

Predicting the potential for zoonotic transmission and host associations for novel viruses

Pranav Pandit (✉ pspandit@ucdavis.edu)

University of California Davis <https://orcid.org/0000-0001-7649-0649>

Simon Anthony

Columbia University

Tracey Goldstein

University of California, Davis <https://orcid.org/0000-0002-1672-7410>

Kevin Olival

EcoHealth Alliance <https://orcid.org/0000-0003-3211-1875>

Megan Doyle

University of California Davis

Nicole Gardner

University of California, Davis <https://orcid.org/0000-0002-4046-9685>

Brian Bird

University of California Davis

Woutrina Smith

University of California Davis

David Wolking

University of California Davis

Kristen Gilardi

University of California Davis

Corina Monagin

University of California Davis

Terra Kelly

University of California Davis

Marcela Uhart

University of California, Davis <https://orcid.org/0000-0003-3525-541X>

Jonathan Epstein

EcoHealth Alliance <https://orcid.org/0000-0002-1373-9301>

Catherine Machalaba

EcoHealth Alliance

Melinda Rostal

EcoHealth Alliance

Patrick Dawson

EcoHealth Alliance

Emily Hagan

EcoHealth Alliance

Ava Sullivan

EcoHealth Alliance

Hongying Li

Ecohealth Alliance

Aleksei Chmura

EcoHealth Alliance

Alice Latinne

EcoHealth Alliance

Christian Lange

Labyrinth Global Health, Inc

Tammie O'Rourke

Labyrinth Global Health, Inc

Sarah Olson

Wildlife Conservation Society <https://orcid.org/0000-0002-8484-9006>

Lucy Keatts

University of California Davis

A. Patricia Mendoza

Wildlife Conservation Society <https://orcid.org/0000-0002-8631-0119>

Alberto Perez

Wildlife Conservation Society

Catia Dejuste de Paula

Wildlife Conservation Society

Dawn Zimmerman

Smithsonian Conservation Biology Institute

Marc Valitutto

Smithsonian's National Zoological Park and Conservation Biology Institute

Matthew LeBreton

Mosaic (Environment, Health, Data, Technology)

David McIver

Metabiota Inc.

Ariful Islam

EcoHealth Alliance

Veasna Duong

Institut Pasteur du Cambodge <https://orcid.org/0000-0003-0353-1678>

Moctar Mouiche

Mosaic/Global Viral Cameroon

Zheng-Li Shi

Wuhan Institute of Virology <https://orcid.org/0000-0001-8089-163X>

Prime Mulembakani

University of Kinshasa

Charles Kumakamba

Metabiota Inc.

Mohamed Ali

Public Health Initiative, Cairo, Egypt.

Nigatu Kebede

Addis Ababa University

Ubald Tamoufe

Metabiota Inc.

Samuel Bel-Nono

Military Veterinarian (Rtd.)

Alpha Camara

Centre de Recherche en Virologie (VRV)

Joko Pamungkas**Julien Kalpy Coulibaly**

Institut Pasteur de Côte d'Ivoire

Ehab Abu-Basha

Jordan University of Science and Technology

Joseph Kamau

University of Nairobi

Soubanh Silithammavong

Metabiota Inc,

James Desmond

EcoHealth Alliance

Tom Hughes

Conservation Medicine <https://orcid.org/0000-0002-5713-9738>

Enkhtuvshin Shiilegdamba

Wildlife Conservation Society, Mongolia Program

Ohnmar Aung

Smithsonian's National Zoological Park and Conservation Biology Institute

Dibesh Karmacharya

Center for Molecular Dynamics

Julius Nziza

Mountain Gorilla Veterinary Project

Daouda Ndiaye

Université Cheikh Anta Diop

Aiah Gbakima

Metabiota, Inc

Zikankuba Sijali

Sokoine University of Agriculture

Supaporn Wacharapluesadee

Chulalongkorn University

Erika Alandia Robles

Wildlife Conservation Society (WCS) – Bolivia Program

Benard Ssebide

Mountain Gorilla Veterinary Project

Gerardo Suzán

Universidad Nacional Autónoma de México

Luis Aguirre

Universidad Mayor de San Simón <https://orcid.org/0000-0002-4194-1523>

Monica Solorio

Universidade Federal do Pará

Tapan Dhole

Sanjay Gandhi Post Graduate Institute of Medical Sciences

Peta Hitchens

Swedish University of Agricultural Sciences <https://orcid.org/0000-0002-7528-7056>

Damien Joly

Metabiota Inc

Karen Saylor

Labyrinth Global Health, Inc

Amanda Fine

Wildlife Conservation Society

Suzan Murray

Smithsonian's National Zoological Park and Conservation Biology Institute

William Karesh

EcoHealth Alliance

Peter Daszak

EcoHealth Alliance <https://orcid.org/0000-0002-2046-5695>

Jonna Mazet

UC Davis

PREDICT Consortium

University of California Davis

Christine Johnson

Karen C. Drayer Wildlife Health Center and EpiCenter for Disease Dynamics, One Health Institute, University of California Davis School of Veterinary Medicine

Article

Keywords:

Posted Date: January 25th, 2022

DOI: <https://doi.org/10.21203/rs.3.rs-846253/v1>

License:   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Version of Record: A version of this preprint was published at Communications Biology on August 19th, 2022. See the published version at <https://doi.org/10.1038/s42003-022-03797-9>.

Predicting the potential for zoonotic transmission and host associations for novel viruses

Authors: P. S. Pandit^{1*}, S.J. Anthony^{2**}, T. Goldstein^{1**}, K. J. Olival³, M. M. Doyle¹, N. R. Gardner¹, B. Bird¹, W. Smith¹, D. Wolking¹, K. Gilardi¹, C. Monagin¹, T. Kelly¹, M. Uhart¹, J. H. Epstein³, C. Machalaba³, M. K. Rostal³, P. Dawson³, E. Hagan³, A. Sullivan³, H. Li³, A. A. Chmura³, A. Latinne³, C. Lange⁴, T. O'Rourke⁴, S. Olson⁵, L. Keatts¹, P. Mendoza⁵, A. Perez⁵, C. Dejuste de Paula⁵, D. Zimmerman⁶, M. Valitutto⁶, M. LeBreton⁷, D. McIver⁸, A. Islam³, V. Duong⁹, M. Mouiche⁷, Z. Shi¹⁰, P. Mulembakani¹¹, C. Kumakamba¹², M. Ali¹³, N. Kebede¹⁴, U. Tamoufe¹⁵, S. Bel-Nono¹⁶, A. Camara¹⁷, J. Pamungkas^{18,19}, K. Coulibaly²⁰, E. Abu-Basha²¹, J. Kamau²², S. Silithammavong⁸, J. Desmond³, T. Hughes^{3,23}, E. Shiilegdamba²⁴, O. Aung⁶, D. Karmacharya²⁵, J. Nziza²⁶, D. Ndiaye²⁷, A. Gbakima²⁸, Z. Sijali²⁹, S. Wacharapluesadee³⁰, E. Alandia Robles³¹, B. Ssebide²⁶, G. Suzán³², L. F. Aguirre³³, M. R. Solorio³⁴, T. N. Dhole³⁵, P. L. Hitchens³⁶, D. O. Joly³⁷, K. Saylor⁴, A. Fine⁶, S. Murray⁷, W. Karesh³, P. Daszak³, J. A. K. Mazet¹, PREDICT Consortium, & C. K. Johnson^{1*}

*Corresponding author: P.S. Pandit, pspandit@ucdavis.edu and C. K. Johnson, ckjohnson@ucdavis.edu

Affiliations:

¹One Health Institute, School of Veterinary Medicine, University of California, Davis, Davis, CA, 95616, USA.

²Center for Infection and Immunity, Columbia University, New York, NY, 10032, USA.

³EcoHealth Alliance, 520 Eighth Avenue, New York, NY, 10018, USA.

⁴Labyrinth Global Health, Inc., 546 15th Ave NE, St Petersburg, FL 33704, USA.

⁵Wildlife Conservation Society, Bronx, NY, 10460, USA.

⁶Global Health Program, Smithsonian's National Zoological Park and Conservation Biology Institute, Washington, District of Columbia, United States of America

⁷Mosaic/Global Viral Cameroon, Yaoundé, Cameroon

⁸Metabiota Inc, Nanaimo, Canada

⁹Institut Pasteur du Cambodge, 5 Monivong Blvd, PO Box 983, Phnom Penh, 12201, Cambodia.

¹⁰CAS Key Laboratory of Special Pathogens, Wuhan Institute of Virology, Center for Biosafety Mega-Science, Chinese Academy of Sciences, Wuhan, China.

¹¹Kinshasa School of Public Health, University of Kinshasa, Kinshasa, Democratic Republic of the Congo.

¹²Metabiota Inc., Kinshasa, Democratic Republic of the Congo

¹³Egypt National Research Centre, 12311 Dokki, Giza, Egypt

¹⁴Addis Ababa University, Aklilu Lemma Institute of Pathobiology, P.O.Box 1176, Addis Ababa, Ethiopia

- 37 ¹⁵Metabiota Inc. Cameroon.
- 38 ¹⁶Military Veterinarian (Rtd.), P.O. Box CT2585, Accra, Ghana.
- 39 ¹⁷Centre de Recherche en Virologie (VRV) Projet Fievres Hemoraiques en Guinée; BP:5680
40 Nongo/Contéya-Commune de Ratoma, Guinea.
- 41 ¹⁸Primate Research Center, Bogor Agricultural University, Bogor 16151, Indonesia.
- 42 ¹⁹Faculty of Veterinary Medicine, Bogor Agricultural University, Darmaga Campus, Bogor
43 16680, Indonesia.
- 44 ²⁰Department Environment and Health, Institut Pasteur de Côte d’Ivoire, PO BOX 490 Abidjan
45 01, Ivory Coast.
- 46 ²¹Department of Basic Medical Veterinary Sciences, College of Veterinary Medicine, Jordan
47 University of Science and Technology, Jordan.
- 48 ²²Molecular Biology Laboratory, Institute of Primate Research, Nairobi, Kenya, Department of
49 Biochemistry, University of Nairobi, Nairobi, Kenya.
- 50 ²³Conservation Medicine, Selangor, Malaysia.
- 51 ²⁴Wildlife Conservation Society (WCS) - Vietnam Program, 1302, 57 Lang Ha, Hanoi, Vietnam
- 52 ²⁵Center for Molecular Dynamics Nepal (CMDN), Thapathali -11, Kathmandu, Nepal.
- 53 ²⁶Regional Headquarters, Mountain Gorilla Veterinary Project, Musanze, Rwanda.
- 54 ²⁷Université Cheikh Anta Diop, BP 5005, Dakar, Sénégal.
- 55 ²⁸Metabiota, Inc. Sierra Leone, Freetown, Sierra Leone.
- 56 ²⁹Department of Veterinary Medicine and Public Health, College of Veterinary Medicine and
57 Biomedical Sciences, Sokoine University of Agriculture, Morogoro, Tanzania.
- 58 ³⁰Thai Red Cross Emerging Infectious Diseases Health Science Centre, WHO Collaborating
59 Centre for Research and Training on Viral Zoonoses, King Chulalongkorn Memorial Hospital,
60 Faculty of Medicine, Chulalongkorn University, Bangkok, Thailand.
- 61 ³¹Wildlife Conservation Society (WCS) – Bolivia Program, c. Gabino Villanueva 340, La Paz,
62 Bolivia.
- 63 ³²Facultad de Medicina Veterinaria y Zootecnia, Universidad Nacional Autónoma de México,
64 México City, 04510, Mexico.
- 65 ³³Centro de Biodiversidad y Genética, Universidad Mayor de San Simón, Cochabamba, Bolivia.
- 66 ³⁴Laboratório de Epidemiologia e Geoprocessamento (EpiGeo), Instituto de Medicina
67 Veterinária (IMV) Universidade Federal do Pará (UFPA), BR-316 Km 31, Castanhal, PA 69746-
68 360, Brazil.
- 69 ³⁵Department of Microbiology, Sanjay Gandhi Post Graduate Institute of Medical Sciences, Uttar
70 Pradesh, India.
- 71 ³⁶Melbourne Veterinary School, Faculty of Veterinary and Agricultural Sciences, University of
72 Melbourne, Werribee, VIC 3030, Australia
- 73 ³⁷Nyati Health Consulting, 2175 Dodds Road, Nanaimo V9X0A4, Canada.

74

75 *Correspondence to: pspandit@ucdavis.edu, ckjohnson@ucdavis.edu

76 **Contributed equally to the manuscript

77 **Abstract:** Host-virus associations have co-evolved under ecological and evolutionary selection
78 pressures that shape cross-species transmission and spillover to humans. Observed virus-host
79 associations provide relevant context for newly discovered wildlife viruses to assess knowledge
80 gaps in host range and estimate pathways for potential human infection. Using models to predict
81 virus-host networks, we predicted the likelihood of humans as host for 513 newly discovered
82 viruses detected by large scale wildlife surveillance at high-risk animal-human interfaces in
83 Africa, Asia, and Latin America. Predictions indicated that novel coronaviruses are likely to
84 infect a greater number of host species than viruses from other families. Our models further
85 characterize novel viruses through prioritization scores and directly inform surveillance targets to
86 identify host ranges for newly discovered viruses.

87 **One Sentence Summary:** Potential host range and spillover risk for novel viruses can be
88 predicted using a network informed by known virus-host associations.

89 **Main**

90 Identifying zoonotic virus emergence events at the earliest possible stage is key to mitigating
91 outbreaks and preventing future epidemic and pandemic threats. By the time novel viruses are
92 recognized in humans, often within the context of a cluster of unusual cases, public health
93 interventions to prevent or contain an epidemic face major challenge. However, determining the
94 potential zoonotic transmission for newly discovered animal viruses, in the absence of
95 documented human infection, is currently a major scientific challenge. New approaches are
96 needed to evaluate and characterize risk of zoonotic transmission of newly discovered animal
97 viruses in the face of very limited data. Here we analyze human, domesticated animal, and wild
98 animal surveillance and viral discovery data collected from 2009-2019, as part of a consortium
99 led One Health project aimed at strengthening pandemic threat detection capabilities in Africa,
100 Asia, and Latin America¹. Surveillance efforts resulted in 944 novel monophyletic clusters of
101 virus sequences in wildlife (referred to as novel viruses henceforth) from 18 virus families
102 sampled at high-risk animal human disease transmission interfaces in 34 countries. As none of
103 these viruses have yet been identified in humans, other indices were established to assess
104 potential risk, including host range or plasticity of viruses and integration of virus and ecological
105 characteristics with expert opinion²⁻⁵. Using an analysis of the host-virus network we were able
106 to quantify risk of zoonotic transmission for 531 out of 944 novel animal viruses.

107 Patterns observed across host-virus networks have been used to understand virus sharing among
108 vertebrate species^{3,6,7}, and predict cryptic links between mammalian, and avian hosts and their
109 viruses⁸⁻¹⁰. Host-virus network linkages can be informed by virus traits, virus biogeography, host
110 ecological niches, and propensity for host sharing among viruses^{10,11}. Precedence in viral sharing
111 among species and ecological opportunities for spillover, as characterized by network topology,
112 can inform propensities for newly discovered viruses that lack data². Further exploration of these
113 networks can aid in estimating the host plasticity of viruses, an important characteristic
114 associated with zoonotic potential^{2,3}. Unfortunately, systematically collected surveillance data to
115 parameterize and validate these models have been missing⁴. Here, we apply a network approach
116 to gain ecological insights from viruses that have been shared among species in nature and
117 inform potential virus-host associations and zoonotic risk of novel viruses recently discovered
118 from in wildlife.

119 Using data from the literature, we developed a network that included 269 known zoonotic and
120 307 non-zoonotic viruses infecting 885 avian and mammalian hosts (G_c ; Fig. 1). The network was
121 used to train and validate two gradient boosting decision tree models to predict links and
122 taxonomic orders of missing links generated by sharing of hosts¹². Trained models were used to
123 predict possible host links for 531 novel viruses due to commonalities in host sharing with
124 known viruses and generated a predicted host-virus network ($G_{predicted}$, Fig. 1) formed due to
125 inclusion of novel viruses and their predicted linkages. We also predicted taxonomic order of the
126 probable host shared as a link between two virus nodes of the network and the likelihood of the
127 link to be humans, indicative of viruses' predicted potential to be zoonotic.

128 **Results and discussion**

129 **Virus-host network for known viruses (G_c):** We developed a unipartite network with viruses
130 as nodes and host species as an edge for all species recognized as a host for viruses based on data

131 presented in previous studies and databases, specifically, data shared by Olival et al.⁵, Pandit et
132 al.⁴, and Johnson et al.¹³ and GenBank. In the observed network (G_c), viruses were represented as
133 nodes and a link (edge) was generated if two viruses had been detected in the same host species.
134 The observed network (G_c) included 576 viruses as nodes and 35,838 edges (viruses linked
135 because of shared hosts) representing 352 vertebrate species (Fig.1). Exploration of network
136 characteristics of known viruses revealed differences in host sharing among virus families. The
137 distributions of centrality measures (Fig. 2a, 2b, 2e, 2i) for *Filoviridae*, *Flaviviridae*,
138 *Hantaviridae*, and *Orthomyxoviridae* families were statistically different from the mean
139 distribution (Kolmogorov-Smirnov, $p < 0.05$). Furthermore, after accounting for sampling bias
140 for individual viruses using PubMed hits, we ran a linear regression model with node-level
141 permutations (10,000 permutations to further characterize the distribution of viruses within virus
142 families in the network). Viruses in families *Hantaviridae*, *Filoviridae*, *Flaviviridae*, and
143 *Orthomyxoviridae* had a significantly higher degree ($p < 0.05$) and eigenvector centrality ($p <$
144 0.05), indicating more connections in the host-virus network than other represented virus
145 families. Viruses from the *Flaviviridae* family also had higher betweenness centrality ($p = 0.01$)
146 indicating more connections based on shared host species (Fig S2-S5). Results based on
147 distributions of centrality measures, as well as node level regression models, show similar
148 directionality for *Hantaviridae*, *Filoviridae*, *Flaviviridae*, and *Orthomyxoviridae* families across
149 multiple network topological metrics. Our findings provide further evidence for direct
150 relationship between higher host plasticity and greater zoonotic potential^{3,5}. Viruses from
151 *Nairoviridae* ($p = 0.01$) and *Rhabdoviridae* ($p = 0.01$) families (Fig S6) were significantly more
152 clustered together than viruses from other families.

153 The wildlife surveillance data consisted of tests for 99,375 animals, representing specimens from
154 861 species, mostly bats, rodents, primates, and other mammals
155 (<https://zenodo.org/record/5899054>)¹. To predict associations between novel viruses nodes
156 related to sharing common host species, gradient boosting models were trained using network
157 topological characteristics and families of viruses in the virus pairs to estimate: 1) whether virus
158 pairs have a species host in common; and 2) the taxonomical order of shared hosts (Fig. 1).

159 **Characteristics of predicted network ($G_{predicted}$) and newly discovered viruses:** The binary
160 model performed high performance in predicting the presence of links formed due to sharing of
161 hosts between two virus nodes in the network. The binary model performed well in predicting
162 sharing of viruses (mean positive predictive value = 0.99, sensitivity = 0.96, F-score 0.97, Fig.
163 S6) The distribution of predicted probability for all links using the binary model showed clear
164 bimodal distribution (Fig. S7a). The accuracy scores as a function of precision and recall
165 indicated good model performance beyond 0.15 predicted probability for the binary model (Fig.
166 S8). Hence, as a more conservative approach and to give weightage to the precision, we decided
167 to use 0.7 as an optimum threshold for detecting a positive link between two nodes (viruses). The
168 performance of the multilabel model varied for taxonomical orders, with higher moderate
169 performance for predicting taxonomical orders and groups of ‘humans’ and Cetartiodactyla (Fig
170 S7, Fig S9). For 531 novel viruses, we identified 184,055 possible links to new hosts (based on
171 optimum probability threshold of 0.7 identified for the binary model) to generate the predicted
172 network ($G_{predicted}$, Fig. 1, Fig S7a). For these predicted links, between two viruses, the
173 multiclass model was able to estimate potential taxonomic order of the shared species for
174 175,113 links. For the remaining links, the model was not able to confidently predict a specific
175 taxonomic order. Empirical biological networks are rarely scale-free (network with large hubs

176 and showing a power-law distribution for degree)¹⁴ but recent host-based host-virus unipartite
177 networks have shown scale-free nature where models with power-law distributions showed the
178 best fit for host-parasite networks¹⁵. Similarly, both observed (G_c) and predicted ($G_{predicted}$)
179 networks provided evidence that some viruses shared significantly larger numbers of hosts,
180 creating hubs of preferential attachment and showed weak evidence of scale-free nature
181 (loglikelihood ratio test $p > 0.05$). The predicted network ($G_{predicted}$) had longer tails at network
182 level (Kolmogorov-Smirnov, $p < 0.05$) as well as at virus family level for degree (Fig. 2a, e, f)
183 and betweenness centrality (Fig. 2 b, i, j) distributions than the observed network (G_c). Mean
184 network degree for all virus families reduced significantly with the addition of newly discovered
185 viruses that were predicted to have fewer links than known viruses, indicating lower host
186 plasticity for novel viruses than known viruses or insufficient adjustment of reporting bias (Fig
187 S10).

188 Based on a linear regression model with node-level permutations (10,000 permutations), our
189 adjustment for search effort (PubMed hits) was found to have no effect on the degree ($p = 0.38$,
190 Fig S11) and betweenness centrality ($p = 0.21$, Fig S12), but did significantly affect the
191 eigenvector ($p < 0.05$, Fig S13) and clustering coefficient ($p < 0.05$, Fig S14) of novel viruses.
192 These results indicate that sampling and reporting efforts affect our understanding of the
193 predilection towards certain species as illustrated by clustering in the network, but do not affect
194 the prediction of missing host links quantified by degree centrality within the network. Many of
195 the newly discovered viruses were mostly detected in only one species (mean = 1.32, $SD \pm 0.99$, n
196 = 944). Long tails of centrality distributions generated for the predicted network ($G_{predicted}$) and
197 comparatively lower centrality measures for novel viruses, when compared with known viruses,
198 support a tendency for newly discovered viruses to be more host-specific than previously
199 recognized viruses, a pattern that should be further evaluated with additional sampling effort to
200 identify the full host range for novel viruses.

201 Importantly, a comparison between virus families of novel viruses showed that novel
202 coronaviruses had higher degree ($p < 0.001$, Fig. 2C, Fig S11), betweenness ($p = 0.02$, Fig. 2D,
203 Fig S12), and eigenvector ($p < 0.001$) centralities in the predicted network compared to newly
204 discovered viruses in all sixteen other virus families (Fig. 2 C, D, G). In addition, the raw
205 detection data showed significantly higher host diversity for novel coronaviruses with a mean of
206 2.02 ($SD \pm 2.03$, $n = 114$) unique host species (maximum of 15 species) compared to 1.22 ($SD \pm$
207 0.70, $n = 834$) for other novel viruses detected in this study. This finding raises concern about the
208 ability of novel coronaviruses to infect a greater number of species than viruses from other
209 families. The recently emerged SARS-CoV-2 and the previously emerged SARS-CoV-1, have
210 shown a wide host breadth¹⁶. These predictions for novel coronaviruses highlight their key
211 ecological properties that can influence spillover into humans. Following coronaviruses, novel
212 flaviviruses showed significantly higher betweenness centrality ($p < 0.001$). Host taxonomic
213 order for novel viruses had no significant association with the degree centrality of the virus in the
214 predicted network. Predicted network characteristics not only differentiate virus families based
215 on network characteristics but also predict network characteristics that are key in understanding
216 the ecology of a novel virus and its behavior within the network community of hosts, including
217 the expected breadth of host species most likely to be infected by that novel virus.

218 **Prioritizing novel viruses for further characterization:** For the 531 newly detected viruses,
219 we developed prioritization metrics based on multiclass model predicted human links for known
220 viruses that inform on the ecological and evolutionary tendencies for spillover. Novel viruses

221 from *Herpesviridae*, *Rhabdoviridae*, *Coronaviridae*, *Adenoviridae*, *Astroviridae*, and
222 *Paramyxoviridae* families not only showed a high median probability of sharing human links
223 with known viruses (Fig S15) but also were predicted to have large numbers of human links in
224 the predicted network ($G_{predicted}$). Novel members of the *Picobirnaviridae* and *Rhabdoviridae*
225 families detected here have been speculated to be hyper-parasites infecting bacteria and insects
226 and were identified in mammalian host samples. Hence the predicted associations for these virus
227 families should not be inferred as infection but only as detection in host samples (e.g. potentially
228 insect viruses detected in oral swab samples from bats). Based on Generalized Linear Mixed
229 models, search effort (PubMed hits) was not associated with the predicted number of human
230 links ($p=0.24$, Table S1) nor the mean probability of sharing human links for novel viruses
231 ($p=0.778$, Table S2).

232 For relative comparison of zoonotic risk for the newly detected viruses, a prioritization metric
233 was developed based on the predicted probability of links being human and the number of shared
234 human links in the predicted network for a given virus. To understand the performance of the
235 prioritization score, we compared scores for known zoonotic and non-zoonotic viruses generated
236 by the ensemble of both binary and multi-class models. Results indicated significantly higher
237 prioritization scores for known zoonotic viruses (Fig S 16, $p < 0.001$) compared to known non-
238 zoonotic viruses. Prioritization scores were derived essentially from the prediction of new/yet
239 unobserved network links generated by the virus with another virus formed due to sharing of
240 hosts. However, models were unable to predict new links for well recognized that have
241 numerous hosts, such as Rabies virus and West Nile virus, and consequently resulted in a
242 prioritization score of zero. Fig. 3A-D shows the top ten and bottom five novel viruses from four
243 virus families for relative comparison based on the prioritization metric (Fig S17-23).
244 PREDICT_CoV-15 found in two *Phyllostomidae* bats from South America (*Artibeus lituratus*,
245 *Sturnira lilium*) scored the highest prioritization score in all novel viruses. Other top ten novel
246 coronaviruses based on the prioritization score included viruses detected in *Phyllostomidae* bats
247 (PREDICT_CoV-4, PREDICT_CoV-13, PREDICT_CoV-11, PREDICT_CoV-5). Out of these,
248 PREDICT_CoV-11 was also detected in *Mormoopidae* species (*Pteronotus personatus*) and
249 PREDICT_CoV-5 was found in *Vespertilionidae* species (*Bauerus dubiaquercus*) during the
250 surveillance. These also included coronaviruses detected in South-east Asian *Pteropodidae* bat
251 species such as PREDICT_CoV-16 and PREDICT_CoV-22. PREDICT_CoV-22 was also
252 detected in *Hipposideridae* bat species (*Hipposideros lekaguli*). PREDICT_CoV-78 detected in
253 multiple bat and rodent species of Southeast Asia also showed a high prioritization score. These
254 model outcomes, especially the prioritization score, provide a data driven tool to quantify
255 zoonotic risk for novel viruses. Even though the model is trained on numerous data points for
256 known zoonotic and non-zoonotic viruses, individual predictions for new virus discoveries
257 would only requires the data on hosts and virus families if used within our modeling framework.

258

259 **Prioritizing future surveillance:** The sharing of viruses among hosts is driven by geographical
260 overlap and synergies in ecological niches of hosts, as well as virus-specific characteristics that
261 enable cross-species transmission¹⁰. Novel viruses discovered in rodents, bats, primates, and
262 other mammalian hosts that were sampled from sites in close association with people, or at high-
263 risk interfaces that can facilitate disease transmission in urban and rural settings^{1,13}. Additional
264 surveillance across a broader taxonomic range is essential to gain additional insight on newly
265 detected viruses, further inform spillover risk, and improve model predictions presented here.

266 We used our network model and host taxonomic data in which the novel virus is first detected to
267 prioritize host species (surveillance targets) for further surveillance for newly discovered viruses
268 (Supplementary Data File 1). Moreover, given the recent SARS-CoV-2 pandemic we further
269 explored surveillance targets for novel coronaviruses. Novel coronaviruses were detected in bats,
270 rodents, birds, and primates (Fig. 4a). For novel coronaviruses, that were detected in bats,
271 predicted surveillance targets for bat coronaviruses showed three distinct clusters (Fig. 4b). The
272 first cluster of novel coronaviruses in bats had a higher proportion of predicted species from
273 *Miniopteridae* family (Bent-winged bats) but none from *Natalidae* (Neotropical funnel-eared
274 bats). Another prominent cluster prioritized all 11 chiropteran families, while the third cluster of
275 coronaviruses showed relatively fewer host recommendations from *Miniopteridae* bats.
276 Representation of these surveillance targets through these clusters highlights host predilection of
277 novel coronaviruses and indicates the preferential sharing of hosts by the novel coronaviruses.
278 These clusters also support earlier results related to the scale-free nature of the predicted network
279 ($G_{predicted}$) by creating virus hubs in the virus-host network. Cluster maps for other virus
280 families providing evidence for future surveillance are shown in Fig S24-S33 and supplementary
281 data file 1.

282 Grange et al developed a tool that ranks viruses for animal to human spillover using a risk-based
283 approach validated inputs by various experts from the field of virology, epidemiology and
284 ecology². Our approach, on the other hand, quantifies the risk of spillover agnostically and
285 informs predicted host range solely based on existing data available across the breadth of viruses
286 and natural infections observed in free-ranging mammalian and avian hosts. Although numerous
287 studies have been recently published that predict host-pathogen predictions, our framework
288 quantifies the risk for viruses that have been recently discovered in animal hosts. Network
289 models have shown to perform well with the inclusion of ecological trait data^{10,17} and genome
290 sequences¹⁸, but ,with the limited data available for novel viruses, the approach provided here is
291 an important step towards characterizing zoonotic potential for newly discovered animal viruses
292 in the face of sparse data. Our virus-centric approach (virus as nodes and edges as shared hosts)
293 showed improved performance over previous host-centric models¹⁷. Our network approach
294 presents some limitations specifically for viruses that have been detected in species with limited
295 surveillance effort to date and are thus not part of the training data. For this reason, we were able
296 to generate predictions for only 531 novel viruses out of 944. The remaining 413 novel viruses
297 without predictions were detected in species that were never found positive for any virus, starkly
298 indicating the lack of surveillance in wildlife. Further, model findings should be interpreted as
299 associations between hosts and viruses (based on detection of viruses in samples collected from
300 host species) with these associations requiring further to understand relationship between viruses
301 and hosts that might serve as reservoir, amplifying, or dead-end hosts. Detection of a virus in a
302 host species is not always correlated with that host's ability to produce viremia for further
303 transmission. Similarly, some of the novel viruses from *Picobirnaviridae* and *Rhabdoviridae*
304 have been speculated to be hyperparasites and the interpretation of these detections and predicted
305 host-associations need further investigations.

306 Novel viruses with high scores on the prioritization metrics present a strong eco-evolutionary
307 case for further genetic and *in-vivo* characterization to understand the risk of spillover. The
308 scoring will help streamline in-depth *in-vivo* characterization and develop additional hypotheses
309 related to genetic and ecological mechanisms for cross-species transmission and zoonotic
310 spillover. Nucleotide data associated with novel viruses presented here are short, hence the

311 current model framework of using only host associations provides a key advantage. However,
312 network models have shown to improve prediction capacities when nucleotide data is included as
313 features for prediction¹¹. These tools will improve with the as well as the discovery of new
314 viruses and further surveillance²⁰, ultimately informing our understanding of the mechanisms of
315 zoonotic emergence for viruses from wildlife.

316 **Methods**

317 **Data collection:** Virus-host data was collated from various sources. Major sources for the
318 association databases included data shared by Olival et al.⁵, Pandit et al.⁴, and Johnson et al.¹³. In
319 data provided by Olival et al (assessed September 2019), host-virus associations have been
320 assigned a score, based on detection methods and tests that are specific and more reliable. We
321 used associations that have been identified as the most reliable (stringent data) from Olival et al⁵.
322 In addition, a query in GenBank was run to parse out hosts reported for each GenBank
323 submission for viruses presented in each of these three databases. Initially, for each virus name,
324 taxonomic ID was identified using *entrez.esearch* function in biopython package. The taxonomic
325 ID helped identify ICTV lineage and associated data in PubMed. This included virus genus and
326 family information along with a standard virus name. Host data were aggregated based on the
327 taxonomic ID and associated standard name. Finally, for each virus, a search was completed in
328 PubMed to compile the number of hits related to the virus and their vertebrate hosts using the
329 search terms below. The number of PubMed hits (*PMHI*) were used as a proxy for sampling
330 bias^{4,13}. The virus-host association data source is presented in supplementary code and data files
331 (<https://zenodo.org/record/5899054>).

```
332     searchterm = (+virus_name  
333                 + [Title/Abstract]) AND (host OR hosts OR reservoir OR reservoirs OR  
334                 wild OR wildlife OR domestic OR animal OR animals OR  
335                 mammal OR bird OR birds OR aves OR avian OR avians  
336                 OR vertebrate OR vertebrates OR surveillance OR sylvatic)
```

337 Along with the PubMed terms we also queried the *nucleotide* database on PubMed using the
338 taxonomic ID to find the number of GenBank entries for these viruses (*PMH2*). A correlation
339 analysis between the *PMHI* and *PMH2* showed a high correlation with each other for us to
340 safely use GenBank hits for novel viruses during the prediction stage of the model (Fig S. 31).

341 **Development of G_c**

342 a. Centrality measures of observed network (G_c)

343 To test if centrality measures (degree centrality, betweenness centrality, eigenvector centrality,
344 clustering coefficient) for viral nodes in the observed network (G_c) vary significantly between
345 viral families, we firstly used the Kolmogorov-Smirnov (KS) test. KS test is routinely used to
346 identify distances between cumulative distribution functions of two probability distributions and
347 is largely used to compare degree distributions of networks^{21,22}. For each viral family,
348 distributions of centrality measures (degree centrality, betweenness centrality, and eigenvector
349 centrality) and clustering coefficient within the observed network (G_c) were compared with the

350 distribution of all nodes in the network using the two-tailed KS test. Secondly, a linear regression
351 model with virus family as a categorical variable and the number of PubMed hits as a covariate
352 to adjust for sampling bias were fitted to understand associations of viral families with centrality
353 measures.

354

$$355 \quad \text{centrality measure} = \beta_0 \text{intercept} + \beta_1 \text{Viral family}_{\text{categorical}} + \beta_2 \text{PubMed hits}$$

356

357 After fitting the model, node-level permutations were implemented. For each random
358 permutation, the output variable was randomly assigned to covariate values and the model was
359 re-fitted. Finally, a *p-value* was calculated by comparing the distribution of coefficients from
360 permutations with the original model coefficient.

361 **Network topology feature selection:** Using the observed network (G_c), multiple network
362 topological features for all node pairs were calculated. The following are topographical network
363 features calculated.

364 1. The Jaccard coefficient: a commonly used similarity metric between nodes in information
365 retrieval, is also called an intersection of over the union for two nodes in the network. In the
366 unipartite network generated here, it represents the proportion of common neighbor viruses from
367 the union of neighbor viruses for two nodes. Neighbor viruses are defined as viruses with which
368 the virus shares at least a single host. Higher Jaccard index represents similar host predilection.

369 2. Adamic/Adar (Frequency-Weighted Common Neighbors): Is the sum of inverse logarithmic
370 degree centrality of the neighbors shared by two nodes in the network²³. The concept of the
371 Adamic Adar index is a weighted common neighbors for viruses in the network. Within network
372 prediction, the index assumes that viruses with large neighborhoods have a less significant
373 impact while predicting a connection between two viruses compared with smaller
374 neighborhoods.

375 Both Jaccard and Adamic Adar coefficients have been routinely used for generalized network
376 prediction²⁴.

377 3. Resource allocation: Similarity score of two nodes defined by the weights of common
378 neighbors of two nodes. Resource allocation is another measure to quantify the closeness of two
379 nodes in the network and hence to understand the similarity of hosts they infect.

380 4. Preferential attachment coefficients: The mechanism of preferential attachment can be used to
381 generate evolving scale-free networks, where the probability that a new link is connected to node
382 x is proportional to k^{25} .

383 5. Betweenness centrality: For a node in the network betweenness centrality is the sum of the
384 fraction of all-pairs shortest paths that pass through it. The feature that we used for training the
385 supervised learning model was the absolute difference between of betweenness centralities of
386 two nodes. The difference between the betweenness centrality represents the difference in the
387 sharing observed by two viruses in the pair.

388 6. Degree centrality: The degree centrality for a node v is the fraction of nodes it is connected to.
389 The feature that we used for training the supervised learning model was the absolute difference
390 between degree centralities of two nodes. Unlike the difference in the betweenness centrality, the
391 difference in degree centrality only looks at the difference in the number of observed host
392 sharing.

393 7. Network clustering: All nodes were classified into community clusters using Louvain
394 methods²⁶. A binary feature variable was generated to describe if both the nodes in the pair were

395 part of the same cluster or not. If both viruses are from the same cluster, it represents similar host
396 predilection than when both viruses are not from the same cluster hence accounting for the
397 evolutionary predilection of viruses (or virus families) to infect a certain type of hosts.

398
399 Pearson's correlation coefficients were calculated to identify highly correlated features and for
400 choosing features for model training (Fig. S32). Virological features included in model training
401 were categorical variables describing the virus family of both the nodes in the pair, followed by a
402 binary variable if both the viruses belong to the same virus family. During the model
403 development, PubMed hits generated three predictive features for each pair of viruses on which
404 model training and predictions were conducted. These included two features representing
405 PubMed hits for the two viruses in the pair (PubMed_{v1} , PubMed_{v2}) and the absolute difference
406 between PubMed_{v1} and PubMed_{v2} to account for sampling bias differences between two viruses.

407 **Cross-validation and fitting generalized boosting machine (GBMs) models:** A nested-cross-
408 validation was implemented for the binary model while simple cross-validation was
409 implemented for the multiclass model (multiple output categories). The model parameters of the
410 binary model were first hyper-tuned using a cross-validated grid-search method. Values were
411 tested using a grid search to find the best-performing model parameters that showed the highest
412 sensitivity (recall). The parameters tested for hypertuning and their performance are provided in
413 the supplementary material (supplementary results and Table S5). For further cross-validation of
414 the overall binary model, all the viruses were randomly assigned to five groups. For each fold,
415 the viruses assigned to a group were dropped from the data, and a temporary training network
416 (G_t) was constructed, assuming that this represented the current observed status of the virus-host
417 community. For all possible pairs in G_o (both that sharing and not sharing any hosts) ten
418 topographical and viral characteristics were calculated as training features (Table S4).
419 Categorical features were one-hot-encoded and numeric features were scaled. An XGBClassifier
420 model with binary: logistic family was trained using the feature dataset to predict if virus pairs
421 share hosts (1,0 encoded output). The cross-validation was also used to determine the optimum
422 decision threshold for determining binary classification (Fig S17) and a precision-recall curve
423 was used to identify positive predictive value and sensitivity at the optimum threshold (Fig S8). -

424 The multiclass model was implemented in the same way, creating an observed network
425 (G_c) based on species-level sharing of hosts and randomly dropping viruses to generate a training
426 network (G_t) to train the XGboost model. The output variables were generated based on the
427 taxonomical orders of shared hosts. A pair of viruses can share multiple hosts, hence we trained a
428 multioutput-multiclass model. Humans were considered an independent category taxonomical
429 order (label) and were given a separate label than primates. For fine-tuning the multiclass model,
430 we started with the best performing parameters of the binary model and manually tested 5
431 combinations of model parameters by adjusting values of the learning rate, number of estimators,
432 maximum depth, and minimum child weight (Supplementary code and results).

433 **Missing links for novel viruses, binary and multiclass prediction:** The wildlife surveillance
434 data represented sampling of 99,379 animals (94,723 wildlife, 4,656 domesticated animals)
435 conducted in 34 countries around the world between 2009-2019 (Table S6)¹. Specimens were
436 tested using conventional Rt-PCR, Quantitative PCR, Sanger sequencing, and Next Generation
437 Sequencing protocols to detect viruses from 28 virus families or taxonomic groups (Table S7).
438 Testing resulted in 951 novel monophyletic clusters of virus sequences (referred to as novel

439 viruses henceforth). Within 951 novel viruses, 944 novel viruses had vertebrate hosts that were
440 identified with certainty based on barcoding methods and field identification. Host species
441 identification was confirmed by cytochrome b (cytb) DNA barcoding using DNA extracted from
442 the samples²⁷. We predicted the shared host links between novel viruses and known viruses
443 using binary and multiclass models in the following steps. Out of 944 novel viruses discovered in
444 the last ten years, we were able to generate predictions for 531 novel viruses that were detected
445 in species already classified as hosts within the network. The remaining 413 viruses were the
446 first detection of any virus in that species and thus host associations could not be informed by the
447 observed network (G_C) data.

- 448 1. A new node representing the novel virus was inserted in the network of the observed
449 network (G_C). Using the list of species in which the novel virus was detected, new edges were
450 created with known viruses that are also known to be found in those hosts. This generated a
451 temporary network for the novel virus (G_{temp}). If the novel virus was not able to generate any
452 edges with known viruses, meaning the host in which they have been found were never found
453 positive for any known virus, predictions were not performed.
- 454 2. Using G_{temp} feature values were calculated for the novel virus (betweenness centrality,
455 clustering, and degree). For all possible pairs of the novel virus with known viruses that are not
456 yet connected with each other through an edge in G_{temp} a feature dataset was generated (Jaccard
457 coefficient_(novel virus, known virus), the difference in betweenness centrality of the novel virus and
458 known virus, if the novel virus and known virus were in the same cluster, the difference in
459 degree centrality_(novel virus, known virus), if the novel virus and known virus were from same virus
460 family, the difference in PubMed hits_(novel virus, known virus), PubMed hits for the novel virus,
461 PubMed hits for the known virus). Studies and nucleotide sequences for novel viruses are
462 expected to be published and shared on PubMed's Nucleotide database and in various peer-
463 reviewed publications. Since, at the time of development of the model, data for all viruses was
464 not shared in a format that would reflect on PubMed's database, we decided to use the number of
465 times the virus was detected in the last ten years of wildlife surveillance. These detections will be
466 reflected in PubMed's Nucleotide database eventually, hence we considered them as a proxy for
467 search terms conducted for known viruses. Currently, evaluation of effects of this substitution of
468 PubMed hits with the number of detections for novel viruses is not possible with limited data on
469 novel viruses but needs to be reevaluated as more studies are published on these novel viruses.
- 470 3. Using this dataset for the novel virus, a binary presence of a link between the novel virus
471 and known viruses was predicted using the trained binary model. The taxonomic order of the
472 host link was predicted using the trained multiclass model.
- 473 4. For each possible link, the binary model predicted a probability of sharing link and the
474 multiclass model predicted multivariate outcomes of taxonomic orders and associated
475 probabilities. A threshold of 0.70 for the binary prediction model was used to classify if the link
476 is present or not and only those links were explored for their corresponding multiclass model
477 outputs.
- 478 5. The multiclass model showed higher performance for correctly classifying links as
479 "human" hosts than other numerous avian and mammalian taxonomic orders. Hence, the
480 multiclass model outputs were summarized into either humans or other taxonomic groups. For
481 the novel virus, a list of known viruses with the predicted link was generated. Using the hosts of
482 these known viruses and the taxonomic order in which the novel virus was detected, a list of
483 most likely species was generated based on the overall frequency of the host species. For
484 understanding the likelihood of infecting humans two factors were considered to be of

485 importance. Firstly, the number of links where humans are predicted as shared hosts with known
486 viruses (n) and the average model-predicted probability of those links. A representation was
487 generated incorporating the probability and available model support in terms of number links to
488 reflect the likelihood and compare viruses relative to each other.

489 To test if virus family, the taxonomic order of hosts in which novel viruses were detected,
490 and the number of times the viruses were detected (equivalent to PubMed hits for known viruses)
491 influenced node (virus) level network centrality measures in the predicted network (G_p) a linear
492 regression model was fitted with centrality measures.

493
494 *centrality measure*
495
$$= \beta_0 \text{intercept} + \beta_1 \text{Viral family}_{\text{categorical}} + \beta_2 \text{Host Order}_{\text{categorical}}$$

496
$$+ \beta_3 \text{PubMed hits}$$

497

498 For each of the random 10,000 node-level permutations, the output variable (centrality
499 measure) was randomly assigned to covariate values and the model was re-fitted. A p -value was
500 calculated by comparing the distributions of coefficients with the original model coefficient.
501 These models were fitted for degree centrality, betweenness centrality, eigenvector centrality,
502 and clustering coefficient of novel viruses in the predicted network.

503
504 **Prioritization score for novel viruses:** Generalized Linear Mixed Models were used to
505 understand the association effects of virus family, taxonomic order of the host and PubMed hits
506 on the number of predicted human links and mean probability of the predicted links. The models
507 were fit using *glmmTMB* and *glm* packages in R. For relative comparison of zoonotic risk and for
508 prioritizing novel viruses for further characterization, a prioritization metric was developed based
509 on the predicted probability of sharing the humans as hosts with known viruses ($p_{\text{sharing humans}}$)
510 and the number of predicted shared human links (n_{humans}) in the predicted network for the
511 given virus ($G_{\text{predicted}}$). Distributions for both $p_{\text{sharing humans}}$ and n_{humans} were normalized
512 and multiplied to generate a single score for a virus and for appropriate relative comparisons
513 between novel viruses. To understand the behavior of the prioritization score when predicting the
514 zoonotic risk of novel viruses, we also compared prioritization scores of known zoonotic and
515 non-zoonotic viruses using the Kolmogorov-Smirnov test.

516

517 **References:**

- 518 1 PREDICT Consortium. 2021. PREDICT Emerging Pandemic Threats Project. Dataset.
519 USAID Development Data Library. <https://data.usaid.gov/d/tqea-hwmmr..>
- 520 2 Grange, Z. L. *et al.* Ranking the risk of animal-to-human spillover for newly discovered
521 viruses. *Proceedings of the National Academy of Sciences* **118**, e2002324118,
522 doi:10.1073/pnas.2002324118 (2021).
- 523 3 Kreuder Johnson, C. *et al.* Spillover and pandemic properties of zoonotic viruses with
524 high host plasticity. *Sci Rep* **5**, 14830, doi:10.1038/srep14830 (2015).
- 525 4 Pandit, P. S. *et al.* Predicting wildlife reservoirs and global vulnerability to zoonotic
526 Flaviviruses. *Nat Commun* **9**, 5425, doi:10.1038/s41467-018-07896-2 (2018).
- 527 5 Olival, K. J. *et al.* Host and viral traits predict zoonotic spillover from mammals. *Nature*
528 **546**, 646-650, doi:10.1038/nature22975 (2017).

529 6 Gomez, J. M., Nunn, C. L. & Verdu, M. Centrality in primate-parasite networks reveals
530 the potential for the transmission of emerging infectious diseases to humans. *Proc Natl*
531 *Acad Sci U S A* **110**, 7738-7741, doi:10.1073/pnas.1220716110 (2013).

532 7 Albery, G. F. *et al.* The science of the host–virus network. *Nature Microbiology* **6**, 1483-
533 1492 (2021).

534 8 Walker, J. G., Plein, M., Morgan, E. R. & Vesk, P. A. Uncertain links in host-parasite
535 networks: lessons for parasite transmission in a multi-host system. *Philos Trans R Soc*
536 *Lond B Biol Sci* **372**, doi:10.1098/rstb.2016.0095 (2017).

537 9 Dallas, T., Park, A. W. & Drake, J. M. Predicting cryptic links in host-parasite networks.
538 *PLoS Comput Biol* **13**, e1005557, doi:10.1371/journal.pcbi.1005557 (2017).

539 10 Albery, G. F., Eskew, E. A., Ross, N. & Olival, K. J. Predicting the global mammalian
540 viral sharing network using phylogeography. *Nat Commun* **11**, 2260,
541 doi:10.1038/s41467-020-16153-4 (2020).

542 11 Wardeh, M., Blagrove, M. S., Sharkey, K. J. & Baylis, M. Divide-and-conquer: machine-
543 learning integrates mammalian and viral traits with network features to predict virus-
544 mammal associations. *Nature Communications* **12**, 1-15 (2021).

545 12 Chen, T. & Guestrin, C. in *Proceedings of the 22nd acm sigkdd international conference*
546 *on knowledge discovery and data mining*. 785-794 (ACM).

547 13 Johnson, C. K. *et al.* Global shifts in mammalian population trends reveal key predictors
548 of virus spillover risk. *Proc Biol Sci* **287**, 20192736, doi:10.1098/rspb.2019.2736 (2020).

549 14 Broido, A. D. & Clauset, A. Scale-free networks are rare. *Nat Commun* **10**, 1017,
550 doi:10.1038/s41467-019-08746-5 (2019).

551 15 Carlson, C. J., Zipfel, C. M., Garnier, R. & Bansal, S. Global estimates of mammalian
552 viral diversity accounting for host sharing. *Nat Ecol Evol* **3**, 1070-1075,
553 doi:10.1038/s41559-019-0910-6 (2019).

554 16 Banerjee, A., Mossman, K. & Baker, M. L. Zooanthroponotic potential of SARS-CoV-2
555 and implications of reintroduction into human populations. *Cell Host & Microbe* **29**, 160-
556 164 (2021).

557 17 Becker, D. J. *et al.* Optimizing predictive models to prioritize viral discovery in zoonotic
558 reservoirs. *bioRxiv*, 2020.2005. 2022.111344 (2021).

559 18 Mollentze, N., Babayan, S. & Streicker, D. Identifying and prioritizing potential human-
560 infecting viruses from their genome sequences. *bioRxiv*, 2020.2011. 2021.379917 (2021).

561 19 Mollentze, N. & Streicker, D. G. Viral zoonotic risk is homogenous among taxonomic
562 orders of mammalian and avian reservoir hosts. *Proceedings of the National Academy of*
563 *Sciences* **117**, 9423-9430 (2020).

564 20 Woolhouse, M., Scott, F., Hudson, Z., Howey, R. & Chase-Topping, M. Human viruses:
565 discovery and emergence. *Philos Trans R Soc Lond B Biol Sci* **367**, 2864-2871,
566 doi:10.1098/rstb.2011.0354 (2012).

567 21 Kossinets, G. & Watts, D. J. Empirical analysis of an evolving social network. *Science*
568 **311**, 88-90, doi:10.1126/science.1116869 (2006).

569 22 Muchnik, L. *et al.* Origins of power-law degree distribution in the heterogeneity of
570 human activity in social networks. *Sci Rep* **3**, 1783, doi:10.1038/srep01783 (2013).

571 23 Adamic, L. A. & Adar, E. Friends and neighbors on the web. *Social networks* **25**, 211-
572 230 (2003).

573 24 Lü, L. & Zhou, T. Link prediction in complex networks: A survey. *Physica A: statistical*
574 *mechanics and its applications* **390**, 1150-1170 (2011).

- 575 25 Barabasi, A. L. & Albert, R. Emergence of scaling in random networks. *Science* **286**,
576 509-512, doi:10.1126/science.286.5439.509 (1999).
577 26 Blondel, V. D., Guillaume, J.-L., Lambiotte, R. & Lefebvre, E. Fast unfolding of
578 communities in large networks. *Journal of statistical mechanics: theory and experiment*
579 **2008**, P10008 (2008).
580 27 Irwin, D. M., Kocher, T. D. & Wilson, A. C. Evolution of the cytochrome b gene of
581 mammals. *J Mol Evol* **32**, 128-144, doi:10.1007/BF02515385 (1991).

582

583 **Acknowledgments:**

584 **Funding:** This work was supported by the United States Agency for International Development
585 (USAID) Emerging Pandemic Threat PREDICT program (Cooperative Agreement nos. GHN-A-
586 00-09-00010-00 and AID-OAA-A-14-00102). . P.S.P. is also supported by R supported by the
587 National Institute Of Allergy And Infectious Diseases of the National Institutes of Health under
588 Award Number U01AI151814. The content is solely the responsibility of the authors and does
589 not necessarily represent the official views of the USAID, National Institutes of Health, or the
590 United States Government. We thank the governments of Bangladesh, Bolivia, Brazil,
591 Cambodia, Cameroon, China, DR Congo, Egypt, Ethiopia, Gabon, Ghana, Guinea, India,
592 Indonesia, Ivory Coast, Jordan, Kenya, Lao PDR, Liberia, Malaysia, Mexico, Mongolia,
593 Myanmar, Nepal, Peru, Republic of Congo, Rwanda, Senegal, Sierra Leone, South Sudan,
594 Tanzania, Thailand, Uganda, and Vietnam for permission to conduct this study, and the field
595 teams and collaborating laboratories that performed sample collection and testing.

596

597 **Author contributions:**

598

599 P.S.P. , C.K.J. S.J.A, T.G, K.L.O, and J.A.K.M conceived of the research; P.S.P analyzed the
600 data; P.S.P., C.K.J., S.J.A, T.G, K.L.O, J.A.K.M, M.M.D., N. R. G., B. B., W. S., D. W., K. G.,
601 C. M., T. K., M. U., J. H. E., C. M., M. K. R., P. D., E. H., A. S., H. L., A. A. C., A. L., C. L., T.
602 O'R., S. O., L. K., P. M., A. P., C. D. de P., D. Z., M. V., M. LB., D. MI., A. I., V. D., M. M.,
603 Z. S., P. M., M. A., N. K., U. T., S. B.N., A. C., J. P., K. C., E. A.B., J. K., S. S., J. D., T. H., E.
604 S., O. A., D. K., J. N., D. N., A. G., Z. S., S. W., E. A. R., B. S., G. S., L. F. A., M. R. S., T. N.
605 D., P. L. H., D. O. J., K. S., A. F., S. M., W. K., P. D., J., and PREDICT Consortium collected
606 data, wrote and revised the manuscript.

607

608 **Competing interests:** Author declare no competing interests.

609

610 **Data and materials availability:** Data and code reported in this paper are available at
611 <https://zenodo.org/record/5899054> and [https://data.usaid.gov/Global-Health-Security-in-](https://data.usaid.gov/Global-Health-Security-in-Development-GHSD-/PREDICT-Emerging-Pandemic-Threats-Project/tqea-hwmr)
612 [Development-GHSD-/PREDICT-Emerging-Pandemic-Threats-Project/tqea-hwmr](https://data.usaid.gov/Global-Health-Security-in-Development-GHSD-/PREDICT-Emerging-Pandemic-Threats-Project/tqea-hwmr)

613

614 **Figures and Tables**

615

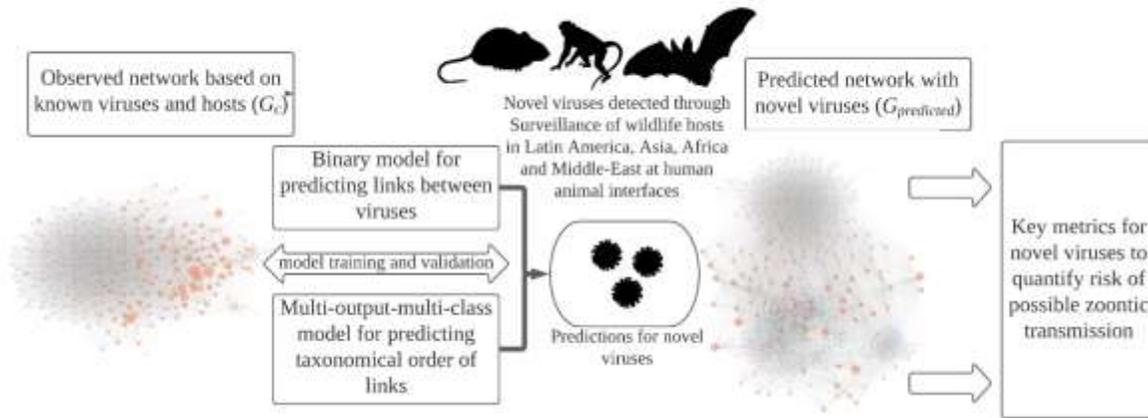
616 **Fig. 1. Modeling workflow:** The figure shows modeling procedure and methods implemented in the study. Orange dot
617 represents a known virus in the observed (G_c) and predicted networks ($G_{predicted}$), blue dots represent novel viruses in the
618 predicted network ($G_{predicted}$). Virus-host networks: G_c , represents a unipartite observed network of known zoonotic and non-
619 zoonotic viruses with nodes representing viruses and edges representing shared hosts. $G_{predicted}$ represents the predicted

620 unipartite network generated after predicting possible linkages between 531 novel viruses (white) and known viruses. The node
 621 size is proportional to the betweenness centrality.

622 **Fig. 2. Predicting missing links between virus-host communities.** Distribution shapes of degree (A) and
 623 betweenness centrality (B) for the observed and predicted network. Degree distributions for virus families in
 624 observed and predicted networks are shown in (E) and (F). Similarly, shapes of betweenness centrality for virus
 625 families in observed and predicted networks are shown in (I) and (J). Right panels show boxplots for novel virus
 626 families describing (C) degree, (D) betweenness centrality, (G) eigenvector centrality, and (H) clustering based on
 627 the predicted network formed by the binary prediction model.

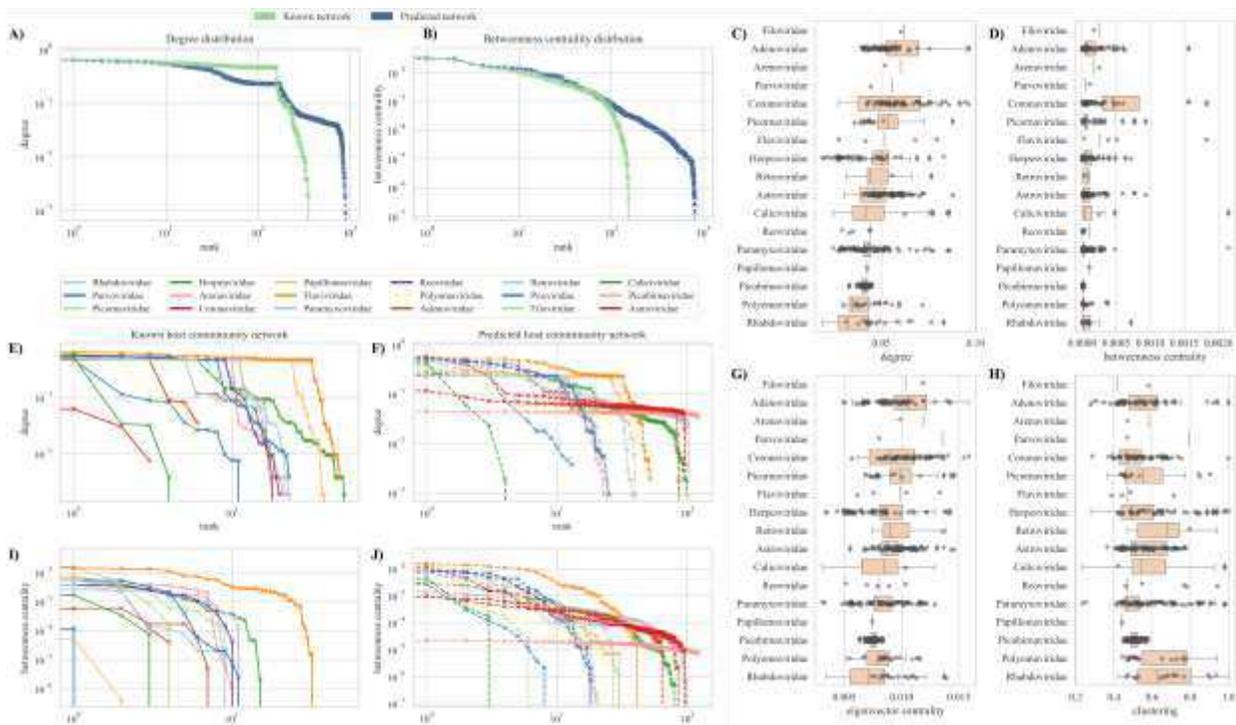
628 **Fig. 3: Prioritization metrics for novel viruses to understand zoonotic risk:** Top ten and bottom five newly
 629 discovered viruses from six virus families (A-F) with the virus prioritization scores based on multiclass model
 630 predictions. Annotations show the score and support represented by number of human links predicted.

631
 632 **Fig. 4: Surveillance targets for novel coronaviruses based on predicted sharing of hosts with known viruses.**
 633 Red color represents the evidence towards species in the taxonomic family (cumulative probability) with darker red
 634 red color indicating higher number of species occurrences from taxonomical families adjusted by model predicted
 635 probability. A) shows clustering of PREDICT coronaviruses by host, and B) focuses on coronaviruses found in bats.
 636 Clustering is based on the Bray-Curtis dissimilarity index.
 637



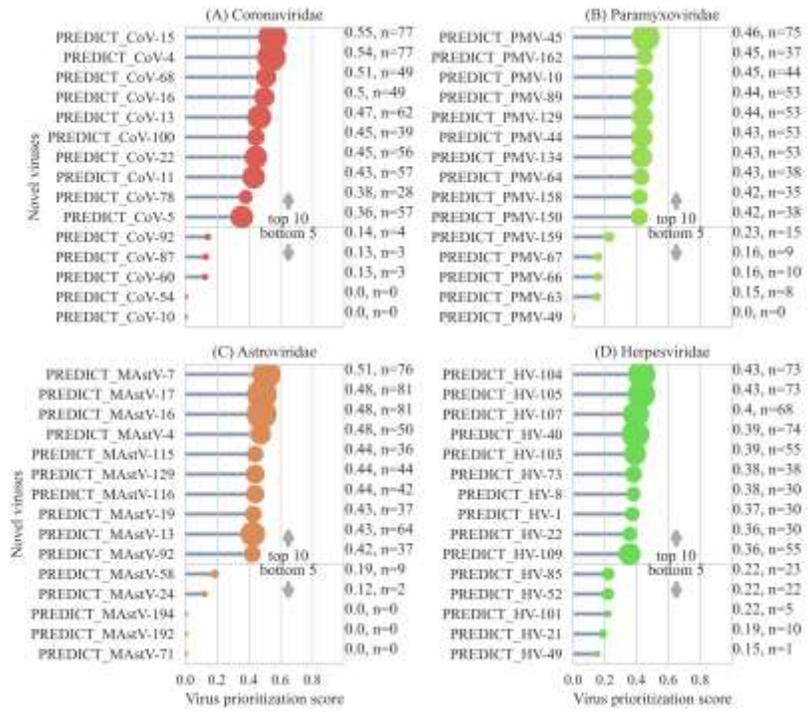
638 *Fig. 1: Model prediction workflow: The figure shows modeling procedure and methods implemented in the study. Orange dot*
 639 *represents a known virus in the observed (G_C) and predicted networks ($G_{predicted}$), blue dots represent novel viruses in the*
 640 *predicted network ($G_{predicted}$). Virus-host networks: G_C , represents a unipartite observed network of known zoonotic and non-*
 641 *zoonotic viruses with nodes representing viruses and edges representing shared hosts. $G_{predicted}$ represents the predicted*
 642 *unipartite network generated after predicting possible linkages between 531 novel viruses (white) and known viruses. The node*
 643 *size is proportional to the betweenness centrality.*

645

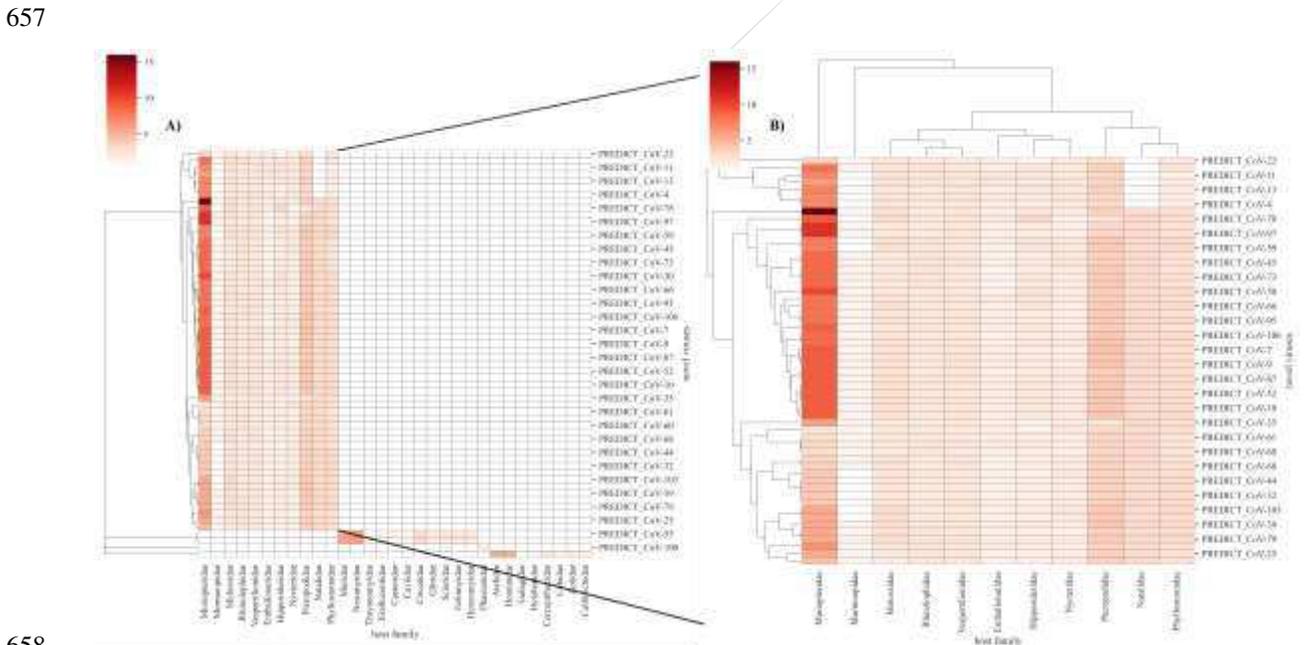


646
 647 *Fig. 2: Predicting missing links between virus-host communities. Distribution shapes of degree (A) and betweenness centrality*
 648 *(B) for the observed and predicted network. Degree distributions for virus families in observed and predicted networks are*
 649 *shown in (E) and (F). Similarly, shapes of betweenness centrality for virus families in observed and predicted networks are*
 650 *shown in (I) and (J). Right panels show boxplots for novel virus families describing (C) degree, (D) betweenness centrality, (G)*
 651 *eigenvector centrality, and (H) clustering based on the predicted network formed by the binary prediction model.*

652



653
 654 *Fig. 3: Prioritization metrics for novel viruses to understand zoonotic risk: Top ten and bottom five newly discovered viruses*
 655 *from six virus families (A-D) with the virus prioritization scores based on multiclass model predictions. Annotations show the*
 656 *score and support represented by number of human links predicted.*



658
 659 *Fig. 4: Surveillance targets for novel coronaviruses based on predicted sharing of hosts with known viruses. Red color represents*
 660 *the evidence towards species in the taxonomic family (cumulative probability) with darker red color indicating higher number of*
 661 *species occurrences from taxonomical families adjusted by model predicted probability. A) shows clustering of newly discovered*
 662 *coronaviruses by host, and B) focuses on coronaviruses found in bats. Clustering is based on the Bray-Curtis dissimilarity index.*

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [SupplementaryMaterialsPanditJohnsonetalr1.docx](#)
- [Supplementraydatafile1.xlsx](#)