# Artificial intelligence redefines RNA virus discovery

Xin Hou[1,15], Yong He[2,15], Pan Fang[2], Shi-Qiang Mei[1], Zan Xu[2], Wei-Chen Wu[1], Jun-Hua Tian[3], Shun Zhang[2], Zhen-Yu Zeng[2], Qin-Yu Gou[1], Gen-Yang Xin[1], Shi-Jia Le[1], Yin-Yue Xia[4], Yu-Lan Zhou[5], Feng-Ming Hui[6,7], Yuan-Fei Pan[8], John-Sebastian Eden[9], Zhao-Hui Yang[10], Chong Han[11], Yue-Long Shu[12], Deyin Guo[13], Jun Li[14], Edward C. Holmes[9], Zhao-Rong Li[2⬚] & Mang Shi[1⬚]

[1]State key laboratory for biocontrol, the Centre for Infection and Immunity Studies, School of Medicine, Shenzhen Campus of Sun Yat-sen University, Sun Yat-sen University, Shenzhen, China;

[2]Industry Research and Development Department, Alibaba Cloud Intelligence, Alibaba Group, Hangzhou, China;

[3]Wuhan Centers for Disease Control and Prevention, Wuhan, China;

[4]Polar Research Institute of China, Shanghai, China;

[5]Department of Nursing, The Fifth Affiliated Hospital, Sun Yat-sen University, Zhuhai, China;

[6]School of Geospatial Engineering and Science, Southern Marine Science and Engineering Guangdong Laboratory (Zhuhai), Sun Yat-sen University, Zhuhai, China;

[7]Key Laboratory of Comprehensive Observation of Polar Environment, Ministry of Education, Sun Yat-sen University, Zhuhai, China;

[8]Ministry of Education Key Laboratory of Biodiversity Science and Ecological Engineering, National Observations and Research Station for Wetland Ecosystems of the Yangtze Estuary, Institute of Biodiversity Science and Institute of Eco-Chongming, School of Life Sciences, Fudan University, Shanghai, China;

[9]Sydney Institute for Infectious Diseases, School of Medical Sciences, The University of Sydney, Sydney, NSW 2006, Australia;

[10]College of Life Sciences, Zhejiang University, Hangzhou, China;

[11]School of life science, Guangzhou University, Guangzhou, China;

[12]School of Public Health (Shenzhen), Shenzhen campus of Sun Yat-sen University, Sun Yat-

30    sen University, Shenzhen, China;

31    [13]Guangzhou National Laboratory, Guangzhou International Bio-Island, Guangzhou, China;

32    [14]Department of Infectious Diseases and Public Health, Jockey Club College of Veterinary

33    Medicine and Life Sciences, City University of Hong Kong, Hong Kong, China.

34

35

36    [15]These authors contributed equally

37

38    [✉]Corresponding authors:

39    Zhao-Rong Li (zhaorong.lzr@alibaba-inc.com)

40    Mang Shi (shim23@mail.sysu.edu.cn)

41 **Abstract**

42 RNA viruses are diverse components of global ecosystems. The metagenomic identification

43 of RNA viruses is currently limited to those with sequence similarity to known viruses, such

44 that highly divergent viruses that comprise the "dark matter" of the virosphere remain

45 challenging to detect. We developed a deep learning algorithm – LucaProt – to search for

46 highly divergent RNA-dependent RNA polymerase (RdRP) sequences in 10,487 global meta-

47 transcriptomes. LucaProt integrates both sequence and structural information to accurately

48 and efficiently detect RdRP sequences. With this approach we identified 180,571 RNA viral

49 species and 180 superclades (viral phyla/classes). This is the broadest diversity of RNA

50 viruses described to date, including many viruses undetectable using BLAST or HMM

51 approaches. The newly identified RNA viruses were present in diverse ecological niches,

52 including the air, hot springs and hydrothermal vents, and both virus diversity and abundance

53 varied substantially among ecological types. We also identified the longest RNA virus

54 genome (nido-like) observed so far, at 47,250 nucleotides, and expanded the diversity of

55 RNA bacteriophage to more than ten phyla/classes. This study marks the beginning of a new

56 era of virus discovery, with the potential to redefine our understanding of the global

57 virosphere and reshape our understanding of virus evolutionary history.

58

59 **Key Words:**

60 RNA virus; Virome; Virus discovery; Deep learning; Meta-transcriptomes; Phylogeny

61

## Introduction

RNA viruses infect a huge array of host species. Despite this ubiquity, their role as an essential component of global ecosystems has only recently been recognized thanks to systematic and large-scale virus discovery projects performed in animals[1,2], plants[3], fungi[4], marine[5], and soil environments[6]. A common feature of these studies is that they were based on the analysis of RNA-dependent RNA polymerase (RdRP) sequences, a canonical component of RNA viruses. Combined, they have resulted in the discovery of tens of thousands of new virus species, leading to at least a ten-fold expansion of the virosphere and the addition of five new phyla of RNA viruses, including the "Taraviricota"[5]. Similarly, data mining exercises that reanalyzed over $10^9$ meta-transcriptomic contigs associated with diverse ecosystems have identified several divergent clades of RNA bacteriophage[7]. Despite such significant progress in filling the gaps of RNA virus diversity through ecological sampling and sequencing, our understanding of the full spectrum of the RNA virosphere is likely limited[8,9]. This is in part because the BLAST-based sequence similarity searching approaches used to discover new RNA virus sequences have limitations in detecting highly divergent RdRPs[10], while the profile alignment (i.e., HMM) based approach misses a significant proportion of viruses due to a high false-negative rate[11]. To efficiently uncover the full range of RNA virus diversity, the development of novel strategies is therefore essential.

Over the past decade, artificial intelligence (AI) related approaches, especially deep learning algorithms, have had a huge impact on various research fields in the life sciences, including molecular docking, compound screening and interaction, protein structure prediction and functional annotation, and the modelling of infectious diseases[12-17]. These advancements can be attributed to the advantages of deep learning algorithms over classic bioinformatic approaches, including greater accuracy, better performance, less feature engineering, flexible models, and self-learning capabilities[18,19]. Recently, deep learning approaches, such as CHEER, VirHunter, Virtifier and RNN-VirSeeker have also been developed and applied to identify viruses from genomic and metagenomic data[20-23]. However, many of these approaches rely on nucleotide sequence information without incorporating protein sequence or structural information, and are hence less likely to identify highly divergent RNA viruses. The transformer architecture was recently developed and applied to

92    sequence-based protein function predictions, outperforming the convolutional neural

93    networks (CNN) and recurrent neural network (RNN) algorithms implemented in previous

94    virus discovery algorithms[24-26]. As a consequence, transformer architecture can be used to

95    design a better tool to uncover the hidden "dark matter" of highly divergent RNA viruses.

96    Herein, we show how AI can be used to accurately and efficiently detect RNA viruses that are

97    too divergent in sequence to be detected by traditional sequence similarity-based methods, in

98    doing so revealing a hidden world of virus diversity.

99

100    **Results**

101    **Deep learning to reveal the dark matter of the RNA virosphere**

102    We performed all-inclusive searches to reveal the entirety of RNA virus diversity present in

103    different ecological systems sampled at global scale (Extended Data Fig. 1, Supplementary

104    Table 1 and 2). Accordingly, a total of 10,487 meta-transcriptomes (51 Tb of sequencing data)

105    were assembled, which resulted in more than 1,368 million contigs and 872 million predicted

106    proteins. Based on this data set, potential viral RdRPs were revealed and cross-validated

107    using two different strategies (Fig. 1, Extended Data Fig. 2-4). The major AI algorithm used

108    here (i.e., "LucaProt") is a deep learning, transformer-based model established based on

109    sequence and structural features of 5,979 well-characterized RdRPs and 229,434 non-RdRPs.

110    LucaProt had high accuracy (0.03% false positives) and specificity (0.20% false negatives)

111    on the test data set (Fig. 1b, Extended Data Fig. 4). Independently to the deep-learning

112    approach, we applied a more conventional approach (i.e., "ClstrSearch") that clustered all

113    proteins based on their sequence homology and then used BLAST or HMM models to

114    identify any resemblance to viral RdRPs or non-RdRP proteins. The latter approach is

115    distinguished from previous BLAST or HMM based approaches because it queries on protein

116    clusters (i.e., alignments) instead of individual sequences, which greatly reduced both the

117    false positive and negative rates of virus identification.

118        By merging the results of the two search strategies we discovered 513,134 RNA viral

119    contigs, representing 180,571 RNA viral species (i.e., > 90% RdRP identity), and 180 RNA

120    viral superclades at the phylum level taxonomic rank (Fig. 1, Supplementary Table 3 and see

121    Methods). Among these, 512,691 viral contigs (0.04% of total contigs) and 157 superclades

122    (87.2%) were revealed by both "LucaProt" and "ClstrSearch", whereas 443 contigs and 23

123    superclades were only predicted by "LucaProt". Both strategies out-performed previous

124    attempts at RNA virus discovery from ocean[5], soil[6], and more diverse ecosystems[7] (Fig. 1c).

125    Indeed, "LucaProt" was able to identify 98.2% ~ 99.9% of RdRPs discovered in these

126    previous studies, even though none were used in either training or testing of the models

127    (Extended Data Fig. 5). To ensure the robustness and introduce innovative findings from the

128    AI approach, we jointly applied the two strategies and merged the results; this enabled us to

129    identify 93,580 viral species and 59 novel superclades, and resulted in a 9-fold expansion in

130    RNA virus diversity (Fig. 1c). This was reflected in the expansion of both existing viral

131    superclades and the identification of new superclades unlikely to be discovered by sequence

132    homology and HMM based approaches alone (Fig. 1d).

133        All the RNA viral sequences discovered here were organized into clusters and

134    superclades without the influence of the current virus classification system[27,28]. These

135    superclades were then placed back onto the classification system at the phylum (such as

136    phylum Lenarviricota in the case of the Narna-Levi superclade) or class (such as the

137    Stelpaviricetes, Alsuviricetes, Flasuviricetes classes for the Astro-Poty, Hepe-Virga, Flavi

138    superclades) levels (Supplementary Table 4)[28]. Notably, however, the virus superclades

139    comprised much greater phylogenetic diversity than their corresponding phyla/classes. Also

140    of note was that our data did not conform to several of the higher taxonomic ranks, such as

141    the phyla Duplornaviricota and Negarnaviricota, which were now too broad to be regarded as

142    single phyla. Indeed, even the Markov cluster algorithm (MCL) approach, on which the

143    existing virus classification scheme is derived[29,30], fails to re-group these expanded classes

144    into the existing phyla[5].

145

**Verification and confirmation of newly identified viral superclades**

147    That the 180 RNA viral superclades identified represented RNA-based organisms was

148    verified by multiple lines of evidence. At the sequence level, two criteria were used to

149    establish a viral superclade: a lack of homology to cellular proteins and the presence of key

150    RdRP motifs (Fig. 2a). Furthermore, the majority (157/180) of the newly identified

151    superclades shared a variable degree of sequence homology with existing RdRPs (i.e.,

152    BLAST *e*-value ≤ 1*E*-3 and/or had HMM model score ≥ 10). The exception were 23

153    superclades that had no detectable homology to viral RdRPs and therefore named as "AI-

154    specific" superclades (Fig. 2a, Extended Data Fig. 6, Supplementary Table 5). To justify the

155    computational prediction, we performed simultaneous DNA and RNA extraction and

156    sequencing to examine whether the viral superclades identified here also exist in DNA form.

157    This analysis revealed that only RNA sequencing reads were mapped to contigs associated

158    with viral RdRPs, whereas both RNA and DNA sequencing reads were mapped to contigs

159    associated with DNA viruses, reverse-transcriptase (RT), and cellular organisms (Fig. 2b,

160    Extended Data Fig. 7-9). These results were further confirmed by a more sensitive RT-PCR

161    approach which showed that none of the sequences encoding viral RdRP were detected in the

162    DNA extractions, suggesting that these viral superclades were *bona fide* RNA organisms (Fig.

163    2c, Extended Data Fig. 7b). Finally, we performed 3D alignment analysis (newly identified

164    viral RdRPs compared with known viral RdRPs, eukaryotic RdRPs, eukaryotic DdRPs and

165    RT) to determine the degree of structure similarity among them (Fig. 2d). The novel viral

166    RdRP superclades (including AI-specific ones) bore at least three signature motifs that gave

167    them much higher structural similarity to known viral RdRPs than their cellular counterparts.

168

169    **Genomic structures reveal modularity and flexibility within the RNA virosphere**

170    We next analyzed the composition and structure of potential RNA virus genomes identified in

171    this study. The length of the RdRP-encoding genomes or genome segments differed markedly

172    within and between viral superclades, although most were centered around 2,569 nt (Fig. 3).

173    Notably, our data set contained some extremely long RNA virus genomes identified from soil

174    that belonged to the Nido-like superclade: the length of one of these, at 47.3 kb, exceeded

175    *Planarian secretory cell nidovirus* (41.2 kb)[31] as the longest RNA virus genome identified to

176    date (Fig. 3c, Extended Data Fig. 10 Supplementary Table 6). In addition to the RdRP, we

177    characterized the remaining proteins encoded by the newly identified virus genomes. While

178    most of these predicted proteins had no homologs in the existing databases, we identified

179    some that were related to structural (i.e., coat, capsid, glycoprotein and envelope proteins,

180    amongst others) and non-structural (i.e., helicase, protease, methyltransferase, movement

181 protein, immune or host-related regulatory proteins, amongst others) proteins from known

182 viruses (Fig. 3d, Extended Data Fig. 11). Importantly, the presence of these additional virus

183 proteins in newly identified supergroups provided further evidence that these were *bona fide*

184 RNA viruses. Furthermore, that the occurrence of these proteins was incongruent with the

185 groupings of RdRPs (Fig. 3e) suggests that RNA virus genomes have a modular-like

186 configuration, transferring proteins across taxonomic groups. This was in line with the

187 dramatic changes in genome structure (genome length, gene organization, ORF numbers, and

188 segmentation) observed among related viruses, such that no prototype genome structures

189 could be defined for each group or supergroup (Fig. 3e).

190

191 **Expanded phylogenetic diversity of RNA viruses**

192 The enormous expansion in the RNA virosphere described here was also reflected in both the

193 growing size of known virus groups and the addition of entirely new groups (Fig. 4). For

194 existing supergroups, the viruses newly described here were distinguished from those

195 identified previously such that they formed unique clusters at more ancestral positions in the

196 phylogenetic trees (Fig. 4). Interestingly, some previously smaller sized viral groups with

197 limited diversity – the Astro-Poty, Hypo, Yan and Cysto – expanded to become large viral

198 groups comprising substantial genetic diversity (Fig. 4). Several newly identified supergroups

199 were also revealed to have high levels of phylogenetic diversity, including SC022 (8,128

200 species), SC024 (3,682 species), and SC37 (1,772 species), highlighting the limitations in

201 previous attempts to identify highly divergent groups of RNA viruses. Following our analysis,

202 the supergroups with the greatest number of species were the Narna-Levi (64,667 species),

203 Picorna-Calici (23,430 species), and Tombus-Nada (16,798 species).

204      In addition to greatly expanding virus genetic diversity, this study identified more virus

205 groups associated with bacterial hosts than the leviviruses, cystoviruses, and the members of

206 Partiti-Picobirna supercluster known previously[7]. Specifically, we identified bacterial viruses

207 within the Narna-Levi, Hepe-Virga, and SC037 supergroups whose sequences were

208 recognized and "recorded" by the bacterial CRISPR system. Furthermore, based on proteins

209 associated with bacterial infection (i.e., Lysis, Prok-E2, and Prok-RING), we inferred

210 potential bacterial RNA viruses in the Tombus-Noda, Yan, and SC022 supergroups

211   (Supplementary Table 7). As a consequence, those RNA viruses associated with bacteria has

212   expanded to ten supergroups, and these numbers are likely to further increase given our

213   limited knowledge of host associations for most of the viruses in this study.

214

215   **Ecological structure of the global RNA virome**

216   Our study investigated the RNA virome of 10,487 ecological samples, revealing the

217   ubiquitous presence of RNA viruses across diverse ecological types (48 categories) and in

218   1,837 locations globally. Despite repeated efforts to uncover the RNA virus diversity from

219   such ecological samples[5-7], a large proportion of the viruses detected here were entirely novel

220   (Fig. 5a). Indeed, the rate of RNA virus discovery did not plateau (Fig. 5b), suggesting that

221   the global space of RNA virus diversity remains largely under-characterized, with a

222   particularly rapid increase in soil (Fig. 5b).

223   We compared alpha diversity (measured by the Shannon index) and abundance levels

224   (measured by the number of reads per million total non-rRNA reads, i.e., RPM) of the RNA

225   virome among diverse ecological types, revealing enormous variation (Fig. 5c,

226   Supplementary Table 8). In general, average alpha diversity was highest in leaf litter, estuary,

227   freshwater, and wetland environments, whereas virus abundance was highest in freshwater,

228   marine sediment, and rhizosphere systems, whose average RPMs were between 12466.9 and

229   26617.3 (Fig. 5c). In contrast, the lowest average diversity and abundance were observed in

230   halite and subsurface environments (Fig. 5c), which as expected as these samples were

231   particularly low in biomass (i.e., host cells). For extreme ecological types such as hot springs

232   and hydrothermal vents, the associated RNA viruses were characterized by low diversity but

233   moderate abundance (1528.9 ~ 3726.9 average RPM) (Fig. 5c). It is also worth noting that

234   the new viral superclades established in this study were mostly identified from aquatic and

235   sediment samples, with few from vertebrate and invertebrate animal samples (Fig. 5c).

236   Our results further revealed the prevalence and abundance levels of single viral species

237   across different ecological types (Fig. 5d), including some that could be considered

238   ecological generalists. For example, members of the Narna-Levi, Partiti-Picobirna and

239   Picorna superclades as well as Superclade022 were among the prominent generalist RNA

240   viruses and found in more than 45 ecological types (Extended Data Fig. 12). Conversely, the

241    majority (85.9%) of the viruses discovered here only occurred in a single ecological type.

242    Finally, we also identified "marker' virus species for each ecotype, which appeared at high

243    prevalence and abundance in one ecological type but not in the others (Fig. 5d). Among these,

244    *Partiti-Picobirna sp.* 4207 and *Partiti-Picobirna sp.* 9871 were associated with hot springs

245    and *Tombus-Noda sp.* 2280 and *Superclade026 sp.* 2292 were associated with hydrothermal

246    vents, suggesting their important role in these ecosystems.

247

248    **Discussion**

249    Our understanding of the genetic diversity of the RNA virosphere, and hence of RNA virus

250    ecology and evolution in general, is greatly hampered by the inability to accurately identify

251    the highly divergent "dark matter" of viruses[32,33]. Indeed, the conventional way to discover

252    RNA viruses has relied heavily on the utility of sequence similarity comparisons and the

253    completeness of sequence databases[11,32]. To address these issues, we developed a data-driven

254    deep learning model (i.e., LucaProt) that overcome these shortcomings while outperforming

255    conventional approaches in accuracy, efficiency, and, most importantly, the scope of diversity.

256    Importantly, LucaProt not only incorporated sequence data but also structural information,

257    which is relevant in predicting protein function (in this case of the RdRP)[34]. Without

258    implementing the structural model, our model had only 41.8% and 94.9% specificity and

259    accuracy, respectively, on the testing data set, and could only detect 44.5% of the predicted

260    RdRP proteins. In addition, the advanced transformer architecture incorporated into our

261    model allowed the parallel processing of larger amino acid sequences[35,36], which can easily

262    capture the relationship between residues from distant parts of sequence space, thereby

263    outperforming the CNN and/or RNN encoders implemented in the CHEER, VirHunter,

264    Virtifier and RNN-VirSeeker RNA virus discovery tools (Extended Data Fig. 13)[20-23].

265    Collectively, we have established an AI framework for large-scale RNA virus discovery,

266    which can be easily extended to the accurate description of any biological dark matter.

267        Despite the large expansion in RNA virus diversity documented here, major gaps remain

268    in our understanding of the ecology and evolution of the newly discovered viruses. In

269    particular, nothing is known about the hosts of the viruses identified, including that with the

270    longest virus genome identified to date. It is possible that the viral clades and superclades

271     identified here were largely associated with diverse microbial eukaryotic hosts, given that the

272     majority of current known RNA viruses infect eukaryotes[37,38] and microbial eukaryotes exist

273     in great abundance and diversity in natural environments[39,40]. Nevertheless, it is also likely

274     that a substantial proportion of the novel viruses discovered are associated with bacterial (and

275     perhaps archaeal) hosts[41-43]. Indeed, based on this and previous studies[7], more than ten

276     superclades contained RNA viruses likely associated with bacteria. Importantly, the presence

277     of RNA bacteriophages in multiple RNA viral superclades underlines the evolutionary

278     connection between RNA viruses from bacterial and eukaryotic hosts. If viewed through the

279     lens of virus-host co-divergence[1,2,44], such a link between bacterial and eukaryotic hosts

280     suggests that the evolutionary history of RNA viruses is at least as long, if not longer, than

281     that of the cellular organisms.

282

283     **Methods**

284     **Samples and data sets**

285     This study comprised the meta-transcriptomic analysis of 10,487 samples for RNA virus

286     discovery. The majority of the samples (n = 10,437) were mined from the NCBI Sequence

287     Read Archive (SRA) database (https://www.ncbi.nlm.nih.gov/sra) between January 16 -

288     August 14, 2020. We targeted samples collected from a wide range of environmental types

289     globally (Extended Data Fig. 1), including: aquatic (such as marine, riverine and lake water),

290     soil (such as sediment, sludge and wetland), host-related (such as biofilm, wood decay, and

291     rhizosphere), and extreme environmental samples (such as hydrothermal vent, hypersaline

292     lake and salt marsh), that were subject to high quality meta-transcriptomics sequencing.

293     Furthermore, the samples included in this study were subject to high-quality short-read

294     sequencing (i.e., utilizing Illumina sequencing platforms), had between 35.1-204.1 Gbp raw

295     sequencing data output, and were not enriched for any specific types of microbial organisms.

296     For highly abundant environmental types, such as "soil" and "marine", representative

297     samples were selected to include as many projects (i.e., independent studies), geographic

298     locations and ecological niches as possible.

299          In addition to data mined from the SRA database, we obtained 50 samples from

300     Antarctica and China for RNA virus discovery and confirmation. The sample types included

301  marine (N = 5), freshwater (N = 12), soil (N = 19), and sediment (N = 14), of which nine

302  sediment samples were collected at the Ross Sea station in Antarctica between January and

303  February 2022, with the others from Zhejiang, Guangdong, Hubei, and Heilongjiang

304  provinces, China between August and October 2022. For each of these samples, DNA and

305  RNA were simultaneously extracted: the soil and sediment samples were extracted using the

306  RNeasy® PowerSoil® Total RNA Kit and RNeasy® PowerSoil® DNA Elution Kit

307  (QIAGEN, Germany), while the marine and freshwater samples were extracted using the

308  DNeasy® PowerWater® Kit and RNeasy® PowerWater® Kit (QIAGEN, Germany). The

309  extracted nucleic acid was then subject to library construction using NEBNext Ultra RNA

310  Library Prep Kit and NEB Next Ultra DNA Library Prep Kit (LTD.NEB, China) for RNA and

311  DNA samples, respectively. Paired-end (150 bp) sequencing of these libraries was performed

312  using the Illumina NovaSeq 6000 platform (Illumina, the United States).

313      For all 10,487 data sets generated and collected for this study, reads were assembled *de*

314  *novo* into contigs using MEGAHIT v1.2.8[45] with default parameters. Potential encoded

315  proteins were predicted from contigs using ORFfinder v0.4.3

316  (https://ftp.ncbi.nlm.nih.gov/genomes/TOOLS/ORFfinder/linux-i64/; parameters, -g 1, -s 2).

317

318  **Identification of RNA viruses based on deep learning**

319  We developed a new deep learning, transformer-based model, termed "Deep Sequential and

320  Structural Information Fusion Network for Protein Function Prediction" (i.e., LucaProt), that

321  takes into account protein sequence composition and structure information to facilitate the

322  accurate identification of viral RdRPs. The model included five modules: Input, Tokenizer,

323  Encoder, Pooling, and Output (Extended Data Fig. 2e).

324      **Input Layer**：Our model uses the amino acid sequence as input.

325      **Tokenizer Layer:** This module consists of two components. One used a frequent

326  substring algorithm[46], which generated subwords from the input sequence, treated co-

327  occurring amino acids as a whole (namely, "words"), and resulted in a vocabulary with

328  20,000 such "words". The other component broke down each protein sequence into a

329  combination of single amino acid characters which were later used in protein structure

330  modeling.

331  **Encoder Layer:** This module processes the two types of input into sequence and

332  structural representation matrices, respectively. In the case of subword processing, an

333  advanced Transformer-Encoder was applied to obtain the sequence representation matrix,

334  while for structural processing, two strategies were considered to calculate the protein

335  structure representation matrix. The first strategy used a structural model (such as

336  RoseTTAFold[47], AlpahFold[15], and ESMFold[48]) to predict 3D protein structure, calculated the

337  distance between the C-atoms (Alpha-C or Beta-C) of all amino acid residues into a Contact

338  Map matrix, and applied Graph Convolutional Network (GCN)[49] to encode the Contact Map

339  into a representation matrix. The second approach was to directly use the intermediate matrix

340  from the structural model and employ it as the structural representation matrix. This method

341  not only addressed the issue of the insufficient number of 3D structures observed in

342  experiments, but also circumvented the need to perform the encoder, resulting in a cost-

343  effective approach suitable for large-scale implementation such as this study. We therefore

344  adopted the second strategy here and used the faster ESMFold[48] for structural representation.

345  **Pooling Layer:** The previous module obtained the sequence and structure representation

346  matrices. A value-level attention pooling (VLAP) approach[50] was then used to transform

347  these two matrices into two vectors.

348  **Output Layer:** A concatenation operator was used to join the two vectors generated by

349  the pooling layer. A fully connected layer and the sigmoid function (Extended Data Fig. 2e)

350  were then used to generate the probability values between 0.0 and 1.0 as a measure of

351  confidence, and a threshold of 0.5 was used to determine whether it represents viral RNA

352  (Extended Data Fig. 4).

353  **Model Building:** We constructed a data set with 235,413 samples for model building,

354  which included 5,979 positive samples of known viral RdRPs (i.e., the well-curated RdRP

355  database described above), and randomly selected 229,434 negative samples of confirmed

356  non-virus RdRPs (as the positive sample accounts for a very small portion of the total data,

357  we constructed the training data set using the conventional 1:40 ratio of positive to negative

358  data). The non-virus RdRP-like sequences contained proteins from the eukaryotic RNA

359    dependent RNA polymerase (Eu RdRP, N = 2,233), the eukaryotic DNA dependent RNA

360    polymerase (Eu DdRP, N = 1,184), reverse transcriptase (RT, N = 48,490), proteins obtained

361    from DNA viruses (N = 1,533), non-RdRP proteins obtained from RNA viruses (N = 1,574),

362    as well as a wide array of cellular proteins from different functional categories (N = 174,420).

363    We randomly divided the data set into training, validation, and testing sets with a ratio of

364    8.5:1:1, which were used for model fitting, model finalization (based on the best F1-score

365    training iteration), and performance reporting (including accuracy, precision, recall, F1-score,

366    and Area under the ROC Curve (AUC)), respectively (Extended Data Fig. 4).

367         LucaProt identified 792,436 putative RdRP signatures from 144,690,558 proteins. These

368    results were first compared with the RdRPs identified based on sequence homology (see

369    below). RdRPs that were identified only by deep learning algorithms were either incorporated

370    into the superclades using the Diamond blastp program v0.9.25.126[51] with an $e$-value

371    threshold of 1$E$-3, or, if they remained unclassified, were subjected to clustering, merging,

372    and manual alignment inspection as described below to form deep learning specific

373    superclades (the case for 23 superclades).

374

375    **Identification of RNA viruses based on homologous clustered proteins**

376    The first approach to identify RNA viruses was based on sequence and structural similarity to

377    previously known RdRP amino acid sequences (Extended Data Fig. 2a). A total of 871.8

378    million amino acid sequences predicted by ORFfinder (see Samples and data sets) were

379    compared against a well-curated RdRP database (N = 5,979) that contained only those

380    derived from reference RNA virus genomes downloaded from the NCBI GenBank database

381    and their close relatives from vertebrate and invertebrate hosts[1,2]. The comparisons were

382    performed using the Diamond blastp program v0.9.25.126[51], with the $e$-value threshold set at

383    1$E$+5 to identify more divergent RdRP proteins (Extended Data Fig. 2a, Extended Data Fig.

384    3a). This process resulted in 75.3 million hits which were further subjected to homology-

385    based and multi-step clustering (three iterations with 90%, 60%, and 20% amino acid identity,

386    respectively) using CD-HIT v4.8.1 (https://github.com/weizhongli/cdhit), which resulted in

387    3,805,584 clusters. False positives and hits to known RdRP proteins were removed by

388    comparing against the NCBI non-redundant (nr) protein database, the NCBI RefSeq protein

389    database and the virus RdRP database (Extended Data Fig. 2b). The remaining unknown

390    protein clusters were subject to viral RdRP domain search using a hidden Markov models

391    (HMMs) built from a manually reviewed profile of known RdRP clusters using the program

392    hmmscan v3.3.2 ($e = 10$, hits $\geq 1$)[52]. Clusters that contain more than one hmmscan hit were

393    subsequently aligned and inspected for the presence of conserved RdRP motifs. Finally, a

394    total of 713 novel RdRP clusters were retained as a result of our rigorous screening and

395    checking steps.

396    To further expand the RdRP collection based on the viruses newly discovered here, we

397    updated the RdRP protein database with the 713 novel RdRP clusters identified here and used

398    it to detect additional RdRP sequences from the original 144.6 million amino acid sequences

399    using the Diamond blastp and an $e$-value threshold of $1E$-3. The newly detected RdRPs were

400    again incorporated into the RdRP database for another round of detection. This process was

401    repeated for ten iterations. The resulting RdRP proteins (21,747,015 in total) were subjected

402    to the homology-based clustering, the removal of false positives, a HMMs-based search, and

403    manual alignment inspection as described above (Extended Data Fig. 2c, Extended Data Fig.

404    3b).

405    Finally, the remaining clusters were merged into superclades using a hierarchical method

406    employing the Girvan–Newman algorithm[53], with the edge betweenness determined based on

407    median $e$-value threshold of $1E$-3 for each pair of clusters (Extended Data Fig. 2d, Extended

408    Data Fig. 3c and 3d). Briefly, the merging of clusters used the following four steps: (i) the

409    betweenness of all edges (median $e$-value between clusters) in the network was calculated; (ii)

410    the edge(s) with the highest betweenness were removed; (iii) the betweenness of all edges

411    affected by the removal was recalculated; (iv) steps ii and iii were repeated until no edges

412    remained. All processes related to merging were performed using igraph package v1.3.5[54]

413    implemented in *R*.

414

415    **Virus verification**

416    To determine whether the newly discovered viral RdRPs belonged to RNA viruses rather than

417    organisms with DNA genomes, we performed two experiments. First, the 50 environmental

418    samples collected in this study were subject to simultaneous RNA and DNA extraction and

419    sequencing. The reads from the DNA sequencing results were mapped against the RdRP

420    sequences to verify that there was no DNA counterpart. Quality control of viral contigs was

421    performed using bbduk.sh (https://sourceforge.net/projects/bbmap/), and the mapping

422    analyses were performed by Bowtie2 v2.4.2[55] with the "end-to-end" setting. Similarly, from

423    our collection of SRA data, we also searched for those studies that performed both RNA and

424    DNA sequencing, and these data were used for mapping analyses to confirm that the viruses

425    discovered had *bona fide* RNA genomes.

426        In addition to read mapping, RT-PCR assays were performed to confirm that the detected

427    viral superclades were RNA organisms. Two pairs of validation primers were designed for

428    each of the representative RdRP sequences from 17 RNA viral superclades, gene sequences

429    from two DNA virus families (i.e., *Podoviridae* and *Siphoviridae*), and RT sequences

430    identified in this study, with a product length of 300-550 bp. For each of the samples, both the

431    reverse-transcribed RNA and the matching DNA underwent simultaneous PCR amplification,

432    and the amplification products were subject to electrophoresis using a 1% agarose gel with

433    GelRed dye, which was subsequently visualized under UV.

434

435    **Structural prediction and comparisons of viral RdRPs and homologous proteins**

436    Three-dimensional structures of newly identified viral RdRPs from diverse RNA viral

437    superclades were predicted from primary sequences using AlphaFold 2 v2.3[15] and visualized

438    using the PyMol software v2.5.4 (http://www.pymol.org/pymol). AlphaFold 2 prediction is a

439    relatively reliable source of structure information as the pLDDT socre of more than 2/3

440    residues it predicted are above 75%. The previously resolved or predicted structures of viral

441    RdRP, eukaryotic RdRP, eukaryotic DdRP and RT proteins were compared using the Super

442    algorithm[56]. Considering that the protein structures have similar molecular weights but

443    substantial variations in their conformations, the "number of aligned atoms after refinement"

444    option was employed to evaluate the similarity between each pair of proteins. Subsequently,

445    networkX (https://networkx.org/) was employed to construct a three-dimensional structure

446    diagram using the "edge-weighted spring embedded" approach, with results then mapped as a

447    scatter plot (depicted in the Fig. 2d). Simultaneously, we visualized four viral RdRP domain

448     proteins using PyMol.

449

450     **Annotation and characterization of virus genomes**

451     Potential open reading frames (ORFs) were predicted from newly identified virus genomes

452     based on two criteria: (i) the predicted amino acid sequences were longer than 200 amino

453     acids in length, and (ii) they were not completely nested within larger ORFs. The annotation

454     of non-RdRP ORFs was mainly based on comparisons of predicted proteins to hidden

455     Markov models (HMMs) collected from the Pfam database (https://pfam-legacy.xfam.org/)

456     using hmmscan implemented in HMMER[52]. For the remaining ORFs, the annotation was

457     carried out by blastp comparisons against the nr protein database with an *e*-value threshold of

458     1*E*-3.

459

460     **Analyses of virome diversity, evolution and ecology**

461     To reveal the diversity of the RNA viruses identified, we used an RdRP identity threshold of

462     90% to define new virus species. Abundance levels were subsequently estimated for every

463     virus species based on the number of non-rRNA reads per million (RPM) within each sample

464     (i.e. sequencing runs) mapped to viral sequences belonging to that species. Virus alpha

465     diversity (measured with the Shannon index) and overall abundance were subsequently

466     estimated and compared across different geographic locations and ecological types, namely;

467     soil, marine, freshwater, wetland, hot spring, salt marsh, and other types. "Marker virus

468     species" that were greatly enriched in certain ecological types were also identified based on

469     virus mapping results. The marker virus species were defined as present only in one

470     ecological subtype with RPM $\geq$ 1 and coverage $\geq$ 20%. To reveal the diversity and

471     evolutionary relationship of RNA viruses within a superclade, RdRP representatives of

472     overall diversity were first selected based on homology-based clustering. These

473     representatives were aligned using L-INS-I algorithm implemented in Mafft v7.475[57].

474     Phylogenetic analyses were performed based on the alignment using a maximum likelihood

475     algorithm, a LG amino acid substitution model, a Subtree Pruning and Regrafting (SPR)

476     branch swapping algorithm, and a Shimodaira–Hasegawa-like procedure implemented in the

477     Phyml program v3.1[58].

478

**Identification of CRISPR spacer hits**

480     A CRISPR-Cas spacer database was compiled from 65,703 genomes of bacteria and archaea

481     downloaded from the GTDB database (https://gtdb.ecogenomic.org/)[59] using a modified

482     version of the CRISPR Recognition Tool (CRT)[60]. This database was supplemented with an

483     additional 11.8 million precompiled CRISPR-Cas spacers obtained from the CrisprOpenDB

484     spacer database (http://crispr.genome.ulaval.ca)[61]. All spacers were queried for exact matches

485     against viral contigs using the BLASTn-short function implemented in the NCBI BLAST

486     v2.9.0+ package[62] with parameters "-evalue $1E$-10, -perc_identity 95, -dust no -word_size 7",

487     allowing only 0-1 mismatches across the entire length of the spacer to minimize the number

488     of false-positive hits.

489

**Data availability**

491     Raw sequence reads newly generated in this study are available at the NCBI Sequence Read

492     Archive (SRA) database under the BioProject accession PRJNA956286 and PRJNA956287

493     (Extended Data Table. 2). All virus sequence data produced in this study are publicly

494     available at http://47.93.21.181/, which includes all RNA virus contigs, RdRP CDS, RdRP

495     proteins, RdRP HMM profiles and phylogenetic tree files. Additionally, this website also

496     includes related data sets for model building and validation, and the trained model of

497     LucaProt.

498

**Code availability**

500     The original codes of ClstrSearch and LucaProt are stored at GitHub repository

501     (https://github.com/alibaba/LucaProt), and the link will be available upon acceptance of the

502     paper. Currently, the codes are provided for the review process only. Any additional

503     information required to reanalyze the data reported in this paper is available from the lead

504     contact upon request.

505

**Acknowledgements**

520

**Author contributions**

522    Conceptualization, X.H., Y.H., E.C.H., Z.-R.L. and M.S.; Methodology, X.H., Y.H., J.-S.E.,

523    J.L., Z.-R.L. and M.S.; Investigation, X.H., Y.H., P.F., S.Q.M., Z.X. and Q.-Y.G.; Writing –

524    Original Draft, X.H., Y.H., E.-C.H. and M.S.; Writing – Review and Editing, All authors.

525    Funding Acquisition, F.-M.H., Y.-L.S., D.-Y.G., Z.-R.L. and M.S.; Resources (sampling),

526    X.H., S.-Q.M., W.-W.C., J.-H.T., G.-Y.X., S.-J.L., Y.-Y.X., Y.-L.Z., F.-M.H., Y.-F.P., Z.-H.Y.

527    and C.H.; Resources (computational), S.Z., Z.-Y.Z. and Z.-R.L.; Supervision, Z.R.L. and M.S.

528

**Competing interests**

530    The authors declare no competing interests.

531

532  1      Shi, M. *et al.* Redefining the invertebrate RNA virosphere. *Nature* **540**, 539-543,
533          doi:10.1038/nature20167 (2016).
534  2      Shi, M. *et al.* The evolutionary history of vertebrate RNA viruses. *Nature* **556**, 197-202,
535          doi:10.1038/s41586-018-0012-7 (2018).
536  3      Rivarez, M. P. S. *et al.* In-depth study of tomato and weed viromes reveals undiscovered plant virus
537          diversity in an agroecosystem. *Microbiome* **11**, 60, doi:10.1186/s40168-023-01500-6 (2023).
538  4      Sutela, S. *et al.* The virome from a collection of endomycorrhizal fungi reveals new viral taxa with
539          unprecedented genome organization. *Virus Evolution* **6**, veaa076, doi:10.1093/ve/veaa076 (2020).
540  5      Zayed, A. A. *et al.* Cryptic and abundant marine viruses at the evolutionary origins of Earth's RNA
541          virome. *Science* **376**, 156-162, doi:10.1126/science.abm5847 (2022).
542  6      Chen, Y. M. *et al.* RNA viromes from terrestrial sites across China expand environmental viral diversity.
543          *Nat Microbiol* **7**, 1312-1323, doi:10.1038/s41564-022-01180-2 (2022).
544  7      Neri, U. *et al.* Expansion of the global RNA virome reveals diverse clades of bacteriophages. *Cell* **185**,
545          4023-4037 e4018, doi:10.1016/j.cell.2022.08.023 (2022).
546  8      Breitbart, M. & Rohwer, F. Here a virus, there a virus, everywhere the same virus? *Trends Microbiol* **13**,
547          278-284, doi:10.1016/j.tim.2005.04.003 (2005).
548  9      Youle, M., Haynes, M. & Rohwer, F. in *Viruses: Essential Agents of Life*    (ed Günther Witzany)
549          61-81 (Springer Netherlands, 2012).
550  10     Rost, B. Twilight zone of protein sequence alignments. *Protein Eng* **12**, 85-94 (1999).
551  11     Chen, J., Guo, M., Wang, X. & Liu, B. A comprehensive review and comparison of different
552          computational methods for protein remote homology detection. *Brief Bioinform* **19**, 231-244,
553          doi:10.1093/bib/bbw108 (2016).
554  12     McNutt, A. T. *et al.* GNINA 1.0: molecular docking with deep learning. *J Cheminform* **13**, 43,
555          doi:10.1186/s13321-021-00522-2 (2021).
556  13     Pham, T. H., Qiu, Y., Zeng, J., Xie, L. & Zhang, P. A deep learning framework for high-throughput
557          mechanism-driven phenotype compound screening and its application to COVID-19 drug repurposing.
558          *Nat Mach Intell* **3**, 247-257, doi:10.1038/s42256-020-00285-9 (2021).
559  14     Du, B.-X. *et al.* Compound-protein interaction prediction by deep learning: Databases, descriptors and
560          models. *Drug Discov Today* **27**, 1350-1366, doi:10.1016/j.drudis.2022.02.023 (2022).
561  15     Jumper, J. *et al.* Highly accurate protein structure prediction with AlphaFold. *Nature* **596**, 583-589,
562          doi:10.1038/s41586-021-03819-2 (2021).
563  16     Gligorijevic, V. *et al.* Structure-based protein function prediction using graph convolutional networks.
564          *Nat Commun* **12**, 3168, doi:10.1038/s41467-021-23303-9 (2021).
565  17     Xu, L., Magar, R. & Barati Farimani, A. Forecasting COVID-19 new cases using deep learning
566          methods. *Comput Biol Med* **144**, 105342, doi:10.1016/j.compbiomed.2022.105342 (2022).
567  18     Deng, L. Deep Learning: Methods and Applications. *Foundations and Trends® in Signal Processing* **7**,
568          197-387, doi:10.1561/2000000039 (2014).
569  19     Sarker, I. H. Deep Learning: A Comprehensive Overview on Techniques, Taxonomy, Applications and
570          Research Directions. *SN Comput Sci* **2**, 420, doi:10.1007/s42979-021-00815-1 (2021).
571  20     Shang, J. & Sun, Y. CHEER: HierarCHical taxonomic classification for viral mEtagEnomic data via
572          deep leaRning. *Methods* **189**, 95-103, doi:10.1016/j.ymeth.2020.05.018 (2021).
573  21     Sukhorukov, G. *et al.* VirHunter: A Deep Learning-Based Method for Detection of Novel RNA Viruses
574          in Plant Sequencing Data. *Front Bioinform* **2**, 867111, doi:10.3389/fbinf.2022.867111 (2022).
575  22     Miao, Y., Liu, F., Hou, T. & Liu, Y. Virtifier: a deep learning-based identifier for viral sequences from

576        metagenomes. *Bioinformatics* **38**, 1216-1222, doi:10.1093/bioinformatics/btab845 (2022).

577   23    Liu, F., Miao, Y., Liu, Y. & Hou, T. RNN-VirSeeker: A Deep Learning Method for Identification of Short Viral Sequences From Metagenomes. *IEEE/ACM Trans Comput Biol Bioinform* **19**, 1840-1849, doi:10.1109/TCBB.2020.3044575 (2022).

580   24    Kabir, A. & Shehu, A. GOProFormer: A Multi-Modal Transformer Method for Gene Ontology Protein Function Prediction. *Biomolecules* **12**, doi:10.3390/biom12111709 (2022).

582   25    Cao, Y. & Shen, Y. TALE: Transformer-based protein function Annotation with joint sequence–Label Embedding. *Bioinformatics* **37**, 2825-2833, doi:10.1093/bioinformatics/btab198 (2021).

584   26    Nambiar, A. *et al.* Transforming the Language of Life: Transformer Neural Networks for Protein Prediction Tasks. *bioRxiv* (2020).

586   27    Wolf, Y. I. *et al.* Origins and Evolution of the Global RNA Virome. *mBio* **9**, doi:10.1128/mBio.02329-18 (2018).

588   28    Koonin, E. V. *et al.* Global Organization and Proposed Megataxonomy of the Virus World. *Microbiol Mol Biol Rev* **84**, doi:10.1128/MMBR.00061-19 (2020).

590   29    Dongen, S. v. Graph clustering by flow simulation. (2000).

591   30    Enright, A. J., Van Dongen, S. & Ouzounis, C. A. An efficient algorithm for large-scale detection of protein families. *Nucleic Acids Res* **30**, 1575-1584 (2002).

593   31    Saberi, A., Gulyaeva, A. A., Brubacher, J. L., Newmark, P. A. & Gorbalenya, A. E. A planarian nidovirus expands the limits of RNA genome size. *PLoS Pathog* **14**, e1007314, doi:10.1371/journal.ppat.1007314 (2018).

596   32    Krishnamurthy, S. R. & Wang, D. Origins and challenges of viral dark matter. *Virus Res* **239**, 136-142, doi:10.1016/j.virusres.2017.02.002 (2017).

598   33    Cobbin, J. C., Charon, J., Harvey, E., Holmes, E. C. & Mahar, J. E. Current challenges to virus discovery by meta-transcriptomics. *Curr Opin Virol* **51**, 48-55, doi:10.1016/j.coviro.2021.09.007 (2021).

601   34    Monttinen, H. A. M., Ravantti, J. J. & Poranen, M. M. Structure Unveils Relationships between RNA Virus Polymerases. *Viruses* **13**, doi:10.3390/v13020313 (2021).

603   35    Vaswani, A. *et al.* Attention Is All You Need. *arXiv* (2017).

604   36    Avsec, Ž. *et al.* Effective gene expression prediction from sequence by integrating long-range interactions. *Nature Methods* **18**, 1196-1203, doi:10.1038/s41592-021-01252-x (2021).

606   37    Wu, R. *et al.* RNA Viruses Linked to Eukaryotic Hosts in Thawed Permafrost. *mSystems* **7**, e0058222, doi:10.1128/msystems.00582-22 (2022).

608   38    Charon, J., Murray, S. & Holmes, E. C. Revealing RNA virus diversity and evolution in unicellular algae transcriptomes. *Virus Evolution* **7**, veab070, doi:10.1093/ve/veab070 (2021).

610   39    Ibarbalz, F. M. *et al.* Global Trends in Marine Plankton Diversity across Kingdoms of Life. *Cell* **179**, 1084-1097 e1021, doi:10.1016/j.cell.2019.10.008 (2019).

612   40    Kalu, E. I., Reyes-Prieto, A. & Barbeau, M. A. Community dynamics of microbial eukaryotes in intertidal mudflats in the hypertidal Bay of Fundy. *ISME Communications* **3**, 21, doi:10.1038/s43705-023-00226-8 (2023).

615   41    Bollback, J. P. & Huelsenbeck, J. P. Phylogeny, genome evolution, and host specificity of single-stranded RNA bacteriophage (family Leviviridae). *J Mol Evol* **52**, 117-128 (2001).

617   42    Poranen, M. M., Mäntynen, S. & Ictv Report, C. ICTV Virus Taxonomy Profile: Cystoviridae. *J Gen Virol* **98**, 2423-2424, doi:10.1099/jgv.0.000928 (2017).

619   43    Callanan, J. *et al.* RNA Phage Biology in a Metagenomic Era. *Viruses* **10**, doi:10.3390/v10070386

620    (2018).

621  44  Sharp, P. M. & Simmonds, P. Evaluating the evidence for virus/host co-evolution. *Curr Opin Virol* **1**,
622      436-441, doi:10.1016/j.coviro.2011.10.018 (2011).

623  45  Li, D., Liu, C.-M., Luo, R., Sadakane, K. & Lam, T.-W. MEGAHIT: an ultra-fast single-node solution
624      for large and complex metagenomics assembly via succinct de Bruijn graph. *Bioinformatics (Oxford,*
625      *England)* **31**, 1674-1676, doi:10.1093/bioinformatics/btv033 (2015).

626  46  Lennox, M., Robertson, N. & Devereux, B. Expanding the Vocabulary of a Protein: Application of
627      Subword Algorithms to Protein Sequence Modelling. *Annu Int Conf IEEE Eng Med Biol Soc* **2020**,
628      2361-2367, doi:10.1109/EMBC44109.2020.9176380 (2020).

629  47  Baek, M. *et al.* Accurate prediction of protein structures and interactions using a three-track neural
630      network. *Science* **373**, 871-876, doi:10.1126/science.abj8754 (2021).

631  48  Lin, Z. *et al.* Evolutionary-scale prediction of atomic-level protein structure with a language model.
632      *Science* **379**, 1123-1130, doi:10.1126/science.ade2574 (2023).

633  49  Thomas N. Kipf, M. W. Semi-Supervised Classification with Graph Convolutional Networks.    (2017).

634  50  Yong He, C. W., Shun Zhang, Nan Li, Zhaorong Li ,    Zhenyu Zeng. KG-MTT-BERT: Knowledge
635      Graph Enhanced BERT for Multi-Type Medical Text Classification. *arXiv* (2022).

636  51  Buchfink, B. A.-O. X., Reuter, K. A.-O. & Drost, H. A.-O. X. Sensitive protein alignments at tree-of-
637      life scale using DIAMOND. *Nat Methods* (2021).

638  52  Potter, S. C. *et al.* HMMER web server: 2018 update. *Nucleic Acids Res* **46**, W200-W204,
639      doi:10.1093/nar/gky448 (2018).

640  53  Girvan, M. & Newman, M. E. J. Community structure in social and biological networks. *Proc Natl*
641      *Acad Sci U S A* **99**, 7821-7826 (2002).

642  54  Newman, M. E. J. & Girvan, M. Finding and evaluating community structure in networks. *Phys Rev E*
643      *Stat Nonlin Soft Matter Phys* **69**, 026113 (2004).

644  55  Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nature methods* **9**, 357-359,
645      doi:10.1038/nmeth.1923 (2012).

646  56  Hasegawa, H. & Holm, L. Advances and pitfalls of protein structural alignment. *Curr Opin Struct Biol*
647      **19**, 341-348, doi:10.1016/j.sbi.2009.04.003 (2009).

648  57  Katoh, K. & Standley, D. M. MAFFT multiple sequence alignment software version 7: improvements
649      in performance and usability. *Molecular biology and evolution* **30**, 772-780,
650      doi:10.1093/molbev/mst010 (2013).

651  58  Guindon, S. & Gascuel, O. A simple, fast, and accurate algorithm to estimate large phylogenies by
652      maximum likelihood. *Systematic biology* **52**, 696-704 (2003).

653  59  Donovan H Parks, M. C., Christian Rinke, Aaron J Mussig, Pierre-Alain Chaumeil, Philip Hugenholtz.
654      GTDB: an ongoing census of bacterial and archaeal diversity through a phylogenetically consistent,
655      rank normalizedand complete genome-based taxonomy. *Nucleic Acids Res* **50**, D785-D794,
656      doi:10.1093/nar/gkab776 (2021).

657  60  Bland, C. *et al.* CRISPR recognition tool (CRT): a tool for automatic detection of clustered regularly
658      interspaced palindromic repeats. *BMC Bioinformatics* **8**, 209 (2007).

659  61  Dion, M. B. *et al.* Streamlining CRISPR spacer-based bacterial host predictions to decipher the viral
660      dark matter. *Nucleic Acids Res* **49**, 3127-3138, doi:10.1093/nar/gkab133 (2021).

661  62  Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local alignment search tool.
662      *Journal of molecular biology* **215**, 403-410 (1990).

663

664     **Figure and Table Legends**

665

666     **Fig. 1. Global diverse RNA virosphere. a**, RNA virus discovery pipeline. The pathway for

667     sequence homolog-based virus discovery is highlighted in blue on the left, including the

668     clustering, expand and merge steps. The RdRP AI modeling pathway is highlighted in orange

669     on the right, including the modeling, clustering and merge steps. **b**, Number of viral

670     superclades discovered using two methods (left), and the detection accuracy of RdRP AI

671     modelling (right). **c**, Venn diagram shows the shared representative viral species between

672     available data from Zayed *et al*., Neri *et al*., Chen *et al*., and this study. The bar graph shows

673     the shared viral superclades between the four studies and the unique viral superclades

674     identified in this study. **d**, Diverse clusters of RNA viruses (dark colored small circle) and

675     RNA virus superclades (light colored large circle). The known viral clusters and superclades

676     are denoted in dark grey and light grey, respectively. The novel viral clusters discovered in

677     this study are denoted with dark blue and dark orange circles, while the novel viral

678     superclades are denoted with light blue and light orange circle. The size of the circle reflects

679     the size of the viral cluster.

680

681     **Fig. 2. Evaluation of authenticity of RNA viral superclades. a**, Distribution of BLAST

682     median *e*-value, HMM score and mean AI modeling probabilities of RNA virus superclade

683     grouping by the sensitivity of three methods, with the primary sequence-identified conserved

684     RdRP motif C of each superclade shown on the left. The known viral superclades show high

685     sensitivity for all three methods and are shown in grey. The novel superclades show declining

686     homology but the relative stable AI probability. **b**, The positive libraries and mean RPM (the

687     number of mapped reads per million non-rRNA reads) of representative viral superclades,

688     DNA viruses, RT, and cell organisms in 50 samples collected in this study. DNA libraries are

689     shown in purple and RNA libraries in yellow, the different groups of RNA viruses and DNA

690     organisms are shown in different colors, and red asterisks refer to those subsequently

691     validated by RT-PCR. **c**, RT-PCR results of first pairs of validation primers for representative

692     RdRP sequences from 17 RNA viral superclades, capsid sequences from two DNA viral

693     families (*Podoviridae* and *Siphoviridae*), and RT sequences. **d**, Three-dimensional (3D)

694    structure homology analysis of representative RdRPs from 180 viral superclades with Eu

695    DdRPs, Eu RdRPs, and RT. Each point stands for a representative structure. The distance

696    between different points represents structure similarity and the greater the distance, the lower

697    the structure similarity. Four RdRP domain structures of the AI-specific superclades are

698    displayed with the A, B and C motifs highlighted.

699

700    **Fig 3. Genomic features of viral superclades**. **a**, Size (the number of contigs) of all novel

701    viral superclades compared to 21 known superclades. **b**, Genome length of all novel viral

702    superclades compared to 21 known superclades. Centre lines in the box plots represent the

703    median bounds. **c**, Histogram of the genome size distribution of RNA viruses from known

704    and novel viral superclades. **d**, The distribution of annotated functional protein in each viral

705    superclade. **e**, Genome structure of representatives from six known superclades, 17 novel

706    superclades and eight AI-specific superclades. Grey stars represent reference virus genomes

707    of known superclades. Domains not commonly found in RNA viruses are shown in yellow

708    and are labeled above their corresponding positions. At the bottom, scale length in

709    nucleotides. Abbreviations: GOLGA2L5: golgin subfamily A member 2-like protein 5;

710    Pentaxin: pentaxin family; Tme5 EGF: thrombomodulin like fifth domain, EGF-like; Mg

711    trans NIPA: magnesium transporter NIPA; NUDIX, nucleoside diphosphate-X hydrolase;

712    RecX: RecX family; TssO: type VI secretion system, TssO; Securin: securin sister-chromatid

713    separation inhibitor; Rax2: cortical protein marker for cell polarity; Abhydrolase: alpha/beta-

714    hydrolase family; OmdA: bacteriocin-protection, YdeI or OmpD-Associated; Blt1 C: Get5

715    carboxyl domain; DnaJ: DnaJ domain; Trypan PARP: procyclic acidic repetitive protein

716    (PARP); SAM KSR: kinase suppressor RAS 1; CBD PlyG: PlyG cell wall binding domain;

717    LydB: LydA-holin antagonist; RelB: RelB antitoxin; T2SSE: type II/IV secretion system

718    protein; PARP regulatory: poly A polymerase regulatory subunit; Pheromone: fungal mating-

719    type pheromone; Y phosphatase2: tyrosine phosphatase family; PseudoU synth: RNA

720    pseudouridylate synthase; Glyco hydro 35: glycosyl hydrolases family 35; TIP: tuftelin

721    interacting protein.

722

723    **Fig 4. Phylogenetic diversity of 32 RNA viral superclades**. Each phylogenetic tree was

724    estimated using a maximum likelihood method based on the conserved RdRP domain. Within

725    each phylogenetic tree, the viruses newly identified here are shaded yellow, those reported

726    previously are shaded green and blue. The name of each superclade is shown on the top of

727    each phylogeny and the names of the families within each superclade are shown on right of

728    the tree. The proteins associated with bacterial hosts are denoted with different shapes on the

729    right side of the corresponding viral sequence. All trees are midpoint-rooted for clarity only,

730    and the scale bar indicates 0.5 amino acid substitutions per site.

731

732    **Fig 5. Ecological dynamics of the global RNA virome**. **a**, Global distribution of RNA

733    viruses identified in this study. Species of known virus superclades are shown in gray and

734    species from novel superclades are shown in magenta. Pie size reflects the number of viral

735    species. **b**, Rarefaction curve of all RNA viral species. Inset, Rarefaction curve of RNA viral

736    species at the ecotype level with colors indicating different ecotypes. **c**, Distribution of alpha

737    diversity, RPM, novel viral species and AI-specific species at different ecological subtypes

738    and colored by their ecotype. The ecological subtypes on the y-axis are ordered from the

739    highest to the lowest alpha diversity for each ecotype. **d**, Viral distribution patterns in

740    environmental and animal samples. The relative abundance of viruses in each library was

741    calculated and normalized by the number of mapped reads per million no-rRNA reads (RPM).

742    Viral species from 11 ecological subtypes are shown and divided into three groups, indicated

743    by the colors on the heatmap.

744

745    **Extended Data Fig. 1 Geographic coverage of the meta-transcriptomic data analyzed in**

746    **this study. a,** Geographical distribution of samples at the ecotype level. Pie size is positively

747    correlated to the number of samples. b, Total number of samples at different ecotypes.

748

749    **Extended Data Fig. 2. Detailed RNA virus discovery pipeline. a**, Schematic diagram of

750    homology-based discovery and RdRP AI modeling. **b**, Protein clustering process; only

751    clusters with more than ten members are retained for viral cluster discovery. **c**, Ten iterations

752    of RdRP expansion by recruiting newly detected RdRP in this process. **d**, RdRP clusters

753    merging into RdRP superclades using BLAST median *e*-value. **e**, RdRP identification by a

754    new deep learning model (i.e., LucaProt), includes five modules: Input, Tokenizer, Encoder,

755    Pooling, and Output.

756

757    **Extended Data Fig. 3. Benchmarking of the threshold at three processes (clustering,**

758    **expand and merge). a**, Number of hits using different *e*-values at the test stage. **b**,

759    Benchmarking of hmmscan bitscore and aligned fraction using the RdRP and non-RdRP data

760    sets (including RT, Eu DdRP and Ed RdRP derived from NCBI GenBank database). **c**,

761    BLAST Median *e*-value within the same known RdRP cluster. **d**, BLAST Median *e*-value

762    between pairwise comparisons of known RdRP clusters, with a 1*E*-3 cut-off used for cluster

763    merging.

764

765    **Extended Data Fig. 4. Benchmarking of the AI RdRP modeling. a,** The sigmoid function

766    of the AI modeling. **b,** Statistics of the data set for AI model building, including the entire

767    data set, training set, validation set, and testing set. **c,** The distribution of AI modeling

768    probabilities of positive data sets, **d,** The AI distribution of AI modeling probabilities of

769    negative data sets, including RT, Eu DdRP and Eu RdRP.

770

771    **Extended Data Fig. 5. Comparisons of RNA virus discovery results between three**

772    **previous studies and the current study. a,** The distribution of representative viral RdRPs of

773    four studies at the superclade level and the study-specific level. **b,** Venn diagram shows the

774    number of RdRP superclades found in each study and those shared between and among four

775    studies. **c,** Venn diagram shows the number of representative RdRPs found in each study and

776    those shared between and among four studies. **d,** Bar graph shows the number of known,

777    novel, AI-specific and study-specific RdRPs of four studies.

778

779    **Extended Data Fig. 6. The distribution of AI modeling probabilities of viral RdRPs. a,**

780    Distribution of AI modeling probabilities for all RdRPs from known viral superclades (first

781    left column) and representative RdRP superclades (right four columns). **b,** Distribution of AI

782    modeling probabilities for all RdRPs from novel viral superclades (first left column) and

783    representative RdRP superclades (right four columns) captured by BLAST, HMM and the

784    deep learning model. **c,** The distribution of AI modeling probabilities for all RdRPs from

785    novel viral superclades (first left column) and representative RdRP superclades (right four

786    columns) captured by both HMM and the deep learning model. **d**, Distribution of AI

787    modeling probabilities for all AI-specific RdRPs (first left column) and representative RdRP

788    superclades (right four columns) that could only be captured by the deep learning model.

789

790    **Extended Data Fig. 7. Expression difference of RNA viruses and DNA organisms in our**

791    **newly sequenced data. a,** Abundance comparisons between 58 RNA viral superclades, four

792    DNA virus families, RT and cell organisms at DNA and RNA libraries. **b**, RT-PCR results of

793    second pairs of validation primers for representative RdRP sequences from 17 RNA viral

794    superclades, capsid sequences from two DNA virus families (*Podoviridae* and *Siphoviridae*),

795    and RT sequences.

796

797    **Extended Data Fig. 8. Genome coverage of representative genome for RNA viruses and**

798    **DNA organisms in our newly sequenced data.** For 42 RNA viral superclades, four DNA

799    virus families, RT, and cell organism, genomes with high abundance in RNA libraries were

800    chosen to check reads coverage in DNA libraries.

801

802    **Extended Data Fig. 9. The coverage and abundance of RNA viruses and DNA organisms.**

803    The coverage of viral sequences shown as rising with rpm.

804

805    **Extended Data Fig. 10. Phylogenetic tree of the Nido-like superclade and the genome**

806    **structure of representatives.** The tree was estimated using a maximum likelihood method

807    based on the conserved RdRP domain. The reference sequences reported previously are

808    shaded grey, the viruses newly identified here are shaded by different colors according to

809    different ecotypes. The names of viral families are shown on right of the tree. The tree was

810    midpoint-rooted for clarity only, and the scale bar of tree indicates 0.2 amino acid

811    substitutions per site. The genome structures of representative viruses are shown on right of

812    the tree. At the bottom, scale indicates the length in nucleotides.

813

814     **Extended Data Fig. 11. Association between RNA viral superclades and other non-RdRP**

815     **protein clusters.** Grey and pink circles denote known and novel superclades, respectively.

816     Blue circles denote non-RdRP protein clusters.

817

818     **Extended Data Fig. 12. Specificity and shareability of RNA viruses. a,** Number of specific

819     viral species ("marker" species) in each ecological subtype. b, Association between RNA

820     viruses and different environmental ecotypes. The size of the colored circles indicates the

821     number of viral species identified by each ecotype, while the thickness of the line indicates

822     the number of viral species shared by each ecotype.

823

824     **Extended Data Fig. 13. Comparison of CHEER, VirHunter, Virtifier, RNN-VirSeeker**

825     **and LucaProt. a,** Positive rate of prediction results for CHEER, VirHunter, Virtifier, RNN-

826     VirSeeker and LucaProt based on the test data set. **b,** False positive rate of prediction results

827     for CHEER, VirHunter, Virtifier, RNN-VirSeeker and LucaProt based on the test data set. **c,**

828     Recall rate of prediction results for CHEER, VirHunter, Virtifier, RNN-VirSeeker and

829     LucaProt based on all RdRPs identified this study. **d,** Number of viral sequences of different

830     groups by contig length identified by CHEER, VirHunter, Virtifier, RNN-VirSeeker and

831     LucaProt. The training machines, training data sets, training strategies, and final model

832     selection of all comparison models are consistent with LucaProt. All comparison models were

833     built using multiple sets of hyperparameters with the best results selected for the comparison.

834

835     **Supplementary Table 1**. Detailed information of 10,437 meta-transcriptomics retrieved from

836     the SRA database.

837

838     **Supplementary Table 2**. Detailed information on the 50 environmental samples collected in

839     this study.

840

841     **Supplementary Table 3.** Information on all the RdRP sequences identified in this study.

842

843     **Supplementary Table 4.** Taxonomic comparison of known viral superclades between this
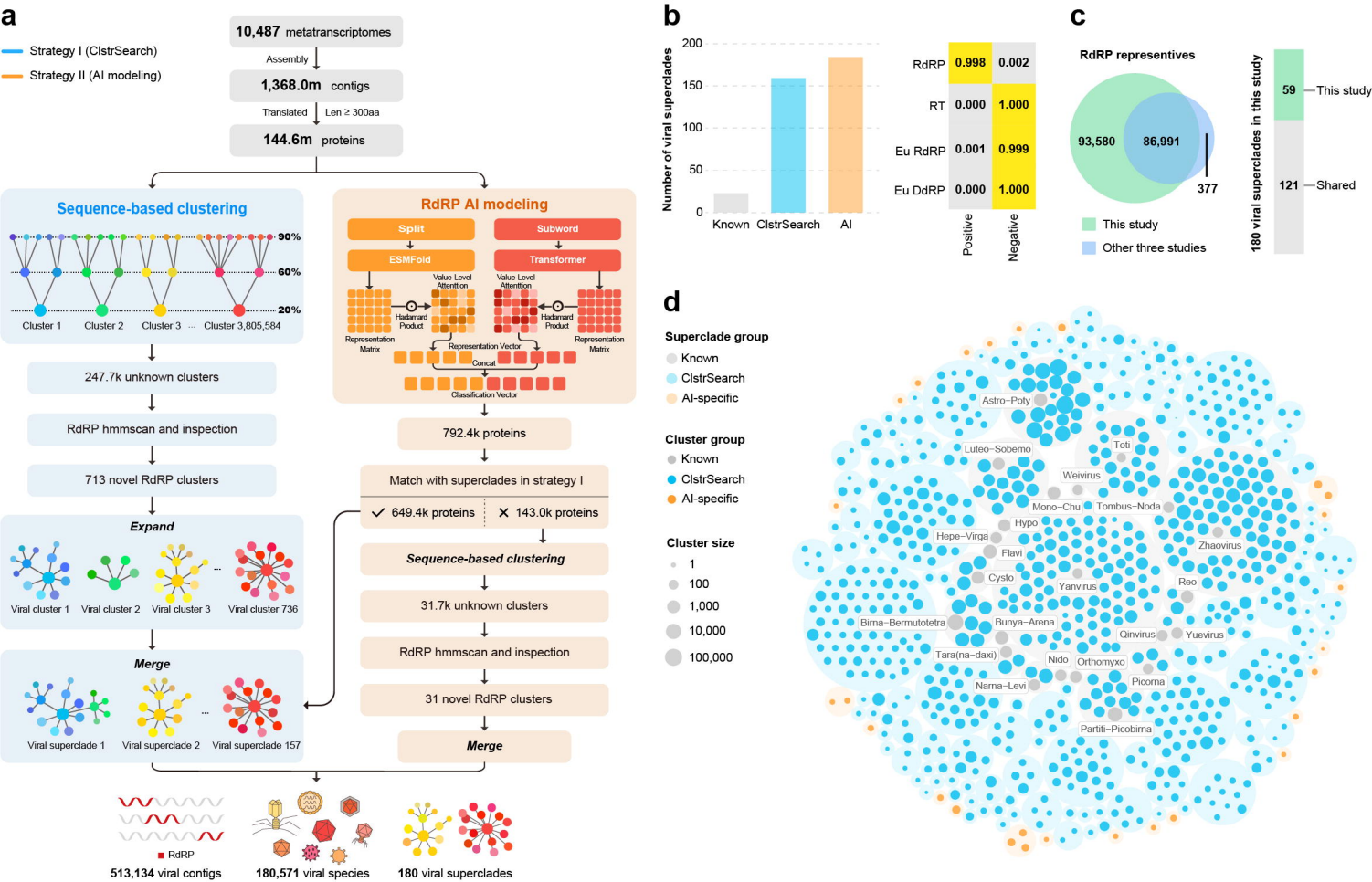
844    study and the current RNA virus classification system. New taxonomies that are incompatible

845    with current viral phyla or classes are shown in orange.

846

847    **Supplementary Table 5.** Predicted results of RdRPs identified in this study using three

848    methods (Threshold: BLAST: e<=1E-3; HMM: score>=10; AI: prob>=0.5).

849

850    **Supplementary Table 6.** Size information (the number of contigs) of all viral superclades.

851

852    **Supplementary Table 7.** Distribution of proteins associated with bacterial hosts in viral

853    superclades.

854

855    **Supplementary Table 8:** Normalized abundance levels (measured by RPM) of each viral

856    species in environmental samples (RPM>=1, coverage>=20%).

**a**

**b**

**c**

**d**

Domain

Superclade

Non-RdRP protein

Structural  Non-structural  Phage-related

Host-related

Superclade

Known  Novel  AI-specific

Hits

Yes  No

**e**

Known

Novel

AI-specific

0 — 5,000bp

Genome length (bp)

Protein domain

RdRP  Coat  Glycoprotein  Helicase  Methyltransferase

Peptidase  Phage coat  Phage integrase  Phage mat-A  Other

Astro-Poty, Flavi, Hepe-Virga, Birna-Permutotetra, Luteo-Sobemo, Tombus-Noda, Toti-Chryso, Picorna-Calici, Reo, Narna-Levi, Mono-Chu, Hypo, Partiti-Picobirna, Yanvirus, Bunya-Arena, Tara (na-daxi), Orthomyxo, Weivirus, Cysto, Qin-Yue, Nido, SC022, SC029, SC037, SC040, SC061, SC077, SC098, SC123, SC090, SC167, SC171

Legend:
- Refseq
- This study
- Partiti (Neri et al., 2022)
- ◆ CRISPR spacer
- ▲ Lysis
- ■ Prok-E2
- ● Prok-RING

**a** Legend: Know (grey), Novel (pink). Circle sizes: [1, 10), [10, 100), [100, 1,000), [1,000, 10,000), [10,000, )

**b** Overall

Number of vial species vs Number of samples

Inset: Ecotype level, Number of samples

Ecotype:
- Aquatic
- Host (microbe)
- Air
- Soil
- Host (plant)
- Invertebrate
- Extremity
- Food
- Vertebrate

**c** Columns: Alpha diversity | RPM (log10) | No. of species in novel SCs | No. of species in AI-specific SCs

Rows: Mine drainage, Marine, Groundwater, Aquatic, Pond, Freshwater, Lake water, Wastewater, Estuary, Halite, Subsurface, Aquifer, Surface, Sludge, Stromatolite, Soil crust, Compost, Freshwater sediment, Sediment, Activated sludge, Marine sediment, Soil, Peat, Wetland, Hot springs, Hydrothermal vent, Salt lake, Antarctic sediment, Hypersaline lake, Salt marsh, Biofilm, Bioreactor, Bioreactor sludge, Anaerobic digester, Microbial mat, Wood decay, Biogas fermenter, Rhizosphere, Leaf litter, Food, Food fermentation, Food production, Air, Invertebrate, Vertebrate

**d** Ecotype (subtype)

Ecotype (subtype):
- Marine
- Freshwater sediment
- Rhizosphere
- Invertebrate
- Freshwater
- Hydrothermal vent
- Leaf litter
- Vertebrate
- Marine sediment
- Salt marsh
- Air

RNA virus:
- Environment-specific
- Animal-specifc
- Share

RPM (log10): 0 1 2 3 4 5