

# A vast world of viroid-like circular RNAs revealed by mining metatranscriptomes

Benjamin D. Lee<sup>1,2</sup>      Uri Neri<sup>3</sup>      Simon Roux<sup>4</sup>  
Yuri I. Wolf<sup>1</sup>      Antonio Pedro Camargo<sup>4</sup>  
Mart Krupovic<sup>5</sup>      RNA Virus Discovery Consortium  
Peter Simmonds<sup>2</sup>      Nikos Kyrpides<sup>4</sup>      Uri Gophna<sup>3</sup>  
Valerian V. Dolja<sup>6</sup>      Eugene V. Koonin<sup>1,\*</sup>

<sup>1</sup> National Center for Biotechnology Information, National Library of  
Medicine, National Institutes of Health, Bethesda, MD 20894, USA

<sup>2</sup> Nuffield Department of Medicine, University of Oxford, Oxford OX3 7BN,  
UK

<sup>3</sup> The Shmunis School of Biomedicine and Cancer Research, Tel Aviv  
University, Tel Aviv 6997801, Israel

<sup>4</sup> Department of Energy Joint Genome Institute, Lawrence Berkeley Na-  
tional Laboratory, Berkeley, CA 94720, USA

<sup>5</sup> Institut Pasteur, Université de Paris, CNRS UMR6047, Archaeal Virology

17 Unit, 75015 Paris, France

18 <sup>6</sup> Department of Botany and Plant Pathology, Oregon State University,  
19 Corvallis, OR 97331, USA

20 \* Correspondence: Eugene V. Koonin <koonin@ncbi.nlm.nih.gov>

## 21 Summary

22 Viroids and viroid-like agents are unique, minimal RNA replicators that typ-  
23 ically encode no proteins and hijack cellular enzymes for their genome repli-  
24 cation. As the extent and diversity of viroid-like agents are poorly under-  
25 stood, we developed a computational pipeline to identify viroid-like cova-  
26 lently closed circular (ccc) RNAs and applied it to 5,131 global metatran-  
27 scriptomes and 1,344 plant transcriptomes. The search resulted in 11,420  
28 viroid-like, ribozyme-containing cccRNAs spanning 4,409 species-level clus-  
29 ters, which is a five-fold increase compared to the previously known set of  
30 viroids and viroid-like RNA agents. Within this diverse collection, we identi-  
31 fied numerous putative novel viroids, satellite RNAs, retrozymes, and ribozy-  
32 like viruses. We also found previously unknown ribozyme combinations and  
33 unusual ribozymes within the cccRNAs. Self-cleaving ribozymes were identi-  
34 fied in both RNA strands of ambiviruses and some mito-like viruses as well  
35 as in capsid-encoding satellite virus-like cccRNAs. The broad presence of  
36 viroid-like cccRNAs in diverse transcriptomes and ecosystems implies that  
37 their host range is not limited to plants, and matches between viroid-like

cccRNAs and CRISPR spacers suggest that some of them might replicate in prokaryotes.

## Introduction

Viroids, which cause several economically important diseases in agricultural plants, are the smallest and simplest among the known infectious agents (Daròs et al., 2006; Diener, 2001; Mascia and Gallitelli, 2017). Viroids are small, covalently closed circular (ccc) RNA molecules of 220 to 450 nucleotides that encode no proteins and consist largely of RNA structures that are required for replication or viroid-host interaction. In contrast to viruses, which hijack the host translation system to produce proteins encoded in virus genes, viroids take advantage of the host transcriptional machinery. Specifically, viroids hijack the host plant's DNA-dependent RNA polymerase II to transcribe their RNA and thus catalyze viroid replication (Mühlbach and Sängler, 1979; Navarro et al., 2000; Schindler and Mühlbach, 1992). Viroids utilize the rolling circle replication (RCR) mechanism, producing multimeric intermediates that are cleaved into genome-size monomers by ribozymes that are present in both the plus and the minus strands of viroid RNA or by recruited host RNases (Branch and Robertson, 1984; Flores et al., 2017). The resulting linear monomers are then ligated by a host DNA ligase to form the mature cccRNA (Nohales et al., 2012a, 2012b).

Since the discovery of viroids in 1971 (Diener, 1971), about fifty distinct

viroid species have been identified in plants and classified into two families, *Avsunviroidae* and *Pospiviroidae*. Members of the *Avsunviroidae* use viroid-encoded autocatalytic hammerhead (HHR) ribozymes, a defining feature of this family, to process replication intermediates into unit length viroid genomes (Di Serio et al., 2017; Wang, 2021). Members of the family *Pospiviroidae* lack ribozymes and instead rely on conserved sequence motifs that serve as recognition and cleavage sites for host RNase III (Branch et al., 1988). The members of the two viroid families also adopt distinct RNA structures: a branched RNA conformation is predominant in avsunviroids (Giguère et al., 2014a), in contrast to the typically rod-shaped conformation of pospiviroids (Giguère et al., 2014b).

In addition to viroids, several other groups of infectious agents also possess genomes consisting of cccRNA (de la Peña et al., 2020). Many plant viruses support the replication of small (about 300 nt) circular satellite RNAs (sometimes called virusoids but hereafter, satRNAs) that closely resemble viroids (Navarro et al., 2017) and also replicate via the rolling circle mechanism (Bruening et al., 1991). The satRNAs differ from viroids in that they are replicated by the RNA-dependent RNA polymerase (RdRP) of the helper virus and are encapsidated in that virus's capsid (Huang et al., 2017; Rao and Kalantidis, 2015). Thus, these satRNAs are effectively encapsidated viroids. Unlike viroids, satRNAs encode both HHRs and hairpin ribozymes (HPR), a distinct ribozyme variety (Ferré-D'Amaré and Scott, 2010).

81 Another notable viroid-like agent is the so-called retroviroid, carnation small  
82 viroid-like RNA (CarSV) which, unlike viroids, does not appear to trans-  
83 mit horizontally among plants (Daròs and Flores, 1995). CarSV, the only  
84 currently known retroviroid, is a cccRNA that is similar to viroids in size  
85 and contains HHRs in both strands. However, in contrast to the viroids, an  
86 extrachromosomal DNA copy of CarSV has been discovered and shown to  
87 integrate into the plant genome with the help of a pararetrovirus (Hegedus  
88 et al., 2004; Vera et al., 2000).

89 A recently discovered distinct group of cccRNA agents are retrozymes, retro-  
90 transposons that propagate via circular RNA intermediates of about 170  
91 to 400 nucleotides. The retrozymes are viroid-like in that they do not en-  
92 code any proteins but contain self-cleaving HHRs (Cervera et al., 2016; de  
93 la Peña and Cervera, 2017). However, unlike viroids, the retrozymes are  
94 neither infectious nor autonomous, but rather, hijack the replication machin-  
95 ery of autonomous retrotransposons. Resembling satRNAs and avsunviroids,  
96 retrozyme cccRNAs also adopt a branched conformation.

97 A distinct group of viroid-like agents is the viral realm *Ribozyviria* (Hepojoki  
98 et al., 2020) that includes deltaviruses, such as hepatitis delta virus (HDV),  
99 an important human pathogen. Similarly to pospiviroids, ribozyviruses pos-  
100 sess rod-shaped cccRNA genomes that replicate via the rolling circle mecha-  
101 nism and encode distinct ribozymes, unrelated to those of viroids, that auto-  
102 catalytically process multimeric replication intermediates (Kos et al., 1986;

103 Modahl et al., 2000; Sureau and Negro, 2016). Ribozviruses have substan-  
 104 tially larger genomes than viroids (about 1.7 kb), encode their own nucleo-  
 105 capsid protein, and rely for reproduction on a helper virus (hepatitis B virus  
 106 in the case of HDV), which provides the envelope protein for ribozvirus viri-  
 107 ons. For years, HDV remained the only known deltavirus. Recently, however,  
 108 viruses more distantly related to HDV have been discovered in diverse ver-  
 109 tebrates and invertebrates including rodents (Paraskevopoulou et al., 2020),  
 110 bats (Bergner et al., 2021), snakes (Hetzel et al., 2019), birds (Wille et al.,  
 111 2018), fish (Chang et al., 2019) and termites (Chang et al., 2019), suggesting  
 112 a considerable uncharacterized diversity of ribozviruses.

113 Viroids and viroid-like cccRNAs comprise a fundamentally distinct type of  
 114 minimal replicators, or ultimate parasites, that lack genes and effectively con-  
 115 sist only of RNA structures required for replication. This extreme simplicity  
 116 of viroids triggered speculation on a potential direct descent of viroids  
 117 from primordial RNA replicators (Diener, 2016; Flores et al., 2022). However,  
 118 the apparent narrow host range of viroids, which so far have been reported  
 119 only in plants, does not appear to be easily compatible with such an evolu-  
 120 tionary scenario. Instead, given the major similarities between retrozymes  
 121 and avsunviroids, it has been suggested that avsunviroids descended from  
 122 retrozymes (de la Peña and Cervera, 2017; Lee and Koonin, 2022).

123 Given the ultimate structural simplicity of viroids and related cccRNAs and  
 124 the universality of the DNA-dependent RNA polymerases involved in their

125 replication across life forms, the current narrow spread and limited diversity  
 126 of parasitic cccRNAs appear puzzling. Furthermore, this apparent paucity  
 127 of viroid-like agents is in stark contrast to the burgeoning diversity of RNA  
 128 viruses, many thousands of which including numerous, distinct new groups  
 129 have been discovered by metatranscriptome analyses (Edgar et al., 2022; Neri  
 130 et al., 2022; Wolf et al., 2020; Zayed et al., 2022). At present, there are at  
 131 least three orders of magnitude more known RNA viruses than there are  
 132 viroids and viroid-like cccRNAs.

133 We were interested in investigating the global diversity of viroids and viroid-  
 134 like agents. To this end, we performed an exhaustive search for cccRNAs  
 135 in a collection of about 5,131 diverse metatranscriptomes that have been  
 136 recently employed for massive RNA virus discovery (Neri et al., 2022), and  
 137 additionally searched 1,341 plant transcriptomes (One Thousand Plant Tran-  
 138 scriptomes Initiative, 2019). This search yielded 10,183,455 putative cccR-  
 139 NAs within the size range of 100–10,000 nt, of which 11,378 were classified  
 140 as viroid-like based on the presence of predicted self-cleaving ribozymes or  
 141 direct similarity to reference sequences. This set of cccRNAs represents an  
 142 about five-fold increase of the known diversity of viroid-like agents. Further  
 143 analysis of these cccRNAs led to the identification of numerous putative novel  
 144 viroids, satRNAs, retrozymes and ribozy-like viruses.

# Results

## Computational approach for the discovery of viroid-like cccRNAs

We first developed an integrated, scalable computational pipeline for the *de novo* discovery and analysis of viroids and viroid-like cccRNAs directly from assembled transcriptomes and metatranscriptomes (Figure 1). The pipeline starts with the reference-free and *de novo* identification of cccRNAs or RCR intermediates. This method depends on the fact both complete circular monomers and multimeric linear intermediates will be assembled containing detectable head-to-tail repeats (Qin et al., 2020). Once identified, the sequences are then cleaved to unit length and deduplicated, taking circularity into account. Starting from the set of detected cccRNAs, the pipeline performs both alignment-free and alignment-based searches. The primary approach for the identification of viroid-like agents among the cccRNAs is the prediction of self-cleaving ribozymes using RNA sequence and secondary structure covariance models (Nawrocki and Eddy, 2013). Under the premise that the diversity of ribozymes in viroid-like RNAs could be greater than so far uncovered, we curated a database of known self-cleaving ribozyme models from Rfam (Kalvari et al., 2021), of which only a minority have been described in viroids and viroid-like RNAs. We supplemented this model database with the pospiviroid RY motif (Gozmanova, 2003) to enable the search to detect potential pospiviroids, which do not contain ribozymes.



167 The pipeline also performs direct sequence similarity searches against refer-  
168 ence databases such as ViroidDB (Lee et al., 2021).

169 Ribozyme-containing cccRNA sequences were classified as symmetric or  
170 asymmetric depending on whether they contained predicted ribozymes in  
171 both or only one RNA polarity, respectively, reflecting the RCR mode  
172 these cccRNAs are likely to undergo. However, it cannot be ruled out that  
173 some apparently asymmetric cccRNAs actually contain a second ribozyme  
174 distinct from the currently known ones.

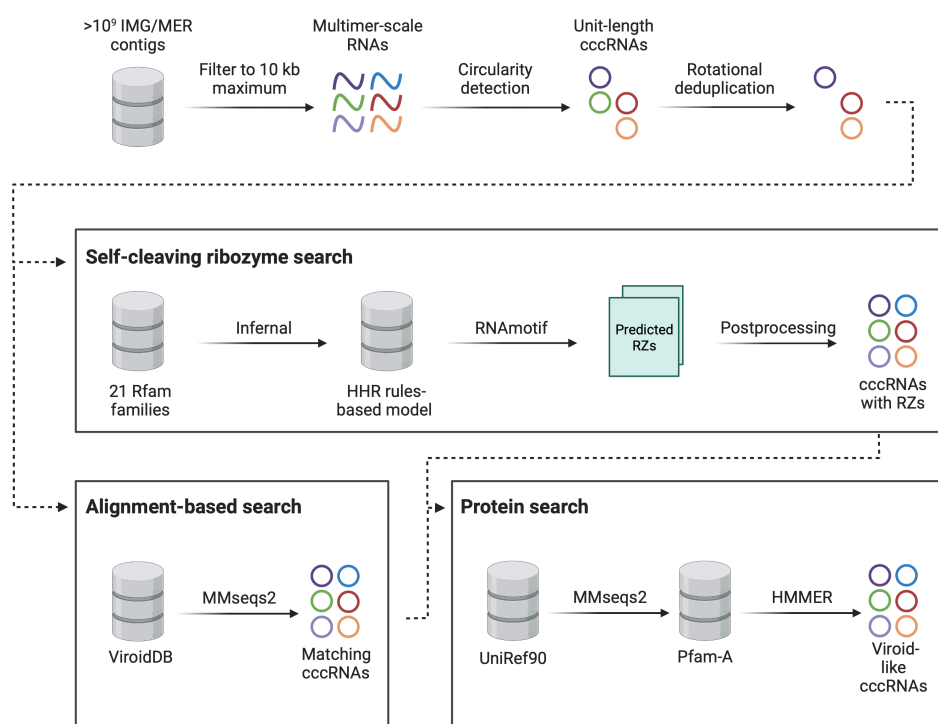


Figure 1: Viroid-like cccRNA detection pipeline

175 We validated this method by demonstrating its ability to recover known  
 176 viroid-like RNAs in both transcriptomes and metatranscriptomes (Table S1).  
 177 For the transcriptomic validation, we processed and searched the 1,000 Plant  
 178 transcriptome (1KP) data set (One Thousand Plant Transcriptomes Initia-  
 179 tive, 2019). We chose this data set due to the known presence of all type  
 180 of viroid-like cccRNAs except ribozymoviruses in plant transcriptomes. Assem-  
 181 bling the raw reads of 1,344 transcriptomes resulted in 103,139,086 contigs,  
 182 of which 163,970 were predicted to be circular. Of these putative cccRNAs,  
 183 42 were identified as viroid-like via ribozyme search (15 sequences), sequence  
 184 search against ViroidDB (33 sequences), or both (6 sequences).

185 To verify the efficacy of the detection method, we performed a direct search  
 186 of all contigs against ViroidDB and identified 12 contigs that matched a *bona*  
 187 *fide* viroid sequence with at least 50% target coverage. The detection pipeline  
 188 found four of these potentially complete viroid contigs. Of the rejected con-  
 189 tigs, four were much larger than typical viroids (>1000 nt) and contained  
 190 major ambiguous regions. The other four were low-coverage fragments that  
 191 were rejected due to being smaller than unit length and therefore unable to  
 192 be verified as circular. Iresine viroid 1, Citrus exocortis viroid, and a Coleus  
 193 blumei viroid (CbVd) were successfully retrieved. While the first two viroids  
 194 were nearly identical to the corresponding reference sequences, the CbVd-like  
 195 sequence was not. At 350 nt, this sequence differed in length from all known  
 196 coleviroids (CbVd 1-6 are 248-51, 301, 361-364, 295, 274, and 340-343 nt  
 197 long, respectively) and in the terminal conserved region, which was identical

198 to that of Dahlia latent viroid, suggesting an origin of this novel viroid by  
 199 recombination, as has been reported for other CbVd species (Hou et al., 2009;  
 200 Nie and Singh, 2017). At 85% identity, this CbVd-like sequence falls below  
 201 the species membership threshold for coleviroids (Chiumenti et al., 2021).

## 202 **A five-fold expansion of the known diversity of viroid-** 203 **like cccRNAs**

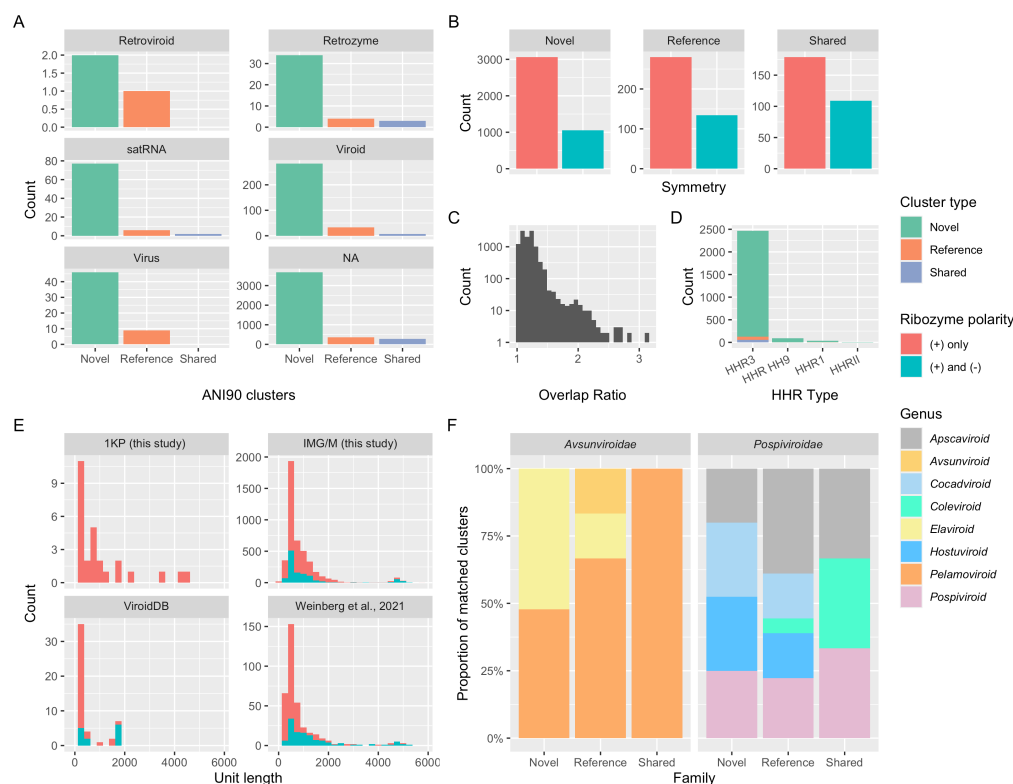
204 After testing the pipeline on plant transcriptomes, we applied it to a set  
 205 of 5,131 diverse metatranscriptomes totalling 1.5 billion metatranscriptomic  
 206 contigs (708 Gbp) after size filtration. The pipeline processed the entire set  
 207 of transcriptomes using 720 CPU hours. We identified 10,183,455 putative  
 208 cccRNAs with a median contig length of 269 nt. After removing overlapping  
 209 regions and eliminating rotationally identical sequences, the median length  
 210 of the 8,748,001 resultant monomers was 165 nt. Of these, 2,791,251 were  
 211 within the known size range of viroids (200–400 nt), including 11,378 we  
 212 classified as viroid-like because they contained a confidently predicted self-  
 213 cleaving ribozyme in at least one RNA polarity. No metatranscriptomic cc-  
 214 cRNAs matched the pospiviroid RY motif. Among the viroid-like cccRNAs,  
 215 10,181 were detected entirely by alignment-free methods, that is, showed  
 216 no detectable direct sequence similarity (including within the ribozyme re-  
 217 gions) to known viroids. The remaining 1,197 sequences shared significant  
 218 nucleotide sequence similarity with known viroid-like RNAs spanning the en-

219 tire gamut from viroids to satRNAs to retrozymes (Table 1). 907 sequences  
 220 were identified as viroid-like by both the ribozyme detection and alignment-  
 221 based search approaches. Among the 10,181 ribozyme-containing viroid-like  
 222 cccRNAs, 3,434 were symmetric, that is, contained predicted ribozymes in  
 223 both polarities. Of the 5,131 metatranscriptomes searched, 1,841 contained  
 224 at least one viroid-like cccRNA.

Table 1: Metatranscriptomic cccRNAs and clusters with sequence matches to viroid-like agents

	Unique	ANI90	
Agent type	cccRNAs	clusters	Most frequent match
Viroid	859	303	Eggplant latent viroid
Satellite RNA	196	91	Lucerne transient streak virus satRNA
Retrozyme	91	29	<i>Fragaria</i> × <i>ananassa</i> retrozyme
Ribozyvirus	49	40	Hepatitis delta virus
Retroviroid	2	2	Carnation small viroid-like RNA

225 Of the sequences aligning to viroid-like agents, the majority only contained  
 226 short (<40 nt) alignable regions, generally localized to the ribozyme motifs.  
 227 However, 33 sequences yielded long (>100 nt) alignments. These cccRNAs  
 228 aligned to Tobacco ringspot virus satRNA (satTRSV), Lucerne transient  
 229 streak virus satRNA (satLTSV), citrus dwarfing viroid (CDVd), and two  
 230 retrozymes. For satTRSV and satLTSV, the range of identity between the



**Figure 2: Viroid-like cccRNAs identified in metatranscriptomes.** A. Number of ANI90 clusters with most significant matches to given viroid-like cccRNA agent types which are either “novel” (derived exclusively from transcriptome and metatranscriptome analysis in this work), “reference” (no novel members), or “shared” (containing at least one of both types of sequence). B. Comparative distribution of inferred ribozyme architectures by cluster type. C. Plot of overlap ratios in cccRNAs, defined as the assembled length divided by the monomer length, from IMG and 1KP. D. Counts of HHR types in representative clusters. E. Length distributions of cluster representatives in the present analysis (transcriptomes and metatranscriptomes), ViroidDB, and a previous study Weinberg et al. (2021). F. Relative abundance of clusters matching different genera within each viroid family by cluster type

recovered cccRNAs and the reference sequence ranged 80%–98% and 81%–99%, respectively. The match to CDVd was 80% identical to the nearest reference sequence. In all cases, the cccRNAs were similar in length and structure to the reference sequences.

We clustered the viroid-like cccRNAs identified here in order to estimate the increase in diversity compared to previously known viroid-like RNAs (Figure 2). Aligning cccRNAs poses a challenge due to the variation in the rotation of the sequences. Two identical cccRNAs could appear to have only half the bases aligning if rotated completely out of phase. Therefore, we took special care to compensate for the circularity of the sequences during the postprocessing of the pairwise nucleotide search results (see Methods). To validate our clustering method, we tested it on ViroidDB. Previously, we identified 458 clusters at the average nucleotide identity (ANI) 90% level in ViroidDB using a method that was not circular-aware (Lee et al., 2021). We identified 50 clusters in ViroidDB using the improved method, generally corresponding to individual species.

In the combined metatranscriptomic, transcriptomic, and reference datasets, we identified 4,823 ANI90 clusters of which 4,121 did not include any sequences from the reference datasets and thus were considered novel, which comprises a 5.9 fold increase in viroid-like RNA diversity. Of the remaining 702 clusters containing at least one known sequence, 288 (41%) were expanded by at least one novel sequence. Notably, 39 novel clusters were

253 represented in plant transcriptomes, of which 8 were symmetric.

254 The relative abundance of HHR types in the cccRNAs varied significantly

255 from what would be expected given the sequence and species count. Within

256 Rfam, HHR1 swamps HHR3 by two orders of magnitude by sequence count

257 (190,679 vs 538 sequences). The same is true for the other three HHR types.

258 However, among the cccRNA cluster representatives, the situation was in-

259 verted: HHR3 was found to be two orders of magnitude more common than

260 HHR1 (1,952 vs 32). Similarly, HHR3 is present in less than one fifth of the

261 HHR-containing species in Rfam whereas among the HHR-containing clus-

262 ter representatives, 94% (1,952/2,074 clusters) contained at least one HHR3.

263 Given the dominance of HHR3 in known viroids (Flores et al., 2017; Navarro

264 et al., 2017), this overabundance of HHR3 is suggestive of the presence of

265 numerous viroid-like cccRNAs.

## 266 **Putative novel viroid-like cccRNAs**

267 We briefly describe the 5 largest novel ANI90 clusters (denoted 1 to 5, in

268 the descending order of the cluster size) derived from the metatranscriptomic

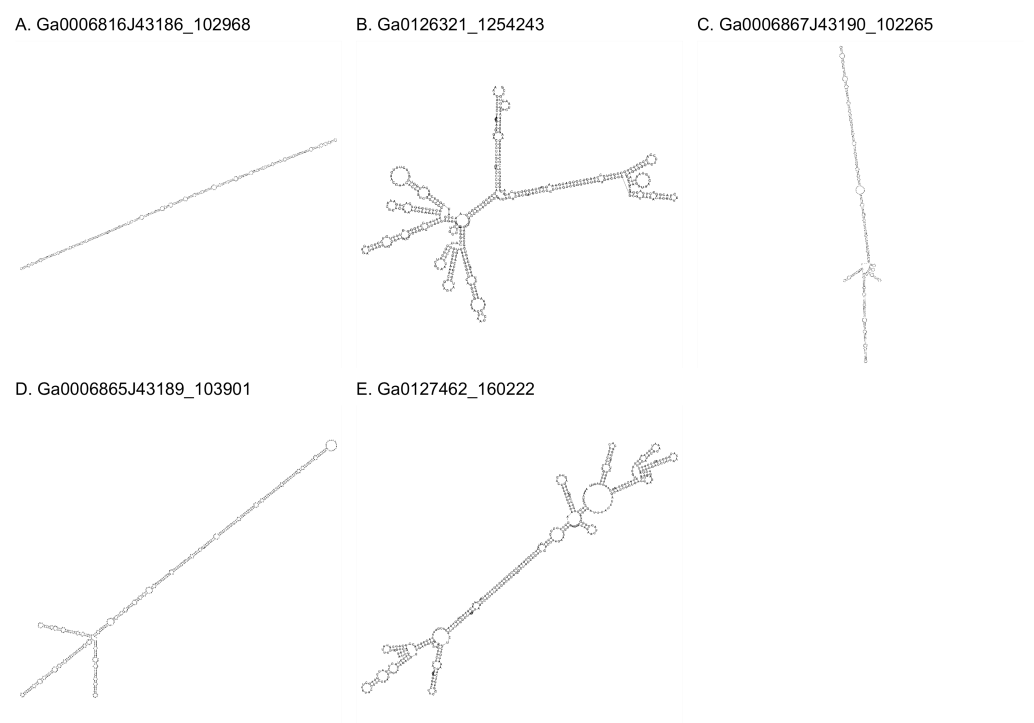
269 data to exemplify the type of findings obtained in this work. All these clusters

270 included members with symmetric, matched ribozymes. The cccRNAs in four

271 of these clusters contained matched HHR3s, whereas those in the fifth one

272 contained twister-P1 ribozymes.

273 The largest, cluster 1, consisted of 149 sequences with a mean length of



**Figure 3: Predicted secondary structures of representatives of the five largest clusters of novel viroid-like cccRNAs.** Structures were predicted using ViennaRNA's RNAfold program configured to operate on circular sequences. Sequence data and metadata are available in Table S1.



274 562 nt ( $\pm 9.0$  nt). During circularity detection, an average of 18% of the  
 275 monomer was removed, although one member of the cluster yielded a contig  
 276 with 60% (341 nt) overlap. The cccRNAs in this cluster are predicted to  
 277 adopt a rod-shaped conformation with 74% of the bases (on average) paired  
 278 in both polarities. Most members of this cluster ( $n=137$ ) contain symmetric  
 279 HHR3s, whereas for the remaining 12 members, only one ribozyme was pre-  
 280 dicted, suggesting the presence of a divergent HHR3. Among the members of  
 281 this cluster, 38 sequences yielded a short (26–37 nt) alignment to the HHR3  
 282 of Eggplant latent viroid or Grapevine hammerhead viroid-like RNA. How-  
 283 ever, the cccRNAs comprising this cluster are substantially longer than the  
 284 respective viroids (334 and 375 nt, respectively). Among the members of this  
 285 cluster, the majority ( $n=133$ ) were found in terrestrial metatranscriptomes  
 286 from 11 distinct locations, 4 members were identified in 3 distinct freshwater  
 287 locations, and one member was found in a spruce rhizosphere sample.

288 The cccRNAs in cluster 2 (68 members) differed in that they contain  
 289 twister-P1 ribozymes in both polarities. All but 4 of these cccRNAs are  
 290 494 nt in length and form a branched structure with a mean of 70% of  
 291 self-complementary bases. These agents were found in 10 unique locations in  
 292 terrestrial ecosystems, such as soil and plant litter. Nearly half ( $n=31$ ) of the  
 293 members were found in switchgrass phyllosphere samples. Cluster members  
 294 were repeatedly found independently during sampling of the switchgrass  
 295 phyllosphere over the course of a year at a sample site in Michigan, USA.

Like cluster 1, cluster 3 (61 members) consisted of comparatively large (605 nt), rod-shaped cccRNA with symmetric HHR3s. However, unlike the largest cluster, no member shows a detectable direct nucleotide match to any ViroidDB sequence, ribozyme or otherwise. In the majority ( $n=40$ ) of the members symmetric ribozymes were not detected, but remaining ones were symmetric, again, suggesting the presence of divergent HHRs. Between 73% and 83% of the bases in these cccRNAs were paired in the (+) polarity and between 74% and 82% of bases were paired in the (-) polarity. On average, 28% of the monomer's length was cleaved during circularity detection although one member was sequenced at 2.78 unit length (almost a complete head-to-tail trimer). The largest two members of the cluster (1,693 and 1,640 nt originally, 1,224 nt after cleavage) were not correctly monomerized by the circularity detection procedure due to mismatches in the seed sequence (see Methods). Manual monomerization of these sequences showed they both were 612 nt, resulting in overlap ratios of 2.76 and 2.67, respectively. Alignment results in 99.0% and 99.5% identity between monomers within each dimer, approximately the error rate of RNA polymerase II when using an RNA template (Gago et al., 2009; López-Carrasco et al., 2017; Wu and Bisaro, 2020). Almost all members of this cluster were identified in soil samples from 6 locations around the world (including Colombia, Czech Republic, Germany, and USA), and two members were found in creek freshwater samples.

Cluster 4 ( $n=61$ ) is similar to clusters 1 and 3 in that its members consist of 615 nt and are predicted to form rods containing HHR3s in both polarities.

319 The secondary structure results in a high mean self-complementarity of 75%  
 320 of bases. However, these cccRNAs contain no alignable regions to the other  
 321 large clusters or to any ViroidDB sequence. As with many other HHR3 rods,  
 322 the majority ( $n=39$ ) contain only one HHR3 above the significance threshold.  
 323 All but two members of the cluster were derived from eight soil locations, and  
 324 the remaining ones were found in the spruce rhizosphere.

325 Unique among the largest clusters we examined, the cccRNAs in cluster 5  
 326 (55 members) are quasi-rod shaped, with 64% bases paired on average in  
 327 the predicted secondary structures. At 454 nt mean length ( $\pm 8.9$  nt), these  
 328 cccRNAs are the smallest among the 5 largest clusters. In this case, all  
 329 members contained two significant HHR3s, but no nucleotide matches to  
 330 ViroidDB were detected. Most members were found in terrestrial samples  
 331 including soil ( $n=38$ ) and plant litter and peat ( $n=2$ ). As in the case of  
 332 cluster 2, 14 members of this cluster were also found over the span of a year  
 333 in the phyllosphere of switchgrass, and one member was found in a freshwater  
 334 sample. Altogether, nine distinct locations contained members of this cluster.

335 In summary, analysis of these 5 largest clusters showed that they consisted  
 336 of cccRNAs endowed with all the hallmarks of viroids including symmetric ri-  
 337 bozymes (HHR3s, with one notable exception), extensive (quasi) rod-shaped  
 338 or branched secondary structure, and evidence of multimeric intermediates.  
 339 Furthermore, the members of these clusters were independently identified in  
 340 diverse samples, primarily, those from soil, indicating they are widespread

341 and consistent with these cccRNAs being infectious agents.

## 342 **Virus-like elements blurring lines between riboviruses,** 343 **ribozyviruses and viroids**

344 Among the cccRNAs containing symmetric HHR3s, we identified rod-shaped  
345 sequences up to 4,705 nt, far outside the length range of viroids and HDV-  
346 like viruses. We hypothesized that these cccRNAs could be novel ribozy-  
347 like viruses. To perform a comprehensive search for potential ribozy-like  
348 viruses, all open reading frames (ORFs) longer than 75 codons present in  
349 cccRNAs were translated, and the resulting sequences of putative proteins  
350 were clustered by amino acid sequence similarity and compared to protein  
351 sequence databases using PSI-BLAST and HHPred (Table S2). Additionally,  
352 we searched the full set of translated cccRNAs sequences against profiles from  
353 the Pfam-A (Mistry et al., 2021) database and protein sequences from the  
354 UniRef90 (Supek et al., 2015) database.

355 Notably, almost all reliable matches were to virus proteins, mostly, to capsid  
356 proteins (CP) of circoviruses, plant satellites RNA viruses, such as Satellite  
357 Tobacco Necrosis Virus, and tombus-like viruses (Table S2). One cluster  
358 of predicted proteins showed significant sequence similarity to the predicted  
359 RdRPs of a unique group of ssRNA viruses, ambiviruses, that were recently  
360 discovered in fungal isolates and transcriptomes (Forgia et al., 2021; Lin-  
361 nakoski et al., 2021; Sutela et al., 2020). Ambiviruses have RNA genomes

of approximately 4 kb which encompass bidirectional ORFs, one of which encodes a predicted RdRP. To date, ambiviruses have not been reported to be circular (Forgia et al., 2021). In the IMG metatranscriptomes, 163 ANI90 clusters (274 cccRNAs total) were predicted to encode ambivirus-like RdRP (E-values between  $1.3e-229$  and  $8.7e-04$ ). Notably, these clusters of cccRNAs were also predicted to contain HHR3, HPR-meta1, and CEPB3 ribozymes, including symmetric sequences with different ribozymes in the two RNA polarities. Furthermore, all these sequences were predicted to adopt a rod-like structure in which the two ORFs encoding, respectively, the RdRP and an uncharacterized protein are arranged along the rod in the opposite strands (Figure 4). These sequences showed varying degrees of terminal overlap, with a median trimmed repeat length of 123 nt. Three representative sequences were recovered with >2000 nt overlaps, of which one was an almost-complete dimer, suggestive of replication via RCR. Three of the ambi-like clusters were detected at very low levels in 10 plant transcriptomes.

We then ran the detection pipeline on the 30 ambivirus and ambivirus-like sequences from GenBank and detected significant ribozyme matches in 15 of these sequences, of which 13 contained two predicted ribozymes. Of the remaining 15 sequences, 11 showed ribozyme matches in the expected locations that failed to pass the significance threshold (Table S3). As in the IMG data, the HHR3 and HPR-meta1 ribozymes are present in both matched and mismatched combinations. Similarly, three of the published genomes (MT354566.2, MN793994.1, and MT354567.1) contain terminal overlaps of

160-250 nt, suggestive of circularity. Furthermore, all known ambivirus and ambivirus-like sequences were predicted to adopt a rod-shaped conformation. Taken together, these observations strongly suggest that ambiviruses comprise a distinct group of ribozy-like viruses that encode a RdRP homologous to the RdRPs of riboviruses. There was no high similarity between the RdRP sequences of ambiviruses and any specific group of riboviruses. However, the HHPred search initiated with the alignment of ambivirus RdRP sequences yielded the top match with the RdRP profile of mitoviruses, suggesting a potential relationship with this group of capsid-less derivatives of RNA bacteriophages (leviviruses) that belongs to the phylum *Lenarviricota* in the realm *Riboviria*.

Three cccRNA clusters with significant mitovirus RdRP matches were detected, including two with symmetric ribozymes. The symmetric singletons are 3,283 and 3,058 nt in size and contain matched twister-P1 ribozymes and an HHR3/twister-P1 combination, respectively. The HHR3 aligns to ELVd with 96% identity. A third cccRNA cluster with three members encoding putative mitovirus-like RdRP, of 3,363 nt, contains a similar match to ELVd (including the HHR conserved core) that was not identified as an HHR by either detection method. This cccRNA lacks the HHR core in the opposite polarity but shows weak similarity (E-value = 0.19) to twister-P1. All cccRRNAs were detected with >100 nt overlaps. These sequences are predicted to adopt a branched conformation with between 63% and 66% of

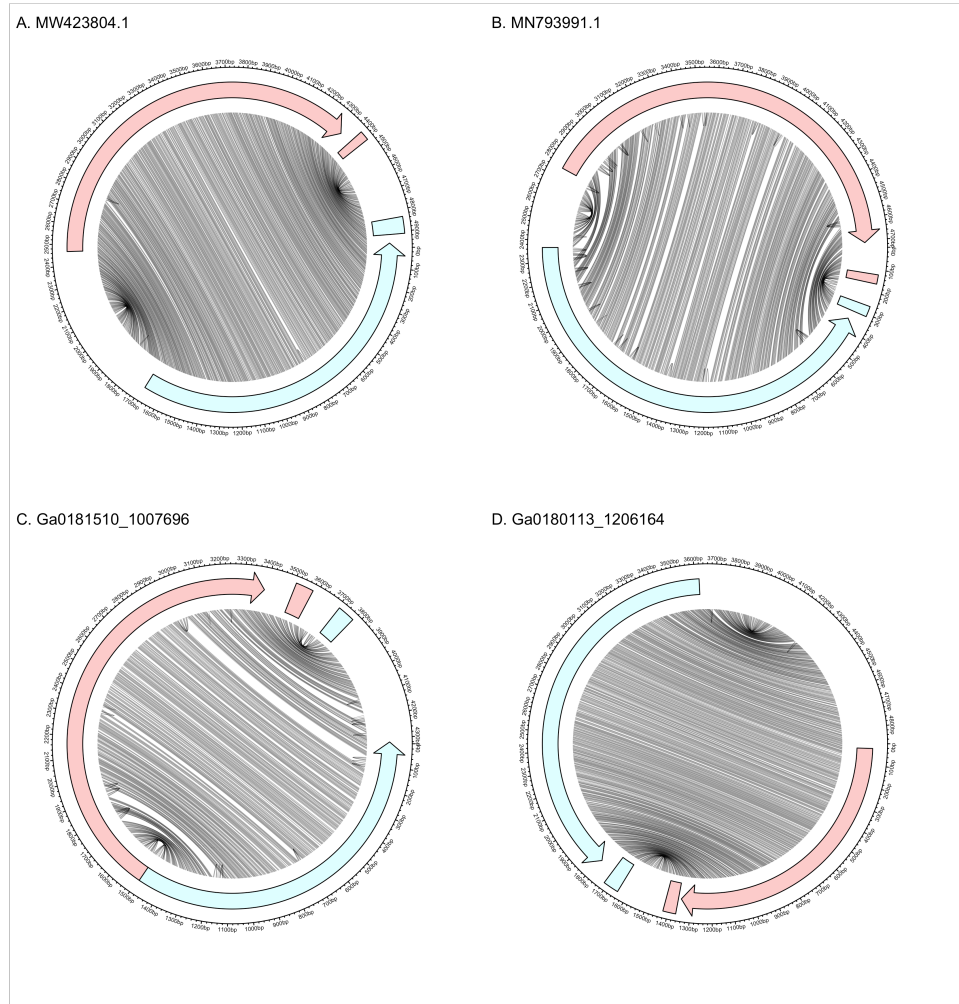


Figure 4: **Genomic and secondary structure** of *Armillaria borealis* ambi-like virus 1 (A), *Tulasnella ambivirus* 1 (B), and novel ambi-like sequences (C, D). Red and blue denote (+) and (-) polarities, respectively. Lines connect bases in the genome that are paired in the predicted secondary structure. Arrows represent ORFs and rectangles represent self-cleaving ribozymes. The (+) and (-) ribozymes are HHR3 and HPR-meta1 in (A), HHR3 and HH3 in (B), CPEB3 and HHR3 in (C), and HHR3 and HPR-meta1 in (D). In all cases, the ribozymes are located outside the ORFs at the end of the rod.

407 bases paired in both polarities. In all three clusters, the putative ribozymes  
 408 are separated by approximately 200 nt. These three genomes have a low (36–  
 409 40%) GC content, a hallmark of mitoviruses (Marais et al., 2017). Search-  
 410 ing the predicted RdRPs against the protein sequence databases yielded the  
 411 most significant matches for the three cccRNAs to Grapevine-associated mi-  
 412 tovirus 13, Grapevine-associated mitovirus 14, and *Fusarium asiaticum* mi-  
 413 tovirus 8. Upon closer inspection of Grapevine-associated mitoviruses 11  
 414 and 13, we found that they also contained HHR3s in both polarities and a  
 415 twister-P1/HHR3 combination, respectively. Ribozyme searches of all avail-  
 416 able *Lenarviricota* (taxid 2732407) and unclassified *Riboviria* (taxid 2585030)  
 417 sequences did not yield other matches besides the ambiviruses and these few  
 418 mitoviruses.

419 Apart from the RdRps, we identified 135 sequences comprising 53 ANI90  
 420 clusters with significant similarity to capsid proteins of single-stranded (ss)  
 421 DNA viruses, in particular, CRESS viruses (Krupovic et al., 2020). The  
 422 sequences in 50 of these clusters contained predicted ribozymes, and 13 con-  
 423 tained HHR3s in both polarities. Two clusters contained paired HHR3 and  
 424 twister-P1 ribozymes, whereas two other clusters contained symmetric HP-  
 425 meta1 ribozymes. 26 clusters contained a single HHR3, three a twister-P1,  
 426 and four a HPR-meta1. 21 clusters, including all three without ribozyme pro-  
 427 file matches, produced a nucleotide alignment to a known viroid’s ribozyme.  
 428 These alignments ranged in length from 25 to 50 nt at 83–96% identity. The  
 429 cccRNAs in these clusters ranged between 1,092 and 1,632 nt in length, with



a mean of 1,317 nt and GC content between 35% and 51%, with the mean of 44%. Four cccRNAs were sequenced as complete head-to-tail dimers. The secondary structure of these cccRNAs showed extensive self complementarity, with 66% of the bases paired. Given the strong evidence of circularity, extensive self-complementarity resulting in predicted branched structure and confident prediction of ribozymes, these cccRNAs most likely represent a novel class of ribozy-like satellite viruses.

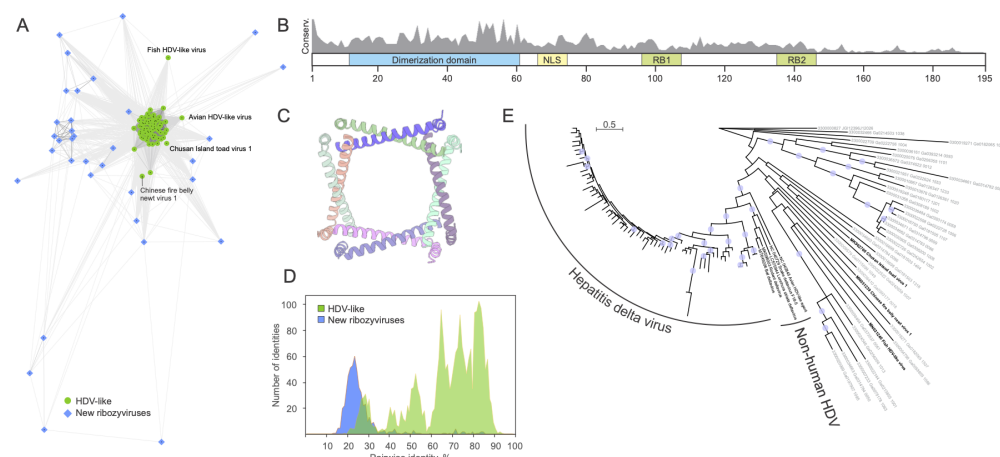
## **Novel ribozy-like viruses**

Apart from the viruses that resembled ribozyviruses conceptually, that is, were identified as protein-coding viroid-like cccRNA but encoded proteins unrelated to HDV antigen (HDVAg), we searched for actual relatives of HDV. To this end, the metatranscriptomes were searched for cccRNAs encoding HDVAg homologs by comparing clusters of ORFs from cccRNAs to the sequences of the HDVAg and its homologs from other known ribozyviruses. A total of 12 ORF clusters were identified above the HDVAg Pfam profile's gathering threshold; additional 21 representative ORFs were significant at the E-value < 1e-03 level, and 34 at E-value < 1e-02. Of these clusters, only one showed a significant nucleotide alignment to a known ribozyvirus. The other clusters were found in a variety of environments ranging from soil to wastewater to coastal wetland sediment. Samples with matching cccRNAs were collected from as far north as Alaska, USA, to as far south as Florida, USA. All 69 members of the clusters encompassing ORFs with HDVAg pro-

file matches (E-value < 1e-02) were predicted to adopt a rod shape in both polarities. The genome size of ribozyviruses in ViroidDB ranges from 1,547 to 1,735 nt. However, among the HDV-like clusters identified in metatranscriptomes, the size ranged from 1,019 nt to 1,757 nt, with a median length of 1,317 nt.

Clustering of the HDVAg sequences, their homologs from other animals and the homologs from metatranscriptomes with a permissive threshold using CLANS showed that all previously known HDVAg homologs formed a single tight cluster whereas the metatranscriptome sequences formed multiple smaller clusters and singletons distant from each other and from HDV (Figure 5A). The conservation profile of the multiple alignment of the HDVAg homologs showed that the dimerization region and one of the RNA-binding regions were prominently conserved whereas the second RNA-binding region was not (Figure 5B,C). The sequences of the distant HDVAg homologs from metatranscriptomes showed low sequence similarity to the previously known HDVAg, far below the similarity among the latter, with the distributions of percent identities almost non-overlapping (Figure 5D). Finally, in the phylogenetic tree of the HDVAg homologs, all previously known sequences formed one compact clade, whereas the homologs from metatranscriptomes identified here comprised several remaining clades, with a much greater phylogenetic depth (Figure 5E).

The nucleotide sequences of these HDV-like cccRNAs formed 26 ANI90



**Figure 5: Diversity of HDV antigen-like proteins in known and newly discovered ribozyviruses.** A. Clustering of the HDV antigen (Ag)-like protein homologs based on their sequence similarity. Lines connect nodes (sequences) with P-value < 1e-05. Reference HDVAg-like sequences from GenBank are shown as green circles, whereas those detected in metatranscriptomic datasets as blue diamonds. Some of the divergent reference sequences are labeled. B. Schematic representation of the HDVAg with functionally important regions indicated with colored boxes. RB1 and RB2, RNA-binding sites 1 and 2, respectively; NLS, nuclear localization signal. Gray histogram shows the sequence conservation (percent identity) of HDVAg-like sequences from metatranscriptomic datasets. C. Octameric structure of the conserved dimerization domain of HDVAg. PDB ID: 1A92 (Zuccola et al., 1998). Each protein molecule is shown with a different color. D. Comparison of the sequence conservation among reference HDVAg from GenBank (green) and those from metatranscriptomic datasets (blue). E. Maximum likelihood phylogeny of HDVAg-like sequences. The tree was constructed with IQ-TREE (Nguyen et al., 2015). Circles at the nodes represent SH-aLRT support higher than 90%. The scale bar represents the number of substitutions per site.

clusters, none which contained confidently predicted self-cleaving ribozymes above the respective gathering thresholds. However, 13 of these clusters produced weak ribozyme matches (E-value < 1e-01), and 8 of these were symmetric. Both HHR-like ( $n=36$ ) and HDV-like ( $n=26$ ) ribozymes were detected although no clusters contained ribozymes of both types. Of the HDV-like ribozymes detected, only five most closely matched the canonical HDV ribozyme. Ten putative ribozymes showed the strongest similarity to the HDV ribozyme (HDVR) found in the genome of *F. prausnitzii* (Webb et al., 2009), seven were most similar to the HDV-like ribozyme found in the genome *A. gambiae* (Webb et al., 2009), and four were most similar to the mammalian CPEB3 ribozyme (Chadalavada et al., 2010; Salehi-Ashtiani et al., 2006).

The limited number of significant ribozyme matches among the HDV-like sequences posed an opportunity for detecting novel ribozymes or diverged variants of known ones. For example, we examined an HDV-like cccRNA cluster (representative member 3300009579\_Ga0115599\_1049451) with no predicted ribozymes. However, upon closer examination, sequences from this cluster were shown to contain the conserved HHR core in both polarities in the expected locations, a recently discovered ribozyme configuration (de la Peña et al., 2021). Some clusters entirely lacked the HHR core in either polarity, suggesting the use of alternative, yet unknown ribozymes.

Clustering the HDV-like sequences in combination with the known ri-

bozyviruses in ViroidDB produces no ANI80 clusters with both reference and novel members. Each of the 26 HDV-like clusters falls below the species demarcation criterion for ribozyviruses (80% nucleotide identity) (Hepojoki et al., 2020). Clustering the ORFs from both the detected HDV-like sequences and reference ribozyviruses with 60% minimum identity (the genus demarcation criterion) using CD-HIT resulted in 36 clusters, of which 10 consisted entirely of reference sequences whereas 26 were entirely novel.

## **Novel ribozyme combinations and unusual self-cleaving ribozymes**

Almost all viroid-like RNAs described to date contain the same type of ribozyme in both polarities, with the exception of some satRNAs. Surprisingly, many viroid-like cccRNAs identified in this work were predicted to contain self-cleaving ribozyme combinations that have not been so far reported in replicating cccRNAs (Figure 6).

Specifically, we identified numerous cccRNAs containing twister ribozymes, a recently described ribozyme motif that so far has only been found in combination with the HP-metal ribozyme. Both symmetric ( $n=381$ ) and asymmetric ( $n=930$ ) variants are present in the metatranscriptomic cccRNA clusters. Most symmetric twister clusters contained matched twister ribozymes (218 clusters) in both polarities, a novel combination. In 87 clusters including mitovirus-like and satellite-like cccRNAs, we found another novel combi-

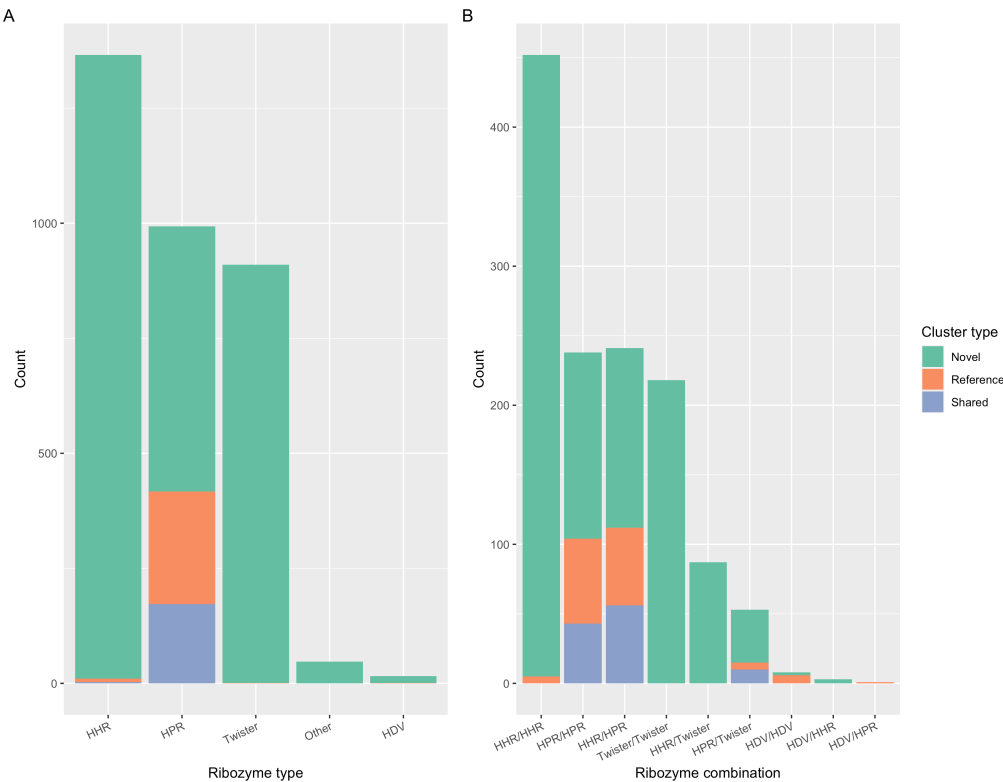


Figure 6: **Ribozyme diversity in viroid-like cccRNAs.** A. Distribution of ribozyme types in asymmetric clusters. B. Ribozyme co-occurrence within the symmetric viroid-like cccRNA cluster representatives derived from metatranscriptomes.

517 nation of ribozymes, with HHRs opposite twister ribozymes. The unusual  
518 twister ribozyme is widespread in plant transcriptomes, with 59% of the tran-  
519 scriptomes containing reads mapping to a twister-bearing cccRNA. Indeed,  
520 we recovered three asymmetric cccRNA clusters from plants that contained  
521 a twister-P1 ribozyme.

522 In addition to the twister ribozyme combinations, we identified several other  
523 novel ribozyme combinations in symmetric cccRNAs. Previously, HHR3s  
524 have been found in conjunction with HDVAg (de la Peña et al., 2021), but  
525 have not been reported to be paired with HDV ribozymes. We identified three  
526 clusters in which HHR3s were paired with HDV-type ribozymes, namely,  
527 CPEB3 and HDVR *F. prausnitzii*. The novel CPEB3-HHR3 combination  
528 was found in an ambivirus-like sequence, and the two HDVR *F. prausnitzii*/-  
529 HHR3 clusters were both predicted to adopt rod-shaped structures. One  
530 of these, a 1,052 nt singleton, did not match any sequences in ViroidDB,  
531 nt, or UniRef90, but in the other cluster (978 nt, two members), the HHR  
532 was closely similar to that of *Cryphonectria parasitica* ambivirus 1 (44/50 nt  
533 identical), whereas the HDVR *F. prausnitzii* motif (32/33 nt and 41/46 nt)  
534 aligned to two chromosomes of the *Vanessa atalanta* butterfly.

535 Among the asymmetric cccRNAs, we identified two additional types of self-  
536 cleaving ribozymes that so far have not been reported among viroid-like  
537 RNAs. The hatchet ribozyme was found in 34 ANI90 clusters that ranged  
538 from 357 nt to 567 nt in length and came primarily from aquatic metatran-

539 scriptomes, in contrast to the general trend among the cccRNAs that derived  
540 from soil metatranscriptomes (Figure 7B). For example, the most diverse of  
541 these clusters (440 nt) contained 22 members from aquatic (almost all fresh-  
542 water) metatranscriptomes sampled from around the United States. For the  
543 hatchet clusters, the Rfam profile matches were the strongest among all de-  
544 tected ribozymes (median E-value 1.4e-08) and, unusually for viroid-like cc-  
545 cRNAs, the GC content was low (median 35%). Among ViroidDB sequences,  
546 only avocado sunblotch viroid and one retrozyme are below 40% GC. The  
547 predicted structures of these sequences varied from branched to quasi-rod  
548 shaped, with a mean of 62% of the bases paired.

549 The other self-cleaving ribozyme previously not known to exist in viroid-like  
550 agents, the pistol ribozyme, was identified in asymmetric cccRNAs. Like the  
551 clusters containing the hatchet ribozyme, clusters with the pistol ribozyme  
552 were found primarily (9/13 clusters) in marine metatranscriptomes ranging  
553 from the Antarctic Ocean to the Baltic Sea. The clusters have a slightly  
554 lower median profile match E-value (E-value = 6.4e-05) compared to the  
555 hatchet ribozyme but, unlike the hatchet ribozyme, have GC content ranging  
556 from 33% to 59% (median 49%) more characteristic of viroid-like RNAs.  
557 The predicted secondary structures of the pistol-containing cccRNAs were  
558 branched, often with several long hairpin structures. At least one pistol-  
559 containing cccRNA (E-value = 4.3e-11) was actually a genomic segment from  
560 a bequatrovirus (a dsDNA bacteriophage) with a short terminal repeat. The  
561 presence of a ribozyme in a bequatrovirus is surprising and represents the



562 first report of a self-cleaving ribozyme within this genus as well as the first  
563 report of a pistol ribozyme in a virus. Further investigation is required to  
564 determine whether this virus genomic segment actually encodes a viroid-like  
565 cccRNA.

## 566 **CRISPR spacers matching cccRNAs**

567 CRISPR spacer matches provide one of the most reliable means for assign-  
568 ing hosts to viruses and other mobile genetic elements in prokaryotes, and to  
569 differentiate prokaryote-infecting from eukaryote-infecting viruses (Munson-  
570 McGee et al., 2018; Paez-Espino et al., 2016). For instance, our recent  
571 search for riboviruses in the same set of metatranscriptomes that is ana-  
572 lyzed here identified multiple spacers matching RNA viruses, resulting in  
573 the assignment of several groups of viruses to bacterial hosts including sev-  
574 eral previously thought to infect eukaryotes (Neri et al., 2022). To identify  
575 viroid-like agents that potentially might replicate in prokaryotes, we searched  
576 the viroid-like cccRNA sequences identified here against the IMG CRISPR  
577 spacer database, built from CRISPR arrays predicted in IMG genomes and  
578 metagenomes (Chen et al., 2021), and detected 89 spacers with significant  
579 matches to viroid-like cccRNAs from 9 clusters.

580 One spacer was an identical match of 37 nt to a member of a cluster of 16  
581 cccRNAs with prominent viroid-like features. The cccRNAs of this cluster  
582 are 315 nt long, contain symmetric HHR3s, and are predicted to adopt a

rod shape, with 73% of the bases paired. Unusual for viroids, the cccR-  
 NAs comprising this cluster were found in Mushroom Spring, a hot spring  
 at Yellowstone National Park. The spacer match was also detected in a  
 60° C hot spring, Great Boiling Spring, albeit more than 800 km away, in  
 northwestern Nevada (Thomas et al., 2019). The repeats in this CRISPR  
 locus were identical to those in the type III CRISPR locus of *Roseiflexus*  
 sp. RS-1, an anoxygenic filamentous bacterium of the phylum *Chloroflexota*  
 that was itself identified in Yellowstone hot springs (Madigan et al., 2005;  
 van der Meer et al., 2010). Searching the same IMG metagenomic spacer  
 database for matches to all 16 cluster members with more relaxed criteria  
 (*i.e.*, more than 1 mismatch but with E-value < 1e-05), we identified a further  
 13 nearly identical matching spacers from 8 Yellowstone hot springs samples  
 collected between 2007 and 2017. One spacer (35/37 identities, E-value =  
 3e-06) included a precise match to the HHR core motif. The repeats from  
 this expanded set of loci all matched those from *Roseiflexus* sp. RS-1. In our  
 previous study, we identified multiple spacers in *Roseiflexus* sp. RS-1 that  
 matched a group of partiti-like riboviruses that were accordingly assigned  
 to this bacterial host (Neri et al., 2022). Apparently, the type III CRISPR  
 system of *Roseiflexus* sp. that encompasses a reverse transcriptase actively  
 incorporates spacers from multiple RNA replicons.

The cluster with the most spacer matches—57 matches spanning 26  
 metagenomes, largely from sludge and bioreactor samples—includes a single  
 cccRNA of 606 nt (recovered as 841 nt) with an asymmetric twister-P5

606 ribozyme. Eight spacers, detected in 7 metagenomes, covered the predicted  
607 ribozyme region. This sequence was predicted to adopt a branched con-  
608 formation with 63% of bases paired. Nucleotide and translated nucleotide  
609 searches did not produce any matches to this cccRNA.

610 The second most frequently matched cluster, also a singleton, was associated  
611 with 13 spacers from 10 metagenomes, all from the same location in the  
612 Southern Indian Ocean. This cccRNA is 286 nt long (recovered with a 123  
613 nt overlap) and contains a predicted HHRII in one polarity only. Like the  
614 most matched singleton, this sequence also is predicted to adopt a branched  
615 conformation, with 57% of bases paired. The spacers collectively covered  
616 33% of the sequence but do not include the HHR region. No homologs of  
617 this sequence were detected in nucleotide or protein databases.

618 We also identified a 29 nt spacer that was identical to the HHR3 core of two  
619 distinct asymmetric cccRNA clusters. One of these is a 1,121 nt singleton  
620 and the other is a set of three 535 nt cccRNAs; both are predicted to form a  
621 branched conformation, with 68% and 65% of the bases paired, respectively.  
622 All four cccRNAs in these clusters also matched the HHR3 of Velvet mottle to-  
623 bacco virus satRNA with 27/29 (93%) and 37/43 (86%) identical nucleotides,  
624 respectively. Furthermore, an identical match to the spacer was detected in  
625 the HHR of an uncultured organism found in the cerebrospinal fluid virome  
626 of a patient with unexplained encephalitis (Jimenez et al., 2011). 26 nt sub-  
627 sequences of the spacer also aligned with 100% identity to Chicory yellow

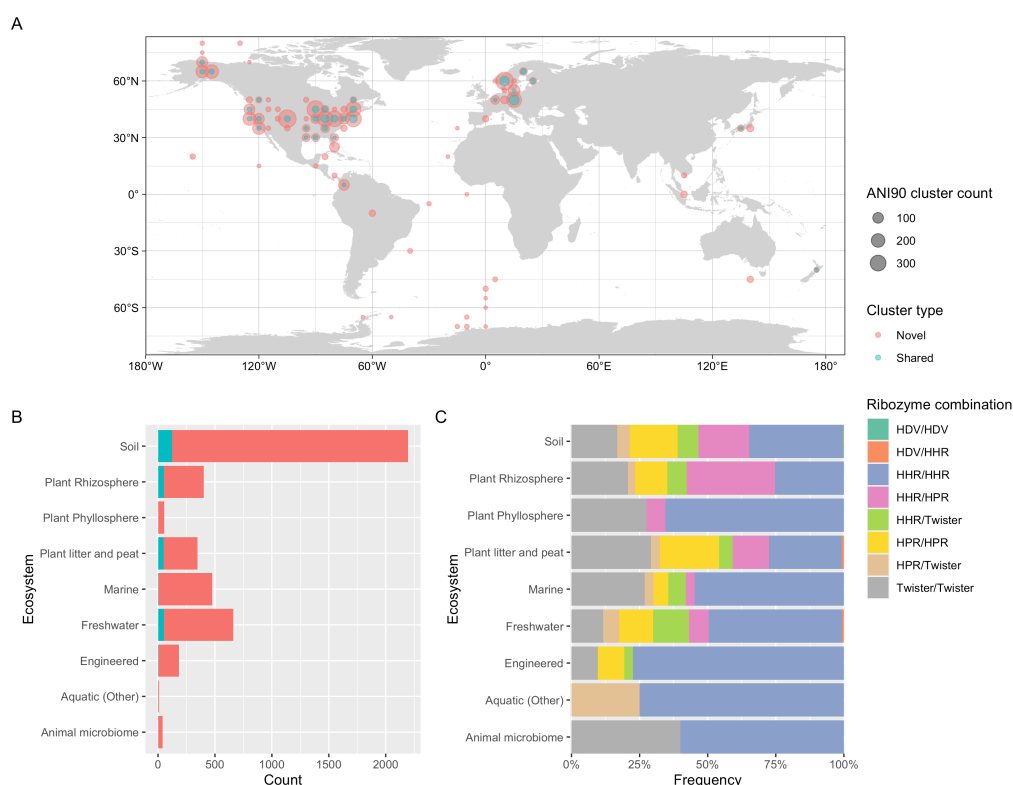
628 mottle virus satRNA and Cereal yellow dwarf virus-RPA satRNA. Four of  
 629 the five genes within the spacer's scaffold (including Cas1 and reverse tran-  
 630 scriptase) were most similar to a type III CRISPR locus of a purple sulfur  
 631 bacterium, *Thiorhodococcus drewsii*. This scaffold originated from a 29° C  
 632 desert spring sample from southern Nevada, whereas the cccRNA clusters  
 633 were found in peatland samples from Minnesota and in spruce litter from the  
 634 Bohemian Forest, Czech Republic.

## 635 **Geographic and ecological distribution of viroid-like cc-** 636 **cRNAs**

637 We examined the global distribution of the cccRNA clusters. Novel clusters  
 638 were found throughout the world (Figure 7A) and in all types of ecosystems  
 639 (Figure 7B). Soil samples were the primary source of both novel and shared  
 640 clusters, reflecting both the greater number of such samples (twice as many  
 641 as the next most common sample type) and the apparent greater sequence  
 642 diversity in soils.

643 Surprisingly, the viroid-like cccRNAs displayed non-uniform ribozyme distri-  
 644 bution among ecosystems (Figure 7B). Mismatched HPR/HHR ribozymes  
 645 were especially prevalent among samples from plant rhizospheres, whereas  
 646 matched HHRs were significantly more abundant in engineered ecosystems,  
 647 such as bioreactors, than in soil environments.

648 Among the 10 most geographically dispersed novel clusters (that is, clus-



**Figure 7: Global distribution and habitats of viroid-like cccRNAs found in metatranscriptomes.** A. Map of sample locations from which viroid-like cccRNAs were detected. The size of each circle corresponds to the number of unique clusters identified in each location (grouped to the nearest five degrees of latitude and longitude) while the color represents the novelty of the clusters. B. The distribution of clusters within each type of ecosystem by novelty. C. The relative frequency of ribozyme combinations within symmetric cccRNA clusters in each ecosystem type.

ters with the most members found in distinct latitudes and longitudes after  
rounding to the nearest degree), 8 included symmetric members, of which  
6 were matched HHR3s. The other two symmetric clusters contain HPR-  
metal/twister-P1 ribozymes and the two asymmetric clusters contain HHR3s.  
These widely dispersed clusters ranged in length from 372 nt to 1,039 nt and  
were predicted to adopt either a rod-like shape (the 6 HHR3 clusters) or a  
branched conformations.

Identifying potential hosts of the viroid-like agents within these ecosystems  
remains a challenge. Based on the IMG annotation pipeline, the majority of  
the analyzed metatranscriptomes were dominated by prokaryotic sequences  
(Neri et al., 2022), but still contained at least 1% of contigs affiliated to  
eukaryotes (Table S4). Nonetheless, 187 metatranscriptomes in which viroid-  
like cccRNAs were detected contained  $< 0.1\%$  of eukaryote-affiliated contigs,  
suggesting that these elements replicate in either rare and undetected eu-  
karyotic hosts or in some of the much more abundant prokaryotic hosts.  
Notable among these datasets were the hot spring metatranscriptomes in  
which putative CRISPR targeting of viroid-like cccRNAs were identified (see  
above). The apparent lack of eukaryotic RNA in these samples strengthen  
the hypothesis of a prokaryotic host. Similarly, we found 104 symmetric  
clusters in marine samples which are far beyond the habitation range of the  
known hosts of viroids and satRNAs. To date, only ribozyviruses have been  
reported in aquatic hosts (Chang et al., 2019) although it has to be taken  
into account that viroids are highly stable and have been reported to re-

main infectious even after up to 7 weeks in the water (Mehle et al., 2014).  
The existence of these viroid-like clusters combined with the clusters from  
prokaryote-dominated samples suggests that the actual ecological and host  
range of viroid-like agents is far broader than currently appreciated.

## Discussion

Viroids and viroid-like cccRNAs, such as satRNAs and ribozymes, are the  
smallest, simplest known replicators that hijack either a host DNA-dependent  
RNA polymerase or a virus RdRP for their replication. Given the universality  
of the cellular transcription machinery across all life forms and the widespread  
of RdRP-encoding riboviruses, the narrow diversity and host range of the  
known viroid-like agents appeared puzzling. We suspected that viroid-like  
agents actually could be far more common than presently known and, with  
this motivation, searched a collection of about 5,000 metatranscriptomes for  
viroid-like cccRNAs.

We were able to identify millions of putative cccRNAs by identifying sig-  
natures of circularity or RCR, namely, the presence of head-to-tail repeats  
in assembled contigs. Because reads spanning the origin cannot be recon-  
ciled with a linear sequence, the assembler produces contigs with the same  
subsequence repeated at both the end and the beginning (Qin et al., 2020).  
Alternatively, when linear replication intermediates containing head-to-tail  
repeats are sequenced, the ends of the sequences are also repeated. After com-

693 compensating for the low-fidelity of RNA polymerase II by allowing for up to 5%  
694 mismatches in the repeated regions, testing this method on assembled plant  
695 transcriptomes demonstrated that known viroids were reliably recovered in  
696 the absence of major assembly errors. However, the secondary structure of  
697 many viroid-like cccRNAs and the use of poly-A enrichment during RNA  
698 isolation prior to sequencing (Johnson et al., 2012) makes it likely that many  
699 viroid-like cccRNAs either were not sequenced at all or were grossly missam-  
700 bled and thus could not be recognized as circular. Even under a conservative  
701 approach, where only predicted cccRNAs containing confidently identified  
702 ribozymes counted as “viroid-like”, this search resulted in an approximately  
703 five-fold increase in the diversity of viroid-like agents. This is most likely a  
704 substantial underestimate of the true span of the viroid-like domain of the  
705 replicator space because among the millions of the predicted cccRNAs, in  
706 which no ribozymes were confidently identified, some, and possibly, many  
707 could be viroid-like agents containing novel ribozymes or lacking ribozymes  
708 altogether like pospiviroids. Although perhaps only a coincidence, it is worth  
709 noting that a recent analysis of the same collection of metatranscriptomes  
710 also yielded an approximately five-fold increase in the diversity of riboviruses  
711 (Neri et al., 2022).

712 The majority of the detected viroid-like cccRNAs possessed characteristic fea-  
713 tures of viroids including the presence of HHR, often in both polarities, and  
714 predicted rod-like or extensive branched conformation. However, the search  
715 resulted not only in quantitative but also in considerable qualitative novelty,



716 in particular, with regard to previously undetected ribozyme combinations,  
717 such as those including twister, and detection of ribozymes not previously  
718 known to exist in viroid-like RNAs, hatchet and pistol. There seems to be  
719 little doubt that additional ribozyme combinations and novel ribozymes in  
720 viroid-like cccRNAs remain to be discovered.

721 Another key finding is the discovery of novel groups of ribozy-like viruses.  
722 The substantial expansion of ribozyviruses themselves, that is, viroid-like cc-  
723 cRNAs encoding homologs of the HDV antigen, probably, should have been  
724 expected. It is nevertheless notable that the diversity of the ribozyvirus  
725 sequences discovered in metatranscriptomes far exceeds that of the previ-  
726 ously known HDV relatives including recently discovered non-mammalian  
727 ones. Moreover, many of the newly identified ribozyviruses lack the HDV  
728 ribozyme or even any known ribozymes, suggesting novel replication mecha-  
729 nisms. The host range of the new ribozyviruses remains to be explored but,  
730 probably, includes non-animal hosts (see discussion below).

731 In contrast, the demonstration that ambiviruses are actually viroid-like  
732 agents and the discovery of viroid-like mitoviruses and satellite viruses was  
733 surprising. These findings amend our understanding of the relationships  
734 between viroids and riboviruses. The latter three groups of viroid-like agents  
735 resemble ribozyviruses in that these are relatively large, protein-coding  
736 viroid-like cccRNAs. However, unlike HDV and its relatives, these viroid-like  
737 agents are clearly linked to riboviruses through the RdRPs encoded by

738 ambiviruses and mitoviruses, and capsid proteins encoded by satellite  
 739 viruses. These findings show that combinations of viroid-like cccRNA  
 740 and protein-coding genes emerged on multiple occasions during evolution.  
 741 It appears likely that these ribozy-like (but unrelated to HDV) viruses  
 742 originated through recombination between typical riboviruses and viroids.  
 743 The implications for virus taxonomy, in particular, whether such viruses  
 744 should be classified into the existing realm *Ribozyviria* or into the respective  
 745 divisions of the realm *Riboviria* (Koonin et al., 2020), or into completely  
 746 new taxa, remain to be sorted out.

747 One of the most interesting but also most challenging problems is the host  
 748 range of the expanded diversity of viroid-like agents. There is currently no  
 749 direct computational approach for connecting viroid-like RNA with specific  
 750 hosts. Nevertheless, it appears exceedingly unlikely that all or even the  
 751 majority of the viroid-like cccRNAs discovered in metatranscriptomes are  
 752 parasites of plants. Indeed, firstly, we identified orders of magnitude more  
 753 viroid-like cccRNAs in metatranscriptomes than in plant transcriptomes, and  
 754 secondly, most of the analyzed metatranscriptomes are dominated by bacte-  
 755 ria followed by unicellular eukaryotes. Furthermore, ambiviruses were iso-  
 756 lated from fungi (Forgia et al., 2021; Linnakoski et al., 2021; Sutela et al.,  
 757 2020), and the demonstration that these are viroid-like agents expands the  
 758 host range of the latter. A potential prokaryotic connection of viroid-like  
 759 cccRNAs through CRISPR spacer matches is of special interest given that  
 760 so far the host range of such agents appeared to be exclusively eukaryotic.

761 The detected spacer matches were not numerous but reliable, in particular,  
 762 because multiple spacers matching the same cccRNA were identified in di-  
 763 verse metagenomes. At least the typical viroid-like cccRNAs that matched  
 764 spacers from the reverse-transcribing type III CRISPR system of *Roseiflexus*  
 765 sp. appear to be strong candidates for novel bacterial parasites. These  
 766 and other viroid-like cccRNAs matched by CRISPR spacers seem to merit  
 767 further, dedicated metatranscriptome and metagenome searches as well as ex-  
 768 perimental investigation. These findings, even if preliminary, echo the recent  
 769 expansion of the bacterial RNA virome through the search of the same meta-  
 770 transcriptome collection and suggest that bacteria might support a much  
 771 greater diversity of RNA replicators than previously suspected (Neri et al.,  
 772 2022).

## 773 Limitations of the Study

774 This work deals with a relatively low hanging fruit in the search for viroid-  
 775 like agents, being limited to the cccRNAs that contain reliably identifiable,  
 776 known ribozymes or align directly to known viroid-like agents. This con-  
 777 servative approach was adopted purposefully, in order to avoid potential  
 778 artifacts resulting from erroneous identification of cccRNA, contamination  
 779 with DNA-encoded sequences or other sources. A potential opportunity for  
 780 the discovery of a far greater diversity of viroid-like agents and a challenge  
 781 for further research is a comprehensive analysis of the massive set of pre-  
 782 dicted cccRNAs that lack known ribozymes. Evidently, the computational

783 approaches applied in this work only identify candidates for viroid-like agents.  
784 Experimental validation is needed and is especially important in the case of  
785 putative cccRNAs lacking known ribozymes.

## 786 **STAR Methods**

### 787 **Resource availability**

#### 788 **Lead contact**

789 Further information and requests for resources and data should be directed  
790 to and will be fulfilled by the lead contact, Benjamin Lee (benjamin.lee@chc  
791 h.ox.ac.uk).

### 792 **Materials Availability**

793 This study did not generate new unique reagents. All data used as inputs  
794 for the pipeline are listed in the “data acquisition” section. All output of  
795 analyses are listed in the “data and code availability” section below.

### 796 **Data and code availability**

- 797 • This paper analyzes existing, publicly available data. These accession  
798 numbers for the datasets are listed in the key resources table. Data  
799 generated during downstream analysis have been deposited at Zenodo  
800 and are publicly available as of the date of publication. DOIs are listed

- 801 in the key resources table.
- 802 • All original code has been deposited at Zenodo and is publicly available
  - 803 as of the date of publication. DOIs are listed in the key resources table.
  - 804 • Any additional information required to reanalyze the data reported in
  - 805 this paper is available from the lead contact upon request.

## 806 **Method Details**

### 807 **Data acquisition**

808 The search for novel cccRNAs was performed on a collection of 5,131 assem-  
 809 bled metatranscriptomes sourced from the IMG/MER database. In addition  
 810 to the IMG metatranscriptomes, we searched the complete set of transcrip-  
 811 tomes of the 1,000 Plants (1KP) project, totalling 1,344 paired-end sequenc-  
 812 ing runs. Before applying the search pipeline to 1KP, we filtered the raw  
 813 reads for quality using fastp (Chen et al., 2018) and assembled them using  
 814 rnaSPAdes (Bushmanova et al., 2019) using default parameters. We also  
 815 included the 2021-09-07 release of ViroidDB and the set of HPR sequences  
 816 identified by Weinberg et al. (2021).

### 817 **Circularity detection**

818 We identified cccRNAs using a modified and improved version of the  
 819 reference-free and *de novo* Cirit algorithm (Gao et al., 2015; Qin et al.,  
 820 2020) implemented in the Nim programming language. This method relies

821 upon assembly errors for circular sequences resulting in terminal repeats.  
 822 Detecting cccRNAs via this method requires searching forward for the last  
 823 several bases of the sequence (the seed) and, if a match is found, comparing  
 824 backwards to the start of the sequence. If the start and end of the  
 825 sequence “overlap”, this repetitive region is then trimmed off. However, the  
 826 existing implementation was unable to monomerize multimeric transcripts  
 827 resulting from rolling circle replication among known viroid-like agents  
 828 due to its single pass design and requirement of exact sequence identity  
 829 within repetitive regions. Our reimplementaion solves these problems by  
 830 reiteratively attempting to monomerize putative cccRNAs while allowing  
 831 for a configurable minimum identity within repeats. For this study, we  
 832 required a minimum of 95% identity with no insertions or deletions within  
 833 the overlapping region. In addition, the ratio of the length of the contig and  
 834 computed monomers is reported. Sequences with monomer lengths below a  
 835 threshold of 100 nt were excluded.

## 836 **cccRNA deduplication**

837 Standard approaches to sequence deduplication are insufficient for cccRNAs.  
 838 Most modern approaches rely on hashing for memory efficiency. Such ap-  
 839 proaches are effective for linear sequences but circular sequences pose a chal-  
 840 lenge due to the arbitrariness of their start position. To enable deduplication  
 841 of putative cccRNAs, we define a sequence’s canonical representation as the  
 842 alphabetically earlier of the lexicographically minimal rotations of the se-

843 quence and its reverse complement. This approach, drawn from k-mer count-  
844 ing methods (Marçais and Kingsford, 2011; Melsted and Pritchard, 2011)  
845 and, if further optimized, would be able to be computed in linear time and  
846 with constant memory for even greater scalability (Duval, 1983).

## 847 **Ribozyme-based filtering**

848 To identify sequences likely to replicate via ribozyme-catalyzed rolling cir-  
849 cle replication, we searched the unique cccRNAs for the presence of known  
850 self-cleaving ribozymes using Infernal (Nawrocki and Eddy, 2013). In each  
851 polarity, we identified ribozymes above Rfam’s curated gathering cutoff or  
852 with E-values  $< 0.1$ . Sequences with ribozymes in both polarities that met  
853 these criteria were considered viroid-like. Alternatively, we considered se-  
854 quences as viroid-like with one ribozyme with an E-value  $< 0.01$  or a score  
855 above the gathering cutoff. For each polarity, we considered only the most  
856 significant (by E-value) ribozyme. To identify more divergent ribozymes  
857 that were not detected using Infernal, we searched sequences containing one  
858 significant (E-value  $< 0.01$ ) using RNAmotif (Macke, 2001).

## 859 **Sequence search**

860 We searched all unique cccRNAs against ViroidDB using MMseqs2 easy-  
861 search (version 13.45111) (Steinegger and Söding, 2017) with the highest  
862 available sensitivity (`-s 7.5`). For each sequence, we considered only the  
863 most significant match as determined by bit score. In addition to searching

ViroidDB, we also searched the novel cccRNAs identified by metatranscriptome mining against the set of cccRNAs recovered from plant transcriptomes using the same method.

## RNA secondary structure prediction

We predicted the secondary structures of all viroid-like cccRNAs for both polarities using the ViennaRNA package (Lorenz et al., 2011). For each predicted structure, we computed the percentage of bases paired and the number of hairpins present. We used a temperature of 25° C and the circular prediction mode.

## Clustering

We performed several types of clustering including both alignment-based and alignment-free methods. To produce the alignment-based clustering, we first performed an all-versus-all search using MMseqs2 on the sequences of ViroidDB, the HPR dataset of Weinberg et al. (2021), and this study. For this method, each sequence was first concatenated to itself to compensate for potential variation in the sequence relative to otherwise-similar sequences due to their circular nature. Next, we executed MMseqs2 (v13.45111, `easy-search -s 7.5 --min-seq-id 0.40 --search-type 3 -e 0.001 -k 5 --max-seqs 1000000`) and computed the pairwise average nucleotide identity (ANI) between sequences by taking the alignment identity of the best hit for each pair. We computed the ANI for two self-concatenated



885 sequences by first taking the length of the smaller sequence and dividing  
886 by two. We then cap the computed alignment length the the length of  
887 the now-monomerized smaller sequence. The ANI is then defined as  
888 the percent identity within the aligned region times the alignment length  
889 divided by the smaller sequence monomer length. Similarly, we defined the  
890 alignment fraction as the smaller of the doubled query coverage, doubled  
891 target coverage, or one.

892 To cluster the viroids based on their pairwise ANI, we first build a graph  
893 by connecting pairs of sequences where the alignment covers at least 25%  
894 of the shorter sequence with 40% identity within the alignment. We then  
895 weighted the connections between the sequences by the ANI and employed  
896 the Leiden algorithm (as implemented in the igraph Python library, version  
897 0.9.10) to delineate communities of similar sequences (Traag et al., 2019).  
898 The clustering granularity was optimized by iterating over the resolution  
899 parameter space until the difference between average intra-cluster ANI and  
900 the target ANI began to increase.

## 901 **ORF prediction**

902 To find ORFs present within the sequences, we used orfipy (Singh and  
903 Wurtele, 2021) configured to operate on circular sequences. Specifically, we  
904 searched sequences concatenated to themselves to ensure ORFs spanning  
905 the origin were detected. Only unique ORFs longer than 100 amino acids  
906 and using the standard genetic code were considered for each cccRNA.

## 907 **Protein searches**

908 We searched all viroid-like cccRNAs for matches to known proteins. The  
909 primary search method we used was by performing translated searches  
910 (BLASTX-style) against the UniRef90 protein database (Suzek et al., 2015)  
911 using MMseqs2 (Steinegger and Söding, 2017). For each cccRNA, we  
912 considered only the best match by E-value.

913 As a second approach, we also searched the ORFs from all cccRNAs, viroid-  
914 like or not, using HMMER (Eddy, 2011). We searched both the full Pfam-A  
915 profile database using hmmscan as well as a curated subset (the profiles for  
916 RdRP clan combined with the HDVAg profile) using hmmsearch.

## 917 **HDV antigen analysis**

918 Sequences were clustered using CLANS with BLASTP option (BLOSUM62  
919 matrix, E-value cutoff of 1e-03) (Frickey and Lupas, 2004). Sequence similar-  
920 ity among reference HDVAg from GenBank and those from metatranscrip-  
921 tomic datasets was analyzed with the Sequence Demarcation Tool (Muhire  
922 et al., 2014). For phylogenetic analysis, HDVAg-like sequences were aligned  
923 using PROMALS3D (Pei and Grishin, 2014). Due to the short length of the  
924 sequences, the alignment was not further processed. Maximum likelihood  
925 phylogenetic analysis was performed using IQ-TREE (Nguyen et al., 2015).  
926 The best fitting model was selected by IQ-TREE and was VT+F+R4. The  
927 tree was visualized with iTOL (Letunic and Bork, 2021).

## 928 **Read mapping**

929 We used bowtie2 to perform read mapping from the 1KP transcriptomes  
930 to the entire viroid-like cccRNA data set in parallel. We configured bowtie  
931 to use its most sensitive setting (`--very-sensitive`) and ignore unaligned  
932 reads.

## 933 **CRISPR spacer analysis**

934 Viroid-like sequences were compared to predicted CRISPR spacers from  
935 prokaryotic (meta)genomes to identify potential cases of spacer acquisition  
936 from, and possible defense against, viroids by prokaryotes. The full set  
937 of 22,109 viroid and viroid-like sequences, including all reference and  
938 novel sequences, was compared to 1,961,109 CRISPR spacers predicted  
939 from whole genomes of bacteria and archaea (v.June2022) and 61,658,467  
940 CRISPR spacers predicted from metagenomes in the IMG database (Chen  
941 et al., 2021) using BLASTN v2.9.0 with options `-dust no -word_size 7`.  
942 To minimize the number of false-positive hits due to low-complexity and/or  
943 repeat sequences, CRISPR spacers were excluded from this analysis if (i)  
944 they were encoded in a predicted CRISPR array including 2 spacers or less,  
945 (ii) less than 66% of the predicted repeats were 100% identical to each other,  
946 (iii) the spacers were at most 20 bp, or (iv) they included a low-complexity  
947 or repeat sequence as detected by dustmasker (v1.0.0) (Morgulis et al., 2006)  
948 (options `-window 20 -level 10`) or a direct repeat of at least 4 bp detected  
949 with etandem (Rice et al., 2000) (options `-minrepeat 4 -maxrepeat 15`

950 -threshold 2). To initially link viroid-like sequences to CRISPR spacers,  
951 only hits with 0 or 1 mismatch over the entire spacer were considered  
952 (Table S5). To find additional spacer matches, we searched all members  
953 of the clusters with a spacer match against the IMG public metagenomic  
954 spacer data (dated 2022-06-18) set using IMG's workspace BLAST with a  
955 minimum E-value of 1e-05. We also extracted the repeats matching loci  
956 using MinCED (Bland et al., 2007) and searched them against nt using  
957 BLASTN v2.13.0 (Altschul et al., 1990; Camacho et al., 2009).

958 **Key Resources**

959 All the code and data generated in this work are freely available through  
960 public portals listed below.

---

REAGENT or		
RESOURCE	SOURCE	IDENTIFIER

---

**Deposited data**

All original code and data for this paper	This paper	<a href="https://doi.org/10.5281/zenodo.6859104">https://doi.org/10.5281/zenodo.6859104</a>
The viroid-like cccRNA detection pipeline	This paper	<a href="https://github.com/Benjamin-Lee/vdsearch">https://github.com/Benjamin-Lee/vdsearch</a>

**Software and  
algorithms**

REAGENT or		
RESOURCE	SOURCE	IDENTIFIER
CD-HIT v4.8.1	Li and Godzik (2006)	<a href="https://github.com/weizhongli/cdhit">https://github.com/weizhongli/cdhit</a>
HMMER v3.3.2	Eddy (2011)	<a href="https://github.com/EddyRivasLab/hmmer">https://github.com/EddyRivasLab/hmmer</a>
Infernal v1.1.4	Nawrocki and Eddy (2013)	<a href="https://github.com/EddyRivasLab/infernal">https://github.com/EddyRivasLab/infernal</a>
python-igraph v0.9.10	Gábor and Nepusz (2005)	<a href="https://github.com/igraph/python-igraph">https://github.com/igraph/python-igraph</a>
scikit-bio v0.5.6	Caporaso et al. (2010); Knight et al. (2007)	<a href="https://github.com/biocore/scikit-bio">https://github.com/biocore/scikit-bio</a>
MMseqs2 v1.3.45111	Steinegger and Söding (2017)	<a href="https://github.com/soedinglab/mmseqs2">https://github.com/soedinglab/mmseqs2</a>
R v4.2.0	R Foundation for Statistical Computing	<a href="https://www.r-project.org">https://www.r-project.org</a>
Python v3.8.3	Python Software Foundation	<a href="https://www.python.org">https://www.python.org</a>
Nim v1.6.2	Rumpf (2022)	<a href="https://nim-lang.org">https://nim-lang.org</a>
ViennaRNA v2.5.0	Lorenz et al. (2011)	<a href="https://github.com/ViennaRNA/ViennaRNA">https://github.com/ViennaRNA/ViennaRNA</a>
SeqKit v2.1.0	Shen et al. (2016)	<a href="https://github.com/shenwei356/seqkit">https://github.com/shenwei356/seqkit</a>

REAGENT or RESOURCE	SOURCE	IDENTIFIER
Pandas v1.2.0	McKinney (2010)	<a href="https://github.com/pandas-dev/pandas/">https://github.com/pandas-dev/pandas/</a>
ggplot2 v3.3.6	Wickham (2016)	<a href="https://github.com/tidyverse/ggplot2/">https://github.com/tidyverse/ggplot2/</a>
orfipy v0.0.4	Singh and Wurtele (2021)	<a href="https://github.com/urmi-21/orfipy">https://github.com/urmi-21/orfipy</a>
fastp v0.20.1	Chen et al. (2018)	<a href="https://github.com/OpenGene/fastp">https://github.com/OpenGene/fastp</a>
Snakemake v6.10.0	Mölder et al. (2021)	<a href="https://github.com/snakemake/snakemake">https://github.com/snakemake/snakemake</a>
circlize v0.4.15	Gu et al. (2014)	<a href="https://github.com/jokergoo/circlize">https://github.com/jokergoo/circlize</a>
RNAmotif v3.1.1	Macke (2001)	<a href="https://github.com/dacase/rnamotif">https://github.com/dacase/rnamotif</a>
rnaSPAdes v3.14.1	Bushmanova et al. (2019)	<a href="https://github.com/ablab/spades">https://github.com/ablab/spades</a>
MinCED v0.4.2	Bland et al. (2007)	<a href="https://github.com/ctSkennerton/minced">https://github.com/ctSkennerton/minced</a>
IQ-TREE v1.6.12	Nguyen et al. (2015)	<a href="https://github.com/iqtree/iqtree2">https://github.com/iqtree/iqtree2</a>
CLANS	Frickey and Lupas (2004)	<a href="http://protevo.eb.tuebingen.mpg.de/download">http://protevo.eb.tuebingen.mpg.de/download</a>
Sequence Demarcation Tool (SDT) v1.2	Muhire et al. (2014)	<a href="https://github.com/brejnev/SDTv1.2">https://github.com/brejnev/SDTv1.2</a>
PROMALS3D	Pei and Grishin (2014)	<a href="http://prodata.swmed.edu/PROMALS3D">http://prodata.swmed.edu/PROMALS3D</a>

REAGENT or		
RESOURCE	SOURCE	IDENTIFIER
iTOL v6.5.8	Letunic and Bork (2021)	<a href="https://itol.embl.de">https://itol.embl.de</a>
bowtie2 v2.4.2	Langmead and Salzberg (2012)	<a href="http://bowtie-bio.sourceforge.net/bowtie2/index.shtml">http://bowtie- bio.sourceforge.net/bowtie2/index.shtml</a>
BLAST+ suite	Altschul et al. (1990); Camacho et al. (2009).	<a href="https://blast.ncbi.nlm.nih.gov/">https://blast.ncbi.nlm.nih.gov/</a>
EMBOSS etandem v6.0.0	Rice et al. (2000)	<a href="http://emboss.sourceforge.net">http://emboss.sourceforge.net</a>

## 961 The RNA Virus Discovery Consortium

962 Adrienne B. Narrowe, Alexander J. Probst, Alexander Sczyrba, Annegret  
 963 Kohler, Armand Séguin, Ashley Shade, Barbara J. Campbell, Björn D. Lin-  
 964 dahl, Brandi Kiel Reese, Breanna M. Roque, Chris DeRito, Colin Averill,  
 965 Daniel Cullen, David A. C. Beck, David A. Walsh, David M. Ward, Donald  
 966 A. Bryant, Dongying Wu, Emiley Eloë-Fadrosh, Eoin L. Brodie, Erica B.  
 967 Young, Erik A. Lilleskov, Federico J. Castillo, Francis M. Martin, Gary R.  
 968 LeCleir, Graeme T. Attwood , Hinsby Cadillo-Quiroz, Holly M. Simon, Ian  
 969 Hewson, Igor V. Grigoriev, James M. Tiedje, Janet K. Jansson, Janey Lee,

970 Jean S. VanderGheynst, Jeff Dangel, Jeff S. Bowman, Jeffrey L. Blanchard,  
 971 Jennifer L. Bowen, Jiangbing Xu, Jillian F. Banfield, Jody W Deming, Joel  
 972 E. Kostka, John M. Gladden, Josephine Z Rapp, Joshua Sharpe, Katherine  
 973 D. McMahon, Kathleen K. Treseder, Kay D. Bidle, Kelly C. Wrighton, Kim-  
 974 berlee Thamatrakoln, Klaus Nusslein, Laura K. Meredith, Lucia Ramirez,  
 975 Marc Buee, Marcel Huntemann, Marina G. Kalyuzhnaya, Mark P Waldrop,  
 976 Matthew B Sullivan, Matthew O. Schrenk, Matthias Hess, Michael A. Vega,  
 977 Michelle A. O'Malley, Monica Medina, Naomi E. Gilbert, Nathalie Delherbe,  
 978 Olivia U. Mason, Paul Dijkstra, Peter F. Chuckran, Petr Baldrian, Philippe  
 979 Constant, Ramunas Stepanauskas, Rebecca A. Daly, Regina Lamendella,  
 980 Robert J Gruninger, Robert M. McKay, Samuel Hylander, Sarah L. Lebeis,  
 981 Sarah P Esser, Silvia G. Acinas, Steven S. Wilhelm, Steven W. Singer, Su-  
 982 sannah S. Tringe, Tanja Woyke, TBK Reddy, Terrence H. Bell, Thomas  
 983 Mock, Tim McAllister, Vera Thiel, Vincent J. Denef, Wen-Tso Liu, Willm  
 984 Martens-Habbena, Xiao-Jun Allen Liu, Zachary S. Cooper, Zhong Wang

## 985 **Acknowledgements**

986 The authors would like to thank Samuel Wilder for his support while devel-  
 987 oping the software pipeline and Caleb Oh for advice on software architecture.  
 988 This work utilized the computational resources of the NIH HPC Biowulf clus-  
 989 ter (<http://hpc.nih.gov>). Figure 1 was created with BioRender.com. B.D.L.  
 990 was supported by a fellowship from the National Institutes of Health Oxford-



991 Cambridge Scholars Program. Y.I.W. and E.V.K. are supported through  
 992 the Intramural Research Program of the US National Institutes of Health  
 993 (National Library of Medicine). U.G. and U.N. are supported by the Euro-  
 994 pean Research Council (ERC-AdG 787514). U.N. is partially supported by  
 995 a fellowship from the Edmond J. Safra Center for Bioinformatics at Tel Aviv  
 996 University. V.V.D. was partially supported by NIH/NLM/NCBI Visiting Sci-  
 997 entist Fellowship. The work of the U.S. Department of Energy Joint Genome  
 998 Institute (S.R., N.K. and all JGI co-authors), a DOE Office of Science User  
 999 Facility, is supported by the Office of Science of the U.S. Department of  
 1000 Energy under contract no. DE-AC02-05CH11231. M.K. was supported by  
 1001 l'Agence Nationale de la Recherche grants ANR-20-CE20-009-02 and ANR-  
 1002 21-CE11-0001-01.

## 1003 **Author Contributions**

1004 E.V.K. incepted the project; B.D.L. and E.V.K. designed research; B.D.L.  
 1005 and U.N. compiled the datasets; B.D.L., U.N., S.R., A.P.C., Y.I.W and M.K.  
 1006 analyzed the data; P.S., U.G., N.K., V.V.D. and E.V.K. supervised research;  
 1007 BDL and E.V.K. wrote the manuscript that was read, edited and approved  
 1008 by all authors.

# Supplementary Information

- Table S1: Summary of each viroid-like cccRNA from plant transcriptomes, metatranscriptomes, ViroidDB, and Weinberg et al. (2021)
- Table S2: Predicted proteins in cccRNAs
- Table S3: Predicted self-cleaving ribozymes in known ambi- and ambi-like viruses
- Table S4: Sample preparation method and taxonomic composition for each metatranscriptome
- Table S5: CRISPR spacer matches to cccRNAs

# References

- Altschul, S.F., Gish, W., Miller, W., Myers, E.W., and Lipman, D.J. (1990). Basic local alignment search tool. *J Mol Biol* *215*, 403–410. [https://doi.org/10.1016/S0022-2836\(05\)80360-2](https://doi.org/10.1016/S0022-2836(05)80360-2).
- Bergner, L.M., Orton, R.J., Broos, A., Tello, C., Becker, D.J., Carrera, J.E., Patel, A.H., Biek, R., and Streicker, D.G. (2021). Diversification of mammalian deltaviruses by host shifting. *Proc Natl Acad Sci USA* *118*, e2019907118. <https://doi.org/10.1073/pnas.2019907118>.
- Bland, C., Ramsey, T.L., Sabree, F., Lowe, M., Brown, K., Kyrpides, N.C., and Hugenholtz, P. (2007). CRISPR Recognition Tool (CRT): A tool for automatic detection of clustered regularly interspaced palindromic repeats. *BMC Bioinformatics* *8*, 209. <https://doi.org/10.1186/1471-2105-8-209>.

1030 Branch, A.D., and Robertson, H.D. (1984). A replication cycle for viroids  
1031 and other small infectious RNA's. *Science* *223*, 450–455. [https://doi.org/10](https://doi.org/10.1126/science.6197756)  
1032 [.1126/science.6197756](https://doi.org/10.1126/science.6197756).

1033 Branch, A.D., Benenfeld, B.J., and Robertson, H.D. (1988). Evidence for  
1034 a single rolling circle in the replication of potato spindle tuber viroid. *Proc*  
1035 *Natl Acad Sci USA* *85*, 9128–9132. <https://doi.org/10.1073/pnas.85.23.9128>.

1036 Bruening, G., Passmore, B.K., van Tol, H., Buzayan, J.M., and Feldstein,  
1037 P.A. (1991). Replication of a Plant Virus Satellite RNA: Evidence Favors  
1038 Transcription of Circular Templates of Both Polarities. *MPMI* *4*, 219–225.  
1039 <https://doi.org/10.1094/MPMI-4-219>.

1040 Bushmanova, E., Antipov, D., Lapidus, A., and Prjibelski, A.D. (2019). rnaS-  
1041 PAdes: A de novo transcriptome assembler and its application to RNA-Seq  
1042 data. *GigaScience* *8*, giz100. <https://doi.org/10.1093/gigascience/giz100>.

1043 Camacho, C., Coulouris, G., Avagyan, V., Ma, N., Papadopoulos, J., Bealer,  
1044 K., and Madden, T.L. (2009). BLAST+: Architecture and applications.  
1045 *BMC Bioinformatics* *10*, 421. <https://doi.org/10.1186/1471-2105-10-421>.

1046 Caporaso, J.G., Kuczynski, J., Stombaugh, J., Bittinger, K., Bushman, F.D.,  
1047 Costello, E.K., Fierer, N., Peña, A.G., Goodrich, J.K., Gordon, J.I., et al.  
1048 (2010). QIIME allows analysis of high-throughput community sequencing  
1049 data. *Nat Methods* *7*, 335–336. <https://doi.org/10.1038/nmeth.f.303>.

1050 Cervera, A., Urbina, D., and de la Peña, M. (2016). Retrozymes are a unique  
1051 family of non-autonomous retrotransposons with hammerhead ribozymes  
1052 that propagate in plants through circular RNAs. *Genome Biol* *17*, 135.

1053 <https://doi.org/10.1186/s13059-016-1002-4>.

1054 Chadalavada, D.M., Gratton, E.A., and Bevilacqua, P.C. (2010). The human  
1055 HDV-like CPEB3 ribozyme is intrinsically fast-reacting. *Biochemistry* *49*,  
1056 5321–5330. <https://doi.org/10.1021/bi100434c>.

1057 Chang, W.-S., Pettersson, J.H.-O., Le Lay, C., Shi, M., Lo, N., Wille, M.,  
1058 Eden, J.-S., and Holmes, E.C. (2019). Novel hepatitis D-like agents in verte-  
1059 brates and invertebrates. *Virus Evol* *5*, vez021. <https://doi.org/10.1093/ve>  
1060 [/vez021](https://doi.org/10.1093/ve/vez021).

1061 Chen, I.-M.A., Chu, K., Palaniappan, K., Ratner, A., Huang, J., Hunte-  
1062 mann, M., Hajek, P., Ritter, S., Varghese, N., Seshadri, R., et al. (2021).  
1063 The IMG/M data management and analysis system v.6.0: New tools and  
1064 advanced capabilities. *Nucleic Acids Research* *49*, D751–D763. [https://doi.](https://doi.org/10.1093/nar/gkaa939)  
1065 [org/10.1093/nar/gkaa939](https://doi.org/10.1093/nar/gkaa939).

1066 Chen, S., Zhou, Y., Chen, Y., and Gu, J. (2018). Fastp: An ultra-fast all-in-  
1067 one FASTQ preprocessor. *Bioinformatics* *34*, i884–i890. <https://doi.org/10>  
1068 [.1093/bioinformatics/bty560](https://doi.org/10.1093/bioinformatics/bty560).

1069 Chiumenti, M., Navarro, B., Candresse, T., Flores, R., and Di Serio, F.  
1070 (2021). Reassessing species demarcation criteria in viroid taxonomy by pair-  
1071 wise identity matrices. *Virus Evol* *7*, veab001. <https://doi.org/10.1093/ve>  
1072 [/veab001](https://doi.org/10.1093/ve/veab001).

1073 Daròs, J.-A., and Flores, R. (1995). Identification of a retroviroid-like ele-  
1074 ment from plants. *Proc Natl Acad Sci USA* *92*, 6856–6860. [https://doi.org/](https://doi.org/10.1073/pnas.92.15.6856)  
1075 [10.1073/pnas.92.15.6856](https://doi.org/10.1073/pnas.92.15.6856).

1076 Daròs, J.-A., Elena, S.F., and Flores, R. (2006). Viroids: An Ariadne's  
1077 thread into the RNA labyrinth. *EMBO Rep* 7, 593–598. [https://doi.org/10](https://doi.org/10.1038/sj.embor.7400706)  
1078 .1038/sj.embor.7400706.

1079 de la Peña, M., and Cervera, A. (2017). Circular RNAs with hammerhead ri-  
1080 bozymes encoded in eukaryotic genomes: The enemy at home. *RNA Biology*  
1081 14, 985–991. <https://doi.org/10.1080/15476286.2017.1321730>.

1082 de la Peña, M., Ceprián, R., and Cervera, A. (2020). A Singular and  
1083 Widespread Group of Mobile Genetic Elements: RNA Circles with Auto-  
1084 catalytic Ribozymes. *Cells* 9, E2555. <https://doi.org/10.3390/cells9122555>.

1085 de la Peña, M., Ceprián, R., Casey, J.L., and Cervera, A. (2021). Hepati-  
1086 tis delta virus-like circular RNAs from diverse metazoans encode conserved  
1087 hammerhead ribozymes. *Virus Evol* 7, veab016. <https://doi.org/10.1093/ve>  
1088 /veab016.

1089 Di Serio, F., Li, S.-F., Pallás, V., Owens, R.A., Randles, J.W., Sano, T.,  
1090 Verhoeven, J.Th.J., Vidalakis, G., and Flores, R. (2017). Viroid Taxonomy.  
1091 In *Viroids and Satellites*, (Elsevier), pp. 135–146.

1092 Diener, T.O. (1971). Potato spindle tuber “virus.” *Virology* 45, 411–428.  
1093 [https://doi.org/10.1016/0042-6822\(71\)90342-4](https://doi.org/10.1016/0042-6822(71)90342-4).

1094 Diener, T.O. (2001). The viroid: Biological oddity or evolutionary fossil? In  
1095 *Advances in Virus Research*, (Elsevier), pp. 137–184.

1096 Diener, T.O. (2016). Viroids: “Living fossils” of primordial RNAs? *Biol*  
1097 *Direct* 11, 15. <https://doi.org/10.1186/s13062-016-0116-7>.

1098 Duval, J.P. (1983). Factorizing words over an ordered alphabet. *Journal of*

1099 Algorithms 4, 363–381. [https://doi.org/10.1016/0196-6774\(83\)90017-2](https://doi.org/10.1016/0196-6774(83)90017-2).

1100 Eddy, S.R. (2011). Accelerated Profile HMM Searches. PLoS Comput Biol

1101 7, e1002195. <https://doi.org/10.1371/journal.pcbi.1002195>.

1102 Edgar, R.C., Taylor, J., Lin, V., Altman, T., Barbera, P., Meleshko, D., Lohr,

1103 D., Novakovsky, G., Buchfink, B., Al-Shayeb, B., et al. (2022). Petabase-

1104 scale sequence alignment catalyses viral discovery. Nature 602, 142–147. <https://doi.org/10.1038/s41586-021-04332-2>.

1105

1106 Ferré-D’Amaré, A.R., and Scott, W.G. (2010). Small self-cleaving ribozymes.

1107 Cold Spring Harb Perspect Biol 2, a003574. <https://doi.org/10.1101/cshperspect.a003574>.

1108

1109 Flores, R., Minoia, S., López-Carrasco, A., Delgado, S., Martínez de Alba, Á.-

1110 E., and Kalantidis, K. (2017). Viroid Replication. In Viroids and Satellites,

1111 (Elsevier), pp. 71–81.

1112 Flores, R., Navarro, B., Serra, P., and Di Serio, F. (2022). A scenario for the

1113 emergence of protoviroids in the RNA world and for their further evolution

1114 into viroids and viroid-like RNAs by modular recombinations and mutations.

1115 Virus Evolution 8, veab107. <https://doi.org/10.1093/ve/veab107>.

1116 Forgia, M., Isgandarli, E., Aghayeva, D.N., Huseynova, I., and Turina, M.

1117 (2021). Virome characterization of Cryphonectria parasitica isolates from

1118 Azerbaijan unveiled a new myomonavirus and a putative new RNA virus un-

1119 related to described viral sequences. Virology 553, 51–61. <https://doi.org/10.1016/j.virol.2020.10.008>.

1120

1121 Frickey, T., and Lupas, A. (2004). CLANS: A Java application for visualizing

1122 protein families based on pairwise similarity. *Bioinformatics* *20*, 3702–3704.  
1123 <https://doi.org/10.1093/bioinformatics/bth444>.

1124 Gábor, C., and Nepusz, T. (2005). The igraph software package for complex  
1125 network research. *InterJournal Complex Systems* 1695.

1126 Gago, S., Elena, S.F., Flores, R., and Sanjuán, R. (2009). Extremely High  
1127 Mutation Rate of a Hammerhead Viroid. *Science* *323*, 1308–1308. <https://doi.org/10.1126/science.1169202>.  
1128

1129 Gao, Y., Wang, J., and Zhao, F. (2015). CIRI: An efficient and unbiased  
1130 algorithm for de novo circular RNA identification. *Genome Biol* *16*, 4. <https://doi.org/10.1186/s13059-014-0571-3>.  
1131

1132 Giguère, T., Adkar-Purushothama, C.R., and Perreault, J.-P. (2014b).  
1133 Comprehensive secondary structure elucidation of four genera of the family  
1134 Pospiviroidae. *PLoS ONE* *9*, e98655. [https://doi.org/10.1371/journal.pone](https://doi.org/10.1371/journal.pone.0098655)  
1135 [.0098655](https://doi.org/10.1371/journal.pone.0098655).

1136 Giguère, T., Adkar-Purushothama, C.R., Bolduc, F., and Perreault, J.-P.  
1137 (2014a). Elucidation of the structures of all members of the Avsunviroidae  
1138 family. *Molecular Plant Pathology* *15*, 767–779. [https://doi.org/10.1111/](https://doi.org/10.1111/mpp.12130)  
1139 [mpp.12130](https://doi.org/10.1111/mpp.12130).

1140 Gozmanova, M. (2003). Characterization of the RNA motif responsible for  
1141 the specific interaction of potato spindle tuber viroid RNA (PSTVd) and  
1142 the tomato protein Virp1. *Nucleic Acids Research* *31*, 5534–5543. <https://doi.org/10.1093/nar/gkg777>.  
1143

1144 Gu, Z., Gu, L., Eils, R., Schlesner, M., and Brors, B. (2014). Circlize imple-

1145 ments and enhances circular visualization in R. *Bioinformatics* *30*, 2811–2812.  
1146 <https://doi.org/10.1093/bioinformatics/btu393>.

1147 Hegedus, K., Dallmann, G., and Balázs, E. (2004). The DNA form of a  
1148 retroviroid-like element is involved in recombination events with itself and  
1149 with the plant genome. *Virology* *325*, 277–286. <https://doi.org/10.1016/j.virol.2004.04.035>.

1151 Hepojoki, J., Hetzel, U., Paraskevopoulou, S., Drosten, C., Harrach, B.,  
1152 Zerbini, F.M., Koonin, E.V., Krupovic, M., Dolja, V.V., and Kuhn, J.H.  
1153 (2020). Create one new realm (Ribozviria) including one new family  
1154 (Kolmioviridae) including genus Deltavirus and seven new genera for a total  
1155 of 15 species (International Committee on Taxonomy of Viruses).

1156 Hetzel, U., Szilovics, L., Smura, T., Prähauser, B., Vapalahti, O., Kipar,  
1157 A., and Hepojoki, J. (2019). Identification of a Novel Deltavirus in Boa  
1158 Constrictors. *mBio* *10*, e00014–19. <https://doi.org/10.1128/mBio.00014-19>.

1159 Hou, W.-Y., Li, S.-F., Wu, Z.-J., Jiang, D.-M., and Sano, T. (2009). *Coleus*  
1160 *blumei* viroid 6: A new tentative member of the genus Coleviroid derived  
1161 from natural genome shuffling. *Arch Virol* *154*, 993–997. <https://doi.org/10.1007/s00705-009-0388-7>.

1163 Huang, Y.-W., Hu, C.-C., Hsu, Y.-H., and Lin, N.-S. (2017). Replication of  
1164 Satellites. In *Viroids and Satellites*, (Elsevier), pp. 577–586.

1165 Jimenez, R.M., Delwart, E., and Lupták, A. (2011). Structure-based Search  
1166 Reveals Hammerhead Ribozymes in the Human Microbiome. *Journal of*  
1167 *Biological Chemistry* *286*, 7737–7743. <https://doi.org/10.1074/jbc.C110.209>



1168 288.

1169 Johnson, M.T.J., Carpenter, E.J., Tian, Z., Bruskewich, R., Burris, J.N.,  
1170 Carrigan, C.T., Chase, M.W., Clarke, N.D., Covshoff, S., dePamphilis, C.W.,  
1171 et al. (2012). Evaluating Methods for Isolating Total RNA and Predicting  
1172 the Success of Sequencing Phylogenetically Diverse Plant Transcriptomes.  
1173 PLoS ONE 7, e50226. <https://doi.org/10.1371/journal.pone.0050226>.

1174 Kalvari, I., Nawrocki, E.P., Ontiveros-Palacios, N., Argasinska, J.,  
1175 Lamkiewicz, K., Marz, M., Griffiths-Jones, S., Toffano-Nioche, C., Gau-  
1176 theret, D., Weinberg, Z., et al. (2021). Rfam 14: Expanded coverage of  
1177 metagenomic, viral and microRNA families. Nucleic Acids Research 49,  
1178 D192–D200. <https://doi.org/10.1093/nar/gkaa1047>.

1179 Knight, R., Maxwell, P., Birmingham, A., Carnes, J., Caporaso, J.G., Easton,  
1180 B.C., Eaton, M., Hamady, M., Lindsay, H., Liu, Z., et al. (2007). PyCogent:  
1181 A toolkit for making sense from sequence. Genome Biol 8, R171. <https://doi.org/10.1186/gb-2007-8-8-r171>.

1183 Koonin, E.V., Dolja, V.V., Krupovic, M., Varsani, A., Wolf, Y.I., Yutin, N.,  
1184 Zerbini, F.M., and Kuhn, J.H. (2020). Global Organization and Proposed  
1185 Megataxonomy of the Virus World. Microbiol Mol Biol Rev 84, e00061–19.  
1186 <https://doi.org/10.1128/MMBR.00061-19>.

1187 Kos, A., Dijkema, R., Arnberg, A.C., van der Meide, P.H., and Schellekens,  
1188 H. (1986). The hepatitis delta ( $\delta$ ) virus possesses a circular RNA. Nature  
1189 323, 558–560. <https://doi.org/10.1038/323558a0>.

1190 Krupovic, M., Varsani, A., Kazlauskas, D., Breitbart, M., Delwart, E.,

1191 Rosario, K., Yutin, N., Wolf, Y.I., Harrach, B., Zerbini, F.M., et al.  
1192 (2020). Cressdnaviricota: A Virus Phylum Unifying Seven Families of  
1193 Rep-Encoding Viruses with Single-Stranded, Circular DNA Genomes. *J*  
1194 *Virol* *94*, e00582–20. <https://doi.org/10.1128/JVI.00582-20>.  
1195 Langmead, B., and Salzberg, S.L. (2012). Fast gapped-read alignment with  
1196 Bowtie 2. *Nat Methods* *9*, 357–359. <https://doi.org/10.1038/nmeth.1923>.  
1197 Lee, B.D., and Koonin, E.V. (2022). Viroids and Viroid-like Circular RNAs:  
1198 Do They Descend from Primordial Replicators? *Life* *12*, 103. <https://doi.org/10.3390/life12010103>.  
1200 Lee, B.D., Neri, U., Oh, C.J., Simmonds, P., and Koonin, E.V. (2021). Vi-  
1201 roidDB: A database of viroids and viroid-like circular RNAs. *Nucleic Acids*  
1202 *Research* gkab974. <https://doi.org/10.1093/nar/gkab974>.  
1203 Letunic, I., and Bork, P. (2021). Interactive Tree Of Life (iTOL) v5: An  
1204 online tool for phylogenetic tree display and annotation. *Nucleic Acids Re-*  
1205 *search* *49*, W293–W296. <https://doi.org/10.1093/nar/gkab301>.  
1206 Li, W., and Godzik, A. (2006). Cd-hit: A fast program for clustering and  
1207 comparing large sets of protein or nucleotide sequences. *Bioinformatics* *22*,  
1208 1658–1659. <https://doi.org/10.1093/bioinformatics/btl158>.  
1209 Linnakoski, R., Sutela, S., Coetzee, M.P.A., Duong, T.A., Pavlov, I.N.,  
1210 Litovka, Y.A., Hantula, J., Wingfield, B.D., and Vainio, E.J. (2021). Armil-  
1211 laria root rot fungi host single-stranded RNA viruses. *Sci Rep* *11*, 7336.  
1212 <https://doi.org/10.1038/s41598-021-86343-7>.  
1213 López-Carrasco, A., Ballesteros, C., Sentandreu, V., Delgado, S., Gago-

1214 Zachert, S., Flores, R., and Sanjuán, R. (2017). Different rates of spon-  
1215 taneous mutation of chloroplastic and nuclear viroids as determined by high-  
1216 fidelity ultra-deep sequencing. *PLoS Pathog* 13, e1006547. [https://doi.org/](https://doi.org/10.1371/journal.ppat.1006547)  
1217 10.1371/journal.ppat.1006547.

1218 Lorenz, R., Bernhart, S.H., Höner zu Siederdissen, C., Tafer, H., Flamm, C.,  
1219 Stadler, P.F., and Hofacker, I.L. (2011). ViennaRNA Package 2.0. *Algo-*  
1220 *rithms Mol Biol* 6, 26. <https://doi.org/10.1186/1748-7188-6-26>.

1221 Macke, T.J. (2001). RNAMotif, an RNA secondary structure definition and  
1222 search algorithm. *Nucleic Acids Research* 29, 4724–4735. [https://doi.org/10](https://doi.org/10.1093/nar/29.22.4724)  
1223 .1093/nar/29.22.4724.

1224 Madigan, M.T., Jung, D.O., Karr, E.A., Sattley, W.M., Achenbach, L.A.,  
1225 and van der Meer, M.T.J. (2005). Diversity of anoxygenic phototrophs in  
1226 contrasting extreme environments. In *Geothermal Biology and Geochem-*  
1227 *istry in Yellowstone National Park*, W.P. Inskeep, and T.R. McDermott, eds.  
1228 (Bozeman, MT: Thermal Biology Institute), pp. 203–219.

1229 Marais, A., Nivault, A., Faure, C., Theil, S., Comont, G., Candresse, T., and  
1230 Corio-Costet, M.-F. (2017). Determination of the complete genomic sequence  
1231 of Neofusicoccum luteum mitovirus 1 (NLMV1), a novel mitovirus associated  
1232 with a phytopathogenic Botryosphaeriaceae. *Arch Virol* 162, 2477–2480. [ht](https://doi.org/10.1007/s00705-017-3338-9)  
1233 [tps://doi.org/10.1007/s00705-017-3338-9](https://doi.org/10.1007/s00705-017-3338-9).

1234 Marçais, G., and Kingsford, C. (2011). A fast, lock-free approach for efficient  
1235 parallel counting of occurrences of k-mers. *Bioinformatics* 27, 764–770. [https:](https://doi.org/10.1093/bioinformatics/btr011)  
1236 [//doi.org/10.1093/bioinformatics/btr011](https://doi.org/10.1093/bioinformatics/btr011).

1237 Mascia, T., and Gallitelli, D. (2017). Economic Significance of Satellites. In  
1238 Viroids and Satellites, (Elsevier), pp. 555–563.

1239 McKinney, W. (2010). Data Structures for Statistical Computing in Python.  
1240 (Austin, Texas), pp. 56–61.

1241 Mehle, N., Gutiérrez-Aguirre, I., Prezelj, N., Delić, D., Vidic, U., and  
1242 Ravnikar, M. (2014). Survival and Transmission of Potato Virus Y, Pepino  
1243 Mosaic Virus, and Potato Spindle Tuber Viroid in Water. *Appl Environ*  
1244 *Microbiol* *80*, 1455–1462. <https://doi.org/10.1128/AEM.03349-13>.

1245 Melsted, P., and Pritchard, J.K. (2011). Efficient counting of k-mers in  
1246 DNA sequences using a bloom filter. *BMC Bioinformatics* *12*, 333. <https://doi.org/10.1186/1471-2105-12-333>.

1247

1248 Mistry, J., Chuguransky, S., Williams, L., Qureshi, M., Salazar, G.A.,  
1249 Sonnhammer, E.L.L., Tosatto, S.C.E., Paladin, L., Raj, S., Richardson, L.J.,  
1250 et al. (2021). Pfam: The protein families database in 2021. *Nucleic Acids*  
1251 *Research* *49*, D412–D419. <https://doi.org/10.1093/nar/gkaa913>.

1252 Modahl, L.E., Macnaughton, T.B., Zhu, N., Johnson, D.L., and Lai, M.M.C.  
1253 (2000). RNA-Dependent Replication and Transcription of Hepatitis Delta  
1254 Virus RNA Involve Distinct Cellular RNA Polymerases. *Mol. Cell. Biol.* *20*,  
1255 6030–6039. <https://doi.org/10.1128/MCB.20.16.6030-6039.2000>.

1256 Mölder, F., Jablonski, K.P., Letcher, B., Hall, M.B., Tomkins-Tinch, C.H.,  
1257 Sochat, V., Forster, J., Lee, S., Twardziok, S.O., Kanitz, A., et al. (2021).  
1258 Sustainable data analysis with Snakemake. *F1000Res* *10*, 33. <https://doi.org/10.12688/f1000research.29032.1>.

1259

1260 Morgulis, A., Gertz, E.M., Schäffer, A.A., and Agarwala, R. (2006). A Fast  
1261 and Symmetric DUST Implementation to Mask Low-Complexity DNA Se-  
1262 quences. *Journal of Computational Biology* 13, 1028–1040. [https://doi.org/](https://doi.org/10.1089/cmb.2006.13.1028)  
1263 10.1089/cmb.2006.13.1028.

1264 Muhire, B.M., Varsani, A., and Martin, D.P. (2014). SDT: A Virus Classifi-  
1265 cation Tool Based on Pairwise Sequence Alignment and Identity Calculation.  
1266 PLoS ONE 9, e108277. <https://doi.org/10.1371/journal.pone.0108277>.

1267 Mühlbach, H.-P., and Sängler, H.L. (1979). Viroid replication is inhibited by  
1268  $\alpha$ -amanitin. *Nature* 278, 185–188. <https://doi.org/10.1038/278185a0>.

1269 Munson-McGee, J.H., Peng, S., Dewerff, S., Stepanauskas, R., Whitaker,  
1270 R.J., Weitz, J.S., and Young, M.J. (2018). A virus or more in (nearly) every  
1271 cell: Ubiquitous networks of virus–host interactions in extreme environments.  
1272 *ISME J* 12, 1706–1714. <https://doi.org/10.1038/s41396-018-0071-7>.

1273 Navarro, B., Rubino, L., and Di Serio, F. (2017). Small Circular Satellite  
1274 RNAs. In *Viroids and Satellites*, (Elsevier), pp. 659–669.

1275 Navarro, J.-A., Vera, A., and Flores, R. (2000). A Chloroplastic RNA Poly-  
1276 merase Resistant to Tagetitoxin Is Involved in Replication of Avocado Sun-  
1277 blotch Viroid. *Virology* 268, 218–225. <https://doi.org/10.1006/viro.1999.01>  
1278 61.

1279 Nawrocki, E.P., and Eddy, S.R. (2013). Infernal 1.1: 100-fold faster RNA  
1280 homology searches. *Bioinformatics* 29, 2933–2935. [https://doi.org/10.1093/](https://doi.org/10.1093/bioinformatics/btt509)  
1281 [bioinformatics/btt509](https://doi.org/10.1093/bioinformatics/btt509).

1282 Neri, U., Wolf, Y.I., Roux, S., Camargo, A.P., Lee, B., Kazlauskas, D., Chen,

1283 I.M., Ivanova, N., Allen, L.Z., Paez-Espino, D., et al. (2022). A five-fold  
1284 expansion of the global RNA virome reveals multiple new clades of RNA  
1285 bacteriophages (bioRxiv).

1286 Nguyen, L.-T., Schmidt, H.A., von Haeseler, A., and Minh, B.Q. (2015). IQ-  
1287 TREE: A Fast and Effective Stochastic Algorithm for Estimating Maximum-  
1288 Likelihood Phylogenies. *Molecular Biology and Evolution* 32, 268–274. <https://doi.org/10.1093/molbev/msu300>.  
1289

1290 Nie, X., and Singh, R.P. (2017). *Coleus Blumei* Viroids. In *Viroids and*  
1291 *Satellites*, (Elsevier), pp. 289–295.

1292 Nohales, M.-Á., Flores, R., and Daròs, J.-A. (2012b). Viroid RNA redirects  
1293 host DNA ligase 1 to act as an RNA ligase. *Proc Natl Acad Sci USA* 109,  
1294 13805–13810. <https://doi.org/10.1073/pnas.1206187109>.

1295 Nohales, M.-Á., Molina-Serrano, D., Flores, R., and Daròs, J.-A. (2012a).  
1296 Involvement of the Chloroplastic Isoform of tRNA Ligase in the Replication  
1297 of Viroids Belonging to the Family Avsunviroidae. *Journal of Virology* 86,  
1298 8269–8276. <https://doi.org/10.1128/JVI.00629-12>.

1299 One Thousand Plant Transcriptomes Initiative (2019). One thousand plant  
1300 transcriptomes and the phylogenomics of green plants. *Nature* 574, 679–685.  
1301 <https://doi.org/10.1038/s41586-019-1693-2>.

1302 Paez-Espino, D., Elie-Fadrosh, E.A., Pavlopoulos, G.A., Thomas, A.D.,  
1303 Huntemann, M., Mikhailova, N., Rubin, E., Ivanova, N.N., and Kyrpides,  
1304 N.C. (2016). Uncovering Earth’s virome. *Nature* 536, 425–430. <https://doi.org/10.1038/nature19094>.  
1305

1306 Paraskevopoulou, S., Pirzer, F., Goldmann, N., Schmid, J., Corman, V.M.,  
1307 Gottula, L.T., Schroeder, S., Rasche, A., Muth, D., Drexler, J.F., et al.  
1308 (2020). Mammalian deltavirus without hepadnavirus coinfection in the  
1309 neotropical rodent *Proechimys semispinosus*. *Proc Natl Acad Sci USA* *117*,  
1310 17977–17983. <https://doi.org/10.1073/pnas.2006750117>.  
1311 Pei, J., and Grishin, N.V. (2014). PROMALS3D: Multiple Protein Sequence  
1312 Alignment Enhanced with Evolutionary and Three-Dimensional Structural  
1313 Information. In *Multiple Sequence Alignment Methods*, D.J. Russell, ed.  
1314 (Totowa, NJ: Humana Press), pp. 263–271.  
1315 Qin, Y., Xu, T., Lin, W., Jia, Q., He, Q., Liu, K., Du, J., Chen, L., Yang, X.,  
1316 Du, F., et al. (2020). Reference-free and de novo Identification of Circular  
1317 RNAs (Bioinformatics).  
1318 Rao, A.L.N., and Kalantidis, K. (2015). Virus-associated small satellite  
1319 RNAs and viroids display similarities in their replication strategies. *Virology*  
1320 *479–480*, 627–636. <https://doi.org/10.1016/j.virol.2015.02.018>.  
1321 Rice, P., Longden, I., and Bleasby, A. (2000). EMBOSS: The European  
1322 Molecular Biology Open Software Suite. *Trends Genet* *16*, 276–277. [https://doi.org/10.1016/s0168-9525\(00\)00204-2](https://doi.org/10.1016/s0168-9525(00)00204-2).  
1323  
1324 Rumpf, A. (2022). Mastering Nim: A complete guide to the programming  
1325 language.  
1326 Salehi-Ashtiani, K., Lupták, A., Litovchick, A., and Szostak, J.W. (2006).  
1327 A Genomewide Search for Ribozymes Reveals an HDV-Like Sequence in the  
1328 Human CPEB3 Gene. *Science* *313*, 1788–1792. <https://doi.org/10.1126/sc>

ience.1129308.

Schindler, I.-M., and Mühlbach, H.-P. (1992). Involvement of nuclear DNA-dependent RNA polymerases in potato spindle tuber viroid replication: A reevaluation. *Plant Science* *84*, 221–229. [https://doi.org/10.1016/0168-9452\(92\)90138-C](https://doi.org/10.1016/0168-9452(92)90138-C).

Shen, W., Le, S., Li, Y., and Hu, F. (2016). SeqKit: A Cross-Platform and Ultrafast Toolkit for FASTA/Q File Manipulation. *PLoS ONE* *11*, e0163962. <https://doi.org/10.1371/journal.pone.0163962>.

Singh, U., and Wurtele, E.S. (2021). Orfipy: A fast and flexible tool for extracting ORFs. *Bioinformatics* *37*, 3019–3020. <https://doi.org/10.1093/bioinformatics/btab090>.

Steinegger, M., and Söding, J. (2017). MMseqs2 enables sensitive protein sequence searching for the analysis of massive data sets. *Nat Biotechnol* *35*, 1026–1028. <https://doi.org/10.1038/nbt.3988>.

Sureau, C., and Negro, F. (2016). The hepatitis delta virus: Replication and pathogenesis. *Journal of Hepatology* *64*, S102–S116. <https://doi.org/10.1016/j.jhep.2016.02.013>.

Sutela, S., Forgia, M., Vainio, E.J., Chiapello, M., Daghino, S., Vallino, M., Martino, E., Girlanda, M., Perotto, S., and Turina, M. (2020). The virome from a collection of endomycorrhizal fungi reveals new viral taxa with unprecedented genome organization. *Virus Evolution* *6*, veaa076. <https://doi.org/10.1093/ve/veaa076>.

Suzek, B.E., Wang, Y., Huang, H., McGarvey, P.B., Wu, C.H., and the



1352 UniProt Consortium (2015). UniRef clusters: A comprehensive and scalable  
1353 alternative for improving sequence similarity searches. *Bioinformatics* *31*,  
1354 926–932. <https://doi.org/10.1093/bioinformatics/btu739>.

1355 Thomas, S.C., Tamadonfar, K.O., Seymour, C.O., Lai, D., Dodsworth, J.A.,  
1356 Murugapiran, S.K., Eloë-Fadrosch, E.A., Dijkstra, P., and Hedlund, B.P.  
1357 (2019). Position-Specific Metabolic Probing and Metagenomics of Microbial  
1358 Communities Reveal Conserved Central Carbon Metabolic Network Activi-  
1359 ties at High Temperatures. *Front. Microbiol.* *10*, 1427. [https://doi.org/10](https://doi.org/10.3389/fmicb.2019.01427)  
1360 [.3389/fmicb.2019.01427](https://doi.org/10.3389/fmicb.2019.01427).

1361 Traag, V.A., Waltman, L., and van Eck, N.J. (2019). From Louvain to  
1362 Leiden: Guaranteeing well-connected communities. *Sci Rep* *9*, 5233. [https:](https://doi.org/10.1038/s41598-019-41695-z)  
1363 [//doi.org/10.1038/s41598-019-41695-z](https://doi.org/10.1038/s41598-019-41695-z).

1364 van der Meer, M.T.J., Klatt, C.G., Wood, J., Bryant, D.A., Bateson, M.M.,  
1365 Lammerts, L., Schouten, S., Sinninghe Damsté, J.S., Madigan, M.T., and  
1366 Ward, D.M. (2010). Cultivation and Genomic, Nutritional, and Lipid  
1367 Biomarker Characterization of *Roseiflexus* Strains Closely Related to Pre-  
1368 dominant In Situ Populations Inhabiting Yellowstone Hot Spring Microbial  
1369 Mats. *J Bacteriol* *192*, 3033–3042. <https://doi.org/10.1128/JB.01610-09>.

1370 Vera, A., Daròs, J.A., Flores, R., and Hernández, C. (2000). The DNA of  
1371 a plant retroviroid-like element is fused to different sites in the genome of  
1372 a plant pararetrovirus and shows multiple forms with sequence deletions. *J*  
1373 *Virol* *74*, 10390–10400. <https://doi.org/10.1128/jvi.74.22.10390-10400.2000>.

1374 Wang, Y. (2021). Current view and perspectives in viroid replication. *Curr*

Opin Virol 47, 32–37. <https://doi.org/10.1016/j.coviro.2020.12.004>.

Webb, C.-H.T., Riccitelli, N.J., Ruminski, D.J., and Lupták, A. (2009). Widespread occurrence of self-cleaving ribozymes. *Science* 326, 953. <https://doi.org/10.1126/science.1178084>.

Weinberg, C.E., Olzog, V.J., Eckert, I., and Weinberg, Z. (2021). Identification of over 200-fold more hairpin ribozymes than previously known in diverse circular RNAs. *Nucleic Acids Research* 49, 6375–6388. <https://doi.org/10.1093/nar/gkab454>.

Wickham, H. (2016). *Ggplot2* (Cham: Springer International Publishing).

Wille, M., Netter, H., Littlejohn, M., Yuen, L., Shi, M., Eden, J.-S., Klaassen, M., Holmes, E., and Hurt, A. (2018). A Divergent Hepatitis D-Like Agent in Birds. *Viruses* 10, 720. <https://doi.org/10.3390/v10120720>.

Wolf, Y.I., Silas, S., Wang, Y., Wu, S., Bocek, M., Kazlauskas, D., Krupovic, M., Fire, A., Dolja, V.V., and Koonin, E.V. (2020). Doubling of the known set of RNA viruses by metagenomic analysis of an aquatic virome. *Nat Microbiol* <https://doi.org/10.1038/s41564-020-0755-4>.

Wu, J., and Bisaro, D.M. (2020). Biased Pol II fidelity contributes to conservation of functional domains in the Potato spindle tuber viroid genome. *PLoS Pathog* 16, e1009144. <https://doi.org/10.1371/journal.ppat.1009144>.

Zayed, A.A., Wainaina, J.M., Dominguez-Huerta, G., Pelletier, E., Guo, J., Mohssen, M., Tian, F., Pratama, A.A., Bolduc, B., Zablocki, O., et al. (2022). Cryptic and abundant marine viruses at the evolutionary origins of Earth’s RNA virome. *Science* 376, 156–162. <https://doi.org/10.1126/science.abm5>

1398 847.

1399 Zuccola, H.J., Rozzelle, J.E., Lemon, S.M., Erickson, B.W., and Hogle, J.M.

1400 (1998). Structural basis of the oligomerization of hepatitis delta antigen.

1401 Structure 6, 821–830. [https://doi.org/10.1016/S0969-2126\(98\)00084-7](https://doi.org/10.1016/S0969-2126(98)00084-7).