiliations is still unknown.
31
32
33
34

35
36 In other algal groups, such as the brown algae (Phaeophyceae), double-stranded DNA
37
38 viruses from *Phycodnaviridae* can persist as integrated proviruses with fragments of the viral
39
40 genome, interspersed as repeats and pseudogenes throughout the host genome. The viral genome
41
42 remains latent until it becomes active in reproductive cells after light and temperature stimulation
43
44 (Bräutigam et al. 1995: 1; Müller et al. 1998; Delaroque and Boland 2008; McKeown et al.
45
46 2017). In *Ectocarpus* (strain Ec 32), the genome of a virus was found to be integrated in the host
47
48 genome however, the viral genes were not expressed and viral particles were not produced (Cock
49
50 et al. 2010). This is similar to observations in *Feldmannia*, which has viral fragments in its
51
52 genome but does not show evidence of viral production (Lee et al. 1998). In the red alga
53
54
55
56
57
58
59
60

1
2
3 (Rhodophyta) *Chondrus crispus*, multiple copies of a partial RNA viral genome were found to be
4 expressed, including a copy of the capsid gene, but there was no apparent virus production or
5 host symptoms (Rousvoal et al. 2016). In green algae (Chlorophyta), although some of their
6 genomes contain numerous giant “endogenous viral elements” (EVEs) with spliceosomal introns
7 and segmental gene duplications (Moniruzzaman et al. 2020), the only evidence thus far of viral
8 latency is from the chlorophyte *Cylindrocapsa geminella*, with production of VLPs and lytic
9 infection observed after heat-shock treatment (Hoffman and Stanker 1976).

10
11
12
13
14
15
16
17
18
19 Recently, Irwin et al. (2022) reported that at least 23 genes in Symbiodiniaceae (Fkaw,
20 Sgor and Smic) show strong support for HGTs with viruses such as Nucleocytoviricota giant
21 viruses (although their analysis did not focus on Symbiodiniaceae), with possible roles in
22 catalysis of ATP-dependent conversion of ribonucleotides, HNH endonuclease, among others. In
23 addition, despite the fact that viral sequences account for < 3% of the *Symbiodinium*
24 *microadriaticum* genome (Aranda et al. 2016) and ~0.1% of the *Cladocopium goreau* genome
25 (Liu et al. 2018), their origin, taxonomic breadth, and metabolic potential has not been explored in
26 depth and there has been no previous evidence of viral genome integration into Symbiodiniaceae
27 nuclear DNA. The latter finding would strongly support the proviral hypothesis (Wilson et al.
28 2001).

29
30
31
32
33
34
35
36
37
38
39
40
41
42 Therefore, to address questions about the occurrence of integrated viral sequences in
43 Symbiodiniaceae genomes, and to uncover their possible origins and functional profiles, we
44 searched predicted proteins from the 13 available Symbiodiniaceae genomes (Shoguchi et al.
45 2013; Lin et al. 2015; Liu et al. 2018; Shoguchi et al. 2018; Chen et al. 2020; González-Pech et
46 al. 2021) for sequences of putative viral origin. We identified 59 candidate viral HGT sequences
47 that formed 12 phylogenies, across 10 genomes. We describe for each of these sequences their
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 taxonomic affiliation with viruses and the composition, structure, and genomic context of their
4 associated genes. Moreover, we searched for corroborative evidence that would support these
5 sequences being HGT-derived and for the presence of viral genome integrations.
6
7
8
9

10 11 12 **2. Methods:**

13 14 15 **2.1 Screening for viral HGT events in Symbiodiniaceae protein datasets**

16
17 Each predicted protein from the 13 Symbiodiniaceae genomes (445,306 total) was
18 searched against the RefSeq viral database using DIAMOND (v0.9.22; BLASTp, sensitive
19 mode) (Buchfink et al. 2015) with an e -value cutoff of $< 1e^{-5}$, retaining just the best viral hits per
20 query sequence. To exclude false positive hits, we perform another DIAMOND (BLASTp)
21 search against the complete Genbank non-redundant database, using only the Symbiodiniaceae
22 protein queries with hits to the RefSeq viral database. For each query, the top 400 hits were
23 retained and filtered using an e -value of $< 1e^{-5}$. For each hit, the full taxonomy of the subject
24 sequence was retrieved for the top 50 filtered hits using a custom perl script (available at:
25 https://github.com/LFelipe-B/Symbiodiniaceae_vHGT_scripts) which uses the subject accession
26 number to fetch the whole taxonomic tree from the Genbank taxonomy suit. Query sequences
27 that had at least one of their top 50 hits to a viral sequence were retained for downstream
28 analysis. Genomes are abbreviated as follows: *Breviolum minutum* (Bmin); *Cladocopium* sp.
29 C92 (CC92); *Cladocopium goreau* (Cgor); *Fugacium kawagutii* (Fkaw); *Symbiodinium*
30 *linucheae* CCMP2456 (Slin_CCMP2456); *S. microadriaticum* (Smic); *S. microadriaticum* 04-
31 503SCI.03 (Smic_04-503SCI.03); *S. microadriaticum* CassKB8 (Smic_CassKB8); *S. natans*
32 CCMP2548 (Snat_CCMP2548); *S. necroappetens* CCMP2469 (Snec_CCMP2469); *S. pilosum*
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 CCMP2461 (Spil_CCMP2461); *S. tridacnidorum* (Stri); and *S. tridacnidorum* CCMP2592
4
5 (Stri_CCMP2592).
6
7
8
9

10 **2.2 Inferring pre-candidates for viral HGTs**

11
12 To determine the scale of viral HGT in Symbiodiniaceae genomes we collected
13
14 sequences that showed significant similarity to viral sequences by calculating four independent
15
16 metrics using custom python scripts and the full taxonomic annotations retrieved using the
17
18 subject accession numbers from the DIAMOND outputs. Each hit to the complete Genbank non-
19
20 redundant database was categorized as “non-target” if the subject sequence was from a virus and
21
22 as “target” if the subject sequence was from a eukaryote; hits to Symbiodiniaceae sequences
23
24 were omitted from the subsequent calculations to avoid hits to sequences from the same or
25
26 similar species already submitted to GenBank from biasing the analysis. The first metric
27
28 calculated was the Alien index (AI) (Gladyshev et al. 2008; Rancurel et al. 2017), which assesses
29
30 how many orders of magnitude there is between the *e*-values of the top target and non-target hits;
31
32 query sequences were retained if non-target sequences scored $AI > 30$. The second metric
33
34 calculated was the HGT index (hU) (Boschetti et al. 2012), defined as the difference
35
36 between the bitscores of the top non-target and target hits; query
37
38 sequences with a $Hu \geq 30$ were retained. The third metric was
39
40 calculated by counting the total number of target and non-target
41
42 hits (collapsing repeated taxa) associated with each query; query
43
44 sequences were retained if the sum of non-target hits was $> 50\%$ of
45
46 the total number of hits. Finally, the fourth metric was calculated by
47
48 taking the sum of the non-target hit bitscores and comparing it with
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 the sum of the target hit bitscores, retaining query sequences which
4
5 had higher non-target bitscore totals compared to their target
6
7 bitscore totals. Query sequences with scores above the cut offs
8
9 associated with each of the four metrics were considered pre-
10
11 candidates of viral HGTs (pre-vHGTs).
12
13
14
15
16
17

18 **2.3 Orthogroup clustering of pre-vHGTs and Phylogenetic reconstruction**

19
20 All Symbiodiniaceae proteins were clustered into orthogroups using OrthoFinder v2.3.3
21
22 (Emms and Kelly 2019), with groups containing at least one pre-vHGT sequence extracted for
23
24 further analysis. We downloaded the subject sequences associated with each of the pre-vHGT
25
26 DIAMOND hits to the RefSeq viral and Genbank non-redundant databases (accession numbers
27
28 in Supplementary Table S11). The Symbiodiniaceae proteins in each pre-vHGT-containing
29
30 orthogroup were combined with the downloaded subject sequences and complemented with
31
32 subject hits downloaded from a BLASTp search (<https://blast.ncbi.nlm.nih.gov>) to retrieve the
33
34 maximum number of subject hits associated with the pre-vHGTs for each group. The combined
35
36 set of proteins associated with each group were aligned using MAFFT v7.305b (L-INS-i
37
38 algorithm: --localpair --maxiterate 1000) (Kato and Standley 2013). The resulting alignments
39
40 were trimmed using TrimAI v1.2 in automated mode (-automated1) (Capella-Gutiérrez et al.
41
42 2009). A Maximum Likelihood phylogenetic tree was inferred from each group alignment using
43
44 IQ-TREE v2.0.3 (Nguyen et al. 2015), with the best evolutionary model selected by
45
46 ModelFinder (Kalyaanamoorthy et al. 2017) (-m TEST) and 10,000 ultrafast bootstrap replicates
47
48 calculated using UFBoot (Minh et al. 2013). Phylogenetic trees were midpoint rooted and
49
50 screened for topologies in which Symbiodiniaceae query sequences (target) were clustered or
51
52
53
54
55
56
57
58
59
60

1
2
3 nested with viral sequences (non--target) in a clade with > 90% ultrafast bootstrap support using
4
5 PhySortR (Stephens et al. 2016) (clade.exclusivity = 0.95, min.prop.target = 0.7). Phylogenies
6
7 were further processed using the packages ggtree (Yu et al. 2017), phangorn (Schliep 2011), and
8
9 phytools (Revell 2012) in the R environment (version 3.4.2). A final manual curation step was
10
11 implemented, discarding orthogroups with ambiguous phylogenetic signals: cases where there
12
13 was an over representation of non-viral sequences or excessively long branch lengths were
14
15 considered as having weak evidence of viral HGT and were discarded. The phylogenetic groups
16
17 that remained after filtering were considered to contain Symbiodiniaceae sequences (hereinafter
18
19 referred to as vHGTs) with strong evidence of having originated from viral HGT.
20
21
22
23
24
25

26 **2.4 Functional analysis of Symbiodiniaceae vHGTs**

27
28 To gain insights into the functional landscape of viral acquired sequences, we analyzed
29
30 protein family membership and gene ontology (GO) terms for the vHGTs by re-annotating
31
32 individual sequences using the InterProScan online suite (<https://www.ebi.ac.uk/interpro/>) and
33
34 QuickGO (<https://www.ebi.ac.uk/QuickGO/annotations>). In some instances, we also subjected
35
36 the vHGTs to a BLASTp search against the UniProtKB database
37
38 (<https://www.uniprot.org/uniprot/>).
39
40
41
42
43
44

45 **2.5 Compositional and structural analysis of Symbiodiniaceae vHGTs**

46
47 We retrieved information about the GC-content, length of coding sequences (CDSs),
48
49 presence of introns, intron length, and location along the scaffold of the gene models associated
50
51 with the vHGTs. For each of these features, the difference between the mean values calculated
52
53 for the vHGTs was compared against a background of all Symbiodiniaceae genes using the
54
55
56
57
58
59
60

1
2
3 Welch Two Sample t-test in R. The presence of a dinoflagellate spliced leader (DinoSL)
4
5 sequence upstream of the first exon of the vHGT genes in the genome was assessed to identify if
6
7 mRNA recycling (Slamovits and Keeling 2008) has played a role in the evolution of these genes.
8
9
10 DinoSL and relic DinoSL sequences were identified by an approach similar to that used by
11
12 (Stephens et al. 2020). Briefly, the DinoSL sequence (CCGTAGCCATTTTGGCTCAAG) and
13
14 relic DinoSL sequences (which are composed of two or more DinoSL sequences joined together
15
16 at their canonical splice sites) were searched against each Symbiodiniaceae genome using
17
18 BLASTn v2.10.1 (-max_target_seqs 1000000 -task blastn-short -evaluate 1000). Hits were
19
20 retained if they started ≤ 5 bp from the 5'-end of the query sequence (which is the position of the
21
22 conserved canonical splice site) and ≤ 2 bp from the 3'-end of the query. The proximity of the
23
24 identified DinoSL and relic DinoSL sequences to vHGT genes within the genome was assessed
25
26 using bedtools v2.29.2 ("bedtools closest -s -D a -id"). Only cases where a DinoSL or relic
27
28 DinoSL were identified on the same strand as the vHGT gene, were upstream or overlapping
29
30 with the vHGT, and were < 2 Kbp from the first (most upstream) exon of the vHGT gene were
31
32 reported.
33
34
35
36
37
38
39

40 **2.6 Gene expression analysis and gene model confirmation of Symbiodiniaceae vHGTs**

41
42 We retrieved RNA-seq files from the Sequence Read Archive (SRA) database
43
44 (<https://www.ncbi.nlm.nih.gov/sra>) (accession list in Supplementary Table S10) to study the
45
46 expression of vHGTs and to validate the vHGTs gene models. When possible (with the
47
48 exception of Snec, which lacks RNA-seq data), we used RNA-seq data generated from the same
49
50 isolate to verify the gene modes in the genome. The RNA-seq data was downloaded from NCBI
51
52 using Fastq-dump v2.9.6
53
54
55
56
57
58
59
60

1
2
3 (https://trace.ncbi.nlm.nih.gov/Traces/sra/sra.cgi?view=toolkit_doc&f=fastq-dump), trimmed
4 using Cutadapt v1.18 (Martin 2011) (quality-cutoff 20 and minimum-length 25), and mapped
5 against the associated reference genome using HISAT2 v2.2.1 (Kim et al. 2019). The resulting
6 alignment files were sorted and indexed using SAMtools v1.11 (Li et al. 2009), before being
7 visualized using the Integrative Genomics Viewer (IGV) v2.12.3 tool (Robinson et al. 2011).
8
9

10
11
12
13
14
15 The gene models associated with the putative sites of *Symbiodinium* +ssRNA virus
16 integration into the host genome were further analyzed for internal stop codons and frameshift
17 mutation to determine if these sequences had become pseudogenes or if they are still potentially
18 functional. For each viral +ssRNA gene of interest, the genome sequence, starting 2 kbp
19 upstream and ending 2 kbp downstream, was extracted using seqkit v0.15.0 (-u 2000 -d 2000)
20 (Shen et al. 2016). These sequences were used in alignments that included known RNA-
21 dependent RNA polymerase (YP_009337004.1 and YP_009342067.1; top hits using online
22 BLASTP against the viral gene models) and major capsid (AOG17586.1) proteins using
23 Exonerate v2.3.0 (--model protein2genome --exhaustive yes) (Slater and Birney 2005).
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39

3. Results:

3.1 Screening Symbiodiniaceae sequences using four different sequence similarity HGT

scoring metrics

40
41
42
43
44 The four HGT metrics (see Methods) were applied to Symbiodiniaceae proteins with hits
45 to viral sequences in RefSeq to produce the list of viral HGT pre-candidates (pre-vHGTs). We
46 identified additional viral HGT candidates by constructing orthogroups and taking all sequences
47 that were in an orthogroup with a pre-vHGT. The Symbiodiniaceae sequences from each group,
48 or individually if single sequences, were combined with the sequences retrieved from the top hits
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 of each pre-vHGT against a taxonomically diverse database and expanded to contain the
4
5 maximum number of subject hits with BLASTp; for each group of combined sequences a
6
7 multiple sequence alignment was created that was then used for phylogeny reconstruction. The
8
9 resulting phylogenies were filtered, retaining only those in which viruses were overrepresented
10
11 and in which viral sequences were nested with Symbiodiniaceae sequences in a clade that had
12
13 UFBoot branch support > 90%. Our workflow resulted in the identification of 59 vHGT
14
15 sequences (that comprise 12 phylogenies) that are spread across the 10 Symbiodiniaceae
16
17 genomes, ranging from n = 14 sequences in Smic to n = 1 in the Cgor and Slin_CCMP2456
18
19 genomes. The majority of vHGTs were identified in the genomes of species from the
20
21 *Symbiodinium* genus, and were absent from the genomes of CC92, Fkaw and the free living
22
23 *Spil_CCMP2461* species.
24
25
26
27
28
29
30

3.2 Taxonomic distribution of vHGTs sequences in Symbiodiniaceae

31
32
33 The taxonomic distribution of the vHGTs was assessed at all classification levels using
34
35 the top viral hit associated with each sequence and associated taxonomic information in the
36
37 NCBI database. The number of viral sequences grouped at the order level in each of the genomes
38
39 is shown in Fig. 1. The vHGTs were derived from both DNA and RNA viruses, however, the
40
41 majority were annotated as Unclassified RNA viruses (NCBI Taxonomy ID: 1922348; n = 44
42
43 sequences). The genome with the most Unclassified RNA viruses vHGTs was *Symbiodinium*
44
45 *microadriaticum* (n = 14), followed by the isolates *S. microadriaticum* KB8 (n = 11), and *S.*
46
47 *microadriaticum* 04-503SCI.03 (n = 6). The second most abundant group of vHGTs (n = 11) was
48
49 similar to giant viruses from the order Pimascovirales (Phylum Nucleotycoviricota; formerly
50
51
52
53
54
55
56
57
58
59
60

1
2
3 NCLDV). Finally, the third group of vHGTs ($n = 4$) was annotated to negative-strand RNA
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

viruses from the Mononegavirales (full annotation in Supplementary Table S2).

3.3 Functional annotation and expression of vHGTs sequences in Symbiodiniaceae

We found that the most abundant predicted functional annotation associated with the vHGTs is the DNA/RNA polymerase superfamily ($n = 12$), followed by RNA-directed RNA polymerase ($n = 11$), and viral RNA helicase ($n = 3$). Eleven vHGTs had RNA-directed RNA polymerase annotations (GO:0003968), followed by ATP binding (GO:0005524), viral RNA genome replication (GO:0039694), and mRNA capping (GO:0006370) as the most abundant GO terms. There were no protein family membership or gene ontology terms (GO terms) predicted for ten of the vHGTs (Fig. 1) (full annotation in Supplementary Table S2).

We also noted six vHGTs occurring in *Smic*, *Smic_04-503SCI.03*, *Smic_CassKB8* and *Stri_CCMP2592* that have two predicted domains, one of viral origin and the other of eukaryotic origin, suggesting the possibility that they might be chimeric genes (Méheust et al. 2018) (Supplementary Table S9). At the 5'-end of these proteins is a predicted Clavaminase synthase-like or a Winged helix DNA-binding domain and at the 3'-end is a DNA/RNA polymerase domain. Aligned RNA-seq reads were visualized and compared manually against all the predicted vHGT gene models to assess if these sequences are being expressed with at least one read mapping to the gene model. This analysis showed that only 13/59 vHGTs had evidence of read mapping to the gene model. This analysis showed that only 13/59 vHGTs had evidence of read mapping (Supplementary Table S2). This could be explained by the experimental conditions used to generate the RNA-seq data, or alternatively, the vHGTs are non-functional. Because there was also no expression of the putative chimeric domain genes we could not validate if these sequences are true chimeras or artifacts of gene prediction.

3.4 Phylogenetic profile of vHGTs sequences in Symbiodiniaceae

3.4.1 Unclassified RNA viruses

The majority of vHGTs were grouped into nine phylogenies (n = 44 sequences); these sequences are all putatively related to unclassified RNA viruses (Shi et al. 2016) (Fig. 2). The taxonomy of top viral hits contained the Wenzhou weivirus-like virus (n = 15), Beihai sobemo-like virus (n = 7), Beihai weivirus-like virus (n = 7), Beihai narna-like virus (n = 4) and Hubei Beny-like virus (n = 3). Re-annotation of these sequences with BLASTp found evidence of putative homology with hallmark RNA viral genes such as the RNA-dependent RNA polymerase (RdRp; n = 23 sequences), putative major capsid protein (n = 4), viral RNA helicase (n = 4), and polyproteins encoding replicases, including a RNA-dependent RNA polymerase region (n = 5). Moreover, these vHGT sequences had hits to the previously characterized *Symbiodinium* +ssRNA virus (accession KX538960 - KX787934) (Levin et al. 2017).

3.4.1.1 RdRp-like phylogenies:

The Symbiodiniaceae vHGT-derived RdRps-like sequences were grouped into four phylogenies that matched two different viral groups, indicating that multiple independent vHGTs had occurred. The first viral group (RdRp-like group 1), comprising the sequences from phylogeny A (n = 15 sequences), B (n = 8) and C (n = 3) (Fig. 2 and Supplementary Fig. S4), have hits to Weivirus-like viruses and the *Symbiodinium* +ssRNA virus. In the phylogenetic tree associated with each group, all Symbiodiniaceae vHGTs are clustered into a separate clade or within Weivirus-like viral sequences, distant from the *Symbiodinium* +ssRNA virus. Weivirus-

1
2
3 like viruses were first identified in the transcriptomes of molluscs and are related to alveolate
4
5 EVEs (Shi et al. 2016), and were later found in the Porifera *Halichondria panacea* (breadcrumb
6
7 sponge) (Waldron et al. 2018).
8
9

10 The second viral group (RdRp-like group 2) comprised the phylogeny D (n = 1), and
11
12 was a single Symbiodiniaceae sequence from the Bmin genome. This
13
14 sequence has hits with RdRps from narna-like viruses and is located
15
16 in a clade containing crustacean and Porifera associated viruses
17
18 (Beihai and Wenling narna-like viruses and Barns Ness breadcrumb
19
20 sponge narna-like virus) (Shi et al. 2016). Narnavirus genomes often contain a single
21
22 open reading frame coding for RdRp. This virus replicates in the cytoplasm of the host and is
23
24 suggested to be transmitted vertically through cell division, or horizontally during the mating
25
26 cycle of its fungal host (Dinan et al. 2020). Narnaviruses are associated with trypanosomatids (in
27
28 which a genomic narnavirus-like element was described (Lye et al.), diatoms (Charon et al.
29
30 2021), and red and brown algae (Dinan et al. 2020). Hits were also observed with plant viruses
31
32 (Ourmiavirus), which have putative origins *via* a chimeric fusion between a fungal narnavirus
33
34 and other plant viruses (Rastgou et al. 2009).
35
36
37
38
39
40
41
42
43

44 **3.4.1.2 MCP-like phylogenies:**

45
46 The MCP-like vHGT sequences formed three phylogenies: E (n = 1), F (n = 3) and G (n
47
48 = 2) (Fig. 2 and Supplementary Fig. S4). These sequences matched MCPs from the
49
50 dinoflagellate HcRNAV Dinornavirus, the invertebrate associated
51
52 Weivirus-like virus, Sobemo-like virus, narna-like virus, and MCPs
53
54
55
56
57
58
59
60

1
2
3 from the *Symbiodinium* +ssRNA virus. In the phylogenies E and G, vHGTs are clustered with
4
5 *Symbiodinium* +ssRNA virus; the sequences from phylogeny F are clustered at the base of the
6
7 tree. These patterns also suggest multiple independent vHGTs, with high divergence of
8
9 Symbiodiniaceae MCPs-like sequences, given that viral sequence hits and the tree topology is
10
11 approximately the same in all three phylogenies.
12
13

14
15 The remaining vHGTs with unclassified viral hits, from phylogenies H (OG7-
16
17 OG0009590; n = 6) and I (OG3-OG0006249; n = 5) , were re-annotated as viral RNA helicases
18
19 and "polyprotein coding for replicases including RNA-dependent RNA polymerase region"
20
21 (respectively). In the phylogeny of the polyprotein-like vHGTs (Phylogeny H), one vHGT
22
23 sequence (Stri.gene18597) is in a clade with fungal and plant Barnaviruses (Nibert et al. 2018)
24
25 and Solemoviruses from plant hosts (Sõmera et al. 2021), whereas the other vHGT sequences are
26
27 at the base of the tree. Interestingly, one sequence in this tree is from the plant *Poinsettia*, which
28
29 is not only a latent virus, but is suggested to be a hybrid between polerovirus and sobemovirus
30
31 (aus dem Siepen et al. 2005), which seems to be common in these types of RNA viruses. In the
32
33 phylogeny I, all Symbiodiniaceae vHGT-derived RNA helicase-like sequences) clustered with a
34
35 plant virus from the genus *Higrevirus* (Hibiscus green spot virus 2), and with a fungal virus,
36
37 *Agaricus bisporus* virus 8, forming a bigger clade which includes the aggressive plant disease
38
39 causing viruses from the genus *Benyvirus* (Gilmer et al. 2017), and Beny-like viruses.
40
41
42
43
44
45
46

47 **3.4.2 Pimascovirales**

48
49 The second group of viruses with similarities to vHGTs (n = 11) were all contained
50
51 within one phylogeny: J (n = 11 sequences) (Fig. 2 and Supplementary Fig. S4). These vHGTs
52
53 are similar to sequences from the Brazilian marseillevirus (n = 2), Noumeavirus (n = 2), Cannes
54
55
56
57
58
59
60

1
2
3 8 virus (n = 3), and Kurlavirus BKC-1 (n = 2); all of these taxa are from the Marseilleviridae
4 family of large DNA viruses that infect *Acanthamoeba* protists (Aherfi et al. 2014). Re-
5
6 annotation of these sequences with BLASTp showed that these sequences are restriction
7
8 endonucleases, which are suggested to be hotspots for mutations and footprints of mobile
9
10 elements in other Marseilleviridae (Doutre et al. 2014; Mueller et al. 2017); they also play a role
11
12 in degradation of host DNA during the early stages of Chlorovirus infection (Agarkova et al.
13
14 2006) (Phycodnaviridae; another large DNA virus infecting green algae). Although, in this
15
16 phylogeny, giant and large viral sequences are overrepresented, such as Marseilleviridae,
17
18 Mimiviridae, and Phycodnaviridae, we found that vHGTs are clustered at the base of the tree, or
19
20 with another eukaryote, *Pelagomonas calceolata* (Pelagophyceae). We also find that these
21
22 sequences occur in the genomes of the dinoflagellate *Polarella glacialis*. This suggests a more
23
24 complex scenario involving multiple independent and possibly ancient acquisitions, given that
25
26 almost all genomes from the *Symbiodinium* genus, which is the most ancient in relation to the
27
28 whole family, contained these sequences.
29
30
31
32
33
34
35
36
37

38 3.4.3 Mononegavirales

39
40 The final group of vHGTs have similarity to Mononegavirales, an order of ssRNA
41
42 negative-strand viruses. These vHGTs were contained within two phylogenies, K and L (Fig. 2
43
44 and Supplementary Fig. S4); four of the vHGTs have top hits to the RNA-directed RNA
45
46 polymerase catalytic domain. Whereas the best hits of the vHGTs were affiliated with metazoan
47
48 viruses (Tacheng Tick Virus 6 [n = 3] and Wuchan romanomermis nematode virus 2 [n = 1]), the
49
50 majority of sequences in the tree were from Rhabdoviridae or Rhabdo-like viruses, which are
51
52 complex and diverse viruses that include several plant pathogens such as Strawberry crinkle
53
54
55
56
57
58
59
60

1
2
3 virus and Citrus leprosis virus among others (Dietzgen et al. 2017). Whereas in phylogeny K
4 (OG6-OG0008326) the only vHGT present is outside the main group containing metazoan
5
6 (OG6-OG0008326) the only vHGT present is outside the main group containing metazoan
7
8 viruses and the plant soybean leaf-associated negative virus, the phylogeny L (OG1-OG0001147;
9
10 n = 3), the vHGTs are positioned at the bottom of the tree, clustered with another eukaryote (the
11
12 perkinsid *Perkinsus olseni*). Because we found hits with the same viral sequences that are split
13
14 into two phylogenies, it is possible that these were two independent HGT events: one involving
15
16 the *S. microadriaticum*, Smic, Smic_04-503SCI.03 and Smic_CassKB8 genomes, and the other
17
18 involving the Snat_CCMP2548 genome.
19
20
21
22
23

24 **3.5 Sequence composition of Symbiodiniaceae vHGTs sequences**

25
26 We calculated the coding sequence (CDS) GC-content for vHGTs and for all
27
28 Symbiodiniaceae CDSs predicted in the genomes that contain vHGTs. That later served as a
29
30 background for comparative analysis (Fig. 1). Whereas the mean GC-content of the
31
32 (background) Symbiodiniaceae CDSs ranged from 59.06 % (in *S. microadriaticum* CassKB8) to
33
34 51.38 % (in *B. minutum*), the mean GC-content of vHGTs ranged from 60.13 % (in *S. natans*
35
36 CCMP2548) to 51.29 % (in *B. minutum*). Although the GC-content of the vHGTs was
37
38 marginally lower than the background CDSs, this difference was statistically significant ($P =$
39
40 3.6E-07). However, when these two sets are compared genome-by-genome (i.e., comparing the
41
42 vHGTs and background CDSs from the same genome) we only found significant GC-content
43
44 differences in the Smic_CassKB8 ($P = 1.4E-06$), Smic ($P = 2.4E-05$), Snec_CCMP2469 ($P =$
45
46 4.4E-05), and Smic_04-503SCI.03 ($P = 3.9E-04$) genomes. This comparison was not performed
47
48 for the Cgor, and Slin_CCMP2456 genomes because the number of vHGTs was too small for
49
50 statistical analysis. There were no significant differences in the length of CDSs between the
51
52
53
54
55
56
57
58
59
60

1
2
3 vHGTs and the background CDSs when comparing all sets together ($P = 8.0E-02$) or
4
5 individually, although vHGTs are slightly longer (Supplementary Table S3 and S4).
6
7
8
9

10 **3.6 Presence of introns and relic DinoSL motifs in Symbiodiniaceae vHGTs**

11
12 The vHGTs are predominantly encoded by genes with multiple exons, with only 4/59
13
14 present as single-exon genes. The average number of introns per gene in the vHGT gene sets is
15
16 either roughly equal to (e.g., in *S. microadriaticum* CassKB8), or lower than (e.g., in *S.*
17
18 *microadriaticum*) the background gene set for each genome (Fig. 1) (Supplementary Table S6).
19
20 The average intron length of the vHGT genes tends to differ from that of background genes in
21
22 each genome, however, this could be a result of the small number of viral HGT genes being
23
24 analyzed, however, this could be a result of the small number of viral HGT genes being
25
26 analyzed.
27

28
29 Of the 59 vHGT-derived genes, 13 have a DinoSL sequence that is upstream (within 2
30
31 kbp) of the first exon (Supplementary Table S7); there was no evidence of relic DinoSL
32
33 sequences downstream of the identified DinoSL, in close proximity to the vHGT genes.
34
35
36
37

38 **3.7 Distribution of vHGT sequences in Symbiodiniaceae genome scaffolds**

39
40 We evaluated if there was enrichment of vHGT genes in any of the genome scaffolds in
41
42 the 10 studied taxa. Our reasoning was that multiple vHGTs in a single scaffold from the same
43
44 taxonomic source may indicate that they were derived from an exogenous virus that was
45
46 included during sequencing (i.e., a contaminant in the assembly), or it could represent a complete
47
48 or near-complete viral genome that had integrated into host DNA. The majority of vHGTs (53
49
50 out of 59) were not on the same scaffold as another vHGT. However, in two of the *S.*
51
52 *microadriaticum* genomes we found cases of two adjacent vHGTs; these genes were in
53
54
55
56
57
58
59
60

1
2
3 scaffold3501 from Smic_04-503SCI.03 (gene15955-15956) and in scaffold792 and scaffold67
4 from Smic (gene19249 -19250 and gene3240-3241, respectively) (Supplementary Table S2). In
5 all three cases, one of the proteins in each pair has similarity to viral RNA-dependent RNA
6 polymerase proteins and the other to viral major capsid proteins, which are the two proteins that
7 comprise the +ssRNA *Symbiodinium* virus genome.
8
9
10
11
12
13

14 15 16 17 **3.8 *Symbiodinium* +ssRNA virus similarities in Symbiodiniaceae genomes**

18
19 We found that the majority of vHGTs (n = 34) had similarity with the most well
20 characterized virus known to infect this group, the *Symbiodinium* +ssRNA virus (TR74740
21 c13_g1_i1 and _i2, accession KX538960 - KX787934) (Levin et al. 2017; Montalvo-Proaño et
22 al. 2017) that have been implicated in coral holobiont health (Grupstra et al. 2022), and also with
23 the *Heterocapsa circularisquama* RNA virus (HcRNAV; another dinoflagellate infecting-virus).
24 In addition, in two of the *S. microadriaticum* genomes we found adjacent vHGTs re-annotated
25 using BLASTp as the two open reading frames (ORFs) that compose the complete genomes of
26 these two viruses: RNA-dependent RNA polymerase (RdRp) (n = 21), and putative major capsid
27 protein (MCP; n = 4) (Supplementary Table S2). Given these findings, we compared the
28 composition and genomic context of the MCPs-like and RdRp-like vHGTs, with respective
29 sequences from the +ssRNA *Symbiodinium* and HcRNAV viruses. We found that the mean GC-
30 content of Symbiodiniaceae vHGTs MCPs-like and RdRps-like genes was higher (52.95% and
31 53.76%, respectively) than in the putative +ssRNA *Symbiodinium* virus homologs (48.10% and
32 43.95%, respectively), lower than the HcRNAV homologs (58.06% and 54.03%, respectively),
33 and lower than in the full Symbiodiniaceae genomes (56.79%). In addition, the mean CDS length
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 of MCPs-like genes was slightly lower than their viral counterparts, whereas the RdRps-like
4
5 genes were slightly longer (Supplementary Table S4 and S5).
6
7

8 Whereas these proteins do not appear to be expressed (i.e., they did not have any RNA-
9
10 seq reads aligned using HISAT2; Supplementary Fig. 1A, 2A, and 3A), they were on scaffolds
11
12 with other expressed multi-exon genes. Realignment of selected RdRp-like and MCP-like
13
14 proteins against the regions of the host genome that encoded these genes (see Methods)
15
16 demonstrated that all six of these proteins are significantly degraded (Supplementary Fig. S1-3).
17
18 These genes often had similarity to only part of the query viral protein (e.g., Supplementary Fig.
19
20 S1B), in-frame stop codons, or frame-shift mutations (e.g., Supplementary Fig. S2C) that would
21
22 suggest that they are no longer under selection and have become pseudogenes. Furthermore, two
23
24 of the putative integrated viral genomes, one from *S. microadriaticum* (gene3240-3241) and one
25
26 from *S. microadriaticum* 04-503SCI.03 (gene15955-15956) appear to be homologous
27
28 (Supplementary Fig. S1 and S3), potentially having arisen from an infection in the common
29
30 ancestor of the two isolates. The third putative integrated genome (gene19249 -19250) only
31
32 appears in the *S. microadriaticum* genome.
33
34
35
36
37
38
39

40 **3.8 Genomic context of Symbiodiniaceae vHGTs**

41
42 Through analysis of the position of these sequences in the genome scaffolds, it was
43
44 observed that the majority of vHGTs encoding MCP-like and RdRp-like genes were located at
45
46 the start of their scaffolds, or in single gene scaffolds. The CDSs of genes of non-vHGTs up- and
47
48 downstream of vHGTs (if available) were annotated using the BLASTp online suite; we found
49
50 that all vHGTs on scaffolds with other genes were located within non-viral host sequences
51
52 (Supplementary Table S8). In these scaffolds, we found the presence of protein sequences such
53
54
55
56
57
58
59
60

1
2
3 as the CCHC-type domain-containing protein (n = 2 sequences in total), that was shown in other
4
5 eukaryotes as either an anti- or pro-viral, targeting several viruses such as RNA viruses
6
7 (Hajikhezri et al. 2020), that could also be “hijacked” by nuclear-replicating viruses to promote
8
9 viral production and are also induced after heat shock stress conditions (Younis et al. 2018).
10
11 Others genes located close to vHGTs are putative E3 ubiquitin-protein ligases (n = 5), RING
12
13 finger proteins (n = 3; which have roles in viral evasion of host innate immunity) (Xu et al.
14
15 2017), F-box proteins (n = 4; which is a component of E3 ubiquitin–ligase complex that targets
16
17 proteins for ubiquitination and degradation) (R.L. Correa et al. 2013), and the ankyrin repeat
18
19 domain-containing proteins (n = 6; which are known to regulate virus-host interactions) (Than et
20
21 al. 2016).
22
23
24
25

26
27 Furthermore, some vHGTs were found to be located in close proximity to transposons,
28
29 retrotransposons, and repetitive sequences, such as the Retrovirus-related Pol polyproteins (n =
30
31 8), Transposon Ty2 Gag-Pol polyprotein (n = 2), pentatricopeptide repeat (n = 2), LINE-1
32
33 retrotransposable element ORF2 protein (n = 1), Copia protein (n = 5), and Reverse transcriptase
34
35 domain-containing protein (n = 2). Finally, three scaffolds were identified in *S. microadriaticum*
36
37 genomes (Fig. 3) in which the MCPs-like and RdRps-like genes were located adjacent to each
38
39 other on the same scaffold (scaffold67 and scaffold792 for Smic and scaffold3501 for Smic 04-
40
41 503SCI.03), where in Smic sequences were in an orientation that was inverted relative to the
42
43 +ssRNA *Symbiodinium* RNA virus genome (i.e., RdRp upstream of MCP). On scaffold67 these
44
45 sequences are flanked by a downstream unannotated protein and by a Retrovirus-related Pol
46
47 polyprotein (from a type-1 retrotransposable element R2). Finally, on scaffold792, there is a
48
49 reverse transcriptase domain-containing protein at the end of this scaffold.
50
51
52
53
54
55
56
57
58
59
60

4 Discussion

In this study, we describe the function, genomic features, and taxonomic distribution of 59 high-confidence viral-derived HGT genes present in 10/13 Symbiodiniaceae genomes. These genes were identified using a combination of HGT scoring metrics, phylogeny-based detection, and significant manual curation. Due to the stringent nature of the filtering that was applied and the lack of a comprehensive database consisting of dinoflagellate-infecting virus genomes, the number of vHGTs we found likely represents a lower bound of the actual number of viral HGT events in these genomes. Nevertheless, by describing a more conservative set of viral HGT candidates in Symbiodiniaceae, we open-up the possibility for additional experimental evaluation to be undertaken to advance functional analysis of these genes. There were between 1-14 vHGTs identified in each genome, comprising ~0.2% of the total Symbiodiniaceae gene repertoire. Whereas HGT score metrics, such as AI and HGT indexes, are able to identify sequences that are highly similar to foreign sources, phylogenetics can identify more distantly related sequences. Thus, by combining the two approaches we have a higher likelihood of identifying genuine viral HGT events (Fan et al. 2020).

Previous reports of viral genes in Symbiodiniaceae genomes showed that in *S. microadriaticum*, genes with similarity to known viral sequences accounted for < 3% of the total gene inventory (Aranda et al. 2016), and in *C. goreau* this number was only ~0.1% (Liu et al. 2018). In the latter case, regions were identified that match viruses known to infect *C. goreau* (Weynberg et al. 2017) although the origin in these genomes was not clear. Shoguchi et al. (2018) suggested that DNA methylation in some Symbiodiniaceae genomes (clade SymA and SymC) was related to the presence and expansion of endogenous retroviruses such as retrovirus-

1
2
3 related genes. In addition, in the *S. microadriaticum* genome, recently expanded Ty1-*cop*ia-type
4 LTR retrotransposons are activated under acute temperature increase (Chen et al. 2018).
5
6

7
8 Our findings expand upon these studies by demonstrating the presence of candidate viral
9 sequences from giant virus sources (Pimascovirales) that are potentially functioning as
10 endonucleases, as well as polymerase sequences originating from RNA viruses that could
11 function in viral replication. Exploration of vHGT sequence composition and structural features
12 showed that their GC-contents are slightly lower than for Symbiodiniaceae background CDSs,
13 but higher than homologous viral sequences. The sizes of the coding sequences do not differ
14 from the background Symbiodiniaceae CDS repertoire. This could indicate that some of these
15 transfers were recent in the evolutionary history of this group and are still in the process of being
16 ameliorated into the recipient host genome (Lawrence and Ochman 1997). Moreover, we found
17 that viral sequences are primarily localized to different scaffolds in the genome but tend to be
18 near eukaryotic genes with known roles in virus-host interaction. Finally, we found these viral
19 sequences near mobile genetic elements, such as retrotransposons and transposons, that may act
20 not only as preferential integration sites, but as a mechanism for their integration and expansion.
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37

38 Currently, the best characterized virus that infects Symbiodiniaceae is a positive single-
39 stranded RNA virus (+ssRNA) (NCBI:txid1909300) (Nagasaki et al. 2005; A.M.S. Correa et al.
40 2013; Montalvo-Proañó et al. 2017; Weynberg et al. 2017), first described in stressed
41 *Montastraea cavernosa* coral viromes (A.M.S. Correa et al. 2013), that is related to the
42 dinoflagellate virus *Heterocapsa circularisquama* (HcRNAV) from the *Alvernaviridae* family.
43 Its complete genome was described by Levin et al (2016); they found that thermally stressed
44 *Cladocopium* (formerly *Symbiodinium* clade C1) populations experienced an extreme +ssRNA
45 virus infection, suggesting that viral infections may influence *Symbiodinium* thermal sensitivity
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 and consequently, coral bleaching. Nevertheless, these authors were unable to amplify MCP
4
5 sequences from genomic DNA (only from cDNA), suggesting that these sequences were not
6
7 endogenous in the host genome, but rather exogenous. Moreover, they reported the presence of a
8
9 DinoSL at the 5'-end of the amplified viral transcripts, proposing that this motif was
10
11 incorporated into the viral genome as a form of “molecular mimicry” to evade the host immune
12
13 response and to drive efficient translation of the viral genome by the host cellular machinery.
14
15

16
17 It has been suggested that viral infections in Symbiodiniaceae are latent, given that VLPs
18
19 and expressed sequences were found in asymptomatic and healthy cultures subjected to stress
20
21 conditions (Wilson et al. 2001; Lawrence et al. 2014; Brüwer et al. 2017). Lawrence et al. (2017)
22
23 showed that in Symbiodiniaceae cultures exposed to UV stress, there was upregulation of virus-
24
25 like genes associated with viral genome replication and protein production, such as ubiquitin-
26
27 encoding genes and homing endonucleases that could be involved in the excision of the
28
29 integrated viral genome. This pattern is consistent with latent virus replication and agrees with
30
31 our results.
32
33

34
35 Latent viruses are well described in plants and often do not have any apparent symptoms
36
37 in wild populations, however, they can cause disease in crops or due to experimentally created
38
39 stresses such as mechanical wounding (Roossinck and García-Arenal 2015). Mixed viral
40
41 infections and changes in environmental conditions trigger the switch from latent to active
42
43 infections and diseases (Takahashi et al. 2019). In plants, the integration of a virus genome is not
44
45 required for latent viral replication, although partial genome integration and endogenization of
46
47 entire viral genomes may play a major role in latency (Takahashi et al. 2019). Latent infections
48
49 and EVEs present major risks for cultivated crops if they become active (Bertsch et al. 2009).
50
51 This is the case for the hybrid *Musa* (Iskra-Caruana et al. 2010), where the endogenous *banana*
52
53
54
55
56
57
58
59
60

1
2
3 *streak virus* (eBSV), which has multiple copies integrated into the host genome, becomes
4
5 activated under stress, recombining, producing episomes, viral particles, and disease symptoms
6
7 such as streaks and necrosis in plant leaves. Recent work using CRISPR/Cas9-mediated editing
8
9 of eBSV has been successful in preventing activation of these viruses (Tripathi et al. 2019). In
10
11 contrast, benefits of latent viruses and EVEs are heat and drought tolerance, protection of host
12
13 plants against virulent viral strains (Agüero et al. 2018), and suppression of heterologous viral
14
15 infections (Staginnus et al. 2007).
16
17
18
19
20
21

22 **4.1 Proposed models for vHGTs in Symbiodiniaceae genomes**

23
24 Given our findings, we propose a model with two hypothetical scenarios to explain the
25
26 occurrence of viral sequences in the genomes of Symbiodiniaceae and the underlying
27
28 mechanisms for their integration (Fig. 4). In this model, whereas several viral groups infect
29
30 Symbiodiniaceae, independently entering the host cell and subjugating the host machinery to
31
32 produce viral mRNAs, some virus sequences become endogenous viral elements (EVEs) by a
33
34 series of incidental integration events into Symbiodiniaceae genomes (A). For this to occur, viral
35
36 mRNAs need to be reverse transcribed before they can be integrated into the host genome. This
37
38 process might occur *via* the mRNA recycling mechanism that is prevalent in dinoflagellates
39
40 (Stephens et al. 2020), whereby spliced-leader-containing transcripts are reverse transcribed
41
42 (possibly mediated by retrotransposons (Lee et al. 2014; Song et al. 2017) into DNA before
43
44 being integrated back into the genome through non-homologous recombination (Slamovits and
45
46 Keeling 2008). The infecting viral elements must escape degradation by host defenses. Because
47
48 viral mRNAs are highly expressed during an active infection, the host spliceosomal machinery
49
50 may erroneously add a DinoSLs motif onto viral transcripts, making the transcripts appear
51
52
53
54
55
56
57
58
59
60

1
2
3 “native”. This may have allowed the viral genome to better evade host defenses and providing
4 the virus with a dinoflagellate TATA-box-like promoter sequence; dinoflagellates appear to use a
5 TTTT motif (which is also present in the DinoSL) in place of the canonical TATA motif in other
6 eukaryotes (Guillebault et al. 2002). This is supported by the presence of a DinoSL in the
7
8 *Symbiodinium* RNA virus genome (Levin et al. 2017) and suggests that this motif could have
9
10 been acquired by the virus from the host either in the common ancestor of these viruses or during
11 each new infection cycle. Whereas convergence is also possible (i.e., independent evolution of a
12 DinoSL-like motif), it is less likely given the high level of sequence similarity between the viral
13 and host DinoSL sequences. Furthermore, the occurrence of DinoSL sequences upstream of 13
14 vHGTs, plus the localization of these sequences near retrovirus related polyproteins, copia
15 proteins, and proteins containing reverse transcriptase domains in the Symbiodiniaceae scaffolds,
16 indirectly support this mechanism of integration.
17
18
19
20
21
22
23
24
25
26
27
28
29

30
31 However, given that we found complete viral genomes with similarity to the +ssRNA
32 *Symbiodinium* virus, we propose a second scenario (B) whereby vHGTs from +ssRNA
33 *Symbiodinium* viruses resulted from a previously hidden (cryptic) life stage of these viruses. In
34 this scenario, copies of viral genomes become endogenized and deactivated by the host genome,
35 decaying as pseudogenes. Alternatively, they remain intact, transcriptionally silent (as
36 proviruses), and may eventually become activated in response to stress, shifting from a proviral
37 silent life-stage to being infectious. These RNA viruses could actively exploit the dinoflagellate
38 mRNA recycling process to facilitate their integration into host DNA although this would require
39 the virus to also have a system for permanence (survival) and activation, both of which are
40 unknown at this time. However, given that the genes associated with the three putative integrated
41 *Symbiodinium* +ssRNA virus genomes, identified in two of the *S. microadriaticum* isolates,
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 appear to have become pseudogenized, it seems likely that many of the viral genes that we
4 identified in the Symbiodiniaceae genomes are (or are becoming) pseudogenes. This is an
5 expected result given that there is strong selective pressure on the host to inactivate endogenous
6 viruses and that we only expect to see remnants of inactivated endogenous viruses originating
7 from past (failed) infections in the genomes of extant lineages (i.e., active endogenous viral
8 elements would likely result in an active viral infection that would cause the lineage to perish).
9
10 In addition, the two vHGT genes that are part of the putative *Symbiodinium* +ssRNA viruses
11 identified in the *S. microadriaticum* (Smic) genome are encoded in the reverse order to what is
12 observed in the available virus reference genome. This may result from the host genome
13 inactivating the invasive viruses by shuffling their gene order, preventing the formation of
14 infective viral particles. Furthermore, the putative inactivation of these elements by the host
15 suggests that they were active when they were integrated into the genome, lending support to the
16 second scenario (which produces viral elements that are “dead on arrival”). Moreover, the
17 multiple, complete and incomplete copies of these viral genomes in Symbiodiniaceae, from
18 different isolates and genera, suggests that viral endogenization is an ongoing process in these
19 organisms, with a potentially rapid turnover of invasion and pseudogenization or activation and
20 infection, similar to Mavirus virophage infection-integration cycles (Hackl et al. 2021).
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41

42 Under this scenario, endogenization does not need to be an obligatory step in the life
43 cycle of the virus (i.e., it is a transient stage) and does not need to become fixed in all hosts
44 occurring, for instance, only when there is a low number of susceptible hosts. Becoming
45 endogenized when the density of host cells is low is advantageous for the virus, given that a copy
46 of the virus genome is always transmitted to host progeny, effectively preserving the virus until
47 conditions promote the lytic life stage. The virus life cycle may therefore mirror that of the host,
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 remaining inactive during phases where the density of susceptible cells is low, such as during the
4 host coccoid endosymbiont life stage, and becoming active when the density of susceptible cells
5 is high, such as during the host mastigote, free-living stage. From this perspective, the more
6 susceptible stage would be the free-living state with no endogenous provirus in its genome. From
7 the perspective of the Symbiodiniaceae host, it could be argued that having an endogenous
8 provirus would confer superinfection immunity against other exogenous viruses *via* cross
9 protection and may benefit the cell under certain conditions (Mette et al. 2002). A hidden
10 transient viral life cycle, if confirmed for these viruses, would unify the current knowledge of
11 viruses in Symbiodiniaceae, explaining the observations of: (1) latent viral infections in healthy
12 algal cultures after stress induction, (2) exogenous RNA viruses that infect Symbiodiniaceae, and
13 (3) our findings of endogenous versions of these exogenous *Symbiodinium* RNA virus genomes.
14
15
16
17
18
19
20
21
22
23
24
25
26
27

28 As previously noted, the triggers that cause the switch from latent to active viral
29 infections and diseases in plants and algae (including cultured Symbiodiniaceae) overlaps with
30 those that trigger and influence coral bleaching and other coral diseases (Harvell et al. 1999;
31 Lesser and Farrell 2004; Randall and van Woesik 2015; Lawrence et al. 2017). It has been
32 recently proposed that breakdown of the Symbiodiniaceae-coral symbiosis due to heat stress
33 (coral bleaching) is caused by increased proliferation of the symbiont (Rädecker et al. 2021). The
34 observed increase in viral particle production under thermal stress may coincide with the
35 increased proliferation of the symbionts during the early stages of coral bleaching. This result
36 suggests that the virus life cycle could be directly linked to that of the host.
37
38
39
40
41
42
43
44
45
46
47
48

49 Consequently, the main implication of our work is that if EVEs and latent proviruses that
50 are hidden in Symbiodiniaceae genomes could become activated by environmental stressors, that
51 they may cause local outbreaks and possibly contribute to the onset of some coral diseases. To
52
53
54
55
56
57
58
59
60

1
2
3 conclude, by describing a variety of acquired viral genes originating from multiple HGT events,
4
5 we provide strong evidence that viruses have invaded and introduced novel genetic material into
6
7 Symbiodiniaceae genomes. Our work suggests the possibility that endogenous and latent viruses
8
9 could modulate and participate in the life cycle, evolution, and potentially breakdown of the
10
11 Symbiodiniaceae-host symbiosis.
12
13

14 15 16 17 **Data availability**

18
19 All Symbiodiniaceae sequence data used in this study are available on
20
21 <https://espace.library.uq.edu.au/view/UQ:f1b3a11> and
22
23 <https://espace.library.uq.edu.au/view/UQ:8279c9a>. Code to generate the HGT calculations is
24
25 available at https://github.com/LFelipe-B/Symbiodiniaceae_vHGT_scripts.
26
27
28
29

30 31 **Supplementary data**

32
33 Supplementary data is available at Virus Evolution online.
34
35
36
37

38 39 **Acknowledgements**

40 LFB wish to thank Paulo Sérgio Salomon (Universidade Federal do Rio de Janeiro - UFRJ) and
41
42 coordinators from “Programa de Pós-graduação em Biodiversidade e Biologia Evolutiva”
43
44 (PPGBBE - UFRJ), for the initial conceptual support in early stage of these ideas; Lílian Caesar
45
46 for critical reading, discussions, love and support during all stages of this work; François
47
48 Bucchini for support in initial codes; GenoToul bioinformatics platform for access to computing
49
50 cluster facilities during early stages of this work; Girish Beedessee and Cheong Xin Chan for
51
52 guidance in data retrieval.
53
54
55
56
57
58
59
60

Funding

This work was supported by an appointment to the NASA Postdoctoral Program at Rutgers University, New Brunswick, administered by Oak Ridge Associated Universities under contract with NASA, to LFB. DB and TGS were supported by a NASA grant (80NSSC19K0462) to DB and a NIFA-USDA Hatch grant (NJ01180) to DB.

Conflict of interest

The authors declare no conflict of interest.

References

- Agarkova IV, Dunigan DD, Etten JLV. 2006. Virion-Associated Restriction Endonucleases of Chloroviruses. *J. Virol.* 80:8114.
- Agüero J, Gómez-Aix C, Sempere RN, García-Villalba J, García-Núñez J, Hernando Y, Aranda MA. 2018. Stable and Broad Spectrum Cross-Protection Against Pepino Mosaic Virus Attained by Mixed Infection. *Front. Plant Sci.* [Internet] 9. Available from: <https://www.frontiersin.org/article/10.3389/fpls.2018.01810>
- Aherfi S, La Scola B, Pagnier I, Raoult D, Colson P. 2014. The expanding family Marseilleviridae. *Virology* 466–467:27–37.
- Aranda M, Li Y, Liew YJ, Baumgarten S, Simakov O, Wilson MC, Piel J, Ashoor H, Bougouffa S, Bajic VB, et al. 2016. Genomes of coral dinoflagellate symbionts highlight evolutionary adaptations conducive to a symbiotic lifestyle. *Sci. Rep.* 6:39734.
- Aranda M, Li Y, Liew YJ, Baumgarten S, Simakov O, Wilson MC, Piel J, Ashoor H,

1
2
3 Bougouffa S, Bajic VB, et al. 2016. Genomes of coral dinoflagellate symbionts highlight
4 evolutionary adaptations conducive to a symbiotic lifestyle. *Sci. Rep.* 6:39734.

5
6
7 aus dem Siepen M, Pohl JO, Koo B-J, Wege C, Jeske H. 2005. Poinsettia latent virus is
8 not a cryptic virus, but a natural polerovirus–sobemovirus hybrid. *Virology* 336:240–250.

9
10
11 Benites LF, Silva-Lima AW, da Silva-Neto ID, Salomon PS. 2018. Megaviridae-like
12 particles associated with Symbiodinium spp. from the endemic coral *Mussismilia*
13 *braziliensis*. *Symbiosis* 76:303–311.

14
15
16
17
18 Bertsch C, Beuve M, Dolja VV, Wirth M, Pelsy F, Herrbach E, Lemaire O. 2009.
19 Retention of the virus-derived sequences in the nuclear genome of grapevine as a
20 potential pathway to virus resistance. *Biol. Direct* 4:21.

21
22
23
24
25 Boschetti C, Carr A, Crisp A, Eyres I, Wang-Koh Y, Lubzens E, Barraclough TG,
26 Micklem G, Tunnacliffe A. 2012. Biochemical diversification through foreign gene
27 expression in bdelloid rotifers. *PLoS Genet.* 8:e1003035.

28
29
30
31
32 Bräutigam M, Klein M, Knippers R, Müller DG. 1995. Inheritance and Meiotic
33 Elimination of a Virus Genome in the Host *Ectocarpus Szlzculosus* (phaeophyceae)1. *J.*
34 *Phycol.* 31:823–827.

35
36
37
38
39 Brüwer JD, Agrawal S, Liew YJ, Aranda M, Voolstra CR. 2017. Association of coral
40 algal symbionts with a diverse viral community responsive to heat shock. *BMC*
41 *Microbiol.* 17:174.

42
43
44
45
46 Buchfink B, Xie C, Huson DH. 2015. Fast and sensitive protein alignment using
47 DIAMOND. *Nat. Methods* 12:59–60.

48
49
50
51
52 Capella-Gutiérrez S, Silla-Martínez JM, Gabaldón T. 2009. trimAl: a tool for automated
53 alignment trimming in large-scale phylogenetic analyses. *Bioinformatics* 25:1972–1973.

1
2
3 Charon J, Murray S, Holmes EC. 2021. Revealing RNA virus diversity and evolution in
4 unicellular algae transcriptomes. *Virus Evol.* [Internet]. Available from:

5
6
7 <https://doi.org/10.1093/ve/veab070>
8
9

10 Chen JE, Cui G, Wang X, Liew YJ, Aranda M. 2018. Recent expansion of heat-activated
11 retrotransposons in the coral symbiont *Symbiodinium microadriaticum*. *ISME J.* 12:639–
12
13 643.
14
15

16
17 Chen Y, González-Pech RA, Stephens TG, Bhattacharya D, Chan CX. 2020. Evidence
18 That Inconsistent Gene Prediction Can Mislead Analysis of Dinoflagellate Genomes. *J.*
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
Phycol. 56:6.

Cock JM, Sterck L, Rouzé P, Scornet D, Allen AE, Amoutzias G, Anthouard V,
Artiguenave F, Aury J-M, Badger JH, et al. 2010. The *Ectocarpus* genome and the
independent evolution of multicellularity in brown algae. *Nature* 465:617–621.

Correa AMS, Ainsworth TD, Rosales SM, Thurber AR, Butler CR, Vega Thurber RL.
2016. Viral Outbreak in Corals Associated with an In Situ Bleaching Event: Atypical
Herpes-Like Viruses and a New Megavirus Infecting *Symbiodinium*. *Front. Microbiol.*
7:127.

Correa AMS, Welsh RM, Vega Thurber RL. 2013. Unique nucleocytoplasmic dsDNA
and +ssRNA viruses are associated with the dinoflagellate endosymbionts of corals.
ISME J. 7:13–27.

Correa AMS, Welsh RM, Vega Thurber RL. 2013. Unique nucleocytoplasmic dsDNA
and +ssRNA viruses are associated with the dinoflagellate endosymbionts of corals.
ISME J. 7:13–27.

Correa RL, Bruckner FP, Cascardo R de S, Alfenas-Zerbini P. 2013. The Role of F-Box

- 1
2
3 Proteins during Viral Infection. *Int. J. Mol. Sci.* 14:4030.
4
5 Davy SK, Burchett SG, Dale AL, Davies P, Davy JE, Muncke C, Hoegh-Guldberg O,
6
7 Wilson WH. 2006. Viruses: agents of coral disease? *Dis. Aquat. Organ.* 69:101–110.
8
9 Delaroque N, Boland W. 2008. The genome of the brown alga *Ectocarpus*
10
11 *siliculosus* contains a series of viral DNA pieces, suggesting an ancient association with
12
13 large dsDNA viruses. *BMC Evol. Biol.* 8:110.
14
15 Dietzgen RG, Kondo H, Goodin MM, Kurath G, Vasilakis N. 2017. The family
16
17 Rhabdoviridae: mono- and bipartite negative-sense RNA viruses with diverse genome
18
19 organization and common evolutionary origins. *Virus Res.* 227:158.
20
21 Dinan AM, Lukhovitskaya NI, Olendraite I, Firth AE. 2020. A case for a negative-strand
22
23 coding sequence in a group of positive-sense RNA viruses. *Virus Evol.* [Internet] 6.
24
25 Available from: <https://www.ncbi.nlm.nih.gov/labs/pmc/articles/PMC7010960/>
26
27
28 Dautre G, Philippe N, Abergel C, Claverie J-M. 2014. Genome Analysis of the First
29
30 Marseilleviridae Representative from Australia Indicates that Most of Its Genes
31
32 Contribute to Virus Fitness. *J. Virol.* 88:14340.
33
34
35 Emms DM, Kelly S. 2019. OrthoFinder: phylogenetic orthology inference for
36
37 comparative genomics. *Genome Biol.* 20:238.
38
39
40
41 Fan X, Qiu H, Han W, Wang Y, Xu D, Zhang X, Bhattacharya D, Ye N. 2020.
42
43 Phytoplankton pangenome reveals extensive prokaryotic horizontal gene transfer of
44
45 diverse functions. *Sci. Adv.* 6:eaba0111.
46
47
48 Gilmer D, Ratti C, ICTV Report Consortium YR 2017. ICTV Virus Taxonomy Profile:
49
50 Benyviridae. *J. Gen. Virol.* 98:1571–1572.
51
52
53 Gladyshev EA, Meselson M, Arkhipova IR. 2008. Massive horizontal gene transfer in
54
55
56
57
58
59
60

- 1
2
3 bdelloid rotifers. *Science* 320:1210–1213.
- 4
5 González-Pech RA, Stephens TG, Chen Y, Mohamed AR, Cheng Y, Shah S, Dougan
6
7 KE, Fortuin MDA, Lagorce R, Burt DW, et al. 2021. Comparison of 15 dinoflagellate
8
9 genomes reveals extensive sequence and structural divergence in family Symbiodiniaceae
10
11 and genus Symbiodinium. *BMC Biol.* 19:73.
- 12
13
14 Gornik SG, Ford KL, Mulhern TD, Bacic A, McFadden GI, Waller RF. 2012. Loss of
15
16 nucleosomal DNA condensation coincides with appearance of a novel nuclear protein in
17
18 dinoflagellates. *Curr. Biol. CB* 22:2303–2312.
- 19
20
21 Gornik SG, Hu I, Lassadi I, Waller RF. 2019. The Biochemistry and Evolution of the
22
23 Dinoflagellate Nucleus. *Microorganisms* 7.
- 24
25
26 Grupstra CGB, Howe-Kerr LI, Veglia AJ, Bryant RL, Coy SR, Blackwelder PL, Correa
27
28 AMS. 2022. Thermal stress triggers productive viral infection of a key coral reef
29
30 symbiont. *ISME J.*
- 31
32
33 Hackl T, Duponchel S, Barenhoff K, Weinmann A, Fischer MG. 2021. Virophages and
34
35 retrotransposons colonize the genomes of a heterotrophic flagellate. *eLife* 10:e72674.
- 36
37
38 Hajikhezri Z, Darweesh M, Akusjärvi G, Punga T. 2020. Role of CCCH-Type Zinc
39
40 Finger Proteins in Human Adenovirus Infections. *Viruses* [Internet] 12. Available from:
41
42 <https://www.ncbi.nlm.nih.gov/labs/pmc/articles/PMC7698620/>
- 43
44
45 Harvell CD, Kim K, Burkholder JM, Colwell RR, Epstein PR, Grimes DJ, Hofmann EE,
46
47 Lipp EK, Osterhaus ADME, Overstreet RM, et al. 1999. Emerging Marine Diseases--
48
49 Climate Links and Anthropogenic Factors. *Science* 285:1505–1510.
- 50
51
52 Hoffman LR, Stanker LH. 1976. Virus-like particles in the green alga *Cylindrocapsa*.
53
54 *Can. J. Bot.* 54:2827–2841.
- 55
56
57
58
59
60

- 1
2
3 Hyman P, Abedon ST. 2012. Smaller Fleas: Viruses of Microorganisms. *Scientifica*
4
5 2012:e734023.
6
7
8 Irwin NAT, Martin BJE, Young BP, Browne MJG, Flaus A, Loewen CJR, Keeling PJ,
9
10 Howe LJ. 2018. Viral proteins as a potential driver of histone depletion in dinoflagellates.
11
12 *Nat. Commun.* 9:1535.
13
14
15 Irwin NAT, Pittis AA, Richards TA, Keeling PJ. 2022. Systematic evaluation of
16
17 horizontal gene transfer between eukaryotes and viruses. *Nat. Microbiol.* 7:327–336.
18
19
20 Iskra-Caruana M-L, Baurens F-C, Gayral P, Chabannes M. 2010. A four-partner plant–
21
22 virus interaction: enemies can also come from within. *Mol. Plant-Microbe Interact.*
23
24 *MPMI* 23:1394–1402.
25
26
27 Janouškovec J, Gavelis GS, Burki F, Dinh D, Bachvaroff TR, Gornik SG, Bright KJ,
28
29 Imanian B, Strom SL, Delwiche CF, et al. 2017. Major transitions in dinoflagellate
30
31 evolution unveiled by phylotranscriptomics. *Proc. Natl. Acad. Sci. U. S. A.* 114:E171–
32
33 E180.
34
35
36 Kalyaanamoorthy S, Minh BQ, Wong TKF, von Haeseler A, Jermini LS. 2017.
37
38 ModelFinder: fast model selection for accurate phylogenetic estimates. *Nat. Methods*
39
40 14:587–589.
41
42
43 Katoh K, Standley DM. 2013. MAFFT Multiple Sequence Alignment Software Version
44
45 7: Improvements in Performance and Usability. *Mol. Biol. Evol.* 30:772–780.
46
47
48 Kim D, Paggi JM, Park C, Bennett C, Salzberg SL. 2019. Graph-based genome
49
50 alignment and genotyping with HISAT2 and HISAT-genotype. *Nat. Biotechnol.* 37:907–
51
52 915.
53
54
55 LaJeunesse TC, Parkinson JE, Gabrielson PW, Jeong HJ, Reimer JD, Voolstra CR,
56
57
58
59
60

1
2
3 Santos SR. 2018. Systematic Revision of Symbiodiniaceae Highlights the Antiquity and
4 Diversity of Coral Endosymbionts. *Curr. Biol. CB* 28:2570-2580.e6.
5

6
7 Lawrence JG, Ochman H. 1997. Amelioration of bacterial genomes: rates of change and
8 exchange. *J. Mol. Evol.* 44:383–397.
9

10
11 Lawrence SA, Floge SA, Davy JE, Davy SK, Wilson WH. 2017. Exploratory analysis of
12 Symbiodinium transcriptomes reveals potential latent infection by large dsDNA viruses.
13
14
15
16
17
18 *Environ. Microbiol.* 19:3909–3919.

19
20 Lawrence SA, Wilson WH, Davy JE, Davy SK. 2014. Latent virus-like infections are
21 present in a diverse range of Symbiodinium spp. (Dinophyta). *J. Phycol.* 50:984–997.
22

23
24 Lee AM, Ivey RG, Meints RH. 1998. Repetitive DNA insertion in a protein kinase ORF
25 of a latent FSV (Feldmannia sp. virus) genome. *Virology* 248:35–45.
26

27
28 Lee R, Lai H, Malik SB, Saldarriaga JF, Keeling PJ, Slamovits CH. 2014. Analysis of
29 EST data of the marine protist *Oxyrrhis marina*, an emerging model for alveolate biology
30 and evolution. *BMC Genomics* 15:122.
31
32
33

34
35 Lesser MP, Farrell JH. 2004. Exposure to solar radiation increases damage to both host
36 tissues and algal symbionts of corals during thermal stress. *Coral Reefs* 23:367–377.
37

38
39 Levin RA, Voolstra CR, Weynberg KD, van Oppen MJH. 2017. Evidence for a role of
40 viruses in the thermal sensitivity of coral photosymbionts. *ISME J.* 11:808–812.
41
42

43
44 Levin RA, Voolstra CR, Weynberg KD, van Oppen MJH. 2017. Evidence for a role of
45 viruses in the thermal sensitivity of coral photosymbionts. *ISME J.* 11:808–812.
46
47

48
49 Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G,
50

51 Durbin R, 1000 Genome Project Data Processing Subgroup. 2009. The Sequence
52 Alignment/Map format and SAMtools. *Bioinformatics* 25:2078–2079.
53
54
55

- 1
2
3 Lin S, Cheng S, Song B, Zhong X, Lin X, Li W, Li L, Zhang Y, Zhang H, Ji Z, et al.
4
5 2015. The Symbiodinium kawagutii genome illuminates dinoflagellate gene expression
6
7 and coral symbiosis. *Science* 350:691–694.
8
9
10 Liu H, Stephens TG, González-Pech RA, Beltran VH, Lapeyre B, Bongaerts P, Cooke I,
11
12 Aranda M, Bourne DG, Forêt S, et al. 2018. Symbiodinium genomes reveal adaptive
13
14 evolution of functions related to coral-dinoflagellate symbiosis. *Commun. Biol.* 1:1–11.
15
16
17 Lohr J, Munn CB, Wilson WH. 2007. Characterization of a Latent Virus-Like Infection
18
19 of Symbiotic Zooxanthellae. *Appl. Environ. Microbiol.* 73:2976–2981.
20
21
22 Lye L-F, Akopyants NS, Dobson DE, Beverley SM. A Narnavirus-Like Element from the
23
24 Trypanosomatid Protozoan Parasite *Leptomonas seymouri*. *Genome Announc.* 4:e00713-
25
26 16.
27
28
29 Martin M. 2011. Cutadapt removes adapter sequences from high-throughput sequencing
30
31 reads. *EMBnet.journal* 17:10–12.
32
33
34 McKeown DA, Stevens K, Peters AF, Bond P, Harper GM, Brownlee C, Brown MT,
35
36 Schroeder DC. 2017. Phaeoviruses discovered in kelp (Laminariales). *ISME J.* 11:2869–
37
38 2873.
39
40
41 Méheust R, Bhattacharya D, Pathmanathan JS, McInerney JO, Lopez P, Bapteste E.
42
43 2018. Formation of chimeric genes with essential functions at the origin of eukaryotes.
44
45 *BMC Biol.* 16:30.
46
47
48 Messyasz A, Rosales SM, Mueller RS, Sawyer T, Correa AMS, Thurber AR, Vega
49
50 Thurber R. 2020. Coral Bleaching Phenotypes Associated With Differential Abundances
51
52 of Nucleocytoplasmic Large DNA Viruses. *Front. Mar. Sci.* 7:789.
53
54 Mette MF, Kanno T, Aufsatz W, Jakowitsch J, van der Winden J, Matzke MA, Matzke AJM.
55
56
57
58
59
60

1
2
3 2002. Endogenous viral sequences and their potential contribution to heritable virus resistance in
4 plants. *EMBO J.* 21:461–469.

7
8 Minh BQ, Nguyen MAT, von Haeseler A. 2013. Ultrafast Approximation for
9
10 Phylogenetic Bootstrap. *Mol. Biol. Evol.* 30:1188–1195.

11
12 Moniruzzaman M, Weinheimer AR, Martinez-Gutierrez CA, Aylward FO. 2020.
13
14 Widespread endogenization of giant viruses shapes genomes of green algae. *Nature*
15
16 588:141–145.

17
18
19 Montalvo-Proañó J, Buerger P, Weynberg KD, van Oppen MJH. 2017. A PCR-Based
20
21 Assay Targeting the Major Capsid Protein Gene of a Dinornia-Like ssRNA Virus That
22
23 Infects Coral Photosymbionts. *Front. Microbiol.* 8:1665.

24
25
26 Mueller L, Bertelli C, Pillonel T, Salamin N, Greub G. 2017. One Year Genome
27
28 Evolution of Lausannevirus in Allopatric versus Sympatric Conditions. *Genome Biol.*
29
30 *Evol.* 9.

31
32
33 Müller DG, Kapp M, Knippers R. 1998. Viruses in Marine Brown Algae. In:
34
35 Maramorosch K, Murphy FA, Shatkin AJ, editors. *Advances in Virus Research.* Vol. 50.
36
37 Academic Press. p. 49–67. Available from:
38
39 <http://www.sciencedirect.com/science/article/pii/S0065352708608052>

40
41
42 Nagasaki K, Shirai Y, Takao Y, Mizumoto H, Nishida K, Tomaru Y. 2005. Comparison
43
44 of genome sequences of single-stranded RNA viruses infecting the bivalve-killing
45
46 dinoflagellate *Heterocapsa circularisquama*. *Appl. Environ. Microbiol.* 71:8888–8894.

47
48
49 Nand A, Zhan Y, Salazar OR, Aranda M, Voolstra CR, Dekker J. 2021. Genetic and
50
51 spatial organization of the unusual chromosomes of the dinoflagellate *Symbiodinium*
52
53 *microadriaticum*. *Nat. Genet.* 53:618–629.

- 1
2
3 Nguyen L-T, Schmidt HA, von Haeseler A, Minh BQ. 2015. IQ-TREE: A Fast and
4 Effective Stochastic Algorithm for Estimating Maximum-Likelihood Phylogenies. *Mol.*
5
6 *Biol. Evol.* 32:268–274.
7
8
9
10 Nibert ML, Manny AR, Debat HJ, Firth AE, Bertini L, Caruso C. 2018. A barnavirus
11 sequence mined from a transcriptome of the Antarctic pearlwort *Colobanthus quitensis*.
12
13 *Arch. Virol.* 163:1921.
14
15
16 Rådecker N, Pogoreutz C, Gegner HM, Cárdenas A, Roth F, Bougoure J, Guagliardo P,
17
18 Wild C, Pernice M, Raina J-B, et al. 2021. Heat stress destabilizes symbiotic nutrient
19 cycling in corals. *Proc. Natl. Acad. Sci.* 118:e2022653118.
20
21
22 Rancurel C, Legrand L, Danchin EGJ. 2017. Alieness: Rapid Detection of Candidate
23
24 Horizontal Gene Transfers across the Tree of Life. *Genes* [Internet] 8. Available from:
25
26 <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5664098/>
27
28
29
30 Randall CJ, van Woesik R. 2015. Contemporary white-band disease in Caribbean corals
31 driven by climate change. *Nat. Clim. Change* 5:375–379.
32
33
34 Rastgou M, Habibi MK, Izadpanah K, Masenga V, Milne RG, Wolf YI, Koonin EV,
35
36 Turina M. 2009. Molecular characterization of the plant virus genus Ourmiavirus and
37
38 evidence of inter-kingdom reassortment of viral genome segments as its possible route of
39
40 origin. *J. Gen. Virol.* 90:2525.
41
42
43
44 Revell LJ. 2012. phytools: an R package for phylogenetic comparative biology (and other
45
46 things). *Methods Ecol. Evol.* 3:217–223.
47
48
49 Robinson JT, Thorvaldsdóttir H, Winckler W, Guttman M, Lander ES, Getz G, Mesirov
50
51 JP. 2011. Integrative genomics viewer. *Nat. Biotechnol.* 29:24–26.
52
53
54 Roossinck MJ, García-Arenal F. 2015. Ecosystem simplification, biodiversity loss and
55
56
57
58
59
60

1
2
3 plant virus emergence. *Curr. Opin. Virol.* 10:56–62.

4
5 Rousvoal S, Bouyer B, López-Cristoffanini C, Boyen C, Collén J. 2016. Mutant swarms
6
7 of a totivirus-like entities are present in the red macroalga *Chondrus crispus* and have
8
9 been partially transferred to the nuclear genome. *J. Phycol.* 52:493–504.

10
11 Schliep KP. 2011. phangorn: phylogenetic analysis in R. *Bioinforma. Oxf. Engl.* 27:592–
12
13 593.

14
15 Shen W, Le S, Li Y, Hu F. 2016. SeqKit: A Cross-Platform and Ultrafast Toolkit for
16
17 FASTA/Q File Manipulation. *PLOS ONE* 11:e0163962.

18
19 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,
20
21 et al. 2016. Redefining the invertebrate RNA virosphere. *Nature* 540:539–543.

22
23 Shoguchi E, Beedessee G, Tada I, Hisata K, Kawashima T, Takeuchi T, Arakaki N, Fujie
24
25 M, Koyanagi R, Roy MC, et al. 2018. Two divergent Symbiodinium genomes reveal
26
27 conservation of a gene cluster for sunscreen biosynthesis and recently lost genes. *BMC*
28
29 *Genomics* 19:458.

30
31 Shoguchi E, Shinzato C, Kawashima T, Gyoja F, Mungpakdee S, Koyanagi R, Takeuchi
32
33 T, Hisata K, Tanaka M, Fujiwara M, et al. 2013. Draft Assembly of the Symbiodinium
34
35 minutum Nuclear Genome Reveals Dinoflagellate Gene Structure. *Curr. Biol.* 23:1399–
36
37 1408.

38
39 Slamovits C, Keeling P. 2008. Widespread recycling of processed cDNAs in
40
41 dinoflagellates. *Curr. Biol. CB* 18:R550-2.

42
43 Slamovits C, Keeling P. 2008. Widespread recycling of processed cDNAs in
44
45 dinoflagellates. *Curr. Biol. CB* 18:R550-2.

46
47 Slater GSC, Birney E. 2005. Automated generation of heuristics for biological sequence
48
49

1
2
3 comparison. *BMC Bioinformatics* 6:31.

4
5 Sömera M, Fargette D, Hébrard E, Sarmiento C, Consortium IR. 2021. ICTV Virus
6
7 Taxonomy Profile: Solemoviridae 2021. *J. Gen. Virol.* [Internet] 102. Available from:
8
9 <https://www.ncbi.nlm.nih.gov/labs/pmc/articles/PMC8744267/>

10
11 Song B, Morse D, Song Y, Fu Y, Lin X, Wang W, Cheng S, Chen W, Liu X, Lin S.
12
13 2017. Comparative Genomics Reveals Two Major Bouts of Gene Retroposition
14
15 Coinciding with Crucial Periods of Symbiodinium Evolution. *Genome Biol. Evol.* 9:2037.

16
17 Song B, Morse D, Song Y, Fu Y, Lin X, Wang W, Cheng S, Chen W, Liu X, Lin S.
18
19 2017. Comparative Genomics Reveals Two Major Bouts of Gene Retroposition
20
21 Coinciding with Crucial Periods of Symbiodinium Evolution. *Genome Biol. Evol.* 9:2037.

22
23 Staginuss C, Gregor W, Mette MF, Teo CH, Borroto-Fernández EG, Machado ML da C,
24
25 Matzke M, Schwarzacher T. 2007. Endogenous pararetroviral sequences in tomato
26
27 (Solanum lycopersicum) and related species. *BMC Plant Biol.* 7:24.

28
29 Stephens TG, Bhattacharya D, Ragan MA, Chan CX. 2016. PhySortR: a fast, flexible
30
31 tool for sorting phylogenetic trees in R. *PeerJ* 4:e2038.

32
33 Stephens TG, González-Pech RA, Cheng Y, Mohamed AR, Burt DW, Bhattacharya D,
34
35 Ragan MA, Chan CX. 2020. Genomes of the dinoflagellate *Polarella glacialis* encode
36
37 tandemly repeated single-exon genes with adaptive functions. *BMC Biol.* 18:56.

38
39 Takahashi H, Fukuhara T, Kitazawa H, Kormelink R. 2019. Virus Latency and the
40
41 Impact on Plants. *Front. Microbiol.* 10:2764.

42
43 Than TT, Tran GVQ, Son K, Park E-M, Kim S, Lim Y-S, Hwang SB. 2016. Ankyrin
44
45 Repeat Domain 1 is Up-regulated During Hepatitis C Virus Infection and Regulates
46
47 Hepatitis C Virus Entry. *Sci. Rep.* 6:20819.

- 1
2
3 Thurber RLV, Correa AMS. 2011. Viruses of reef-building scleractinian corals. *J. Exp.*
4
5 *Mar. Biol. Ecol.* 408:102–113.
6
7
8 Tripathi JN, Ntui VO, Ron M, Muiruri SK, Britt A, Tripathi L. 2019. CRISPR/Cas9
9
10 editing of endogenous banana streak virus in the B genome of *Musa* spp. overcomes a
11
12 major challenge in banana breeding. *Commun. Biol.* 2:1–11.
13
14
15 Waldron FM, Stone GN, Obbard DJ. 2018. Metagenomic sequencing suggests a diversity
16
17 of RNA interference-like responses to viruses across multicellular eukaryotes. *PLOS*
18
19 *Genet.* 14:e1007533.
20
21
22 Weynberg KD, Neave M, Clode PL, Voolstra CR, Brownlee C, Laffy P, Webster NS,
23
24 Levin RA, Wood-Charlson EM, van Oppen MJH. 2017. Prevalent and persistent viral
25
26 infection in cultures of the coral algal endosymbiont *Symbiodinium*. *Coral Reefs* 36:773–
27
28 784.
29
30
31 Weynberg KD, Wood-Charlson EM, Suttle CA, van Oppen MJH. 2014. Generating viral
32
33 metagenomes from the coral holobiont. *Front. Microbiol.* 5:206.
34
35
36 Wilson WH, Dale AL, Davy JE, Davy SK. 2005. An enemy within? Observations of
37
38 virus-like particles in reef corals. *Coral Reefs* 24:145–148.
39
40
41 Wilson WH, Francis I, Ryan K, Davy SK. 2001. Temperature induction of viruses in
42
43 symbiotic dinoflagellates. *Aquat. Microb. Ecol.* 25:99–102.
44
45
46 Wisecaver JH, Brosnahan ML, Hackett JD. 2013. Horizontal gene transfer is a significant
47
48 driver of gene innovation in dinoflagellates. *Genome Biol. Evol.* 5:2368–2381.
49
50
51 Wisecaver JH, Brosnahan ML, Hackett JD. 2013. Horizontal gene transfer is a significant
52
53 driver of gene innovation in dinoflagellates. *Genome Biol. Evol.* 5:2368–2381.
54
55
56 Wood-Charlson EM, Weynberg KD, Suttle CA, Roux S, van Oppen MJH. 2015.
57
58
59
60

- 1
2
3 Metagenomic characterization of viral communities in corals: mining biological signal
4 from methodological noise. *Environ. Microbiol.* 17:3440–3449.
- 5
6
7 Xu Q, Zhu N, Chen S, Zhao P, Ren H, Zhu S, Tang H, Zhu Y, Qi Z. 2017. E3 Ubiquitin
8 Ligase Nedd4 Promotes Japanese Encephalitis Virus Replication by Suppressing
9 Autophagy in Human Neuroblastoma Cells. *Sci. Rep.* 7:45375.
- 10
11
12 Younis S, Kamel W, Falkeborn T, Wang H, Yu D, Daniels R, Essand M, Hinkula J,
13 Akusjärvi G, Andersson L. 2018. Multiple nuclear-replicating viruses require the stress-
14 induced protein ZC3H11A for efficient growth. *Proc. Natl. Acad. Sci. U. S. A.*
15 115:E3808.
- 16
17
18 Yu G, Smith DK, Zhu H, Guan Y, Lam TT-Y. 2017. ggtree: an r package for
19 visualization and annotation of phylogenetic trees with their covariates and other
20 associated data. *Methods Ecol. Evol.* 8:28–36.
- 21
22
23
24
25
26
27
28
29
30
31
32
33
34

35 Figure legends

36
37
38 Fig. 1: Main characteristics of genes in Symbiodiniaceae genomes that were acquired by viral
39 horizontal gene transfer (vHGT). Cladogram depicting phylogenetic relationship of
40 Symbiodiniaceae genomes in which vHGTs were identified (left); the hosts associated with each
41 Symbiodiniaceae species are represented by colored circles (top left key). The total number of
42 vHGTs and their putative taxonomic origins (at the order level; center key), the coding sequence
43 (CDS) percent Guanine-Cytosine (GC%) content, the average number of introns, and the total
44 number of associated protein domains are shown for the vHGTs identified in each genome. The
45 CDS GC% content and average number of introns is shown separately for the vHGTs and the
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 background Symbiodiniaceae sequences (i.e., all Symbiodiniaceae CDS or genes [respectively]
4 excluding the vHGTs). Abbreviations: *Breviolum minutum* (Bmin); *Cladocopium* sp. C92
5 (CC92); *Cladocopium goreau* (Cgor); *Fugacium kawagutii* (Fkaw); *Symbiodinium linucheae*
6 (Slin_CCMP2456); *S. microadriaticum* (Smic); *S. microadriaticum* 04-503SCI.03
7 (Smic_04-503SCI.03); *S. microadriaticum* CassKB8 (Smic_CassKB8); *S. natans* CCMP2548
8 (Snat_CCMP2548); *S. necroappetens* CCMP2469 (Snec_CCMP2469); *S. pilosum* CCMP2461
9 (Spil_CCMP2461); *S. tridacnidorum* (Stri); *S. tridacnidorum* CCMP2592 (Stri_CCMP2592);
10 DNA/RNA pol (DNA/RNA polymerase superfamily); DUF4116 (Domain of unknown function
11 DUF4116); None (None predicted); RdRp (RNA dependent RNA polymerase); RNA hel core
12 (RNA virus helicase core domain); R-dRP (RNA-directed RNA polymerase); R-dRp core (RNA-
13 directed RNA polymerase catalytic domain); RNA hel (Viral (Superfamily 1) RNA helicase);
14 Vir hel (Viral helicase1).

15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33 Fig. 2: Phylogenetic profile of vHGT's grouped first by taxonomic affiliation (viral order) and
34 second by putative functional annotation. Unclassified RNA viruses: RNA dependent RNA
35 polymerase (RdRps-like group 1 [A, B, C] and RdRp-like group 2 [D]), Major capsid protein
36 (MCPs-like [E, F, G]), viral RNA helicase (H) and polyprotein coding for replicases including
37 RNA-dependent RNA polymerase region (I); Pimascovirales: restriction endonuclease
38 (Endonuclease) (J); and Mononegavirales: RNA-directed RNA polymerase catalytic domain
39 (RdRP catalytic domain) (K and L). Sequence names are not shown in the trees to enhance
40 readability; the taxonomic affiliation of each sequence in the trees is represented by a colored
41 circle, with the colors described in the legend at the bottom of the figure. Symbiodiniaceae
42 vHGTs are highlighted with thicker borders; only bootstrap support values > 95% (from 10000
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 bootstrap replicates) are shown. The full phylogenies with sequence names and branch support
4 values are included in Supplementary Fig. S4.
5
6
7
8
9

10 Fig. 3: Genome maps comparing the coding sequences of the dinoflagellate virus (dinornavirus)
11 HcRNAV34 (*Heterocapsa circularisquama* RNA virus - accession: AB218608.1) and Symbiod.
12 RNA virus (*Symbiodinium* +ssRNA virus - accession KX538960.1) with the putative integrated
13 viral genomes in *S. microadriaticum* (Smic and Smic04). The putative viral elements (vHGTs)
14 annotated as RdRp-like and MCP-like proteins are shown, as are the genes surrounding these
15 elements. Each feature points in the direction that it is encoded and a scale bar is shown in the
16 top right corner. Slash bars (/) denote intergenic regions. Abbreviations: scaf. (scaffold); DinoSL
17 (dinoflagellate splice leader; annotated as small arrows); RdRp and RdRp-like (RNA dependent
18 RNA polymerase); MCP and MCP-like (major capsid protein); Retrovirus-related Pol
19 polyprotein R2 (Retrovirus-related Pol polyprotein from type-1 retrotransposable element R2);
20 Amino acid permease YfnA (putative amino acid permease YfnA).
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37

38 Fig. 4: hypothetical model with two scenarios (A and B) explaining the origin of vHGTs in
39 Symbiodiniaceae genomes. During an acute infection, viral particles enter the host cell and viral
40 mRNAs accumulate in the cytoplasm to a high abundance. Host machinery can accidentally
41 incorporate viral mRNA into its genome *via* reverse transcription and non-homologous
42 recombination; this may use the same (or similar) mechanisms as the mRNA recycling process
43 that is prevalent in dinoflagellates. These viral sequences, known as endogenous viral elements
44 (EVEs), would be “dead on arrival” and are expected to decay into pseudogenes (which we
45 would then observe in the genome as vHGTs). In the second scenario (B), integration of the virus
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 genome into the host genome is part of a previously unobserved transient life stage of these
4
5 viruses (specifically the *Symbiodinium* +ssRNA virus). The virus infects and is incorporated into
6
7 the host genome *via* the same mechanisms as the first scenario (A) however, instead of this
8
9 process resulting in non-functional EVEs it produces active proviral elements. These integrated
10
11 viral genomes can remain silent during times of low host abundance (e.g., during the free-living
12
13 stage of facultative Symbiodiniaceae lifecycles) and can become activated by environmental
14
15 cues (such as stress or high Symbiodiniaceae cell density) that result in viral lysis, escape from
16
17 the host cell and new infections in susceptible hosts. The proviral elements can also become
18
19 deactivated, either through random mutation or host driven process, resulting in EVEs that are
20
21 expected to decay into pseudogenes, which we would then observe in the genome as vHGTs.
22
23
24
25
26
27

28
29 Supplementary Fig. S1: Putative integration of an +ssRNA *Symbiodinium* virus genome into the
30
31 genome of *S. microadriaticum* 04-503SCI.03 (scaffold: Smic_04-503SCI.03.scaffold3501;
32
33 genes: gene15955 and gene15956). (A) Image taken from IGV showing the region of the
34
35 genome encoding the two vHGT genes (represented as blue bars; solid bars represent exons and
36
37 lines represent introns) along with aligned RNA-seq data (gray discontinuous bars). (B)
38
39 Alignment (generated using exonerate) of a RNA-dependent RNA polymerase protein
40
41 (YP_009337004.1) against a region (Smic_04-503SCI.03.scaffold3501:13547-15709) of the *S.*
42
43 *microadriaticum* 04-503SCI.03 genome corresponding to the vHGT gene15955. (C) Alignment
44
45 (generated using exonerate) of a major capsid protein (AOG17586.1) against a region (Smic_04-
46
47 503SCI.03.scaffold3501:16298-17334) of the *S. microadriaticum* 04-503SCI.03 genome
48
49 corresponding to the vHGT gene15956. In B and C, a predicted intron is highlighted using a blue
50
51 box and in-frame stop codons using red boxes.
52
53
54
55
56
57
58
59
60

1
2
3
4
5
6 Supplementary Fig. S2: Putative integration of an +ssRNA *Symbiodinium* virus genome into the
7
8 genome of *S. microadriaticum* (scaffold: scaffold792; genes: gene19249 and gene19250). (A)
9
10 Image taken from IGV showing the region of the genome encoding the two vHGT genes
11
12 (represented as blue bars; solid bars represent exons and lines represent introns) along with
13
14 aligned RNA-seq data (gray discontinuous bars). (B) Alignment (generated using exonerate) of a
15
16 major capsid protein (AOG17586.1) against a region (Smic.scaffold792:89753-90508) of the *S.*
17
18 *microadriaticum* genome corresponding to the vHGT gene19249. (C) Alignment (generated
19
20 using exonerate) of a RNA-dependent RNA polymerase protein (YP_009342067.1) against a
21
22 region (Smic.scaffold792:96851-98156) of the *S. microadriaticum* genome corresponding to the
23
24 vHGT gene19250. In B and C, a predicted intron is highlighted using a blue box, in-frame stop
25
26 codons using red boxes, and frame-shifts using orange boxes.
27
28
29
30
31
32

33 Supplementary Fig. S3: Putative integration of an +ssRNA *Symbiodinium* virus genome into the
34
35 genome of *S. microadriaticum* (scaffold: scaffold67; genes: gene3240 and gene3241). (A) Image
36
37 taken from IGV showing the region of the genome encoding the two vHGT genes (represented
38
39 as blue bars; solid bars represent exons and lines represent introns) along with aligned RNA-seq
40
41 data (gray discontinuous bars). (B) Alignment (generated using exonerate) of a major capsid
42
43 protein (AOG17586.1) against a region (Smic.scaffold67:1208081-1209117) of the *S.*
44
45 *microadriaticum* genome corresponding to the vHGT gene3240. (C) Alignment (generated using
46
47 exonerate) of a RNA-dependent RNA polymerase protein (YP_009337004.1) against a region
48
49 (Smic.scaffold67:1209706-1211868) of the *S. microadriaticum* genome corresponding to the
50
51
52
53
54
55
56
57
58
59
60

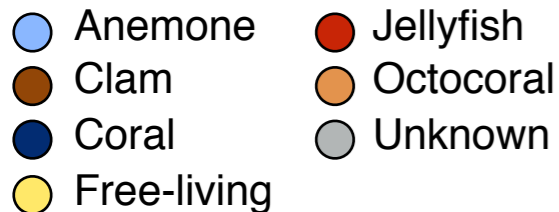
1
2
3 vHGT gene3241. In B and C, a predicted intron is highlighted using a blue box, in-frame stop
4
5 codons using red boxes, and frame-shifts using orange boxes.
6
7
8
9

10 Supplementary Fig. S4: Phylogenies organized by grouping (A-L), viral order (Unclassified
11
12 RNA viruses, Pimascovirales and Mononegavirales) of subject hits and predicted annotation,
13
14 showing complete sequence names and annotation with branch labels and associated ultrafast
15
16 bootstrap values.
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

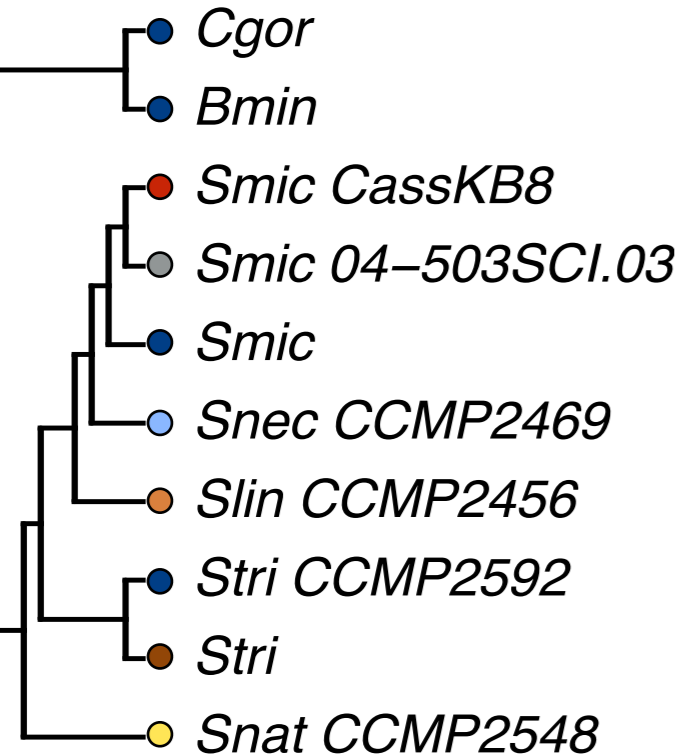
For Review Only

Symbiodiniaceae Hosts

vHGT Viral order



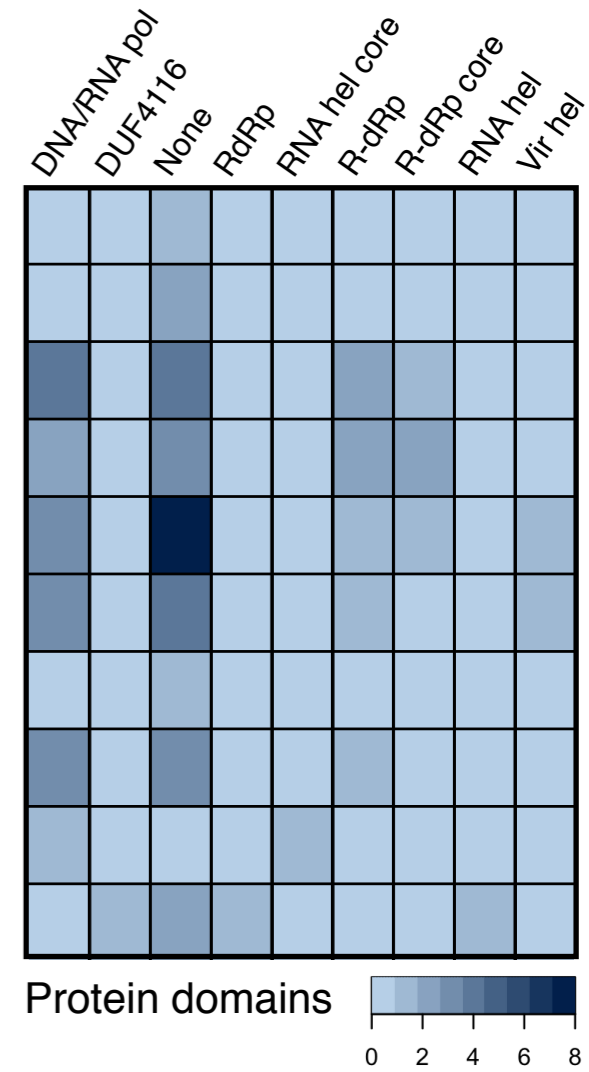
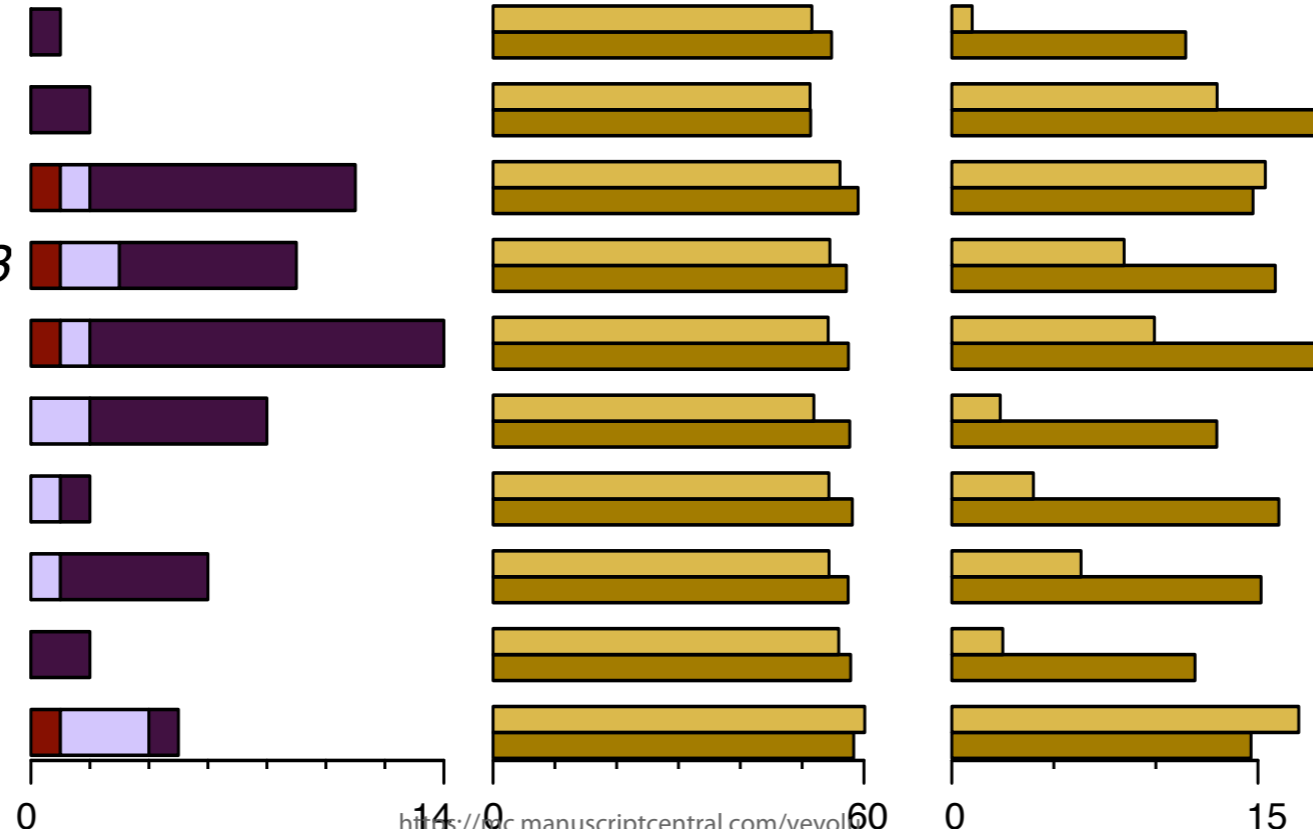
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42

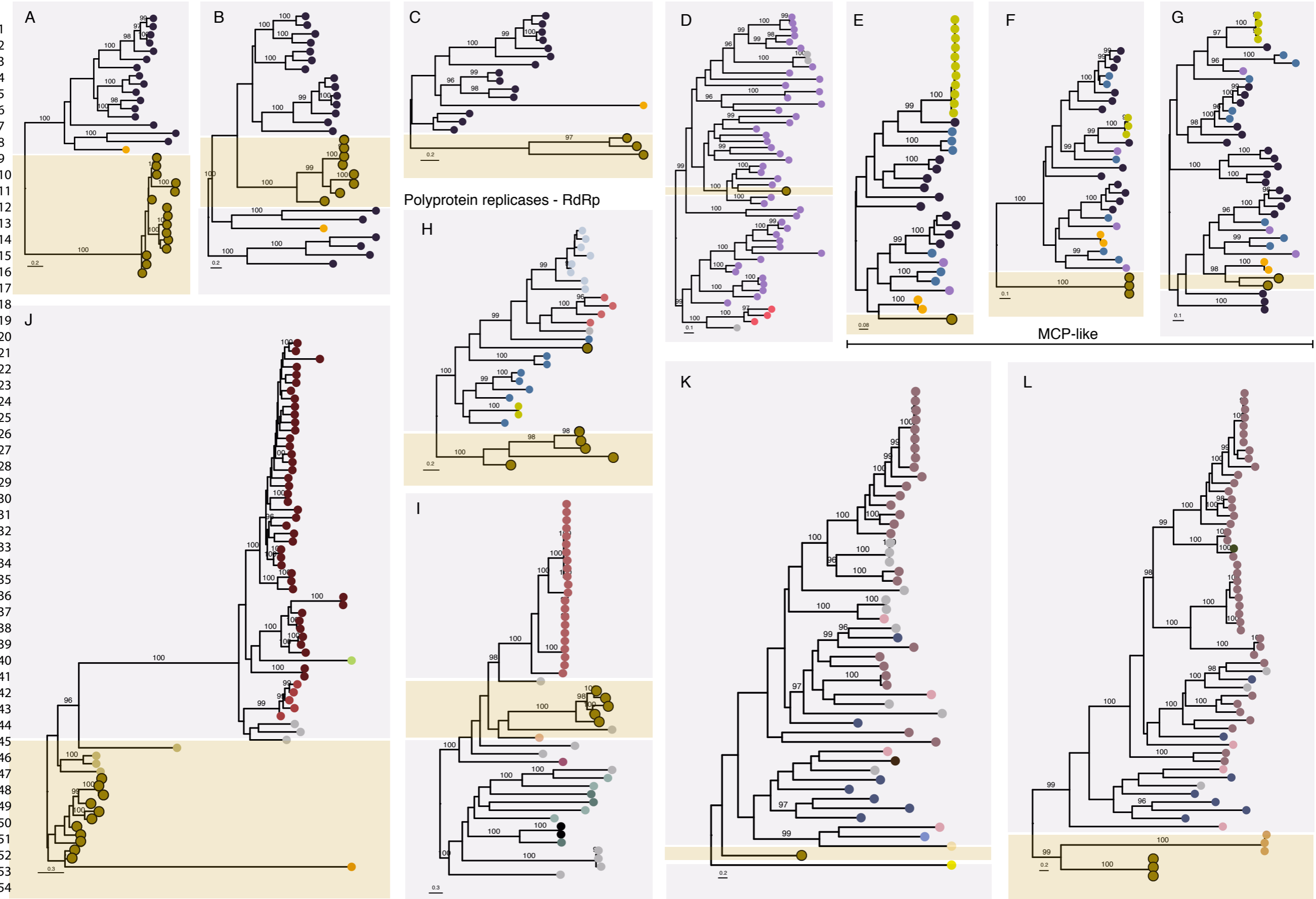


Number of vHGTs

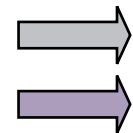
GC (%)

Number of introns





- 55 Endonucleases
- 56
- 57
- 58 Ascomycota ● Astroviridae ● Barnavirus ● Beny-like virus ● Benyvirus ● Caulimoviridae ● Dinophyceae ● Dinornavirus ● Filoviridae ● Hepe-like virus
- 59 Hepeviridae ● Higrevirus ● Insecta ● Lispiviridae ● Marseilleviridae ● Mimiviridae ● Mymonaviridae ● Narna-like virus ● Ourmiavirus ● Pelagophyceae
- 60 Perkinsidae ● Phycodnaviridae ● Rhabdo-like virus ● Rhabdoviridae ● Sobemo-like virus ● Solemoviridae ● Symbiodiniaceae ● Symbiodiniaceae RNA virus ● Unclassified viruses ● Weivirus-like virus
- <https://mc.manuscriptcentral.com/vevolu>



Eukaryotic genes



Major capsid protein (MCP)



Retroelements

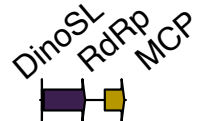


RNA-dependent RNA polymerase (RdRp)

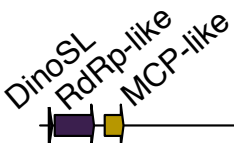
HcRNAV34



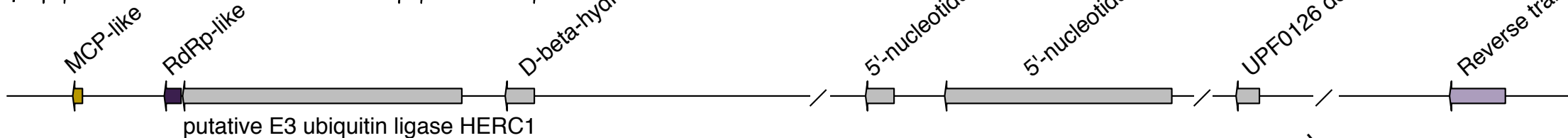
Symbiod. RNA virus



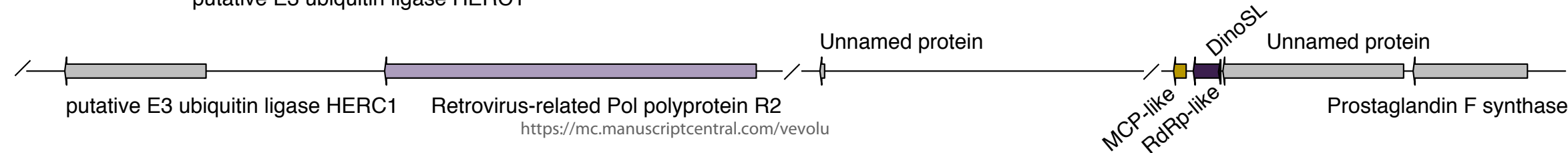
Smic04 (scaf. 3501)



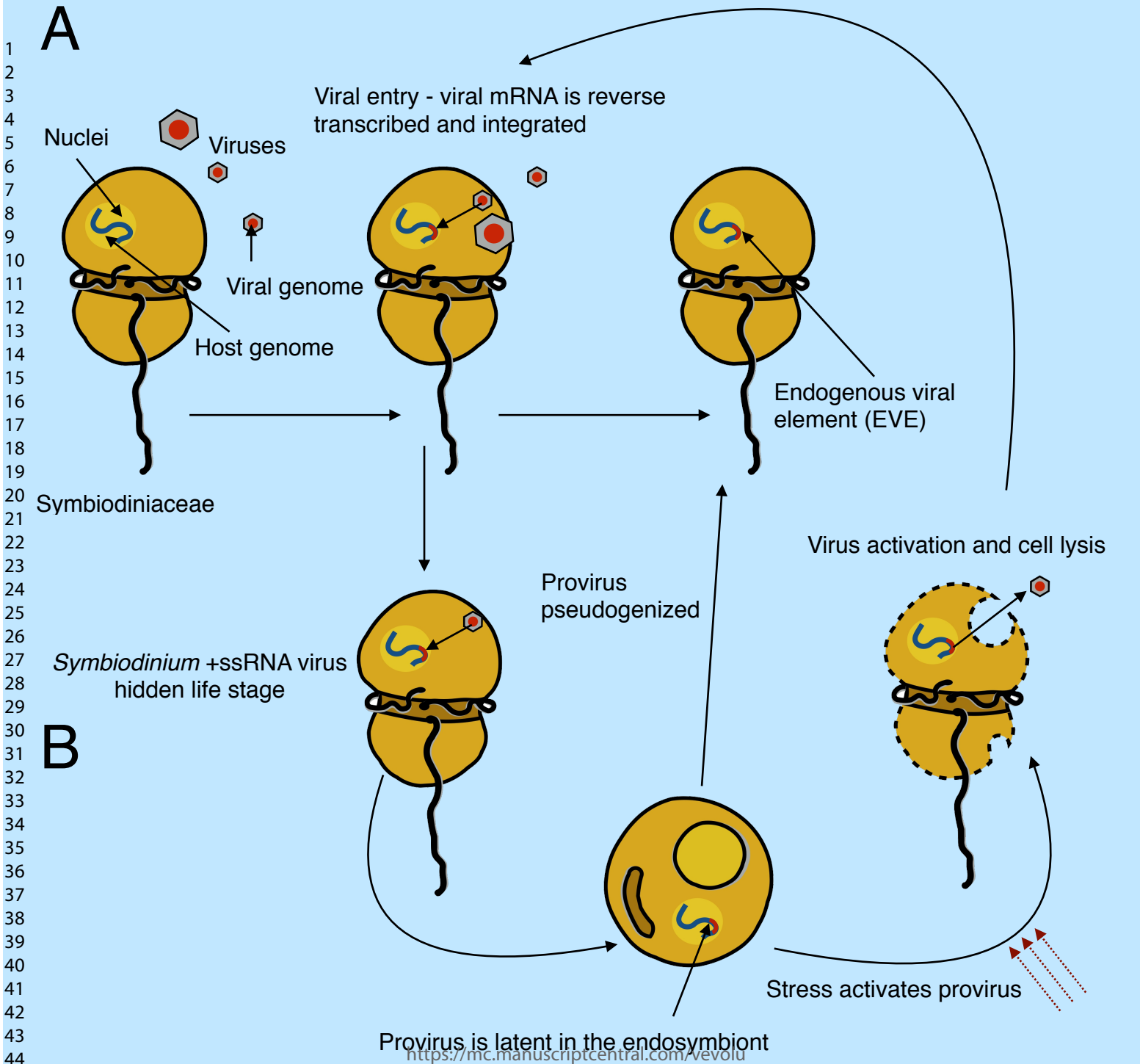
Smic (scaf. 792)



Smic (scaf. 67)



25 kb





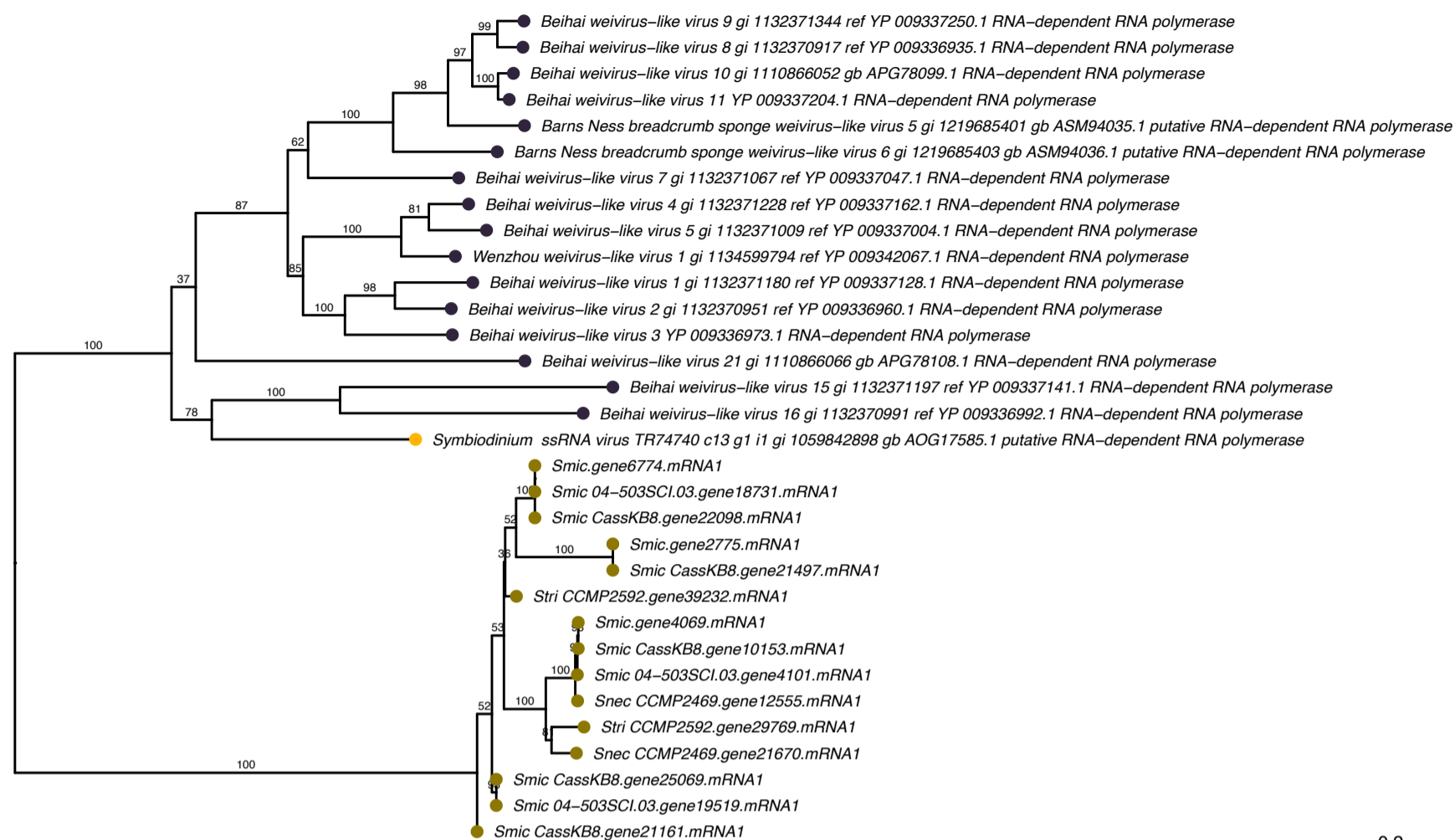
B Query: YP_009337004.1 RNA-dependent RNA polymerase [Beihai weivirus-like virus 5]
Target: Smic_04-503SCI.03.scaffold3501:13547-15709:+ (Smic_04-503SCI.03.gene15955)
Model: protein2genome:local
Raw score: 545
Query range: 357 -> 767 (Length:820aa; 50% query coverage)
Target range: 2639 -> 3938

358 : AspLysArgCysGlyValPheAlaLysHisArgIleGluGluTrpAlaIleAlaHisPheAs : 378
:!!::!!::!! !!!!!!!!:!!::!!::!! !!::!! !! !
21 AsnArgLysTyrHisValPheSerLysGluGlnIleAspAlaGluIleValSerIlePheHi
22 2640 : AACCGAAGTATACAGTCTTCAGCAAAGAGCAGATTGTGCTGAAATTGTGAGCATAATCCA : 2700
23
24 379 : pLeuGluGluCysLysSerGlyLysTrpSerIleGluArgPheArgGlySerLeuGluAsnL : 399
:!! !:!! ! ! !!!!!!!! :!!!!!!!..!! !!!!!..!! !
25 sValProGlnPheAlaSerLysLysTrpThrGlnLysArgPheGluAsnMetPheAsnHisL
26 2701 : TGTGCCACAATTCGCCTCCAAGAAATGGACACAGAAAAGATTTGAGAACATGTTCAATCATC : 2763
27
28 400 : euTyrAlaLysGluHisProThrPheSerPheLysAlaAspValLysTyrGlu<-><-> : 417
|| ! !:!! ! !!!!! !!!!! !!!!! !!!!! !!!!! ||
29 euLeuLysGlnValAspProArgPheLysSerLysAlaLysIleLysLeuGluAlaMetGly
30 2764 : TCCTCAACAGGTTGGATCCAGATTCAAGAGCAAGCCAAAATAAGCTAGAAGCAATGGGC : 2823
31
32 418 : CysMetProGluGly<-><->LysAlaProArgMetLeuIleAlaAspGlyAspGluGlyGl : 435
!!!! !!!!!:!!! || !!!!!!!!:!!!!!!
33 CysLysProAspGlySerProLysProProArgLeuLeuIleAlaAspGlyAspGluGlyGl
34 2824 : TGCAACCCGGATGGCTCGCAAAGCCACCCCGTCTGCTTATAGCAGACGGGGATGAAGGGCA : 2886
35
36 436 : nLeuMetAlaLeuAlaValValLysCysPheGluGluLeuLeuPheSerHisPheGluThrL : 456
|:!!!!:!!!! ! !:!! ! !!!!!:!!!!!! ! ! !!!!!:!!
37 nIleMetSerLeuLeuAspIleAlaIlePheGluLysLeuLeuPheArgLysPheHisSerA
38 2887 : AATCATGTCGCTCCGACATAGCAATTTTGGAGAAATGCTCTCCGGAGTCCACAGTA : 2949
39
40 457 : ysSerIleLysHisLeuAlaLysArgAspAlaIleAspArgValLeuLysGluLeuArgAla : 476
:!!!!!! ! !!!!!:!!!!!!:!!:!!:!! !!!!!:!!:!! !!!!!!!
41 rgSerIleLysGlyArgSerArgGlnValGluGlnAspValGluTyrLeuArgPro
42 2950 : GGAGCATAAAGCCGCTCAAGGAGACAAGTGTGCAGGACGTGGTTCAGTACTTACGTCCC : 3009
43
44 477 : ProGlyAlaLys >>>> Target Intron 1 >>>> AlaValGluGlyAspGlySer : 487
! ! ! ! 33 bp !!!!!!! !!!!!!!
45 SerLysLysHis+- -+MetValGluGlyHisGlySer
46 3010 : AGCAAGAACATga.....cgATGGTGAAGGACATGGTTCC : 3075
47
48 488 : AlaTrpAspThrThrCysAsnValLeuIleArgGlyLeuValGluAsnProValLeuArgHi : 508
!!!!!! !!!!!: ! !!!!! !!!!!!! !!!!!!!
49 AlaTrpAsp---CysCysSerLysGluLeuArgAspMetValGluAsnProValLeuArgHi
50 3076 : GCTTGGGAC---TGCTGCTCGAAAGAGCTTAGAGACATGGTAGAAAATCCGGTACTTCGACA : 3135
51
52 509 : sIleThrThrValLeuCysAsnPheGlyValIleProSerThrTrpMetGluGluHisGlnA : 529
|!!!:!!!! !!!!!: !!!!!:!! ! !!!!!:!! ! !!!!!:!! !
53 sIleAlaThrHisLeuMetAspTyrTyrLeuValProProGlnTrpGluGlnGluHisAlaA
54 3136 : CATCGCGACACATCTGATGGATTACTACCTCGTCCCGCTCAATGGGAGCAAGAACACGCGA : 3198
55
56 530 : rgAlaCysGluLysLysThrLeuArgLeuPhePheSerAsnLysPheGluThrMetSerThr : 549
|!!! ! ! !!!!!:!!!! !!!!!:!! ! !!!!!:!! ! !!!!!:!! !
57 rgThrAsnThrAlaAspArgTyrAsnLeuIlePheArgAspLysLeuMetThrTyrPheVal
58 3199 : GGACCAACACCCGGGACCGCTACACCTCATCTCCGGGACAAGCTGATGACCTACTTTGTA : 3258
59
60 550 : SerIleAspAlaIleArgArgSerGlyHisArgGlyThrSerCysLeuAsnTrpTrpIleAs : 570
..! !!!!! !!!!!:!!!!!! !!!!!:!! !!!!!:!! !!!!!:!! !
3259 : GACTAGAAAGAACACGTCGATCAGCCATAGAGGTACGTCCTGCTGAACTGGTGGTCAA : 3321
571 : nPheValLeuTrpValSerSerValPheLysGluProGluArgPheLeuAspValAlaValA : 591
!!!!!! !!!!!:!! !!!!!:!! !!!!!:!! !!!!!:!! ! !!!!!:!! !
3322 : CTTCGTGTTGGAGCGCATCAGTATCCATCAACCTTGGGTATTGCTCTATGCCAACAGG : 3384
592 : rgAsnGlyThrAspLeuThrGlyArgSerArgTrpTrpAsnGlyCysPheGluGlyAspAsp : 611
! !!!!! !!!!!:!!!!!! !!!!! !!!!! !!!!! !!!!! !!!!! !!!!!
3385 : ACTTTTGAAGGATATAAATGGGAAGAGGCTTTGGATCAGGATCGTCTGGAGGAGACGAT : 3444
612 : SerLeu<-><->CysThrMetArgProProMetValGluGlyAspAlaLeuCysGlnValPh : 630
!!!! ! ! !!!!! ! !!!!! !!!!! !!!!! !!!!! !!!!! !!!!! !!!!!
3445 : TCCTGGAGGGATGTCGCCGAGGCTCGCCGATCTGAAGAGGACCCGAGAACAAAGTACTT : 3507
631 : eLeuAlaPheTrpLysSerAlaGlyPheAsnMetLysIleValPheCysLysThrArg<-><-> : 650
!!!! !!!!!:!! !!!!!:!! !!!!!:!! !!!!!:!! !!!!!:!! !!!!!:!! !
3508 : TCTCGACTATTGGAAACGCCAGGGTTTCGATATGAAAATCCGCCAATGTGGGGTTCGTCCAG : 3570
651 : -><-><-><-><->AlaThrPheValGlyTrpHisValGlyCysThrAspGly<->GluLeu : 664
|| !!!!!:!!!! !!!!! !!!!! !!!!! !!!!! !!!!! !!!!! !!!!!
3571 : ACACGAAGCCATAGGCCATGTTTCATTGGGACTCATTTCATGTAGATGACCACTGGATCTC : 3630
665 : AsnAsnPheArgCysProGluLeuProArgAlaLeuAlaAsnSerGlyValSerValSerPr : 685
!....! ! ! !!!!!:!!!!!! !!!!!:!! !!!!!:!! !!!!!:!! !!!!!:!! !
3631 : GAGGCCACCTTTGTGCCGACCTGCCAGGAATCTGACGAACAATTGG---TCAACAACCCC : 3690
686 : oGluAlaIleLysAlaAlaLysAspMetAsnArgSerAlaValAsnValLeuAlaAlaAlaS : 706
|| !!!!!:!!!! !!!!!:!! !!!!!:!! !!!!!:!! !!!!!:!! !!!!!:!! !
3691 : GGGATGATACAGTGTACGAACAGCAGAAATCCACCTGGTGAGGCAACACGCCAGCAGCAG : 3753
707 : erAlaLeuAlaArgAlaSerAspPheAlaGlyIleLeuProSerValSerValLysTyrMet : 726
!!!!!! !!!!!:!! !!!!!:!! !!!!!:!! !!!!!:!! !!!!!:!! !!!!!:!! !
3754 : CAGCAGTTCACGCTGCTTGGATTATGCAGGCATACTACCGATGCTGCTGATGAAATATGTC : 3813
727 : AspPheAlaGluSerValSerArgThrAspPheSerAspArgGluMetSerIleArgAlaPh : 747
..!:!!!! ! ! ! ! !!!!! !!!!!:!! !!!!!:!! !!!!!:!! !!!!!:!! !
3814 : CAATATGCCATCAATGCTACGCTGGTACTCCACAATGAAGATCTCTGATCTTGCCAC : 3876
748 : eGlyGluAspGlyPheSerAlaAsnAlaValArgThrGlnIleMetGluArgAsnIleGly : 767
!!!!!! || !!!!!:!!!!!! !!!!!:!! !!!!!:!! !!!!!:!! !!!!!:!! !
3877 : AGGTGAAGCAGGTACCACCTCACAGCTGTCATTTCCAAAGTACCAAGGCTGAATGTTGGG : 3938

C Query: AOG17586.1 putative major capsid protein [Symbiodinium +ssRNA virus TR74740 c13_g1_i1]
Target: Smic_04-503SCI.03.scaffold3501:16298-17334:+ (Smic_04-503SCI.03.gene15956)
Model: protein2genome:local
Raw score: 167
Query range: 32 -> 146 (Length:358aa; 32% query coverage)
Target range: 2083 -> 2428

33 : ArgArgSerArgAlaArgProAlaAsnGlyGlySerGlnLeuMetGlyIleLysGlnGlyVa : 53
!:!!!!:!! ! !!!!! !!!!! !!!!! !!!!! !!!!! !!!!! !!!!! !!!!!
2084 : AAGCGCAACACGATCGCAGAGAACTAAGTCTGACCAACTTGTGGGATTAGACAAGGAGT : 2144
54 : lGlyAlaIleThrArgGlyProPheGlySerAsnValAlaTyrGlyProHisCysPheAsnA : 74
!!!!!! !!!!!:!! !!!!!:!! !!!!!:!! !!!!!:!! !!!!!:!! !!!!!:!! !
2145 : CGGAGCCGCCCAAAACGGGGATACGGTTCAAGTGGCAAGTATCATTCGCGCCGCTCGACG : 2207
75 : laPheAsnTyrCysHisMetProLeuProArgAlaIleGlyProTyrThrValIleArgThr : 94
!!!!!!:!! !!!!!:!!!!!! !!!!!:!! !!!!!:!! !!!!!:!! !!!!!:!! !!!!!:!! !
2208 : CTTTTGATCTCTGCCATTTGCCTTTACCCAGGCGAGTGGTGGTAAACTGGCATACGGACA : 2267
95 : ThrArgValIleLysSerAsnLeuGlnLeuMetAsnPheGlyThr<->MetTyrAsnGluAr : 114
|| !!!!!!! !!!!!:!! !!!!!:!! !!!!!:!! !!!!!:!! !!!!!:!! !!!!!:!! !
2268 : ACCGTTGTTATCACCAGTATGATCCATTTGCCATCTTCGGACAGATGATGGTGTCTGACAC : 2330
115 : gGlnThrAlaPheAsnSerSerThrTrpSerAsnValCysAlaTrpGlyThrAsnAsnLeuA : 135
! !!!!! !!!!!:!! !!!!!:!! !!!!!:!! !!!!!:!! !!!!!:!! !!!!!:!! !
2331 : TGTGCAACAATTCAGGCGCAGAGGAGTGGAGTCAAGTTCATGCTACTCCATGCCAGATGCTG : 2393
136 : laAsnProMetAsnGlyAlaAlaAsnAlaThrArg : 146
!....! !!!!!:!!!!!! !!!!!:!! !!!!!:!! !!!!!:!! !!!!!:!! !!!!!:!! !
2394 : GCCAGCTCATCACTGGGACCAACGCTGTCACTAAG : 2428

RdRp-like group 1 (Phylogeny A)

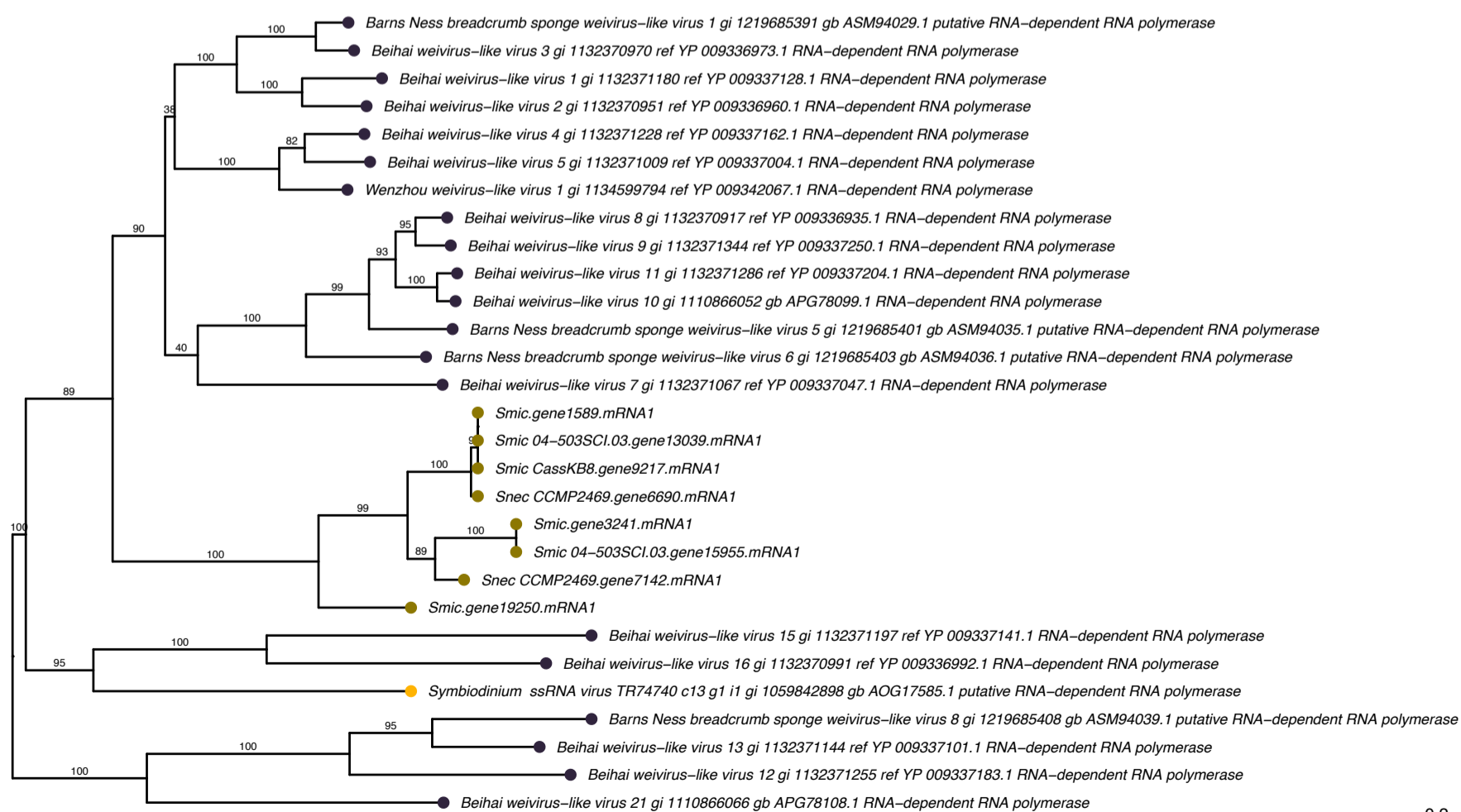


0.2

Taxa

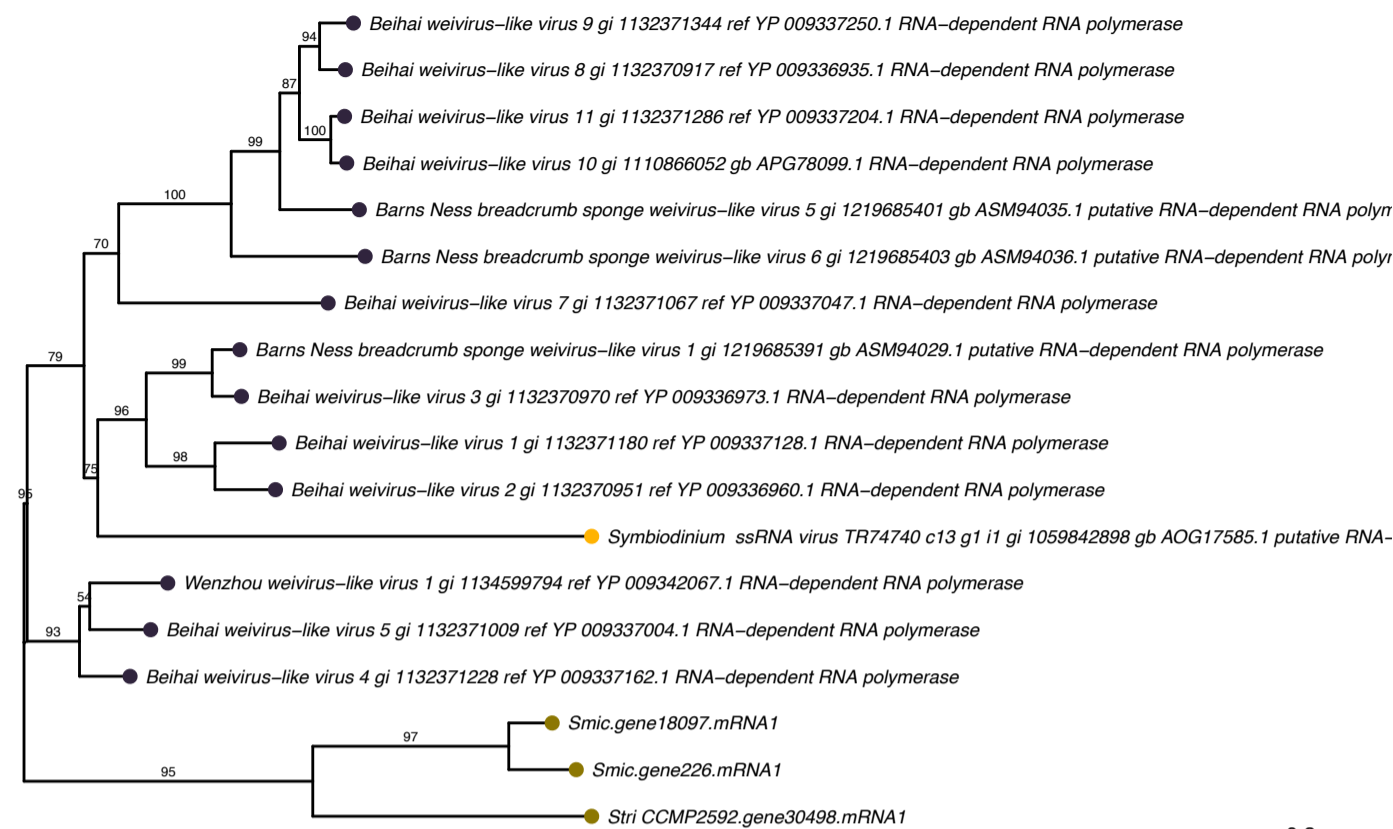
- Symbiodiniaceae
- Symbiodiniaceae RNA virus
- Weivirus-like virus

RdRp-like group 1 (Phylogeny B)



0.2

RdRp-like group 1 (Phylogeny C)

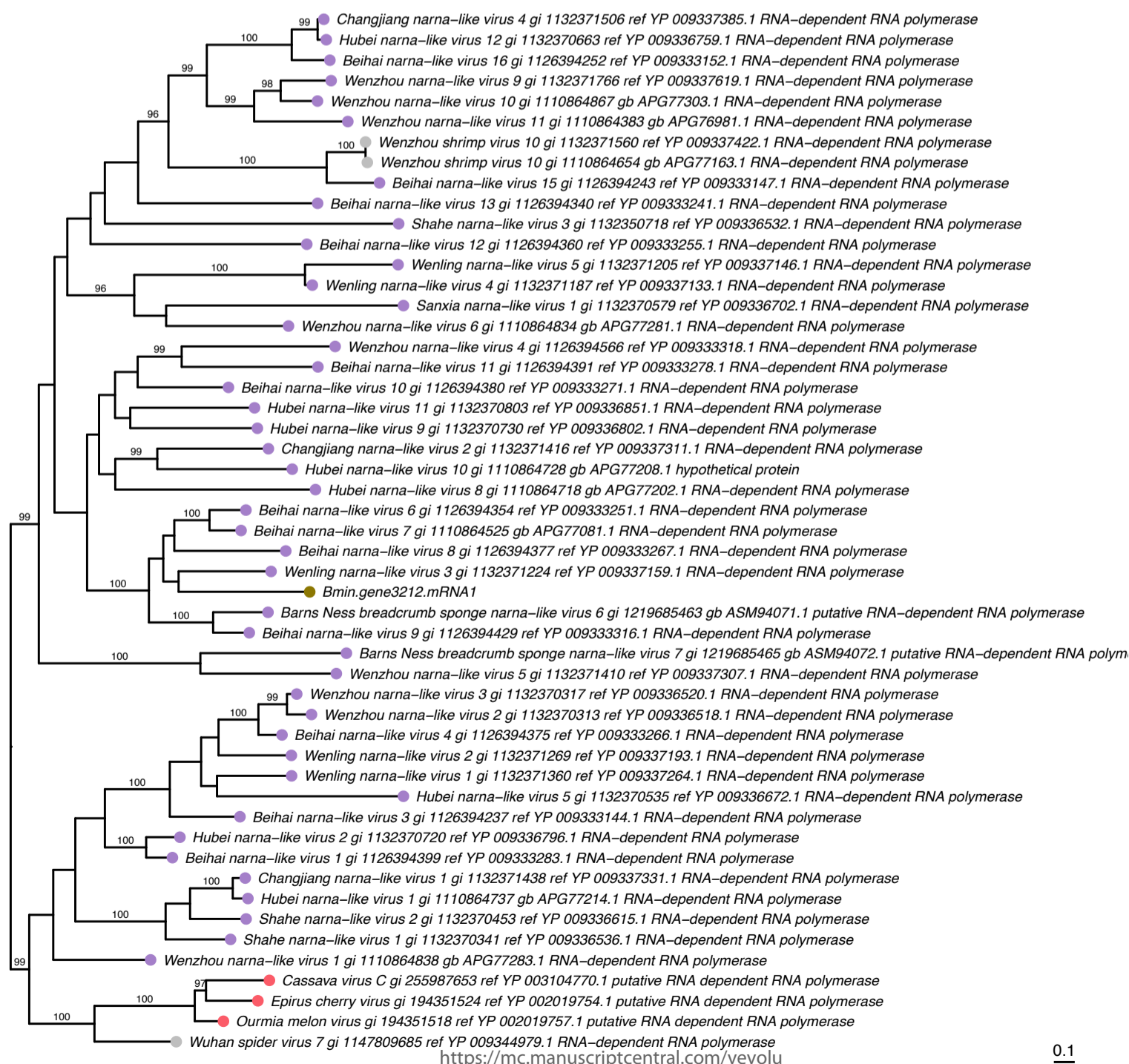


Taxa

- Symbiodiniaceae
- Symbiodiniaceae RNA virus
- Weivirus-like virus

0.2

RdRp-like group 2 (Phylogeny D)

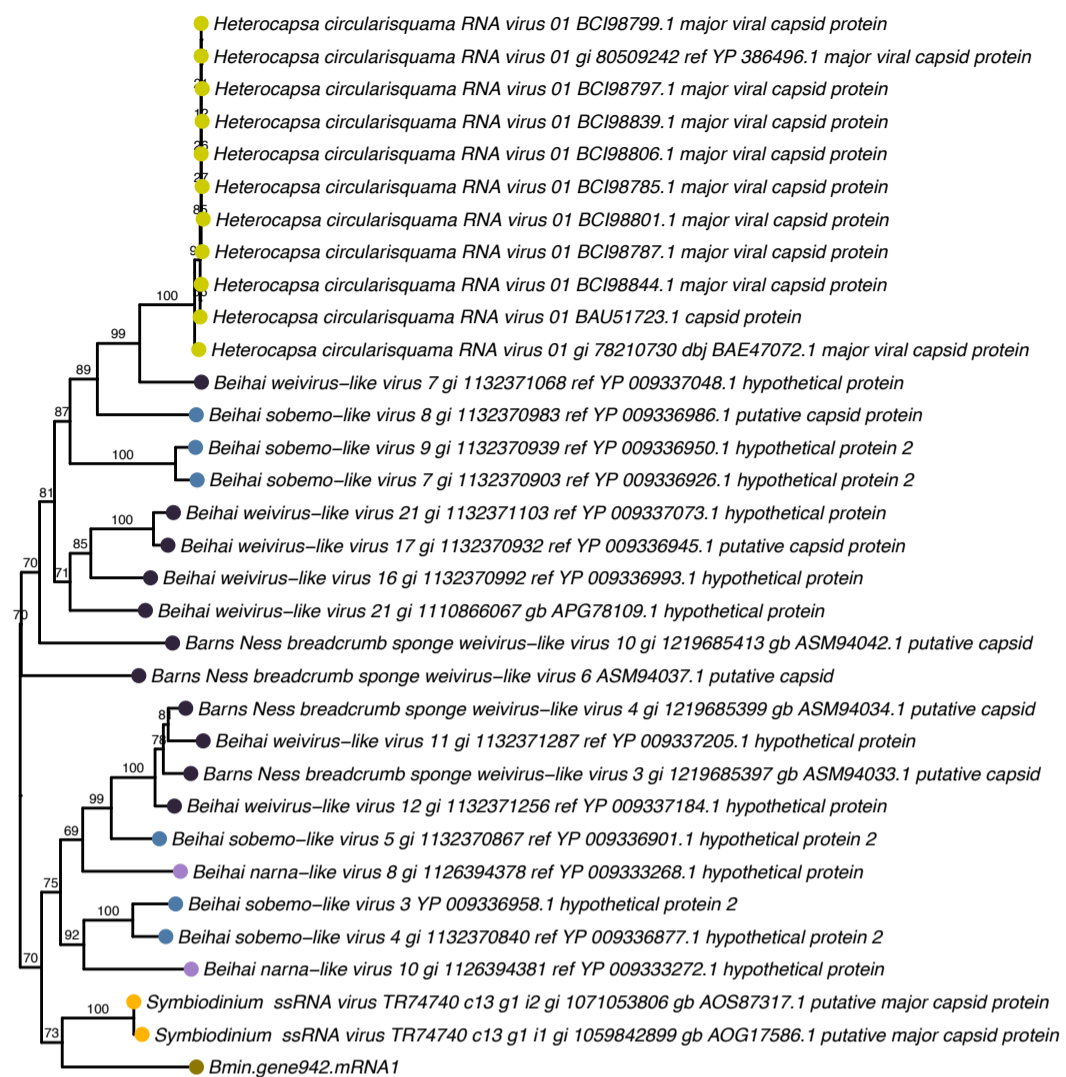


Taxa

- Narna-like virus
- Ourmiavirus
- Symbiodiniaceae
- Unclassified viruses

0.1

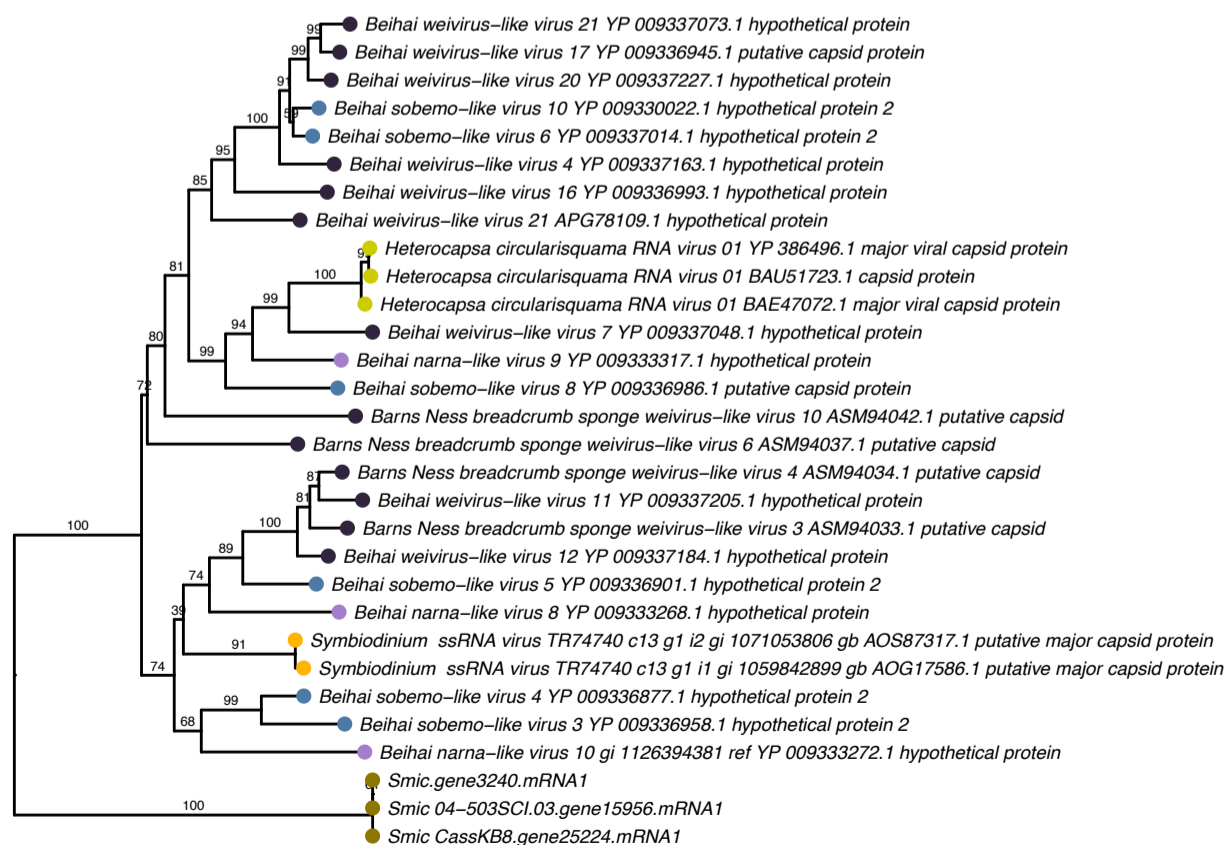
MCP-like (Phylogeny E)



Taxa

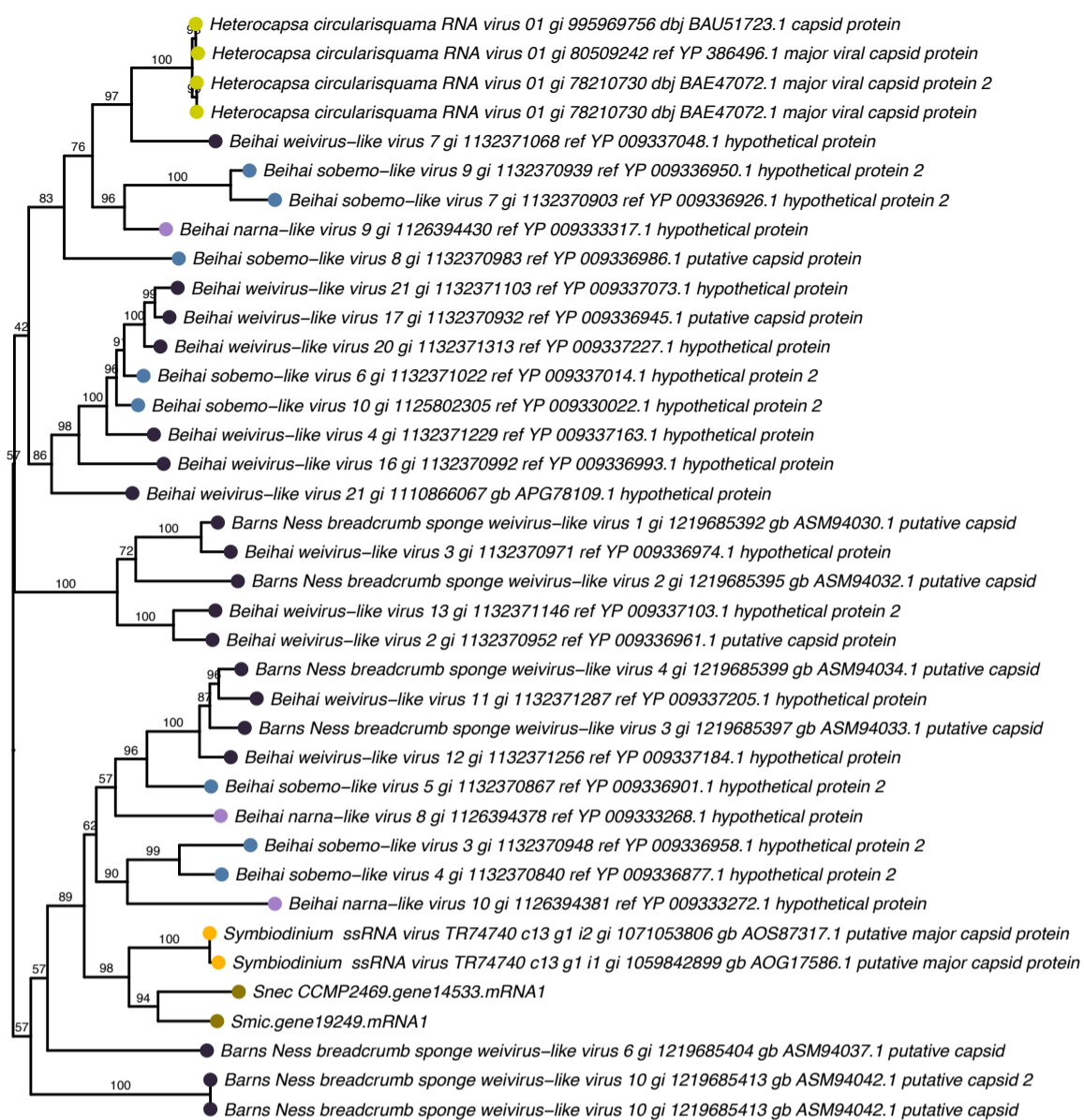
- Dinornavirus
- Narna-like virus
- Sobemo-like virus
- Symbiodiniaceae
- Symbiodiniaceae RNA virus
- Weivirus-like virus

MCP-like (Phylogeny F)



0.1

MCP-like (Phylogeny G)

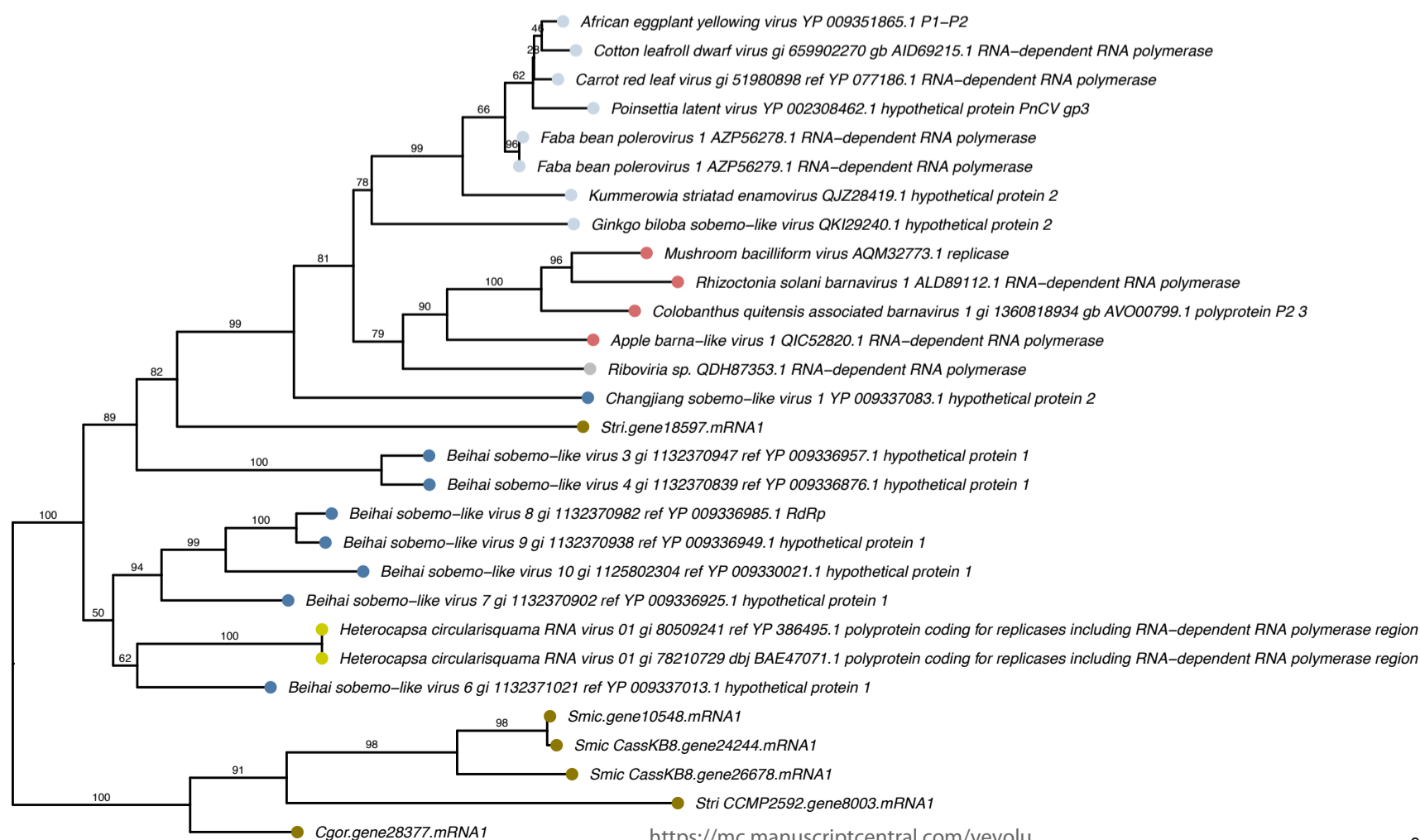


Taxa

- Dinornavirus
- Narna-like virus
- Sobemo-like virus
- Symbiodiniaceae
- Symbiodiniaceae RNA virus
- Weivirus-like virus

0.1

Polyprotein replicases - RdRp (Phylogeny H)



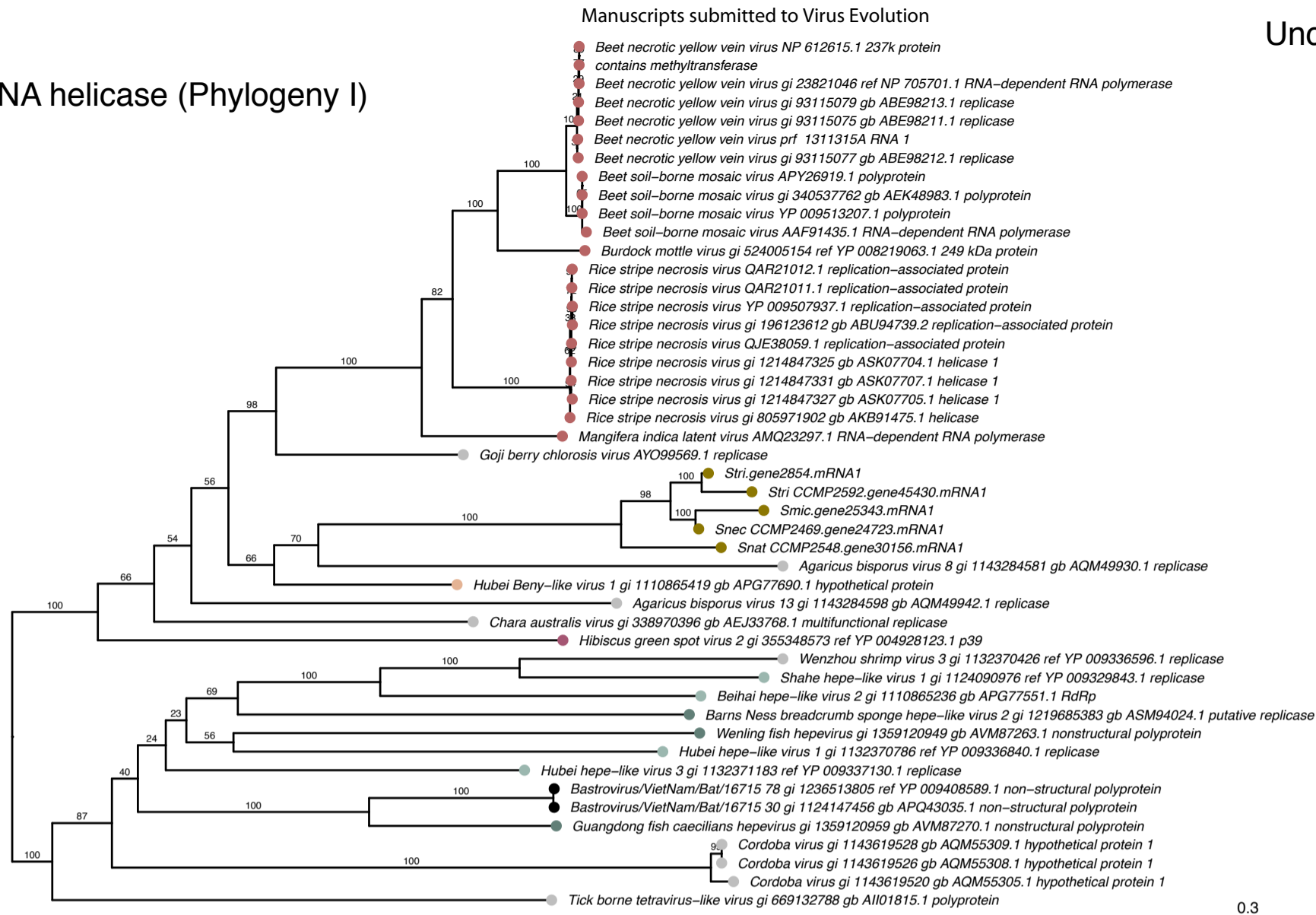
Taxa

- Barnavirus
- Dinornavirus
- Sobemo-like virus
- Solemoviridae
- Symbiodiniaceae
- Unclassified viruses

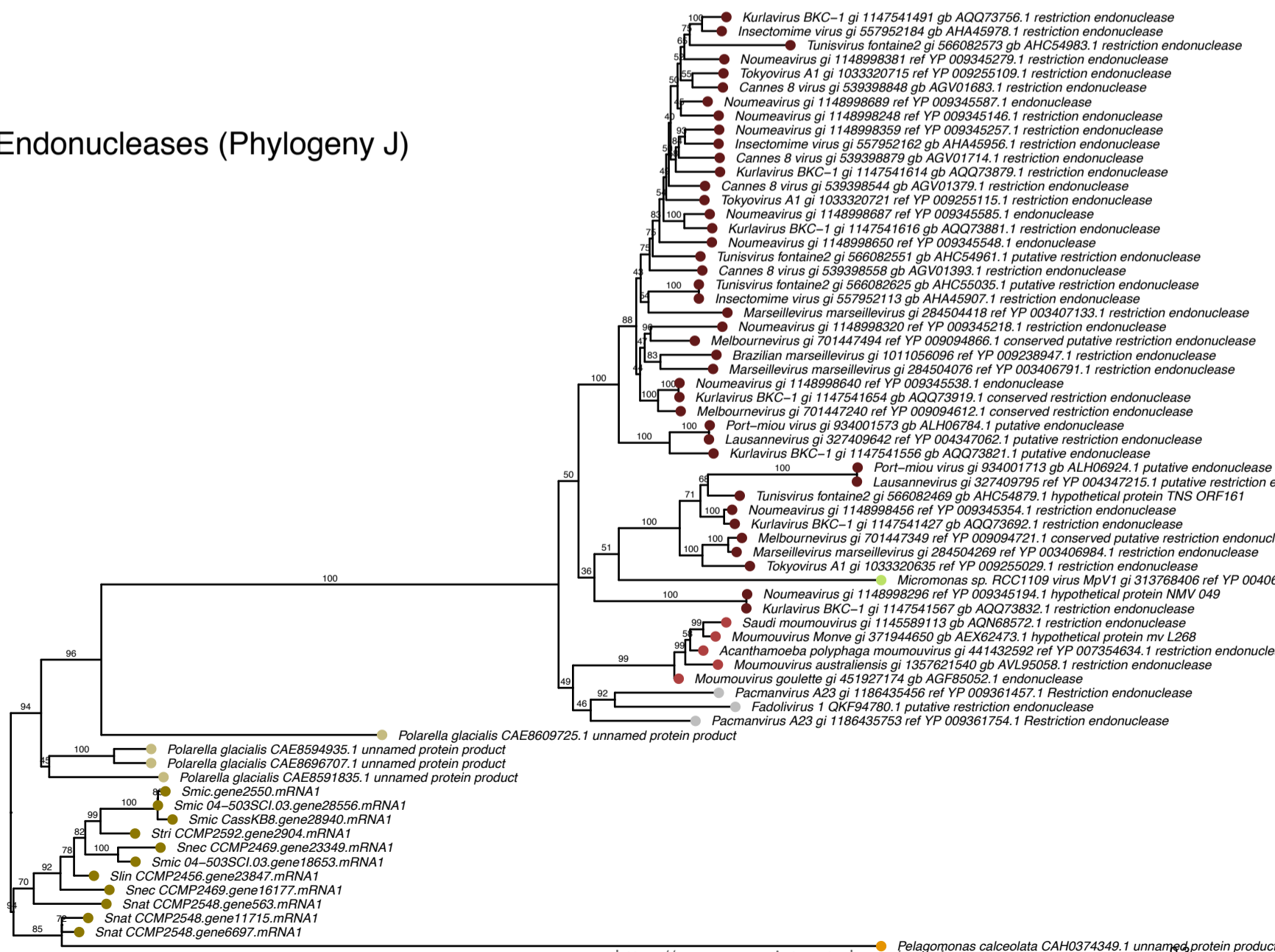
<https://mc.manuscriptcentral.com/vevolu>

0.2

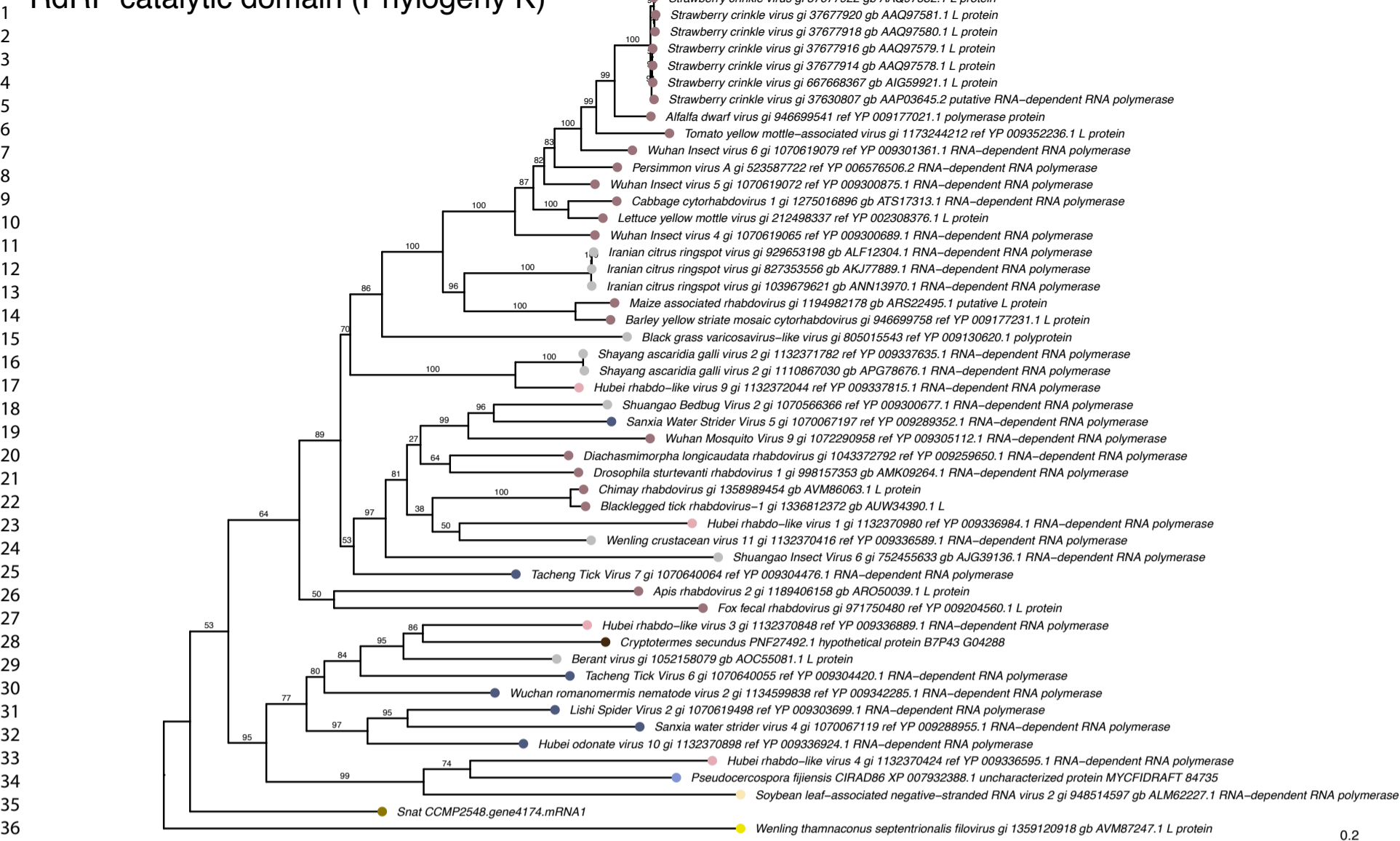
Viral RNA helicase (Phylogeny I)



Endonucleases (Phylogeny J)

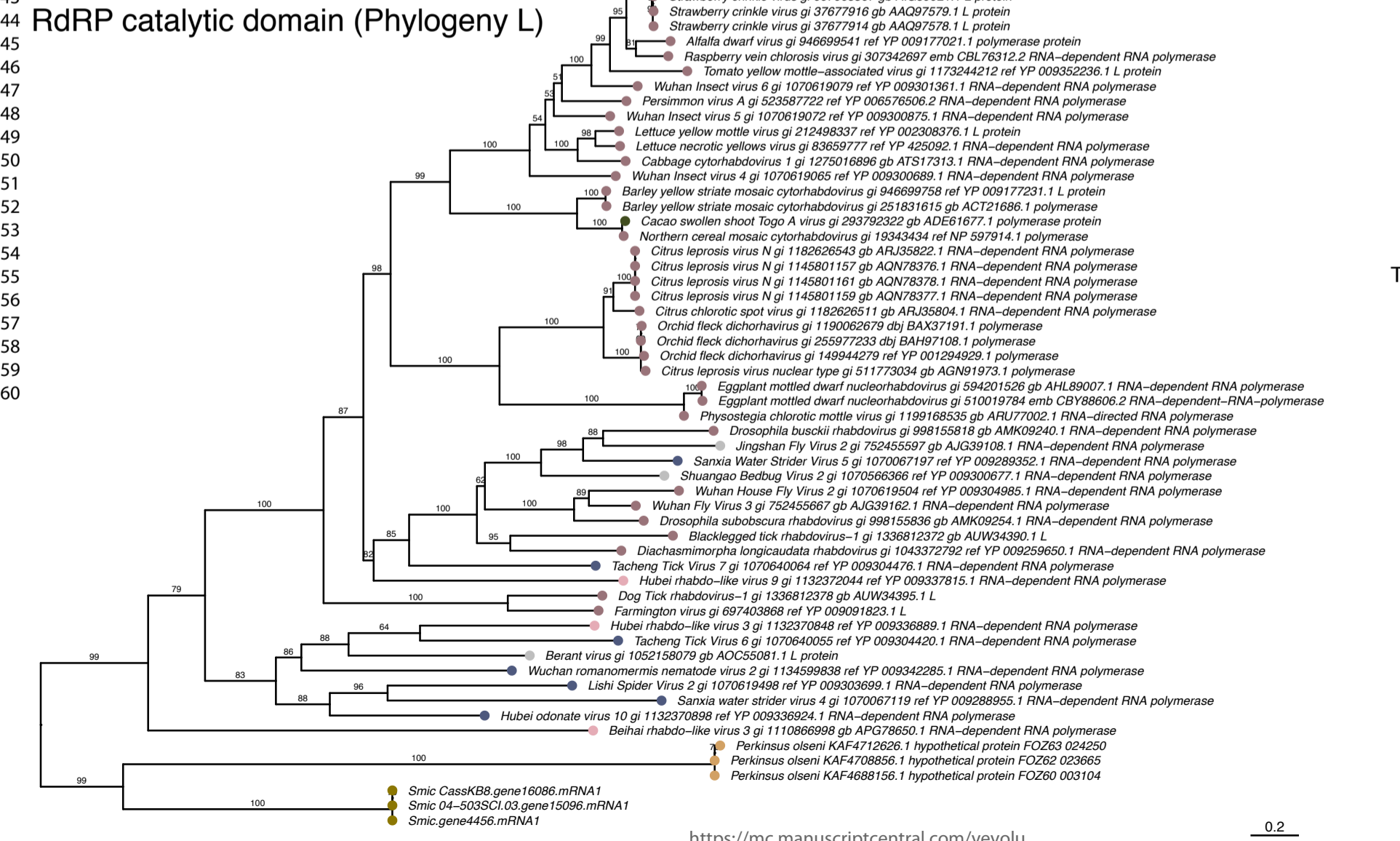


RdRP catalytic domain (Phylogeny K)



- Taxa**
- Ascomycota
 - Filoviridae
 - Insecta
 - Lispiviridae
 - Mymonaviridae
 - Rhabdo-like virus
 - Rhabdoviridae
 - Symbiodiniaceae
 - Unclassified viruses

RdRP catalytic domain (Phylogeny L)



- Taxa**
- Caulimoviridae
 - Lispiviridae
 - Perkinsidae
 - Rhabdo-like virus
 - Rhabdoviridae
 - Symbiodiniaceae
 - Unclassified viruses

Genome abbreviation	Species (isolate) name	Protein numbers	pre-vHGTs	vHGTs-2022	% vHGTs-	Host species	Host type
Bmin	Breviolum minutum	32803	30	2	0.006	<i>Orbicella faveolata</i> (former	Coral (Hexacorallia)
CC92	Cladocopium sp. C92	33421	23	0	0.000	<i>Fragum sp.</i>	Clam (Bivalvia)
Cgor	Cladocopium goreau	39006	24	1	0.003	<i>Acropora tenuis</i>	Coral (Hexacorallia)
Fkaw	Fugacium kawagutii	31520	17	0	0.000	Isolated from <i>Montipora verrucosa</i>	Unknown
Slin_CCMP2456	Symbiodinium linucheae CCMP2456	32053	31	1	0.003	<i>Plexaura homamalla</i>	Octocoral
Smic	Symbiodinium microadriaticum	29728	28	14	0.047	<i>Stylophora pistillata</i>	Coral (Hexacorallia)
Smic_04-503SCI.03	Symbiodinium microadriaticum 04-	38462	39	9	0.023	Unknown	Unknown
Smic_CassKB8	Symbiodinium microadriaticum CassKB8	42652	40	11	0.026	<i>Cassiopea sp.</i>	Jellyfish (Scyphozoa)
Snat_CCMP2548	Symbiodinium natans CCMP2548	35270	29	5	0.014	Free-living	Free-living
Snec_CCMP2469	Symbiodinium necroappetens CCMP2469	35672	33	8	0.022	<i>Condylactis gigantea</i>	Anemone
Spil_CCMP2461	Symbiodinium pilosum CCMP2461	23437	20	0	0.000	<i>Zoanthus sociatus</i>	Zoanthids
Stri	Symbiodinium tridacnidorum	25808	31	2	0.008	<i>Tridacna crocea</i>	Clam (Bivalvia)
Stri_CCMP2592	Symbiodinium tridacnidorum CCMP2592	45474	49	6	0.013	<i>Heliofungia actiniformis</i>	Coral (Hexacorallia)
Total count (sum)	13	445306	394	59	0.166		

For Review Only

Genome	Protein	WH17a	Genome CDB length	Genome CDB length	WH17a-2022 length	WH17a-2022 length	Genome CDB X WH17a-2022 mean Length (μ)	Genome (CDB) GC	Genome (CDB) GC	WH17a-2022 GC	WH17a-2022 GC	Genome CDB X WH17a-2022 mean GC (μ)
Btlnh	32803	2	1829.794427	10.82024303	1888.5	196.5	0.868167902396248	51.37612503	0.0153835697082351	51.29071211	0.699695563	0.922673331
Cpdr	39006	1	1625.915806	8.497688799	444	NA	NA	54.76078942	0.0194836725897129	51.57657658	NA	NA
SHL_OCMF2406	32053	1	1560.219	8.208636	1581	NA	NA	58.12416	0.01960339	54.3327	NA	NA
Btmo	29728	14	2217.576732	12.72658751	2112.857143	354.520447	0.772465892628945	57.48950778	0.0170804400100531	54.21847764	0.512893973471282	2.4160789793009e-05
Btlnh_OA	38462	9	1814.093231	9.5547768	2250.111111	470.694469	0.381442376630971	57.16127731	0.0204591987085494	54.4877883	0.458226877344755	0.000391896952288187
Btlnh_Cnem1Q8	42652	11	1836.738139	9.117236782	2441.909091	380.6169502	0.142992653950441	59.05804217	0.0232012009002331	56.13108859	0.291590139994192	1.43047047341923e-06
SHL_OCMF2546	35270	5	1660.039529	8.674021033	2316.6	717.1489803	0.411737537199503	58.34446068	0.0190850079742343	60.12620774	0.727413048	0.0704593703673462
Btmo_CCMF2406	35672	8	1485.112719	8.052933748	1443.875	255.9522443	0.87860646521478	57.704944	0.0189547506501442	51.89649044	0.649232722	4.4008241308554e-05
Btlnh	25808	2	1424.609897	8.244818801	4185	2702	0.495508406182338	57.84946773	0.0206391872762781	55.91560915	0.680426332649998	0.215103686496037
SHL_OCMF2582	45474	6	2033.565101	9.553814759	1539.166667	357.6453939	0.225462206553619	57.44785905	0.015095646672105	54.30502807	1.374054837	0.0738923795476341

For Review Only

All genomes CDS length	All genome CDS length (stderr)	All virseqs CDS length (mean)	All virseqs CDS length (stderr)	All genomes CDS X All virseqs CDS (p-value)
1774.998182	2.937872618	2088.016949	175.9260984	8.0E-02
All genomes GC (mean)	All CDS genomes GC (stderr)	All virseqs GC (mean)	All CDS virseqs GC (stderr)	All CDS genomesX vHGTs-2022 (p-value)
56.7864614	0.00679069403987519	54.73083571	0.357844201006567	3.6E-07
All genomes CDS length	All genome CDS length (stderr)	All virseqs RNAPOL CDS length	All virseqs RNAPOL CDS length	All genomes CDS X All virseqs RNAPOL CDS (p-
1774.998182	2.937872618			
All genomes GC (mean)	All CDS genomes GC (stderr)			
56.7864614	0.00679069403987519			

For Review Only

Sequence_id	GC%		Length	
Symb_vir_RNA_pol_KX538960.1	43.95		2280.00	
Symb_vir_MCP_KX538960.1	48.10		1077.00	
HcRNAV34_RNA-pol_AB218608.1:19-3018	54.03		3000.00	
HcRNAV34_MCP_AB218608.1:3182-4261	58.06		1080.00	
	GC% (mean)	stderror GC	Length (mean)	stderror Length
Symbio_vir_genome_RNAPOL-like	53.76	0.46	2374.19	249.66
Symbio_vir_genome_MCP-like	52.95	1.06	1072.00	230.53
Symbio_background_genome_all_CDS 2022 (only same genomes as RNAPOL-like	57.80	0.01	1872.46	4.39
Symbio_background_genome_all_CDS 2022 (only same genomes as MCP-like virseqs)	56.7159937	0.010867106268647	1842.160157	4.487962579

For Review Only

Genome	All CDSs						Viral CDSs						
	Total no. of genes	No. genes with introns	No. introns	Average no. of introns	Total intron length (bp)	Average intron length (bp)	Total no. viral genes	No. genes with introns	No. genes without introns (single-exon genes)	No. introns	Average no. of introns	Total intron length (bp)	Average intron length (bp)
<i>Breviolum minutum</i>	32803	30749	591597	18.03484437	267003468	451.327	2	2	0	26	13	23459	902.269
<i>Cladocopium goreauli</i>	39006	37444	447075	11.4616962	265351575	593.528	1	1	0	1	1	94	94
<i>Cladocopium</i> sp. C92	33421	32632	589620	17.64220101	367817944	623.822	0	0	0	NA		0	0
<i>Fugacium kawagutii</i>	31520	30259	335019	10.62877538	220578588	658.406	1	1	0	1	1	774	774
<i>Symbiodinium linucheae</i> CCMP2456	32053	29498	513645	16.02486507	234193267	455.944	2	1	1	8	4	4183	522.875
<i>Symbiodinium microadriaticum</i> O4-	38462	35728	609648	15.85065779	246520488	404.365	9	7	2	76	8.444444444	55698	732.868
<i>Symbiodinium microadriaticum</i> CasaKB8	42652	39578	629550	14.78015193	259462674	412.14	11	11	0	169	15.36363636	128821	762.254
<i>Symbiodinium natans</i> CCMP2548	35270	30171	517113	14.66155373	251116222	485.612	5	5	0	85	17	61823	727.329
<i>Symbiodinium neoappetens</i>	35672	32485	463202	12.98503028	120122982	259.332	8	8	0	19	2.375	3512	184.842
<i>Symbiodinium pilosum</i> CCMP2461	23437	22006	275531	11.75624013	129873189	471.356	0	0	0	0	NA	0	0
<i>Symbiodinium tridacnidorum</i> CCMP2592	45474	40282	689089	15.15347231	391733376	568.48	6	5	1	38	6.333333333	41973	1104.55
<i>Symbiodinium microadriaticum</i>	29728	28449	541336	18.20963402	209996798	387.923	14	14	0	139	9.928571429	124767	897.604
<i>Symbiodinium tridacnidorum</i>	25808	22775	307547	11.91673125	130134686	423.138	2	2	0	5	2.5	388	77.6
Total		412056	6509972		3093805277	475.2563109	61	57	4	567		445492	785.7001764

For Review Only

gene_scaffold	gene_start	gene_stop	gene_ID	gene_score	gene_strand	DinoSL_scaffold	DinoSL_start	DinoSL_stop	DinoSL_info	DinoSL_score	DinoSL_strand	distance_between_gene_and_DinoSL
Smic_04-503SCL03.scaffold3501	13545	15709	Smic_04-503SCL03.gene15955.mRNA1	0	+	Smic_04-503SCL03.scaffold3501	12942	12960	DinoSL100.0000-3-	0	+	-586
Smic_04-503SCL03.scaffold3509	146636	150480	Smic_04-503SCL03.gene4101.mRNA1	0	+	Smic_04-503SCL03.scaffold3509	148112	148132	DinoSL100.0000-1-	0	+	0
Smic_CassKB8.scaffold1834	661	25265	Smic_CassKB8.gene10153.mRNA1	0	-	Smic_CassKB8.scaffold1834	9470	9490	DinoSL100.0000-1-	0	-	0
Smic_CassKB8.scaffold8947	7360	10122	Smic_CassKB8.gene21161.mRNA1	0	+	Smic_CassKB8.scaffold8947	8009	8027	DinoSL100.0000-3-	0	+	0
Smic_CassKB8.scaffold7092	1465	27295	Smic_CassKB8.gene21497.mRNA1	0	+	Smic_CassKB8.scaffold7092	10508	10528	DinoSL100.0000-1-	0	+	0
Smat_CCMP2548.scaffold1176	20901	21539	Smat_CCMP2548.gene30156.mRNA1	0	-	Smat_CCMP2548.scaffold1176	22801	22817	DinoSL100.0000-5-	0	-	-1263
Smec_CCMP2469.scaffold2078	966	2617	Smec_CCMP2469.gene1742.mRNA1	0	+	Smec_CCMP2469.scaffold2078	620	636	DinoSL100.0000-5-	0	+	-331
Smec_CCMP2469.scaffold5477	7815	8640	Smec_CCMP2469.gene12555.mRNA1	0	-	Smec_CCMP2469.scaffold5477	9624	9644	DinoSL100.0000-1-	0	-	-685
Smec_CCMP2469.scaffold7240	9106	11417	Smec_CCMP2469.gene14533.mRNA1	0	-	Smec_CCMP2469.scaffold7240	11771	11791	DinoSL95.238-1-2 21	0	-	-355
Smic.scaffold956	1019745	1047674	Smic.gene2775.mRNA1	0	+	Smic.scaffold956	1027828	1027848	DinoSL100.0000-1-	0	+	0
Smic.scaffold987	1209704	1211868	Smic.gene3241.mRNA1	0	-	Smic.scaffold987	1212454	1212472	DinoSL100.0000-3-	0	-	-587
Smic.scaffold991	455265	475826	Smic.gene4069.mRNA1	0	-	Smic.scaffold991	464065	464085	DinoSL100.0000-1-	0	-	0
Sti.scaffold302	138207	145121	Sti.gene2854.mRNA1	0	-	Sti.scaffold302	145242	145262	DinoSL100.0000-1-	0	-	-122

For Review Only

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



For Review Only

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

The table is a complex grid of data, likely representing a phylogenetic tree or a list of sequences. It consists of many columns and rows, with some cells containing bolded text or specific identifiers. The overall structure is highly detailed and spans most of the page's height.

For Review Only

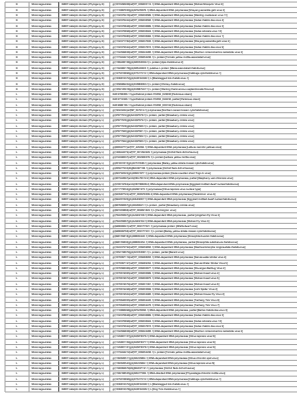
Symbio_seqs	Interpro_3.25					
Smic_D4-	SUPERFAMILY	SSF51197	Clavaminate synthase-like	111	286	4E-17
	SUPERFAMILY	SSF56672	DNA/RNA polymerases	636	828	9E-10
Smic_CassKB8.gene21497.mRNA1	SUPERFAMILY	SSF56672	DNA/RNA polymerases	573	762	8E-06
	SUPERFAMILY	SSF46785	Winged helix DNA-binding	271	330	5E-08
Smic_CassKB8.gene9217.mRNA1	SUPERFAMILY	SSF56672	DNA/RNA polymerases	668	860	9E-10
	SUPERFAMILY	SSF51197	Clavaminate synthase-like	143	318	2E-17
Smic.gene1589.mRNA1	SUPERFAMILY	SSF51197	Clavaminate synthase-like	48	218	6E-16
	SUPERFAMILY	SSF56672	DNA/RNA polymerases	568	760	
Smic.gene2775.mRNA1	SUPERFAMILY	SSF46785	Winged helix DNA-binding	476	535	6E-08
	SUPERFAMILY	SSF56672	DNA/RNA polymerases	778	967	1E-05
Stri_CCMP2592.gene39232.mRNA1	SUPERFAMILY	SSF56672	DNA/RNA polymerases	590	742	2E-05
	SUPERFAMILY	SSF46785	Winged helix DNA-binding	276	331	4E-08

For Review Only

Genome abbreviation	Species (strain) name	bioproject	SOURCE	SRA	Run	OBS
Bmin	Breviolum minutum		TRANSCRIPTOMI	OIST	OIST	https://marinegenomics.oist.jp/symb/viewer/download?project_id=21
CC92	Cladocopium sp. C92	Symbiodinium sp. Clade C Y103 PRJDB3243	TRANSCRIPTOMI	DRX082293	DRR088470	Illumina HiSeq 2000
CC92	Cladocopium sp. C92	Symbiodinium sp. Clade C Y103 PRJDB3243	TRANSCRIPTOMI	DRX082292	DRR088469	Illumina HiSeq 2000
CC92	Cladocopium sp. C92	Symbiodinium sp. Clade C Y103 PRJDB3243	TRANSCRIPTOMI	DRX082291	DRR088468	Illumina HiSeq 2000
CC92	Cladocopium sp. C92	Symbiodinium sp. Clade C Y103 PRJDB3243	TRANSCRIPTOMI	DRX082290	DRR088467	Illumina HiSeq 2000
CC92	Cladocopium sp. C92	Symbiodinium sp. Clade C Y103 PRJDB3243	TRANSCRIPTOMI	DRX082289	DRR088466	Illumina HiSeq 2000
Cgor	Cladocopium goreau	PRJEB20399	TRANSCRIPTOMI	SRX1205243	SRR2298883	GSM1872186: MI_day 9_32C_sample 6; Symbiodinium sp. C1; RNA-
Cgor	Cladocopium goreau	PRJNA295075	TRANSCRIPTOMI	SRX1205242	SRR2298882	
Cgor	Cladocopium goreau	PRJNA295075	TRANSCRIPTOMI	SRX1205238	SRR2298878	
Cgor	Cladocopium goreau	PRJNA723630	TRANSCRIPTOMI	SRX1397118	SRR17809162	polyA-enriched Stranded TruSeq
Cgor	Cladocopium goreau	PRJNA723630	TRANSCRIPTOMI	SRX1397118	SRR17809163	polyA-enriched Stranded TruSeq
Fkaw	Fugacium kawaguti	PRJNA412800	TRANSCRIPTOMI	SRX3236099	SRR6123446	Not the same as the genome
Fkaw	Fugacium kawaguti	PRJNA412800	TRANSCRIPTOMI	SRX3236098	SRR6123445	Not the same as the genome
Fkaw	Fugacium kawaguti	PRJNA248394	TRANSCRIPTOMI	SRX554110	SRR1300302	
Slin_CCMP2456	Symbiodinium linucheae CCMP2456	PRJEB34894	TRANSCRIPTOMI	ERX3714737	ERR3711290	cDNA Illumina TruSeq RNA / https://www.ebi.ac.uk/ena/browser/view/ERR3711290?show=reads
Smlc	Symbiodinium microadriaticum Smic_CCMP2467	S. microadriaticum (strain CCMP2467) -	TRANSCRIPTOMI	SRX1681571	SRR3337505	cDNA Illumina HiSeq 2000
Smlc	Symbiodinium microadriaticum Smic_CCMP2467	S. microadriaticum (strain CCMP2467) -	TRANSCRIPTOMI	SRX1681570	SRR3337504	cDNA Illumina HiSeq 2000
Smlc	Symbiodinium microadriaticum Smic_CCMP2467	S. microadriaticum (strain CCMP2467) -	TRANSCRIPTOMI	SRX1681569	SRR3337503	
Smlc	Symbiodinium microadriaticum Smic_CCMP2467	S. microadriaticum (strain CCMP2467) -	TRANSCRIPTOMI	SRX1681568	SRR3337502	
Smlc	Symbiodinium microadriaticum Smic_CCMP2467	S. microadriaticum (strain CCMP2467) -	TRANSCRIPTOMI	SRX1681567	SRR3337501	
Smlc	Symbiodinium microadriaticum Smic_CCMP2467	S. microadriaticum (strain CCMP2467) -	TRANSCRIPTOMI	SRX1681566	SRR3337500	
Smlc	Symbiodinium microadriaticum Smic_CCMP2467	S. microadriaticum (strain CCMP2467) -	TRANSCRIPTOMI	SRX1681565	SRR3337498	
Smlc	Symbiodinium microadriaticum Smic_CCMP2467	S. microadriaticum (strain CCMP2467) -	TRANSCRIPTOMI	SRX1681564	SRR3337497	
Smlc	Symbiodinium microadriaticum Smic_CCMP2467	S. microadriaticum (strain CCMP2467) -	TRANSCRIPTOMI	SRX1681563	SRR3337496	
Smlc	Symbiodinium microadriaticum Smic_CCMP2467	S. microadriaticum (strain CCMP2467) -	TRANSCRIPTOMI	SRX1681562	SRR3337495	
Smlc	Symbiodinium microadriaticum Smic_CCMP2467	S. microadriaticum (strain CCMP2467) -	TRANSCRIPTOMI	SRX1681561	SRR3337494	
Smlc	Symbiodinium microadriaticum Smic_CCMP2467	S. microadriaticum (strain CCMP2467) -	TRANSCRIPTOMI	SRX1681560	SRR3337493	
Smlc_04-5038Cl.03	Symbiodinium microadriaticum 04-5038Cl.03	PRJEB34894	TRANSCRIPTOMI	ERX3714702	ERR3711255	cDNA Illumina TruSeq RNA
Smlc_Case8KB8	Symbiodinium microadriaticum Case8KB8	PRJEB34894	TRANSCRIPTOMI	ERX3714714	ERR3711267	cDNA Illumina TruSeq RNA
Smlc_Case8KB8	Symbiodinium microadriaticum Case8KB8	PRJNA80083	TRANSCRIPTOMI	SRX076696	SRR278693	454 GS FLX Titanium cDNA
Sna_CCMP2548	Symbiodinium natans CCMP2548	PRJEB34894	TRANSCRIPTOMI	ERX4439268	ERR4501197	
Stf	Symbiodinium tridacnidorum	Symbiodinium sp. Clade A Y106 - PRJDB3242	TRANSCRIPTOMI	DRX073090	DRR079246	
Stf	Symbiodinium tridacnidorum	Symbiodinium sp. Clade A Y106 - PRJDB3242	TRANSCRIPTOMI	DRX073089	DRR079245	
Stf	Symbiodinium tridacnidorum	Symbiodinium sp. Clade A Y106 - PRJDB3242	TRANSCRIPTOMI	DRX073088	DRR079244	
Stf	Symbiodinium tridacnidorum	Symbiodinium sp. Clade A Y106 - PRJDB3242	TRANSCRIPTOMI	DRX073087	DRR079243	
Stf	Symbiodinium tridacnidorum	Symbiodinium sp. Clade A Y106 - PRJDB3242	TRANSCRIPTOMI	DRX073086	DRR079242	
Stf_CCMP2592	Symbiodinium tridacnidorum CCMP2592	PRJEB34894	TRANSCRIPTOMI	ERX4439230	ERR4501159	

Review Only

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



For Review Only