# Origins and Evolution of the Global RNA Virome

1   Yuri I. Wolf[a], Darius Kazlauskas[b,c], Jaime Iranzo[a], Adriana Lucía-Sanz[a,d], Jens H.

2   Kuhn[e], Mart Krupovic[c], Valerian V. Dolja[f,#], Eugene V. Koonin[a]

3

[a]National Center for Biotechnology Information, National Library of Medicine, National Institutes of Health, Bethesda, Maryland, USA

[b] Vilniaus universitetas biotechnologijos institutas, Vilnius, Lithuania

[c] Département de Microbiologie, Institut Pasteur, Paris, France

[d]Centro Nacional de Biotecnología, Madrid, Spain

[e]Integrated Research Facility at Fort Detrick, National Institute of Allergy and Infectious Diseases, National Institutes of Health, Frederick, Maryland, USA

[f]Department of Botany and Plant Pathology, Oregon State University, Corvallis, Oregon, USA

#Address correspondence to Valerian V. Dolja, doljav@oregonstate.edu

**Running title:** Global RNA Virome

19  **ABSTRACT**

20  Viruses with RNA genomes dominate the eukaryotic virome, reaching enormous diversity in

21  animals and plants. The recent advances of metaviromics prompted us to perform a detailed

22  phylogenomic reconstruction of the evolution of the dramatically expanded global RNA virome.

23  The only universal gene among RNA viruses is the RNA-dependent RNA polymerase (RdRp).

24  We developed an iterative computational procedure that alternates the RdRp phylogenetic tree

25  construction with refinement of the underlying multiple sequence alignments. The resulting tree

26  encompasses 4,617 RNA virus RdRps and consists of 5 major branches, 2 of which include

27  positive-sense RNA viruses, 1 is a mix of positive-sense (+) RNA and double-stranded (ds) RNA

28  viruses, and 2 consist of dsRNA and negative-sense (-) RNA viruses, respectively. This tree

29  topology implies that dsRNA viruses evolved from +RNA viruses on at least two independent

30  occasions, whereas -RNA viruses evolved from dsRNA viruses. Reconstruction of RNA virus

31  evolution using the RdRp tree as the scaffold suggests that the last common ancestors of the

32  major branches of +RNA viruses encoded only the RdRp and a single jelly-roll capsid protein.

33  Subsequent evolution involved independent capture of additional genes, particularly, those

34  encoding distinct RNA helicases, enabling replication of larger RNA genomes and facilitating

35  virus genome expression and virus-host interactions. Phylogenomic analysis reveals extensive

36  gene module exchange among diverse viruses and horizontal virus transfer between distantly

37  related hosts. Although the network of evolutionary relationships within the RNA virome is

38  bound to further expand, the present results call for a thorough reevaluation of the RNA virus

39  taxonomy.

40 **IMPORTANCE**

41 The majority of the diverse viruses infecting eukaryotes have RNA genomes, including

42 numerous human, animal, and plant pathogens. Recent advances of metagenomics have led to

43 the discovery of many new groups of RNA viruses in a wide range of hosts. These findings

44 enable a far more complete reconstruction of the evolution of RNA viruses than what was

45 attainable previously. This reconstruction reveals the relationships between different Baltimore

46 Classes of viruses and indicates extensive transfer of viruses between distantly related hosts,

47 such as plants and animals. These results call for a major revision of the existing taxonomy of

48 RNA viruses.

## INTRODUCTION

Early evolution of life is widely believed to have first involved RNA molecules that functioned both as information storage devices and as catalysts (ribozymes) (1, 2). Subsequent evolution involved the emergence of DNA, the dedicated genomic material, and proteins, the ultimate operational molecules. RNA molecules remained central for translating information from genes to proteins (mRNA), the functioning of the translation machinery itself (rRNA and tRNA), and for a variety of regulatory functions (various classes of non-coding RNA that are increasingly discovered in all life forms) (3). Viruses with RNA genomes ("RNA viruses" for short) that do not involve DNA in their genome replication and expression cycles (4, 5) can be considered the closest extant recapitulation, and possibly, a relic of the primordial RNA world.

RNA viruses comprise 3 of the 7 so-called Baltimore Classes of viruses that differ with respect to the nature of the genome (i.e., the nucleic acid form that is packaged into virions) and correspond to distinct strategies of genome replication and expression: positive-sense (+) RNA viruses, double-stranded (ds) RNA viruses and negative-sense (-) RNA viruses (6). The +RNA viruses use the simplest possible strategy of replication and expression as the same molecule functions as both genome and mRNA (7). Most likely, the first replicators to emerge in the RNA world, after the evolution of translation, resembled +RNA viruses (8). Because the +RNA released from a virion can be directly used for translation to produce viral proteins, virions of +RNA viruses contain only structural proteins, in addition to the genome. In contrast, -RNA and dsRNA viruses package their transcription and replication machineries into their virions because these are necessary to initiate the virus reproduction cycles (9, 10).

RNA viruses comprise a major part of the global virome. In prokaryotes, the known representation of RNA viruses is narrow. Only one family of +RNA viruses (*Leviviridae*) and

4

72    one family of dsRNA viruses (*Cystoviridae*) are formally recognized, and furthermore, their

73    members have limited host ranges. No -RNA viruses have been isolated from prokaryotes (4, 11,

74    12). Although recent metagenomic studies suggest that genetic diversity and host range of

75    prokaryotic RNA viruses could be substantially underestimated (13, 14), it appears that the scope

76    of the prokaryotic RNA virome is incomparably less than that of the DNA virome. As discussed

77    previously, potential causes of the vast expansion of the RNA virome in eukaryotes might

78    include the emergence of the compartmentalized cytosol that provided a hospitable, protective

79    environment for RNA replication that is known to be associated with the endoplasmic reticulum

80    and other membrane compartments (11). Conversely, the nuclear envelope could be a barrier that

81    prevents the access of DNA viruses to the host replication and transcription machineries and thus

82    partially relieves the stiff competition that DNA viruses of prokaryotes represent for RNA

83    viruses.

84        In a sharp contrast, the 3 Baltimore Classes of RNA viruses dominate the eukaryotic

85    virosphere (11, 15). Eukaryotes from all major taxa are hosts to RNA viruses, and particularly in

86    plants and invertebrates, these viruses are enormously abundant and diverse (15-17). Until

87    recently, the study of RNA viromes had been heavily skewed towards viruses infecting humans,

88    livestock, and agricultural plants. Because of these limitations and biases, the attempts to

89    reconstruct the evolutionary history of RNA viruses were bound to be incomplete. Nonetheless,

90    these studies have yielded important generalizations. A single gene encoding an RNA-dependent

91    RNA polymerase (RdRp) is universal among RNA viruses, including capsid-less RNA replicons

92    but excluding some satellite viruses (18). Even within each of the 3 Baltimore Classes, virus

93    genomes do not include fully conserved genes other than those encoding RdRps (15). However,

94    several hallmark genes are shared by broad ranges of RNA viruses including, most notably, those

5

95    encoding capsid proteins of icosahedral and helical virions of +RNA viruses ([19], [20]), and key

96    enzymes involved in virus replication such as distinct helicases and capping enzymes ([15]).

97         The RdRp gene and encoded protein are natural targets of evolutionary analysis because they

98    are the only universal gene (protein) in RNA viruses. However, obtaining strongly supported

99    phylogenies for RdRps is a difficult task due to the extensive sequence divergence, apart from

100   several conserved motifs that are required for polymerase activity ([21-23]). RdRps belong to the

101   expansive class of polymerases containing so-called Palm catalytic domains along with the

102   accessory Fingers and Thumb domains ([24], [25]). In addition to viral RdRps, Palm domain

103   polymerases include reverse transcriptases (RTs) of retroelements and reverse-transcribing

104   viruses, and DNA polymerases that are responsible for genome replication in cellular organisms

105   and diverse DNA viruses. Within the Palm domain class of proteins, RT and the +RNA virus

106   RdRps are significantly similar in sequence and structure, and appear to comprise a

107   monophyletic group ([22], [24-26]). More specifically, the highest similarity is observed between

108   the +RNA virus RdRps and the RTs of group II introns. These introns are wide-spread

109   retrotransposons in prokaryotes that are thought to be ancestral to the RTs of all other

110   retrotransposons as well as retroviruses and pararetroviruses (recently jointly classified as the

111   order *Ortervirales*) of eukaryotes ([27-30]).

112        A phylogenetic analysis of the +RNA virus RdRps revealed only a distant relationship

113   between the leviviruses and the bulk of the eukaryotic +RNA viruses, leaving the ancestral

114   relationships uncertain ([31]). The origin of eukaryotic +RNA viruses from their prokaryotic

115   counterparts is an obvious possibility. However, given the dramatically greater prevalence of

116   +RNA viruses among eukaryotes compared to the narrow spread of leviviruses and "levi-like

117   viruses" in bacteria, an alternative scenario has been proposed. Under this scenario, RdRps of the

118     prokaryotic and eukaryotic +RNA viruses independently descended from distinct RTs (11, 31).

119     Among eukaryotic +RNA viruses, phylogenetic analysis of the RdRps strongly supports the

120     existence of picornavirus and alphavirus "supergroups", which are further validated by additional

121     signature genes (7, 15). However, both the exact compositions of these supergroups and the

122     evolutionary relationships among many additional groups of viruses remain uncertain. Some

123     RdRp phylogenies suggest a third supergroup combining animal "flavi-like viruses" and plant

124     tombusviruses but this unification is not supported by additional shared genes and thus remains

125     tenuous (7, 15, 21).

126     The similarity of RdRps among dsRNA viruses is limited but these RdRps are similar to

127     varying degree to +RNA virus RdRps. Therefore, different groups of dsRNA viruses might have

128     evolved from different +RNA viruses independently, on multiple occasions (15, 32). Although it

129     is not entirely clear how prokaryotic dsRNA viruses fit into this concept, evidence of an

130     evolutionary affinity between cystoviruses and reoviruses has been presented (33, 34).

131     For a long time, the evolutionary provenance of -RNA virus RdRps remained uncertain due

132     to their low sequence similarity to other RdRps and RTs (23). However, recent protein structure

133     comparisons point to a striking similarity between the RdRps of -RNA orthomyxoviruses and

134     those of +RNA flaviviruses and dsRNA cystoviruses (35). All these findings notwithstanding,

135     the overall evolutionary relationships among the RdRps of +RNA, -RNA, and dsRNA viruses,

136     and RTs remain unresolved. In particular, whether the RdRps of +RNA and -RNA viruses are

137     mono- or polyphyletic, is unclear.

138     Many deep evolutionary connections between RNA virus groups that originally were thought

139     to be unrelated have been delineated using the results of pre-metagenomic era evolutionary

140     studies. These discoveries culminated in the establishment of RNA virus supergroups (7, 9, 36).

141    However, the evolutionary provenance of many other RNA virus groups remained unclear, and

142    so did the relationships between the RNA viruses of the 3 Baltimore Classes and retroelements,

143    and their ultimate origins. The prospects of substantial progress appeared dim because of the

144    extreme sequence divergence among RNA viruses, which could amount to irrevocable loss of

145    evolutionary information.

146        Recent revolutionary developments in virus metagenomics (metaviromics) dramatically

147    expanded the known diversity of RNA viruses and provided an unprecedented amount of

148    sequence data for informed investigation into RNA virus evolution (11, 17, 37). The foremost

149    development was the massive expansion of the known invertebrate virome that was achieved

150    primarily through meta-transcriptome sequencing of various holobionts. The subsequent

151    phylogenetic analysis revealed previously unknown lineages of +RNA and -RNA viruses and

152    prompted reconsideration of high rank virus unifications, such as +RNA virus supergroups (14,

153    38-41). The RNA viromes of fungi and prokaryotes also underwent notable expansion albeit not

154    as massive as that of invertebrates (13, 42, 43).

155        Here we reexamine the evolutionary relationships among and within the 3 Baltimore Classes

156    of RNA viruses through a comprehensive analysis of the available genomic and metagenomic

157    sequences. In particular, to build a phylogenetic tree of thousands of viral RdRps, we designed

158    an iterative computational procedure that alternates phylogeny construction with refinement of

159    the underlying multiple alignments. Although RNA viruses have relatively short genomes (~3–

160    41 kb), the combined gene repertoire (pangenome) of these viruses includes numerous genes that

161    are shared, to varying degrees, by related subsets of RNA viruses. To obtain further insight into

162    virus evolution, we therefore attempted to reconstruct the history of gain and loss of conserved

163    proteins and domains in different virus lineages. We also investigated evolution of the single

164  jelly-roll capsid protein (SJR-CP), the dominant type of capsid protein among +RNA viruses.

165  Our analysis revealed patterns that are generally congruent with the RdRp phylogeny and

166  provide further insights into the evolution of different branches of RNA viruses. Finally, we

167  analyzed a bipartite network in which RNA virus genomes are connected via nodes representing

168  virus genes (44, 45) to identify distinct modules in the RNA virosphere. The results shed light on

169  the evolution of RNA viruses, revealing, in particular, the monophyly of -RNA viruses and their

170  apparent origin from dsRNA viruses, which seem to have evolved from distinct branches of

171  +RNA viruses on at least two independent occasions.

172

173  **RESULTS**

174  **Comprehensive phylogeny of RNA virus RNA-dependent RNA polymerases: overall**

175  **structure of the tree and the 5 major branches**

176  Amino acid sequences of RdRps and RTs were collected from the non-redundant National

177  Center for Biotechnology Information (NCBI) database and analyzed using an iterative

178  clustering-alignment-phylogeny procedure (Figure S1; see Materials and Methods for details).

179  This procedure ultimately yielded a single multiple alignment of the complete set of 4,617 virus

180  RdRp sequences (Data set S1) and 1,028 RT sequences  organized in 50+2 clusters (50 clusters

181  of RdRps and 2 clusters of RTs; see Materials and Methods for details). This sequence set did

182  not include RdRps of members of the families *Birnaviridae* or *Permutatetraviridae*, distinct

183  groups of RNA viruses that encompass a circular permutation within the RdRp Palm domain

184  (46) and therefore could not be confidently included in the alignment over their entire lengths.

185    The phylogenetic tree of RdRps and RTs (Data set S2 and Figure 1) was assembled from a

186    set of trees that represent three hierarchical levels of relationships. At the lowest level, full

187    complements of sequences from each cluster were used to construct cluster-specific trees. At the

188    intermediate level, up to 15 representatives from each cluster were selected to elucidate

189    supergroup-level phylogeny. At the highest level, up to 5 representatives from each cluster were

190    taken to resolve global relationships (Data set S3A and S4). The final tree (Data set S2) was

191    assembled by replacing the cluster representatives with the trees from the previous steps.

192    The large number and immense diversity of the viruses included in our analysis create

193    serious challenges for a systematic, phylogeny-based nomenclature of the identified evolutionary

194    lineages of RNA viruses. Many such lineages consist of viruses newly discovered by

195    metaviromics, are not yet formally classified by the International Committee on Taxonomy of

196    Viruses (ICTV), and therefore, cannot be assigned formal names. For the purpose of the present

197    work, we adopted a semi-arbitrary naming scheme using the following approach. 1) We use

198    taxon names that have been fully accepted by the ICTV as of March 2018 (47) whenever

199    possible. These names are recognizable through their capitalization and italicization and rank-

200    specific suffixes (e.g., *-virales* for orders, *-viridae* for families). As is the common practice in

201    virus taxonomy, the officially classified members of each ICTV-approved taxon are referred to

202    via vernaculars (recognizable through their lack of capitalization and italicization). For instance,

203    the members of the ICTV-approved order *Bunyavirales* are called bunyaviruses, whereas those

204    of the family *Tombusviridae* are called tombusviruses. However, in this work, both taxon and

205    vernacular terms are to be understood *sensu lato*: if our analysis indicates certain viruses to be

206    members of or very closely related to an ICTV-established taxon, we consider them members of

207    that taxon despite the lack of current ICTV recognition. As a result, in our analysis, the order

10

208    *Bunyavirales* has more members than in the official taxonomy. 2) we use vernacular names in

209    quotation marks for viruses/lineages that are clearly distinct from those covered by the official

210    ICTV framework. Whenever possible, we use names that circulate in the literature (e.g.,

211    "hepeliviruses", "statoviruses"). In the absence of such unofficial names, we name the lineage

212    reminiscent of the next closely related lineage (e.g., "levi-like viruses" are a clearly distinct sister

213    group to *Leviviridae*/leviviruses). 3) Monophyletic clusters that transcend the currently highest

214    ICTV-accepted rank (i.e., order) are labeled according to terms circulating in the literature (i.e.,

215    "alphavirus supergroup", "flavivirus supergroup", and "picornavirus supergroup"). 4) Lineages

216    represented by a single virus are labeled with the respective virus name.

217        Rooting the phylogenetic tree, generated using PhyML (48), between the RTs and RdRps

218    resulted in a well-resolved topology of RNA viruses which the tree splits into 5 major branches,

219    each including substantial diversity of viruses (Figure 1).

220        Branch 1: leviviruses and their eukaryotic relatives, namely, "mitoviruses", "narnaviruses"

221    and "ourmiaviruses" (the latter three terms are placed in quotation marks as our analysis

222    contradicts the current ICTV framework, which considers mitoviruses and narnaviruses members

223    of one family, *Narnaviridae*, and ourmiaviruses members of a free-floating genus *Ourmiavirus*);

224        Branch 2 ("picornavirus supergroup"): a large assemblage of +RNA viruses of eukaryotes, in

225    particular, those of the orders *Picornavirales* and *Nidovirales*, the families *Caliciviridae*,

226    *Potyviridae, Astroviridae*, and *Solemoviridae*, a lineage of dsRNA viruses including

227    partitiviruses and picobirnaviruses; and several other, smaller groups of +RNA and dsRNA

228    viruses.

11

229    Branch 3: a distinct subset of +RNA viruses including the "alphavirus supergroup" along

230    with "flavivirus supergroup", nodaviruses, and tombusviruses, the "statovirus", "wèivirus",

231    "yànvirus", "zhàovirus" groups; and several additional, smaller groups.

232    Branch 4: dsRNA viruses including cystoviruses, reoviruses, and totivirues, and several

233    additional families.

234    Branch 5: -RNA viruses.

235    Each of these 5 major branches of the tree is strongly supported by bootstrap replications

236    (Figure 1). Assuming the RT rooting of the tree, Branch 1, which consists of leviviruses and their

237    relatives infecting eukaryotes, is a sister group to the rest of RNA viruses; this position is highly

238    robust to the choice of the phylogenetic method and parameters. This tree topology is compatible

239    with the monophyly of the RdRps and, by inference, of RNA viruses and the origin of eukaryotic

240    RNA viruses from a prokaryotic RNA virus ancestor shared with leviviruses. The deeper history

241    remains murky. We have no information on the nature of the common ancestor of retroelements

242    and RNA viruses, let alone whether the ancestor was an RNA virus or a retroelement. However,

243    parsimony considerations suggest that a retroelement ancestor is more likely given that capsids

244    first appeared in the virus part of the tree rather than having been lost in retroelements.

245    The next split in the tree occurs between Branch 2 and the short stem that formally joins

246    Branches 3, 4, and 5. However, the unification of Branch 3 with Branches 4 and 5 is weakly

247    supported and might not reflect actual common ancestry.

248    Arguably, the most striking feature of the RNA virus tree topology is the paraphyly of +RNA

249    viruses relative to dsRNA and -RNA viruses. Indeed, according to this phylogeny, -RNA viruses

250    evolved from within dsRNA viruses, whereas dsRNA viruses are polyphyletic (Figure 1). One

251    major group of dsRNA viruses that includes partitiviruses and picobirnaviruses is firmly

12

252 embedded within the +RNA virus Branch 2, whereas another, larger dsRNA virus group

253 including cystoviruses, reoviruses, totiviruses and viruses from several other families comprise

254 the distinct Branch 4 that might be related to +RNA virus Branch 3 (Figure 1). This placement of

255 the two branches of dsRNA viruses is conceptually compatible with the previous evolutionary

256 scenarios of independent origins from +RNA viruses. However, the presence of a strongly

257 supported branch combining 3 lineages of dsRNA viruses that infect both prokaryotes and

258 eukaryotes suggests a lesser extent of polyphyly in the evolution of dsRNA viruses than

259 originally proposed (15, 32).

260    An alternative phylogenetic analysis of the same RdRp alignment using RAxML yielded the

261 same 5 main branches, albeit some with weak support (Data set S3B). Furthermore, although the

262 dsRNA viruses are split the same way using RAxML as in the PhyML tree, the nested tree

263 topology, in which Branch 4 (the bulk of dsRNA viruses) is lodged deep within +RNA viruses,

264 and Branch 5 (-RNA viruses) is located inside Branch 4, is not reproduced (Data set S3B).

265 Instead, Branches 4 and 5 are separate and positioned deep in the tree, right above the split

266 between Branch 1 and the rest of the RdRps. Given the poor resolution of the RAxML tree and a

267 strong biological argument, namely, the absence of identified -RNA viruses in prokaryotes or

268 protists (with the exception of the "leishbuviruses" infecting kinetoplastids (49, 50); see also

269 discussion below), we believe that the tree topology in Figure 1 carries more credence than that

270 shown in Data set S3B. Nevertheless, these discrepancies emphasize that utmost caution is due

271 when biological interpretation of deep branching in trees of highly divergent proteins is

272 attempted.

273

13

274 **Evolution of the 5 major branches of RNA viruses and reconstruction of gene gain and loss**

275 **events**

276     ***Reconstruction of gene gains and losses***. The RdRp is the only universal protein of RNA

277 viruses. Accordingly, other viral genes must have been gained and/or lost at different stages of

278 evolution. Thus, after performing the phylogenetic analysis of RdRps, we assigned the proteins

279 and domains shared by diverse viruses to the branches of the RdRp tree. Multiple alignments and

280 hidden Markov model (hmm) profiles were constructed for 16,814 proteins and domains

281 encoded by RNA viruses, and a computational pipeline was developed to map these domains on

282 the viral genomes (Data set S5). The resulting patterns of domain presence/absence in the

283 branches of the RdRp tree were used to reconstruct the history of the gains and losses of RNA

284 virus genes (or proteins and domains, for simplicity), both formally, using the ML-based Gloome

285 method (51), and informally, from parsimony considerations.

286     These reconstructions reveal a high level of branch-specificity in the evolution of the gene

287 repertoire of RNA viruses. The only protein that is likely to have been gained at the base of the

288 eukaryotic virus subtree (Branches 2 to 5) (that also includes the bacterial cystoviruses) is the

289 single jelly-roll capsid protein (SJR-CP) (Figure S2A). Retroelements lack capsid proteins, and

290 therefore, there is no indication that SJR-CP was present in the hypothetical element that

291 encoded the common ancestor of RdRps and RTs. Furthermore, reconstruction of the evolution

292 of Branch 1 (leviviruses and their relatives) argues against the ancestral status of SJR-CP in this

293 branch.

294

295     ***Branch 1: leviviruses and their descendants***. The current RdRp tree topology combined

296 with gene gain-loss reconstruction suggests the following evolutionary scenario for Branch 1

14

297    (Figure 2A): a levivirus-like ancestor that, like the extant members of the *Leviviridae*, possessed

298    a capsid protein unrelated to SJR-CP (19, 52) gave rise to naked eukaryotic RNA replicons

299    known as "mitoviruses" and "narnaviruses". These replicons consist of a single RdRp gene (Fig.

300    2B) and replicate either in mitochondria or in the cytosol of the host cells, respectively, of fungal

301    and invertebrate hosts (the latter hosts were identified in metaviromic holobiont analyses) (14,

302    53). Recently, the existence of plant "mitoviruses" has been reported although it is not known

303    whether these viruses reproduce in the mitochondria (54). The "narnavirus" RdRp is also the

304    ancestor of the RdRp of the expanding group of "ourmiaviruses" (Figure 2A). "Ourmiaviruses"

305    were originally identified in a narrow range of plants, and genomic analysis revealed the

306    chimeric nature of these viruses, with a "narnavirus"-like RdRp but SJR-CPs and movement

307    proteins (MPs) which were apparently acquired from "picorna-like" and "tombus-like" viruses,

308    respectively (55). Use of metaviromics has led to the identification of numerous related viruses

309    associated with invertebrates, many of which encode distinct SJR-CP variants and some of which

310    acquired an RNA helicase (Figure 2B and S2) (14). Thus, the evolution of this branch apparently

311    involved the loss of the structural module of leviviruses, which yielded naked RNA replicons

312    that reproduced in the mitochondria of early eukaryotes. A group of these replicons subsequently

313    escaped to the cytosol which was followed by the reacquisition of unrelated structural modules

314    from distinct lineages of eukaryotic viruses inhabiting the same environment (Figure 2B). This

315    complex evolutionary scenario emphasizes the key role of modular gene exchange in the

316    evolution of RNA viruses.

317

318    ***Branch 2: Picornavirus supergroup***. The expansive Branch 2 generally corresponds to the

319    previously described "picornavirus supergroup" (Figures 1 and 3) (7, 31). Some of the virus

15

320  groups that were previously considered peripheral members of this supergroup, such as

321  totiviruses and nodaviruses, were relocated to different branches in the present tree (Branches 4

322  and 3, respectively), whereas the viruses of the order *Nidovirales* were moved inside Branch 2

323  from an uncertain position in the tree. Nevertheless, the core of the supergroup remains coherent,

324  suggestive of common ancestry. Within Branch 2, 3 major clades are strongly supported (Figure

325  3A); however, many of the internal branches are less reliable, so that the relative positions of

326  partitivirus-picobirnavirus, potyvirus-astrovirus and nidovirus clades within Branch 2 remain

327  uncertain.

328      The largest and most coherent of the Branch 2 clades includes the cornerstone of the

329  picornavirus supergroup, the ~826 viruses-strong order *Picornavirales* (56)*,* expanded with

330  caliciviruses, solinviviruses and a multitude of unclassified viruses infecting invertebrates,

331  vertebrates, fungi, protists and undefined hosts (for viruses discovered by metaviromics) (Figure

332  3) (11, 14, 17, 57-59). The second largest, deep-branching clade consists of two lineages that

333  include, respectively, +RNA and dsRNA viruses (Figure 3). The +RNA virus lineage combines

334  astroviruses and potyviruses, the evolutionary affinity of which is well recognized (31, 60). The

335  dsRNA lineage includes the members of the families *Amalgaviridae, Hypoviridae,*

336  *Partitiviridae*, and *Picobirnaviridae,* with each of these families greatly expanded by

337  unclassified affiliates. Finally, the 'middle' clade is smaller and less diverse: it encompasses

338  nidoviruses, including the longest of all +RNA virus genomes (61), and solemoviruses with

339  much shorter genomes (Figure 3). Notably, some members of the family *Luteoviridae* and

340  Heterocapsa circularisquama RNA virus, the only known alvernavirus (62), are nested within the

341  solemovirus clade. Given the lack of support beyond the phylogenetic affinity of the RdRps and

342  the dramatic differences in the genomic architectures of nidoviruses and solemoviruses, the

16

343   possibility that this unification is caused by a tree construction artifact is difficult to rule out (the

344   branch support notwithstanding).

345        Hypoviruses, a group of fungal capsid-less RNA replicons, have been traditionally viewed as

346   dsRNA viruses. However, comparison of genome architectures and phylogenetic analysis

347   suggested that hypoviruses are derivatives of potyviruses that have lost the capsid protein (63,

348   64). In the current RdRp tree, hypoviruses cluster with the dsRNA viruses of the partitivirus-

349   picobirnavirus clade rather than with potyviruses (Figure 3). Whether this position is an artifact

350   of tree construction or whether hypoviruses actually share the RdRps with dsRNA viruses is

351   unclear.

352        The partitivirus-picobirnavirus clade within Branch 2 represents a transition to the *bona fide*

353   dsRNA Baltimore Class (Figure 3). Typical partitiviruses and picobirnaviruses have minimalist

354   genomes that consist of two dsRNA segments encapsidated separately into distinct 120-subunit

355   T=1 capsids (65-68). These genome segments encode, respectively, RdRps and CPs that are

356   clearly homologous between the two families. The CPs of the partitivirus-picobirnavirus clade

357   have been suggested to be distantly related to those of other dsRNA viruses that belong to

358   Branch 4 (33, 69). Notably, this clade also includes some naked RNA replicons that reproduce in

359   algal mitochondria or chloroplasts, use a mitochondrial genetic code and, in terms of lifestyle,

360   resemble "mitoviruses" (14, 70, 71). By analogy, the origin of the partitivirus-picobirnavirus

361   group from an as-yet undiscovered lineage of prokaryotic RNA viruses seems likely. More

362   specifically, this group of dsRNA viruses could have evolved through reassortment of genomic

363   segments encoding, respectively, a +RNA virus RdRp of Branch 2 (possibly a naked RNA

364   replicon) and a dsRNA virus capsid protein related to those of Branch 4 viruses. The most

365   recently evolved branch of partitiviruses is characterized by larger, 4–6-partite genomes, in

17

366    contrast to mono- or bipartite genomes in the deeper branches (14). This observation emphasizes

367    a major tendency in virus evolution: increase in genome complexity via gradual acquisition of

368    accessory genes (72).

369    Apart from the SJR-CP, an apparently ancestral protein that is likely to be a shared derived

370    character (synapomorphy) of Branch 2 is a serine protease that is present in members of the order

371    *Picornavirales* (with the diagnostic substitution of cysteine for the catalytic serine), members of

372    the potyvirus-astrovirus clade, solemoviruses, alvernavirus, and nidoviruses (Figure 3B and

373    S2B). As demonstrated previously, this viral protease derives from a distinct bacterial protease,

374    probably of mitochondrial origin, which is compatible with an early origin of Branch 2 in

375    eukaryotic evolution (31).

376    The reconstruction of protein gain-loss, together with the comparison of genome

377    architectures in this branch, reveal extensive rearrangements as well as gene and module

378    displacement (Figure 3B and S2B). Branch 2 includes viruses with relatively long genomes and

379    complex gene repertoires (nidoviruses, potyviruses and many members of *Picornavirales*) along

380    with viruses with much shorter genomes and minimal sets of genes (astroviruses and

381    solemoviruses). Clearly, evolution of Branch 2 viruses involved multiple gene gains. Of special

382    note is the gain of 3 distinct helicases in 3 clades within this branch: superfamily 3 helicases

383    (S3H) in members of *Picornavirales*, superfamily 2 helicases (S2H) in potyviruses, and

384    superfamily 1 helicases (S1H) in nidoviruses (Figure 3B and S2B). This independent, convergent

385    gain of distinct helicases reflects the trend noticed early in the study of RNA virus evolution,

386    namely, that most viruses with genomes longer than ~6 kb encode helicases, whereas smaller

387    ones do not. This difference conceivably exists because helicase activity is required for the

388    replication of longer RNA genomes (73). Another notable feature is the change of virion

18

389    morphology among potyviruses (replacement of ancestral SJR-CP by an unrelated CP forming

390    filamentous virions) and nidoviruses (displacement by a distinct nucleocapsid protein). The

391    dramatic change in virion morphology and mode of genome encapsidation might have been

392    necessitated by the inability of SJR-CP-based icosahedral capsids to accommodate the larger

393    genomes of the ancestral potyviruses and nidoviruses. In addition, nidoviruses gained capping

394    enzymes (CapE; Figure 3B and S2B) that most likely were acquired independently of the

395    capping enzymes of other RNA viruses. Nidoviruses also gained ribonucleases and other

396    accessory proteins that are involved in genome replication, virulence and other aspects of the

397    infection cycle of these largest known RNA viruses (61, 74-76).

398

399    ***Branch 3: "Alphavirus supergroup", "flavivirus supergroup" and the extensive diversity of***

400    ***"tombus-like viruses"***. Branch 3 is the part of the +RNA virus RdRp tree that underwent the

401    most dramatic rearrangements compared to previous versions. This branch consists of two

402    strongly supported clades of +RNA viruses: i) the assemblage that originally was defined as the

403    "alphavirus supergroup" (7, 15) joined by several additional groups of viruses; and ii)

404    flaviviruses and related viruses ("flavivirus supergroup"; Figure 1 and 4). In the former clade, the

405    alphavirus supergroup encompasses an enormous diversity of plant, fungal, and animal +RNA

406    viruses,       and       consists       of       3       well-supported       lineages:       tymoviruses,

407    virgaviruses/alphaviruses/endornaviruses, and hepeviruses/benyviruses, each accompanied by

408    related viruses that often form yet unclassified lineages (Figure 4A). Within the "alphavirus

409    supergroup" alone, the genome lengths range from ~6 to ~20 kb. Despite this length variation, all

410    supergroup members harbor a conserved RNA replication gene module encoding a CapE, S1H,

19

411   and RdRp; the conservation of this module attests to the monophyly of the supergroup (Figure

412   4B).

413       In contrast, virion architectures vary dramatically even within each of the three lineages of

414   the "alphavirus supergroup". The major structural themes include: variants of icosahedral capsids

415   formed by SJR-CP (e.g., bromoviruses, tymoviruses); unrelated icosahedral capsids enveloped in

416   a lipoprotein bilayer (togaviruses); flexuous filamentous capsids formed by a distinct type of CP

417   (alphaflexiviruses, betaflexiviruses, gammaflexiviruses, closteroviruses); and rigid rod-shaped

418   capsids assembled from another distinct CP (benyiviruses, virgaviruses). It was traditionally

419   thought that the latter capsid type is specific to viruses of flowering plants (20). However, the

420   recent discovery of a virgavirus-like CP in invertebrate viruses (e.g., Běihǎi charybdis crab virus

421   1 in Figure 4B) (14) suggests that the emergence of this unique CP fold antedates land

422   colonization by plants by ~100 Mya. Yet another 'structural' theme is offered by endornaviruses,

423   naked RNA replicons which, similarly to hypoviruses in Branch 2 (see above), originally were

424   classified as dsRNA viruses. However, endornaviruses possess all the hallmarks of the

425   alphavirus supergroup and clearly are derived from +RNA viruses of this group. They seem to

426   have been mislabeled dsRNA viruses due to the accumulation of dsRNA replication

427   intermediates in infected cells (18, 77). A parallel loss of the CP genes apparently occurred in

428   deltaflexiviruses which, in RdRp phylogenies, form a sister group to the flexible filamentous

429   gammaflexiviruses (78), and in umbraviruses that are included in the family *Tombusviridae*

430   based on the RdRp phylogeny. Notably, unlike most other capsid-less viruses that are vertically

431   inherited, umbraviruses can hijack capsids of co-infecting luteoviruses for aphid transmission

432   (79).

20

433    Within Branch 3, the phylogenetically compact alphavirus supergroup is embedded within

434    the radiation of diverse virus groups including the well-known tombusviruses and nodaviruses,

435    along with several newcomers discovered via metaviromics, such as the "statovirus", "wèivirus",

436    "yànvirus", and "zhàovirus" groups (14, 80, 81) (Figure 4A). Our RdRp analysis revealed

437    remarkable phylogenetic heterogeneity within and among these groups and split "tombus-like

438    viruses" into 5 lineages with distinct evolutionary affinities (groups 'Uncl. inv.', and subsets of

439    "tombus-like viruses" and "nodaviruses" in Figure 4A). This subdivision is also supported by the

440    analysis of the CPs of these viruses (see section on SJR-CP evolution; Fig. 7). Therefore, in

441    contrast to the alphavirus supergroup, nodaviruses or flavivirus supergroup, the term 'tombus-

442    like' loses its evolutionary and taxonomic coherence. Accordingly, we use the term

443    "tombusviruses" (without quotation marks) only for one lineage that includes the members of the

444    current family *Tombusviridae* along with a broad variety of related plant and invertebrate

445    holobiont viruses (14).

446    The previously suggested, tenuous *Flaviviridae-Tombusviridae* affinity is gone in the present

447    tree although members of both families belong to the same major Branch 3. Plant tombusviruses

448    (and members of closely related plant virus genera), the only group of "tombus-like viruses" that

449    was available at the time of previous analyses (7, 21), now form but a small twig deep within the

450    large assemblage we refer to as tombusviruses. Tombusviruses are affiliated with "statoviruses"

451    (80) and a subset of unclassified viruses from invertebrate holobionts rather than with

452    flaviviruses (Figure 4A). Flaviviruses now form a separate clade within Branch 3, the flavivirus

453    supergroup that includes members of four recognized flaviviral genera (*Pegivirus, Hepacivirus,*

454    *Flavivirus*, and *Pestivirus*), the newly discovered "jīngménviruses" with segmented genomes

455    (38) and a variety of unclassified, extremely divergent "flavi-like viruses" of animals and plants.

21

456    This clade is split into two well-supported lineages, one of which includes pegiviruses and

457    hepaciviruses, and the other one consists of the rest of flaviviruses (Figure 4A). Flaviviral virions

458    are enveloped, with the envelope proteins forming an external icosahedral shell, whereas the core

459    nucleocapsid is apparently disordered; the evolutionary provenance of the core protein, with its

460    distinct fold, is unclear (19, 82, 83). Notably, flaviviral envelope proteins are class II fusion

461    proteins that are closely related to alphavirus envelope E1 proteins (84). The theme of gene

462    swapping between these distantly related virus groups of Branch 3 is further emphasized by the

463    homology between alphavirus CPs which form icosahedral capsids under the lipid envelopes,

464    and flavivirus non-structural NS3 proteases that share a chymotrypsin-like fold (84). Because the

465    RdRp tree topology implies that the alphavirus ancestor is more recent than the ancestor of

466    flaviviruses (Figure 1), such adoption of the NS3 protease for a structural role is suggestive of

467    emerging alphaviruses borrowing their structural module from preexisting flaviviruses (19).

468        The hallmark of Branch 3 is the capping enzyme (CapE), which is present in the entire

469    "alphavirus supergroup" and in flaviviruses (Figure 4 and S2B). A highly divergent version of

470    CapE has been identified in nodaviruses (85) and, in our present analysis, in the additional subset

471    of viruses that grouped with nodaviruses, as well as a few viruses scattered throughout the clade.

472    Formally, CapE is inferred to be ancestral in the entire Branch 3. However, CapEs of "alphavirus

473    supergroup" members, nodaviruses, and flaviviruses are only distantly related to one another,

474    and at least the latter have closer eukaryotic homologs, the FtsJ family methyltransferases (86,

475    87). Furthermore, tombusviruses, statoviruses, yànviruses, zhàoviruses, wèiviruses, and members

476    of the *Pegivirus-Hepacivirus* lineage of flaviviruses lack CapE, putting into question its presence

477    in the ancestor of this branch (Figure 4B, S2B). The most credible evolutionary scenario seems

478    to involve convergent acquisition of CapEs on at least 3 independent occasions, recapitulating

479   the apparent history of helicases in Branch 2 (see above; Figure 3B and S2B). The trend of the

480   capture of helicases by +RNA viruses with longer genomes also holds in Branch 3 and includes

481   the acquisition of S1H at the base of the alphavirus supergroup and S2H by the ancestral

482   flavivirus (Figure 4B, S2B).

483       To an even greater extent than in Branch 2, the apparent routes of virus evolution in Branch 3

484   involve lineage-specific gene capture resulting in evolution of complex genome architectures

485   (Figure 4B). The most notable cases are closteroviruses and divergent flaviviruses that have

486   genomes of up to 20–26 kb, rivalling coronaviruses in terms of genome length and the

487   complexity of the gene repertoire (38, 88-90).

488       The lack of genes assigned to the common ancestor of Branch 3 (with the obvious exception

489   of the RdRp) prevents development of a coherent evolutionary scenario for the entire branch. In

490   the case of the clade encompassing the "alphavirus supergroup" and related viruses, a potential

491   common ancestor could be a simple virus that encoded only an RdRp and a SJR-CP, a CP fold

492   most broadly represented in this clade including diverse tombusviruses, nodaviruses and

493   members of *Bromoviridae*, and *Tymoviridae* within the "alphavirus supergroup". Proposing such

494   an ancestor for the flavivirus clade is challenged by the lack of viruses with short and simple

495   genomes among flaviviruses. Indeed, the lengths of the genomes in this clade vary from ~9 kb to

496   ~26 kb, with even the shortest ones encoding at least three of the flavivirus signature genes

497   (serine protease [Spro], S2H, and RdRp). One potential clue, however, is provided by

498   "jīngménviruses" with tetra-partite genomes in which the protease-helicase modules and the

499   RdRps are encoded by separate genome segments; two other segments apparently encode

500   structural proteins of unclear provenance (38). This genome architecture could hint at an

501   ancestral flavivirus genome that was assembled from genes borrowed from pre-existing viruses,

23

502    one of which possessed a divergent "tombus-like virus" RdRp. Although the origins of Branch 3

503    are murky, major trends in its subsequent evolution clearly included lineage-specific gene

504    capture, starting with helicases and CapEs in the ancestors of the major lineages and followed by

505    diverse genes in smaller groups (Figure 4B).

506

507    **Branch 4: dsRNA viruses**. Branch 4, which joins Branch 3 with weak support, includes the

508    bulk of the dsRNA viruses (Figure 1 and 5). All dsRNA viruses in this branch share a unique

509    virion organization and encode homologous CPs. In particular, the specialized icosahedral

510    capsids of these viruses, involved in transcription and replication of the dsRNA genome inside

511    cells, are constructed from 60 homo- or heterodimers of CP subunits organized on an unusual

512    $T=1$ (also known as pseudo-$T=2$) lattice (69, 91). The only exception are the chrysoviruses which

513    encode large CPs corresponding to the CP dimers of other dsRNA viruses and form genuine $T=1$

514    capsids (92). Icosahedral capsids of partitiviruses and picobirnaviruses, which encode RdRps

515    belonging to Branch 2, are also constructed from 60 homodimers (66, 67, 93) and have been

516    suggested to be evolutionarily related to those of the dsRNA viruses from RdRp Branch 4 (94)

517    despite little structural similarity between the corresponding CPs. Totiviruses, many of which

518    have "minimal" genomes encoding only RdRps and CPs, comprise one of the two major clades

519    in Branch 4, whereas cystoviruses, the only known prokaryotic dsRNA viruses, together with the

520    vast family *Reoviridae*, which consists of multi-segmented dsRNA viruses infecting diverse

521    eukaryotes, comprise the second clade (Figure 5). The closer phylogenetic affinity between

522    cystoviruses and reoviruses appears to be corroborated by the fact that the inner $T=1$ icosahedral

523    capsid is uniquely encased by the outer icosahedral shell constructed on a $T=13$ lattice in both

524    families (34). Both cystoviruses and reoviruses appear to have gained many clade-specific genes,

24

525  in particular, RecA-like packaging ATPases of the former (95) and the CapEs of the latter that

526  are only distantly related to CapEs of other RNA viruses and likely were acquired independently

527  (96, 97) (Figure 5, S2B).

528

529  ***Branch 5: -RNA viruses***. Branch 5, the 100% supported lineage combining all -RNA

530  viruses, is lodged within Branch 4 as the sister group of reoviruses, and this position is upheld by

531  two strongly supported internal branches in the RdRp tree (Figure 1 and 6). The -RNA branch

532  splits into 2 strongly supported clades. The first clade encompasses the 348 viruses-strong

533  membership of the order *Mononegavirales* (98), along with the members of the distantly related

534  family *Aspiviridae* (99), 3 groups of -RNA viruses discovered through metaviromics

535  ("chǔviruses", "qínviruses", "yuèviruses") (14, 39), and a group of unclassified fungal viruses

536  (Figure 6A) (42, 100). In contrast to the members of the *Mononegavirales*, most of which

537  possess unsegmented genomes, the remainder of this clade is characterized by bi-, tri-, or even

538  tetrasegmented genomes (Figure 6B). The second clade combines the family *Orthomyxoviridae*,

539  the genus *Tilapinevirus* (101) and the large order *Bunyavirales* (394 viruses) (102). The latter

540  order consists of two branches, one of which is the sister group to the

541  orthomyxovirus/tilapinevirus clade (albeit with weak support). The number of negative-sense or

542  ambisense genome segments in this clade varies from 2–3 in most of bunyaviruses to 8 in viruses

543  of the genus *Emaravirus* and the orthomyxovirus/tilapinevirus group (10, 99, 101, 103). A

544  notable acquisition in the first clade is a CapE, whereas members of the second clade share "cap-

545  snatching" endonucleases (En) (10).

546

547  **Patterns of the single jelly-roll capsid protein evolution**

25

548    The SJR-CP is the dominant type of CP among +RNA viruses and is also found in members of

549    one family of dsRNA viruses (*Birnaviridae*). Structural comparisons indicate that SJR-CPs of

550    RNA viruses form a monophyletic group and likely have been recruited from cellular SJR

551    proteins on a single occasion during the evolution of RNA viruses (19). The short length and

552    high divergence of SJR-CPs preclude adequate resolution in phylogenetic analysis; thus, we

553    performed profile-profile sequence comparisons and clustering of all viral SJR-CP sequences in

554    our dataset (see Methods for details). Analysis of the resulting network revealed patterns that are

555    generally congruent with the RdRp phylogeny and provide further insights into the evolution of

556    different branches of RNA viruses.

557        At conservative P value thresholds ($P<1e^{-10}$), the majority of SJR-CPs segregated into two

558    large clusters both of which contained representatives from RdRp Branch 2. Cluster 1 included

559    the members of *Picornavirales*, *Caliciviridae*, and diverse "picorna-like viruses" of

560    invertebrates, whereas cluster 2 consisted of the members of the families *Astroviridae*,

561    *Luteoviridae*, and *Solemoviridae* and "solemo-like viruses" (Figure S3). In addition, cluster 2

562    contained members of several families from RdRp Branch 3, namely, *Tombusviridae* (and

563    diverse "tombus-like viruses"), *Hepeviridae*, a subgroup of *Nodaviridae* and "statoviruses".

564        At less restrictive P value thresholds ($P<1e^{-03}$), all SJR-CPs were interconnected, largely,

565    through making contacts to the core of cluster 2. Only "ourmiaviruses" had stronger affinity to

566    picornaviruses in cluster 1 (Figure 7). This pattern of connectivity is consistent with the radiation

567    of SJR-CPs from a common ancestor, likely resembling sequences from cluster 2 of Branch 2.

568    This analysis also revealed high CP sequence divergence among members of some families (e.g.,

569    *Bromoviridae*) and numerous cases of apparent CP gene replacement. For instance, the CPs of

570    nodaviruses fall into two groups: one is related to the turreted CPs of tetraviruses and the other is

26

571    similar to CPs of tombusviruses, mirroring the RdRp phylogeny (Figure 7). At a greater

572    phylogenetic distance, CPs of astroviruses and hepeviruses are closely related despite them being

573    affiliated to Branches 2 and 3, respectively, suggesting CP gene replacement in the ancestor of

574    one of the two families. Given that the CPs of hepeviruses connect to SJR-CPs of other viruses

575    through astroviruses, CP gene replacement most likely occurred in the ancestor of hepeviruses

576    (Figure S3). Notably, the CPs of "zhàoviruses", "wèiviruses", and "tombus-like" and "solemo-

577    like viruses" (diverse virus assemblage within solemovirus branch, to the exclusion of bona fide

578    *Solemoviridae*) did not form discrete clusters but rather were affiliated with diverse virus groups,

579    suggesting extensive recombination in these viruses, with multiple CP gene exchanges (Figure

580    7). In the case of unclassified "narnaviruses" and "ourmiaviruses", the CP genes apparently have

581    been acquired on more than 3 independent occasions from different groups of viruses,

582    emphasizing the impact of recombination and gene shuffling in the evolution of RNA viruses.

583    Previously, a similar extent of chimerism has been also observed among ssDNA viruses (104,

584    105), highlighting the evolutionary and functional plasticity of short viral genomes.

585

586    **The modular gene-sharing network of RNA viruses: gene transfer and module shuffling**

587    The pronounced structural and functional modularity of virus proteomes and pervasive shuffling

588    of the genomic regions encoding distinct protein modules are key features of virus evolution (11,

589    15, 17). Therefore, a productive approach to the study of the virosphere that complements

590    phylogenetics is the construction and analysis of networks of gene sharing. Bipartite networks, in

591    which one type of nodes corresponds to genes and the other one to genomes, have been

592    employed to investigate the dsDNA domain of the virosphere (45). This analysis revealed

593    hierarchical modular organization of the network, with several modules including non-obvious

27

594    connections between disparate groups of viruses (44, 106). Although, for RNA viruses, this type

595    of analysis is less informative due to the small number of proteins encoded in each viral genome,

596    the "pan-proteome" of RNA viruses is large (Data set S5), prompting us to experiment with

597    bipartite gene sharing networks for RNA viruses. The initial search for statistically significant

598    modularity identified 54 distinct modules, most of which included a single virus family (Figure

599    8A, B). Remarkably, the family *Reoviridae* has been split into 5 modules, highlighting the vast

600    diversity of this family, comparable to that in order-level taxa. Among the exceptions, the most

601    expansive module included the viruses of the order *Picornavirales* (module 29), together with

602    the family *Caliciviridae* (module 47), that are linked through the conserved suit of genes

603    including SJR-CP, chymotrypsin-like protease, S3H and the RdRp (Figure 8A and C). Viruses of

604    the order *Bunyavirales* were also recovered in a single module that is characterized by the

605    presence of a conserved nucleocapsid (with the exception of the families *Nairoviridae* and

606    *Arenaviridae*) and the cap-snatching endonuclease (module 51; Figure 8A and C).

607    The next stage of the network analysis aims at detecting supermodules that are formed from

608    the primary modules via connecting genes. The supermodules of RNA viruses failed to attain

609    statistical significance due to the small number of shared genes but nevertheless, some notable

610    connections are revealed by this analysis. Specifically, 8 overlapping supermodules were

611    identified (Figure 8C). The largest and, arguably, most remarkable is a supermodule that

612    combines +RNA, dsRNA and -RNA viruses that share the capping enzymes, S1H (with the

613    exception of *Reoviridae* and *Mononegavirales*) and additional connector genes (e.g., the OTU

614    family protease) that link some of the constituent modules. The second largest supermodule

615    combines large subsets of viruses from the RdRp Branches 2 and 3 that are connected through

616    the SJR-CP and the chymotrypsin-like protease. Another supermodule encompasses enveloped

28

617     +RNA viruses of the families *Flaviviridae* and *Togaviridae*, and -RNA viruses of the order

618     *Bunyavirales* (except for *Arenaviridae*) that share homologous envelope glycoproteins (except

619     for flaviviruses) and class II fusion proteins.

620     Information from gene sharing is inherently limited for RNA viruses due to the small number

621     of genes in each genome. Nevertheless, the bipartite network analysis reveals prominent

622     "horizontal" connections that are underlain either by actual gene exchange or by parallel

623     acquisition of homologous genes by distinct RNA viruses.

624

**Host ranges of RNA viruses: evolutionary implications and horizontal virus transfer**

626     RNA viruses have been identified in representatives of all major divisions of eukaryotes,

627     whereas in prokaryotes, members of two families of RNA viruses are known to infect only a

628     limited range of hosts (11, 13, 15, 31) . For Branch 1 in our phylogenetic tree of RdRps, the route

629     of evolution from leviviruses infecting prokaryotes to eukaryotic "ourmiaviruses" of plants and

630     invertebrates is readily traceable and involves a merger between a levivirus-derived naked RNA

631     replicon that eukaryotes most likely inherited from the mitochondrial endosymbiont with the

632     SJR-CP of a eukaryotic "picorna-like virus". Notably, such a merger seems to have occurred on

633     at least three other independent occasions in Branch 1 because several groups of invertebrate

634     holobiont "narnaviruses" and some "ourmiaviruses" encode distantly related SJR-CPs that

635     apparently were acquired from different groups of plant and animal viruses (Figure 7).

636     The case of cystoviruses is less clear given that this clade is sandwiched between eukaryotic

637     viruses in Branch 4 and therefore does not seem to be a good candidate for the ancestor of this

638     branch. It appears more likely that the ancestor was a toti-like virus, whereas cystoviruses are

639     derived forms, which implies virus transfer from eukaryotes to prokaryotes. However, an

29

640    alternative scenario might be considered. No known prokaryotic viruses are classified in Branch

641    2 but it has been proposed that picobirnaviruses, for which no hosts have been reliably identified,

642    actually are prokaryotic viruses. This proposal is based on the conspicuous conservation of

643    functional, bacterial-type, ribosome-binding sites (Shine-Dalgarno sequences) in picobirnavirus

644    genomes (107, 108). Should that be the case, viruses of prokaryotes might be lurking among

645    totiviruses as well. Then, Branch 4 would stem from a prokaryotic ancestor avoiding the need to

646    invoke virus transfer from eukaryotes to prokaryotes to explain the origin of the cystoviruses.

647    We made an attempt to quantify the potential horizontal virus transfer (HVT) events in RNA

648    viruses that represent the 5 major branches of the RdRp tree. The leaves of the tree were labeled

649    with the known hosts, the entropy of the host ranges for each subtree was calculated, and the

650    resulting values were plotted against the distance from the root (Figure 9). By design, for all

651    branches, entropy (host diversity) drops from the maximum values at the root to zero at the

652    leaves. All branches show substantial host range diversity such that, for example, at half-distance

653    from root to leaves, all branches, except for Branch 1, retain at least half of the diversity (Figure

654    9). Furthermore, differences between the branches are substantial, with the highest entropy

655    observed in branches 4 (dsRNA viruses) and 5 (-RNA viruses). With all the caveats due to

656    potential errors and ambiguities in host assignment, this analysis strongly suggests that HVT

657    played an important role in the evolution of all major groups of RNA viruses.

658

659    **DISCUSSION**

660    **RNA virus evolution coming into focus**

30

661     This work was prompted by the advances of metaviromics, which have dramatically increased

662     the known diversity of RNA viruses ([11](#), [13](#), [17](#), [37](#), [57](#)). We reasoned that this expansion of the

663     RNA virosphere could provide for an improved understanding of virus evolution. Although

664     further progress of metaviromics and enhanced phylogenomic methods will undoubtedly change

665     current ideas, we believe that some key aspects of RNA virus evolution are indeed coming into

666     focus.

667          The expanded diversity of RNA viruses combined with the iterative procedure for

668     phylogenetic analysis allowed us to obtain a tree of all RdRps and the most closely related RTs

669     in which the main branches are strongly supported and thus appear to be reliable (Figure 1). To

670     our knowledge, the picture of RNA virus evolution emerging from the tree has not been

671     presented previously. The tree seems to clarify the relationships between the 3 Baltimore Classes

672     of RNA viruses by revealing the nested tree structure in which dsRNA viruses evolved, on at

673     least two occasions, from +RNA viruses, whereas -RNA viruses evolved from a distinct group of

674     dsRNA viruses.

675          The derivation of -RNA viruses from dsRNA viruses is, arguably, the most unexpected

676     outcome of the present analysis, considering the lack of genes (other than the RdRp) shared by

677     these virus classes. Clearly, given that the primary evidence behind the derivation of -RNA

678     viruses from within dsRNA viruses comes from deep phylogeny, extreme caution is due in the

679     interpretation of this observation. However, the pronounced similarity between the 3D structures

680     of the RdRps of the -RNA influenza virus A and bacteriophage φ6 dsRNA cystovirus ([35](#)) is

681     compatible with our findings. Further, because virtually no -RNA viruses are known in

682     prokaryotes or unicellular eukaryotes [with the single exception of "leischbuviruses" in parasitic

31

683    trypanosomatids (50) that were likely acquired from the animal hosts of these protists], their later

684    origin from a preexisting group of +RNA or dsRNA viruses appears most likely.

685        The +RNA to dsRNA to -RNA scenario of RNA virus genome evolution also makes sense in

686    terms of the molecular logic of genome replication-expression strategies. Indeed, +RNA viruses

687    use the simplest genomic strategy and, in all likelihood, represent the primary pool of RNA

688    viruses. The dsRNA viruses, conceivably, evolved when a +RNA virus switched to

689    encapsidating a replicative intermediate (dsRNA) together with the RdRp. Naked replicons

690    similar to "mitoviruses", hypoviruses, and endornaviruses might have been evolutionary

691    intermediates in this process. This switch does not seem to be as "easy" and common as

692    previously suspected (15, 32) but, nevertheless, appears to have occurred at least twice during

693    the evolution of RNA viruses. The origin of -RNA viruses is the next step during which the plus-

694    strand is discarded from the virions, perhaps simplifying the processes of transcription and

695    replication. Conceivably, the evolution of dsRNA and -RNA viruses, in which transcription and

696    replication of the viral genomes are confined to the interior of virions or nucleocapsid

697    transcription/replication complexes and no dsRNA accumulates in the infected cells, was driven

698    by the advantage of escaping some of the host defense mechanisms, in particular, RNA

699    interference (109, 110). The membrane-associated replication complexes of +RNA viruses could

700    represent an initial step in this direction (9).

701        Obviously, evolution of the RdRp does not equal evolution of viruses: other genes, in

702    particular those encoding capsid and other structural proteins, are crucial for virus reproduction,

703    and these genes often have different histories. The reconstruction of gene gain and loss sheds

704    some light on these aspects of RNA virus evolution. The ancestors of each of the major branches

705    of RNA viruses except for Branch 1 appear to have been simply organized +RNA viruses

32

706      resembling tombusviruses (Figure 10). Thus, these types of viruses encoding RdRps and SJR-

707      CPs might have been ancestral to the bulk of eukaryotic RNA viruses (apart from those in

708      Branch 1 that derive directly from prokaryotic leviviruses). The subsequent parallel capture of

709      different helicases enabled evolution of increasingly complex genomes via accumulation of

710      additional genes (Figure 10). Notably, similar levels of complexity, with the complete coding

711      capacity of 20 to 40 kb, were reached independently in 4 branches of the RdRp tree, namely,

712      Branch 2 (coronaviruses), Branch 3 (closteroviruses and flaviviruses), Branch 4 (reoviruses), and

713      Branch 5 (filoviruses and paramyxoviruses) (Figure 3-6).

714         Gene exchange and shuffling of gene modules are important factors of RNA virus evolution.

715      For example, it appears that all dsRNA viruses in the two major clades within Branches 2 and 4

716      share homologous structural modules that combine with distinct RdRps. At least in the case of

717      partiti-picobirnaviruses in Branch 2, the dsRNA virus ("toti-like virus") CP apparently displaced

718      the ancestral SJR-CP. However, this particular protein structure does not seem to be essential to

719      encapsidate a dsRNA genome: birnaviruses, whose provenance is uncertain due to the

720      permutation in their RdRps, retain SJR-CP, which is most closely related to SJR-CP of

721      nodaviruses and tetraviruses (19). An even more striking example of module shuffling is

722      presented by amalgaviruses, dsRNA viruses that group with partitiviruses in the RdRp tree but

723      encode a distant homolog of the nucleocapsid protein of -RNA bunyaviruses (111-114). More

724      generally, structural and replication modules have been repeatedly shuffled during the evolution

725      of +RNA viruses. Examples include displacement of the ancestral SJR-CP by a filamentous CP

726      in potyviruses and by a helical nucleocapsid protein in nidoviruses, and multiple cases of

727      displacement with rod-shaped-like CP and unique nucleocapsid proteins in Branch 3. Thus,

728     exchange of genes and gene modules among RNA viruses is pervasive and can cross the

729     boundaries of Baltimore Classes.

730       Another recurring trend in RNA virus evolution is the loss of the structural module resulting

731     in the emergence of naked RNA replicons such as "narnaviruses" and "mitoviruses" in Branch 1,

732     hypoviruses in Branch 2, and endornaviruses, umbraviruses and deltaflexiviruses in Branch 3

733     (18). On some occasions, broad horizontal spread of a gene leads to a major shift in the lifestyle

734     of viruses, such as adaptation of viruses to a new type of hosts. The primary cases in point are

735     the movement proteins (MPs) of plant viruses that are represented in all 5 branches of the RdRp

736     tree and, outside of the RNA part of the virosphere, in plant caulimoviruses and badnaviruses,

737     and ssDNA viruses (115).

738

739     **The prevalence of HVT and the overall course of RNA virus evolution**

740     Arguably, the most striking realization brought about by metaviromics is the diverse host range

741     of numerous groups of viruses, even tight ones that occupy positions near the tips of the RdRp

742     tree. Extending early observations on then highly unexpected similarities between viruses of

743     animals and plants, recent metaviromic analyses reveal numerous clusters of indisputably related

744     viruses infecting animals and plants, plants and fungi, and in some cases, animals or plants and

745     protists. The evolutionary relationships between viruses with distinct host ranges are supported

746     not only by the phylogeny of the RdRp, but also by the fact that these viruses share additional

747     conserved domains, such as, for example, SJR1, S3H and 3C-Pro in the case of *Picornavirales*

748     members infecting protists, plants, fungi, invertebrates and vertebrates (Fig. 3).

749       Invertebrates are particularly promiscuous hosts for viruses, often sharing the same virus

750     group with distantly related organisms. Certainly, much caution is due in the interpretation of

751  host range assignments from metaviromics, especially, holobiont studies. Viruses identified in

752  holobiont samples of, say, invertebrates could actually infect protists associated with these

753  animals (50) or could represent contamination from fungal, plant or even prokaryotic sources.

754  These uncertainties notwithstanding, the extensive diversity of hosts even within small branches

755  of the RdRp tree is undeniable. A coevolution scenario in which the ancestors of all these viruses

756  originated from the common ancestor of the respective groups of eukaryotes and coevolved with

757  the hosts implies an enormous diversity of RNA viruses in early eukaryotes. This scenario

758  appears to be highly unlikely given the apparent paucity of RNA viruses in the extant protists

759  (although new metaviromic studies might substantially expand the range of protist viruses). The

760  pervasive HVT alternative seems much more plausible, especially given that arthropods,

761  nematodes, and other invertebrates are well known as virus vectors and thus fit the role of RNA

762  virus reservoirs (11) .

763  In addition to invertebrates that appear to be dominant HVT agents, fungi could also play an

764  important role in HVT within the global RNA virome. Indeed, fungi that are tightly associated

765  with plants and insects often share closely related viruses with these organisms (116-118).

766  Furthermore, an indisputable case of cross-kingdom transfer of an insect iflavirus to an

767  entomopathogenic fungus has been recently described (119).

768  These findings appear to be best compatible with a grand evolutionary scenario (Figure 10)

769  in which the ancestor of the eukaryotic RNA virome was a levi/narnavirus-like naked RNA

770  replicon that originally reproduced in mitochondria and combined with a host carbohydrate-

771  binding SJR protein (19) or a preexisting SJR-CP from a DNA virus to form a simple ancestral

772  virus. Given that viruses of Branches 2, 3, and 4 are present in modern protists, it appears likely

773  that these branches emerged in early eukaryotes. However, because of the apparent dominance of

35

774    the viruses of the RdRp Branch 2 ("picornavirus supergroup" in general and "aquatic picorna-

775    like viruses" clade in particular) in protists, this branch likely diversified first, whereas the

776    diversification of Branches 3 and 4 occurred later, after ancestral protist viruses were transferred

777    to marine invertebrates during the Cambrian explosion. The recent analysis of the viromes of

778    ctenophores, sponges and cnidarians suggests that substantial diversification of RNA viruses

779    occurred already in these deeply branching metazoa (120). Invertebrates brought their already

780    highly diverse RNA virome to land at terrestrialization and subsequently inoculated land plants.

781    In land plants, RNA viruses, particularly those of Branch 3, dramatically expanded, in part,

782    perhaps, because of the exclusion of competing large DNA viruses. Finally, it seems plausible

783    that, given the high prevalence of -RNA viruses in metazoa and their virtual absence in protists

784    [with the exception of the recently discovered "leishbuviruses" that likely invaded their parasitic

785    protist hosts via HVT from an animal host (49, 50)], these viruses that comprise the RdRp

786    Branch 5 evolved in animals via mixing and matching genes from reovirus-like and flavivirus-

787    like ancestors.

788

789    **The impending overhaul of RNA virus taxonomy**

790    The expansion of the global RNA virome thanks to the advances of metaviromics, combined

791    with the phylogenomics results, seem to call for an overhaul of the current virus taxonomy on

792    multiple levels. Most importantly, creation of a coherent, hierarchical system with multiple

793    taxonomic ranks seems to be imminent. This process has already started with the proposal of a

794    phylum rank for -RNA viruses, for which monophyly is unequivocally supported by the present

795    analysis (Figures 1 and 6). This phylum could consist of two subphyla with multiple classes and

796    orders. At least 4 additional phyla of RNA viruses can be confidently predicted to emerge,

36

797   including, respectively, the dsRNA viruses of Branch 4, and +RNA viruses of Branches 1, 2, and

798   3. Each of these phyla will undoubtedly have a rich internal structure. In addition, some of the

799   current families do not seem to be compatible with the expansive RdRp trees present here and in

800   a previous analysis (14). For instance, the families *Coronaviridae, Togaviridae* and

801   *Rhabdoviridae* are likely to be split into two families each.

802   While the present study was in preparation, a major attempt on a comprehensive virus

803   classification has been published (121). This work analyzed a dendrogram that was produced

804   from distance matrices between viruses derived from sequence similarity scores combined with

805   measures of gene composition similarity. . Unlike our present analysis, this approach pre-

806   supposes monophyly of each of the Baltimore classes. Furthermore, given that their analysis is

807   based on measures of similarity rather than on phylogenetic analysis proper, this approach is best

808   regarded as producing a phenetic classification of viruses rather that an evolutionary

809   reconstruction as such. Some of the groups delineated by this method, particularly, among +RNA

810   viruses, are closely similar to those reported here. Others, however, are widely different: for

811   instance, the order *Mononegavirales* does not come across as monophyletic in their

812   dendrograms. We did not attempt a complete comparison; such an exercise could be useful in the

813   future, for better understanding the routes of RNA virus evolution.

814

## CONCLUDING REMARKS

815

816   Through metaviromics, many aspects of the global RNA virome evolution can be clarified.

817   Certainly, reconstruction of the deepest events in this evolutionary history is bound to remain

818   tentative, especially, because the RdRp is the only universal gene of the RNA viruses, and hence,

819   the only one that can serve as the template for evolutionary reconstructions. At the depth of

820    divergence characteristic of RdRps, the relationship between the major branches in the tree

821    cannot be established with confidence. Nevertheless, monophyly of several expansive groups, in

822    particular the 5 main branches in the RdRp tree, is strongly supported. Because of the stability of

823    these branches, biologically plausible scenarios of evolution emerge under which dsRNA viruses

824    evolved from different groups of +RNA viruses, whereas -RNA viruses evolved from dsRNA

825    viruses.

826    Evolutionary reconstructions suggest that the last common ancestors of each major lineage of

827    eukaryotic +RNA viruses were simple viruses that encoded only the RdRp and a CP, most likely,

828    of the SJR fold. Subsequent evolution involved independent capture of distinct helicases which

829    apparently facilitate replication of larger, more complex +RNA genomes. The helicase-assisted

830    replication of +RNA genomes created the opportunities for parallel acquisition of additional

831    genes encoding proteins involved in polyprotein processing and virus genome expression, such

832    as proteases and capping enzymes, respectively, and proteins involved in virus-host interactions,

833    such as MPs or RNAi suppressors of plant viruses. In addition to these processes of vertical

834    evolution of RNA viruses, phylogenomic analysis reveals multiple cases of gene module

835    exchange among diverse viruses and pervasive HVT, often between distantly related hosts, such

836    as animals and plants. Together, these processes have shaped a complex network of evolutionary

837    relationships among RNA viruses.

838    The much anticipated comprehensive exploration of the RNA viromes of prokaryotes and

839    unicellular eukaryotes, such as free-living excavates, chromalveolates, rhizaria, amoebozoa and

840    choanoflagellates, as well as deeply rooted metazoa, will undoubtedly help in developing better

841    supported evolutionary scenarios for each of the 5 major branches of the RNA virus tree.

842 Nevertheless, it is already clear that the current taxonomy of RNA viruses is due for a complete

843 overhaul.

844

845 **MATERIALS AND METHODS**

846 **Phylogeny of RNA-dependent RNA polymerases.** Protein sequences belonging to RNA

847 viruses, excluding retroviruses, and unclassified viruses were downloaded from the NCBI

848 GenBank database in April 2017 (122). Initial screening for RdRp domains was performed using

849 PSI-BLAST (123)(e-value of 0.01, effective database size of $2x10^8$) with position-specific

850 scoring matrices (PSSMs) produced from the available RdRp alignments. The sources included

851 group-specific alignments for +RNA viruses and dsRNA viruses (12, 31) and the 4 PFAM

852 alignments for the -RNA viruses (pfam00602, pfam00946, pfam04196, and pfam06317) from

853 the NCBI conserved domain database (CDD) (124). Additionally, a set of RTs from group-II

854 introns and non-long terminal repeat (LTR) retrotransposons was extracted from GenBank as an

855 outgroup.

856 Extracted RdRp footprints were filtered for the fraction of unresolved amino acids (at most

857 10%) and clustered using UCLUST (125) with a similarity threshold of 0.9. One representative

858 from each cluster was selected for further analysis. The resulting set contained 4,640 virus

859 RdRps. This set went through several rounds of semi-manual curation whereby sequences were

860 clustered using UCLUST, aligned using MUSCLE (126), and cross-searched against each other

861 and their parent sequences (often, complete viral polyproteins) using PSI-BLAST and

862 HHSEARCH (127). Upon the results of these searches, the boundaries of the RdRp domain were

863 expanded or trimmed to improve their compatibility with each other.

864    The RdRp and RT sequences were subjected to an iterative clustering and aligned procedure,

865    organized as follows: Initially, sequences were clustered using UCLUST with a similarity

866    threshold of 0.5; clustered sequences were aligned using MUSCLE, and singletons were

867    converted to pseudo-alignments consisting of just one sequence. Sites containing more than 67%

868    of gaps were temporarily removed from alignments and pairwise similarity scores were obtained

869    for clusters using HHSEARCH. Scores for a pair of clusters were converted to distances (the

870    $d_{A,B} = -\log(s_{A,B}/\min(s_{A,A}, s_{B,B})$ formula, in which $d_{A,B}$ is the distance between cluster A and B and

871    $s_{A,B}$ is the HHSEARCH score for the comparison of these clusters, was used to convert scores $s$

872    to distances $d$). The matrix of pairwise distances was used to make an unweighted pair group

873    method with arithmetic (UPGMA) mean ([128]). Shallow tips of the tree were used as the guide

874    tree for a progressive pairwise alignment of the clusters at the tree leaves using HHALIGN

875    ([127]), resulting in larger clusters. This procedure was reiterative, ultimately resulting in the

876    single alignment of the whole set of 4,640 virus RdRp sequences and 1,028 RT sequences.

877    During the clustering procedure, 50 virus RdRp clusters, consisting of 1 to 545 sequences,

878    were defined. These clusters represent either well-established groups of related viruses (roughly

879    comparable to the ICTV family rank) or, in case of poorly characterized and unclassified viruses,

880    groups of well-aligned RdRps that are clearly distinct from others. In all cases, uncertainties

881    were treated conservatively, i.e., in favor of placing sequences with questionable relatedness into

882    separate clusters. Additionally, RT sequences were placed into two clusters consisting of group-

883    II intron RTs and non-LTR retrotransposon RTs.

884    For each cluster consisting of more than two sequences, an approximate maximum likelihood

885    (ML) phylogenetic tree was constructed from the cluster-specific alignment using the FastTree

886    program ([129]) (Whelan and Goldman [WAG] evolutionary model with gamma-distributed site

887   rates) with sites, containing more than 50% of gaps removed from the alignment. Trees were

888   rooted using a variant of the mid-point rooting procedure such that the difference between the

889   (weighted) average root-to-tip distances in the subtrees on the opposite sides of the root is

890   minimized.

891       To resolve the structure of the global relationships, up to five representatives from each

892   cluster were selected using the within-cluster trees to ensure the diversity of the selected

893   sequences. This procedure resulted in a set of 228 virus RdRps and 10 RTs. The alignments of

894   the selected sequences were extracted from the master alignment and filtered for sites containing

895   more than 50% of gaps. A ML phylogenetic tree was reconstructed for the resulting alignment

896   using PhyML (48) (Le Gascuel [LG] evolutionary model with gamma-distributed site rates and

897   empirical amino acid frequencies; aBayes support values). Another form of branch support,

898   bootstrap support by transfer (BOOSTER) phylogenetic bootstrap (130), was also used to assess

899   the reliability of the major tree divisions. Alternatively, the same RdRp alignment was used as

900   the input for ML phylogenetic analysis using RAxML (LG evolutionary model with gamma-

901   distributed site rates and empirical amino acid frequencies).

902       RdRps in the global tree were divided into 5 major branches (supergroups). Up to 15

903   representatives from each cluster were selected to form supergroup-level alignments. Respective

904   trees were reconstructed from these alignments using the same procedure (PhyML tree with LG

905   evolutionary model with gamma-distributed site rates and empirical amino acid frequencies).

906       The overall tree was assembled manually by first replacing the supergroup representatives in

907   the global tree with the supergroup trees, and then, replacing the cluster representatives with the

908   cluster trees. The lower-level trees were rooted according to the arrangement of the

909   representatives in the upper-level tree.

41

910     **Identification of protein domains.** Protein domains were identified using a representative set of

911     RNA virus genomes, including representative members of all ICTV-approved virus families and

912     unclassified virus groups. This set was annotated manually using sensitive profile-profile

913     comparisons with the HHsuite package (127), and hmm profiles for annotated proteins or their

914     domains were generated by running one iteration of HHblits against the latest (October 2017)

915     uniclust30 database (131). Each annotated profile was assigned to a functional category (e.g.,

916     "capsid protein_jelly-roll", "chymotrypsin-like protease"). These profiles were used to annotate

917     the genomes of all viruses included in the RdRp phylogenetic analysis and for which complete

918     (or near-complete) genome sequences were available. Profiles for the latter proteins were

919     generated by running one iteration of Jackhmmer (132) against UniRef50 database. Protein

920     regions that did not have significant hits were extracted and clustered with cluster analysis of

921     sequences (CLANS) (133), and groups containing at least three members were identified,

922     annotated (if possible), and added to the manually annotated profile database. The last step of the

923     domain identification procedure was then repeated using the updated RNA virus profile

924     database. Highly similar viruses (having identical domain organization and >94% identical

925     RdRps) were removed from the dataset. In addition, only one representative genome was

926     retained for some members of over-represented species (e.g., *Hepacivirus C*). Many genomes

927     encoded readthrough proteins (e.g., alphaviruses, tombusviruses, nidoviruses), resulting in

928     domains from the shorter protein contained within the longer readthrough protein. Such

929     redundancies were removed by filtering out the shorter of the two proteins sharing >80%

930     identity. The final set of annotated genomes included 2,839 viruses.

931     **Reconstruction of the history of gene gain and loss.** The protein domains identified in the

932     2,839 virus genomes were mapped onto the composite tree of virus RdRps. A set of 500

42

933    representatives was chosen among the (mostly complete) genomes. The domain complement

934    data and the respective tree were analyzed using the GLOOME program (51). Domain gains and

935    losses were inferred at each tree branch using the difference of posterior probabilities for the

936    domain occurrence in the nodes at the proximal and the distal ends of the branch (difference $>0.5$

937    implies gain; difference $<-0.5$ implies loss).

938    **Clustering analysis.** The SJR CP network was generated by performing all-against-all

939    comparison of the SJR CP profiles. To generate hmm profiles, SJR CP sequences were extracted

940    from the total proteome of RNA viruses and two iterations of HHblits were performed against

941    the uniprot20_2016_02 database. The resultant profiles were compared to each other with

942    HHsearch (127). Their similarity P-values were extracted from the result files and used as an

943    input for CLANS program (133). Clusters were identified using a network-based algorithm

944    implemented in CLANS. Resulting clusters were manually inspected and refined.

945    **Analysis of bipartite gene-genome networks.** A bipartite network was built to study the

946    patterns of gene sharing among viral genomes (44, 106). After removing genomes with less than

947    2 domains and domains that appear in less than 3 genomes, the network consisted of 2,829

948    nodes, of which 2,515 correspond to genomes and 314 to domains. Genome and domain nodes

949    were connected by links whenever a domain is present in a genome. A consensus community

950    detection approach was used to identify the modules of the network (134, 135). First, we ran 500

951    replicas with Infomap (136) (bipartite setting, using domains as factors) and built a similarity

952    matrix by assigning to each pair of nodes a similarity value equal to the fraction of replicas in

953    which both nodes were placed in the same module. Then, hierarchical clustering was performed

954    on the similarity matrix, setting the number of clusters equal to the median number of modules

955    obtained in the 500 replicas. The order statistics local optimization method (Oslom) software

956 (137) was subsequently used to filter significant modules with a p-value threshold equal to 0.05.

957 To detect higher-order (super)modules, we first identified connector domains as those present in

958 at least 2 modules with prevalence greater than 0.65. The 2nd-order network composed of 54

959 modules and 34 connector domains was searched for 2nd-order modules with Infomap (500

960 replicas, bipartite setting with connector domains as factors). Due to the small size of the 2nd-

961 order network, consensus community detection did not qualitatively improve the results of the

962 search, and therefore we took the replica with the best Infomap score and skipped hierarchical

963 clustering for this and subsequent steps of the higher-order module search. After assessing

964 statistical significance with Oslom, 4 2nd-order modules were recovered, encompassing 12 of the

965 original modules associated with closely related virus families. A 3rd-order network was built by

966 pooling these 4 2nd-order modules withthe 46 modules that remained unmerged. The 3rd-order

967 network was analyzed in the same manner to obtain the 9 supermodules. None of these

968 supermodules was assessed as significant by Oslom.

969 **Quantification of virus host range diversity**. Known hosts or, in case of viruses, isolated from

970 holobionts, virus sources were identified for 3,456 viruses. The host taxonomy (to the phylum

971 level) was mapped to the leaves of the combined tree. The tree was ultrametrized and the weights

972 were computed for all leaves as described in (138). Each internal node of the tree, therefore,

973 defines a set of leaves with assigned hosts; the distribution of (weighted) relative frequencies of

974 hosts can be characterized by its Shannon entropy. Each leaf has one host defined, so the entropy

975 of the host range distribution is 0, whereas at the root the host diversity is maximal. The entropy

976 generally declines along the tree because viruses that belong to relatively shallow branches tend

977 to share their host ranges. The node depth vs node host range entropy data was averaged for each

44

978    major branch separately using a Gaussian kernel with the bandwidth equal to 10% of the total

979    tree depth.

980

45

## REFERENCES

992

993    1.    **Bernhardt HS.** 2012. The RNA world hypothesis: the worst theory of the early evolution

994          of life (except for all the others)(a). Biol Direct **7:**23.

995    2.    **Gilbert W.** 1986. Origin of Life - the RNA World. Nature **319:**618-618.

996    3.    **Nelson JW, Breaker RR.** 2017. The lost language of the RNA World. Sci Signal

997          **10:**eaam8812.

998    4.    **Koonin EV, Senkevich TG, Dolja VV.** 2006. The ancient Virus World and evolution of

999          cells. Biol Direct **1:**29.

1000   5.    **Baltimore D.** 1971. Expression of animal virus genomes. Bacteriol Rev **35:**235-241.

1001   6.    **Baltimore D.** 1971. Viral genetic systems. Trans N Y Acad Sci **33:**327-332.

1002   7.    **Koonin EV, Dolja VV.** 1993. Evolution and taxonomy of positive-strand RNA viruses:

1003          implications of comparative analysis of amino acid sequences. Crit Rev Biochem Mol

1004          Biol **28:**375-430.

1005   8.    **Koonin EV, Dolja VV.** 2013. A virocentric perspective on the evolution of life. Curr

1006          Opin Virol **3:**546-557.

1007   9.    **Ahlquist P.** 2006. Parallels among positive-strand RNA viruses, reverse-transcribing

1008          viruses and double-stranded RNA viruses. Nat Rev Microbiol **4:**371-382.

1009   10.   **Reguera J, Gerlach P, Cusack S.** 2016. Towards a structural understanding of RNA

1010          synthesis by negative strand RNA viral polymerases. Curr Opin Struct Biol **36:**75-84.

1011   11.   **Dolja VV, Koonin EV.** 2018. Metagenomics reshapes the concepts of RNA virus

1012          evolution by revealing extensive horizontal virus transfer. Virus Res **244:**36-52.

1013    12.    **Bolduc B, Shaughnessy DP, Wolf YI, Koonin EV, Roberto FF, Young M.** 2012.

1014           Identification of novel positive-strand RNA viruses by metagenomic analysis of archaea-

1015           dominated Yellowstone hot springs. J Virol **86:**5562-5573.

1016    13.    **Krishnamurthy SR, Janowski AB, Zhao G, Barouch D, Wang D.** 2016.

1017           Hyperexpansion of RNA bacteriophage diversity. PLoS Biol **14:**e1002409.

1018    14.    **Shi M, Lin X-D, Tian J-H, Chen L-J, Chen X, Li C-X, Qin X-C, Li J, Cao J-P, Eden**

1019           **J-S, Buchmann J, Wang W, Xu J, Holmes EC, Zhang Y-Z.** 2016. Redefining the

1020           invertebrate RNA virosphere. Nature **540:**539-543.

1021    15.    **Koonin EV, Dolja VV, Krupovic M.** 2015. Origins and evolution of viruses of

1022           eukaryotes: the ultimate modularity. Virology **479-480:**2-25.

1023    16.    **Dolja VV, Koonin EV.** 2011. Common origins and host-dependent diversity of plant and

1024           animal viromes. Curr Opin Virol **1:**322-331.

1025    17.    **Shi M, Zhang Y-Z, Holmes EC.** 2018. Meta-transcriptomics and The Evolutionary

1026           Biology of RNA Viruses. Virus Res **243:**83-90.

1027    18.    **Koonin EV, Dolja VV.** 2014. Virus world as an evolutionary network of viruses and

1028           capsidless selfish elements. Microbiol Mol Biol Rev **78:**278-303.

1029    19.    **Krupovic M, Koonin EV.** 2017. Multiple origins of viral capsid proteins from cellular

1030           ancestors. Proc Natl Acad Sci U S A **114:**E2401-E2410.

1031    20.    **Dolja VV, Boyko VP, Agranovsky AA, Koonin EV.** 1991. Phylogeny of capsid

1032           proteins of rod-shaped and filamentous RNA plant viruses: two families with distinct

1033           patterns of sequence and probably structure conservation. Virology **184:**79-86.

1034    21.    **Koonin EV.** 1991. The phylogeny of RNA-dependent RNA polymerases of positive-

1035           strand RNA viruses. J Gen Virol **72 ( Pt 9):**2197-2206.

1036    22.    **Xiong Y, Eickbush TH.** 1990. Origin and evolution of retroelements based upon their

1037          reverse transcriptase sequences. EMBO J **9:**3353-3362.

1038    23.    **Poch O, Sauvaget I, Delarue M, Tordo N.** 1989. Identification of four conserved motifs

1039          among the RNA-dependent polymerase encoding elements. EMBO J **8:**3867-3874.

1040    24.    **Ng KK-S, Arnold JJ, Cameron CE.** 2008. Structure-function relationships among

1041          RNA-dependent RNA polymerases. Curr Top Microbiol Immunol **320:**137-156.

1042    25.    **te Velthuis AJW.** 2014. Common and unique features of viral RNA-dependent

1043          polymerases. Cell Mol Life Sci **71:**4403-4420.

1044    26.    **Eickbush TH, Jamburuthugoda VK.** 2008. The diversity of retrotransposons and the

1045          properties of their reverse transcriptases. Virus Res **134:**221-234.

1046    27.    **Gladyshev EA, Arkhipova IR.** 2011. A widespread class of reverse transcriptase-related

1047          cellular genes. Proc Natl Acad Sci U S A **108:**20311-20316.

1048    28.    **Lambowitz AM, Belfort M.** 2015. Mobile bacterial group II introns at the crux of

1049          eukaryotic evolution. Microbiol Spectr **3:**MDNA3-0050-2014.

1050    29.    **Novikova O, Belfort M.** 2017. Mobile group II introns as ancestral eukaryotic elements.

1051          Trends Genet **33:**773-783.

1052    30.    **Krupovic M, Blomberg J, Coffin JM, Dasgupta I, Fan H, Geering AD, Gifford R,**

1053          **Harrach B, Hull R, Johnson W, Kreuze JF, Lindemann D, Llorens C, Lockhart B,**

1054          **Mayer J, Muller E, Olszewski NE, Pappu HR, Pooggin MM, Richert-Pöggeler KR,**

1055          **Sabanadzovic S, Sanfaçon H, Schoelz JE, Seal S, Stavolone L, Stoye JP, Teycheney**

1056          **P-Y, Tristem M, Koonin EV, Kuhn JH.** 2018. *Ortervirales*: new virus order unifying

1057          five families of reverse-transcribing viruses. J Virol **92:**e00515-00518.

1058    31.    **Koonin EV, Wolf YI, Nagasaki K, Dolja VV.** 2008. The Big Bang of picorna-like virus

1059           evolution antedates the radiation of eukaryotic supergroups. Nat Rev Microbiol **6:**925-

1060           939.

1061    32.    **Koonin EV.** 1992. Evolution of double-stranded RNA viruses: a case for polyphyletic

1062           origin from different groups of positive-stranded RNA viruses. Semin Virol **3:**327-339.

1063    33.    **El Omari K, Sutton G, Ravantti JJ, Zhang H, Walter TS, Grimes JM, Bamford DH,**

1064           **Stuart DI, Mancini EJ.** 2013. Plate tectonics of virus shell assembly and reorganization

1065           in phage Φ8, a distant relative of mammalian reoviruses. Structure **21:**1384-1395.

1066    34.    **Poranen MM, Bamford DH.** 2012. Assembly of large icosahedral double-stranded

1067           RNA viruses. Adv Exp Med Biol **726:**379-402.

1068    35.    **Pflug A, Guilligay D, Reich S, Cusack S.** 2014. Structure of influenza A polymerase

1069           bound to the viral RNA promoter. Nature **516:**355-360.

1070    36.    **Goldbach R, Wellink J.** 1988. Evolution of plus-strand RNA viruses. Intervirology

1071           **29:**260-267.

1072    37.    **Greninger AL.** 2018. A decade of RNA virus metagenomics is (not) enough. Virus Res

1073           **244:**218-229.

1074    38.    **Shi M, Lin X-D, Vasilakis N, Tian J-H, Li C-X, Chen L-J, Eastwood G, Diao X-N,**

1075           **Chen M-H, Chen X, Qin X-C, Widen SG, Wood TG, Tesh RB, Xu J, Holmes EC,**

1076           **Zhang Y-Z.** 2016. Divergent viruses discovered in arthropods and vertebrates revise the

1077           evolutionary history of the *Flaviviridae* and related viruses. J Virol **90:**659-669.

1078    39.    **Li C-X, Shi M, Tian J-H, Lin X-D, Kang Y-J, Chen L-J, Qin X-C, Xu J, Holmes EC,**

1079           **Zhang Y-Z.** 2015. Unprecedented genomic diversity of RNA viruses in arthropods

1080           reveals the ancestry of negative-sense RNA viruses. Elife **4:**e05378.

49

1081 40. **Webster CL, Longdon B, Lewis SH, Obbard DJ.** 2016. Twenty-five new viruses

1082 associated with the Drosophilidae (Diptera). Evol Bioinform Online **12:**13-25.

1083 41. **Fauver JR, Grubaugh ND, Krajacich BJ, Weger-Lucarelli J, Lakin SM, Fakoli LS,**

1084 **3rd, Bolay FK, Diclaro JW, 2nd, Dabire KR, Foy BD, Brackney DE, Ebel GD,**

1085 **Stenglein MD.** 2016. West African *Anopheles gambiae* mosquitoes harbor a

1086 taxonomically diverse virome including new insect-specific flaviviruses,

1087 mononegaviruses, and totiviruses. Virology **498:**288-299.

1088 42. **Marzano S-YL, Nelson BD, Ajayi-Oyetunde O, Bradley CA, Hughes TJ, Hartman**

1089 **GL, Eastburn DM, Domier LL.** 2016. Identification of diverse mycoviruses through

1090 metatranscriptomics characterization of the viromes of five major fungal plant pathogens.

1091 J Virol **90:**6846-6863.

1092 43. **Deakin G, Dobbs E, Bennett JM, Jones IM, Grogan HM, Burton KS.** 2017. Multiple

1093 viral infections in *Agaricus bisporus* - Characterisation of 18 unique RNA viruses and 8

1094 ORFans identified by deep sequencing. Sci Rep **7:**2469.

1095 44. **Iranzo J, Krupovic M, Koonin EV.** 2016. The double-stranded DNA virosphere as a

1096 modular hierarchical network of gene sharing. MBio **7:**e00978-00916.

1097 45. **Iranzo J, Krupovic M, Koonin EV.** 2017. A network perspective on the virus world.

1098 Commun Integr Biol **10:**e1296614.

1099 46. **Gorbalenya AE, Pringle FM, Zeddam J-L, Luke BT, Cameron CE, Kalmakoff J,**

1100 **Hanzlik TN, Gordon KHJ, Ward VK.** 2002. The palm subdomain-based active site is

1101 internally permuted in viral RNA-dependent RNA polymerases of an ancient lineage. J

1102 Mol Biol **324:**47-62.

1103  47.  **King AMQ, Lefkowitz EJ, Mushegian AR, Adams MJ, Dutilh BE, Gorbalenya AE,**

1104  **Harrach B, Harrison RL, Junglen S, Knowles NJ, Kropinski AM, Krupovic M,**

1105  **Kuhn JH, Nibert ML, Rubino L, Sabanadzovic S, Sanfaçon H, Siddell SG,**

1106  **Simmonds P, Varsani A, Zerbini FM, Davison AJ.** 2018. Changes to taxonomy and

1107  the International Code of Virus Classification and Nomenclature ratified by the

1108  International Committee on Taxonomy of Viruses (2018). Arch Virol **163:**2601-2631.

1109  48.  **Guindon S, Dufayard J-F, Lefort V, Anisimova M, Hordijk W, Gascuel O.** 2010.

1110  New algorithms and methods to estimate maximum-likelihood phylogenies: assessing the

1111  performance of PhyML 3.0. Syst Biol **59:**307-321.

1112  49.  **Akopyants NS, Lye L-F, Dobson DE, Lukeš J, Beverley SM.** 2016. A novel

1113  bunyavirus-like virus of trypanosomatid protist parasites. Genome Announc **4:**e00715-

1114  00716.

1115  50.  **Grybchuk D, Akopyants NS, Kostygov AY, Konovalovas A, Lye L-F, Dobson DE,**

1116  **Zangger H, Fasel N, Butenko A, Frolov AO, Votýpka J, d'Avila-Levy CM, Kulich P,**

1117  **Moravcová J, Plevka P, Rogozin IB, Serva S, Lukeš J, Beverley SM, Yurchenko V.**

1118  2018. Viral discovery and diversity in trypanosomatid protozoa with a focus on relatives

1119  of the human parasite *Leishmania*. Proc Natl Acad Sci U S A **115:**E506-E515.

1120  51.  **Cohen O, Ashkenazy H, Belinky F, Huchon D, Pupko T.** 2010. GLOOME: gain loss

1121  mapping engine. Bioinformatics **26:**2914-2915.

1122  52.  **Golmohammadi R, Valegård K, Fridborg K, Liljas L.** 1993. The refined structure of

1123  bacteriophage MS2 at 2.8 Å resolution. J Mol Biol **234:**620-639.

1124  53.  **Hillman BI, Cai G.** 2013. The family *Narnaviridae*: simplest of RNA viruses. Adv Virus

1125  Res **86:**149-176.

1126    54.    **Nibert ML, Vong M, Fugate KK, Debat HJ.** 2018. Evidence for contemporary plant

1127            mitoviruses. Virology **518:**14-24.

1128    55.    **Rastgou M, Habibi MK, Izadpanah K, Masenga V, Milne RG, Wolf YI, Koonin EV,**

1129            **Turina M.** 2009. Molecular characterization of the plant virus genus *Ourmiavirus* and

1130            evidence of inter-kingdom reassortment of viral genome segments as its possible route of

1131            origin. J Gen Virol **90:**2525-2535.

1132    56.    **Le Gall O, Christian P, Fauquet CM, King AMQ, Knowles NJ, Nakashima N,**

1133            **Stanway G, Gorbalenya AE.** 2008. *Picornavirales*, a proposed order of positive-sense

1134            single-stranded RNA viruses with a pseudo-T = 3 virion architecture. Arch Virol

1135            **153:**715-727.

1136    57.    **Culley A.** 2018. New insight into the RNA aquatic virosphere via viromics. Virus

1137            Research **244:**84-89.

1138    58.    **Ng TFF, Marine R, Wang C, Simmonds P, Kapusinszky B, Bodhidatta L, Oderinde**

1139            **BS, Wommack KE, Delwart E.** 2012. High variety of known and new RNA and DNA

1140            viruses of diverse origins in untreated sewage. J Virol **86:**12161-12175.

1141    59.    **Hause BM, Palinski R, Hesse R, Anderson G.** 2016. Highly diverse posaviruses in

1142            swine faeces are aquatic in origin. J Gen Virol **97:**1362-1367.

1143    60.    **Jiang B, Monroe SS, Koonin EV, Stine SE, Glass RI.** 1993. RNA sequence of

1144            astrovirus: distinctive genomic organization and a putative retrovirus-like ribosomal

1145            frameshifting signal that directs the viral replicase synthesis. Proc Natl Acad Sci U S A

1146            **90:**10539-10543.

1147    61.    **Saberi A, Gulyaeva AA, Brubacher JL, Newmark PA, Gorbalenya AE.** 2018. A

1148           planarian nidovirus expands the limits of RNA genome size. bioRxiv

1149           doi:10.1101/299776**:**299776.

1150    62.    **Nagasaki K, Shirai Y, Takao Y, Mizumoto H, Nishida K, Tomaru Y.** 2005.

1151           Comparison of genome sequences of single-stranded RNA viruses infecting the bivalve-

1152           killing dinoflagellate *Heterocapsa circularisquama*. Appl Environ Microbiol **71:**8888-

1153           8894.

1154    63.    **Koonin EV, Choi GH, Nuss DL, Shapira R, Carrington JC.** 1991. Evidence for

1155           common ancestry of a chestnut blight hypovirulence-associated double-stranded RNA

1156           and a group of positive-strand RNA plant viruses. Proc Natl Acad Sci U S A **88:**10647-

1157           10651.

1158    64.    **Dawe AL, Nuss DL.** 2013. Hypovirus molecular biology: from Koch's postulates to host

1159           self-recognition genes that restrict virus transmission. Adv Virus Res **86:**109-147.

1160    65.    **Nibert ML, Ghabrial SA, Maiss E, Lesker T, Vainio EJ, Jiang D, Suzuki N.** 2014.

1161           Taxonomic reorganization of family *Partitiviridae* and other recent progress in

1162           partitivirus research. Virus Res **188:**128-141.

1163    66.    **Tang J, Ochoa WF, Li H, Havens WM, Nibert ML, Ghabrial SA, Baker TS.** 2010.

1164           Structure of Fusarium poae virus 1 shows conserved and variable elements of partitivirus

1165           capsids and evolutionary relationships to picobirnavirus. J Struct Biol **172:**363-371.

1166    67.    **Duquerroy S, Da Costa B, Henry C, Vigouroux A, Libersou S, Lepault J, Navaza J,**

1167           **Delmas B, Rey FA.** 2009. The picobirnavirus crystal structure provides functional

1168           insights into virion assembly and cell entry. EMBO J **28:**1655-1665.

1169　68.　**Nibert ML, Tang J, Xie J, Collier AM, Ghabrial SA, Baker TS, Tao YJ.** 2013. 3D

1170　　　　structures of fungal partitiviruses. Adv Virus Res **86:**59-85.

1171　69.　**Luque D, Gómez-Blanco J, Garriga D, Brilot AF, González JM, Havens WM,**

1172　　　　**Carrascosa JL, Trus BL, Verdaguer N, Ghabrial SA, Castón JR.** 2014. Cryo-EM

1173　　　　near-atomic structure of a dsRNA fungal virus shows ancient structural motifs preserved

1174　　　　in the dsRNA viral lineage. Proc Natl Acad Sci U S A **111:**7641-7646.

1175　70.　**Koga R, Fukuhara T, Nitta T.** 1998. Molecular characterization of a single

1176　　　　mitochondria-associated double-stranded RNA in the green alga Bryopsis. Plant Mol Biol

1177　　　　**36:**717-724.

1178　71.　**Koga R, Horiuchi H, Fukuhara T.** 2003. Double-stranded RNA replicons associated

1179　　　　with chloroplasts of a green alga, *Bryopsis cinicola*. Plant Mol Biol **51:**991-999.

1180　72.　**Koonin EV, Dolja VV.** 2006. Evolution of complexity in the viral world: the dawn of a

1181　　　　new vision. Virus Res **117:**1-4.

1182　73.　**Gorbalenya AE, Koonin EV.** 1989. Viral proteins containing the purine NTP-binding

1183　　　　sequence pattern. Nucleic Acids Res **17:**8413-8440.

1184　74.　**Snijder EJ, Decroly E, Ziebuhr J.** 2016. The nonstructural proteins directing

1185　　　　coronavirus RNA synthesis and processing. Adv Virus Res **96:**59-126.

1186　75.　**Enjuanes L, Zuñiga S, Castaño-Rodriguez C, Gutierrez-Alvarez J, Canton J, Sola I.**

1187　　　　2016. Molecular basis of coronavirus virulence and vaccine development. Adv Virus Res

1188　　　　**96:**245-286.

1189　76.　**Sola I, Almazán F, Zúñiga S, Enjuanes L.** 2015. Continuous and discontinuous RNA

1190　　　　synthesis in coronaviruses. Annu Rev Virol **2:**265-288.

1191    77.    **Roossinck MJ, Sabanadzovic S, Okada R, Valverde RA.** 2011. The remarkable

1192           evolutionary history of endornaviruses. J Gen Virol **92:**2674-2678.

1193    78.    **Li K, Zheng D, Cheng J, Chen T, Fu Y, Jiang D, Xie J.** 2016. Characterization of a

1194           novel *Sclerotinia sclerotiorum* RNA virus as the prototype of a new proposed family

1195           within the order *Tymovirales*. Virus Res **219:**92-99.

1196    79.    **Syller J.** 2002. Umbraviruses - the unique plant viruses that do not encode a capsid

1197           protein. Acta Microbiol Pol **51:**99-113.

1198    80.    **Janowski AB, Krishnamurthy SR, Lim ES, Zhao G, Brenchley JM, Barouch DH,**

1199           **Thakwalakwa C, Manary MJ, Holtz LR, Wang D.** 2017. Statoviruses, a novel taxon

1200           of RNA viruses present in the gastrointestinal tracts of diverse mammals. Virology

1201           **504:**36-44.

1202    81.    **Greninger AL, DeRisi JL.** 2015. Draft genome sequence of tombunodavirus UC1.

1203           Genome Announc **3:**e00655-00615.

1204    82.    **Dokland T, Walsh M, Mackenzie JM, Khromykh AA, Ee K-H, Wang S.** 2004. West

1205           Nile virus core protein; tetramer structure and ribbon formation. Structure **12:**1157-1163.

1206    83.    **Ma L, Jones CT, Groesch TD, Kuhn RJ, Post CB.** 2004. Solution structure of dengue

1207           virus capsid protein reveals another fold. Proc Natl Acad Sci U S A **101:**3414-3419.

1208    84.    **Kuhn RJ, Rossmann MG.** 2005. Structure and assembly of icosahedral enveloped RNA

1209           viruses. Adv Virus Res **64:**263-284.

1210    85.    **Ahola T, Karlin DG.** 2015. Sequence analysis reveals a conserved extension in the

1211           capping enzyme of the alphavirus supergroup, and a homologous domain in nodaviruses.

1212           Biol Direct **10:**16.

1213    86.    **Koonin EV.** 1993. Computer-assisted identification of a putative methyltransferase

1214        domain in NS5 protein of flaviviruses and λ2 protein of reovirus. J Gen Virol **74 ( Pt**

1215        **4):**733-740.

1216    87.    **Liu L, Dong H, Chen H, Zhang J, Ling H, Li Z, Shi P-Y, Li H.** 2010. Flavivirus RNA

1217        cap methyltransferase: structure, function, and inhibition. Front Biol (Beijing) **5:**286-303.

1218    88.    **Dolja VV, Kreuze JF, Valkonen JPT.** 2006. Comparative and functional genomics of

1219        closteroviruses. Virus Res **117:**38-51.

1220    89.    **Kobayashi K, Atsumi G, Iwadate Y, Tomita R, Chiba K-i, Akasaka S, Nishihara M,**

1221        **Takahashi H, Yamaoka N, Nishiguchi M, Sekine K-T.** 2013. Gentian Kobu-sho-

1222        associated virus: a tentative, novel double-stranded RNA virus that is relevant to gentian

1223        Kobu-sho syndrome. J Gen Plant Pathol **79:**56-63.

1224    90.    **Teixeira M, Sela N, Ng J, Casteel CL, Peng H-C, Bekal S, Girke T, Ghanim M,**

1225        **Kaloshian I.** 2016. A novel virus from *Macrosiphum euphorbiae* with similarities to

1226        members of the family *Flaviviridae*. J Gen Virol **97:**1261-1271.

1227    91.    **Mata CP, Luque D, Gómez-Blanco J, Rodríguez JM, González JM, Suzuki N,**

1228        **Ghabrial SA, Carrascosa JL, Trus BL, Castón JR.** 2017. Acquisition of functions on

1229        the outer capsid surface during evolution of double-stranded RNA fungal viruses. PLoS

1230        Pathog **13:**e1006755.

1231    92.    **Castón JR, Luque D, Gómez-Blanco J, Ghabrial SA.** 2013. Chrysovirus structure:

1232        repeated helical core as evidence of gene duplication. Adv Virus Res **86:**87-108.

1233    93.    **Pan J, Dong L, Lin L, Ochoa WF, Sinkovits RS, Havens WM, Nibert ML, Baker TS,**

1234        **Ghabrial SA, Tao YJ.** 2009. Atomic structure reveals the unique capsid organization of

1235        a dsRNA virus. Proc Natl Acad Sci U S A **106:**4225-4230.

1236  94.  **Abrescia NGA, Bamford DH, Grimes JM, Stuart DI.** 2012. Structure unifies the viral

1237        universe. Annu Rev Biochem **81:**795-822.

1238  95.  **El Omari K, Meier C, Kainov D, Sutton G, Grimes JM, Poranen MM, Bamford DH,**

1239        **Tuma R, Stuart DI, Mancini EJ.** 2013. Tracking in atomic detail the functional

1240        specializations in viral RecA helicases that occur during evolution. Nucleic Acids Res

1241        **41:**9396-9410.

1242  96.  **Sutton G, Grimes JM, Stuart DI, Roy P.** 2007. Bluetongue virus VP4 is an RNA-

1243        capping assembly line. Nat Struct Mol Biol **14:**449-451.

1244  97.  **Yu X, Jiang J, Sun J, Zhou ZH.** 2015. A putative ATPase mediates RNA transcription

1245        and capping in a dsRNA virus. Elife **4:**e07901.

1246  98.  **Amarasinghe GK, Ceballos NGA, Banyard AC, Basler CF, Bavari S, Bennett AJ,**

1247        **Blasdell KR, Briese T, Bukreyev A, Caì Y, Calisher CH, Lawson CC, Chandran K,**

1248        **Chapman CA, Chiu CY, Choi K-S, Collins PL, Dietzgen RG, Dolja VV, Dolnik O,**

1249        **Domier LL, Dürrwald R, Dye JM, Easton AJ, Ebihara H, Echevarría JE, Fooks AR,**

1250        **Formenty PBH, Fouchier RAM, Freuling CM, Ghedin E, Goldberg TL, Hewson R,**

1251        **Horie M, Hyndman TH, Jiāng D, Kityo R, Kobinger GP, Kondō H, Koonin EV,**

1252        **Krupovic M, Kurath G, Lamb RA, Lee B, Leroy EM, Maes P, Maisner A, Marston**

1253        **DA, Mor SK, Müller T, et al.** 2018. Taxonomy of the order *Mononegavirales*: update

1254        2018. Arch Virol **163:**2283-2294.

1255  99.  **Kormelink R, Garcia ML, Goodin M, Sasaya T, Haenni A-L.** 2011. Negative-strand

1256        RNA viruses: the plant-infecting counterparts. Virus Res **162:**184-202.

1257  100.  **Osaki H, Sasaki A, Nomiyama K, Tomioka K.** 2016. Multiple virus infection in a

1258        single strain of *Fusarium poae* shown by deep sequencing. Virus Genes **52:**835-847.

1259    101.    **Bacharach E, Mishra N, Briese T, Zody MC, Kembou Tsofack JE, Zamostiano R,**

1260            **Berkowitz A, Ng J, Nitido A, Corvelo A, Toussaint NC, Abel Nielsen SC, Hornig M,**

1261            **Del Pozo J, Bloom T, Ferguson H, Eldar A, Lipkin WI.** 2016. Characterization of a

1262            novel orthomyxo-like virus causing mass die-offs of tilapia. MBio **7:**e00431-00416.

1263    102.    **Maes P, Alkhovsky SV, Bào Y, Beer M, Birkhead M, Briese T, Buchmeier MJ,**

1264            **Calisher CH, Charrel RN, Choi IR, Clegg CS, Torre JCdl, Delwart E, DeRisi JL,**

1265            **Bello PLD, Serio FD, Digiaro M, Dolja VV, Drosten C, Druciarek TZ, Du J,**

1266            **Ebihara H, Elbeaino T, Gergerich RC, Gillis AN, Gonzalez J-PJ, Haenni A-L,**

1267            **Hepojoki J, Hetzel U, Hồ T, Hóng N, Jain RK, Vuren PJv, Jin Q, Jonson MG,**

1268            **Junglen S, Keller KE, Kemp A, Kipar A, Kondov NO, Koonin EV, Kormelink R,**

1269            **Korzyukov Y, Krupovic M, Lambert AJ, Laney AG, LeBreton M, Lukashevic IS,**

1270            **Marklewitz M, Markotter W, et al.** 2018. Taxonomy of the family *Arenaviridae* and

1271            the order *Bunyavirales*: update 2018. Arch Virol **163:**2295-2310.

1272    103.    **Mielke-Ehret N, Mühlbach H-P.** 2012. *Emaravirus*: a novel genus of multipartite,

1273            negative strand RNA plant viruses. Viruses **4:**1515-1536.

1274    104.    **Krupovic M.** 2013. Networks of evolutionary interactions underlying the polyphyletic

1275            origin of ssDNA viruses. Curr Opin Virol **3:**578-586.

1276    105.    **Kazlauskas D, Varsani A, Krupovic M.** 2018. Pervasive chimerism in the replication-

1277            associated proteins of uncultured single-stranded DNA viruses. Viruses **10:**187.

1278    106.    **Iranzo J, Koonin EV, Prangishvili D, Krupovic M.** 2016. Bipartite network analysis of

1279            the archaeal virosphere: evolutionary connections between viruses and capsidless mobile

1280            elements. J Virol **90:**11043-11055.

1281  107.  **Krishnamurthy SR, Wang D.** 2018. Extensive conservation of prokaryotic ribosomal

1282  binding sites in known and novel picobirnaviruses. Virology **516:**108-114.

1283  108.  **Boros Á, Polgár B, Pankovics P, Fenyvesi H, Engelmann P, Phan TG, Delwart E,**

1284  **Reuter G.** 2018. Multiple divergent picobirnaviruses with functional prokaryotic Shine-

1285  Dalgarno ribosome binding sites present in cloacal sample of a diarrheic chicken.

1286  Virology **525:**62-72.

1287  109.  **Jinek M, Doudna JA.** 2009. A three-dimensional view of the molecular machinery of

1288  RNA interference. Nature **457:**405-412.

1289  110.  **Wu Q, Wang X, Ding S-W.** 2010. Viral suppressors of RNA-based viral immunity: host

1290  targets. Cell Host Microbe **8:**12-15.

1291  111.  **Krupovic M, Dolja VV, Koonin EV.** 2015. Plant viruses of the *Amalgaviridae* family

1292  evolved via recombination between viruses with double-stranded and negative-strand

1293  RNA genomes. Biol Direct **10:**12.

1294  112.  **Martin RR, Zhou J, Tzanetakis IE.** 2011. Blueberry latent virus: an amalgam of the

1295  *Partitiviridae* and *Totiviridae*. Virus Res **155:**175-180.

1296  113.  **Sabanadzovic S, Valverde RA, Brown JK, Martin RR, Tzanetakis IE.** 2009.

1297  Southern tomato virus: the link between the families *Totiviridae* and *Partitiviridae*. Virus

1298  Res **140:**130-137.

1299  114.  **Sabanadzovic S, Abou Ghanem-Sabanadzovic N, Valverde RA.** 2010. A novel

1300  monopartite dsRNA virus from rhododendron. Arch Virol **155:**1859-1863.

1301  115.  **Mushegian AR, Koonin EV.** 1993. Cell-to-cell movement of plant viruses. Insights

1302  from amino acid sequence comparisons of movement proteins and from analogies with

1303  cellular transport systems. Arch Virol **133:**239-257.

59

1304    116.    **Ghabrial SA, Castón JR, Jiang D, Nibert ML, Suzuki N.** 2015. 50-plus years of

1305            fungal viruses. Virology **479-480:**356-368.

1306    117.    **Hillman BI, Annisa A, Suzuki N.** 2018. Viruses of plant-interacting fungi. Adv Virus

1307            Res **100:**99-116.

1308    118.    **Kotta-Loizou I, Coutts RHA.** 2017. Studies on the virome of the entomopathogenic

1309            fungus *Beauveria bassiana* reveal novel dsRNA elements and mild hypervirulence. PLoS

1310            Pathog **13:**e1006183.

1311    119.    **Coyle MC, Elya CN, Bronski M, Eisen MB.** 2018. Entomophthovirus: an insect-

1312            derived iflavirus that infects a behavior manipulating fungal pathogen of dipterans.

1313            bioRxiv doi:10.1101/371526**:**371526.

1314    120.    **Waldron FM, Stone GN, Obbard DJ.** 2018. Metagenomic sequencing suggests a

1315            diversity of RNA interference-like responses to viruses across multicellular eukaryotes.

1316            PLoS Genet **14:**e1007533.

1317    121.    **Aiewsakun P, Simmonds P.** 2018. The genomic underpinnings of eukaryotic virus

1318            taxonomy: creating a sequence-based framework for family-level virus classification.

1319            Microbiome **6:**38.

1320    122.    **Benson DA, Cavanaugh M, Clark K, Karsch-Mizrachi I, Ostell J, Pruitt KD, Sayers**

1321            **EW.** 2018. GenBank. Nucleic Acids Res **46:**D41-D47.

1322    123.    **Altschul SF, Madden TL, Schäffer AA, Zhang J, Zhang Z, Miller W, Lipman DJ.**

1323            1997. Gapped BLAST and PSI-BLAST: a new generation of protein database search

1324            programs. Nucleic Acids Res **25:**3389-3402.

1325    124.    **Marchler-Bauer A, Derbyshire MK, Gonzales NR, Lu S, Chitsaz F, Geer LY, Geer**

1326            **RC, He J, Gwadz M, Hurwitz DI, Lanczycki CJ, Lu F, Marchler GH, Song JS,**

1327          **Thanki N, Wang Z, Yamashita RA, Zhang D, Zheng C, Bryant SH.** 2015. CDD:

1328          NCBI's conserved domain database. Nucleic Acids Res **43:**D222-226.

1329    125.   **Edgar RC.** 2010. Search and clustering orders of magnitude faster than BLAST.

1330          Bioinformatics **26:**2460-2461.

1331    126.   **Edgar RC.** 2004. MUSCLE: a multiple sequence alignment method with reduced time

1332          and space complexity. BMC Bioinformatics **5:**113.

1333    127.   **Söding J.** 2005. Protein homology detection by HMM-HMM comparison.

1334          Bioinformatics **21:**951-960.

1335    128.   **Sokal RB, Michener CD.** 1958. A statistical method for evaluating systematic

1336          relationships. Univ Kansas Sci Bull **XXXVIII:**1409-1438.

1337    129.   **Price MN, Dehal PS, Arkin AP.** 2010. FastTree 2 - approximately maximum-likelihood

1338          trees for large alignments. PLoS One **5:**e9490.

1339    130.   **Lemoine F, Domelevo Entfellner J-B, Wilkinson E, Correia D, Dávila Felipe M, De

1340          Oliveira T, Gascuel O.** 2018. Renewing Felsenstein's phylogenetic bootstrap in the era

1341          of big data. Nature **556:**452-456.

1342    131.   **Mirdita M, von den Driesch L, Galiez C, Martin MJ, Söding J, Steinegger M.** 2017.

1343          Uniclust databases of clustered and deeply annotated protein sequences and alignments.

1344          Nucleic Acids Res **45:**D170-D176.

1345    132.   **Eddy SR.** 2011. Accelerated profile HMM searches. PLoS Comput Biol **7:**e1002195.

1346    133.   **Frickey T, Lupas A.** 2004. CLANS: a Java application for visualizing protein families

1347          based on pairwise similarity. Bioinformatics **20:**3702-3704.

1348    134.   **Lancichinetti A, Fortunato S.** 2012. Consensus clustering in complex networks. Sci

1349          Rep **2:**336.

1350    135.    **Fortunato S, Hric D.** 2016. Community detection in networks: A user guide. Physics

1351            Rep Rev Sect Physics Lett **659:**1-44.

1352    136.    **Rosvall M, Bergstrom CT.** 2008. Maps of random walks on complex networks reveal

1353            community structure. Proc Natl Acad Sci U S A **105:**1118-1123.

1354    137.    **Lancichinetti A, Radicchi F, Ramasco JJ, Fortunato S.** 2011. Finding statistically

1355            significant communities in networks. PLoS One **6:**e18961.

1356    138.    **Petitjean C, Makarova KS, Wolf YI, Koonin EV.** 2017. Extreme deviations from

1357            expected evolutionary rates in archaeal protein families. Genome Biol Evol **9:**2791-2811.

1358

1359 **Figure legends**

1360

1361 **Figure 1.** Phylogeny of RNA virus RNA-dependent RNA polymerases (RdRps) and reverse

1362 transcriptases (RTs): the main branches (1-5). Each branch represents collapsed sequences of the

1363 respective set of RdRps. The 5 main branches discussed in the text are labeled accordingly. The

1364 bootstrap support obtained by (numerator) and by (denominator) is shown for each internal

1365 branch. LTR, long-terminal repeat.

1366

1367 **Figure 2**. Branch 1 of the RNA virus RNA-dependent RNA polymerases (RdRps): leviviruses

1368 and their relatives. (A) Phylogenetic tree of the virus RdRps showing ICTV-accepted virus taxa

1369 and other major groups of viruses. Approximate numbers of distinct virus RdRps present in each

1370 branch are shown in parentheses. Symbols to the right of the parentheses summarize the

1371 presumed virus host spectrum of a lineage. Green dots represent well-supported branches (≥0.7)

1372 whereas yellow dot corresponds to a weakly supported branch. (B) Genome maps of a

1373 representative set of Branch 1 viruses (drawn to scale) showing major, color-coded conserved

1374 domains. When a conserved domain comprises only a part of the larger protein, the rest of this

1375 protein is shown in light gray. The locations of such domains are approximated (indicated by

1376 fuzzy boundaries). CP, capsid protein; MP, movement protein; S3H, superfamily 3 helicase;

1377 SJR1 and SJR2, single jelly-roll capsid proteins of type 1 and 2 (see Figure 7).

1378

1379 **Figure 3.** Branch 2 of the RNA virus RNA-dependent RNA polymerases (RdRps): "picornavirus

1380 supergroup" of the +RNA viruses expanded to include nidoviruses and two groups of dsRNA

1381 viruses, partitiviruses and picobirnaviruses. (A) Phylogenetic tree of the virus RdRps showing

63

1382    ICTV accepted virus taxa and other major groups of viruses. Approximate numbers of distinct

1383    virus RdRps present in each branch are shown in parentheses. Symbols to the right of the

1384    parentheses summarize the presumed virus host spectrum of a lineage. Green dots represent well-

1385    supported branches (≥0.7). Inv., viruses of invertebrates (many found in holobionts making host

1386    assignment uncertain); myco., mycoviruses; uncl., unclassified; vert., vertebrate. (B) Genome

1387    maps of a representative set of Branch 2 viruses (drawn to scale) showing major, color-coded

1388    conserved domains. When a conserved domain comprises only a part of the larger protein, the

1389    rest of this protein is shown in light gray. The locations of such domains are approximated

1390    (indicated by fuzzy boundaries). $3C^{pro}$, 3C chymotrypsin-like protease; CP, capsid protein; E,

1391    envelope protein; En, nidoviral uridylate-specific endoribonuclease (NendoU); Exo, 3'-to-5'

1392    exoribonuclease domain; fCP, capsid protein forming filamentous virions; M, membrane protein;

1393    MD macro domain; MP, movement protein; MT, ribose-2-$O$-methyltransferase domain; N,

1394    nucleocapsid protein; N7, guanine-N7-methyltransferase; Ppro, papain-like protease; SJR1 and

1395    SJR2, single jelly-roll capsid proteins of type 1 and 2; spike, spike protein; S1H, superfamily 1

1396    helicase; S2H, superfamily 2 helicase; S3H, superfamily 3 helicase; VP2, virion protein 2; Z, Zn-

1397    finger domain; Spro, serine protease; P3, protein 3. Distinct hues of same color (e.g., green for

1398    MPs) are used to indicate the case when proteins that share analogous function are not

1399    homologous.

1400

1401    **Figure 4**. Branch 3 of the RNA virus RNA-dependent RNA polymerases (RdRps): Alphavirus

1402    superfamily, radiation of related tombusviruses, nodaviruses and unclassified viruses, and

1403    flavivirus supergroup. (A) Phylogenetic tree of the virus RdRps showing ICTV-accepted virus

1404    taxa and other major groups of viruses. Approximate numbers of distinct virus RdRps present in

1405    each branch are shown in parentheses. Symbols to the right of the parentheses summarize the

1406    presumed virus host spectrum of a lineage. Green dots represent well-supported branches ($\geq$0.7),

1407    whereas yellow dots correspond to weakly supported branches. Inv., viruses of invertebrates

1408    (many found in holobionts making host assignment uncertain); myco., mycoviruses; uncl.,

1409    unclassified.. (B) Genome maps of a representative set of Branch 3 viruses (drawn to scale)

1410    showing major, color-coded conserved domains. When a conserved domain comprises only a

1411    part of the larger protein, the rest of this protein is shown in light gray. The locations of such

1412    domains are approximated (indicated by fuzzy boundaries). C, nucleocapsid protein; CapE,

1413    capping enzyme; CP-Spro, capsid protein-serine protease; E, envelope protein; fCP, divergent

1414    copies of the capsid protein forming filamentous virions; Hsp70h, Hsp70 homolog; MP,

1415    movement protein; NS, nonstructural protein; nsP2–3, non-structural proteins; Ppro, papain-like

1416    protease; prM, precursor of membrane protein; rCP, capsid protein forming rod-shaped virions;

1417    RiS, RNA interference suppressor; S1H, superfamily 1 helicase; S2H, superfamily 2 helicase;

1418    SJR2, single jelly-roll capsid proteins of type 2; Spro, serine protease; vOTU, virus OTU-like

1419    protease; NS, nonstructural protein. Distinct hues of same color (e.g., green for MPs) are used to

1420    indicate the cases when proteins that share analogous function are not homologous.

1421

1422    **Figure 5.** Branch 4 of the RNA virus RNA-dependent RNA polymerases (RdRps): dsRNA

1423    viruses of eukaryotes and prokaryotes. (A) Phylogenetic tree of the virus RdRps showing ICTV-

1424    accepted virus taxa and other major groups of viruses. Approximate numbers of distinct virus

1425    RdRps present in each branch are shown in parentheses. Symbols to the right of the parentheses

1426    summarize the presumed virus host spectrum of a lineage. Inv., viruses of invertebrates (many

1427    found in holobionts making host assignment uncertain); myco., mycoviruses; uncl., unclassified..

1428    Green dots represent well-supported branches (≥0.7). (B) Genome maps of a representative set of

1429    Branch 4 viruses (drawn to scale) showing major, color-coded conserved domains. When a

1430    conserved domain comprises only a part of the larger protein, the rest of this protein is shown in

1431    light gray. The locations of such domains are approximated (indicated by fuzzy boundaries).

1432    CapE, capping enzyme; CP, capsid protein; iCP, internal capsid protein; NS, non-structural

1433    protein; NTPase, nucleotide triphosphatase; oCP, outer capsid protein; P, protein; phytoreoS7,

1434    homolog of S7 domain of phytoreoviruses; pHel, packaging helicase; vOTU, virus OTU-like

1435    protease; VP, viral protein; The CPs of totiviruses and chrysoviruses are homologous to iCPs of

1436    reoviruses and cystoviruses (black rectangles).

1437

1438    **Figure 6.** Branch 5 of the RNA virus RNA-dependent RNA polymerases (RdRps): -RNA

1439    viruses. (A) Phylogenetic tree of the virus RdRps showing ICTV-accepted virus taxa and other

1440    major groups of viruses. Approximate numbers of distinct virus RdRps present in each branch

1441    are shown in parentheses. Symbols to the right of the parentheses summarize the presumed virus

1442    host spectrum of a lineage. Green dots represent well-supported branches (≥0.7), whereas yellow

1443    dots correspond to weakly supported branches. Inv., viruses of invertebrates (many found in

1444    holobionts making host assignment uncertain); uncl., unclassified. (B) Genome maps of a

1445    representative set of Branch 5 viruses (drawn to scale) showing major, color-coded conserved

1446    domains. When a conserved domain comprises only a part of the larger protein, the rest of this

1447    protein is shown in light gray. The locations of such domains are approximated (indicated by

1448    fuzzy boundaries). CapE, capping enzyme; CP, capsid protein; EN, "cap-snatching"

1449    endonuclease; GP, glycoprotein; GPC; glycoprotein precursor; HA; hemagglutinin; M, matrix

1450    protein; MP, movement protein; NA; neuraminidase; NP, nucleoprotein; NS, nonstructural

1451    protein; $NS_M$, medium nonstructural protein; NSs, small nonstructural protein; PA, polymerase

1452    acidic protein; PB, polymerase basic protein; vOTU; virus OTU-like protease; VP; viral protein;

1453    Z, zinc finger protein.

1454

1455    **Figure 7.** Sequence similarity networks of SJR-CPs. Protein sequences were clustered by the

1456    pairwise similarity of their hmm profiles using CLANS. Different groups of SJR-CPs are shown

1457    as clouds of differentially colored circles, with the corresponding subgroups labeled as indicated

1458    in the figure. Edges connect sequences with CLANS P-value of $\leq$ 1e-03.

1459

1460    **Figure 8.** Bipartite network of gene sharing in RNA viruses. (A) Groups of related viruses were

1461    identified as the modules of the bipartite genome-gene network (not shown), whereas connector

1462    genes were defined as those genes present in two or more modules with prevalence greater than

1463    65%. The network in A shows viral modules as colored circles (blue, +RNA viruses; green,

1464    dsRNA viruses; red, -RNA viruses), linked to the connector genes (black dots) that are present in

1465    each module. The size of the circles is proportional to the number of genomes in each module.

1466    Shaded ovals indicate statistically significant, 1st-order supermodules that join modules from

1467    taxonomically related groups. (B) Taxonomic analysis of the network modules confirms that

1468    most modules contain viruses from a single family and that families do not tend to split among

1469    modules. (C) High-order supermodules of the RNA virus network, obtained by iteratively

1470    applying a community detection algorithm on the bipartite network of (super)modules and

1471    connector genes. GP, glycoprotein; GT, guanylyltransferase; MT, methyltransferase; NCP,

1472    nucleocapsid; RdRp, RNA-dependent RNA polymerase; SF, superfamily; SJR-CP, single jelly

1473    roll capsid protein.

1474

1475 **Figure 9.** Quantitative analysis of the host range diversity of RNA viruses. The entropy of host

1476 ranges is plotted against the ultrameterized tree depth for the 5 main branches of the RdRp

1477 phylogeny (see Figure 1).

1478

1479 **Figure 10.** A general scenario of RNA virus evolution. The figure is a rough scheme of the key

1480 steps of RNA virus evolution inferred in this work. The main branches from the phylogenetic

1481 tree of the RdRps are denoted 1 to 5 as in Figure 1. Only the RdRp, CP genes and helicase (S1H,

1482 S2H and S3H for the helicases of superfamilies 1, 2 and 3, respectively) are shown

1483 systematically. The helicases appear to have been captured independently and in parallel in 3

1484 branches of +RNA viruses, facilitating the evolution of larger, more complex genomes.

1485 Additional genes, namely, Endo and Exo (for endonuclease and exonuclease, respectively) and

1486 Hsp70h (heat shock protein 70 homolog), are shown selectively, to emphasize the increased

1487 genome complexity, respectively, in *Nidovirales* and in *Closteroviridae*. The virion architecture

1488 is shown schematically for each included group of viruses. Icosahedral capsids composed of

1489 unelated CPs are shown by different colors (see text for details). The question mark at the

1490 hypothetical ancestral eukaryotic RNA virus indicates the uncertainty with regard to the nature

1491 of the host (prokaryotic or eukaryotic) of this ancestral form. The block arrow at the bottom

1492 shows the time flow and the complexification trend in RNA virus evolution.

1493 **SUPPLEMENTAL MATERIAL**

1494

1495 **Supplemental Figure S1**. Iterative clustering-alignment-phylogeny procedure.

1496

1497 **Supplemental Figure S2**. (A) Gain and loss of capsid (nucleocapsid) proteins and movement

1498 proteins. (B) Gain and loss of key enzymes. Each branch is shown by a triangle which represents

1499 collapsed sequences of the respective set of RdRps. The 5 main branches discussed in the text

1500 are indicated. The genes (domains) are denoted by distinct shapes shown at the bottom of each

1501 panel. Different colors show distinct families of proteins with similar functions (helicases,

1502 proteases, capping enzymes) and different hues of the same color show distinct subfamilies of

1503 the same family that are likely to have been acauired independently. The symbols at internal

1504 nodes denoted the inferred point of origin (gain) of the respective gene. Empty shapes show loss

1505 of the respective genes. Shapes placed inside triangles indicate the presence of the respective

1506 gene in a subset of the viruses in the respective branch.

1507

1508 **Supplemental Figure S3.** Sequence similarity networks of SJR-CPs. Protein sequences were

1509 clustered by the pairwise similarity of their hmm profiles using CLANS. Different groups of

1510 SJR-CPs are shown as clouds of differentially colored circles, with the corresponding subgroups

1511 labeled as indicated in the figure. Edges connect sequences with CLANS P-value ≤ 1e-10.

1512

1513 **Supplemental Data set S1.** List of viruses.

1514

1515 **Supplemental Data set S2**. Phylogenetic tree for the global set of RdRps; Newick format.

1516

1517    **Supplemental Data set S3**. Phylogenetic trees for the global set of RdRp representatives;

1518    Newick format. A. Reconstructed using PhyML. B. Reconstructed using RAxML.

1519

1520    **Supplemental Data set S4**. Phylogenetic trees for RdRp representatives; Newick format.

1521    A. Branch 1 representatives.

1522    B. Branch 2 representatives.

1523    C. Branch 3 representatives.

1524    D. Branch 4 representatives.

1525    E. Branch 5 representatives

1526

1527    **Supplemental Data set S5.** Proteins domains encoded by the virus genomes.

1528

1529    Additional data can be found at the following directory:

1530    ftp://ftp.ncbi.nlm.nih.gov/pub/wolf/_suppl/rnavir18

**A  Branch 1 (299)**

"Mitoviruses" (98)
"Ourmiaviruses" (75)
"Narnaviruses" (31)
*Leviviridae* (69)
"Levi-like viruses" (26)

— Positive-sense RNA viruses

Invertebrate   Fungi   Plants   Bacteria   Protists

**B**

Genome length (kb)

Cryphonectria mitovirus 1 (*Mitovirus*)
RdRp

Běihǎi narna-like virus 8 ("ourmiavirus")
RdRp  SJR2

Wēnzhōu narna-like virus 5 ("ourmiavirus")
RdRp  SJR2

Wēnzhōu narna-like virus 9 ("ourmiavirus")
S3H  RdRp

Ourmia melon virus (*Ourmiavirus*)
RdRp  RNA 1  MP  RNA 2  SJR1  RNA 3

Saccharomyces 23S RNA narnavirus (*Narnavirus*)
RdRp

Escherichia phage Qβ (*Leviviridae: Allolevivirus*)
A Protein  CP  RdRp

AVE000 (metagenomic, "levi-like virus")
Maturation  RdRp

**A**  **Branch 2: "Picornavirus supergroup" (1901)**



Legend:
— Double-stranded RNA viruses
— Positive-sense RNA viruses

🕷 Invertebrate  🐕 Vertebrate  🍄 Fungi  🌿 Plants  🦠 Bacteria  Protists

**B**

**A** Branch 3 (1,208)

**B**

**A** Branch 4 (346)

Branch 5

*Reoviridae*
- *Spinareovirinae* (22)
- *Sedoreovirinae* (82)

*Cystoviridae* (7)

Uncl. myco., inv. (29)
Diatom colony-associated dsRNA virus 16 **?**
*Megabirnaviridae* (11)
*Chrysoviridae* (28)
"Botybirnaviruses" (15)
*Totiviridae* / *Quadriviridae* (102)
"Alternaviruses" (3)
*Giardiavirus* / Uncl. inv. (46)

— Negative-sense RNA viruses
— Double-stranded RNA viruses

Invertebrate   Vertebrate   Fungi   Plants   Bacteria   Protists

**B**

Genome length (kb)

Golden shiner virus
(*Reoviridae: Spinareovirinae: Aquareovirus*)
CapE — RNA 1
RdRp — RNA 2
VP3 (iCP) — RNA 3
NS1 — RNA 4
VP5 (NTPase) — RNA 5
VP4 (oCP) — RNA 6
NS5 NS4 — RNA 7
VP6 — RNA 8
NS2 — RNA 9
VP7 — RNA 10
NS3 — RNA 11

Pseudomonas phage phi6
(*Cystoviridae: Cystovirus*)
P8 (oCP)
P12 P9 P5a P5b — S segment
P10
P6 P3 (spike) — M segment
P13
P7 P2 RdRp pHel — L segment
P1 (iCP)

Penicillium chrysogenum virus
(*Chrysoviridae: Chrysovirus*)
RdRp — RNA 1
CP — RNA 2
P3 phytoreo S7 — RNA 3
P4 vOTU — RNA 4

Saccharomyces cerevisiae virus L-A
(*Totiviridae: Totivirus*)
CP RdRp

**A** Branch 5 (859)



*Bunyavirales* (394)

*Mononegavirales* (348)

0.5

— Negative-sense RNA viruses

Invertebrate   Vertebrate   Fungi   Plants   Protists

**B**



Genome length (kb)

Cluster 1

Cluster 2

*Ourmiavirus*

*Tymoviridae*

*Astroviridae*

"sobemo-like"

*Hepeviridae*

*Solemoviridae*

"tombus-like"

*Luteoviridae*

"narna-like"

*Picornavirales*

"tombus-like"

"zhàoviruses"

"sobemo-like"

*Alvernaviridae*

*Barnaviridae*

"narna-like"

*Caliciviridae*

"wèiviruses"

"statoviruses"

*Nodaviridae*

*Tombusviridae*

"tombus-like"

tetraviruses

*Nodaviridae*

"tombus-like"

"sobemo-like"

"narna-like"

"sobemo-like"

"zhàoviruses"

"wèiviruses"

"yuèviruses"

*Bromoviridae*

**Prokaryotic hosts**

**Eukaryotic hosts**

RT → RdRp

① Leviviruses

Mitoviruses → Narnaviruses

Mitoviruses → Ourmiaviruses

? SJR — ancestral eukaryotic RNA virus

② Picorna-like (S3H) | Picobirna-like (S2H) | Poty-like (S2H) | Nido-like (S1H, Endo Exo)

③ Tombus-/noda-like | Flavi-like (S2H) | Alpha-like (S1H) | Clostero-like (S1H Hsp70h)

④ Toti-like | Reo-like → ⑤ Negarnaviricota

**TIME, COMPLEXIFICATION**

Evolution of the ancestral RdRp from a reverse transcriptase → Emergence of eukaryotes → Origin of multicellularity → Diversification of invertebrates; terrestrialization → HVT from invertebrates to plants and vertebrates → Diversification and complexification in different virus lineages